

43362

Translations of
**MATHEMATICAL
MONOGRAPHS**

Volume 229

Lectures in
Mathematical
Statistics

Parts 1 and 2

Yu. N. Lin'kov



American Mathematical Society

Lectures in
Mathematical
Statistics

Parts 1 and 2

Translations of
**MATHEMATICAL
MONOGRAPHS**

Volume 229

Lectures in
Mathematical
Statistics

Parts 1 and 2

Yu. N. Lin'kov

Translated by Oleg Klesov and Vladimir Zayats



American Mathematical Society
Providence, Rhode Island

EDITORIAL COMMITTEE

AMS Subcommittee

Robert D. MacPherson Grigorii A. Margulis James D. Stasheff (Chair)

ASL Subcommittee Steffen Lemp (Chair)

IMS Subcommittee Mark I. Freidlin (Chair)

Ю. Н. Линьков

ЛЕКЦИИ ПО МАТЕМАТИЧЕСКОЙ СТАТИСТИКЕ

“ИСТОКИ”, ДОНЕЦК, 2001

This work was originally published in Russian by Istoki, Donetsk under the title “Лекции по математической статистике, части 1, 2” © Ю. Н. Линьков, 1999. The present translation was created under license for the American Mathematical Society and is published by permission.

Translated from the Russian by Oleg Klesov and Vladimir Zayats.

2000 *Mathematics Subject Classification*. Primary 62-01.

For additional information and updates on this book, visit
www.ams.org/bookpages/mmono-229

Library of Congress Cataloging-in-Publication Data

Lin'kov, Iu. N.

[Leksii po matematicheskoi statistike. English]

Lectures in mathematical statistics : parts 1 and 2 / Yu. N. Lin'kov ; translated by Oleg Klesov and Vladimir Zayats.

p. cm. — (Translations of mathematical monographs, ISSN 0065-9282 ; v. 229)

Includes bibliographical references and index.

ISBN 0-8218-3732-X (alk. paper)

1. Mathematical statistics. I. Title II. Series

QA276 .16 .L5513 2005

519 .5—dc22

2005052661

Copying and reprinting. Individual readers of this publication, and nonprofit libraries acting for them, are permitted to make fair use of the material, such as to copy a chapter for use in teaching or research. Permission is granted to quote brief passages from this publication in reviews, provided the customary acknowledgment of the source is given.

Republication, systematic copying, or multiple reproduction of any material in this publication is permitted only under license from the American Mathematical Society. Requests for such permission should be addressed to the Acquisitions Department, American Mathematical Society, 201 Charles Street, Providence, Rhode Island 02904-2294, USA. Requests can also be made by e-mail to reprint-permission@ams.org.

© 2005 by the American Mathematical Society. All rights reserved.

The American Mathematical Society retains all rights
except those granted to the United States Government.

Printed in the United States of America.

∞ The paper used in this book is acid-free and falls within the guidelines
established to ensure permanence and durability.

Visit the AMS home page at <http://www.ams.org/>

10 9 8 7 6 5 4 3 2 1 10 09 08 07 06 05

Contents

Foreword to the English Translation	vii
Part 1	1
Preface to Part 1	3
Chapter 1. Samples from One-Dimensional Distributions	5
1.1. Empirical distribution function and its asymptotic behavior	5
1.2. Sample characteristics and their properties	8
1.3. Order statistics and their properties	13
1.4. The distributions of some functions of Gaussian random vectors	20
Chapter 2. Samples from Multidimensional Distributions	25
2.1. Empirical distribution function, sampling moments, and their properties	25
2.2. Sampling regression and its properties	31
Chapter 3. Estimation of Unknown Parameters of Distributions	39
3.1. Statistical estimators and their quality measures	39
3.2. Estimation of a location parameter	49
3.3. Estimation of a scale parameter	56
3.4. The Cramér–Rao inequality and efficient estimators	61
3.5. The Cramér–Rao inequality for a multidimensional parameter	80
3.6. Integral inequalities of Cramér–Rao type	88
Chapter 4. Sufficient Statistics	99
4.1. Sufficient statistics and a theorem on factorization	99
4.2. Sufficient statistics and optimal estimators	113
Chapter 5. General Methods for Constructing Estimators	131
5.1. Method of moments	131
5.2. The maximum likelihood method	133
5.3. Bayes and minimax methods	142
5.4. Confidence intervals and regions	147
References to Part 1	153

Part 2	155
Preface to Part 2	157
Chapter 1. General Theory of Hypotheses Testing	159
1.1. Testing two simple hypotheses	159
1.2. Distinguishing a finite number of simple hypotheses	173
1.3. Distinguishing composite hypotheses	182
Chapter 2. Asymptotic Distinguishability of Simple Hypotheses	203
2.1. Statistical hypotheses and tests	203
2.2. Types of the asymptotic distinguishability of families of hypotheses. The characterization of types	205
2.3. Complete asymptotic distinguishability under the strong law of large numbers	218
2.4. Complete asymptotic distinguishability under the weak convergence	238
2.5. Contiguous families of hypotheses	248
Chapter 3. Goodness-of-Fit Tests	263
3.1. The setting of the problem. Kolmogorov test	263
3.2. The Pearson test	266
3.3. Smirnov test	275
3.4. Other goodness-of-fit tests	282
Chapter 4. Sequential Tests	293
4.1. Bayes sequential tests of hypotheses	293
4.2. Wald sequential tests	300
4.3. The optimality of a sequential Wald test	310
References to Part 2	317
Index	319

Foreword to the English Translation

Parts 1 and 2 of “Lectures in Mathematical Statistics” by Yu. N. Lin'kov were originally published in Russian as two separate books. For the English translation, the two parts are combined into one book. Each part has its own preface and list of references, with chapters, sections, theorems, etc., numbered independently in each part.

Part 1

Preface to Part 1

The author's idea was that this textbook should be aimed at students of mathematics having a background in general university courses in probability theory and mathematical statistics. The textbook was written based on the courses given by the author for students of mathematical departments at Volgograd University, Volgograd, Russia, and Donetsk University, Donetsk, Ukraine.

Among the books which may be used as a first reading in mathematical statistics, we mention the books by Cramér [9] and van der Waerden [34], which have already become the cornerstones in statistics. These books are still an authority, and many generations of experts have been brought up with these books. Elements of mathematical statistics are an essential ingredient of other general courses on probability theory. Let us mention the textbooks by Gnedenko [12], Gikhman, Skorokhod, and Yadrenko [11], Rozanov [27], Sevast'yanov [29], Tutubalin [33], and Shiryaev [30]. The textbooks by Shmetterer [31], Ivchenko and Medvedev [14], and Kozlov and Prokhorov [19] can be regarded as thoroughly developed introductions into mathematical statistics. The books on mathematical statistics by Borovkov [5], [6] take a special rank among textbooks for undergraduate and postgraduate students.

In writing this book, the author has used Russian and foreign literature on mathematical statistics, as well as the experience and traditions of teaching probability at Volgograd University and Donetsk University. Let us mention here the books by Rao [26], Cox and Hinkley [8], van der Waerden [34], and Bickel and Doksum [4] that thoroughly work out, each in its own way, problems for teaching mathematical statistics.

Part 1 of this book begins with a presentation of sampling using one-dimensional samples (Chapter 1) and multidimensional samples (Chapter 2) as an example. The basic sample characteristics are introduced and their asymptotic and nonasymptotic properties are studied. Main distributions related to the multidimensional Gaussian distribution are defined.

Chapter 3 deals with the estimation of parameters of distributions. In this chapter, measures of quality of statistical estimators are introduced and some optimality criteria are given. Optimal estimation of a scale parameter and a location parameter is studied. For regular families of distributions, approaches leading to effective estimators based on the Cramér–Rao inequality are given.

Chapter 4 deals with the theory of sufficient statistics and its applications to the construction of optimal estimators of unknown parameters and parametric functions.

In Chapter 5, general methods for constructing statistical estimators of parameters of distributions are considered and the main properties of the corresponding estimators are established.

The limited size of the book did not allow us to include some important statistical procedures or to consider other topics in the theory of parametric estimation. Part 2 of the textbook will deal with problems related to testing statistical hypotheses. The author hopes that this textbook will enable the reader to work independently, using other sources, on the topics we only touch upon here. We would recommend the books by Wilks [35] and Lehmann [21] and the three-volume monograph by Kendall and Stuart [16]–[18]. Our textbook can be used in preparation for general courses on mathematical statistics as well as specialized courses on the subject.

The list of references at the end of Part 1 includes only references available for students in Russia and Ukraine and is by no means complete.

In the textbook, we use the common notational conventions: P and P_θ for probabilities; E and E_θ for mathematical expectations; D and D_θ for variances, etc. We use triple notation for theorems, lemmas, formulas, etc. Therefore, for example, Theorem 4.1.2 refers to Theorem 2 in Section 1 of Chapter 4. Sections are enumerated by double numbers: Section 1.4 stands for Section 4 in Chapter 1. The sign \square marks the end of a proof.

Samples from One-Dimensional Distributions

1.1. Empirical distribution function and its asymptotic behavior

Empirical distribution function. Order statistics. Let ξ be a real-valued random variable with the distribution function

$$F(x) = P\{\xi < x\}, \quad x \in \mathbf{R} = (-\infty, \infty).$$

Let $\xi_1, \xi_2, \dots, \xi_n$ be n independent observations of the random variable ξ . Therefore $\xi_1, \xi_2, \dots, \xi_n$ are independent identically distributed random variables whose distribution function coincides with that of the random variable ξ , that is,

$$P\{\xi_i < x\} = F(x)$$

for all $i = 1, 2, \dots, n$. Denote by $\xi^{(n)} = (\xi_1, \xi_2, \dots, \xi_n)$ the vector of observations (also called a *sample*).

Given $x \in \mathbf{R}$, introduce the random variable

$$\nu_n(x) = \sum_{i=1}^n I_{(-\infty, x)}(\xi_i)$$

where $I_A(x)$ is the indicator of a set A . The function

$$(1.1.1) \quad F_n(x) = \nu_n(x)/n, \quad x \in \mathbf{R},$$

is called the *empirical distribution function*.

We rearrange the observations $\xi_1, \xi_2, \dots, \xi_n$ in ascending order and denote the resulting random variables by

$$(1.1.2) \quad \zeta_{n,1} \leq \zeta_{n,2} \leq \dots \leq \zeta_{n,n}.$$

The terms of this sequence are called *order statistics*.

Note that the empirical distribution function F_n possesses all the properties of regular distribution functions, namely it

- (1) assumes values in the interval $[0, 1]$,
- (2) does not decrease, and
- (3) is left-continuous.

Note also that F_n is a step function whose jumps are at the points $\zeta_{n,1}, \dots, \zeta_{n,n}$. If all the observations of a sample $\xi^{(n)}$ are different (in which case all the inequalities in (1.1.2) are strict), then $F_n(x)$ has n jumps whose heights are $1/n$. In the general case, equalities may appear in (1.1.2) (in which case the function $F_n(x)$ may have less than n jumps; however the jumps are proportional to $1/n$).

Glivenko's theorem. Definition (1.1.1) of the empirical distribution function implies that $F_n(x)$ is the frequency of the random event $\{\xi < x\}$ in n independent observations. Given x the probability of $\{\xi < x\}$ is constant and equals $F(x)$. By the Bernoulli theorem (the law of large numbers for Bernoulli trials) the empirical distribution function $F_n(x)$ tends in probability to $F(x)$ as $n \rightarrow \infty$, that is,

$$(1.1.3) \quad \lim_{n \rightarrow \infty} \mathbb{P}\{|F_n(x) - F(x)| > \varepsilon\} = 0 \quad \text{for all } \varepsilon > 0.$$

Moreover, by the Borel theorem (the strong law of large numbers for Bernoulli trials) $F_n(x)$ tends with probability 1 to $F(x)$ as $n \rightarrow \infty$, that is,

$$(1.1.4) \quad \mathbb{P}\left\{\lim_{n \rightarrow \infty} F_n(x) = F(x)\right\} = 1.$$

The convergence in relations (1.1.3) and (1.1.4) holds for every fixed $x \in \mathbf{R}$. However the following (stronger) Glivenko (1933) result claims that the convergence of $F_n(x)$ to $F(x)$ is, in fact, uniform with respect to x .

THEOREM 1.1.1 (Glivenko).

$$(1.1.5) \quad \mathbb{P}\left\{\lim_{n \rightarrow \infty} \sup_{x \in \mathbf{R}} |F_n(x) - F(x)| = 0\right\} = 1.$$

PROOF. Let $x_{r,k}$ be the minimal number x for which

$$(1.1.6) \quad F(x) \leq \frac{k}{r} \leq F(x+0)$$

where $r = 1, 2, \dots$ and $k = 0, 1, 2, \dots, r$. If the system of inequalities (1.1.6) does not have solutions for $k = 0$, then we put $x_{r,0} = -\infty$. Similarly, if (1.1.6) does not have solutions for $k = r$, then we put $x_{r,r} = \infty$. Consider random events

$$E_k^r = \left\{\lim_{n \rightarrow \infty} |F_n(x_{r,k}) - F(x_{r,k})| \vee |F_n(x_{r,k} + 0) - F(x_{r,k} + 0)| = 0\right\}$$

where $a \vee b$ stands for the maximum of two numbers a and b . We also put

$$E^r = \bigcap_{k=0}^r E_k^r.$$

It is clear that

$$E^r = \left\{\lim_{n \rightarrow \infty} \max_{0 \leq k \leq r} (|F_n(x_{r,k}) - F(x_{r,k})| \vee |F_n(x_{r,k} + 0) - F(x_{r,k} + 0)|) = 0\right\}.$$

We have $\mathbb{P}(E_k^r) = 1$ for all $k = 0, 1, \dots, r$ by the Borel theorem. Thus $\mathbb{P}(E^r) = 1$. Let $E = \bigcap_{r=1}^{\infty} E^r$. Since $\mathbb{P}(E^r) = 1$ for all $r \geq 1$, we get $\mathbb{P}(E) = 1$.

Now let k be such that $x_{r,k} < x_{r,k+1}$ and $x \in (x_{r,k}, x_{r,k+1}]$. Then

$$(1.1.7) \quad F_n(x_{r,k} + 0) \leq F_n(x) \leq F_n(x_{r,k+1}),$$

$$(1.1.8) \quad F(x_{r,k} + 0) \leq F(x) \leq F(x_{r,k+1}).$$

It is clear that

$$(1.1.9) \quad F(x_{r,k+1}) - F(x_{r,k} + 0) \leq \frac{1}{r}.$$

Inequalities (1.1.7) and (1.1.8) together with (1.1.9) yield

$$\begin{aligned}
 (1.1.10) \quad F_n(x) - F(x) &\leq F_n(x_{r,k+1}) - F(x_{r,k} + 0) \\
 &= [F_n(x_{r,k+1}) - F(x_{r,k+1})] + [F(x_{r,k+1}) - F(x_{r,k} + 0)] \\
 &\leq \max_{0 \leq k \leq r} (|F_n(x_{r,k}) - F(x_{r,k})| \vee |F_n(x_{r,k} + 0) - F(x_{r,k} + 0)|) + \frac{1}{r}.
 \end{aligned}$$

Using the same argument we derive from (1.1.7)–(1.1.9) that

$$\begin{aligned}
 (1.1.11) \quad F_n(x) - F(x) \\
 \geq - \max_{0 \leq k \leq r} (|F_n(x_{r,k}) - F(x_{r,k})| \vee |F_n(x_{r,k} + 0) - F(x_{r,k} + 0)|) - \frac{1}{r}.
 \end{aligned}$$

Combining (1.1.10) and (1.1.11) we obtain for all $x \in (x_{r,k}, x_{r,k+1}]$

$$\begin{aligned}
 (1.1.12) \quad |F_n(x) - F(x)| \\
 \leq \max_{0 \leq k \leq r} (|F_n(x_{r,k}) - F(x_{r,k})| \vee |F_n(x_{r,k} + 0) - F(x_{r,k} + 0)|) + \frac{1}{r}.
 \end{aligned}$$

Since the right-hand side of inequality (1.1.12) does not depend on k , it holds for all $x \in \mathbf{R}$. Thus

$$\begin{aligned}
 (1.1.13) \quad \sup_{x \in \mathbf{R}} |F_n(x) - F(x)| \\
 \leq \max_{0 \leq k \leq r} (|F_n(x_{r,k}) - F(x_{r,k})| \vee |F_n(x_{r,k} + 0) - F(x_{r,k} + 0)|) + \frac{1}{r}
 \end{aligned}$$

for all $r \geq 1$.

Since inequality (1.1.13) holds for all $r \geq 1$, we obtain

$$E \subset \left\{ \limsup_{n \rightarrow \infty} \sup_{x \in \mathbf{R}} |F_n(x) - F(x)| = 0 \right\},$$

whence (1.1.5) follows in view of $P(E) = 1$. □

Relations (1.1.3) and (1.1.4) as well as the Glivenko theorem indicate that the empirical distribution function $F_n(x)$ may serve as an approximation of the original distribution function $F(x)$.

Asymptotic normality of the empirical distribution function and Kolmogorov's theorem. According to definition (1.1.1) the empirical distribution function $F_n(x)$ for a fixed x is a random variable assuming values k/n for $k = 0, 1, 2, \dots, n$. Moreover

$$P\{F_n(x) = k/n\} = \binom{n}{k} F^k(x) (1 - F(x))^{n-k},$$

whence

$$EF_n(x) = F(x), \quad DF_n(x) = F(x)(1 - F(x))/n.$$

We say that a sequence of random variables η_n , $n = 1, 2, \dots$, is *asymptotically normal with parameters* (A_n, B_n^2) if

$$(1.1.14) \quad \lim_{n \rightarrow \infty} P \left\{ \frac{\eta_n - A_n}{B_n} < x \right\} = \Phi(x) \quad \text{for all } x \in \mathbf{R}$$

where

$$(1.1.15) \quad \Phi(x) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^x e^{-z^2/2} dz$$

is the distribution function of the standard normal law $\mathcal{N}(0, 1)$, that is, the normal law with mean 0 and variance 1.

By the central limit theorem for Bernoulli trials we obtain the following assertion on the asymptotic normality of the empirical distribution function.

THEOREM 1.1.2. *For every fixed $x \in \mathbf{R}$, the sequence of empirical distribution functions $F_n(x)$, $n = 1, 2, \dots$, is asymptotically normal with parameters*

$$F(x) \quad \text{and} \quad \frac{F(x)(1 - F(x))}{n}.$$

Consider the random variable

$$D_n = \sup_{x \in \mathbf{R}} |F_n(x) - F(x)|$$

measuring the deviation between the empirical distribution function $F_n(x)$ and the distribution function $F(x)$ in the uniform metric.

The following result by Kolmogorov (1933) allows one to estimate, for large n , that the probability D_n differs from zero.

THEOREM 1.1.3 (Kolmogorov). *If the distribution function $F(x)$ is continuous, then*

$$\lim_{n \rightarrow \infty} \mathbf{P} \{ \sqrt{n} D_n < z \} = K(z) = \sum_{j=-\infty}^{\infty} (-1)^j \exp \{ -2j^2 z^2 \}$$

for all $z > 0$.

The function $K(z)$ is called the *Kolmogorov distribution function*.

1.2. Sample characteristics and their properties

Sample moments. Let $\xi^{(n)} = (\xi_1, \dots, \xi_n)$ be a sample, that is, ξ_1, \dots, ξ_n are independent observations of a random variable ξ with the distribution function $F(x)$. Denote by α_k the k -th moment of the random variable ξ (in other words, the k -th moment of the distribution function $F(x)$), that is, $\alpha_k = E\xi^k$. By μ_k we denote the k -th central moment of the random variable ξ (in other words, the k -th central moment of the distribution function $F(x)$), that is, $\mu_k = E(\xi - \alpha_1)^k$. Note that α_1 is the expectation (or *mean value*) of the random variable ξ , while μ_2 is its variance. We also note that $\mu_1 = 0$ and $\mu_2 = \alpha_2 - \alpha_1^2$. Moreover, the moments and central moments are related to each other as follows:

$$(1.2.1) \quad \mu_k = \sum_{j=0}^k \binom{k}{j} (-1)^j \alpha_1^j \alpha_{k-j}.$$

Similar characteristics can be introduced for the empirical distribution function $F_n(x)$ constructed from the sample $\xi^{(n)} = (\xi_1, \xi_2, \dots, \xi_n)$. The k -th moment of the

empirical distribution function $F_n(x)$ is called the k -th *sampling moment*, that is,

$$(1.2.2) \quad a_k = \int x^k dF_n(x) = \frac{1}{n} \sum_{i=1}^n \xi_i^k.$$

The k -th central moment of the empirical distribution function $F_n(x)$ is called the k -th *sampling central moment*, that is,

$$(1.2.3) \quad m_k = \int (x - a_1)^k dF_n(x) = \frac{1}{n} \sum_{i=1}^n (\xi_i - a_1)^k.$$

From (1.2.2) and (1.2.3) we obtain a relation between sampling moments and sampling central moments:

$$(1.2.4) \quad m_k = \sum_{j=0}^k \binom{k}{j} (-1)^j a_1^j a_{k-j}.$$

Expectation and variance of sampling moments. It is clear that, for all k ,

$$(1.2.5) \quad E a_k = \frac{1}{n} \sum_{i=1}^n E \xi_i^k = \alpha_k,$$

$$(1.2.6) \quad D a_k = \frac{1}{n^2} \sum_{i=1}^n D \xi_i^k = \frac{1}{n} (\alpha_{2k} - \alpha_k^2)$$

if the corresponding moments exist. The evaluation of the expectation and variance of higher sampling central moments is a more complicated problem.

THEOREM 1.2.1. *If $\alpha_{2k} < \infty$, then*

$$(1.2.7) \quad E m_k = \mu_k + O\left(\frac{1}{n}\right),$$

$$(1.2.8) \quad D m_k = \frac{1}{n} (\mu_{2k} - 2k\mu_{k-1}\mu_{k+1} - \mu_k^2 + k^2\mu_2\mu_{k-1}^2) + O\left(\frac{1}{n^2}\right).$$

PROOF. Consider the random variables $\tilde{\xi}_i = \xi_i - a_1$, $i = 1, 2, \dots, n$, and put

$$\tilde{a}_j = \frac{1}{n} \sum_{i=1}^n \tilde{\xi}_i^j.$$

Note that $E\tilde{\xi}_i = 0$, $i = 1, 2, \dots, n$, $E\tilde{a}_1 = 0$, and $E\tilde{a}_j = \mu_j$. Applying (1.2.4) we get

$$(1.2.9) \quad m_k = \tilde{a}_k - k\tilde{a}_1\tilde{a}_{k-1} + \sum_{j=2}^k \binom{k}{j} (-1)^j \tilde{a}_1^j \tilde{a}_{k-j}.$$

Since the random variables $\tilde{\xi}_1, \tilde{\xi}_2, \dots, \tilde{\xi}_n$ are independent and $E\tilde{\xi}_i = 0$, we obtain

$$(1.2.10) \quad E\tilde{a}_1\tilde{a}_{k-1} = \frac{1}{n^2} \sum_{i,j=1}^n E\tilde{\xi}_i\tilde{\xi}_j^{k-1} = \frac{1}{n} \mu_k.$$

By the Cauchy–Bunyakovskiĭ inequality we get that

$$(1.2.11) \quad \left| \mathbb{E} \tilde{a}_1^j \tilde{a}_{k-j} \right| \leq \left(\mathbb{E} \tilde{a}_1^{2j} \mathbb{E} \tilde{a}_{k-j}^2 \right)^{1/2}$$

for $j \geq 2$. On the other hand,

$$\begin{aligned} \mathbb{E} \tilde{a}_{k-j}^2 &= \frac{1}{n^2} \sum_{i=1}^n \mathbb{E} \tilde{\xi}_i^{2(k-j)} + \frac{1}{n^2} \sum_{i \neq l} \mathbb{E} \tilde{\xi}_i^{k-j} \tilde{\xi}_l^{k-j} \\ &= \frac{1}{n} \mu_{2(k-j)} + \frac{n-1}{n} \mu_{k-j}^2 = \mu_{k-j}^2 + \frac{1}{n} (\mu_{2(k-j)} - \mu_{k-j}^2), \end{aligned}$$

whence $\mathbb{E} \tilde{a}_{k-j}^2 \leq \mu_{2(k-j)}$ in view of $\mu_{2(k-j)} - \mu_{k-j}^2 \geq 0$. Thus inequality (1.2.11) can be rewritten for an arbitrary $j \geq 2$ as follows:

$$(1.2.12) \quad \left| \mathbb{E} \tilde{a}_1^j \tilde{a}_{k-j} \right| \leq \left(\mu_{2(k-j)} \mathbb{E} \tilde{a}_1^{2j} \right)^{1/2}.$$

Further,

$$(1.2.13) \quad \mathbb{E} \tilde{a}_1^{2j} = \frac{1}{n^{2j}} \sum_{i_1=1}^n \sum_{i_2=1}^n \cdots \sum_{i_{2j}=1}^n \mathbb{E} \tilde{\xi}_{i_1} \tilde{\xi}_{i_2} \cdots \tilde{\xi}_{i_{2j}}.$$

Consider the terms on the right-hand side of (1.2.13) containing a factor $\tilde{\xi}_{i_l}$ whose index i_l differs from those of other factors. All such terms vanish because $\mathbb{E} \tilde{\xi}_i = 0$ and the random variables $\tilde{\xi}_1, \tilde{\xi}_2, \dots, \tilde{\xi}_n$ are independent.

Now consider the terms on the right-hand side of (1.2.13) whose indices i_1, i_2, \dots, i_{2j} fall into j pairs of equal numbers. The number of elements of this set is $N_1 N_2$ where N_1 is the number of ways to split the set $\{1, 2, \dots, 2j\}$ into j pairs and N_2 is the number of ways to choose different j numbers from the set $\{1, 2, \dots, n\}$. Obviously

$$N_1 N_2 = \prod_{k=0}^j (2j - 2k + 1) \cdot \prod_{l=0}^{j-1} (n - l) = O(n^j), \quad n \rightarrow \infty.$$

Note also that the cardinality of the set of all other terms in (1.2.13) can also be represented as a polynomial of n whose degree is less than j . Thus we obtain from (1.2.13) that

$$(1.2.14) \quad \mathbb{E} \tilde{a}_1^{2j} = O(n^{-j})$$

for all $j \geq 2$. Combining (1.2.9), (1.2.10), (1.2.12), and (1.2.14) and taking into account that $\mathbb{E} \tilde{a}_k = \mu_k$, we prove (1.2.7).

Equality (1.2.7) implies that

$$(1.2.15) \quad Dm_k = \mathbb{E}(m_k - \mu_k)^2 + O(n^{-2}).$$

It follows from (1.2.9) that

$$(1.2.16) \quad \begin{aligned} \mathbb{E}(m_k - \mu_k)^2 &= \mathbb{E}(\tilde{a}_k - \mu_k)^2 + 2 \sum_{j=1}^k \binom{k}{j} (-1)^j \mathbb{E} \tilde{a}_1^j \tilde{a}_{k-j} (\tilde{a}_k - \mu_k) \\ &\quad + \sum_{i,j=1}^k \binom{k}{i} \binom{k}{j} (-1)^{i+j} \mathbb{E} \tilde{a}_1^{i+j} \tilde{a}_{k-i} \tilde{a}_{k-j}. \end{aligned}$$

Applying (1.2.6) to $\tilde{\xi}_1, \dots, \tilde{\xi}_n$ we obtain

$$(1.2.17) \quad E(\tilde{a}_k - \mu_k)^2 = D\tilde{a}_k = n^{-1} (\mu_{2k} - \mu_k^2),$$

since $E\tilde{a}_k = \mu_k$. Using the same arguments as those applied for the evaluation of the right-hand side of (1.2.13) we prove that

$$(1.2.18) \quad \begin{aligned} E\tilde{a}_1\tilde{a}_{k-1}(\tilde{a}_k - \mu_k) &= \frac{1}{n^3} \sum_{i_1=1}^n \sum_{i_2=1}^n \sum_{i_3=1}^n E\tilde{\xi}_{i_1}\tilde{\xi}_{i_2}^{k-1}(\tilde{\xi}_{i_3} - \mu_k) \\ &= \frac{1}{n} \mu_{k-1}\mu_{k+1} + O\left(\frac{1}{n^2}\right), \end{aligned}$$

$$(1.2.19) \quad \begin{aligned} E\tilde{a}_1^2\tilde{a}_{k-1}^2 &= \frac{1}{n^4} \sum_{i_1=1}^n \sum_{i_2=1}^n \sum_{i_3=1}^n \sum_{i_4=1}^n E\tilde{\xi}_{i_1}\tilde{\xi}_{i_2}\tilde{\xi}_{i_3}^{k-1}\tilde{\xi}_{i_4}^{k-1} \\ &= \frac{1}{n} \mu_2\mu_{k-1}^2 + O\left(\frac{1}{n^2}\right). \end{aligned}$$

The same method shows that

$$(1.2.20) \quad \sum_{j=2}^k \binom{k}{j} (-1)^j E\tilde{a}_1^j \tilde{a}_{k-j} (\tilde{a}_k - \mu_k) = O\left(\frac{1}{n^2}\right),$$

$$(1.2.21) \quad \sum_{i+j \geq 3} \binom{k}{i} \binom{k}{j} (-1)^{i+j} E\tilde{a}_1^{i+j} \tilde{a}_{k-i} \tilde{a}_{k-j} = O\left(\frac{1}{n^2}\right).$$

Combining (1.2.15)–(1.2.21) we easily obtain relation (1.2.8). \square

Convergence in probability of sampling moments. We study the asymptotic behavior (as $n \rightarrow \infty$) of sampling moments a_k and m_k defined by (1.2.2) and (1.2.3), respectively. To indicate the dependence of moments a_k and m_k on the size of the sample we write a_{nk} and m_{nk} , respectively. Using (1.2.6) and the Chebyshev inequality we prove that $a_{nk} \rightarrow \alpha_k$ in probability as $n \rightarrow \infty$. A similar assertion holds for sampling central moments and even for arbitrary continuous functions of a finite number of sampling moments a_{nk} (the sampling central moment m_{nk} is a polynomial of moments $a_{n1}, a_{n2}, \dots, a_{nk}$ in view of (1.2.4)). The following result contains a precise statement of the latter assertion.

THEOREM 1.2.2. *Let random variables $\zeta_n^{(1)}, \zeta_n^{(2)}, \dots, \zeta_n^{(k)}$ converge in probability to some constants c_1, c_2, \dots, c_k , respectively, as $n \rightarrow \infty$. Let a function $f(z_1, z_2, \dots, z_k)$ be continuous in a neighborhood of the point (c_1, c_2, \dots, c_k) . Then the random variables $\eta_n = f(\zeta_n^{(1)}, \zeta_n^{(2)}, \dots, \zeta_n^{(k)})$ converge to $f(c_1, c_2, \dots, c_k)$ in probability as $n \rightarrow \infty$.*

PROOF. Let $\varepsilon > 0$ be an arbitrary number. Since $f(z_1, z_2, \dots, z_k)$ is continuous in a neighborhood of the point (c_1, c_2, \dots, c_k) , there is a number $\delta = \delta(\varepsilon)$ such that

$$|f(z_1, z_2, \dots, z_k) - f(c_1, c_2, \dots, c_k)| < \varepsilon$$

for $|z_i - c_i| < \delta$, $i = 1, 2, \dots, k$.

Consider the random events $B_i = \{|\zeta_n^{(i)} - c_i| < \delta\}$, $i = 1, 2, \dots, k$. Then $B \subset C$ where

$$B = \bigcap_{i=1}^k B_i, \quad C = \{|\eta_n - f(c_1, c_2, \dots, c_k)| < \varepsilon\}.$$

Thus

$$(1.2.22) \quad P(C) \geq P(B) = 1 - P\left(\bigcup_{i=1}^k \bar{B}_i\right) \geq 1 - \sum_{i=1}^k P(\bar{B}_i).$$

Since $\zeta_n^{(i)}$ converges in probability to c_i , for a given $\delta > 0$ and all $\gamma > 0$ there is $n_i = n_i(\gamma)$ such that $P(\bar{B}_i) < \gamma/k$ for $n \geq n_i$. Then $P(\bar{B}_i) < \gamma/k$ for

$$n \geq n_0 = \max(n_1, \dots, n_k)$$

and all $i = 1, 2, \dots, k$. Therefore $P(C) \geq 1 - \gamma$ by (1.2.22) if $n \geq n_0$. \square

Consider another application of Theorem 1.2.2. For continuous random variables we define the *skewness* γ_1 and *excess* γ_2 by

$$(1.2.23) \quad \gamma_1 = \mu_3 \mu_2^{-3/2}, \quad \gamma_2 = \mu_4 \mu_2^{-2} - 3.$$

If the density of a distribution is symmetric, then $\gamma_1 = 0$. Moreover, $\gamma_2 = 0$ in the case of the Gaussian distribution. Starting from (1.2.23) we construct the sampling skewness $g_1 = g_{n1}$ and sampling excess $g_2 = g_{n2}$ from the sample $\xi^{(n)} = (\xi_1, \xi_2, \dots, \xi_n)$:

$$(1.2.24) \quad g_{n1} = m_{n3} m_{n2}^{-3/2}, \quad g_{n2} = m_{n4} m_{n2}^{-2} - 3.$$

Applying Theorem 1.2.2 and equality (1.2.4), we prove that the sampling skewness and excess defined by (1.2.24) converge in probability as $n \rightarrow \infty$ to the corresponding skewness and excess defined by (1.2.23) for a given random variable.

Asymptotic normality of sampling moments. We introduce the following notation. The *law of distribution* of a random variable ξ is denoted by $\mathcal{L}(\xi)$. The law of distribution of the normal random variable ξ with expectation $\alpha = E\xi = \alpha_1$ and variance $\sigma^2 = D\xi = \mu_2$ is denoted by $\mathcal{L}(\xi) = \mathcal{N}(\alpha, \sigma^2)$.

We say that a sequence of random variables η_n , $n = 1, 2, \dots$, *weakly converges* as $n \rightarrow \infty$ to a random variable η if $\mathcal{L}(\eta_n) \rightarrow \mathcal{L}(\eta)$ as $n \rightarrow \infty$ (the convergence of laws $\mathcal{L}(\eta_n) \rightarrow \mathcal{L}(\eta)$ is understood as the convergence of distribution functions $P\{\eta_n < x\}$ to a distribution function $P\{\eta < x\}$ at all points of continuity of the function $P\{\eta < x\}$). In particular, the asymptotic normality of a sequence η_n with parameters (A_n, B_n^2) defined by (1.1.14) and (1.1.15) means that

$$\mathcal{L}((\eta_n - A_n)/B_n) \rightarrow \mathcal{N}(0, 1), \quad n \rightarrow \infty.$$

In the latter case we also say that a sequence η_n is $\mathcal{N}(A_n, B_n^2)$ asymptotically normal and occasionally write $\mathcal{L}(\eta_n) \sim \mathcal{N}(A_n, B_n^2)$.

The sampling moment a_{nk} is the sum of n independent identically distributed random variables (see (1.2.2)).

Applying the central limit theorem, we obtain the following result.

THEOREM 1.2.3. *If $\alpha_{2k} < \infty$, then the sequence a_{nk} of order k sampling moments is $\mathcal{N}(\alpha_k, (\alpha_{2k} - \alpha_k^2)/n)$ asymptotically normal.*

PROOF. It follows from (1.2.2) that

$$\sqrt{\frac{n}{\alpha_{2k} - \alpha_k^2}}(a_{nk} - \alpha_k) = \frac{1}{\sqrt{n(\alpha_{2k} - \alpha_k^2)}} \left(\sum_{i=1}^n \xi_i^k - n\alpha_k \right) = \eta_n.$$

Taking into account equalities (1.2.5) and (1.2.6) and applying the central limit theorem to the sum $\sum_{i=1}^n \xi_i^k$, we obtain $\mathcal{L}(\eta_n) \rightarrow \mathcal{N}(0, 1)$. Therefore the sequence a_{nk} is $\mathcal{N}(\alpha_k, (\alpha_{2k} - \alpha_k^2)/n)$ asymptotically normal. \square

Similarly to Theorem 1.2.3 one can prove the asymptotic normality of a continuous function of a finite number of sampling moments a_{nk} . In particular,

THEOREM 1.2.4. *If $\alpha_{2k} < \infty$, then the sequence m_{nk} of order k sampling central moments is*

$$\mathcal{N}(\mu_k, (\mu_{2k} - 2k\mu_{k-1}\mu_{k+1} - \mu_k^2 + k^2\mu_2\mu_{k-1}^2)/n)$$

asymptotically normal.

1.3. Order statistics and their properties

The distribution of order statistics. Let ξ be a random variable with the distribution function $F(x)$, let $\xi^{(n)} = (\xi_1, \xi_2, \dots, \xi_n)$ be observations of ξ , and let $\zeta_{n,1}, \zeta_{n,2}, \dots, \zeta_{n,n}$ be order statistics constructed from the sample ξ^n defined by (1.1.2). We study the distribution $F_{n,k}(x) = P\{\zeta_{n,k} < x\}$ of the k -th order statistic $\zeta_{n,k}$. It is clear that $\{\zeta_{n,k} < x\} = \{\nu_n(x) \geq k\}$ where $\nu_n(x)$ is the number of random variables in the sequence $\xi_1, \xi_2, \dots, \xi_n$ such that $\{\xi_k < x\}$. One can treat $\nu_n(x)$ as the number of occurrences of the event $\{\xi < x\}$ in n independent Bernoulli trials (see Section 1.1). Since $P\{\xi < x\} = F(x)$, the binomial distribution of $\nu_n(x)$ shows that

$$(1.3.1) \quad F_{n,k}(x) = P\{\nu_n(x) \geq k\} = \sum_{j=k}^n \binom{n}{j} F^j(x)(1 - F(x))^{n-j}.$$

The following result on the integral representation of the function $F_{n,k}$ is helpful for studies of its asymptotic properties.

THEOREM 1.3.1. *The distribution function $F_{n,k}(x)$ admits the following representation:*

$$(1.3.2) \quad F_{n,k}(x) = n \binom{n-1}{k-1} \int_0^{F(x)} t^{k-1} (1-t)^{n-k} dt.$$

In particular, if F has the density $f(x) = F'(x)$, then so does $F_{n,k}(x)$. Denote the density of $F_{n,k}(x)$ by $f_{n,k}(x)$ if the density f exists. Then

$$(1.3.3) \quad f_{n,k}(x) = n \binom{n-1}{k-1} F^{k-1}(x)(1 - F(x))^{n-k} f(x).$$

PROOF. Evaluating the integral in (1.3.2) by parts, we obtain

$$(1.3.4) \quad \int_0^{F(x)} t^{k-1}(1-t)^{n-k} dt \\ = \frac{1}{k} F^k(x)(1-F(x))^{n-k} + \frac{n-k}{n} \int_0^{F(x)} t^k(1-t)^{n-k-1} dt.$$

Substituting (1.3.4) into (1.3.2) we get

$$(1.3.5) \quad n \binom{n-1}{k-1} \int_0^{F(x)} t^{k-1}(1-t)^{n-k} dt \\ = \binom{n}{k} F^k(x)(1-F(x))^{n-k} + n \binom{n-1}{k} \int_0^{F(x)} t^k(1-t)^{n-k-1} dt.$$

Evaluating the integral on the right-hand side of (1.3.5) by parts, we conclude that

$$n \binom{n-1}{k-1} \int_0^{F(x)} t^{k-1}(1-t)^{n-k} dt = \sum_{j=k}^n \binom{n}{j} F^j(x)(1-F(x))^{n-j},$$

whence (1.3.2) follows by (1.3.1).

If the density $f(x) = F'(x)$ exists, then (1.3.3) follows from (1.3.2). \square

The joint distribution of two order statistics, say $\zeta_{n,k}$ and $\zeta_{n,m}$ with $k < m$, is also easy to evaluate. In particular, if the density $f(x) = F'(x)$ exists, then the density of the joint distribution of $\zeta_{n,k}$ and $\zeta_{n,m}$ also exists. Denoting it by $f_{n;k,m}(x,y)$ we have for $x < y$

$$(1.3.6) \quad f_{n;k,m}(x,y) = n(n-1) \binom{n-2}{k-1} \binom{n-k-1}{m-1} F^{k-1}(x) \\ \times (1-F(y))^{m-1} (F(y)-F(x))^{n-(k+m)} f(x)f(y).$$

To prove (1.3.6), we consider two disjoint intervals $[x, x + \Delta x]$ and $[y, y + \Delta y]$ for small Δx and Δy . Then we evaluate the probability of the event that exactly $k-1$ random variables of the sequence $\xi_1, \xi_2, \dots, \xi_n$ belong to the interval $(-\infty, x)$; only one random variable belongs to $[x, x + \Delta x]$; $m-1$ random variables belong to $[y, y + \Delta y, \infty)$; only one random variable belongs to $[y, y + \Delta y]$; and all other random variables belong to $[x + \Delta x, y)$. The probability of this event equals the right-hand side of (1.3.6) multiplied by $\Delta x \Delta y$ with a remainder term of a higher order with respect to $\Delta x \Delta y$. It is not hard to show that the probabilities of other events favorable to $\{\zeta_{n,k} \in [x, x + \Delta x], \zeta_{n,m} \in [y, y + \Delta y]\}$ are of higher orders with respect to $\Delta x \Delta y$ as compared to the probability of the event discussed above.

Limit theorems for extreme order statistics. Consider the k -th order statistic $\zeta_{n,k}$ whose index $k = k(n)$ depends on n in such a way that $k(n)/n$ approaches either 0 or 1 as $n \rightarrow \infty$. Those order statistics are called *extremes* or *extreme order statistics*. Below we study the cases of $k(n) = k = \text{const}$ and $k(n) = n - m + 1$ where m does not depend on n . In other words, we study the k -th order statistic from the left and m -th order statistic from the right for fixed constants k and m . Consider the limit behavior of extremes $\zeta_{n,k}$ and $\zeta_{n,n-m+1}$.

Let

$$(1.3.7) \quad \eta_n = nF(\zeta_{n,k}), \quad \kappa_n = n[1 - F(\zeta_{n,n-m+1})].$$

The following characteristics of a distribution function play an important role in the theorems below:

$$\bar{x}_0 = \sup\{x: F(x) = 0\}, \quad \underline{x}_1 = \inf\{x: F(x) = 1\}.$$

We agree that $\sup(\emptyset) = -\infty$ and $\inf(\emptyset) = \infty$.

THEOREM 1.3.2. *Assume that there exists $x' > \bar{x}_0$ such that $F(x)$ is continuous in the interval $(-\infty, x')$ and increases in the interval (\bar{x}_0, x') . Then*

$$(1.3.8) \quad \lim_{n \rightarrow \infty} P\{\eta_n < y\} = \Gamma_k(y)$$

for all $y \in \mathbf{R}$ where the random variables η_n are defined by (1.3.7), while $\Gamma_k(y)$ is given by

$$(1.3.9) \quad \Gamma_k(y) = \begin{cases} \frac{1}{(k-1)!} \int_0^y z^{k-1} e^{-z} dz, & y > 0, \\ 0, & y \leq 0. \end{cases}$$

PROOF. Equality (1.3.8) is obvious for $y \leq 0$. Fix a number $y > 0$ and let n be such that $y/n < F(x')$. Then the inverse function $F^{-1}(y)$ exists in the interval (\bar{x}_0, x') . Applying (1.3.2) we get

$$P\{\eta_n < y\} = P\left\{\zeta_{n,k} < F^{-1}\left(\frac{y}{n}\right)\right\} = n \binom{n-1}{k-1} \int_0^{y/n} t^{k-1} (1-t)^{n-k} dt.$$

Changing the variable $z = nt$ we obtain

$$P\{\eta_n < y\} = \binom{n-1}{k-1} \int_0^y \left(\frac{z}{n}\right)^{k-1} \left(1 - \frac{z}{n}\right)^{n-k} dz.$$

Note that $n^{-(k-1)} \binom{n-1}{k-1} \rightarrow 1/(k-1)!$ and $(1 - z/n)^{n-k} \rightarrow e^{-z}$ as $n \rightarrow \infty$, and moreover the convergence in the second relation is uniform with respect to $z \in (0, y)$ for an arbitrary finite y . This implies relation (1.3.8) by the Lebesgue dominated convergence theorem. \square

THEOREM 1.3.3. *Assume that there exists $x'' < \underline{x}_1$ such that $F(x)$ is continuous in the interval (x'', ∞) and increases in the interval (x'', \underline{x}_1) . Then*

$$\lim_{n \rightarrow \infty} P\{\kappa_n < y\} = \Gamma_m(y)$$

for all $y \in \mathbf{R}$ where the random variables κ_n are defined by (1.3.7), while $\Gamma_m(y)$ is defined by (1.3.9).

PROOF. Note that

$$\begin{aligned} P\{\kappa_n < y\} &= P\left\{\zeta_{n,n-m+1} > F^{-1}\left(1 - \frac{y}{n}\right)\right\} \\ &= n \binom{n-1}{m-1} \int_{1-y/n}^1 t^{n-m} (1-t)^{m-1} dt \\ &= n \binom{n-1}{m-1} \int_0^{y/n} z^{m-1} (1-z)^{n-m} dz \end{aligned}$$

if n is sufficiently large. The rest of the proof is the same as that of Theorem 1.3.2. \square

The function $\Gamma_k(y)$ defined by (1.3.9) is the so-called Gamma distribution function with parameter k .

The asymptotic behavior of extremes is a complicated problem in the case of a general k . A rather complete solution of this problem is given by Gnedenko (1943) and Smirnov (1949) (see, for example, [32]). Below we briefly discuss some of their results.

Consider the random variables $\tilde{\zeta}_{n,k} = (\zeta_{n,k} - A_n)/B_n$ where $k = \text{const}$ and A_n and $B_n > 0$ are appropriate constants depending on n . The possible limit distributions for $\tilde{\zeta}_{n,k}$ can only be of the following three types:

$$\psi_{k,\alpha}^{(1)}(x) = \begin{cases} 0, & x \leq 0, \\ \Gamma_k(x^\alpha), & x > 0, \end{cases} \quad \psi_{k,\alpha}^{(2)}(x) = \begin{cases} \Gamma_k(|x|^{-\alpha}), & x \leq 0, \\ 1, & x > 0, \end{cases}$$

$$\psi_k^{(3)}(x) = \Gamma_k(e^x), \quad -\infty < x < \infty,$$

where $\alpha > 0$. Necessary and sufficient conditions for convergence to any of the three types of limit laws are known in terms of the distribution function $F(x)$. Similar results are also obtained for extremes $\zeta_{n,n-m+1}$ with $m = \text{const}$.

To this end, we note that one can obtain the limit distribution for the pair of random variables η_n and κ_n defined by (1.3.7). Below is the corresponding result.

THEOREM 1.3.4. *Assume that all the assumptions of Theorems 1.3.2 and 1.3.3 hold. Then*

$$\lim_{n \rightarrow \infty} P\{\eta_n < x, \kappa_n < y\} = \Gamma_k(x)\Gamma_m(y)$$

for all $x < y$ where the function $\Gamma_k(x)$ is defined by (1.3.9).

Central order statistics and sampling quantiles. If $k = k(n)$ depends on n in such a way that $k(n)/n \rightarrow p$ as $n \rightarrow \infty$ and $0 < p < 1$, then the order statistic $\zeta_{n,k(n)}$ is called *central*. Sampling quantiles of a distribution can be expressed in terms of central order statistics.

Let $p \in (0, 1)$. Any number x_p such that

$$(1.3.10) \quad F(x_p) \leq p \quad \text{and} \quad F(x_p + 0) \geq p$$

is called a *p-quantile of a distribution F(x)*. It is clear that the system of inequalities (1.3.10) has at least one solution. A *p-quantile of the empirical distribution function F_n(x)* is called a *sampling p-quantile* and is denoted by \hat{x}_p .

THEOREM 1.3.5. *For all $p \in (0, 1)$, the sampling p-quantile can be represented as follows:*

$$\hat{x}_p = \begin{cases} \zeta_{n,[np]+1}, & \text{if } np \text{ is not an integer,} \\ \text{any number of the interval } [\zeta_{n,np}, \zeta_{n,np+1}], & \text{if } np \text{ is an integer,} \end{cases}$$

where $[a]$ stands for the integer part of a number a .

PROOF. Assume that np is not an integer. Generally speaking, there are order statistics $\zeta_{n,k}$ with $k \leq [np] + 1$ and such that $\zeta_{n,k} = \zeta_{n,[np]+1}$, whence

$$F_n(\zeta_{n,[np]+1}) \leq [np]/n < p.$$

On the other hand, there are order statistics $\zeta_{n,k}$ with $k \geq [np] + 1$ and such that $\zeta_{n,k} = \zeta_{n,[np]+1}$, thus $F_n(\zeta_{n,[np]+1} + 0) \geq ([np] + 1)/n > p$. Therefore $\hat{x}_p = \zeta_{n,[np]+1}$.

Now let np be an integer and let \hat{x}_p be any number of the interval $[\zeta_{n,np}, \zeta_{n,np+1}]$. As above $F_n(\hat{x}_p) \leq np/n = p$ and $F_n(\hat{x}_p + 0) \geq np/n = p$, that is, \hat{x}_p is a sampling p -quantile. \square

Below we use the notation $\hat{x}_{n,p}$ for \hat{x}_p to highlight that the sampling p -quantile depends on n .

Convergence in probability of the sampling p -quantiles. The consideration below excludes the cases of $p = 0$ and $p = 1$. These two cases require special treatment. If $p = 0$, then $\hat{x}_{n,0} \in (-\infty, \zeta_{n,1}]$ and therefore either $-\infty$ is the unique p -quantile or there are infinitely many p -quantiles.

Below we consider sampling p -quantiles $\hat{x}_{n,p}$ for $p \in (0, 1)$. The following result contains conditions for the convergence in probability of sampling quantiles to the corresponding quantiles x_p .

THEOREM 1.3.6. *If a p -quantile x_p is unique, then the sampling p -quantile $\hat{x}_{n,p}$ converges in probability to x_p as $n \rightarrow \infty$.*

PROOF. It is obvious that $F(x_p + \varepsilon) \geq F(x_p + 0) \geq p$ for all $\varepsilon > 0$. Since a p -quantile x_p is unique, $F(x_p + \varepsilon) > p$ for all $\varepsilon > 0$. The definition of the sampling p -quantile implies that $\{F_n(x_p + \varepsilon) > p\} \subset \{\hat{x}_{n,p} \leq x_p + \varepsilon\}$. Therefore

$$(1.3.11) \quad \mathbb{P}\{F_n(x_p + \varepsilon) > p\} \leq \mathbb{P}\{\hat{x}_{n,p} \leq x_p + \varepsilon\}.$$

By the law of large numbers (1.1.3)

$$\mathbb{P}\{F_n(x_p + \varepsilon) > p\} = \mathbb{P}\{F_n(x_p + \varepsilon) - F(x_p + \varepsilon) > -(F(x_p + \varepsilon) - p)\} \rightarrow 1$$

as $n \rightarrow \infty$. Using inequality (1.3.11) we get for all $\varepsilon > 0$ that

$$(1.3.12) \quad \lim_{n \rightarrow \infty} \mathbb{P}\{\hat{x}_{n,p} \leq x_p + \varepsilon\} = 1.$$

Using again the uniqueness of a p -quantile x_p and the same argument we obtain for all $\varepsilon > 0$

$$(1.3.13) \quad \lim_{n \rightarrow \infty} \mathbb{P}\{\hat{x}_{n,p} \geq x_p - \varepsilon\} = 1.$$

Relations (1.3.12) and (1.3.13) mean that $\hat{x}_{n,p}$ converges in probability to x_p as $n \rightarrow \infty$. \square

REMARK 1.3.1. In fact, $\hat{x}_{n,p}$ converges with probability 1 to x_p as $n \rightarrow \infty$ under the assumptions of Theorem 1.3.6.

REMARK 1.3.2. Let $\zeta_{n,k(n)}$ be a central order statistic such that $k(n)/n \rightarrow p$ as $n \rightarrow \infty$. Assume that a p -quantile x_p is unique. Then one can show by using Theorem 1.3.6 that $\zeta_{n,k(n)} \rightarrow x_p$ in probability as $n \rightarrow \infty$.

Asymptotic normality of sampling quantiles. Below we provide conditions for the asymptotic normality of a sampling p -quantile for $0 < p < 1$. First we give the following central limit theorem for independent identically distributed random variables in the scheme of series.

THEOREM 1.3.7 (LINDBERG). *Let $\xi_{n1}, \xi_{n2}, \dots, \xi_{nn}$ be independent identically distributed random variables such that*

$$\mathbb{E}\xi_{ni} = 0, \quad i = 1, \dots, n, \quad \sum_{i=1}^n \mathbb{E}\xi_{ni}^2 = 1, \quad n \geq 1.$$

Then the sequence of random variables $\sum_{i=1}^n \xi_{ni}$ is $\mathcal{N}(0, 1)$ asymptotically normal if and only if

$$(1.3.14) \quad \lim_{n \rightarrow \infty} \sum_{i=1}^n \mathbb{E}\xi_{ni}^2 I(|\xi_{ni}| > \tau) = 0$$

for all $\tau > 0$ where $I(A)$ is the indicator of a random event A .

The proof of Theorem 1.3.7 can be found in [23], p. 292.

THEOREM 1.3.8. *Let a distribution function F be continuous and let the equation $F(x) = p$ have a unique solution x_p . Moreover let the function $F(x)$ be differentiable at the point x_p and $F'(x_p) = f(x_p) > 0$. Then the sampling p -quantile $\hat{x}_{n,p}$ is $\mathcal{N}(x_p, n^{-1}pqf^{-2}(x_p))$ normal as $n \rightarrow \infty$ where $q = 1 - p$.*

PROOF. In view of Theorem 1.3.5, one can restrict consideration to the case $\hat{x}_{n,p} = \zeta_{n,k(n,p)}$ where $k(n,p) = [np] + 1$. It is obvious that it is sufficient to prove the $\mathcal{N}(0, 1)$ asymptotic normality for random variables

$$\eta_n = f(x_p) \sqrt{n/(pq)} (\hat{x}_{n,p} - x_p), \quad n = 1, 2, \dots$$

Note that $\{\hat{x}_{n,p} < x\} = \{\nu_n(x) \geq k(n,p)\}$ where $\nu_n(x) = \sum_{i=1}^n I_{(-\infty, x)}(\xi_i)$. Then

$$(1.3.15) \quad \begin{aligned} \mathbb{P}\{\eta_n < x\} &= \mathbb{P}\left\{\hat{x}_{n,p} < x_p + \sqrt{\frac{pq}{n}} \frac{x}{f(x_p)}\right\} \\ &= \mathbb{P}\left\{\nu_n\left(x_p + \sqrt{\frac{pq}{n}} \frac{x}{f(x_p)}\right) \geq k(n,p)\right\}. \end{aligned}$$

Consider random variables

$$\mu_{n,j} = I\left(\xi_j < x_p + x \sqrt{pq/n}/f(x_p)\right), \quad j = 1, 2, \dots, n.$$

It follows from (1.3.15) that

$$(1.3.16) \quad \begin{aligned} \mathbb{P}\{\eta_n < x\} &= \mathbb{P}\left\{\sum_{j=1}^n \mu_{n,j} \geq k(n,p)\right\} \\ &= \mathbb{P}\left\{\frac{1}{\sigma_n \sqrt{n}} \sum_{j=1}^n (\mu_{n,j} - a_n) \geq \frac{k(n,p) - na_n}{\sigma_n \sqrt{n}}\right\} \end{aligned}$$

where

$$\begin{aligned} a_n &= \mathbb{E}\mu_{n,j} = F\left(x_p + \sqrt{pq/n}x/f(x_p)\right), \\ \sigma_n^2 &= \mathbb{D}\mu_{n,j} = a_n(1 - a_n). \end{aligned}$$

Since $F(x_p) = p$, the Taylor expansion in a neighborhood of the point x_p shows that

$$(1.3.17) \quad F\left(x_p + \sqrt{\frac{pq}{n}} \frac{x}{f(x_p)}\right) = p + \sqrt{\frac{pq}{n}} x + o\left(\frac{1}{\sqrt{n}}\right)$$

as $n \rightarrow \infty$. Thus $a_n \rightarrow p$ and $\sigma_n^2 \rightarrow pq$ as $n \rightarrow \infty$, whence

$$I\left(|\mu_{n,i} - a_n|/\sqrt{n\sigma_n^2} > \tau\right) = 0$$

for all $\tau > 0$ and all sufficiently large n . This means that condition (1.3.14) holds for $\xi_{n,j} = (\mu_{n,j} - a_n)/\sqrt{n\sigma_n^2}$. Therefore the assumptions of Theorem 1.3.7 are satisfied for the random variables $(\mu_{n,j} - a_n)/\sqrt{n\sigma_n^2}$, $j = 1, 2, \dots, n$.

Applying Theorem 1.3.7 we derive from equality (1.3.16) that

$$(1.3.18) \quad P\{\eta_n < x\} = 1 - \Phi\left(\frac{k(n,p) - na_n}{\sigma_n\sqrt{n}}\right) + o(1)$$

as $n \rightarrow \infty$. Since $k(n,p) = [np] + 1 = np + 1 + r_n$ where $|r_n| < 1$, we obtain from equality (1.3.17) that

$$\frac{k(n,p) - na_n}{\sigma_n\sqrt{n}} = \frac{np + 1 + r_n - np - \sqrt{npq}x + o(\sqrt{n})}{\sqrt{npq}(1 + o(1))} = -x + o(1)$$

as $n \rightarrow \infty$.

This together with relation (1.3.18) implies that

$$P\{\eta_n < x\} = 1 - \Phi(-x) + o(1) = \Phi(x) + o(1)$$

as $n \rightarrow \infty$, that is, the sequence of random variables η_n is asymptotically $\mathcal{N}(0, 1)$ normal. \square

In particular, Theorem 1.3.8 implies that the central order statistic $\zeta_{n,[np]+1}$ is asymptotically normal.

The study of the asymptotic behavior of the central order statistic $\zeta_{n,k(n)}$ is a complicated problem for general $k(n)$. This problem is solved by Smirnov (1949) (see [32]) who showed in particular that if $k(n) = np + o(\sqrt{n})$, then the limit distributions for the sequence of random variables $\tilde{\zeta}_{n,k(n)} = (\zeta_{n,k(n)} - A_n)/B_n$ can only be of the following four types:

$$\begin{aligned} \Phi_\alpha^{(1)}(x) &= \begin{cases} \Phi(cx^\alpha), & x > 0, \\ 0, & x \leq 0, \end{cases} & \Phi_\alpha^{(2)}(x) &= \begin{cases} \Phi(-c|x|^\alpha), & x \leq 0, \\ 1, & x > 0, \end{cases} \\ \Phi_\alpha^{(3)}(x) &= \begin{cases} \Phi(-c_1|x|^\alpha), & x \leq 0, \\ \Phi(c_2x^\alpha), & x > 0, \end{cases} & \Phi^{(4)}(x) &= \begin{cases} 0, & x \leq -1, \\ \frac{1}{2}, & -1 < x \leq 1, \\ 1, & x > 1, \end{cases} \end{aligned}$$

where A_n and $B_n > 0$ are some appropriate constants depending on n and p ; α , c , c_1 , and c_2 are some positive constants; and $\Phi(x)$ is the standard Gaussian distribution function. Smirnov also obtained necessary and sufficient conditions on the distribution function $F(x)$ for the convergence to a given type of the limit laws.

REMARK 1.3.3. More details about order statistics and further references can be found in the book by David [10].

1.4. The distributions of some functions of Gaussian random vectors

We consider in this section the distributions of some functions of Gaussian random vectors that are widely used in various topics of mathematics.

The normal distribution. Let $X = (X_1, \dots, X_n)'$ be a random column-vector (here and throughout the symbol $'$ stands for the transposition of matrices and vectors). By $a = (a_1, a_2, \dots, a_n)'$ we denote the vector of its expectations of X , that is, $a_i = EX_i$, $i = 1, 2, \dots, n$, and by $\Lambda = (\lambda_{ij})$ we denote the $n \times n$ matrix of mixed central moments $\lambda_{ij} = E(X_i - a_i)(X_j - a_j)$, $i, j = 1, 2, \dots, n$. Note that the matrix Λ is symmetric and nonnegative definite. A random vector X is called *normal* (or *Gaussian*) if its characteristic function is of the form

$$(1.4.1) \quad \phi(t) = Ee^{iX't} = \exp \left\{ ia't - \frac{1}{2}t'\Lambda t \right\}$$

where $t = (t_1, t_2, \dots, t_n)'$. If X is a normal vector whose characteristic function is given by (1.4.1), then we write $\mathcal{L}(X) = \mathcal{N}(a, \Lambda)$, which means that X has a normal distribution. The distribution $\mathcal{L}(X) = \mathcal{N}(0, I_n)$ where I_n is the $n \times n$ unit matrix is called the *standard normal distribution*. The coordinates of the standard normal vector X are independent random variables whose distribution is $\mathcal{N}(0, 1)$. If the matrix Λ is nonsingular, then the normal distribution is called *proper* (or *nondegenerate*), in which case the distribution possesses the density

$$(1.4.2) \quad f(x) = (2\pi)^{-n/2}(\det\Lambda)^{-1/2} \exp \left\{ -\frac{1}{2}(x - a)'\Lambda^{-1}(x - a) \right\}$$

where $x = (x_1, x_2, \dots, x_n)'$ and $\det\Lambda$ is the determinant of the matrix Λ .

Linear transformations of Gaussian vectors are again Gaussian vectors. The precise statement is as follows.

LEMMA 1.4.1. *Let $Y = AX$ where $\mathcal{L}(X) = \mathcal{N}(a, \Lambda)$ and A is a $k \times n$ matrix. Then $\mathcal{L}(Y) = \mathcal{N}(b, B)$ for $b = Aa$ and $B = A\Lambda A'$.*

PROOF. Let $\phi_Y(u)$ and $\phi_X(t)$ be the characteristic functions of vectors Y and X , respectively. Then

$$\phi_Y(u) = E \exp \{iY'u\} = E \exp \{iX'A'u\} = \phi_X(A'u).$$

Since $\mathcal{L}(X) = \mathcal{N}(a, \Lambda)$, equality (1.4.1) implies that

$$\phi_Y(u) = \exp \left\{ ia'A'u - \frac{1}{2}(A'u)'\Lambda(A'u) \right\} = \exp \left\{ i(Aa)'u - \frac{1}{2}u'(A\Lambda A')u \right\}.$$

Thus $\mathcal{L}(Y) = \mathcal{N}(b, B)$ for $b = Aa$ and $B = A\Lambda A'$. □

Equality (1.4.1) implies for a diagonal matrix Λ that the coordinates of the vector X are independent. If the matrix Λ is not diagonal, then there is a linear transformation $Y = AX$ such that the coordinates of the vector Y are independent. Indeed, by Lemma 1.4.1, as a matrix A one can take an orthogonal matrix (this means that $AA' = I_n$) such that $A\Lambda A'$ is diagonal. This implies that if Λ is nonsingular, then there exists a nonsingular matrix A such that the vector $Y = AX$ has the standard normal distribution $\mathcal{N}(0, I_n)$ if $a = 0$.

Chi-square distribution and its properties. Let $X = (X_1, X_2, \dots, X_n)'$ and $\mathcal{L}(X) = \mathcal{N}(0, I_n)$. The distribution of the random variable $\chi_n^2 = \sum_{i=1}^n X_i^2$ is called the *chi-square distribution with n degrees of freedom*. Put $\mathcal{L}(\chi_n^2) = \chi^2(n)$ and let us find the density of $\chi^2(n)$. Applying (1.4.2) as $\Delta r \rightarrow 0$ we get

$$P\{\chi_n^2 \in [r, r + \Delta r)\} = P\left\{r \leq \sum_{i=1}^n X_i^2 < r + \Delta r\right\} = ke^{-r/2} \Delta (V_{S(\sqrt{r})}) + o(\Delta r)$$

where V_S is the volume of the ball $S(r) = \{x \in \mathbf{R}^n: |x| \leq r\}$ of radius r . Since $V_{S(\sqrt{r})} = Cr^{n/2}$, we have $\Delta(V_{S(\sqrt{r})}) = C'r^{n/2-1}\Delta r + o(\Delta r)$. Thus the density of the distribution $\chi^2(n)$ is given by

$$(1.4.3) \quad k_n(x) = K_n x^{n/2-1} e^{-x/2}, \quad x > 0,$$

where $K_n = (2^{n/2}\Gamma(n/2))^{-1}$ and $\Gamma(\cdot)$ is the Gamma function.

The characteristic function of the distribution $\chi^2(n)$ is

$$\phi(t; n) = E \exp\{it\chi_n^2\} = K_n \int_0^\infty x^{n/2-1} \exp\{-x(1-2it)/2\} dx.$$

Differentiating with respect to t we obtain

$$(1.4.4) \quad \phi'(t; n) = \frac{in}{1-2it} \phi(t; n).$$

Solving equation (1.4.4) subject to the condition $\phi(0; n) = 1$ yields

$$(1.4.5) \quad \phi(t; n) = (1-2it)^{-n/2}.$$

This equality allows one to find the moments of the distribution $\chi^2(n)$:

$$(1.4.6) \quad E\chi_n^2 = \frac{1}{i}\phi'(0; n) = n, \quad D\chi_n^2 = \frac{1}{i^2}\phi''(0; n) - (E\chi_n^2)^2 = 2n.$$

We also mention the following important property of the distribution $\chi^2(n)$. Let random variables $\chi_{n_1}^2$ and $\chi_{n_2}^2$ be independent and let $\mathcal{L}(\chi_{n_i}^2) = \chi^2(n_i)$, $i = 1, 2$. In view of equality (1.4.5) the characteristic function of the sum $\chi_{n_1}^2 + \chi_{n_2}^2$ is $\phi(t; n_1 + n_2)$, that is, $\mathcal{L}(\chi_{n_1}^2 + \chi_{n_2}^2) = \chi^2(n_1 + n_2)$. This means that the sum of independent chi-square random variables is again a chi-square random variable and its degree of freedom is equal to the sum of degrees of freedom of terms.

Linear and quadratic forms of normal random variables. Let

$$X = (X_1, X_2, \dots, X_n)'$$

be a random vector with the standard normal $\mathcal{L}(0, I_n)$ distribution. Consider a quadratic form

$$Q = \sum_{i,j=1}^n a_{ij} X_i X_j = X' A X$$

where $A = (a_{ij})$ and $A' = A$. We also consider m linear forms

$$Y_k = \sum_{i=1}^n b_{ki} X_i, \quad k = 1, 2, \dots, m.$$

Using matrix notation we rewrite the latter relations in a compact form as $Y = B X$ where B is a rectangle $m \times n$ matrix and $Y = (Y_1, \dots, Y_m)'$. By O we denote the

matrix with zero entries. The following result contains conditions for the independence of functions Q and Y .

LEMMA 1.4.2. *If $BA = O$, then the functions Q and Y are independent.*

PROOF. Since the real matrix A is symmetric, there exists an orthogonal matrix U such that $U'AU = D$ where D is a diagonal matrix with diagonal entries $\lambda_i \geq 0$, $i = 1, 2, \dots, n$. The numbers $\lambda_1, \dots, \lambda_n$ are characteristic numbers of the matrix A , that is, they are the roots of the characteristic equation $\det(A - \lambda I_n) = 0$. The columns u_k of the matrix $U = \|u_1 \dots u_n\|$ are eigenvectors of the matrix A , that is, $Au_k = \lambda_k u_k$, $k = 1, 2, \dots, n$.

Let r be the rank of the matrix A and let $\lambda_1, \dots, \lambda_r$ be nonzero characteristic numbers. The equality $A = UDU'$ can be viewed as the matrix form of the spectral representation of the matrix A , namely

$$(1.4.7) \quad A = \sum_{k=1}^r \lambda_k u_k u_k'.$$

By the assumptions of the lemma, $O = BA = \sum_{k=1}^n \lambda_k B u_k u_k'$. Multiplying this equality on the right by the vector u_s we get

$$(1.4.8) \quad B u_s = 0, \quad s = 1, 2, \dots, r,$$

since the vectors u_j are orthogonal. Put $Z = (Y_1, \dots, Y_m, u_1'X, \dots, u_r'X)$. It is clear that $Z = CX$ for some matrix C and thus Lemma 1.4.1 implies that the distribution of the vector Z is normal with $EZ = 0$. According to representation (1.4.7)

$$Q = \sum_{k=1}^n \lambda_k (X' u_k)(u_k' X) = \sum_{k=1}^n \lambda_k (u_k' X)^2.$$

Thus the equalities $EY_i u_s' X = 0$, $i = 1, 2, \dots, m$, $s = 1, 2, \dots, r$, complete the proof of the lemma, since they mean that Y_i and $u_s' X$ are independent in view of the normal distribution of the vector Z . To prove the above equalities we denote by b_i' the rows of the matrix B , $i = 1, 2, \dots, m$. Then we have by (1.4.8)

$$EY_i u_s' X = E b_i' X u_s' X = E b_i' X X' u_s = b_i' (E X X') u_s = b_i' I_n u_s = 0. \quad \square$$

Consider two quadratic forms $Q_1 = X'AX$ and $Q_2 = X'BX$.

LEMMA 1.4.3. *If $AB = BA = O$, then Q_1 and Q_2 are independent.*

PROOF. Let the matrix A admit the representation (1.4.7), and let the spectral representation of the matrix B be $B = \sum_{i=1}^s \nu_i v_i v_i'$ where $s = \text{rank } B$ is the rank of the matrix B . By the assumptions, $O = AB = \sum_{k,i} \lambda_k \nu_i u_k (u_k' v_i) v_i'$. Multiplying this equality on the left by u_i' and on the right by v_j , we get $u_i' v_j = 0$, $i = 1, \dots, r$, $j = 1, \dots, s$. Since the joint distribution of random variables $u_i' X$ and $v_j' X$ is normal, we prove as above that these random variables are independent. Now it follows from $Q_1 = \sum_{k=1}^r \lambda_k (u_k' X)^2$ and $Q_2 = \sum_{i=1}^s \nu_i (v_i' X)^2$ that the random variables Q_1 and Q_2 are independent. \square

The distributions of quadratic forms of normal random variables. By $\text{tr } A$ we denote the *trace* of a quadratic matrix A , that is, the sum of its diagonal entries.

LEMMA 1.4.4. Let $Q = X'AX$ where $\mathcal{L}(X) = \mathcal{N}(0, I_n)$ and $\text{rank } A = r \leq n$. If the matrix A is idempotent, that is, $A^2 = A$, then $\mathcal{L}(Q) = \chi^2(r)$ and $r = \text{tr } A$.

PROOF. Let the matrix A admit the representation (1.4.7). Since A is symmetric and idempotent, $\lambda_1 = \dots = \lambda_r = 1$. Thus $Q = \sum_{k=1}^r (u'_k X)^2$. The vectors u_k are orthonormal, thus the random variables $u'_k X$, $k = 1, 2, \dots, r$, are independent and $\mathcal{N}(0, 1)$ normal. Hence $\mathcal{L}(Q) = \chi^2(r)$. Since $A = UDU'$, we obtain from $\text{tr}(BA) = \text{tr}(AB)$ that $\text{tr } A = \text{tr}(U'UD) = \text{tr } D = \lambda_1 + \dots + \lambda_r = r$. \square

The following result, which is a corollary of the preceding assertions, is of its own interest as well.

THEOREM 1.4.1. Let an n -dimensional vector Y have a nondegenerate $\mathcal{N}(\mu, \Sigma)$ distribution. Then the distribution of the quadratic form $Q = (Y - \mu)' \Sigma^{-1} (Y - \mu)$ is $\chi^2(n)$.

PROOF. Let U be an orthogonal matrix such that $U' \Sigma U = D$ where D is a diagonal matrix. Since Σ is nondegenerate, all diagonal entries λ_k of the matrix D are positive. Thus $D^{-1/2}$ is well defined as the diagonal matrix with diagonal entries $\lambda_k^{-1/2}$. Consider the random vector $Z = D^{-1/2} U' (Y - \mu)$. By Lemma 1.4.1

$$\mathcal{L}(Z) = \mathcal{N}(0, I_n).$$

On the other hand, $Y - \mu = UD^{1/2}Z$. Thus $Q = Z' D^{1/2} U' \Sigma^{-1} U D^{1/2} Z = Z' Z$, whence $\mathcal{L}(Q) = \chi^2(n)$. \square

The following important result of the sampling theory is proved by Fisher (1925).

THEOREM 1.4.2. Let $\xi^{(n)} = (\xi_1, \dots, \xi_n)$ be a sample from the $\mathcal{N}(\mu, \sigma^2)$ distribution. Then the sampling moments a_1 and m_2 are independent. Moreover $\mathcal{L}(\sqrt{n}(a_1 - \mu)/\sigma) = \mathcal{N}(0, 1)$ and $\mathcal{L}(nm_2/\sigma^2) = \chi^2(n - 1)$.

PROOF. Consider a sample $\tilde{\xi}^{(n)} = (\tilde{\xi}_1, \dots, \tilde{\xi}_n)$ where

$$\tilde{\xi}_i = (\xi_i - \mu)/\sigma, \quad i = 1, 2, \dots, n.$$

Put

$$\tilde{a}_1 = \frac{1}{n} \sum_{i=1}^n \tilde{\xi}_i, \quad \tilde{m}_2 = \frac{1}{n} \sum_{i=1}^n (\tilde{\xi}_i - \tilde{a}_1)^2.$$

Then $\tilde{a}_1 = (a_1 - \mu)/\sigma$ and $\tilde{m}_2 = m_2/\sigma$. Thus it is sufficient to prove that \tilde{a}_1 and \tilde{m}_2 are independent, since $\mathcal{L}(\sqrt{n}\tilde{a}_1) = \mathcal{N}(0, 1)$ and $\mathcal{L}(n\tilde{m}_2) = \chi^2(n - 1)$. Consider an n -dimensional vector-column $b = (1/n, \dots, 1/n)'$ and $n \times n$ matrix $B = \|b \cdots b\|$. It is clear that $\tilde{a}_1 = b' \tilde{\xi}^{(n)}$ and $n\tilde{m}_2 = (\tilde{\xi}^{(n)} - B\tilde{\xi}^{(n)})' (\tilde{\xi}^{(n)} - B\tilde{\xi}^{(n)}) = (\tilde{\xi}^{(n)})' A \tilde{\xi}^{(n)}$ where $A = I_n - B$. Since $b'A = b' - b'B = b' - b' = 0$, the random variables \tilde{a}_1 and \tilde{m}_2 are independent by Lemma 1.4.2.

It is obvious that the distribution of a_1 is normal. Note that the matrix A is idempotent and $\text{tr } A = \text{tr } I_n - \text{tr } B = n - 1$. Then $\mathcal{L}(n\tilde{m}_2) = \chi^2(n - 1)$ by Theorem 1.4.1. \square

Student and Snedekor distributions. Let two random variables ξ and χ^2 be independent and let $\mathcal{L}(\xi) = \mathcal{N}(0, 1)$ and $\mathcal{L}(\chi^2) = \chi^2(n)$. Then the distribution of the random variable $t = \xi/\sqrt{\chi^2/n}$ is called the *Student distribution with n degrees of freedom* and is denoted by $S(n)$. The density $s_n(x)$ of the distribution $S(n)$ is

$$s_n(x) = \frac{1}{\sqrt{\pi n}} \frac{\Gamma((n+1)/2)}{\Gamma(n/2)} \frac{1}{(1+x^2/n)^{(n+1)/2}}, \quad x \in \mathbf{R}.$$

Let random variables χ_1^2 and χ_2^2 be independent and let $\mathcal{L}(\chi_i^2) = \chi^2(n_i)$, $i = 1, 2$. Then the distribution of the ratio

$$F = (\chi_1^2/n_1) / (\chi_2^2/n_2)$$

is called the *Snedekor distribution with n_1 and n_2 degrees of freedom* and is denoted by $S(n_1, n_2)$. This distribution is sometimes called the *F-distribution* or *Fisher distribution*. The density $s_{n_1, n_2}(x)$ of the distribution $S(n_1, n_2)$ is

$$s_{n_1, n_2}(x) = \left(\frac{n_1}{n_2}\right)^{n_1/2} \Gamma\left(\frac{n_1+n_2}{2}\right) \left(\Gamma\left(\frac{n_1}{2}\right)\Gamma\left(\frac{n_2}{2}\right)\right)^{-1} \\ \times x^{n_1/2-1} \left(1 + \frac{n_1 x}{n_2}\right)^{-(n_1+n_2)/2}, \quad x > 0.$$

The distributions $S(n)$ and $S(n_1, n_2)$ play an important role in the sampling theory.

REMARK 1.4.1. Properties of normal distributions and those related to normal distributions are treated in many textbooks on probability theory. A comprehensive text on properties and applications of normal distributions in statistical problems can be found in [1, 28].

Samples from Multidimensional Distributions

2.1. Empirical distribution function, sampling moments, and their properties

Empirical distribution function and its properties. Let (ξ, η) be a two-dimensional random vector with real coordinates ξ and η . Denote its distribution function by $F(x, y) = P\{\xi < x, \eta < y\}$, $x \in \mathbf{R}$, $y \in \mathbf{R}$. Assume that there are n independent observations of the vector (ξ, η) :

$$(2.1.1) \quad (\xi_1, \eta_1), (\xi_2, \eta_2), \dots, (\xi_n, \eta_n).$$

The set of observations is called a *sample from the two-dimensional distribution* $F(x, y)$.

For fixed $x \in \mathbf{R}$ and $y \in \mathbf{R}$ consider the following random variables:

$$(2.1.2) \quad \nu_n(x, y) = \sum_{i=1}^n I(\{\xi_i < x, \eta_i < y\}).$$

Then

$$(2.1.3) \quad F_n(x, y) = \frac{1}{n} \nu_n(x, y), \quad x \in \mathbf{R}, y \in \mathbf{R},$$

is called the *empirical distribution function of the sample* (2.1.1). Note that the empirical distribution function $F_n(x, y)$ possesses all the properties of regular two-dimensional distribution functions.

Equality (2.1.2) implies that $\nu_n(x, y)$ is the total number of occurrences of the event

$$\{\xi < x, \eta < y\}$$

in n independent trials, while (2.1.3) shows that the empirical distribution function $F_n(x, y)$ is the relative frequency of the event $\{\xi < x, \eta < y\}$ in n independent trials. Like the one-dimensional case, the Bernoulli law of large numbers implies that the empirical distribution function $F_n(x, y)$ approaches $F(x, y)$ in probability as $n \rightarrow \infty$ for all $x \in \mathbf{R}$ and $y \in \mathbf{R}$, that is,

$$\lim_{n \rightarrow \infty} P\{|F_n(x, y) - F(x, y)| > \varepsilon\} = 0 \quad \text{for all } \varepsilon > 0.$$

Moreover, the Borel strong law of large numbers implies that $F_n(x, y)$ approaches $F(x, y)$ with probability 1 as $n \rightarrow \infty$ for all $x \in \mathbf{R}$ and $y \in \mathbf{R}$, that is,

$$P\left\{\lim_{n \rightarrow \infty} F_n(x, y) = F(x, y)\right\} = 1.$$

Therefore the empirical distribution function $F_n(x, y)$ may serve as an approximation of the distribution function $F(x, y)$.

According to definition (2.1.3) the empirical distribution function $F_n(x, y)$ is a random variable for all fixed x and y . It assumes values k/n , $k = 0, 1, 2, \dots, n$, and moreover

$$P \left\{ F_n(x, y) = \frac{k}{n} \right\} = \binom{n}{k} F^k(x, y) (1 - F(x, y))^{n-k}.$$

Therefore

$$EF_n(x, y) = F(x, y), \quad DF_n(x, y) = F(x, y)(1 - F(x, y))/n.$$

Applying the De Moivre–Laplace central limit theorem we obtain the following result on the asymptotic normality of the empirical distribution function $F_n(x, y)$.

THEOREM 2.1.1. *The sequence of empirical distribution functions*

$$F_n(x, y), \quad n = 1, 2, \dots,$$

is asymptotically normal with parameters $(F(x, y), F(x, y)(1 - F(x, y))/n)$ for all fixed $x \in \mathbf{R}$ and $y \in \mathbf{R}$.

Moments of two-dimensional distributions. Let (ξ, η) be a real two-dimensional random vector. The number $\alpha_{ij} = E\xi^i\eta^j$ is called the *mixed moment of order $i + j$* (or, *$(i + j)$ -th mixed moment*) of the random vector (ξ, η) . The number $\mu_{ij} = E(\xi - \alpha_{i0})^i(\eta - \alpha_{0j})^j$ is called the *mixed central moment of order $i + j$* (or *$(i + j)$ -th mixed central moment*) of the random vector (ξ, η) . It is easy to see that

$$(2.1.4) \quad \mu_{ij} = \sum_{k=0}^i \sum_{l=0}^j \binom{i}{k} \binom{j}{l} (-1)^{k+l} \alpha_{i0}^k \alpha_{01}^l \alpha_{i-k, j-l}.$$

Note that α_{i0} is the i -th moment of the random variable ξ , while α_{0j} is the j -th moment of the random variable η . Analogously, μ_{i0} is the i -th central moment of the random variable ξ , while μ_{0j} is the j -th central moment of the random variable η . Note further that μ_{20} is the variance of ξ , while μ_{02} is the variance of η . We often use the notation $\sigma_1^2 = \mu_{20}$ and $\sigma_2^2 = \mu_{02}$. It is clear that $\mu_{20} = \alpha_{20} - \alpha_{10}^2$, $\mu_{02} = \alpha_{02} - \alpha_{01}^2$, and $\mu_{11} = \alpha_{11} - \alpha_{10}\alpha_{01}$.

If $\mu_{11} = 0$, then the random variables ξ and η are called *uncorrelated*. In this case $\alpha_{11} = \alpha_{10}\alpha_{01}$, that is, $E\xi\eta = E\xi E\eta$. If ξ and η are independent, then $\mu_{11} = 0$, that is, independent random variables are uncorrelated. The converse is, in general, false. In a particular case where the vector (ξ, η) has a normal distribution, the random variables ξ and η are independent if and only if they are uncorrelated.

Let $z = (t, u)'$ where t and u are real numbers. Consider a quadratic form

$$(2.1.5) \quad Q(z) = E[t(\xi - \alpha_{10}) + u(\eta - \alpha_{01})]^2 = \mu_{20}t^2 + 2\mu_{11}tu + \mu_{02}u^2.$$

Since $Q(z)$ is the expectation of a square of a random variable, $Q(z) \geq 0$ for all vectors z , whence it follows that the quadratic form $Q(z)$ is nonnegative definite. Definition (2.1.5) implies that $Q(z) = z'Mz$ where M is the matrix of central moments of second order:

$$M = \begin{pmatrix} \mu_{20} & \mu_{11} \\ \mu_{11} & \mu_{02} \end{pmatrix}.$$

The matrix M also is nonnegative definite (this follows from the same property of the quadratic form $Q(z)$). Hence

$$(2.1.6) \quad \mu_{02}\mu_{20} - \mu_{11}^2 \geq 0.$$

Rewriting (2.1.6) we obtain $\mu_{11}^2 \leq \mu_{02}\mu_{20}$, which is known as the *Cauchy-Bunyakovskii inequality*.

Denote the rank of M by r . The possible values of r are 0, 1, or 2. If $r = 2$, then (2.1.6) becomes a strict inequality, while (2.1.6) becomes an equality if $r = 0$ or $r = 1$. The following result contains some simple properties of the distribution of the vector (ξ, η) related to the rank r .

THEOREM 2.1.2. *The following are true:*

- 1) $r = 0$ if and only if the distribution of the vector (ξ, η) assigns unit mass to a single point;
- 2) $r = 1$ if and only if the distribution of the vector (ξ, η) is concentrated on a straight line, but not at a single point;
- 3) $r = 2$ if and only if the support of the distribution of the vector (ξ, η) does not coincide with a straight line or with a single point.

PROOF. We consider only the first two cases when $r = 0$ and $r = 1$. The case when $r = 2$ follows from 1) and 2).

If $r = 0$, then $\mu_{20} = \mu_{02} = 0$, so that the distribution of the random variables ξ and η is concentrated at the points α_{10} and α_{01} , respectively. Then the distribution of the vector (ξ, η) is concentrated at the single point $(\alpha_{10}, \alpha_{01})$. Conversely, if the distribution of the vector (ξ, η) is concentrated at a single point, then $\mu_{20} = \mu_{02} = 0$, so that $\mu_{11} = 0$ by (2.1.6). Therefore the rank of the matrix M is equal to zero.

If $r = 1$, then the quadratic form $Q(z)$ is not positive definite. Thus there is a vector $z_0 = (t_0, u_0) \neq 0$ such that $Q(z_0) = 0$. It means by (2.1.5) that with probability 1

$$(2.1.7) \quad t_0(\xi - \alpha_{10}) + u_0(\eta - \alpha_{01}) = 0,$$

which implies that the distribution of the vector (ξ, η) is concentrated at a straight line $t_0(x - \alpha_{10}) + u_0(y - \alpha_{01}) = 0$. Conversely, let the distribution of the vector (ξ, η) be concentrated at a straight line but not be concentrated at a single point. Obviously this line passes the point $(\alpha_{10}, \alpha_{01})$, whence it follows that equality (2.1.7) holds with probability 1 for some constants t_0 and u_0 . By (2.1.5) we have $Q(z_0) = 0$ for $z_0 = (t_0, u_0)'$, that is, the quadratic form $Q(z)$ is not positive definite. Since the distribution is not concentrated at a single point of the line (2.1.7), at least one of the numbers μ_{20} or μ_{02} is nonzero. Thus the rank of the matrix M equals 1. \square

If $\mu_{20} \neq 0$ and $\mu_{02} \neq 0$, then the rank of the matrix M equals either 1 or 2. Let

$$(2.1.8) \quad \rho = \frac{\mu_{11}}{\sqrt{\mu_{20}\mu_{02}}} = \frac{\mu_{11}}{\sigma_1\sigma_2}.$$

We have by (2.1.6) that $\rho^2 \leq 1$, that is, $|\rho| \leq 1$. It is clear that $|\rho| = 1$ if and only if the rank of the matrix M is equal to 1, that is, if and only if the support of the distribution of the vector (ξ, η) belongs to a straight line. In particular, if ξ and η are independent, then $\mu_{11} = 0$, whence $\rho = 0$. On the other hand, the equality $\rho = 0$ means that the random variables ξ and η are uncorrelated.

The number ρ defined by (2.1.8) is called the *coefficient of correlation* (or simply *correlation*) of random variables ξ and η .

Sampling moments. Consider a sample (2.1.1) consisting of n independent observations of a vector (ξ, η) and denote its distribution function by $F(x, y)$. Let $F_n(x, y)$ be the empirical distribution function of the sample (2.1.1) defined by (2.1.3). The $(i + j)$ -th mixed moment of the empirical distribution function $F_n(x, y)$ is called the $(i + j)$ -th *sampling mixed moment*, that is,

$$(2.1.9) \quad a_{ij} = \int \int x^i y^j dF_n(x, y) = \frac{1}{n} \sum_{k=1}^n \xi_k^i \eta_k^j.$$

The $(i + j)$ -th mixed central moment of the empirical distribution function $F_n(x, y)$ is called the $(i + j)$ -th *sampling mixed central moment*, that is,

$$(2.1.10) \quad \begin{aligned} m_{ij} &= \int \int (x - a_{10})^i (y - a_{01})^j dF_n(x, y) \\ &= \frac{1}{n} \sum_{k=1}^n (\xi_k - a_{10})^i (\eta_k - a_{01})^j. \end{aligned}$$

It follows from (2.1.9) and (2.1.10) that an analog of (2.1.4) holds for sampling mixed moments, namely

$$(2.1.11) \quad m_{ij} = \sum_{k=0}^i \sum_{l=0}^j \binom{i}{k} \binom{j}{l} (-1)^{k+l} a_{10}^k a_{01}^l a_{i-k, j-l}.$$

It follows from (2.1.9) that

$$(2.1.12) \quad E a_{ij} = \frac{1}{n} \sum_{k=1}^n E \xi_k^i \eta_k^j = \alpha_{ij},$$

$$(2.1.13) \quad D a_{ij} = \frac{1}{n^2} \sum_{k=1}^n D \xi_k^i \eta_k^j = \frac{1}{n} (\alpha_{2i, 2j} - \alpha_{ij}^2).$$

As in the proof of Theorem 1.2.1, we derive from (2.1.11) that

$$(2.1.14) \quad E m_{ij} = \mu_{ij} + O\left(\frac{1}{n}\right),$$

$$(2.1.15) \quad D m_{ij} = \frac{1}{n} R + O\left(\frac{1}{n^2}\right)$$

where R is a constant depending on mixed central moments μ_{kl} . Relations (2.1.12)–(2.1.15) hold if all expectations involved are finite.

The most important mixed moments are of order less than or equal to 2. We often use the notation $m_{20} = s_1^2$ and $m_{02} = s_2^2$. For sampling mixed central moments m_{ij} of order $i + j = 2$ we easily obtain that

$$E m_{ij} = \frac{n-1}{n} \mu_{ij}, \quad D m_{ij} = \frac{\mu_{2i, 2j} - \mu_{ij}^2}{n} + O\left(\frac{1}{n^2}\right),$$

which is a refinement of (2.1.14) and (2.1.15).

The number

$$r = \frac{m_{11}}{\sqrt{m_{20} m_{02}}} = \frac{m_{11}}{s_1 s_2}$$

is called the *sampling coefficient of correlation* (or simply *sampling correlation*). Since r is the coefficient of correlation of the distribution function $F_n(x, y)$, we have $|r| \leq 1$. The sampling coefficient of correlation r attains values ± 1 if and only if all the sampling points $(\xi_1, \eta_1), (\xi_2, \eta_2), \dots, (\xi_n, \eta_n)$ lie on a straight line. One can show that

$$(2.1.16) \quad Er = \rho + O\left(\frac{1}{n}\right),$$

$$(2.1.17) \quad Dr = \frac{\rho^2}{4n} \left(\frac{\mu_{40}}{\mu_{20}^2} + \frac{\mu_{04}}{\mu_{02}^2} + \frac{2\mu_{22}}{\mu_{20}\mu_{02}} + \frac{4\mu_{22}}{\mu_{11}^2} - \frac{4\mu_{31}}{\mu_{11}\mu_{20}} - \frac{4\mu_{13}}{\mu_{11}\mu_{02}} \right) + O\left(n^{-3/2}\right).$$

By the law of large numbers, the moments a_{ij} approach α_{ij} in probability as $n \rightarrow \infty$. The latter result also follows from (2.1.12) and (2.1.13) by the Chebyshev inequality if the moment $\alpha_{2i,2j}$ exists. Now we apply Theorem 1.2.2 to prove that $m_{ij} \rightarrow \mu_{ij}$ and $r \rightarrow \rho$ in probability as $n \rightarrow \infty$.

Provided $\alpha_{2i,2j} < \infty$ we use relations (2.1.12) and (2.1.13) and the central limit theorem for sums of independent identically distributed random variables $\xi_k^i \eta_k^j$, $k = 1, 2, \dots, n$, and prove that sampling mixed moments a_{ij} are asymptotically $\mathcal{N}(\alpha_{ij}, (\alpha_{2i,2j} - \alpha_{ij}^2)/n)$ normal. Applying (2.1.4) and (2.1.5) one can show that the sampling mixed central moments m_{ij} are asymptotically $\mathcal{N}(\mu_{ij}, R/n)$ normal under appropriate assumptions where R is a function of central moments involved with asymptotic equality (2.1.15). One can also see from (2.1.16) and (2.1.17) that the sampling correlation coefficient r is asymptotically $\mathcal{N}(\rho, C^2(\rho)/n)$ normal where

$$C^2(\rho) = \frac{\rho^2}{4} \left(\frac{\mu_{40}}{\mu_{20}^2} + \frac{\mu_{04}}{\mu_{02}^2} + \frac{2\mu_{22}}{\mu_{20}\mu_{02}} + \frac{4\mu_{22}}{\mu_{11}^2} - \frac{4\mu_{31}}{\mu_{11}\mu_{20}} - \frac{4\mu_{13}}{\mu_{11}\mu_{02}} \right).$$

REMARK 2.1.1. More details about two-dimensional sampling vectors can be found in the classical book by Cramér [9].

Samples from k -dimensional distributions for $k > 2$. Let $(\xi_1, \xi_2, \dots, \xi_k)$ be a k -dimensional vector with real coordinates $\xi_1, \xi_2, \dots, \xi_k$ and the distribution function

$$F(x_1, x_2, \dots, x_k) = P\{\xi_1 < x_1, \xi_2 < x_2, \dots, \xi_k < x_k\}.$$

The *moments* of this distribution are defined by

$$\alpha_{i_1 i_2 \dots i_k} = E\xi_1^{i_1} \xi_2^{i_2} \dots \xi_k^{i_k}.$$

The number $i_1 + i_2 + \dots + i_k$ is called the *order* of the moment. We use the notation $\alpha_1^{(i)} = E\xi_i^i$, $i = 1, 2, \dots, k$, for moments of first order. In particular, $(\alpha_1^{(1)}, \dots, \alpha_1^{(k)}) = (E\xi_1, \dots, E\xi_k)$. The *central moments* are defined by

$$\mu_{i_1 i_2 \dots i_k} = E\left(\xi_1 - \alpha_1^{(1)}\right)^{i_1} \left(\xi_2 - \alpha_1^{(2)}\right)^{i_2} \dots \left(\xi_k - \alpha_1^{(k)}\right)^{i_k}$$

where $i_1 + i_2 + \dots + i_k$ is the *order* of the moment. The general notation is inconvenient for use if $k > 2$. Thus we sometimes use another notation for moments of second order, namely

$$\lambda_{ii} = \sigma_i^2 = E\left(\xi_i - \alpha_1^{(i)}\right)^2, \quad \lambda_{ij} = E\left(\xi_i - \alpha_1^{(i)}\right)\left(\xi_j - \alpha_1^{(j)}\right) = \rho_{ij}\sigma_i\sigma_j.$$

Here σ_i^2 is the variance of the random variable ξ_i and λ_{ij} is the mixed central moment of second order of random variables ξ_i and ξ_j . The coefficient of correlation ρ_{ij} for random variables ξ_i and ξ_j is well defined if $\sigma_i \neq 0$ and $\sigma_j \neq 0$. Moreover $\rho_{ij} = \lambda_{ij}/(\sigma_i\sigma_j)$. The *matrix of central moments of second order* $\Lambda = (\lambda_{ij})$ is nonnegative definite. The matrix of coefficients of correlation

$$\rho = (\rho_{ij})$$

(well defined if all σ_i are positive) also is nonnegative definite. There is a relationship between matrices Λ and ρ , namely

$$\Lambda = \Sigma\rho\Sigma$$

where Σ is a diagonal matrix with diagonal entries $\sigma_1, \dots, \sigma_n$.

In particular, if $\lambda_{ij} = 0$ for all $i \neq j$, then the random variables $\xi_1, \xi_2, \dots, \xi_k$ are uncorrelated and the matrix Λ is diagonal, thus $\det(\Lambda) = \lambda_{11}\lambda_{22}\dots\lambda_{kk}$. If additionally all numbers σ_i are positive, then the matrix ρ is well defined and moreover $\rho = I_k$.

Consider n independent observations of a random vector (ξ_1, \dots, ξ_k) . This means that there is a sample $(\xi_{1i}, \dots, \xi_{ki}), i = 1, \dots, n$. Denote by $F_n(x_1, \dots, x_k)$ the empirical distribution function of this sample defined in the same way as in the case of two-dimensional vectors. The *sampling moments* are defined in this case by

$$a_{i_1 i_2 \dots i_k} = \frac{1}{n} \sum_{j=1}^n \xi_{1j}^{i_1} \xi_{2j}^{i_2} \dots \xi_{kj}^{i_k}.$$

The number $i_1 + \dots + i_k$ is called the *order* of the moment. We use the notation $a_1^{(i)} = n^{-1} \sum_{j=1}^n \xi_{ij}, i = 1, 2, \dots, k$, for sampling moments of first order. The *central sampling moments* are defined by

$$m_{i_1 i_2 \dots i_k} = \frac{1}{n} \sum_{j=1}^n \left(\xi_{1j} - a_1^{(1)} \right)^{i_1} \left(\xi_{2j} - a_1^{(2)} \right)^{i_2} \dots \left(\xi_{kj} - a_1^{(k)} \right)^{i_k}.$$

The number $i_1 + \dots + i_k$ is called the *order* of the moment. For moments of second order we use the simpler notation

$$l_{ii} = s_i^2 = \frac{1}{n} \sum_{j=1}^n \left(\xi_{ij} - a_1^{(i)} \right)^2, \quad i = 1, 2, \dots, k,$$

$$l_{ij} = \frac{1}{n} \sum_{m=1}^n \left(\xi_{im} - a_1^{(i)} \right) \left(\xi_{jm} - a_1^{(j)} \right) = r_{ij} s_i s_j.$$

Here s_i^2 is the sampling variance of the random variable ξ_i constructed from observations of the i -th coordinate of the vector, while $r_{ij} = l_{ij}/(s_i s_j)$ is the sampling coefficient of correlation between random variables ξ_i and ξ_j . Let $L = (l_{ij})$ be the matrix of sampling central moments of second order, and let $R = (r_{ij})$ be the matrix of sampling coefficients of correlation. It is clear that $L = SRS$ for the diagonal matrix S whose diagonal entries are s_1, s_2, \dots, s_k .

Asymptotic behavior of the empirical distribution function, sampling moments, and sampling coefficients of correlation for $k > 2$ are analogous to those in the cases of $k = 2$ and $k = 1$. Further results and other properties can be found in a classical

book by Cramér [9] as well as in other books devoted to the multidimensional case, for example in [1, 4, 18, 28].

2.2. Sampling regression and its properties

General regression. Let ξ and η be two random variables with the joint distribution function $F(x, y)$. We denote the conditional expectations by

$$(2.2.1) \quad m_1(y) = E\{\xi/\eta = y\}, \quad m_2(x) = E\{\eta/\xi = x\}.$$

The function $m_1(y)$ is called the *regression of ξ on η* , while the function $m_2(x)$ is called the *regression of η on ξ* . Regressions (2.2.1) possess an important property of minimality explained in the following result.

THEOREM 2.2.1. *If $E\xi^2 < \infty$, then for any Borel function f*

$$(2.2.2) \quad E(\xi - m_1(\eta))^2 \leq E(\xi - f(\eta))^2.$$

Analogously if $E\eta^2 < \infty$, then for any Borel function g

$$(2.2.3) \quad E(\eta - m_2(\xi))^2 \leq E(\eta - g(\xi))^2.$$

PROOF. We prove inequality (2.2.3), the proof of inequality (2.2.2) is analogous.

Let g be an arbitrary Borel function. Inequality (2.2.3) is trivial if

$$E(\eta - g(\xi))^2 = \infty.$$

Consider the case $E(\eta - g(\xi))^2 < \infty$. Then

$$(2.2.4) \quad \begin{aligned} E(\eta - g(\xi))^2 &= E((\eta - m_2(\xi)) + (m_2(\xi) - g(\xi)))^2 \\ &= E(\eta - m_2(\xi))^2 + 2E(\eta - m_2(\xi))(m_2(\xi) - g(\xi)) \\ &\quad + E(m_2(\xi) - g(\xi))^2. \end{aligned}$$

On the other hand,

$$(2.2.5) \quad \begin{aligned} E(\eta - m_2(\xi))(m_2(\xi) - g(\xi)) &= EE\{(\eta - m_2(\xi))(m_2(\xi) - g(\xi)) / \xi\} \\ &= E(m_2(\xi) - g(\xi))E\{\eta - m_2(\xi) / \xi\} = 0, \end{aligned}$$

since

$$E\{\eta - m_2(\xi) / \xi\} = E\{\eta / \xi\} - m_2(\xi) = m_2(\xi) - m_2(\xi) = 0.$$

It follows from (2.2.4) and (2.2.5) that

$$(2.2.6) \quad E(\eta - g(\xi))^2 = E(\eta - m_2(\xi))^2 + E(m_2(\xi) - g(\xi))^2 \geq E(\eta - m_2(\xi))^2. \quad \square$$

REMARK 2.2.1. Inequality (2.2.6) becomes an equality if and only if

$$E(m_2(\xi) - g(\xi))^2 = 0,$$

that is, if $P\{g(\xi) = m_2(\xi)\} = 1$. Thus inequality (2.2.3) becomes an equality if and only if $P\{g(\xi) = m_2(\xi)\} = 1$. Similarly inequality (2.2.2) becomes an equality if and only if $P\{f(\eta) = m_1(\eta)\} = 1$. This implies that the regressions $m_1(y)$ and $m_2(x)$ can be defined as functions minimizing the right-hand sides of inequalities (2.2.2) and (2.2.3), respectively. More precisely, every Borel function f^* , such that

$$(2.2.7) \quad E(\xi - f^*(\eta))^2 = \min E(\xi - f(\eta))^2$$

where the minimum is taken over all Borel functions f , is the regression of ξ on η . Similarly, every Borel function g^* , such that

$$(2.2.8) \quad E(\eta - g^*(\xi))^2 = \min E(\eta - g(\xi))^2$$

where the minimum is taken over all Borel functions g , is the regression of η on ξ .

Linear regression. We solved problem (2.2.7) in the class of all Borel functions and found a function $f(y)$ such that the random variable $f(\eta)$ as a function of η is the best *mean square approximation* of the random variable ξ . In other words, we found a function $f^*(y)$ minimizing the mean square error $E(\xi - f(\eta))^2$. Problem (2.2.7) is also of interest in the cases where we consider a narrower class of functions $f(y)$ instead of the class of all Borel functions. Say, one can solve problem (2.2.7) in the class of all linear functions or, more generally, in the class of polynomials of a fixed degree, etc. A similar remark concerns problem (2.2.8), too.

Let $L = \{\alpha + \beta x; \alpha, \beta \in (-\infty, \infty)\}$ be the class of all linear functions on \mathbf{R} . A function $g^*(x) = \alpha^* + \beta^*x$ such that

$$(2.2.9) \quad E(\eta - g^*(\xi))^2 = \min_{g \in L} E(\eta - g(\xi))^2$$

is called the *linear regression of η on ξ* .

A function $f^*(y) = \alpha^* + \beta^*y$ such that

$$(2.2.10) \quad E(\xi - f^*(\eta))^2 = \min_{f \in L} E(\xi - f(\eta))^2$$

is called the *linear regression of ξ on η* .

Below we find the linear regression of η on ξ , that is, we find a function

$$g^*(x) = \alpha^* + \beta^*x$$

solving problem (2.2.9). We assume that $\mu_{20} > 0$ and $\mu_{02} > 0$. Therefore we exclude the case of $\mu_{20} = 0$ and $\mu_{02} = 0$ for which the distribution of the vector (ξ, η) is concentrated at the point $(\alpha_{10}, \alpha_{01})$.

Let $G(\alpha, \beta) = E(\eta - g(\xi))^2$ where $g(x) = \alpha + \beta x$ is an arbitrary linear function. Then

$$(2.2.11) \quad \begin{aligned} G(\alpha, \beta) &= E((\eta - \alpha_{01}) - \beta(\xi - \alpha_{10}) + (\alpha_{01} - \beta\alpha_{10} - \alpha))^2 \\ &= \mu_{20}\beta^2 - 2\mu_{11}\beta + \mu_{02} + (\alpha_{01} - \beta\alpha_{10} - \alpha)^2. \end{aligned}$$

To solve the regression problem, it is sufficient to find the minimum of the function $G(\alpha, \beta)$. It follows from (2.2.11) that

$$(2.2.12) \quad \beta^* = \beta_{21} = \frac{\mu_{11}}{\mu_{20}} = \rho \frac{\sigma_2}{\sigma_1}, \quad \alpha^* = \alpha_{01} - \beta^*\alpha_{10}$$

where ρ is the coefficient of correlation between random variables ξ and η defined by (2.1.8). The number β_{21} defined by the first equality in (2.2.12) is called the *coefficient of linear regression of η on ξ* .

Substituting coefficients (2.1.12), the regression equation $y = g^*(x) = \alpha^* + \beta^*x$ becomes of the form

$$(2.2.13) \quad y = \alpha_{01} + \beta_{21}(x - \alpha_{10}).$$

This is the equation determining a straight line passing through the point $(\alpha_{10}, \alpha_{01})$. This equation can also be written in the form

$$(2.2.14) \quad \frac{y - \alpha_{01}}{\sigma_2} = \rho \frac{x - \alpha_{10}}{\sigma_1}.$$

Equation (2.2.14) is called the *canonical equation of the linear regression of η on ξ* .

The number $E(\eta - \alpha^* - \beta^*\xi)^2$ is the minimal mean square error in the problem (2.2.9) and is called the *least variance* of the random variable η . In view of (2.2.11) and (2.2.12) we get

$$(2.2.15) \quad E(\eta - \alpha^* - \beta^*\xi)^2 = \sigma_2^2(1 - \rho^2).$$

It follows from (2.2.15) that $|\rho| = 1$ if and only if $\eta = \alpha^* + \beta^*\xi$ with probability 1 where α^* and β^* are defined by (2.2.12), that is, $\eta = \alpha_{01} + \beta_{21}(\xi - \alpha_{10})$ with probability 1. Therefore we obtained the straight line for the assertion 2) of Theorem 2.1.2, that is, we found the coefficients t_0 and u_0 in equality (2.1.7).

If $\rho = 0$, then it follows from (2.2.14) that the linear regression of η on ξ is of the form $y = \alpha_{01}$. Note that this is the straight line parallel to the x -axis and passing through the point $(\alpha_{10}, \alpha_{01})$. Moreover we obtain from (2.2.15) that $E(\eta - \alpha^* - \beta^*\xi)^2 = \sigma_2^2$, that is, the variance of the random variable η does not decrease after subtracting the linear regression $\alpha^* + \beta^*\xi$.

It is not hard to prove that if the general regression $y = m_2(x)$ of η on ξ is linear, then it coincides with the linear regression given by (2.2.13).

Solving the analogous problem (2.2.10), we find the linear regression $f^*(y) = \alpha^* + \beta^*y$ of ξ on η whose coefficients are given by

$$(2.2.16) \quad \beta^* = \beta_{12} = \frac{\mu_{11}}{\mu_{02}} = \rho \frac{\sigma_1}{\sigma_2}, \quad \alpha^* = \alpha_{10} - \beta^* \alpha_{01}.$$

The number β_{12} is called the *coefficient of the linear regression of ξ on η* . Therefore the equation of the linear regression of ξ on η is of the form

$$(2.2.17) \quad x = \alpha_{10} + \beta_{12}(y - \alpha_{01}).$$

The regression can be rewritten in the canonical form:

$$(2.2.18) \quad \frac{x - \alpha_{10}}{\sigma_1} = \rho \frac{y - \alpha_{01}}{\sigma_2}.$$

The least variance of the random variable ξ is equal in this case to

$$(2.2.19) \quad E(\xi - \alpha^* - \beta^*\eta)^2 = \sigma_1^2(1 - \rho^2).$$

If $\rho = 0$, then it follows from (2.2.16) and (2.2.17) that the linear regression of ξ on η is $x = \alpha_{10}$; this is the straight line passing through the point $(\alpha_{10}, \alpha_{01})$ and parallel to the y -axis. It follows from (2.2.19) that $E(\xi - \alpha^* - \beta^*\eta)^2 = \sigma_1^2$, that is, the least variance of ξ coincides with the variance σ_1^2 .

If $|\rho| = 1$, then we obtain from (2.2.19) that $\xi = \alpha^* + \beta^*\eta$ with probability 1 where α^* and β^* are defined by (2.2.16). Moreover, we obtain from (2.2.14) and (2.2.18) that the linear regression of η on ξ and that of ξ on η coincide.

If $0 < |\rho| < 1$, then linear regression (2.2.18) can be rewritten in the form

$$(2.2.20) \quad \frac{y - \alpha_{01}}{\sigma_2} = \frac{1}{\rho} \frac{x - \alpha_{10}}{\sigma_1}$$

(cf. (2.2.14)). It follows from (2.2.14) and (2.2.20) that the linear regression of η on ξ and that of ξ on η coincide if and only if $|\rho| = 1$. Otherwise they do not coincide; in the case $\rho = 0$ they are perpendicular and each of them is parallel to the corresponding coordinate axis.

Parabolic regression. Let P be the family of polynomials

$$g(x) = c_0 + c_1x + \cdots + c_kx^k$$

of degree k whose coefficients c_0, c_1, \dots, c_k are real. A polynomial

$$g^*(x) = c_0^* + c_1^*x + \cdots + c_k^*x^k$$

such that

$$(2.2.21) \quad E(\eta - g^*(\xi))^2 = \min_{g \in P} E(\eta - g(\xi))^2$$

is called the *parabolic regression of η on ξ* . Assuming that all the moments occurring in (2.2.21) are finite we obtain the following condition for the minimum in (2.2.21):

$$(2.2.22) \quad \frac{1}{2} \frac{\partial G}{\partial c_i} = E(\xi^i(g(\xi) - \eta)) = c_0\alpha_{i,0} + \cdots + c_k\alpha_{i+k,0} - \alpha_{i,1} = 0,$$

$i = 0, 1, \dots, k$ (here $G = G(c_0, c_1, \dots, c_k) = E(\eta - g(\xi))^2$ for $g(x) = c_0 + c_1x + \cdots + c_kx^k$). If the moments α_{ij} are known, the coefficients $c_0^*, c_1^*, \dots, c_k^*$ can be determined from the above $k + 1$ equations. The evaluation of the coefficients can be simplified if the polynomial $g(x)$ is represented as a linear combination of *orthogonal polynomials $p_i(x)$ of degree i related to the distribution of ξ* and such that

$$(2.2.23) \quad E p_m(\xi) p_l(\xi) = \begin{cases} 1, & m = l, \\ 0, & m \neq l. \end{cases}$$

Any polynomial $g(x)$ of degree k can be represented as

$$g(x) = c_0 p_0(x) + c_1 p_1(x) + \cdots + c_k p_k(x)$$

for some coefficients c_0, c_1, \dots, c_k . According to (2.2.23) the condition for the minimum becomes of the form

$$(2.2.24) \quad \frac{1}{2} \frac{\partial G}{\partial c_i} = E(p_i(\xi)(g(\xi) - \eta)) = c_i - E\eta p_i(\xi), \quad i = 0, 1, \dots, k.$$

Using (2.2.24) we determine the coefficients $c_i^* = E\eta p_i(\xi)$, whence the parabolic regression of η on ξ is

$$(2.2.25) \quad g^*(x) = c_0^* p_0(x) + c_1^* p_1(x) + \cdots + c_k^* p_k(x).$$

Note that the coefficients c_i^* do not depend on the degree k of the polynomial $g(x)$. Therefore one can apply the recursion to find the regression. Namely if the parabolic regression $g^*(x)$ of degree k is known in the form (2.2.25), then the parabolic regression of degree $k + 1$ can be obtained in the form

$$c_0^* p_0(x) + \cdots + c_k^* p_k(x) + c_{k+1}^* p_{k+1}(x)$$

by evaluating only one extra number $c_{k+1}^* = E\eta p_{k+1}(\xi)$. By (2.2.23) we get

$$(2.2.26) \quad E(\eta - g^*(\xi))^2 = E\eta^2 - (c_0^*)^2 - \cdots - (c_k^*)^2$$

where $g^*(x) = c_0^*p_0(x) + \dots + c_k^*p_k(x)$. It is seen from (2.2.26) that a larger degree of the parabolic regression results in a smaller error of approximation $E(\eta - g^*(\xi))$.

Note that the above relations hold not only for polynomials $p_i(x)$. In fact, an arbitrary sequence of functions $p_i(x)$ satisfying condition (2.2.23) can be used to construct the function $g(x) = \sum_{i=1}^k c_i p_i(x)$. Note that relations (2.2.24) and (2.2.26) remain true in this case, too.

To this end we note that the parabolic regression of ξ on η can be evaluated in the same way as in the case of problem (2.2.21). This regression possesses the same properties as that of η on ξ .

Sampling linear regression. Let $F_n(x, y)$ be the empirical distribution function constructed from sample (2.1.1) according to (2.1.2) and (2.1.3) where sample (2.1.1) consists of n independent observations of the random vector (ξ, η) with the distribution function $F(x, y)$.

The *sampling linear regression of η on ξ* is called a function $g^*(x) = \alpha^* + \beta^*x$ such that

$$(2.2.27) \quad \sum_{i=1}^n (\eta_i - g^*(\xi_i))^2 = \min_{g \in L} \sum_{i=1}^n (\eta_i - g(\xi_i))^2$$

where L is the family of all linear functions. The *sampling linear regression of ξ on η* is called a function $f^*(y) = \alpha^* + \beta^*y$ such that

$$(2.2.28) \quad \sum_{i=1}^n (\xi_i - f^*(\eta_i))^2 = \min_{f \in L} \sum_{i=1}^n (\xi_i - f(\eta_i))^2.$$

To determine a linear regression $g^*(x) = \alpha^* + \beta^*x$ solving problem (2.2.27) it is sufficient to find $\alpha = \alpha^*$ and $\beta = \beta^*$ for which the function

$$(2.2.29) \quad G(\alpha, \beta) = \frac{1}{n} \sum_{i=1}^n (\eta_i - \alpha - \beta\xi_i)^2$$

attains the minimum. Taking into account (2.1.9) and (2.1.10) and applying the same argument as in the case of (2.2.11) we obtain

$$(2.2.30) \quad G(\alpha, \beta) = m_{20}\beta^2 - 2m_{11}\beta + m_{02} + (a_{01} - \beta a_{10} - \alpha)^2,$$

whence

$$(2.2.31) \quad \beta^* = b_{21} = \frac{m_{11}}{m_{20}} = r \frac{s_2}{s_1}, \quad \alpha^* = a_{01} - \beta^* a_{10}$$

where r is the sampling coefficient of regression (see Section 2.1). The random variable b_{21} defined by (2.2.31) is called the *sampling coefficient of linear regression of η on ξ* . Using (2.2.31) we rewrite the regression equation $y = g^*(x)$ in the form

$$(2.2.32) \quad \frac{y - a_{01}}{s_2} = r \frac{x - a_{10}}{s_1}.$$

Equation (2.2.32) is called the *canonical equation of sampling linear regression of η on ξ* .

As above we get

$$(2.2.33) \quad G(\alpha^*, \beta^*) = \frac{1}{n} \sum_{i=1}^n (\eta_i - \alpha^* - \beta^*\xi_i)^2 = s_2^2 (1 - r^2).$$

Reasoning in the same way as in the case of the least variance (2.2.15) we obtain the same results about $G(\alpha^*, \beta^*)$ defined by (2.2.33). The only difference is that these results involve sampling characteristics instead of the corresponding characteristics of the random vector (ξ, η) .

The method for obtaining sampling linear regressions as solutions of extremal problems (2.2.27) and (2.2.28) is called the *least squares method*. For example, the extremal problem (2.2.27) is to find the minimum of the sum of squares of distances evaluated in the y -direction between sampling points (ξ_i, η_i) and straight lines $y = \alpha + \beta x$. The same is true for the extremal problem (2.2.28) with the difference that the distances are evaluated in the x -direction between sampling points (ξ_i, η_i) and straight lines $y = \alpha + \beta x$.

Solving extremal problem (2.2.28) and using representations, similar to (2.2.29) and (2.2.30), we get the sampling linear regression $f^*(y) = \alpha^* + \beta^* y$ where

$$(2.2.34) \quad \beta^* = b_{12} = \frac{m_{11}}{m_{02}} = r \frac{s_1}{s_2}, \quad \alpha^* = a_{10} - \beta^* a_{01}.$$

The random variable b_{12} defined by (2.2.34) is called the *coefficient of the sampling linear regression of ξ on η* . The canonical equation $x = \alpha^* + \beta^* y$ of the sampling linear regression of ξ on η is

$$(2.2.35) \quad \frac{x - a_{10}}{s_1} = r \frac{y - a_{01}}{s_2}$$

in view of (2.2.34). Reasoning in the same way as in the case of equalities (2.2.18) and (2.2.20) we prove the same results for sampling regressions (2.2.32) and (2.2.35).

The asymptotic behavior of coefficients b_{12} and b_{21} as $n \rightarrow \infty$ is as follows. If $\sigma_1 \neq 0$ and $\sigma_2 \neq 0$, then $r \rightarrow \rho$ in probability as $n \rightarrow \infty$, thus by (2.2.31) and (2.2.34) we get $b_{21} \rightarrow \beta_{21}$ and $b_{12} \rightarrow \beta_{12}$ in probability as $n \rightarrow \infty$. Further applying (2.1.14) and (2.1.16) we obtain

$$Eb_{12} = \beta_{12} + O\left(\frac{1}{n}\right), \quad Eb_{21} = \beta_{21} + O\left(\frac{1}{n}\right)$$

provided that all necessary moments are finite. Moreover one can show that

$$Db_{12} = \frac{c_{12}}{n} + O\left(n^{-3/2}\right), \quad Db_{21} = \frac{c_{21}}{n} + O\left(n^{-3/2}\right)$$

where c_{12} and c_{21} are some positive constants.

One can also show that b_{12} is asymptotically $\mathcal{N}(\beta_{12}n^{-1}c_{12})$ normal, and b_{21} is asymptotically $\mathcal{N}(\beta_{21}, n^{-1}c_{21})$ normal.

Sampling parabolic regression. Let P be a family of polynomials

$$g(x) = \gamma_0 + \gamma_1 x + \dots + \gamma_k x^k$$

of degree k whose coefficients $\gamma_0, \gamma_1, \dots, \gamma_k$ are real. A polynomial

$$g^*(x) = \gamma_0^* + \gamma_1^* x + \dots + \gamma_k^* x^k$$

such that

$$(2.2.36) \quad \sum_{i=1}^n (\eta_i - g^*(\xi_i))^2 = \min_{g \in P} \sum_{i=1}^n (\eta_i - g(\xi_i))^2$$

is called the *sampling parabolic regression of η on ξ* . Let $g(x) = \gamma_0 + \gamma_1 x + \dots + \gamma_k x^k$ and

$$G(\gamma_0, \gamma_1, \dots, \gamma_k) = \frac{1}{n} \sum_{i=1}^n (\eta_i - g(\xi_i))^2 = \frac{1}{n} \sum_{i=1}^n (\eta_i - \gamma_0 - \gamma_1 \xi_i - \dots - \gamma_k \xi_i^k)^2.$$

Then condition (2.2.36) can be rewritten in the form

$$(2.2.37) \quad \frac{1}{2} \frac{\partial G}{\partial \gamma_j} = \frac{1}{n} \sum_{i=1}^n \xi_i^j (g(\xi_i) - \eta_i) = \gamma_0 a_{j,0} + \gamma_1 a_{j+1,0} + \dots + \gamma_k a_{j+k,0} - a_{j,1} \\ = 0, \quad j = 0, 1, \dots, k.$$

Condition (2.2.37) becomes simpler if the polynomial $g(x)$ is represented as a linear combination of orthogonal polynomials $p_m(x)$ of degree m related to the sampling distribution of ξ and such that

$$(2.2.38) \quad \frac{1}{n} \sum_{i=1}^n p_m(\xi_i) p_l(\xi_i) = \begin{cases} 1, & m = l, \\ 0, & m \neq l. \end{cases}$$

Every polynomial $g(x)$ of degree k can be represented as

$$g(x) = c_0 p_0(x) + \dots + c_k p_k(x)$$

with some real coefficients c_0, c_1, \dots, c_k . Using property (2.2.38) we rewrite condition (2.2.37) in the form

$$(2.2.39) \quad \frac{1}{2} \frac{\partial G}{\partial c_j} = \frac{1}{n} \sum_{i=1}^n p_j(\xi_i) (g(\xi_i) - \eta_i) = c_j - \frac{1}{n} \sum_{i=1}^n \eta_i p_j(\xi_i) = 0, \\ j = 0, 1, \dots, k.$$

This system of equations is easy to solve; its solution is given by

$$(2.2.40) \quad c_j^* = \frac{1}{n} \sum_{i=1}^n \eta_i p_j(\xi_i), \quad j = 0, 1, \dots, k.$$

Thus the sampling parabolic regression of η on ξ is

$$(2.2.41) \quad g^*(x) = c_0^* p_0(x) + c_1^* p_1(x) + \dots + c_k^* p_k(x).$$

In view of (2.2.40) the coefficients c_j^* depend only on the polynomial $p_j(x)$. This implies that if the sampling parabolic regression $g^*(x)$ of degree k is represented in the form of (2.2.41) and one needs to obtain the regression of a higher degree, say of the degree $k+1$, that is,

$$g^*(x) = c_0^* p_0(x) + c_1^* p_1(x) + \dots + c_k^* p_k(x) + c_{k+1}^* p_{k+1}(x),$$

then one can use coefficients c_j^* , $j = 0, 1, \dots, k$, known from the preceding regression. The only extra work is to evaluate the coefficient c_{k+1}^* by using (2.2.40) for $j = k+1$. For the regression (2.2.41) we have

$$(2.2.42) \quad \frac{1}{n} \sum_{i=1}^n (\eta_i - g^*(\xi_i))^2 = a_{20} - (c_0^*)^2 - \dots - (c_k^*)^2$$

where a_{20} is the sampling moment of order 2 defined by (2.1.9). It follows from equality (2.2.42) that a larger degree of the regression results in a smaller error of approximation.

To this end we note that the sampling parabolic regression of ξ on η is defined similarly to (2.2.21). This regression possesses all the properties as does that of η on ξ .

REMARK 2.2.2. In this section we considered the regression analysis only for two-dimensional vectors. For the higher-dimensional case see [1, 18, 28].

Estimation of Unknown Parameters of Distributions

3.1. Statistical estimators and their quality measures

Parametric families of distributions and statistical estimators of parameters. Let ξ be an observation that is a random element assuming values in a measurable space (X, \mathcal{B}) . Let the probability distribution of the random element ξ be either unknown or partially known. Let $\{\mathbb{P}_\theta, \theta \in \Theta\}$ be a family of probability measures on (X, \mathcal{B}) and let $\theta = (\theta_1, \theta_2, \dots, \theta_k)$ be a k -dimensional parameter belonging to a subset $\Theta \subset \mathbf{R}^k$, $k \geq 1$. We assume that the distribution of the random element ξ depends on the parameter θ which is unknown for the statistician. Thus the measure \mathbb{P}_θ is the distribution of ξ if the unknown parameter is equal to θ , that is, $\mathbb{P}_\theta\{\xi \in A\} = \mathbb{P}_\theta(A)$ for all $A \in \mathcal{B}$. The problem is to estimate the unknown parameter θ or a function $\phi(\theta)$ of the parameter θ with the help of the observation $\xi = x$.

The space X is called the *sampling space*. Every measurable function $T = T(x)$ mapping (X, \mathcal{B}) onto a measurable space (Y, \mathcal{S}) is called a *statistic*. If Θ is a Borel set of \mathbf{R}^k and $\mathcal{B}(\Theta)$ is the σ -algebra of Borel subsets of Θ , then, in the case of $(Y, \mathcal{S}) = (\Theta, \mathcal{B}(\Theta))$, a statistic $T = T(x)$ is called a *statistical estimator* (or just *estimator*) of an unknown parameter θ constructed from an observation $\xi = x$. In the case of $(Y, \mathcal{S}) = (\mathbf{R}^k, \mathcal{B}^k)$ and $\mathbf{R}^k \neq \Theta$, we sometimes refer to a statistic $T = T(x)$ as an estimator of a parameter θ .

The notion of a statistical estimator of a function of a parameter θ can be introduced in an analogous way. The random variable $T = T(\xi)$ is also called an estimator of a parameter (or, an estimator of a function of a parameter).

Statistical estimators of a parameter θ introduced above are sometimes called *point estimators*. A point estimator T constructed from an observation $\xi = x$ provides a single value $t = T(x)$ which we treat as an approximation of the parameter. However the true value of the parameter can be (and usually is) different from an estimator. Therefore it is very important to know the error of approximation appearing due to a specific estimator. For this purpose, statisticians usually also indicate a region (an interval, if $k = 1$) such that the probability that the true value of a parameter θ belongs to the region is large.

Let $k = 1$; thus θ is a one-dimensional (scalar) parameter. Let $T_1 = T_1(x)$ and $T_2 = T_2(x)$ be two statistics with values in \mathbf{R}^1 . Assume that $T_1 < T_2$ and for a given $\gamma \in (0, 1)$

$$(3.1.1) \quad \mathbb{P}_\theta\{T_1(\xi) < \theta < T_2(\xi)\} \geq \gamma \quad \text{for all } \theta \in \Theta.$$

The interval $(T_1(\xi), T_2(\xi))$ is called a γ -*confidence interval* or a *confidence interval of level γ* for the parameter θ . The number γ is often called the *confidence*

probability or *confidence level*. The numbers $T_1(\xi)$ and $T_2(\xi)$ are called the *lower* and *upper confidence bounds*, respectively.

Now let $k > 1$. Then a parameter θ is a k -dimensional vector and instead of γ -confidence intervals we define γ -confidence regions $G = G(\xi) \subset \mathbf{R}^k$ with the help of a condition similar to (3.1.1), namely

$$(3.1.2) \quad P_\theta\{\theta \in G(\xi)\} \geq \gamma \quad \text{for all } \theta \in \Theta.$$

The γ -confidence intervals and γ -confidence regions defined by (3.1.1) and (3.1.2) can be constructed by using point estimators of the unknown parameter θ . For example, if T is a point estimator of a parameter θ , then, as a confidence interval, one can take $(T - \delta, T + \delta)$ where $\delta > 0$ is found from condition (3.1.1).

In what follows we often treat a sample $\xi^{(n)} = (\xi_1, \xi_2, \dots, \xi_n)$ as an observed random element ξ . For this case, the sampling space is $(\mathbf{R}^n, \mathcal{B}^n)$, while the family of distributions of the sample $\xi^{(n)}$ is $\{P_\theta^n, \theta \in \Theta\}$. The point estimation of parameters for this case is considered in the book by Lehmann [21].

Unbiased and consistent estimators. Let ξ be an observed random element, let $\theta = (\theta_1, \dots, \theta_k)$ be an unknown parameter of the distribution, and let $\hat{\theta} = (\hat{\theta}_1, \dots, \hat{\theta}_k)$ be a statistical estimator of the parameter θ constructed from the observation ξ . An estimator $\hat{\theta}$ is called *unbiased* if

$$(3.1.3) \quad E_\theta \hat{\theta} = \theta \quad \text{for all } \theta \in \Theta$$

where E_θ is the expectation with respect to the distribution P_θ .

In the case of estimation of a function $\phi(\theta)$ of a parameter θ , a statistic $\hat{\phi}$ is called an unbiased estimator of the function $\phi(\theta)$ if

$$(3.1.4) \quad E_\theta \hat{\phi} = \phi(\theta) \quad \text{for all } \theta \in \Theta.$$

By $\hat{\theta}_n = \hat{\theta}(\xi^{(n)})$ we denote statistical estimators of a parameter θ in the case where an observed random element is a sample $\xi^{(n)} = (\xi_1, \xi_2, \dots, \xi_n)$. In such a case we deal with a sequence of estimators $\hat{\theta}_n$, $n = 1, 2, \dots$.

If an estimator $\hat{\theta}$ does not satisfy condition (3.1.3), then we consider the *bias* of the estimator $\hat{\theta}$ defined by

$$b(\theta) = E_\theta \hat{\theta} - \theta.$$

A sequence of estimators $\hat{\theta}_n$, $n = 1, 2, \dots$, is called an *asymptotically unbiased estimator* of a parameter θ if

$$(3.1.5) \quad \lim_{n \rightarrow \infty} E_\theta \hat{\theta}_n = \theta \quad \text{for all } \theta \in \Theta$$

or, in other words,

$$\lim_{n \rightarrow \infty} b_n(\theta) = 0 \quad \text{for all } \theta \in \Theta$$

where $b_n(\theta) = E_\theta \hat{\theta}_n - \theta$ is the bias of the estimator $\hat{\theta}_n$. The notion of asymptotically unbiased estimators of a function $\phi(\theta)$ of a parameter can be introduced in a similar way.

When analyzing data, statisticians often restrict themselves to the case of unbiased estimators, since there exists a simple and useful theory of unbiased estimators where the quality of an estimator is measured by its variance.

On the other hand, there are many cases where the requirement that an estimator should be unbiased is too restrictive. For example, it is possible that unbiased estimators do not exist at all or are useless in practice for a given parametric model. To see this we consider the following examples.

EXAMPLE 3.1.1. Let ξ be a Poisson random variable with parameter $\theta > 0$, that is, $P_\theta\{\xi = x\} = \theta^x e^{-x}/x!$, $x = 0, 1, 2, \dots$. Assume that we want to estimate the function $\phi(\theta) = 1/\theta$ of the parameter θ by an observation ξ . Let $T = T(\xi)$ be an unbiased estimator of $\phi(\theta)$, that is, condition (3.1.4) holds. Then it can be rewritten as

$$\sum_{x=0}^{\infty} T(x) \frac{\theta^x}{x!} e^{-x} = \frac{1}{\theta} \quad \text{for all } \theta \in (0, \infty)$$

or, in other words,

$$(3.1.6) \quad \sum_{x=0}^{\infty} T(x) \frac{\theta^{x+1}}{x!} = e^\theta = \sum_{s=0}^{\infty} \frac{\theta^s}{s!} \quad \text{for all } \theta \in (0, \infty).$$

It is obvious that there is no function $T(x)$ that satisfies condition (3.1.6) for all $\theta \in (0, \infty)$ and does not depend on θ . This means that there is no unbiased estimator of $\phi(\theta) = 1/\theta$.

EXAMPLE 3.1.2. Let ξ have the geometric distribution with parameter

$$\theta \in (0, 1),$$

that is, $P_\theta\{\xi = x\} = \theta^x(1 - \theta)$, $x = 0, 1, 2, \dots$. Assume that we want to estimate the parameter θ . Then the condition that $T = T(\xi)$ is an unbiased estimator is given by

$$\sum_{x=0}^{\infty} T(x) \theta^x = \frac{\theta}{1 - \theta} = \sum_{s=1}^{\infty} \theta^s \quad \text{for all } \theta \in (0, 1).$$

Comparing the coefficients for degrees of θ we see that the only unbiased estimator of θ is the statistic $T(x)$ such that $T(0) = 0$ and $T(x) = 1$ for $x = 1, 2, \dots$. Since this statistic does not belong to the set $\Theta = (0, 1)$ of possible values of the parameter, it gives a wrong approximation of the true value of θ and the estimator is useless for practice.

The following example shows that, at least in some cases, an estimator with a small bias and small mean square error is better than an unbiased estimator with a large variance.

EXAMPLE 3.1.3. Let $\xi^{(n)}$ be a sample from the normal $\mathcal{N}(\theta_1, \theta_2^2)$ distribution, $\xi^{(n)} = (\xi_1, \xi_2, \dots, \xi_n)$. The unknown parameter is $\theta = (\theta_1, \theta_2)$. Consider the problem on estimating the function $\phi(\theta) = \theta_2^2$.

Consider the sampling variance

$$(3.1.7) \quad s^2 = m_2 = \frac{1}{n} \sum_{i=1}^n (\xi_i - a_1)^2$$

as an estimator where a_1 is the first order sampling moment defined by (1.2.2). Applying (1.2.9) for $k = 2$ we get $s^2 = \tilde{a}_2 - \tilde{a}_1^2$ where \tilde{a}_1 and \tilde{a}_2 are defined in the proof of Theorem 1.2.1. Note that $E_\theta \tilde{a}_2 = \theta_2^2$. In view of (1.2.10) for $k = 2$ we have $E_\theta \tilde{a}_1^2 = \theta_2^2/n$, whence $E_\theta s^2 = \theta_2^2(n - 1)/n$, that is, the estimator s^2 is biased

(however it is asymptotically unbiased). This implies that an unbiased estimator is given by

$$(3.1.8) \quad \tilde{s}^2 = \frac{n}{n-1} s^2 = \frac{1}{n-1} \sum_{i=1}^n (\xi_i - a_1)^2.$$

Theorem 1.4.2 implies that $\mathcal{L}((n-1)\tilde{s}^2/\theta_2^2) = \chi^2(n-1)$. This together with (1.4.6) yields

$$(3.1.9) \quad D_\theta \tilde{s}^2 = \frac{2\theta_2^4}{n-1}.$$

Consider the class of estimators $T_\lambda = \lambda \tilde{s}^2$, $\lambda \in (0, \infty)$. Since

$$\mathbf{E}_\theta T_\lambda = \lambda \mathbf{E}_\theta \tilde{s}^2 = \lambda \theta_2^2,$$

there is only one unbiased estimator \tilde{s}^2 of the function $\phi(\theta)$ in this class. The mean square error of the estimator T_λ equals

$$(3.1.10) \quad \mathbf{E}_\theta (T_\lambda - \theta_2^2)^2 = \left(\frac{2\lambda^2}{n-1} + (1-\lambda)^2 \right) \theta_2^4.$$

The right-hand side of (3.1.10) attains its minimum at $\lambda^* = (n-1)/(n+1)$. Taking into account (3.1.9) we obtain

$$\mathbf{E}_\theta (T_{\lambda^*} - \theta_2^2)^2 = \frac{2}{n+1} \theta_2^4 < \frac{2}{n-1} \theta_2^4 = \mathbf{E}_\theta (\tilde{s}^2 - \theta_2^2)^2.$$

Therefore the estimator T_{λ^*} has a smaller mean square error than that of the unbiased estimator \tilde{s}^2 . Since $\mathbf{E}_\theta T_{\lambda^*} = (n-1)\theta_2^2/(n+1)$, the estimator T_{λ^*} is asymptotically unbiased. Note that s^2 also is an asymptotically unbiased estimator, but $s^2 = T_{\lambda'}$ for $\lambda' = (n-1)/n \neq \lambda^*$. This means that the estimator s^2 is worse than T_{λ^*} in the sense of the minimum of the mean square error.

Moreover,

$$(3.1.11) \quad \mathbf{E}_\theta (T_{\lambda^*} - \theta_2^2)^2 < \mathbf{E}_\theta (s^2 - \theta_2^2)^2 < \mathbf{E}_\theta (\tilde{s}^2 - \theta_2^2)^2.$$

Let $\hat{\phi}_n$, $n = 1, 2, \dots$, be a sequence of estimators of a function $\phi(\theta)$ of a parameter θ constructed from a sample $\xi^{(n)} = (\xi_1, \xi_2, \dots, \xi_n)$. A sequence of estimators $\hat{\phi}_n$, $n = 1, 2, \dots$, is called a *consistent sequence of estimators* of a function $\phi(\theta)$ if for all $\varepsilon > 0$

$$(3.1.12) \quad \lim_{n \rightarrow \infty} \mathbf{P}_\theta \{ |\hat{\phi}_n - \phi(\theta)| > \varepsilon \} = 0 \quad \text{for all } \theta \in \Theta.$$

For brevity we also say that an estimator $\hat{\phi}_n$ satisfying condition (3.1.12) is a *consistent estimator* of the function $\phi(\theta)$.

Note that s^2 and \tilde{s}^2 introduced in Example 3.1.3 by formulas (3.1.7) and (3.1.8) respectively, as well as $T_{\lambda^*} = \lambda^* \tilde{s}^2$ are consistent estimators of the function

$$\phi(\theta) = \theta_2^2,$$

since they satisfy condition (3.1.12) in view of the Chebyshev inequality and relations (3.1.9) and (3.1.11).

We say that a sequence of estimators $\widehat{\phi}_n$, $n = 1, 2, \dots$, is a *strongly consistent sequence of estimators of a function $\phi(\theta)$* (alternatively, $\widehat{\phi}_n$ is a *strongly consistent estimator of a function $\phi(\theta)$*) if

$$(3.1.13) \quad P_\theta \left\{ \lim_{n \rightarrow \infty} \widehat{\phi}_n = \phi(\theta) \right\} = 1.$$

For example, it follows from (1.1.4) that the empirical distribution function $F_n(x)$ is a strongly consistent estimator of the distribution function $F(x)$.

Consistent estimators with minimal variance. It is natural to compare unbiased estimators according to their variances. Let ξ be an observed random element with values in a measurable space (X, \mathcal{B}) and with a distribution belonging to a parametric family of probability measures $\{P_\theta, \theta \in \Theta\}$.

Consider the problem of estimating a real function $g(\theta)$. Consider the following classes of estimators:

- (1) the class U_g^θ of unbiased estimators $T = T(\xi)$ of the function $g(\theta)$ at a given point θ and such that $E_\theta T^2 < \infty$,
- (2) the class U_0^θ of unbiased estimators $T = T(\xi)$ of zero at a given point θ and such that $E_\theta T^2 < \infty$.

Therefore

$$U_g^\theta = \{T: E_\theta T = g(\theta), E_\theta T^2 < \infty\},$$

$$U_0^\theta = \{T: E_\theta T = 0, E_\theta T^2 < \infty\}.$$

We also consider the following classes:

$$U_g = \bigcap_{\theta \in \Theta} U_g^\theta, \quad U_0 = \bigcap_{\theta \in \Theta} U_0^\theta.$$

The following result contains necessary and sufficient conditions for an estimator to be optimal in the sense of minimum of the variance in the classes U_g^θ and U_g .

THEOREM 3.1.1. *The variance of an estimator $T \in U_g$ (respectively, $T \in U_g^\theta$) is minimal in the class U_g (respectively, in U_g^θ) if and only if $E_\theta T h = 0$ for all $h \in U_0$ and $\theta \in \Theta$ (respectively, $E_\theta T h = 0$ for a given $\theta \in \Theta$ and for all $h \in U_0^\theta$).*

PROOF. *Necessity.* Let an estimator $T \in U_g^\theta$ have minimal variance in the class U_g^θ for a fixed $\theta \in \Theta$ and $h \in U_0^\theta$. It is clear that $T + \lambda h \in U_g^\theta$ for all constants λ . Then

$$D_\theta(T + \lambda h) = D_\theta T + 2\lambda E_\theta T h + \lambda^2 D_\theta h.$$

If $E_\theta T h \neq 0$, then there exists a number λ such that

$$2\lambda E_\theta T h + \lambda^2 D_\theta h < 0.$$

This implies that $D_\theta(T + \lambda h) < D_\theta T$, which contradicts the assumption that the estimator T has minimal variance in the class U_g^θ . Therefore $E_\theta T h = 0$ for all $h \in U_0^\theta$.

Sufficiency. Let condition $E_\theta T h = 0$ hold for all $h \in U_0^\theta$ and for an estimator $T \in U_g^\theta$ where $\theta \in \Theta$ is fixed. Let $T' \in U_g^\theta$ be another estimator. Then

$$T' - T = h \in U_0^\theta$$

and

$$D_{\theta}T' = D_{\theta}T + 2E_{\theta}Th + D_{\theta}h \geq D_{\theta}T,$$

since $E_{\theta}Th = 0$. This means that the estimator T has minimal variance in the class U_g^{θ} .

The proof for the class U_g is analogous. \square

Theorem 3.1.1 is convenient for applications if a family of distributions contains a sufficiently wide class of unbiased estimators of zero.

Optimal estimators. As in the preceding section we consider the problem of estimating a function $g(\theta)$ from an observation ξ assuming values in (X, \mathcal{B}) and whose distribution belongs to a parametric family $\{P_{\theta}, \theta \in \Theta\}$.

Let $T = T(\xi)$ be an estimator of a function $g(\theta)$ and let $r(T, g)$ be a nonnegative loss function (a loss appears because we approximate $g = g(\theta)$ by an estimator T). A function

$$(3.1.14) \quad R(T; \theta) = E_{\theta}r(T(\xi), g(\theta)), \quad \theta \in \Theta,$$

is called a *risk function of an estimator* $T = T(\xi)$ if the true value of the parameter is θ .

Examples of loss functions are presented by the quadratic function

$$r(T, g) = (T - g)^2$$

used in the preceding section, by the Laplace function $r(T, g) = |T - g|$, and by the function

$$r(T, g) = \begin{cases} 0, & |T - g| \leq b, \\ 1, & |T - g| > b, \end{cases}$$

where $b > 0$. The latter function appears in the interval estimation of parameters.

Sometimes we treat a risk function as a measure of quality of estimators. We consider the general case where an estimator is not necessarily unbiased and the loss function is not necessarily quadratic.

An estimator $T' = T'(\xi)$ belonging to a class \mathcal{K} of estimators of a function $g(\theta)$ is called *admissible* for the class \mathcal{K} with respect to a loss function $r(T, g)$ if there is no estimator $T \in \mathcal{K}$ such that

$$(3.1.15) \quad R(T; \theta) \leq R(T'; \theta) \quad \text{for all } \theta \in \Theta$$

and inequality (3.1.15) is strict for at least one $\theta \in \Theta$. An estimator $T' = T'(\xi)$ that is admissible for the class of all estimators is called an *absolutely admissible* estimator of the function $g(\theta)$.

A statistic $T^* = T^*(\xi) \in \mathcal{K}$ is called an *optimal estimator of a function* $g(\theta)$ in the class \mathcal{K} with respect to a loss function $r(T, g)$ if for all $T \in \mathcal{K}$

$$R(T^*; \theta) \leq R(T; \theta) \quad \text{for all } \theta \in \Theta.$$

A statistic $T' = T'(\xi) \in \mathcal{K}$ is called an *optimal estimator* (or, *locally optimal estimator*) of a function $g(\theta)$ at a point θ_0 in the class \mathcal{K} with respect to a loss function $r(T, g)$ if for all $T \in \mathcal{K}$

$$R(T'; \theta_0) \leq R(T; \theta_0).$$

It is clear that the set of risk functions $R(T; \theta)$ is unordered in the class of all estimators T . For this reason we consider narrower classes of either estimators or

distributions. One of the possible approaches here is to exclude estimators that are not admissible.

EXAMPLE 3.1.4. Let $\xi^{(n)} = (\xi_1, \xi_2, \dots, \xi_n)$ be a sample where $\xi_1, \xi_2, \dots, \xi_n$ are independent identically distributed random variables depending on an unknown parameter θ . Assume that $\theta = E_\theta \xi_1$. Let $T = T(\xi^{(n)}) = \xi_1$. It is obvious that T is an unbiased estimator of the parameter θ . Let $E_\theta \xi_1^2 < \infty$ for all θ and let $r(T, \theta) = |T - \theta|^2$ be the loss function. Then $R(T; \theta) = E_\theta |\xi_1 - \theta|^2 = D_\theta \xi_1$. The estimator $T' = T'(\xi) = n^{-1} \sum_{i=1}^n \xi_i$ also is an unbiased estimator of the parameter θ . In this case

$$R(T'; \theta) = E_\theta |T' - \theta|^2 = \frac{1}{n} D_\theta \xi_1 \leq R(T; \theta).$$

Moreover the inequality becomes strict if $n \geq 2$. Therefore T is not an admissible estimator in the class of unbiased estimators of the parameter θ .

Below we consider the Bayes and minimax approaches allowing one to avoid the problem that the set of risk functions is unordered.

The Bayes and minimax approaches. Let Θ be an open set of \mathbf{R}^k and let \mathbf{Q} be a σ -finite measure on Θ . Without loss of generality we can extend the measure \mathbf{Q} to the whole space \mathbf{R}^k by putting $\mathbf{Q}(\mathbf{R}^k \setminus \Theta) = 0$. Let $T = T(\xi)$ be some estimator of a function $g(\theta)$ and let $R(T; \theta)$ be a risk function of the estimator T defined by (3.1.14). The number

$$(3.1.16) \quad R(T) = \int R(T; \theta) \mathbf{Q}(d\theta)$$

is called the *risk of the estimator* T . The measure \mathbf{Q} is called the *a priori measure*. An estimator T^* is called a *Bayes estimator* of a function $g(\theta)$ with respect to a loss function $r(T, g)$ and the a priori measure \mathbf{Q} if

$$R(T^*) \leq R(T)$$

for all estimators T where the risk $R(T)$ is defined by (3.1.16). In other words,

$$R(T^*) = \min_T R(T)$$

for a Bayes estimator.

Sometimes an estimator is called Bayes only in the case where \mathbf{Q} is a probability measure. Otherwise an estimator minimizing the risk (3.1.16) is called a *generalized Bayes estimator*.

Note that we can think of θ as a random parameter with distribution \mathbf{Q} if \mathbf{Q} is a probability measure. Then all Bayes estimators T^* are of the form

$$T^* = E \{ g(\theta) / \xi \}$$

in the case of $r(T, g) = (T - g)^2$ where the conditional expectation is evaluated with respect to the conditional distribution of the parameter θ subject to ξ . In its turn, the latter distribution can be found by the Bayes formula and this explains why these estimators are called Bayes. In this case, a Bayes estimator minimizes the possibility that the risk is the mean square error

$$R(T) = E(T(\xi) - g(\theta))^2 = ER(T; \theta).$$

Another approach is based on the comparison of maximums of risk functions for estimators $\sup_{\theta \in \Theta} R(T; \theta)$. A statistic $T' = T'(\xi)$ is called a *minimax* estimator of a function $g(\theta)$ with respect to a loss function $r(T, g)$ if for all estimators T

$$\sup_{\theta \in \Theta} R(T'; \theta) \leq \sup_{\theta \in \Theta} R(T; \theta).$$

In other words,

$$\sup_{\theta \in \Theta} R(T'; \theta) = \inf_T \sup_{\theta \in \Theta} R(T; \theta)$$

for a minimax estimator T' . There are many relations between minimax and Bayes estimators; some of them are given below.

THEOREM 3.1.2. *Let T^* be a Bayes estimator of a function $g(\theta)$ with respect to a loss function $r(T, g)$ and the a priori probability measure \mathbf{Q} . If there is an estimator T' such that*

$$(3.1.17) \quad R(T'; \theta) \leq \int R(T^*; t) \mathbf{Q}(dt)$$

for all $\theta \in \Theta$, then the estimator T' is minimax.

PROOF. Let T be an arbitrary estimator of a function $g(\theta)$. Then for all $t \in \Theta$

$$\sup_{\theta \in \Theta} R(T; \theta) \geq \int R(T; t) \mathbf{Q}(dt) \geq \int R(T^*; t) \mathbf{Q}(dt) \geq R(T'; t). \quad \square$$

Assume that the measure \mathbf{Q} possesses the density $q(t)$. Consider the set

$$N_{\mathbf{Q}} = \{t: q(t) > 0\}.$$

Note that inequality (3.1.17) becomes an equality for almost all $\theta \in N_{\mathbf{Q}}$, since otherwise

$$\int R(T'; \theta) q(\theta) d\theta < \int R(T^*; \theta) q(\theta) d\theta,$$

contradicting the assumption that the estimator T^* is Bayes. This remark allows one to obtain the following criterion, which is equivalent to Theorem 3.1.2.

THEOREM 3.1.3. *Assume that an estimator T exists such that*

- 1) T is a Bayes estimator with respect to some probability measure \mathbf{Q} possessing density $q(t)$;
- 2) $R(T; t) = c = \text{const}$ for $t \in N_{\mathbf{Q}}$;
- 3) $R(T; t) \leq c$ for $t \notin N_{\mathbf{Q}}$.

Then the estimator T is minimax.

If an estimator T' is minimax and is Bayes with respect to a probability measure \mathbf{Q} with density $q(t)$, then

$$\sup_{t \in \Theta} R(T'; t) = \int R(T'; t) q(t) dt.$$

Therefore any minimax estimator is a Bayes estimator that smooths the risk function. This means that the a priori measure \mathbf{Q}' related to this estimator suggests that statisticians pay the same attention to all possible parameters θ , instead of the approach suggested by Bayes estimators $T^* = T_{\mathbf{Q}}^*$ that corresponds to other a priori

measures $\mathbf{Q} \neq \mathbf{Q}'$, namely to pay special attention to some (the most probable) values of θ . Therefore

$$\int R(T_{\mathbf{Q}}^*; t) \mathbf{Q}(dt) \leq \int R(T'; t) \mathbf{Q}'(dt).$$

This inequality explains why the distribution \mathbf{Q}' in Theorem 3.1.3 corresponding to a minimax estimator T' is often called the *worse* or *least favorable*.

The least favorable distribution \mathbf{Q}' does not always exist, thus one can use the following criterion for minimax estimators.

THEOREM 3.1.4. *Assume that there are an estimator T' and a sequence of distributions \mathbf{Q}_m , $m = 1, 2, \dots$, possessing the densities $q_m(t)$ such that*

$$(3.1.18) \quad E_t r(T'; g(t)) \leq \limsup_{m \rightarrow \infty} \int E_t(T_m^*; g(t)) q_m(t) dt$$

for all $t \in \Theta$ where T_m^* is a Bayes estimator of a function $g(\theta)$ with respect to a loss function $r(T; g(t))$ and the a priori distribution \mathbf{Q}_m . Then the estimator T' is minimax.

PROOF. For all estimators T of a function $g(\theta)$,

$$\sup_t E_t r(T; g(t)) \geq \int E_t r(T; g(t)) q_m(t) dt \geq \int E_t r(T_m^*; g(t)) q_m(t) dt.$$

According to (3.1.18) this implies that

$$\sup_t E_t r(T; g(t)) \geq \limsup_{m \rightarrow \infty} \int E_t r(T_m^*; g(t)) q_m(t) dt \geq \sup_t E_t r(T'; g(t)),$$

whence it follows that the estimator T' is minimax. \square

EXAMPLE 3.1.5. Let $\xi^{(n)} = (\xi_1, \xi_2, \dots, \xi_n)$ be a sample from the normal $\mathcal{N}(\theta, 1)$ distribution. Let the a priori measure also be normal $\mathbf{Q}_m = \mathcal{N}(0, m)$ distribution where m is the variance. Then the Bayes estimator of the parameter θ with respect to the quadratic loss function and the a priori $\mathcal{N}(0, m)$ distribution coincides with the a posteriori mean $E(\theta / \xi^{(n)} = x) = \theta_m^*(x)$. Simple calculations show that

$$\theta_m^*(x) = \frac{1}{n} \sum_{i=1}^n x_i / \left(1 + \frac{1}{nm} \right).$$

The variance of the a posteriori distribution is

$$D(\theta / \xi^{(n)} = x) = \frac{m}{1 + nm},$$

whence it follows that the mean square error of the estimator θ_m^* is

$$E(\theta_m^* - \theta)^2 = ED(\theta / \xi^{(n)}) = \frac{m}{1 + nm} = \int E_t(\theta_m^* - t) q_m(t) dt$$

where $q_m(t)$ is the density of \mathbf{Q}_m with respect to the Lebesgue measure. This implies for the estimator $\tilde{\theta}_n(x) = n^{-1} \sum_{i=1}^n x_i$ that

$$E_t(\tilde{\theta}_n - t)^2 = \frac{1}{n} = \lim_{m \rightarrow \infty} \int E_t(\theta_m^* - t) q_m(t) dt.$$

Thus the estimator $\tilde{\theta}_n$ is minimax by Theorem 3.1.4. Note also that the least favorable distribution does not exist in this case.

We show in Example 3.2.1 below that $\tilde{\theta}_n$ is a Pitman estimator of the parameter θ . By (3.2.9) this estimator is Bayes with respect to the quadratic loss function and the Lebesgue measure taken as a priori measure.

The following is an example where the least favorable distribution exists.

EXAMPLE 3.1.6. Let $\xi^{(n)} = (\xi_1, \xi_2, \dots, \xi_n)$ be a sample from the binomial distribution with parameter $\theta \in \Theta = [0, 1]$, that is, the random variables ξ_1, \dots, ξ_n are independent, identically distributed, and assume values 1 and 0 with probabilities θ and $1 - \theta$, respectively. For the estimator $\tilde{\theta}_n(x) = n^{-1} \sum_{i=1}^n x_i$,

$$E_{\theta}(\tilde{\theta}_n - \theta)^2 = \theta(1 - \theta)/n.$$

Hence the assumption of Theorem 3.1.3 does not hold for this estimator. Consider another estimator

$$(3.1.19) \quad \theta_n^*(x) = \left(\tilde{\theta}_n(x) + \frac{1}{2\sqrt{n}} \right) \left(1 + \frac{1}{\sqrt{n}} \right)^{-1}$$

for which

$$E_{\theta}(\theta_n^* - \theta)^2 = \left(1 + \frac{1}{\sqrt{n}} \right)^{-2} E_{\theta} \left(\tilde{\theta}_n - \theta + \frac{1}{2\sqrt{n}} - \frac{\theta}{\sqrt{n}} \right)^2 = \frac{1}{4(1 + \sqrt{n})^2},$$

that is, the risk function of the estimator θ_n^* does not depend on θ .

Let B_{λ_1, λ_2} be the beta distribution with density

$$(3.1.20) \quad \frac{\Gamma(\lambda_1 + \lambda_2)}{\Gamma(\lambda_1)\Gamma(\lambda_2)} t^{\lambda_1-1}(1-t)^{\lambda_2-1}, \quad 0 \leq t \leq 1,$$

where $\lambda_1 > 0$ and $\lambda_2 > 0$ are two parameters of the distribution. Let the a priori distribution \mathbf{Q} be a Beta $B_{N+1, N+1}$ distribution. It can be proved in this case that the a priori distribution coincides with the Beta distribution (3.1.20) with parameters $\lambda_1 = N + n\tilde{\theta}_n(x) + 1$ and $\lambda_2 = N + n(1 - \tilde{\theta}_n(x)) + 1$. Since the mean value of the distribution B_{λ_1, λ_2} is $\lambda_1/(\lambda_1 + \lambda_2)$, the Bayes estimator with respect to the a priori distribution $\mathbf{Q} = B_{N+1, N+1}$ and the quadratic loss function is

$$\theta_n^{\mathbf{Q}}(x) = \frac{N + n\tilde{\theta}_n(x) + 1}{2N + n + 2} = \frac{\tilde{\theta}_n(x) + (N + 1)/n}{1 + 2(N + 1)/n}.$$

If $N + 1 = \sqrt{n}/2$, then the latter estimator coincides with the estimator $\theta_n^*(x)$ defined by (3.1.19). By Theorem 3.1.3 θ_n^* is minimax. On the other hand, it is known that this is a Bayes estimator with respect to the a priori distribution

$$\mathbf{Q} = B_{N+1, N+1}$$

for $N = \sqrt{n}/2 - 1$. This means that a priori distribution is the least favorable. If n increases, then the support of this distribution tends to concentrate in a neighborhood of the least favorable value of the parameter $\theta = 1/2$ for which the variance $\theta(1 - \theta)/n = 1/(4n)$ of the estimator $\tilde{\theta}_n$ is maximal. The estimator $\tilde{\theta}_n$ itself is not minimax, since

$$\sup_{\theta} \frac{\theta(1 - \theta)}{n} = \frac{1}{4n} > \frac{1}{4(1 + \sqrt{n})^2}.$$

It is also clear that for all θ outside a small enough neighborhood of the point $\theta = 1/2$ the estimator $\tilde{\theta}_n$ is better than θ_n^Q . The small neighborhood of the point $\theta = 1/2$ mentioned above is determined by the inequality

$$\theta(1 - \theta) < \frac{1}{4(1 + 1/\sqrt{n})^2}.$$

REMARK 3.1.1. In the general case, it is not always possible to give explicit expressions for Bayes and minimax estimators.

3.2. Estimation of a location parameter

In this and the next sections we show how to estimate unknown location and scale parameters of a distribution in the cases where optimal estimators exist.

Location parameters. Equivariant estimators. Let an observed element $\xi^{(n)} = (\xi_1, \dots, \xi_n)$ be a vector whose coordinates are, generally speaking, dependent random variables, let $(\mathbf{R}^n, \mathcal{B}^n)$ be a sampling space, and let $(P_\theta^n, \theta \in \Theta)$ be a family of distributions to which the distribution of the vector $\xi^{(n)}$ belongs. Assume that the parameter θ is one-dimensional and $\Theta = \mathbf{R}^1$. If measures P_θ^n depend on the parameter θ in such a way that

$$(3.2.1) \quad P_\theta^n(A) = P_0^n(A - \theta) \quad \text{for all } A \in \mathcal{B}^n,$$

then θ is called a *location parameter*. We use the notation

$$A - \theta = \{x - \theta = (x_1 - \theta, \dots, x_n - \theta) : x = (x_1, \dots, x_n) \in A\}$$

for all $A \in \mathcal{B}^n$ in equality (3.2.1). One of the models leading to distributions (3.2.1) is the so-called scheme of direct observations where

$$\xi_i = \theta + \varepsilon_i, \quad i = 1, \dots, n,$$

and $\varepsilon_1, \dots, \varepsilon_n$ are, generally speaking, dependent random variables with the joint distribution defined by the measure P_0^n .

Let $P_\theta\{\xi^{(n)} \in A\} = P_\theta^n(A)$ for all $A \in \mathcal{B}^n$. Then condition (3.2.1) can be rewritten in an equivalent form as

$$(3.2.2) \quad P_\theta\{\xi^{(n)} - \theta \in A\} = P_0\{\xi^{(n)} \in A\} \quad \text{for all } A \in \mathcal{B}^n$$

where $\xi^{(n)} - \theta = (\xi_1 - \theta, \dots, \xi_n - \theta)$. Condition (3.2.2) means that if the true value of the parameter is θ , then the vector $\xi^{(n)} - \theta$ has the same distribution as the vector $\xi^{(n)}$ corresponding to the zero value of the parameter.

There is a natural class of estimators used in the estimation of location parameters, namely

$$(3.2.3) \quad \mathcal{T} = \left\{ \tilde{\theta}_n = \tilde{\theta}_n(x) : \tilde{\theta}_n(x + c) = \tilde{\theta}_n(x) + c \text{ for all } x \in \mathbf{R}^n \text{ and all } c \in \mathbf{R}^1 \right\}.$$

Estimators of the class \mathcal{T} are called *equivariant* estimators of a location parameter. Some authors call such estimators "invariant".

Let $r(\tilde{\theta}_n; \theta) = r(\tilde{\theta}_n - \theta)$ be a nonnegative loss function depending on the difference of arguments $\tilde{\theta}_n - \theta$, and let $R(\tilde{\theta}_n; \theta) = E_\theta r(\tilde{\theta}_n - \theta)$ be the risk function of

the estimator $\tilde{\theta}_n$ (here E_θ stands for the expectation with respect to the measure P_θ). If $\tilde{\theta}_n \in \mathcal{T}$, then

$$(3.2.4) \quad R(\tilde{\theta}_n; \theta) = E_\theta r(\tilde{\theta}_n - \theta) = E_\theta r\left(\tilde{\theta}_n\left(\xi^{(n)} - \theta\right)\right) = E_0 r(\tilde{\theta}_n) = \text{const}$$

by (3.2.2), that is, the risk function $R(\tilde{\theta}_n; \theta)$ does not depend on θ . Thus the estimator $\tilde{\theta}_n \in \mathcal{T}$ is either optimal in \mathcal{T} or not admissible in \mathcal{T} for such loss functions.

An optimal in the class \mathcal{T} estimator $\hat{\theta}_n$ is called the *Pitman estimator* of a location parameter θ corresponding to a loss function $r(\hat{\theta}_n - \theta)$ if

$$R(\hat{\theta}_n; \theta) = \min_{\tilde{\theta}_n \in \mathcal{T}} R(\tilde{\theta}_n; \theta) \quad \text{for all } \theta \in \mathbf{R}^1.$$

Below we show that the Pitman estimator exists and find it for some loss functions.

We mention another useful property of equivariant estimators. Let $\tilde{\theta}_n$ and $\tilde{\theta}'_n$ be two equivariant estimators. Then by the definition of equivariant estimators

$$(3.2.5) \quad \tilde{\theta}_n(x) - \tilde{\theta}'_n(x) = \psi(y), \quad y = (x_2 - x_1, x_3 - x_1, \dots, x_n - x_1),$$

where $\psi(y)$ is some measurable function. Indeed,

$$\begin{aligned} \tilde{\theta}_n(x) - \tilde{\theta}'_n(x) &= [\tilde{\theta}_n(x) - x_1] - [\tilde{\theta}'_n(x) - x_1] = \tilde{\theta}_n(x - x_1) - \tilde{\theta}'_n(x - x_1) \\ &= \tilde{\theta}_n(0, x_2 - x_1, \dots, x_n - x_1) - \tilde{\theta}'_n(0, x_2 - x_1, \dots, x_n - x_1), \end{aligned}$$

whence (3.2.5) follows.

The Pitman estimator of a location parameter. In what follows we use the following notation. If $T = T(x)$, $x \in \mathbf{R}^n$, is a statistic, then $E_\theta(T/y)$ denotes the conditional expectation

$$(3.2.6) \quad E_\theta(T/y) = E_\theta \left\{ T\left(\xi^{(n)}\right) / \eta = y \right\}$$

where $\xi^{(n)} = (\xi_1, \xi_2, \dots, \xi_n)$, the vector y is defined by (3.2.5), and

$$(3.2.7) \quad \eta = (\xi_2 - \xi_1, \xi_3 - \xi_1, \dots, \xi_n - \xi_1).$$

The following result establishes the Pitman estimator of a location parameter with respect to the quadratic loss function.

THEOREM 3.2.1. *Assume that $E_0 \xi_i^2 < \infty$ for all $i = 1, 2, \dots, n$. Let*

$$l(x) = \sum_{i=1}^n c_i x_i, \quad x = (x_1, \dots, x_n),$$

be a linear statistic such that $\sum_{i=1}^n c_i = 1$. Then

1) the estimator

$$(3.2.8) \quad \hat{\theta}_n(x) = l(x) - E_0(l/y), \quad x \in \mathbf{R}^n,$$

is the Pitman estimator of a location parameter θ with respect to the quadratic loss function $r(\tilde{\theta}_n, \theta) = |\tilde{\theta}_n - \theta|^2$;

2) if the measure P_0 is absolutely continuous with respect to the Lebesgue measure and its density is $f(x)$, $x \in \mathbf{R}^n$, then the Pitman estimator is of the form

$$(3.2.9) \quad \widehat{\theta}_n(x) = \int_{-\infty}^{\infty} v f(x-v) dv \left(\int_{-\infty}^{\infty} f(x-v) dv \right)^{-1}, \quad x \in \mathbf{R}^n.$$

PROOF. It is clear that $\widehat{\theta}_n \in \mathcal{T}$. Let $\widetilde{\theta}_n$ be an arbitrary equivariant estimator of \mathcal{T} . Then $\widetilde{\theta}_n(x) = \widehat{\theta}_n(x) + \psi(y)$ by (3.2.5). If $E_0 \widetilde{\theta}_n^2 = \infty$, then

$$E_{\theta}(\widehat{\theta}_n - \theta)^2 \leq E_{\theta}(\widetilde{\theta}_n - \theta)^2$$

for all $\theta \in \mathbf{R}^1$. On the other hand, if $E_0 \widetilde{\theta}_n^2 < \infty$, then

$$(3.2.10) \quad \begin{aligned} E_{\theta}(\widetilde{\theta}_n - \theta)^2 &= E_{\theta}(\widetilde{\theta}_n - \widehat{\theta}_n + \widehat{\theta}_n - \theta)^2 \\ &= E_{\theta}(\widehat{\theta}_n - \theta)^2 + 2E_{\theta}(\widetilde{\theta}_n - \widehat{\theta}_n)(\widehat{\theta}_n - \theta) + E_{\theta}(\widetilde{\theta}_n - \widehat{\theta}_n)^2. \end{aligned}$$

It follows from (3.2.2) and (3.2.5) that

$$(3.2.11) \quad \begin{aligned} E_{\theta}(\widetilde{\theta}_n - \widehat{\theta}_n)(\widehat{\theta}_n - \theta) &= E_0(\widetilde{\theta}_n - \widehat{\theta}_n)\widehat{\theta}_n = E_0 E_0 \{ (\widetilde{\theta}_n - \widehat{\theta}_n)\widehat{\theta}_n / \eta \} \\ &= E_0(\widetilde{\theta}_n - \widehat{\theta}_n)E_0 \{ \widehat{\theta}_n / \eta \} = 0 \end{aligned}$$

in view of (3.2.6) and (3.2.7) where $\widehat{\theta}_n$ is defined by (3.2.8). Thus

$$E_0(\widehat{\theta}_n/y) = E_0(l/y) - E_0(l/y) = 0.$$

Relations (3.2.10) and (3.2.11) imply that for all $\theta \in \mathbf{R}^1$ we have

$$E_{\theta}(\widetilde{\theta}_n - \theta)^2 = E_{\theta}(\widehat{\theta}_n - \theta)^2 + E_{\theta}(\widetilde{\theta}_n - \widehat{\theta}_n)^2 \geq E_{\theta}(\widehat{\theta}_n - \theta)^2.$$

Therefore we proved that the estimator $\widehat{\theta}_n$ defined by (3.2.8) is the Pitman estimator of a location parameter with respect to the quadratic loss function.

Now we prove equality (3.2.9). Let $l(x) = x_1$. Then estimator (3.2.8) is of the form

$$(3.2.12) \quad \widehat{\theta}_n(x) = x_1 - E_0(\xi_1/y), \quad x \in \mathbf{R}^n.$$

Consider random variables $\zeta_1 = \xi_1$, $\zeta_2 = \xi_2 - \xi_1, \dots, \zeta_n = \xi_n - \xi_1$. Let $p(z_1, z_2, \dots, z_n)$ be the probability density of the vector $(\zeta_1, \zeta_2, \dots, \zeta_n)$ for $\theta = 0$. It is not hard to show that

$$p(z_1, z_2, \dots, z_n) = f(z_1, z_1 + z_2, \dots, z_1 + z_n)$$

where $f(x)$, $x \in \mathbf{R}^n$, is the density of the distribution P_0 . It is obvious that

$$\begin{aligned}
 E_0(\xi_1/y) &= E_0 \{ \zeta_1 / \zeta_2 = x_2 - x_1, \dots, \zeta_n = x_n - x_1 \} \\
 &= \int_{-\infty}^{\infty} zp(z, x_2 - x_1, \dots, x_n - x_1) dz \\
 &\quad \times \left(\int_{-\infty}^{\infty} p(z, x_2 - x_1, \dots, x_n - x_1) dz \right)^{-1} \\
 (3.2.13) \quad &= \int_{-\infty}^{\infty} zf(z, z + x_2 - x_1, \dots, z + x_n - x_1) dz \\
 &\quad \times \left(\int_{-\infty}^{\infty} f(z, z + x_2 - x_1, \dots, z + x_n - x_1) dz \right)^{-1} \\
 &= x_1 + \int_{-\infty}^{\infty} vf(x_1 - v, x_2 - v, \dots, x_n - v) dv \\
 &\quad \times \left(\int_{-\infty}^{\infty} f(x_1 - v, x_2 - v, \dots, x_n - v) dv \right)^{-1}
 \end{aligned}$$

Combining equalities (3.2.12) and (3.2.13) we obtain representation (3.2.9). \square

REMARK 3.2.1. When proving representation (3.2.9) we put $l(x) = x_1$ in (3.2.8). Note in the general case that

$$E_0(l/y) = l(x) - x_1 + E_0(\xi_1/y), \quad x \in \mathbf{R}^n,$$

where we used the property that $\sum_{i=1}^n c_i = 1$. Thus estimator (3.2.8) becomes of the form $\hat{\theta}_n(x) = x_1 - E_0(\xi_1/y)$ and this result is used in the above proof.

REMARK 3.2.2. If the measure P_0 is absolutely continuous with respect to the Lebesgue measure, then it is seen from equality (3.2.9) that the Pitman estimator is the Bayes estimator with respect to the quadratic loss function and the a priori measure \mathbf{Q} coinciding with the Lebesgue measure. In other words, the Pitman estimator is a generalized Bayes estimator.

REMARK 3.2.3. In fact, statement 1) of Theorem 3.2.1 is a particular case of Theorem 3.1.1. Namely, the optimality of the estimator $\hat{\theta}_n$ in the class \mathcal{T} means that $\hat{\theta}_n$ is orthogonal to unbiased estimators of zero and the latter are of the form $h(x_2 - x_1, \dots, x_n - x_1)$.

The following result establishes the Pitman estimator of a location parameter with respect to the Laplace loss function $r(\hat{\theta}, \theta) = |\hat{\theta} - \theta|$. We use the notation $\text{med}_0(l/y)$ for a median of the conditional distribution $l(\xi)$ in the case of $\theta = 0$ given condition $\eta = y$ for which $\text{med}_0(l/y)$ is a statistic.

THEOREM 3.2.2. Let $E_0|\zeta_i| < \infty$ for all $i = 1, 2, \dots, n$. Let $l(x) = \sum_{i=1}^n c_i x_i$, $x = (x_1, x_2, \dots, x_n)$, be some linear statistic such that $\sum_{i=1}^n c_i = 1$. Then

$$(3.2.14) \quad \bar{\theta}_n(x) = l(x) - \text{med}_0(l/y), \quad x \in \mathbf{R}^n,$$

is the Pitman estimator of a location parameter with respect to the Laplace loss function.

PROOF. Let $\tilde{\theta}_n \in \mathcal{T}$. Then $\tilde{\theta}_n(x) = l(x) + \psi(y)$ by (3.2.5), since $l \in \mathcal{T}$. Thus

$$(3.2.15) \quad \begin{aligned} E_\theta |\tilde{\theta}_n - \theta| &= E_0 |l + \psi(\eta)| = E_0 E_0 (|l + \psi(\eta)|/\eta) \\ &\geq E_0 E_0 (|l - \text{med}_0(l/\eta)|/\eta) = E_0 |\bar{\theta}_n| = E_\theta |\bar{\theta}_n - \theta|, \end{aligned}$$

since for all y

$$(3.2.16) \quad E_0 (|l + \psi(y)|/y) \geq E_0 (|l - \text{med}_0(l/y)|/y).$$

It follows from (3.2.15) that estimator (3.2.14) is the Pitman estimator of a location parameter with respect to the Laplace loss function. \square

REMARK 3.2.4. In the proof of Theorem 3.2.2 we used inequality (3.2.16) which is a well-known property of a median of a distribution (see [9]).

Below are some examples of the evaluation of Pitman estimators.

EXAMPLE 3.2.1. Let $\xi^{(n)} = (\xi_1, \xi_2, \dots, \xi_n)$ where $\xi_1, \xi_2, \dots, \xi_n$ are independent identically distributed $\mathcal{N}(\theta, 1)$ random variables. Then the assumption of statement 2) of Theorem 3.2.1 holds, and moreover

$$f(x) = (2\pi)^{-n/2} \exp \left\{ -2^{-1} \sum_{i=1}^n x_i^2 \right\}, \quad x = (x_1, \dots, x_n).$$

Substituting this density into (3.2.9), we show that the Pitman estimator is of the form $\hat{\theta}_n(x) = n^{-1} \sum_{i=1}^n x_i$. This implies that $D_\theta \hat{\theta}_n = E_\theta (\hat{\theta}_n - \theta)^2 = n^{-1}$.

EXAMPLE 3.2.2. Let $\xi^{(n)} = (\xi_1, \xi_2, \dots, \xi_n)$ where $\xi_1, \xi_2, \dots, \xi_n$ are independent exponential random variables with the density $e^{-x+\theta}$, $x \geq \theta$. Again the assumption of statement 2) of Theorem 3.2.1 holds, and moreover

$$f(x) = \exp \left\{ - \sum_{i=1}^n x_i \right\} I_{[0, \infty)} \left(\min_{1 \leq i \leq n} x_i \right), \quad x = (x_1, \dots, x_n).$$

Substituting this density into (3.2.9) we obtain for the Pitman estimator

$$\hat{\theta}_n(x) = \min_{1 \leq i \leq n} x_i - \frac{1}{n}, \quad x = (x_1, \dots, x_n).$$

Using (1.3.3) for the density of the first order statistic $\zeta_{1,n}$ we obtain that

$$D_\theta \hat{\theta}_n = E_\theta (\hat{\theta}_n - \theta)^2 = 2n^{-2}.$$

This example is remarkable, since the mean square error of the Pitman estimator is of order n^{-2} which is higher than that in the preceding example. This can be explained by the discontinuity of the density with respect to the parameter.

The optimal estimator of a location parameter in the class of linear unbiased estimators. Let $\xi^{(n)} = (\xi_1, \dots, \xi_n)$ where ξ_1, \dots, ξ_n are independent random variables with the distribution functions $F_1(x - \theta), \dots, F_n(x - \theta)$, respectively. Assume that

$$(3.2.17) \quad \int x dF_j(x) = 0, \quad j = 1, \dots, n,$$

$$(3.2.18) \quad 0 < \int x^2 dF_j(x) = \sigma_j^2 < \infty, \quad j = 1, \dots, n.$$

Thus θ is a location parameter and $\theta = E_{\theta}\xi_i$, $i = 1, \dots, n$, by (3.2.17). Let L be the class of linear unbiased estimators l of the parameter θ that are of the form $l(x) = \sum_{i=1}^n c_i x_i$ where $\sum_{i=1}^n c_i = 1$, $x = (x_1, \dots, x_n)$. For any estimator $l \in L$ we have

$$(3.2.19) \quad E_{\theta}(l - \theta)^2 = E_{\theta} \left(\sum_{i=1}^n c_i (\xi_i - \theta) \right)^2 = \sum_{i=1}^n c_i^2 \sigma_i^2.$$

Solving the extremum problem for the function (3.2.19) subject to $\sum_{j=1}^n c_j = 1$ we find that

$$(3.2.20) \quad c_j = c_j^* = \sigma_j^{-2} \left(\sum_{i=1}^n \sigma_i^{-2} \right)^{-1}, \quad j = 1, \dots, n,$$

at the point of extremum.

Thus we have proved the following result.

THEOREM 3.2.3. *Let conditions (3.2.17) and (3.2.18) hold. Then the optimal estimator of a location parameter in the class L of linear unbiased estimators is given by $l^*(x) = \sum_{j=1}^n c_j^* x_j$, $x = (x_1, \dots, x_n)$, where the coefficients c_j^* , $j = 1, \dots, n$, are defined by (3.2.20).*

The following result contains necessary and sufficient conditions that an estimator $l^*(x)$ is admissible in the class of unbiased estimators.

THEOREM 3.2.4. *Assume that ξ_1, \dots, ξ_n , $n \geq 3$, are independent observations with the distribution functions $F_1(x - \theta), \dots, F_n(x - \theta)$, respectively, for which conditions (3.2.17) and (3.2.18) hold. An optimal estimator of a parameter θ in the class of linear unbiased estimators $l^*(x) = \sum_{i=1}^n c_i^* x_i$ is admissible for the quadratic loss function in the class of all unbiased estimators of the parameter θ if and only if all the distribution functions $F_j(x)$ are normal.*

The proof of this theorem can be found in [15], Theorem 7.4.1. The book [15] contains further results on the estimation of a location parameter.

REMARK 3.2.5. If conditions (3.2.17) and (3.2.18) hold and moreover the variances σ_j^2 are equal to each other, that is, $\sigma_1^2 = \sigma_2^2 = \dots = \sigma_n^2 = \sigma^2$, then $c_j^* = n^{-1}$, $j = 1, \dots, n$, in (3.2.20). Thus the optimal estimator of the location parameter in the class L of linear unbiased estimators is of the form $l^*(x) = n^{-1} \sum_{i=1}^n x_i$ and this obviously is the sampling mean.

REMARK 3.2.6. If condition (3.2.18) holds, while condition (3.2.17) does not hold, then one can consider the class of linear unbiased estimators l of the form $l(x) = \sum_{j=1}^n c_j(x_j - b_j)$ where $\sum_{j=1}^n c_j = 1$ and $b_j = \int x dF_j(x)$, $j = 1, \dots, n$. As before we obtain that the optimal estimator is of the form $l^*(x) = \sum_{j=1}^n c_j^*(x_j - b_j)$ where the coefficients c_j^* are defined by (3.2.20).

REMARK 3.2.7. If the assumptions of Theorem 3.2.4 hold, then the optimal estimator θ in the class of linear unbiased estimators $l^*(x) = \sum_{j=1}^n c_j^* x_j$ is absolutely admissible with respect to the quadratic loss function if and only if all the distribution functions $F_j(x)$ are normal (see [15], Theorem 7.4.2). Note that if all the functions $F_j(x)$ are equal to each other, then the optimal estimator is of the form $l^*(x) = n^{-1} \sum_{i=1}^n x_i$ and therefore this estimator is absolutely admissible in the case of Gaussian distributions.

The problem of the confidence estimation of a location parameter. Let $\xi^{(n)} = (\xi_1, \dots, \xi_n)$ be an observation where ξ_1, \dots, ξ_n are independent identically distributed random variables with the distribution function $F(x - \theta)$, $\theta \in \mathbf{R}^1$. Consider the problem of the confidence estimation of the parameter θ with respect to the loss function

$$(3.2.21) \quad r(\tilde{\theta}_n, \theta) = \begin{cases} 0, & |\tilde{\theta}_n - \theta| \leq b, \\ 1, & |\tilde{\theta}_n - \theta| > b, \end{cases}$$

where b is some positive constant. The corresponding risk function of the estimator $\tilde{\theta}_n$ is of the form

$$R(\tilde{\theta}_n; \theta) = E_{\theta} r(\tilde{\theta}_n; \theta) = P_{\theta} \{ |\tilde{\theta}_n - \theta| > b \} = P_{\theta} \{ \theta \notin [\tilde{\theta}_n - b, \tilde{\theta}_n + b] \}.$$

This risk function is the probability of the event that the confidence interval

$$[\tilde{\theta}_n - b, \tilde{\theta}_n + b]$$

does not contain the unknown value θ .

Given an arbitrary $F(x)$, the risk function

$$R(\bar{\xi}^{(n)}; \theta) = 1 - P_{\theta} \{ -b \leq \bar{\xi}^{(n)} - \theta \leq b \} = \gamma, \quad \gamma = \gamma(b),$$

for the estimator $\bar{\xi}^{(n)} = n^{-1} \sum_{i=1}^n \xi_i$ does not depend on θ and therefore

$$[\bar{\xi}^{(n)} - b, \bar{\xi}^{(n)} + b]$$

is a confidence interval of level γ (depending, of course, on $F(x)$).

The following result claims that the sampling mean $\bar{\xi}^{(n)}$ is admissible with respect to the loss function (3.2.21) (in other words, it claims that the confidence interval $[\bar{\xi}^{(n)} - b, \bar{\xi}^{(n)} + b]$ is admissible).

THEOREM 3.2.5. *Let $\xi^{(n)} = (\xi_1, \dots, \xi_n)$, $n \geq 3$, where ξ_1, \dots, ξ_n are independent identically distributed random variables with the distribution function $F(x - \theta)$ whose density $f(x - \theta)$ is bounded. If, for a given sequence of numbers $b_j \rightarrow 0$, the sampling mean $\bar{\xi}^{(n)}$ is an admissible estimator of the parameter $\theta \in \mathbf{R}^1$ with respect to the loss function (3.2.21) for $b = b_j$, then $F(x)$ is the Gaussian distribution function.*

The proof of Theorem 3.2.5 can be found in [31] (see Theorem 7.9.3 therein).

REMARK 3.2.8. If random variables $\xi_1, \xi_2, \dots, \xi_n$ are independent and identically distributed according to the normal $\mathcal{N}(\theta, \sigma^2)$ law, then one can show that the sampling mean $\bar{\xi}^{(n)}$ is an admissible estimator of the location parameter θ with respect to the loss function (3.2.21) for all $b > 0$. Thus the converse statement to Theorem 3.2.5 is also true.

The minimax property of the Pitman estimator of a location parameter. As mentioned above, the Pitman estimator of a location parameter with respect to the quadratic loss function is a Bayes estimator of the location parameter with respect to the Lebesgue a priori measure \mathbf{Q} if there exists the density of the observation. If all the assumptions of Theorem 3.2.1 hold, then the Pitman

estimator $\widehat{\theta}_n$ of the location parameter θ with respect to the quadratic loss function is of the form (3.2.9).

Let the a priori measure \mathbf{Q}_N coincide with the uniform distribution on $[-N, N]$, so that the density is $q_N(t) = (2N)^{-1}I_{[-N, N]}(t)$. Then the Bayes estimator of the parameter θ with respect to the a priori measure \mathbf{Q}_N and the quadratic loss function is given by

$$\theta_{\mathbf{Q}_N}^*(x) = \frac{\int v f(x-v) q_N(v) dv}{\int f(x-v) q_N(v) dv} = \int_{-N}^N v f(x-v) dv \left(\int_{-N}^N f(x-v) dv \right)^{-1},$$

$x \in \mathbf{R}^n$, where $f(x)$ is the density mentioned in statement 2) of Theorem 3.2.1. It is clear that $\widehat{\theta}_n(x) = \lim_{N \rightarrow \infty} \theta_{\mathbf{Q}_N}^*(x)$ for all $x \in \mathbf{R}^n$. One can show that

$$(3.2.22) \quad \lim_{N \rightarrow \infty} E_{\theta} (\theta_{\mathbf{Q}_N}^* - \theta)^2 = E_{\theta} (\widehat{\theta}_n - \theta)^2$$

for all $\theta \in [-N + \sqrt{N}, N - \sqrt{N}]$. Moreover the convergence is uniform in θ in the above interval. Since $E_{\theta} (\widehat{\theta}_n - \theta)^2$ does not depend on θ and convergence (3.2.22) is uniform in the interval $[-N + \sqrt{N}, N - \sqrt{N}]$, we get

$$\begin{aligned} \limsup_{N \rightarrow \infty} \int E_t (\theta_{\mathbf{Q}_N}^* - t)^2 \mathbf{Q}_N(dt) &\geq \limsup_{N \rightarrow \infty} \frac{1}{2N} \int_{-N+\sqrt{N}}^{N-\sqrt{N}} E_t (\theta_{\mathbf{Q}_N}^* - t)^2 dt \\ &\geq \limsup_{N \rightarrow \infty} \frac{2(N-\sqrt{N})}{2N} (E_{\theta} (\widehat{\theta}_n - \theta)^2 - \varepsilon) \\ &= E_{\theta} (\widehat{\theta}_n - \theta)^2 - \varepsilon \end{aligned}$$

for all $\varepsilon > 0$. This implies for all $\theta \in \mathbf{R}^1$ that

$$E_{\theta} (\widehat{\theta}_n - \theta)^2 \leq \limsup_{N \rightarrow \infty} \int E_t (\theta_{\mathbf{Q}_N}^* - t)^2 \mathbf{Q}_N(dt).$$

By Theorem 3.1.4 this means that the Pitman estimator $\widehat{\theta}_n$ is minimax. Thus we proved the following result.

THEOREM 3.2.6. *If all the assumptions of Theorem 3.2.1 hold, then the Pitman estimator $\widehat{\theta}_n$ of a location parameter θ with respect to the quadratic loss function is minimax.*

3.3. Estimation of a scale parameter

Scale parameters. Equivariant estimators. Let $\xi^{(n)} = (\xi_1, \xi_2, \dots, \xi_n)$ be an observed random element assuming values in the space $(\mathbf{R}^k, \mathcal{B}^k)$ and having the distribution belonging to a parametric set of measures $(P_{\sigma}, \sigma \in (0, \infty))$ where σ is an unknown parameter. If the measure P_{σ} depends on the parameter σ such that

$$(3.3.1) \quad P_{\sigma}(A) = P_1(A/\sigma), \quad A \in \mathcal{B}^k,$$

then σ is called a *scale parameter*. We put

$$A/\sigma = \{x/\sigma = (x_1/\sigma, \dots, x_n/\sigma) : x = (x_1, \dots, x_n) \in A\}$$

for all $A \in \mathcal{B}^k$. Condition (3.3.1) can be rewritten as

$$(3.3.2) \quad P_{\sigma} \{ \xi^{(n)} \in A \} = P_{\sigma} \{ \xi^{(n)}/\sigma \in A/\sigma \} = P_1 \{ \xi^{(n)} \in A/\sigma \}, \quad A \in \mathcal{B}^k.$$

Distributions (3.3.1) arise in the case where observations are of the form $\xi_i = \sigma \varepsilon_i$, $i = 1, \dots, n$, for some $\sigma > 0$ and if the vector $(\varepsilon_1, \dots, \varepsilon_n)$ has the distribution defined by a measure P_1 . Generally speaking, the random variables $\varepsilon_1, \dots, \varepsilon_n$ are dependent, thus the random variables ξ_1, \dots, ξ_n are dependent, too.

It is natural to consider the following class of estimators in the case of the estimation of a scale parameter σ :

$$(3.3.3) \quad \Sigma = \{ \tilde{\sigma}_n(x) : \sigma_n(\lambda x) = \lambda \tilde{\sigma}_n(x) \text{ for all } \lambda > 0 \text{ and } x \in \mathbf{R}^n \}$$

where we put $\lambda x = (\lambda x_1, \dots, \lambda x_n)$ for all $x = (x_1, \dots, x_n)$. Estimators of the class Σ are called *equivariant estimators*. Consider a loss function $r(\tilde{\sigma}_n; \sigma)$ such that

$$(3.3.4) \quad r(\tilde{\sigma}_n; \sigma) = r(\tilde{\sigma}_n - \sigma), \quad r(\lambda u) = \lambda^m r(u)$$

for all $\lambda > 0$ and some $m > 0$. Using (3.3.2) and the definition (3.3.3) we get that the risk of any estimator $\tilde{\sigma}_n \in \Sigma$ is such that

$$R(\tilde{\sigma}_n; \sigma) = E_\sigma r(\tilde{\sigma}_n - \sigma) = \sigma^m E_1 r(\tilde{\sigma}_n - 1) = \sigma^m R(\tilde{\sigma}_n; 1).$$

Therefore an estimator $\tilde{\sigma}_n \in \Sigma$ is either optimal in the class Σ or is not admissible in this class provided the loss function satisfies condition (3.3.4).

An optimal estimator $\hat{\sigma}_n \in \Sigma$ of a parameter σ in the class Σ , that is, the one such that

$$R(\hat{\sigma}_n; \sigma) = \min_{\tilde{\sigma}_n \in \Sigma} R(\tilde{\sigma}_n; \sigma) \quad \text{for all } \sigma \in (0, \infty),$$

is called the *Pitman estimator of a scale parameter σ with respect to the loss function* (3.3.4).

In what follows we assume that an observed random element $\xi^{(n)}$ assumes values in the space

$$\mathbf{R}_{+,0}^n = \mathbf{R}_+^n \setminus \{0\} = \{x = (x_1, \dots, x_n) : x_i > 0 \text{ for all } i = 1, \dots, n\}.$$

Alternatively one can think that the sampling space is \mathbf{R}^n and assume that

$$P_\sigma(\mathbf{R}_{+,0}^n) = 1$$

for all $\sigma > 0$.

Equivariant estimators possess the following useful property. If

$$l(x) = \sum_{j=1}^n c_j x_j, \quad x \in \mathbf{R}_{+,0}^n,$$

is a linear statistic with $c_j > 0$ for all $j = 1, 2, \dots, n$, then

$$(3.3.5) \quad \tilde{\sigma}_n(x) = l(x)\psi(y), \quad x = (x_1, \dots, x_n) \in \mathbf{R}_{+,0}^n,$$

for all $\tilde{\sigma}_n \in \Sigma$ where $\psi(y)$ is some measurable function of the vector

$$(3.3.6) \quad y = \left(\frac{x_2}{x_1}, \frac{x_3}{x_1}, \dots, \frac{x_n}{x_1} \right).$$

Using (3.3.3) we also get

$$\tilde{\sigma}_n(x) = l(x)\tilde{\sigma}_n(x/l(x)) = l(x)\tilde{\sigma}_n((x/x_1)/l(x/x_1)),$$

since $l(x) > 0$, $x \in \mathbf{R}_{+,0}^n$, whence (3.3.5) follows for the function $\psi(y)$ specified above.

Pitman estimators. We use the following notation for an arbitrary statistic $T(x)$, $x \in \mathbf{R}^n$:

$$(3.3.7) \quad E_1(T/y) = E_1 \left\{ T \left(\xi^{(n)} \right) / \eta = y \right\}$$

where y is the vector defined by (3.3.6), while the vector η is given by

$$(3.3.8) \quad \eta = \left(\frac{\xi_2}{\xi_1}, \frac{\xi_3}{\xi_1}, \dots, \frac{\xi_n}{\xi_1} \right).$$

The following result provides the explicit form of Pitman estimators of a scale parameter with respect to the quadratic loss function.

THEOREM 3.3.1. *Let $E_1 \xi_j^2 < \infty$ for all $j = 1, 2, \dots, n$ and let $l(x) = \sum_{j=1}^n c_j x_j$, $x \in \mathbf{R}_{+,0}^n$, be some linear statistic with $c_j > 0$, $j = 1, 2, \dots, n$. Then*

1) *the estimator*

$$(3.3.9) \quad \hat{\sigma}_n(x) = l(x) \frac{E_1(l/y)}{E_1(l^2/y)}, \quad x \in \mathbf{R}_{+,0}^n,$$

is the Pitman estimator of a scale parameter σ with respect to the quadratic loss function where y is the vector given by (3.3.6), and $E_1(l/y)$ and $E_1(l^2/y)$ are conditional expectations defined by (3.3.7) and (3.3.8);

2) *if the measure P_1 is absolutely continuous with respect to the Lebesgue measure and its density is $f(x)$, $x \in \mathbf{R}_{+,0}^n$, then the Pitman estimator is of the form*

$$(3.3.10) \quad \hat{\sigma}_n(x) = \int_0^\infty u^n f(ux) du \left(\int_0^\infty u^{n+1} f(ux) du \right)^{-1}, \quad x \in \mathbf{R}_{+,0}^n.$$

PROOF. Let $\tilde{\sigma}_n$ be an arbitrary estimator of the class Σ . Using (3.3.5) we get

$$E_\sigma(\tilde{\sigma}_n - \sigma)^2 = \sigma^2 E_1(l\psi(\eta) - 1)^2 = \sigma^2 E_1 E_1((l\psi(\eta) - 1)^2/\eta).$$

If $\eta = y$ is fixed, then

$$\min_c E_1((lc - 1)^2/y) = E_1((lc^* - 1)^2/y)$$

where

$$c^* = \psi^*(y) = \frac{E_1(l/y)}{E_1(l^2/y)}.$$

Thus

$$\min_\psi E_\sigma(l\psi(\eta) - \sigma)^2 = E_\sigma(l\psi^*(\eta) - \sigma)^2 = E_\sigma(\hat{\sigma}_n - \sigma)^2$$

where $\hat{\sigma}_n$ is the estimator defined by (3.3.9). Therefore estimator (3.3.9) is the Pitman estimator of the scale parameter σ with respect to the quadratic loss function.

Now we prove equality (3.3.10). Let $l(x) = n^{-1} \sum_{i=1}^n x_i = \bar{x}$. Then the Pitman estimator (3.3.9) can be rewritten as follows (here we put $\bar{\xi}^{(n)} = n^{-1} \sum_{i=1}^n \xi_i$):

$$(3.3.11) \quad \begin{aligned} \hat{\sigma}_n(x) &= \bar{x} \frac{E_1(\bar{\xi}^{(n)}/y)}{E_1\left(\left(\bar{\xi}^{(n)}\right)^2/y\right)} = \bar{x} \frac{E_1\left(\xi_1\left(\bar{\xi}^{(n)}/\xi_1\right)/y\right)}{E_1\left(\xi_1^2\left(\bar{\xi}^{(n)}/\xi_1\right)^2/y\right)} \\ &= \bar{x} \frac{(\bar{x}/x_1)E_1(\xi_1/y)}{(\bar{x}/x_1)^2 E_1(\xi_1^2/y)} = x_1 \frac{E_1(\xi_1/y)}{E_1(\xi_1^2/y)}. \end{aligned}$$

Consider the random variables $\beta_1 = \xi_1, \beta_2 = \xi_2/\xi_1, \dots, \beta_n = \xi_n/\xi_1$ and denote by $p(z_1, z_2, \dots, z_n)$ the probability density of the vector $(\beta_1, \beta_2, \dots, \beta_n)$ for $\sigma = 1$. It is not hard to show that

$$p(z_1, z_2, \dots, z_n) = z_1^{n-1} f(z_1, z_1 z_2, \dots, z_1 z_n)$$

where $f(x)$ is the probability density of P_1 with respect to the Lebesgue measure. It is clear that

$$\begin{aligned} & E_1(\xi_1^m/y) \\ &= E_1\left\{\beta_1^m / \beta_2 = \frac{x_2}{x_1}, \dots, \beta_n = \frac{x_n}{x_1}\right\} \\ &= \int_0^\infty z^m p\left(z, \frac{x_2}{x_1}, \dots, \frac{x_n}{x_1}\right) dz \left(\int_0^\infty p\left(z, \frac{x_2}{x_1}, \dots, \frac{x_n}{x_1}\right) dz\right)^{-1} \\ (3.3.12) \quad &= \int_0^\infty z^{m+n-1} f\left(z, z \frac{x_2}{x_1}, \dots, z \frac{x_n}{x_1}\right) dz \\ &\quad \times \left(\int_0^\infty z^{n-1} f\left(z, z \frac{x_2}{x_1}, \dots, z \frac{x_n}{x_1}\right) dz\right)^{-1} \\ &= x_1^m \int_0^\infty u^{m+n-1} f(x_1 u, \dots, x_n u) du \\ &\quad \times \left(\int_0^\infty u^{n-1} f(x_1 u, \dots, x_n u) du\right)^{-1}. \end{aligned}$$

Substituting (3.3.12) into (3.3.11) we obtain (3.3.10). \square

Let Σ_u be the class of unbiased equivariant estimators of a scale parameter σ . It is clear that $\Sigma_u \subset \Sigma$. The following result provides the explicit form of an optimal estimator in the class Σ_u with respect to the quadratic loss function.

THEOREM 3.3.2. *Let all the assumptions of Theorem 3.3.1 hold. An optimal in the class Σ_u estimator of a scale parameter σ with respect to the quadratic loss function is of the form*

$$(3.3.13) \quad \hat{\sigma}_{n,u} = C_u \hat{\sigma}_n$$

where the constant C_u is such that $C_u E_1 \hat{\sigma}_n = 1$.

PROOF. Let $\tilde{\sigma}_n$ be an arbitrary estimator of the class Σ_u . If $E_1 \tilde{\sigma}_n^2 = \infty$, then $E_\sigma \tilde{\sigma}_n^2 = \infty$ for all $\sigma > 0$. Thus the estimator $\tilde{\sigma}_n$ is worse than $\hat{\sigma}_{n,u}$. If $E_1 \tilde{\sigma}_n^2 < \infty$, then

$$\begin{aligned} (3.3.14) \quad E_\sigma(\tilde{\sigma}_n - \sigma)^2 &= \sigma^2 E_1(\tilde{\sigma}_n - 1)^2 = \sigma^2 E_1(\tilde{\sigma}_n - \hat{\sigma}_{n,u} + \hat{\sigma}_{n,u} - 1)^2 \\ &= \sigma^2 \{E_1(\tilde{\sigma}_n - \hat{\sigma}_{n,u})^2 \\ &\quad + 2E_1(\tilde{\sigma}_n - \hat{\sigma}_{n,u})(\hat{\sigma}_{n,u} - 1) + E_1(\hat{\sigma}_{n,u} - 1)^2\}. \end{aligned}$$

Since $\tilde{\sigma}_n \in \Sigma_u$ and $\hat{\sigma}_{n,u} \in \Sigma_u$, we have $E_1(\tilde{\sigma}_n - \hat{\sigma}_{n,u}) = 0$. Since $\hat{\sigma}_n$ is optimal in the class Σ , we obtain $E_1 \hat{\sigma}_n h = 0$ for an arbitrary unbiased estimator of zero $h \in \Sigma$ such that $E_1 h^2 < \infty$ (the proof is the same as that in Theorem 3.1.1). In particular, $E_1 \hat{\sigma}_{n,u}(\tilde{\sigma}_n - \hat{\sigma}_{n,u}) = 0$. Thus (3.3.14) implies

$$E_\sigma(\tilde{\sigma}_n - \sigma)^2 = E_\sigma(\tilde{\sigma}_n - \hat{\sigma}_{n,u})^2 + E_\sigma(\hat{\sigma}_{n,u} - \sigma)^2 \geq E_\sigma(\hat{\sigma}_{n,u} - \sigma)^2. \quad \square$$

Below is an example of the evaluation of a Pitman estimator.

EXAMPLE 3.3.1. Let $\xi^{(n)} = (\xi_1, \dots, \xi_n)$ where ξ_1, \dots, ξ_n are independent identically distributed exponential random variables. Then their common probability density is $\sigma^{-1} \exp(-x/\sigma)$, $x > 0$, $\sigma > 0$. Assumptions of assertion 2) in Theorem 3.3.1 are satisfied and

$$f(x) = \exp\left(-\sum_{i=1}^n x_i\right) I_{[0, \infty)}\left(\min_{1 \leq i \leq n} x_i\right), \quad x = (x_1, \dots, x_n).$$

Substituting this density into (3.3.10), the Pitman estimator of the parameter σ becomes of the form

$$\hat{\sigma}_n(x) = \frac{1}{n+1} \sum_{i=1}^n x_i, \quad x = (x_1, \dots, x_n).$$

Since $E_1 \hat{\sigma}_n = n/(n+1)$, Theorem 3.3.2 implies that the optimal estimator of the parameter σ in the class of the unbiased equivariant estimator Σ_u is of the form

$$\hat{\sigma}_{n,u}(x) = \frac{1}{n} \sum_{i=1}^n x_i, \quad x = (x_1, \dots, x_n).$$

Straightforward calculations yield

$$E_\sigma(\hat{\sigma}_n - \sigma)^2 = \frac{\sigma^2}{n+1} < E_\sigma(\hat{\sigma}_{n,u} - \sigma)^2 = \frac{\sigma^2}{n},$$

that is, the estimator $\hat{\sigma}_n$ is better than $\hat{\sigma}_{n,u}$, and this result is natural.

The optimal estimator of a scale parameter in the class of linear unbiased estimators. Let $\xi^{(n)} = (\xi_1, \dots, \xi_n)$ where ξ_1, \dots, ξ_n are nondegenerate independent random variables with the distribution functions $F_1(x/\sigma), \dots, F_n(x/\sigma)$, respectively. The distribution functions depend on a scale parameter $\sigma \in (0, \infty)$ and are such that

$$(3.3.15) \quad F_j(0+) = 0, \quad E_1 \xi_j^2 < \infty, \quad j = 1, \dots, n.$$

Let L be the class of linear unbiased estimators of the parameter σ of the form $l(x) = \sum_{i=1}^n c_i x_i$, $x = (x_1, \dots, x_n)$, where $c_i > 0$ for all $i = 1, \dots, n$. Put

$$\alpha_{1j} = E_1 \xi_j, \quad \alpha_{2j} = E_1 \xi_j^2, \quad \sigma_j^2 = \alpha_{2j} - \alpha_{1j}^2.$$

Since ξ_j are nondegenerate, it follows that $\sigma_j^2 > 0$. If $l(x) = \sum_{j=1}^n c_j x_j$, then $\sum_{j=1}^n c_j \alpha_{1j} = 1$. Further $E_\sigma(l - \sigma)^2 = \sigma^2 \sum_{j=1}^n c_j^2 \sigma_j^2$. This implies that the optimal estimator of the parameter σ in the class L with respect to the quadratic loss function is of the form $l^*(x) = \sum_{j=1}^n c_j^* x_j$ where the coefficients c_j^* are such that

$$(3.3.16) \quad c_j^* = \frac{\alpha_{1j}}{\sigma_j^2} \left(\sum_{i=1}^n \frac{\alpha_{1i}}{\sigma_i^2} \right)^{-1}, \quad j = 1, \dots, n.$$

Therefore we have proved the following result.

THEOREM 3.3.3. *Let random variables ξ_1, \dots, ξ_n be independent and nondegenerate with the distribution functions $F_1(x/\sigma), \dots, F_n(x/\sigma)$, respectively. If condition (3.3.15) holds, then the optimal linear unbiased estimator of the parameter σ with respect to the quadratic loss function is of the form $l^*(x) = \sum_{j=1}^n c_j^* x_j$ where the coefficients c_j^* are defined by (3.3.16).*

The following result provides necessary and sufficient conditions for an optimal linear unbiased estimator $l^*(x) = \sum_{j=1}^n c_j^* x_j$ of a scale parameter σ to be admissible.

THEOREM 3.3.4. *Let all the assumptions of Theorem 3.3.3 hold. Then an optimal linear unbiased estimator $l^*(x) = \sum_{i=1}^n c_i^* x_i$ of the parameter σ is admissible in the class of unbiased estimators with respect to the quadratic loss function for some two different values of n , say $n = n_1$ and $n = n_2$, $n_2 > n_1 \geq 3$, if and only if the random variables ξ_j have the Gamma distribution*

$$F_j(x) = \frac{\alpha_j^{\gamma_j}}{\Gamma(\gamma_j)} \int_0^x t^{\gamma_j-1} e^{-\alpha_j t} dt, \quad x > 0,$$

for some $\gamma_j > 0$ and $\alpha_j > 0$, $j = 1, \dots, n$.

The proof of Theorem 3.3.4 can be found in [15] (see Theorem 7.12.2 therein).

REMARK 3.3.1. If an estimator $l^*(x)$ is admissible in the class of unbiased estimators of the scale parameter σ , then by Theorem 3.3.2 it is optimal in the class Σ_u . Moreover by Theorem 3.3.4 the distribution of the random variables ξ_j is the Gamma distribution in this case (see Example 3.3.1).

REMARK 3.3.2. Further results concerning the estimation of a scale parameter can be found in [15].

3.4. The Cramér–Rao inequality and efficient estimators

In the preceding sections we studied the quality of statistical estimators of unknown parameters and obtained several qualitative results. Moreover we introduced two classes of parameters, namely the classes of scale and location parameters, for which one can construct optimal estimators in appropriate classes of estimators. In this section we use a somewhat different approach to construct optimal estimators. We also obtain the minimal mean square error of the estimation that can be achieved in an experiment.

Regularity conditions for families of distributions. Let ξ be an observation that is a random element assuming values in a measurable space (X, \mathcal{B}) and whose distribution belongs to a parametric set $\{P_\theta, \theta \in \Theta\}$ where Θ is a subset of \mathbf{R}^k , $k \geq 1$. Throughout this section we assume that for all $\theta \in \Theta$ the measure P_θ is absolutely continuous with respect to some σ -finite measure μ on (X, \mathcal{B}) , that is, $P_\theta \ll \mu$, and that $f(x; \theta)$ is the density of the measure P_θ with respect to the measure μ . In particular, if $(X, \mathcal{B}) = (\mathbf{R}^m, \mathcal{B}^m)$ for some $m \geq 1$, then as the measure μ one can take the Lebesgue measure.

We consider the case of a one-dimensional parameter θ in this section, that is, we consider the case $k = 1$.

Below we use the following set of *regularity conditions*, called (CR):

- (i) Θ is a finite or infinite interval in \mathbf{R}^1 ;

- (ii) the derivative $\partial f(x; \theta)/\partial \theta$ exists and is finite P_θ -almost everywhere for all $\theta \in \Theta$;
- (iii) $|\partial^i f(x; \theta)/\partial \theta^i| \leq g_i(x)$ for all $\theta \in \Theta$ and $i = 1, 2$ where $g_1(x)$ and $g_2(x)$ are nonnegative Borel functions such that $\int g_i(x) \mu(dx) < \infty$, $i = 1, 2$;
- (iv) $0 < E_\theta(\partial \ln f(\xi; \theta)/\partial \theta)^2 < \infty$ for all $\theta \in \Theta$.

The regularity conditions (CR) are also called the *Cramér–Rao regularity conditions*. If conditions (CR) hold, then the family of distributions $\{P_\theta, \theta \in \Theta\}$ is called *CR-regular*.

Let $S(x; \theta) = \partial \ln f(x; \theta)/\partial \theta$, $x \in X$, and put

$$(3.4.1) \quad I(\theta) = E_\theta S^2(\xi; \theta), \quad \theta \in \Theta.$$

The function $I(\theta)$ is called the *Fisher information*. This function is treated as the amount of information about the parameter θ contained in the observation ξ . The notion of the information $I(\theta)$ will become clear after the proof of the Cramér–Rao inequality. Note that condition (iv) above can be rewritten as follows: $0 < I(\theta) < \infty$ for all $\theta \in \Theta$.

First we prove an auxiliary result.

LEMMA 3.4.1. *If regularity conditions (CR) hold, then*

$$(3.4.2) \quad E_\theta S(\xi; \theta) = 0 \quad \text{for all } \theta \in \Theta.$$

PROOF. First we differentiate the identity

$$\int f(x; \theta) \mu(dx) = 1 \quad \text{for all } \theta \in \Theta$$

with respect to θ and obtain

$$(3.4.3) \quad \frac{\partial}{\partial \theta} \int f(x; \theta) \mu(dx) = 0 \quad \text{for all } \theta \in \Theta.$$

Conditions (ii) and (iii) allow one to interchange the differentiation and integration in (3.4.3). Thus

$$(3.4.3') \quad \int \frac{\partial}{\partial \theta} f(x; \theta) \mu(dx) = 0 \quad \text{for all } \theta \in \Theta,$$

whence

$$E_\theta S(\xi; \theta) = \int \frac{\partial}{\partial \theta} f(x; \theta) \mu(dx) = 0 \quad \text{for all } \theta \in \Theta. \quad \square$$

The following result gives another representation for the Fisher information.

LEMMA 3.4.2. *If regularity conditions (CR) hold, then*

$$(3.4.4) \quad I(\theta) = -E_\theta \left(\frac{\partial^2}{\partial \theta^2} \ln f(\xi; \theta) \right)$$

for all $\theta \in \Theta$.

PROOF. It is obvious that

$$(3.4.5) \quad \begin{aligned} E_{\theta} \left(\frac{\partial^2}{\partial \theta^2} \ln f(\xi; \theta) \right) &= E_{\theta} \left(\frac{\partial^2 f(\xi; \theta) / \partial \theta^2}{f(\xi; \theta)} - \left(\frac{\partial f(\xi; \theta) / \partial \theta}{f(\xi; \theta)} \right)^2 \right) \\ &= \int \frac{\partial^2}{\partial \theta^2} f(x; \theta) \mu(dx) - I(\theta). \end{aligned}$$

It follows from (iv) that $I(\theta) < \infty$ for all $\theta \in \Theta$. Condition (iii) implies that

$$(3.4.6) \quad \frac{\partial^2}{\partial \theta^2} \int f(x; \theta) \mu(dx) = \int \frac{\partial^2}{\partial \theta^2} f(x; \theta) \mu(dx) = 0$$

for all $\theta \in \Theta$. Combining (3.4.5) and (3.4.6) we obtain (3.4.4). \square

Sometimes we also consider the following set of regularity conditions, called (R) :

- (i') Θ is a finite or infinite interval in \mathbf{R}^1 ;
- (ii') the function $(f(x; \theta))^{1/2}$ is continuously differentiable with respect to θ for μ -almost all x ;
- (iii') $0 < E_{\theta} S^2(\xi; \theta) = I(\theta) < \infty$ for all $\theta \in \Theta$ and the function $I(\theta)$ is continuous with respect to θ .

In what follows we need an assertion on the continuity of integrals of functions depending on a parameter.

Let $\psi(t, y)$, $t \in \Theta$, be a family of measurable functions defined on a measurable space (Y, \mathcal{B}_Y) equipped with a measure ν . We consider some conditions under which

$$(3.4.7) \quad \int \psi(t, y) \nu(dy) \rightarrow \int \psi(\theta, y) \nu(dy), \quad t \rightarrow \theta.$$

Let $A(t) = A(t, \theta)$, $t \in \Theta$, be a family of sets belonging to \mathcal{B}_Y . Put

$$\bar{A}(t) = Y \setminus A(t).$$

The following result is a generalization of a well-known Lebesgue theorem.

LEMMA 3.4.3. *Let $A(t)$, $t \in \Theta$, be a family such that*

- 1) $\psi(t, y) I_{A(t)}(y) \rightarrow \psi(\theta, y)$ as $t \rightarrow \theta$ for ν -almost all y for which $\psi(\theta, y) \neq 0$;
- 2) $\sup_t |\psi(t, y) I_{A(t)}(y)| \leq \psi(y)$ where $\psi(y)$ is a function integrable with respect to the measure ν , that is, $\int \psi(y) \nu(dy) < \infty$.

Then relation (3.4.7) holds if and only if

$$(3.4.8) \quad \int \psi(t, y) I_{\bar{A}(t)}(y) \nu(dy) \rightarrow 0, \quad t \rightarrow \theta.$$

PROOF. By the Lebesgue theorem

$$\int \psi(t, y) I_{A(t)}(y) \nu(dy) \rightarrow \int \psi(\theta, y) \nu(dy), \quad t \rightarrow \theta.$$

Since

$$\int \psi(t, y) \nu(dy) = \int \psi(t, y) I_{A(t)}(y) \nu(dy) + \int \psi(t, y) I_{\bar{A}(t)}(y) \nu(dy),$$

relation (3.4.7) is equivalent to relation (3.4.8). \square

COROLLARY 3.4.1. Let $T(x)$ be a real measurable bounded function, $T: X \rightarrow \mathbf{R}^1$, and let $f(x; \theta)$ be continuous with respect to θ for μ -almost all $x \in X$. Then the function $E_\theta T(\xi)$ is continuous with respect to θ .

PROOF. We apply Lemma 3.4.3 for $Y = X$, $\nu = \mu$, $\psi(t, x) = T(x)f(x; t)$, and $A(t) = \{x: f(x; t) \leq 2f(x; \theta)\}$. It is obvious that all the assumptions of Lemma 3.4.3 are satisfied. Since $T(x) \equiv 1$ and thus $E_\theta T(\xi) \equiv 1$ is continuous, Lemma 3.4.3 yields

$$\int f(x; t) I_{\bar{A}(t)} \mu(dx) \rightarrow 0$$

as $t \rightarrow \theta$. This together with Lemma 3.4.3 implies that $E_\theta T(\xi)$ is continuous for any bounded function $T(x)$. \square

REMARK 3.4.1. If one seeks a simple sufficient condition for (3.4.7) in the case of $\psi(t, y) \rightarrow \psi(\theta; y)$ as $t \rightarrow \theta$ and ν -almost surely, then an appropriate candidate is the uniform convergence of integrals in (3.4.7). The latter condition can be reformulated as follows: there exists a finite measure λ such that the inequality $\lambda(A) < \delta = \delta(\varepsilon)$ implies $\sup_t \int_A |\psi(t, y)| \nu(dy) < \varepsilon$ for a given $\varepsilon > 0$. Moreover if the function $\psi(y) = \sup_t |\psi(t, y)|$ is integrable, then one can take $\lambda(A) = \int_A \psi(y) \nu(dy)$.

Below we consider some corollaries of conditions (R) that we will use in the proof of the Cramér–Rao inequality.

LEMMA 3.4.4. Let conditions (R) hold. Let $T = T(\xi)$ be an arbitrary real statistic such that $E_\theta T^2 < c < \infty$ for all $\theta \in \Theta$. Then the function $a(\theta) = E_\theta T$ is differentiable with respect to θ , and moreover

$$(3.4.9) \quad a'(\theta) = E_\theta T(\xi) S(\xi; \theta) = \int T(x) \frac{\partial}{\partial \theta} f(x; \theta) \mu(dx)$$

where $S(x; \theta) = \partial \ln f(x; \theta) / \partial \theta$, $x \in X$, and the function $a'(\theta)$ is continuous.

PROOF. In the same manner as in the proof of Lemma 3.4.1 we derive from condition (ii)' that

$$(3.4.10) \quad E_\theta S(\xi; \theta) = \int \frac{\partial}{\partial \theta} f(x; \theta) \mu(dx) = 0, \quad \theta \in \Theta.$$

Note that (3.4.10) also follows from (3.4.9) for $T(x) \equiv 1$. Then

$$(3.4.11) \quad D_\theta S(\xi; \theta) = E_\theta S^2(\xi; \theta) = I(\theta) = 4 \int \left(\frac{\partial}{\partial \theta} \sqrt{f(x; \theta)} \right)^2 \mu(dx).$$

The function $I(\theta)$ is continuous by conditions (R). We apply Lemma 3.4.3 for $Y = X$, $\nu = \mu$, and $\psi(t, x) = (\partial \sqrt{f(x; t)} / \partial \theta)^2$:

$$(3.4.12) \quad A(t) = A_1(\delta) = \left\{ x : \sup_{|v-\theta|<|\delta|} \sqrt{f(x; v)} < 2\sqrt{f(x; \theta)}, \right. \\ \left. \sup_{|v-\theta|<|\delta|} \left| \frac{\partial}{\partial v} \sqrt{f(x; v)} \right| \leq 2 \left| \frac{\partial}{\partial \theta} \sqrt{f(x; \theta)} \right| \right\}$$

where $\delta = t - \theta$. Assumptions of Lemma 3.4.3 hold for $\psi(x) = 4\psi(\theta; x)$, since the functions $\sqrt{f(x; \theta)}$ and $\partial\sqrt{f(x; \theta)}/\partial\theta$ are continuous. Since $I(t)$ converges to $I(\theta)$ as $t \rightarrow \theta$, we prove from (3.4.8) that

$$(3.4.13) \quad \varepsilon(t) = \int \left(\frac{\partial}{\partial\theta} \sqrt{f(x; t)} \right)^2 I_{A_1(\delta)}(x) \mu(dx) \rightarrow 0$$

as $t \rightarrow \theta$.

Similarly to the proof of Corollary 3.4.1 we prove that $\int T(x)\partial f(x; \theta)/\partial\theta \mu(dx)$ is a continuous function. In the proof of this result we apply Lemma 3.4.3 for $Y = X$, $\nu = \mu$, $\psi(t, x) = T(x)\partial f(x; t)/\partial t$, and $A(t) = A_1(\delta)$. Since

$$\partial f(x, \theta)/\partial\theta = 2\sqrt{f(x, \theta)}\partial\sqrt{f(x, \theta)}/\partial\theta,$$

we get $\sup_t |\psi(t, x)|I_{A(t)}(x) \leq \psi(x) = 4|\psi(\theta; x)|$. The Cauchy-Bunyakovskii inequality implies that

$$\int_{A_1(\delta)} |\psi(\theta, x)| \mu(dx) \leq 2 \left(E_\theta T^2(\xi) \int_{A_1(\delta)} (\partial\sqrt{f(x, \theta)}/\partial\theta)^2 \mu(dx) \right)^{1/2}$$

This together with (3.4.13) implies relation (3.4.8), whence it follows that the function $\int T(x)\partial f(x; \theta)/\partial\theta \mu(dx)$ is continuous.

Now we turn to the proof of equality (3.4.9). Note that

$$\begin{aligned} & \frac{1}{\delta} \left(\int T(x)f(x; \theta + \delta) \mu(dx) - \int T(x)f(x; \theta) \mu(dx) \right) \\ &= \int \int_0^1 T(x) \frac{\partial}{\partial\theta} f(x; \theta + u\delta) du \mu(dx) \\ &= \int \int_0^1 2T(x)\sqrt{f(x; \theta + u\delta)} \frac{\partial}{\partial\theta} \sqrt{f(x; \theta + u\delta)} du \mu(dx) \end{aligned}$$

by condition (ii)'. We apply Lemma 3.4.3 again for $Y = R \times X$, $y = (u, x)$, $\nu = \lambda \times \mu$, $\psi(\delta, y) = T(x)\partial f(x; \theta + u\delta)/\partial\theta$, and $A(\delta) = A_1(\delta)$, where λ is the Lebesgue measure, $A_1(\delta)$ is defined by (3.4.12), and $\delta \rightarrow 0$. Since the functions $\sqrt{f(x; \theta)}$ and $\partial\sqrt{f(x; \theta)}/\partial\theta$ are continuous with respect to θ , assumptions 1) and 2) of Lemma 3.4.3 hold, whence

$$\begin{aligned} \psi(\delta, y)I_{A(\delta)}(x) &\rightarrow T(x) \frac{\partial}{\partial\theta} f(x; \theta) = \psi(0; y), \quad \delta \rightarrow 0, \\ \sup_\delta |\psi(\delta, y)I_{A(\delta)}(x)| &\leq 4T(x) \left| \frac{\partial}{\partial\theta} f(x; \theta) \right| \end{aligned}$$

and by the Cauchy-Bunyakovskii inequality

$$\int T(x) \left| \frac{\partial}{\partial\theta} f(x; \theta) \right| \mu(dx) \leq (E_\theta T^2(\xi) \cdot I(\theta))^{1/2} < \infty.$$

Using the Cauchy–Bunyakovskiĭ inequality again we obtain from relation (3.4.13) that

$$\begin{aligned} & \left| \int_{\bar{A}_1(\delta)} \int_0^1 T(x) \sqrt{f(x; \theta + u\delta)} \frac{\partial}{\partial \theta} \sqrt{f(x; \theta + u\delta)} du \mu(dx) \right| \\ & \leq \left[\int_{\bar{A}_1(\delta)} \int_0^1 T^2(x) f(x; \theta + u\delta) du \mu(dx) \right. \\ & \quad \left. \times \int_{\bar{A}_1(\delta)} \int_0^1 \left(\frac{\partial}{\partial \theta} \sqrt{f(x; \theta + u\delta)} \right)^2 du \mu(dx) \right]^{1/2} \\ & \leq \left(c \int_0^1 \varepsilon(\theta + u\delta) du \right)^{1/2} \rightarrow 0 \end{aligned}$$

as $\delta \rightarrow 0$, whence relation (3.4.8) follows. Thus we proved that the derivative $a'(\theta)$ exists and equality (3.4.9) holds. \square

LEMMA 3.4.5. *If the set Θ is compact and the function $\sqrt{f(x; \theta)}$ is continuously differentiable with respect to θ for μ -almost all x , then $I(\theta)$ is continuous if and only if*

$$(3.4.14) \quad \lim_{N \rightarrow \infty} \sup_{\theta} E_{\theta} S^2(\xi; \theta) I(|S(\xi; \theta)| > N) = 0.$$

PROOF. Let the function $I(\theta)$ be continuous but let relation (3.4.14) not hold. Then there is a $\gamma > 0$ and a sequence $t \rightarrow \theta \in \Theta$ such that $N_t \rightarrow \infty$ and

$$(3.4.15) \quad m(t) = E_t S^2(\xi; t) I(|S(\xi; t)| > N_t) > \gamma$$

for all t belonging to this sequence.

Applying Lemma 3.4.3 for $Y = X$ and $\mu = \nu$ we get

$$\begin{aligned} \psi(t, x) &= \left(\frac{\partial}{\partial t} \sqrt{f(x; t)} \right)^2 = 4^{-1} S^2(x; t) f(x; t), \\ A(t) &= \left\{ x: \left| \frac{\partial}{\partial t} \sqrt{f(x; t)} \right| \leq 2 \left| \frac{\partial}{\partial \theta} \sqrt{f(x; \theta)} \right| \right\}. \end{aligned}$$

Since $\partial \sqrt{f(x; \theta)} / \partial \theta$ is continuous, assumptions 1) and 2) of Lemma 3.4.3 hold and the continuity of $I(t)$ implies that

$$m_1(t) = \int_{\bar{A}(t)} \left| \frac{\partial}{\partial t} \sqrt{f(x; t)} \right|^2 \mu(dx) \rightarrow 0$$

as $t \rightarrow \theta$. Note that $m(t) \leq m_1(t) + m_2(t)$ where

$$\begin{aligned} m_2(t) &= \int_{B(t) \cap A(t)} \left(\frac{\partial}{\partial t} \sqrt{f(x; t)} \right)^2 \mu(dx), \\ B(t) &= \left\{ x: 2 \left| \frac{\partial}{\partial t} \sqrt{f(x; t)} \right| > N_t \sqrt{f(x; t)} \right\}. \end{aligned}$$

It follows from the definition of the set $A(t)$ that

$$m_2(t) \leq 4 \int_{B(t)} \left| \frac{\partial}{\partial \theta} \sqrt{f(x; \theta)} \right|^2 \mu(dx).$$

Since $\partial \sqrt{f(x; t)} / \partial t \rightarrow \partial \sqrt{f(x; \theta)} / \partial \theta$ and $\sqrt{f(x; t)} \rightarrow \sqrt{f(x; \theta)}$ as $t \rightarrow \theta$, we prove that $B(t)$ converges to a set whose measure μ is zero. This means that $\mu(B(t)) \rightarrow 0$, $m_2(t) \rightarrow 0$, and $m(t) \rightarrow 0$ as $t \rightarrow \theta$. This contradicts (3.4.15) and thus relation (3.4.14) holds.

Now we prove the converse statement. Let condition (3.4.14) hold. According to Lemma 3.4.3, $I(t)$ is continuous if $m_1(t) \rightarrow 0$ as $t \rightarrow \theta$ on the set $A(t)$ defined as above. Further

$$m_1(t) \leq \int_{|S(x; t)| > N} |S(x; t)|^2 f(x; t) \mu(dx) + N^2 \int_{\bar{A}(t)} f(x; t) \mu(dx)$$

where the first integral is small by (3.4.14) if N is sufficiently large. To estimate the second integral we note that $\mu(\bar{A}(t)) \rightarrow 0$ as $t \rightarrow \theta$ and

$$\int_{\bar{C}(t)} f(x; t) \mu(dx) \rightarrow 0, \quad t \rightarrow \theta,$$

for $C(t) = \{x: f(x; t) \leq 2f(x; \theta)\}$ (see the proof of Corollary 3.4.1). Thus as $t \rightarrow \theta$

$$\int_{\bar{A}(t)} f(x; t) \mu(dx) \leq 2 \int_{\bar{A}(t)} f(x; \theta) \mu(dx) + \int_{\bar{C}(t)} f(x; t) \mu(dx). \quad \square$$

REMARK 3.4.2. If the set Θ is compact and conditions (R) are satisfied, then, due to Lemma 3.4.5, the Fisher information $I(\theta)$ is continuous if and only if condition (3.4.14) holds. It is natural to call the latter condition the uniform convergence condition for the integral $I(\theta) = E_\theta S^2(\xi; \theta)$.

The Cramér–Rao inequality under regularity conditions (CR). The following result contains the lower bound for the variance of an unbiased estimator under the Cramér–Rao conditions (CR).

THEOREM 3.4.1. *Let the distributions P_θ , $\theta \in \Theta$, satisfy the regularity conditions (CR). Assume that $g(\theta)$ is a differentiable real function, $\hat{g} = \hat{g}(x)$ is an unbiased estimator of the function $g(\theta)$ such that the variance of \hat{g} exists, and*

$$\int \left| \hat{g}(x) \frac{\partial}{\partial \theta} \bar{f}(x; \theta) \right| \mu(dx) < \infty \quad \text{for all } \theta \in \Theta.$$

Then

$$(3.4.16) \quad D_\theta \bar{g}(\xi) \geq \frac{(g'(\theta))^2}{I(\theta)} \quad \text{for all } \theta \in \Theta.$$

Inequality (3.4.16) becomes an equality if and only if the density $f(x; \theta)$ is of the form

$$(3.4.17) \quad f(x; \theta) = \exp\{\psi_1(\theta)\hat{g}(x) + \psi_2(\theta) + h(x)\}, \quad x \in X,$$

where $\psi'_1(\theta) \neq 0$.

PROOF. According to the Cauchy–Bunyakovskiĭ inequality

$$(3.4.18) \quad E_{\theta} S(\xi; \theta) (\widehat{g}(\xi) - g(\theta)) \leq (I(\theta) D_{\theta} \widehat{g}(\xi))^{1/2}, \quad \theta \in \Theta.$$

Lemma 3.4.1 implies that

$$(3.4.19) \quad E_{\theta} S(\xi; \theta) (\widehat{g}(\xi) - g(\theta)) = E_{\theta} S(\xi; \theta) \widehat{g}(\xi) = \int \widehat{g}(x) \frac{\partial}{\partial \theta} f(x; \theta) \mu(dx).$$

Since \widehat{g} is an unbiased estimator, we have

$$g(\theta + \delta) - g(\theta) = \int \widehat{g}(x) (f(x; \theta + \delta) - f(x; \theta)) \mu(dx),$$

whence it follows by the regularity condition (iii) that

$$(3.4.20) \quad g'(\theta) = \int \widehat{g}(x) \frac{\partial}{\partial \theta} f(x; \theta) \mu(dx).$$

Thus (3.4.18)–(3.4.20) imply inequality (3.4.16).

It remains to consider the case of equality in (3.4.16). The inequality in (3.4.16) becomes an equality if and only if inequality (3.4.18) becomes an equality. In its turn (3.4.18) becomes an equality if and only if

$$S(x; \theta) = \widehat{g}(x) \psi(\theta) + \widetilde{\psi}(\theta)$$

for all $x \in X$ where $\psi(\theta) \neq 0$. This implies that

$$\ln f(x; \theta) = \widehat{g}(x) \psi_1(\theta) + \psi_2(\theta) + h(x)$$

for all $x \in X$ where $\psi'_1(\theta) \neq 0$. □

Results of the type (3.4.16) are called the *Cramér–Rao inequality*. This inequality gives a lower bound for the variance of an unbiased estimator of the function $g(\theta)$. If $\widehat{\theta}$ is an unbiased estimator of a parameter θ , then the Cramér–Rao inequality becomes of the form

$$(3.4.21) \quad D_{\theta} \widehat{\theta} \geq \frac{1}{I(\theta)}, \quad \theta \in \Theta.$$

Inequality (3.4.16) also gives a lower bound for the variance of estimators $\widehat{\theta}$ of a parameter θ that are not necessarily unbiased. Indeed let $g(\theta) = \theta + b(\theta)$ where $b(\theta)$ is the bias of the estimator $\widehat{\theta}$, that is, $b(\theta) = E_{\theta}(\widehat{\theta} - \theta)$. Assuming that the function $b(\theta)$ is differentiable, we obtain from (3.4.16) that

$$(3.4.22) \quad E_{\theta} (\widehat{\theta} - \theta)^2 \geq b^2(\theta) + \frac{(1 + b'(\theta))^2}{I(\theta)}.$$

Let an observation be a sample $\xi^{(n)} = (\xi_1, \xi_2, \dots, \xi_n)$ and let the random variables $\xi_1, \xi_2, \dots, \xi_n$ have the density $f(x; \theta)$, $x \in \mathbf{R}^1$, where θ is a real unknown parameter, $\theta \in \Theta \subset \mathbf{R}^1$. Assume that the density $f(x; \theta)$ satisfies the Cramér–Rao conditions (CR). Denote by $f_n(x; \theta)$, $x \in \mathbf{R}^n$, the density of the vector $\xi^{(n)}$. Let $I_n(\theta)$ be the Fisher information evaluated with respect to the density $f_n(x; \theta)$, while $I(\theta) = I_1(\theta)$ is the Fisher information evaluated with respect to the density $f(x; \theta)$, that is,

$$I_n(\theta) = E_{\theta} S_n^2(\xi^{(n)}; \theta) \quad \text{and} \quad I(\theta) = E_{\theta} S^2(\xi_1; \theta)$$

where

$$S_n(x; \theta) = \frac{\partial \ln f_n(x; \theta)}{\partial \theta} \quad \text{and} \quad S(x; \theta) = \frac{\partial \ln f(x; \theta)}{\partial \theta}.$$

LEMMA 3.4.6. *If the Cramér-Rao regularity conditions (CR) hold, then*

$$(3.4.23) \quad I_n(\theta) = nI(\theta).$$

PROOF. Since $f_n(x; \theta) = \prod_{k=1}^n f(x_k; \theta)$, $x = (x_1, \dots, x_n)$, and $E_\theta S(\xi_1; \theta) = 0$ by Lemma 3.4.1, we have

$$\begin{aligned} I_n(\theta) &= E_\theta \left(\frac{\partial}{\partial \theta} \ln f_n(\xi^{(n)}; \theta) \right)^2 \\ &= \sum_{k=1}^n E_\theta \left(\frac{\partial \ln f(\xi_k; \theta)}{\partial \theta} \right)^2 + \sum_{k \neq j} E_\theta \frac{\partial \ln f(\xi_k; \theta)}{\partial \theta} E_\theta \frac{\partial \ln f(\xi_j; \theta)}{\partial \theta} \\ &= nI(\theta). \end{aligned} \quad \square$$

COROLLARY 3.4.2. *Assume that all the assumptions of Theorem 3.4.1 hold. If $\widehat{g}_n(x)$ is an unbiased estimator of a function $g(\theta)$, then the Cramér-Rao inequality holds:*

$$(3.4.24) \quad D_\theta \widehat{g}_n(\xi) \geq \frac{(g'(\theta))^2}{nI(\theta)}, \quad \theta \in \Theta.$$

In particular, if $g(\theta) = \theta$ and $\widehat{\theta}_n$ is an unbiased estimator of a parameter θ , then

$$(3.4.25) \quad D_\theta \widehat{\theta}_n \geq \frac{1}{nI(\theta)}, \quad \theta \in \Theta.$$

The proof of Corollary 3.4.2 follows from Theorem 3.4.1 and equality (3.4.23). \square

Not all of the regularity conditions (CR) are used in the proof of the Cramér-Rao inequality. In fact, this result holds under a weaker set of regularity conditions called in what follows the Cramér-Rao conditions (CR)*:

- (i) Θ is a finite or infinite interval in \mathbf{R}^1 ;
- (ii)* the function $f(x; \theta)$ is differentiable with respect to θ for μ -almost all $x \in X$ and

$$0 < I(\theta) = \int \left(\frac{\partial}{\partial \theta} f(x; \theta) \right)^2 (f(x; \theta))^{-1} \mu(dx) < \infty;$$

- (iii)* the following relations are satisfied:

$$(3.4.26) \quad \begin{aligned} \lim_{\Delta \rightarrow 0} \frac{1}{\Delta^2} \int \frac{(f(x; \theta + \Delta) - f(x; \theta))^2}{f(x; \theta)} \mu(dx) \\ = \int \left(\frac{\partial f(x; \theta)}{\partial \theta} \right)^2 \frac{\mu(dx)}{f(x; \theta)}, \end{aligned}$$

$$(3.4.27) \quad \int \frac{\partial f(x; \theta)}{\partial \theta} \mu(dx) = \frac{\partial}{\partial \theta} \int f(x; \theta) \mu(dx).$$

In fact, conditions (3.4.26) and (3.4.27) are necessary to justify the interchange of integration and differentiation. It is easily seen that the regularity conditions $(CR)^*$ are less restrictive than conditions (CR) . Nevertheless the Cramér–Rao inequality can also be proved under conditions $(CR)^*$.

THEOREM 3.4.2. *If regularity conditions $(CR)^*$ hold, then*

$$(3.4.28) \quad D_{\theta}\hat{\theta} \geq 1/I(\theta), \quad \theta \in \Theta,$$

for any unbiased estimator $\hat{\theta} = \hat{\theta}(\xi)$ of a parameter θ .

PROOF. Since $\hat{\theta}$ is an unbiased estimator, we have

$$\Delta = \int (\hat{\theta}(x) - \theta)(f(x; \theta + \Delta) - f(x; \theta)) \mu(dx),$$

whence

$$\begin{aligned} \Delta^2 &= \left(\int (\hat{\theta}(x) - \theta) \sqrt{f(x; \theta)} \frac{f(x; \theta + \Delta) - f(x; \theta)}{\sqrt{f(x; \theta)}} \mu(dx) \right)^2 \\ &\leq D_{\theta}\hat{\theta} \int \frac{(f(x; \theta + \Delta) - f(x; \theta))^2}{f(x; \theta)} \mu(dx) \end{aligned}$$

by the the Cauchy–Bunyakovskiĭ inequality. Thus

$$D_{\theta}\hat{\theta} \geq \left(\frac{1}{\Delta^2} \int \frac{(f(x; \theta + \Delta) - f(x; \theta))^2}{f(x; \theta)} \mu(dx) \right)^{-1}$$

for all Δ . Passing to the limit as $\Delta \rightarrow 0$ we obtain inequality (3.4.28) in view of condition (3.4.26). \square

Let an observation be a sample $\xi^{(n)} = (\xi_1, \dots, \xi_n)$ with the density

$$f_n(x; \theta) = \prod_{i=1}^n f(x_i; \theta), \quad x = (x_1, \dots, x_n) \in \mathbf{R}^n, \quad t \in \Theta.$$

If the density $f(x; \theta)$ satisfies conditions (ii)* and (3.4.27), while the density $f_n(x; \theta)$ satisfies condition (3.4.26), then the Cramér–Rao inequality (3.4.25) holds for any unbiased estimator $\hat{\theta}_n$ of the parameter θ . To prove this result we use equality (3.4.23) that follows from condition (3.4.27).

The Cramér–Rao inequality under the regularity conditions (R) . The following result contains a lower bound for the variance of a biased, generally speaking, estimator of a parameter θ .

THEOREM 3.4.3. *Let regularity conditions (R) hold. Let $\hat{\theta}$ be an estimator of a parameter θ such that $E_{\theta}\hat{\theta}^2 \leq c < \infty$ for all $\theta \in \Theta$. Then*

$$(3.4.29) \quad D_{\theta}\hat{\theta} \geq \frac{(1 + b'(\theta))^2}{I(\theta)}, \quad \theta \in \Theta,$$

where $b(\theta) = E_{\theta}\hat{\theta} - \theta$ is the bias of the estimator $\hat{\theta}$.

If (3.4.9) becomes an equality on some interval $[\theta_1, \theta_2] \subset \Theta$ and $D_\theta \hat{\theta} > 0$ on this interval, then

$$(3.4.30) \quad f(x; \theta) = \exp\{A(\theta)\hat{\theta}(x) + B(\theta)\}h(x), \quad x \in X,$$

for $\theta \in [\theta_1, \theta_2]$ where $A(\theta)$ and $B(\theta)$ do not depend on x .

Conversely if either $\hat{\theta}(x) \equiv \text{const}$ or representation (3.4.30) holds, then inequality (3.4.29) becomes an equality.

PROOF. Let $a(\theta) = E_\theta \hat{\theta}$. Putting $T(x) \equiv 1$ in Lemma 3.4.4 we get from (3.4.9) that

$$(3.4.31) \quad E_\theta S(\xi; \theta) = 0, \quad E_\theta a(\theta)S(\xi; \theta) = 0.$$

Again using Lemma 3.4.4 for $T(x) = \hat{\theta}(x)$ we obtain from (3.4.9) and (3.4.31) that

$$(3.4.32) \quad E_\theta \hat{\theta}(\xi)S(\xi; \theta) = a'(\theta), \quad E_\theta (\hat{\theta}(\xi) - a(\theta))S(\xi; \theta) = a'(\theta).$$

Using the Cauchy-Bunyakovskiï inequality and the second equality in (3.4.32) we get

$$(3.4.33) \quad (a'(\theta))^2 \leq E_\theta (\hat{\theta}(\xi) - a(\theta))^2 E_\theta S^2(\xi; \theta)$$

or, equivalently,

$$(3.4.34) \quad D_\theta \hat{\theta} \geq \frac{(a'(\theta))^2}{E_\theta S^2(\xi; \theta)}, \quad \theta \in \Theta.$$

Since $a(\theta) = \theta + b(\theta)$ and $E_\theta S^2(\xi; \theta) = I(\theta)$, we obtain inequality (3.4.29) from (3.4.34).

Now we prove the second statement of Theorem 3.4.3. For the sake of simplicity we assume that Θ coincides with an interval $[\theta_1, \theta_2]$ and that the measure μ is concentrated on a union of supports of measures P_θ , $\theta \in \Theta$. The equality in (3.4.29) (or, equivalently, in (3.4.33)) means that

$$(3.4.35) \quad \int (\hat{\theta}(x) - a(\theta)) \frac{\partial f(x; \theta)}{\partial \theta} \mu(dx) \\ = \left(\int (\hat{\theta}(x) - a(\theta))^2 f(x; \theta) \mu(dx) \int \frac{(\partial f(x; \theta)/\partial \theta)^2}{f(x; \theta)} \mu(dx) \right)^{1/2}$$

for all $\theta \in \Theta$. Since the first integral on the right-hand side of (3.4.35) is positive by condition, an equality in (3.4.35) is only possible if

$$(3.4.36) \quad \frac{\partial f(x; \theta)/\partial \theta}{\sqrt{f(x; \theta)}} = c(\theta)(\hat{\theta}(x) - a(\theta))\sqrt{f(x; \theta)} \quad (\mu\text{-a.s.}).$$

Let A be the set of $x \in X$ for which (3.4.36) holds and $|\hat{\theta}(x)| < \infty$. Then $\mu(\bar{A}) = 0$ (here $\bar{A} = X \setminus A$ is the complement of the set A). Fix $x \in A$. Since $f(x; \theta)$ is continuous with respect to θ , we have $f(x; t) > 0$ on some interval $(t_1, t_2) \subset \Theta$, and moreover

$$(3.4.37) \quad S(x; t) = c(t)(\hat{\theta}(x) - a(t)), \quad t \in (t_1, t_2),$$

on this interval by (3.4.36). If (3.4.29) becomes an equality, then (3.4.32) and (3.4.37) imply that

$$a'(\theta) = E_{\theta}(\widehat{\theta}(\xi) - a(\theta))S(\xi; \theta) = c(\theta)D_{\theta}\widehat{\theta},$$

$$D_{\theta}\widehat{\theta} = \frac{(a'(\theta))^2}{I(\theta)}, \quad |c(\theta)| = \left(\frac{I(\theta)}{D_{\theta}\widehat{\theta}}\right)^{1/2}$$

This implies that the variance $D_{\theta}\widehat{\theta}$ is continuous with respect to θ as well as $a'(\theta)$, and $I(\theta)$ is continuous with respect to θ while the functions $|c(\theta)|$ and $a(\theta)$ are uniformly bounded on $[\theta_1, \theta_2]$. The derivative $S(x; \theta) = \partial \ln f(x; \theta) / \partial \theta$ in equality (3.4.37) possesses the same property. This means that the function $\ln f(x; \theta)$ is finite and $f(x; \theta) > 0$ for all $\theta \in \Theta = [\theta_1, \theta_2]$, whence (3.4.37) follows for all θ . Integrating equality (3.4.37) with respect to t from θ_1 to θ , we obtain

$$\ln f(x; \theta) = \widehat{\theta}(x) \int_{\theta_1}^{\theta} c(t) dt - \int_{\theta_1}^{\theta} c(t)a(t) dt + \ln f(x; \theta_1)$$

and this is equivalent to (3.4.30) for μ -almost all x . Since the values of $f(x; \theta)$ on a set whose measure μ is zero do not matter, representation (3.4.30) is proved.

Finally, we prove the latter statement of Theorem 3.4.3. If $\widehat{\theta}(x) \equiv \text{const}$, then $b'(\theta) = -1$ and both sides of equality (3.4.29) vanish. Now let representation (3.4.30) hold. Differentiating the function $\ln f(x; \theta)$ with respect to θ we get

$$S(x; \theta) = \widehat{\theta}(x)A'(\theta) + B'(\theta).$$

The first equality in (3.4.31) implies that

$$a(\theta)A'(\theta) + B'(\theta) = 0.$$

Thus

$$S(x; \theta) = A'(\theta)(\widehat{\theta}(x) - a(\theta))$$

and inequality (3.4.29) becomes an equality in view of (3.4.36). \square

REMARK 3.4.3. If $E_{\theta}\widehat{\theta}^2 = \infty$, then $D_{\theta}\widehat{\theta} = \infty$ and inequality (3.4.29) is trivial. In view of (3.4.29), the condition $D_{\theta}\widehat{\theta} > 0$ can be substituted by $1 + b'(\theta) \neq 0$.

COROLLARY 3.4.3. *If all the assumptions of Theorem 3.4.3 are satisfied, then*

$$E_{\theta}(\widehat{\theta} - \theta)^2 \geq \frac{(1 + b'(\theta))^2}{I(\theta)} + b^2(\theta), \quad \theta \in \Theta.$$

For every unbiased estimator $\widehat{\theta}$ of the parameter θ

$$E_{\theta}(\widehat{\theta} - \theta)^2 \geq \frac{1}{I(\theta)}, \quad \theta \in \Theta.$$

Analogs of Theorem 3.4.3 can be proved under other sets of conditions. Below is a set of conditions, called $(R)^*$, which also is sufficient for the the Cramér-Rao inequality:

- (i) Θ is a finite or infinite interval in \mathbf{R}^1 ;
- (ii)'' the function $\sqrt{f(x; \theta)}$ is absolutely continuous with respect to θ for μ -almost all $x \in X$;

(iii)' $0 < E_{\theta} S^2(\xi; \theta) = I(\theta) < \infty$ for all $\theta \in \Theta$ and the function $I(\theta)$ is continuous with respect to θ .

It is obvious that conditions $(R)^*$ are weaker than conditions (R) . Nevertheless the Cramér-Rao inequality holds under conditions $(R)^*$, too.

THEOREM 3.4.4. *Let the regularity conditions $(R)^*$ hold. Let $\hat{\theta}$ be an unbiased estimator of a parameter θ . Then*

$$(3.4.38) \quad E_{\theta}(\hat{\theta} - \theta)^2 \geq \frac{1}{I(\theta)}$$

for all points $\theta \in \Theta$ of continuity of $E_{\theta}(\hat{\theta} - \theta)^2$.

PROOF. Since $\hat{\theta}$ is unbiased,

$$\Delta = \int (\hat{\theta}(x) - \theta) \left(\sqrt{f(x; \theta + \Delta)} + \sqrt{f(x; \theta)} \right) \left(\sqrt{f(x; \theta + \Delta)} - \sqrt{f(x; \theta)} \right) \mu(dx).$$

Applying the Cauchy-Bunyakovskiĭ inequality and then the elementary inequality $(\sqrt{a} + \sqrt{b})^2 \leq 2(a + b)$ for $a > 0$ and $b > 0$, we get

$$(3.4.39) \quad \begin{aligned} \Delta^2 &\leq \int (\hat{\theta}(x) - \theta)^2 \left(\sqrt{f(x; \theta + \Delta)} + \sqrt{f(x; \theta)} \right)^2 \mu(dx) \\ &\quad \times \int \left(\sqrt{f(x; \theta + \Delta)} - \sqrt{f(x; \theta)} \right)^2 \mu(dx) \\ &\leq 2 \int (\hat{\theta}(x) - \theta)^2 (f(x; \theta + \Delta) + f(x; \theta)) \mu(dx) \\ &\quad \times \int \left(\sqrt{f(x; \theta + \Delta)} - \sqrt{f(x; \theta)} \right)^2 \mu(dx). \end{aligned}$$

Then

$$(3.4.40) \quad \int (\hat{\theta}(x) - \theta)^2 f(x; \theta) \mu(dx) = E_{\theta}(\hat{\theta} - \theta)^2,$$

$$(3.4.41) \quad \int (\hat{\theta}(x) - \theta)^2 f(x; \theta + \Delta) \mu(dx) = E_{\theta + \Delta}(\hat{\theta} - \theta - \Delta)^2 + \Delta^2.$$

It follows from (3.4.39)–(3.4.41) that

$$(3.4.42) \quad \begin{aligned} &\frac{E_{\theta}(\hat{\theta} - \theta)^2 + E_{\theta + \Delta}(\hat{\theta} - \theta - \Delta)^2 + \Delta^2}{2} \\ &\geq \left(\frac{4}{\Delta^2} \int \left(\sqrt{f(x; \theta + \Delta)} - \sqrt{f(x; \theta)} \right)^2 \mu(dx) \right)^{-1} \end{aligned}$$

Condition (ii)'' implies that

$$\sqrt{f(x; \theta + \Delta)} - \sqrt{f(x; \theta)} = \int_{\theta}^{\theta + \Delta} \frac{\partial}{\partial u} \sqrt{f(x; u)} du.$$

Applying the Cauchy–Bunyakovskii inequality and Fubini theorem we obtain

$$\begin{aligned}
 & \frac{4}{\Delta^2} \int \left(\sqrt{f(x; \theta + \Delta)} - \sqrt{f(x; \theta)} \right)^2 \mu(dx) \\
 &= \frac{4}{\Delta^2} \int \left(\int_{\theta}^{\theta + \Delta} \frac{\partial \sqrt{f(x; u)}}{\partial u} du \right)^2 \mu(dx) \\
 (3.4.43) \quad & \leq \frac{4}{\Delta} \int \int_{\theta}^{\theta + \Delta} \left(\frac{\partial \sqrt{f(x; u)}}{\partial u} \right)^2 du \mu(dx) \\
 &= \frac{1}{\Delta} \int_{\theta}^{\theta + \Delta} \left(4 \int \left(\frac{\partial \sqrt{f(x; u)}}{\partial u} \right)^2 \mu(dx) \right) du \\
 &= \frac{1}{\Delta} \int_{\theta}^{\theta + \Delta} I(u) du.
 \end{aligned}$$

Substituting (3.4.43) into (3.4.42) and passing to the limit as $\Delta \rightarrow 0$, we complete the proof of (3.4.38) in view of the continuity of the function $I(u)$. \square

The above proof of inequality (3.4.42) does not require any condition posed on the density $f(x; \theta)$ or on an unbiased estimator $\hat{\theta}$. Thus this proof can be used to obtain lower bounds of the variance of an estimator $\hat{\theta}$ even in the case where regularity conditions are not satisfied for $f(x; \theta)$. Below we provide a result of this kind for biased estimators $\hat{\theta}$.

THEOREM 3.4.5. *Let $\theta \in \Theta$ and $\theta + \Delta \in \Theta$ for some $\Delta \neq 0$. Then for all estimators $\hat{\theta}$ of a parameter θ one has*

$$(3.4.44) \quad D_{\theta} \hat{\theta} \geq (\Delta a(\theta))^2 \left(\int \frac{(\Delta f(x; \theta))^2}{f(x; \theta)} \mu(dx) \right)^{-1}$$

where

$$\begin{aligned}
 a(\theta) &= E_{\theta} \hat{\theta}, & \Delta a(\theta) &= a(\theta + \Delta) - a(\theta), \\
 \Delta f(x; \theta) &= f(x; \theta + \Delta) - f(x; \theta).
 \end{aligned}$$

In particular,

$$(3.4.45) \quad D_{\theta} \hat{\theta} \geq \left(\frac{1}{\Delta^2} \int \frac{(\Delta f(x; \theta))^2}{f(x; \theta)} \mu(dx) \right)^{-1}$$

if the estimator $\hat{\theta}$ is unbiased.

PROOF. First let the measure $P_{\theta + \Delta}$ be not absolutely continuous with respect to the measure P_{θ} . Denote by N_{θ} the support of the measure P_{θ} in X and let $N_{\theta} = \{x: f(x; \theta) \neq 0\}$. Then there is a set $A \subset N_{\theta + \Delta}$ such that

$$P_{\theta + \Delta}(A) > 0$$

and $f(x; \theta) = 0$ for all $x \in A$. Thus the integral in (3.4.44) is infinite and inequality (3.4.44) is trivial.

Now let the measure $P_{\theta+\Delta}$ be absolutely continuous with respect to the measure P_θ . Then $N_{\theta+\Delta} \subset N_\theta$. Since $f(x; \theta)$ and $f(x; \theta + \Delta)$ are the densities of measures P_θ and $P_{\theta+\Delta}$, respectively, with respect to μ , we have

$$\int \Delta f(x; \theta) \mu(dx) = 0.$$

Moreover

$$\int \widehat{\theta}(x) \Delta f(x; \theta) \mu(dx) = \Delta a(\theta).$$

This implies that

$$(3.4.46) \quad \int_{N_\theta} (\widehat{\theta}(x) - a(\theta)) \Delta f(x; \theta) \mu(dx) = \Delta a(\theta).$$

The integrand in (3.4.46) can be represented on the set N_θ as

$$(\widehat{\theta}(x) - a(\theta)) \sqrt{f(x; \theta)} \cdot \frac{\Delta f(x; \theta)}{\sqrt{f(x; \theta)}}.$$

Applying the Cauchy-Bunyakovskiĭ inequality we obtain

$$(\Delta a(\theta))^2 \leq \int_{N_\theta} (\widehat{\theta}(x) - a(\theta))^2 f(x; \theta) \mu(dx) \cdot \int_{N_\theta} \frac{(\Delta f(x; \theta))^2}{f(x; \theta)} \mu(dx),$$

whence inequality (3.4.44) follows. Inequality (3.4.45) follows from (3.4.44). \square

COROLLARY 3.4.4. *Assume that an arbitrary set of regularity conditions holds and*

$$\lim_{\Delta \rightarrow 0} \frac{1}{\Delta^2} \int \frac{(\Delta f(x; \theta))^2}{f(x; \theta)} \mu(dx) = I(\theta).$$

Then

$$(3.4.47) \quad D_\theta \widehat{\theta} \geq (a'_+(\theta))^2 / I(\theta)$$

for all estimators $\widehat{\theta}$ of the parameter θ where

$$a'_+(\theta) = \limsup_{\Delta \rightarrow 0} \frac{\Delta a(\theta)}{\Delta}.$$

PROOF. To prove (3.4.47) we pass to the limit in (3.4.44) along a subsequence $\Delta \rightarrow 0$ such that $\Delta a(\theta) / \Delta \rightarrow a'_+(\theta)$. \square

Inequality (3.4.44) is called the *Chapman-Robbins inequality*. Another name for it is the *difference Cramér-Rao type inequality*.

REMARK 3.4.4. If an observation ξ is a sample $\xi^{(n)} = (\xi_1, \dots, \xi_n)$, $f(x; \theta)$ is the density of ξ_1 , and $f_n(x; \theta)$ is the density of the vector $\xi^{(n)}$, then one can obtain analogs of all the above results. If conditions (R) or (R)* hold for the density $f(x; \theta)$, then all the above inequalities hold for estimators $\widehat{\theta}_n$ of the parameter θ . The only exception is that the Fisher information $nI(\theta)$ substitutes the Fisher information $I(\theta)$ where $I(\theta)$ is evaluated with respect to the density $f(x; \theta)$. In particular, inequality (3.4.29) becomes in this case of the form

$$(3.4.48) \quad D_\theta \widehat{\theta}_n \geq \frac{(1 + b'_n(\theta))^2}{nI(\theta)}$$

where $b_n(\theta) = E_\theta \hat{\theta}_n - \theta$ is the bias of the estimator $\hat{\theta}_n$. Further, inequality (3.4.44) can be rewritten in this case as

$$(3.4.49) \quad D_\theta \hat{\theta}_n \geq (\Delta a_n(\theta))^2 \left(\int \frac{(\Delta f_n(x; \theta))^2}{f_n(x; \theta)} \mu(dx) \right)^{-1}$$

where $a_n(\theta) = E_\theta \hat{\theta}_n$.

Efficient and asymptotically efficient estimators of parameters. Let one of the sets of regularity conditions (CR) , $(CR)^*$, (R) , or $(R)^*$ hold. Let ξ be an observation that is a random element assuming values in a measurable space (X, \mathcal{B}) . Let the distribution of ξ belong to a family of probability measures

$$\{P_\theta, \theta \in \Theta\}, \quad \Theta \subset \mathbf{R}^1.$$

Let K_b^g be the class of estimators $\hat{g} = \hat{g}(\xi)$ of a function $g(\theta)$ with a bias $b(\theta)$. Let K_b be the class of estimators $\hat{\theta} = \hat{\theta}(\xi)$ of a parameter θ with a bias $b(\theta)$, that is,

$$K_b^g = \{\hat{g}: E_\theta \hat{g} = g(\theta) + b(\theta)\}, \quad K_b = \{\hat{\theta}: E_\theta \hat{\theta} = \theta + b(\theta)\}.$$

We also consider the class K^g of unbiased estimators \hat{g} of a function $g(\theta)$ and the class K of unbiased estimators $\hat{\theta}$ of a parameter θ , that is,

$$K^g = K_0^g = \{\hat{g}: E_\theta \hat{g} = g(\theta)\}, \quad K = K_0 = \{\hat{\theta}: E_\theta \hat{\theta} = \theta\}.$$

Note that $K^g = U_g$ and $K = U_0$ where U_g and U_0 are the classes of estimators introduced in Section 3.1.

We say that $g^* \in K_b^g$ is an *efficient estimator of a function $g(\theta)$ in the class K_b^g* if the Cramér–Rao inequality for it becomes an equality, that is,

$$(3.4.50) \quad D_\theta g^* = \frac{(g'(\theta) + b'(\theta))^2}{I(\theta)}, \quad \theta \in \Theta.$$

Similarly, $g^* \in K^g$ is called an *efficient estimator of a function $g(\theta)$ in the class K^g* (or, an *efficient estimator of a function $g(\theta)$*) if

$$(3.4.51) \quad D_\theta g^* = \frac{(g'(\theta))^2}{I(\theta)}, \quad \theta \in \Theta.$$

Conditions (3.4.50) and (3.4.51) can be rewritten in the following equivalent form:

$$(3.4.52) \quad E_\theta (g^* - g(\theta))^2 = \frac{(g'(\theta) + b'(\theta))^2}{I(\theta)} + b^2(\theta), \quad \theta \in \Theta,$$

$$(3.4.53) \quad E_\theta (g^* - g(\theta))^2 = \frac{(g'(\theta))^2}{I(\theta)}, \quad \theta \in \Theta.$$

Efficient estimators θ^* of a parameter θ in the classes K_b and K can be introduced similarly to (3.4.50)–(3.4.53) if $g(\theta) = \theta$. In particular, $\theta^* \in K_b$ is called an *efficient estimator of a parameter θ in the class K_b* if the corresponding Cramér–Rao inequality becomes an equality, that is,

$$(3.4.54) \quad E_\theta (\theta^* - \theta)^2 = \frac{(1 + b'(\theta))^2}{I(\theta)} + b^2(\theta), \quad \theta \in \Theta.$$

We say that $\theta^* \in K$ is an *efficient estimator of a parameter θ in the class K* if

$$(3.4.55) \quad E_{\theta}(\theta^* - \theta)^2 = \frac{1}{I(\theta)}, \quad \theta \in \Theta.$$

Conditions (3.4.50)–(3.4.55) suggest a general definition: an *estimator is called efficient in the corresponding class if the Cramér–Rao inequality becomes an equality*.

Efficient estimators exist only in exceptional cases. In other cases one can construct the so-called asymptotically efficient estimators if the size of a sample increases.

Let an observation be a sample $\xi^{(n)} = (\xi_1, \dots, \xi_n)$ with the density

$$f_n(x; \theta) = \prod_{i=1}^n f(x_i; \theta), \quad x = (x_1, \dots, x_n), \quad \theta \in \Theta.$$

A sequence of estimators g_n^* , $n = 1, 2, \dots$, is called an *asymptotically efficient estimator of a function $g(\theta)$* if

$$E_{\theta}(g_n^* - g(\theta))^2 = \frac{(g'(\theta))^2}{nI(\theta)} + o\left(\frac{1}{n}\right), \quad \theta \in \Theta,$$

as $n \rightarrow \infty$. For the sake of brevity we say that g_n^* is an asymptotically efficient estimator of a function $g(\theta)$. Similarly, θ_n^* is called an *asymptotically efficient estimator of a parameter θ* if

$$E_{\theta}(\theta_n^* - \theta)^2 = \frac{1}{nI(\theta)} + o\left(\frac{1}{n}\right), \quad \theta \in \Theta,$$

as $n \rightarrow \infty$.

Another name for asymptotically efficient estimators is *asymptotically efficient estimators in the strong sense*. If a sequence of estimators g_n^* , $n = 1, 2, \dots$, is asymptotically $\mathcal{N}(g(\theta), (g'(\theta))^2/(nI(\theta)))$ normal, then we say that g_n^* is an *asymptotically efficient estimator of a function $g(\theta)$ in the weak sense*. If a sequence of estimators θ_n^* , $n = 1, 2, \dots$, is asymptotically $\mathcal{N}(\theta, 1/(nI(\theta)))$ normal, then the estimator θ_n^* is called an *asymptotically efficient estimator of a parameter θ in the weak sense*.

EXAMPLE 3.4.1. Let $\xi^{(n)} = (\xi_1, \xi_2, \dots, \xi_n)$ be a sample from the normal distribution $\mathcal{N}(\theta, \sigma^2)$. In this case

$$I(\theta) = E_{\theta} S^2(\xi_1; \theta) = E_{\theta} \left(\frac{\xi_1 - \theta}{\sigma^2} \right)^2 = \frac{1}{\sigma^2}.$$

Then $\hat{\theta}_n = n^{-1} \sum_{i=1}^n \xi_i$ is an unbiased estimator of the parameter θ and

$$E_{\theta}(\hat{\theta}_n - \theta)^2 = \frac{\sigma^2}{n} = \frac{1}{nI(\theta)}.$$

Thus $\hat{\theta}_n$ is an efficient estimator of the parameter θ .

EXAMPLE 3.4.2. Let $\xi^{(n)} = (\xi_1, \xi_2, \dots, \xi_n)$ be a sample from the $\mathcal{N}(\alpha, \theta)$ distribution. It is clear that

$$\begin{aligned} I(\theta) &= E_\theta S^2(\xi_1; \theta) = E_\theta \left(\frac{(\xi_1 - \alpha)^2}{2\theta^2} - \frac{1}{2\theta} \right)^2 \\ &= \frac{1}{4\theta^4} E_\theta (\xi_1 - \alpha)^4 - \frac{1}{2\theta^3} E_\theta (\xi_1 - \alpha)^2 + \frac{1}{4\theta^2} = \frac{1}{2\theta^2}, \end{aligned}$$

since $E_\theta (\xi_1 - \alpha)^4 = 3\theta^2$. It is clear that $\widehat{\theta}_n = n^{-1} \sum_{i=1}^n (\xi_i - \alpha)^2$ is an unbiased estimator of the parameter θ and

$$\begin{aligned} E_\theta (\widehat{\theta}_n - \theta)^2 &= E_\theta \left(\frac{1}{n} \sum_{i=1}^n [(\xi_i - \alpha)^2 - \theta] \right)^2 \\ &= \frac{1}{n^2} \sum_{i=1}^n E_\theta [(\xi_i - \alpha)^2 - \theta]^2 \\ &\quad + \frac{1}{n^2} \sum_{j \neq i} E_\theta [(\xi_i - \alpha)^2 - \theta] E_\theta [(\xi_j - \alpha)^2 - \theta] \\ &= \frac{1}{n} E_\theta [(\xi_1 - \alpha)^2 - \theta]^2 = \frac{1}{n} [E_\theta (\xi_1 - \alpha)^4 - 2\theta E_\theta (\xi_1 - \alpha)^2 + \theta^2] \\ &= \frac{2\theta^2}{n} = \frac{1}{nI(\theta)}, \end{aligned}$$

whence it follows that $\widehat{\theta}_n$ is an efficient estimator of the parameter θ .

Consider the estimator $\theta_n = (n-1)^{-1} \sum_{i=1}^n (\xi_i - \bar{\xi})^2$, where $\bar{\xi} = n^{-1} \sum_{i=1}^n \xi_i$. It is clear that $E_\theta \theta_n = \theta$, that is, θ_n is an unbiased estimator of the parameter θ (see Example 3.1.3). According to (3.2.9) we have

$$E_\theta (\theta_n - \theta)^2 = \frac{2\theta^2}{n-1} = \frac{1}{(n-1)I(\theta)} > \frac{1}{nI(\theta)},$$

that is, θ_n is not an efficient estimator of the parameter θ . Nevertheless θ_n is an asymptotically efficient estimator of the parameter in the strong sense.

EXAMPLE 3.4.3. Let $\xi^{(n)} = (\xi_1, \xi_2, \dots, \xi_n)$ be a sample from the Gamma distribution with the density

$$f(x; \theta) = \frac{1}{\Gamma(\theta)} x^{\theta-1} e^{-x} I_{(0, \infty)}(x)$$

where $\theta \in \Theta = (0, \infty)$ and $\Gamma(\theta)$ is the Gamma function. It is obvious that regularity conditions hold in this case. By Lemma 3.4.2 we have

$$(3.4.56) \quad I(\theta) = -E_\theta \left(\frac{\partial^2 \ln f(\xi_1; \theta)}{\partial \theta^2} \right) = \frac{d^2 \ln \Gamma(\theta)}{d\theta^2}.$$

Consider the estimator $\widehat{\theta}_n = n^{-1} \sum_{i=1}^n \xi_i$. It is clear that $E_\theta \widehat{\theta}_n = \theta$, that is, $\widehat{\theta}_n$ is an unbiased estimator of the parameter θ . Moreover,

$$(3.4.57) \quad E_\theta (\widehat{\theta}_n - \theta)^2 = \frac{\theta}{n} = \frac{1}{nI(\theta)} \cdot \kappa(\widehat{\theta}_n; \theta)$$

by (3.4.56) where

$$(3.4.58) \quad \kappa(\hat{\theta}_n; \theta) = \theta \frac{d^2 \ln \Gamma(\theta)}{d\theta^2}.$$

By the Stirling formula (see [9], relation (12.5.3)) we have

$$\ln \Gamma(\theta) = \left(\theta - \frac{1}{2}\right) \ln \theta - \theta + \frac{1}{2} \ln 2\pi + \int_0^\infty \frac{P_1(x)}{\theta + x} dx$$

where $P_1(x)$ is the periodic function with period 1 such that $P_1(x) = -x + \frac{1}{2}$ for $x \in (0, 1)$. This implies

$$\frac{d^2 \ln \Gamma(\theta)}{d\theta^2} = \frac{1}{\theta} + \frac{1}{2\theta^2} + 2 \int_0^\infty \frac{P_1(x)}{(\theta + x)^3} dx.$$

Thus for all $\theta \in (0, \infty)$

$$(3.4.59) \quad \kappa(\hat{\theta}_n; \theta) = 1 + \frac{1}{2\theta} + 2\theta \int_0^\infty \frac{P_1(x)}{(\theta + x)^3} dx > 1.$$

The coefficient $\kappa(\hat{\theta}_n; \theta)$ can be made as large as we want by choosing a sufficiently large θ . Thus (3.4.57)–(3.4.59) imply that the estimator $\hat{\theta}_n$ is not asymptotically efficient whatever the parameter θ is.

EXAMPLE 3.4.4. Let $\xi^{(n)} = (\xi_1, \dots, \xi_n)$ be a sample from the exponential distribution with the density

$$f(x; \theta) = e^{-x+\theta} I_{[\theta, \infty)}(x).$$

Regularity conditions do not hold in this case, since the function $f(x; \theta)$ is discontinuous with respect to θ . Consider the estimator

$$\hat{\theta}_n = \min_{1 \leq i \leq n} \xi_i - \frac{1}{n}.$$

We learned from Example 3.2.2 that $\hat{\theta}_n$ is the Pitman estimator of the parameter θ and moreover

$$E_\theta \hat{\theta}_n = \theta, \quad E_\theta (\hat{\theta}_n - \theta)^2 = \frac{2}{n^2}.$$

This implies that the mean square error $E_\theta (\hat{\theta}_n - \theta)^2$ is of order n^{-2} for large n . This is a higher rate of decay as compared to the one given by the Cramér–Rao lower bound. This phenomenon occurs, since the regularity conditions fail in this case. Other examples of higher rates of decay of $E_\theta (\hat{\theta}_n - \theta)^2$ can be obtained by using the lower bound in the Chapman–Robbins inequality in some other cases where the regularity conditions fail.

REMARK 3.4.5. Further information about the regularity conditions and Cramér–Rao inequalities can be found in [36] and [13].

REMARK 3.4.6. The Cramér–Rao inequalities belong to the family of results, called *information inequalities*, which provide lower bounds for the risk functions or risks of estimators of parameters. See [22], Chapter 5, about the relation between the Cramér–Rao inequalities and for other information about inequalities.

3.5. The Cramér–Rao inequality for a multidimensional parameter

In this section we consider analogs of the Cramér–Rao inequalities for the case of a multidimensional parameter θ .

The Fisher information matrix. Let ξ be an observation that is a random element assuming values in a measurable space (X, \mathcal{B}) . Assume that its distribution belongs to a family of probability measures $\{P_\theta, \theta \in \Theta\}$ where Θ is some subset of \mathbf{R}^k , $k \geq 1$. As in the case of a one-dimensional parameter we assume that for all $\theta \in \Theta$ the measure P_θ is absolutely continuous with respect to some σ -finite measure μ on (X, \mathcal{B}) and that there exists the density $f(x; \theta)$ of the measure P_θ with respect to the measure μ .

Let the derivatives $S_i(x; \theta) = \partial \ln f(x; \theta) / \partial \theta_i$, $i = 1, \dots, k$, exist for μ -almost all $x \in X$. The matrix $I(\theta)$ with the entries $I_{ij}(\theta) = E_\theta S_i(\xi; \theta) S_j(\xi; \theta)$, $i, j = 1, \dots, k$, is called the *Fisher information matrix*. In the case $k = 1$, $I(\theta)$ is the Fisher information.

It is easy to see that the matrix $I(\theta)$ is nonnegative definite. Indeed, for all $\lambda = (\lambda_1, \dots, \lambda_k)' \in \mathbf{R}^k$,

$$\begin{aligned} \lambda' I(\theta) \lambda &= \sum_{i,j=1}^k I_{ij}(\theta) \lambda_i \lambda_j = E_\theta \sum_{i,j=1}^k S_i(\xi; \theta) S_j(\xi; \theta) \lambda_i \lambda_j \\ (3.5.1) \quad &= E_\theta \left(\sum_{i=1}^k S_i(\xi; \theta) \lambda_i \right)^2 = E_\theta (\lambda' S(\xi; \theta))^2 \geq 0 \end{aligned}$$

where $S(\xi; \theta)$ is the vector defined by

$$(3.5.2) \quad S(\xi; \theta) = (S_1(\xi; \theta), S_2(\xi; \theta), \dots, S_k(\xi; \theta))'.$$

Inequality (3.5.1) implies that the Fisher information matrix $I(\theta)$ is nonnegative definite and this explains why we write $I(\theta) \geq 0$ in this case. If the matrix $I(\theta)$ is positive definite, then $\lambda' I(\theta) \lambda > 0$ for all vectors $\lambda \neq 0$. We write $I(\theta) > 0$ in the latter case. We write $A \geq B$ for matrices A and B if $A - B \geq 0$, that is, if the matrix $A - B$ is nonnegative definite.

The Cramér–Rao inequality for unbiased estimators. Let $\hat{\theta}$ be an estimator of a parameter θ constructed from an observation ξ where $\hat{\theta} = (\hat{\theta}_1, \hat{\theta}_2, \dots, \hat{\theta}_k)'$ and $\theta = (\theta_1, \theta_2, \dots, \theta_k)'$. Denote by $R(\hat{\theta}; \theta)$ the matrix with entries

$$E_\theta(\hat{\theta}_i - \theta_i)(\hat{\theta}_j - \theta_j), \quad i, j = 1, 2, \dots, k.$$

In other words, $R(\hat{\theta}; \theta) = E_\theta(\hat{\theta} - \theta)(\hat{\theta} - \theta)'$ is the matrix of mixed moments of the vector $\hat{\theta} - \theta$. It is easy to show that for all vectors $\lambda \in \mathbf{R}^k$

$$(3.5.3) \quad E_\theta((\hat{\theta} - \theta)' \lambda)^2 = \lambda' R(\hat{\theta}; \theta) \lambda,$$

that is, the matrix $R(\hat{\theta}; \theta)$ is nonnegative definite.

If the matrix $I(\theta)$ is nondegenerate, then we will show that $R(\hat{\theta}; \theta) \geq I^{-1}(\theta)$ under some conditions on $f(x; \theta)$, that is, we will show for all vectors $\lambda \in \mathbf{R}^k$ that

$$(3.5.4) \quad \lambda' R(\hat{\theta}; \theta) \lambda \geq \lambda' I^{-1}(\theta) \lambda.$$

The latter is a *matrix analog of the Cramér-Rao inequality*. Inequality (3.5.4) for an unbiased estimator $\widehat{\theta}$ means in view of relations (3.5.1) and (3.5.3) that the variance of the projection of the vector $\widehat{\theta}$ on an arbitrary direction λ in \mathbf{R}^k is greater than or equal to the variance of the projection of the vector $S(\xi; \theta)$ on the same direction.

Let $\mathbf{C}_1^0(\Theta)$ be the class of real functions $\phi(\theta)$ defined on $\Theta \subset \mathbf{R}^k$ that are differentiable in θ almost everywhere with respect to the Lebesgue measure, and such that the function $\phi(\theta + t\lambda)$ is absolutely continuous in t for all θ and λ for which $\theta + t\lambda \in \Theta$ and $0 \leq t \leq 1$. If $\phi \in \mathbf{C}_1^0(\Theta)$ and $\theta + t\lambda \in \Theta$ for $0 \leq t \leq \Delta$, then

$$(3.5.5) \quad \phi(\theta + \Delta\lambda) - \phi(\theta) = \int_0^\Delta \frac{d\phi(\theta + u\lambda)}{du} du = \int_0^\Delta \lambda' \frac{\partial \phi(\theta + u\lambda)}{\partial \theta} du$$

where

$$\frac{\partial \phi(\theta)}{\partial \theta} = \left(\frac{\partial \phi(\theta)}{\partial \theta_1}, \frac{\partial \phi(\theta)}{\partial \theta_2}, \dots, \frac{\partial \phi(\theta)}{\partial \theta_k} \right)'$$

The following result contains sufficient regularity conditions posed on the density $f(x; \theta)$ under which the Cramér-Rao inequality (3.5.4) holds.

THEOREM 3.5.1. *Let $\sqrt{f(x; \theta)} \in \mathbf{C}_1^0(\Theta)$ for μ -almost all $x \in X$. Assume that the matrix $I(\theta)$ is continuous in θ and nondegenerate. If $\widehat{\theta}$ is an unbiased estimator of the parameter θ , then*

$$(3.5.6) \quad R(\widehat{\theta}; \theta) \geq I^{-1}(\theta)$$

for all points θ of continuity of the matrix $R(\widehat{\theta}; \theta)$.

PROOF. Let $\theta \in \Theta$, $\lambda \in \mathbf{R}^k$, and $|\lambda| = 1$. Then $\theta + t\lambda \in \Theta$ for all sufficiently small $\Delta > 0$ and for all $t \in [0, \Delta]$. Since $\widehat{\theta}$ is an unbiased estimator,

$$E_\theta \widehat{\theta} = \theta, \quad E_{\theta + \Delta\lambda} \widehat{\theta} = \theta + \Delta\lambda,$$

whence we obtain

$$\int (\widehat{\theta}(x) - \theta)[f(x; \theta + \Delta\lambda) - f(x; \theta)] \mu(dx) = \Delta\lambda.$$

Multiplying this equality on the left by u' and applying the Cauchy-Bunyakovskiĭ inequality we get

$$(3.5.7) \quad \begin{aligned} \Delta^2 (u'\lambda)^2 &\leq \int \left(u'(\widehat{\theta}(x) - \theta) \right)^2 \left(\sqrt{f(x; \theta + \Delta\lambda)} + \sqrt{f(x; \theta)} \right)^2 \mu(dx) \\ &\quad \times \int \left(\sqrt{f(x; \theta + \Delta\lambda)} - \sqrt{f(x; \theta)} \right)^2 \mu(dx) \\ &\leq 2 \int \left(u'(\widehat{\theta}(x) - \theta) \right)^2 (f(x; \theta + \Delta\lambda) + f(x; \theta)) \mu(dx) \\ &\quad \times \int \left(\sqrt{f(x; \theta + \Delta\lambda)} - \sqrt{f(x; \theta)} \right)^2 \mu(dx). \end{aligned}$$

Relations (3.5.3) and (3.5.7) imply that

$$(3.5.8) \quad \begin{aligned} \Delta^2 (u'\lambda)^2 &\leq 2 \left[u'R(\widehat{\theta}; \theta)u + u'R(\widehat{\theta}; \theta + \Delta\lambda)u + \Delta^2 (u'\lambda)^2 \right] \\ &\quad \times \int \left(\sqrt{f(x; \theta + \Delta\lambda)} - \sqrt{f(x; \theta)} \right)^2 \mu(dx). \end{aligned}$$

Since $\sqrt{f(x; \theta)} \in C_1^0(\Theta)$, we obtain from (3.5.5) that

$$\begin{aligned}
 & \int \left(\sqrt{f(x; \theta + \Delta \lambda)} - \sqrt{f(x; \theta)} \right)^2 \mu(dx) \\
 (3.5.9) \quad & = \int \left(\int_0^\Delta \frac{\lambda' \partial f(x; \theta + t \lambda) / \partial \theta}{2 \sqrt{f(x; \theta + t \lambda)}} dt \right)^2 \mu(dx) \\
 & \leq \frac{\Delta}{4} \int_0^\Delta \lambda' I(\theta + t \lambda) \lambda dt.
 \end{aligned}$$

Let θ be a point of continuity of the matrix $R(\hat{\theta}; \theta)$. Substituting (3.5.9) into (3.5.8) and passing to the limit as $\Delta \rightarrow 0$ we get

$$(3.5.10) \quad (u' R(\hat{\theta}; \theta) u) (\lambda' I(\theta) \lambda) \geq (u' \lambda)^2.$$

Putting $\lambda = I^{-1}(\theta) u$ we derive from (3.5.10) that

$$u' R(\hat{\theta}; \theta) u \geq u' I^{-1}(\theta) u$$

for all vectors $u \in \mathbf{R}^k$. The latter inequality is equivalent to (3.5.6). \square

Let an observation be a sample $\xi^{(n)} = (\xi_1, \dots, \xi_n)$ and let the density of the sample be $f_n(x; \theta) = \prod_{i=1}^n f(x_i; \theta)$ where $x = (x_1, \dots, x_n)$. Let $I_n(\theta)$ and $I(\theta)$ be the information matrices constructed from the densities $f_n(x; \theta)$ and $f(x; \theta)$, respectively. If the regularity conditions of Theorem 3.5.1 hold, then

$$I_n(\theta) = nI(\theta).$$

If $\hat{\theta}_n$ is an unbiased estimator of the parameter θ , then under the conditions of Theorem 3.5.1 we get the following matrix analog of the Cramér–Rao inequality:

$$(3.5.11) \quad R(\hat{\theta}_n; \theta) \geq \frac{1}{n} I^{-1}(\theta).$$

The Cramér–Rao inequality for biased estimators. Let $\hat{\theta} = (\hat{\theta}_1, \dots, \hat{\theta}_k)'$ be an estimator of a parameter $\theta = (\theta_1, \dots, \theta_k)'$ constructed from an observation ξ . Put

$$a(\theta) = E_\theta \hat{\theta} = \theta + b(\theta), \quad b(\theta) = (b_1(\theta), \dots, b_k(\theta))'.$$

Here $b(\theta)$ is the *bias vector of the estimator* $\hat{\theta}$ of the parameter θ .

Consider the *multivariate analog of the regularity conditions (R)* introduced in Section 3.4:

(R) the function $\sqrt{f(x; \theta)}$ is continuously differentiable in θ for μ -almost all x ; the matrix $I(\theta)$ is nondegenerate and continuous in θ .

In what follows we need the following *Cauchy–Bunyakovskii inequality for matrices*.

LEMMA 3.5.1. *Let η and ζ be two random matrices of the same size (they are not necessarily square matrices). Assume that the inverse matrix of $E\eta\eta'$ exists. Then*

$$(3.5.12) \quad E\zeta\zeta' \geq E\zeta\eta'(E\eta\eta')^{-1}E\eta\zeta'.$$

This inequality becomes an inequality if and only if $\zeta = z\eta$ where $z = E\zeta\eta'(E\eta\eta')^{-1}$.

PROOF. Since $AA' \geq 0$ for an arbitrary matrix A (that is, the matrix AA' is nonnegative definite),

$$0 \leq E(\zeta - z\eta)(\zeta - z\eta)' = E\zeta\zeta' - zE\eta\zeta' - E\zeta\eta'z' + zE\eta\eta'z'$$

for a nonrandom square matrix z . Putting $z = E\zeta\eta'(E\eta\eta')^{-1}$ we obtain inequality (3.5.12). The statement concerning the case of an equality in (3.5.12) is obvious. \square

The following result contains the Cramér-Rao inequality for a biased estimator $\hat{\theta}$ under the regularity conditions (R).

THEOREM 3.5.2. *Let conditions (R) hold. Let $D(\hat{\theta}; \theta) = E_{\theta}(\hat{\theta} - a(\theta))(\hat{\theta} - a(\theta))'$ be the matrix of mixed central moments of second order of an estimator $\hat{\theta}$ of a parameter θ . Then*

$$(3.5.13) \quad D(\hat{\theta}; \theta) \geq (I_k + B(\theta))I^{-1}(\theta)(I_k + B(\theta))'$$

where I_k is the unit matrix, $B(\theta) = \|b_{ij}(\theta)\|$, and $b_{ij}(\theta) = \partial b_i(\theta)/\partial \theta_j$.

Let $\det(D(\hat{\theta}; \theta)) > 0$ (or $\det(I_k + B(\theta)) > 0$) for all θ . Then inequality (3.5.13) becomes an equality if and only if the density of the distribution is such that

$$(3.5.14) \quad f(x; \theta) = \exp\{A(\theta)'\hat{\theta}(x) + C(\theta)\}h(x), \quad x \in X,$$

for some scalar functions $C(\theta)$ and $h(x)$ where

$$(3.5.15) \quad \|A_{ij}(\theta)\| = \left\| \frac{\partial A_i(\theta)}{\partial \theta_j} \right\| = ((I_k + B(\theta))^{-1})'I(\theta)$$

is the matrix of derivatives of the vector $A(\theta) = (A_1(\theta), \dots, A_k(\theta))'$.

If $\hat{\theta}$ is an unbiased estimator, then

$$(3.5.16) \quad D(\hat{\theta}; \theta) \geq I^{-1}(\theta).$$

Inequality (3.5.16) becomes an equality if and only if relation (3.5.14) holds where $\|A_{ij}(\theta)\| = I(\theta)$.

PROOF. As in the proof for the one-dimensional case we use the regularity conditions (R) to prove that

$$E_{\theta}S_i(\xi; \theta) = 0, \quad E_{\theta}\hat{\theta}_i S_j(\xi; \theta) = \delta_{ij} + b_{ij}(\theta), \quad i, j = 1, 2, \dots, k$$

(see Lemma 3.4.4), where δ_{ij} is the Kronecker symbol and the functions $b_{ij}(\theta)$ are continuous. The latter condition can be written in matrix form as follows:

$$(3.5.17) \quad E_{\theta}S(\xi; \theta) = 0,$$

$$(3.5.18) \quad E_{\theta}\hat{\theta}S(\xi; \theta)' = I_k + B(\theta)$$

where the matrix $B(\theta)$ is continuous in θ and the vector $S(\xi; \theta)$ is of the form (3.5.2). Equalities (3.5.17) and (3.5.18) imply that

$$(3.5.19) \quad E_{\theta}(\hat{\theta} - a(\theta))S(\xi; \theta)' = I_k + B(\theta).$$

Now we apply the Cauchy–Bunyakovskiĭ inequality for matrices. Put $\zeta = \widehat{\theta} - a(\theta)$ and $\eta = S(\xi; \theta)$ in (3.5.12). Then

$$\begin{aligned} E_{\theta}\zeta\zeta' &= E_{\theta}(\widehat{\theta} - a(\theta))(\widehat{\theta} - a(\theta))' = D(\widehat{\theta}; \theta), \\ E_{\theta}\eta\eta' &= E_{\theta}S(\xi; \theta)S(\xi; \theta)' = I(\theta). \end{aligned}$$

It follows from (3.5.19) that

$$E_{\theta}\zeta\eta' = E_{\theta}(\widehat{\theta} - a(\theta))S(\xi; \theta)' = I_k + B(\theta).$$

This together with (3.5.12) implies inequality (3.5.13).

According to Lemma 3.5.1, inequality (3.5.13) becomes an equality if and only if

$$(\widehat{\theta}(x) - a(\theta)) = (I_k + B(\theta))I^{-1}(\theta)S(x; \theta)$$

for points $(x; \theta)$ such that $f(x; \theta) > 0$. The latter condition is equivalent to

$$(3.5.20) \quad S(x; \theta) = I(\theta)(I_k + B(\theta))^{-1}(\widehat{\theta}(x) - a(\theta)).$$

If inequality (3.5.13) becomes an equality, then

$$\det(I_k + B(\theta))^2 = \det D(\widehat{\theta}; \theta) \det I(\theta).$$

If $\det D(\widehat{\theta}; \theta)$ is far away from zero, then so is $\det(I_k + B(\theta))$, whence it follows that the inverse matrix $(I_k + B(\theta))^{-1}$ exists and is bounded. Thus the derivative $S(x; \theta)$ in (3.5.20) is bounded, $f(x; \theta) > 0$ everywhere on Θ , and equality (3.5.20) holds everywhere on Θ . Let $\theta_0, \theta \in \Theta$ and $\theta_0 + s(\theta - \theta_0) \in \Theta$ for all $s \in [0, 1]$. Then

$$\ln f(x; \theta) = \ln f(x; \theta_0) + \int_0^1 (\theta - \theta_0)' S(x; \theta_0 + s(\theta - \theta_0)) ds$$

in view of conditions (R). Thus

$$(3.5.21) \quad \ln f(x; \theta) = A(\theta)' \widehat{\theta}(x) + C(\theta) + H(x)$$

according to (3.5.20) where $C(\theta)$ and $H(x)$ are some scalar functions, and

$$A(\theta) = (A_1(\theta), \dots, A_k(\theta))'$$

is a column-vector depending only on θ . This means that representation (3.5.14) holds.

If relation (3.5.21) is satisfied, then

$$(3.5.22) \quad S(x; \theta) = \|A_{ij}(\theta)\|' \widehat{\theta}(x) + \partial B(\theta) / \partial \theta$$

and

$$(3.5.23) \quad \partial B(\theta) / \partial \theta = -\|A_{ij}(\theta)\|' a(\theta),$$

since $E_{\theta}S(\xi; \theta) = 0$. It follows from (3.5.22) and (3.5.23) that

$$S(x; \theta) = \|A_{ij}(\theta)\|' (\widehat{\theta}(x) - a(\theta)).$$

Multiplying this equality on the right by $(\widehat{\theta}(x) - a(\theta))'$, we obtain from (3.5.19) that condition (3.5.20) (which is equivalent to the case of equality in (3.5.13)) follows from (3.5.15).

Inequality (3.5.16) follows from (3.5.13), since the matrix $B(\theta)$ is zero if $\widehat{\theta}$ is an unbiased estimator. \square

All the remarks concerning the Cramér-Rao inequality that we made in Section 3.4 for the regularity conditions (R) in the one-dimensional case remain true in the multidimensional case, too.

One can prove the Cramér-Rao inequality for estimators \widehat{g} of a function $g(\theta)$ of a parameter θ in the same way as in the one-dimensional case.

Note that if an observation is a sample $\xi^{(n)} = (\xi_1, \dots, \xi_n)$, then as above

$$I_n(\theta) = nI(\theta)$$

where $I_n(\theta)$ and $I(\theta)$ are the Fisher information matrices constructed from the distribution of the sample $\xi^{(n)}$ and from the distribution of the component ξ_1 , respectively. If the regularity conditions (R) hold, then

$$(3.5.24) \quad D(\widehat{\theta}_n; \theta) \geq \frac{1}{n}(I_k + B_n(\theta))I^{-1}(\theta)(I_k + B_n(\theta))'$$

for all estimators $\widehat{\theta}_n$ constructed from the sample $\xi^{(n)}$ where

$$(3.5.25) \quad D(\widehat{\theta}_n; \theta) = E_\theta(\widehat{\theta}_n - a_n(\theta))(\widehat{\theta}_n - a_n(\theta))',$$

$$(3.5.26) \quad a_n(\theta) = E_\theta \widehat{\theta}_n = \theta + b_n(\theta), \quad b_n(\theta) = (b_n^1(\theta), \dots, b_n^k(\theta))',$$

$$(3.5.27) \quad B_n(\theta) = \|b_n^{ij}(\theta)\|, \quad b_n^{ij}(\theta) = \partial b_n^i / \partial \theta_j.$$

The case of equality in (3.5.24) can be considered by applying Theorem 3.5.2.

Efficient and asymptotically efficient estimators. The definitions of efficient and asymptotically efficient estimators in the case of a multidimensional parameter are similar to those in the case of a one-dimensional parameter. An estimator θ^* is called an *efficient estimator* of a parameter θ if the Cramér-Rao inequality for this estimator becomes an equality. If θ^* is an unbiased estimator of a parameter θ , then θ^* is efficient if (3.5.6) becomes an equality, that is, if

$$(3.5.28) \quad R(\theta^*; \theta) = I^{-1}(\theta), \quad \theta \in \Theta.$$

If θ^* is a biased estimator of a parameter θ , then it is efficient if inequality (3.5.13) becomes an equality, that is, if

$$(3.5.29) \quad D(\theta^*; \theta) = (I_k + B(\theta))I^{-1}(\theta)(I_k + B(\theta))', \quad \theta \in \Theta.$$

If an observation is a sample $\xi^{(n)} = (\xi_1, \dots, \xi_n)$, then conditions (3.5.28) and (3.5.29) become of the form

$$(3.5.30) \quad R(\theta_n^*; \theta) = \frac{1}{n}I^{-1}(\theta), \quad \theta \in \Theta,$$

$$(3.5.31) \quad D(\theta_n^*; \theta) = \frac{1}{n}(I_k + B_n(\theta))I^{-1}(\theta)(I_k + B_n(\theta))', \quad \theta \in \Theta,$$

where $I(\theta)$ is the Fisher information matrix constructed from the distribution of ξ_1 and $B_n(\theta)$ is the matrix defined by (3.5.24)–(3.5.27) for $\widehat{\theta}_n = \theta_n^*$.

Equalities (3.5.30)–(3.5.31) hold, that is, θ_n^* is an efficient estimator, only in exceptional cases. However there exist the so-called asymptotically efficient estimators and conditions for their existence are quite general. An estimator θ_n^* of a parameter θ constructed from a sample $\xi^{(n)}$ is called *asymptotically efficient* if

$$(3.5.32) \quad R(\theta_n^*; \theta) = \frac{1 + o(1)}{n}I^{-1}(\theta), \quad \theta \in \Theta,$$

where $R(\theta_n^*; \theta) = E_\theta(\theta_n^* - \theta)(\theta_n^* - \theta)'$. The estimator θ_n^* is, generally speaking, biased.

EXAMPLE 3.5.1. Let $\xi^{(n)} = (\xi_1, \xi_2, \dots, \xi_n)$ be a sample from the normal distribution $\mathcal{N}(\theta_1, \theta_2)$. Thus $\theta = (\theta_1, \theta_2)'$ and the density $f(x; \theta)$ is of the form

$$(3.5.33) \quad f(x; \theta) = \frac{1}{\sqrt{2\pi}} \exp \left\{ -\frac{x^2}{2\theta_2} + \frac{x\theta_1}{\theta_2} - \frac{\theta_1^2}{2\theta_2} - \frac{1}{2} \ln \theta_2 \right\}.$$

Consider the estimator $\hat{\theta}_n = (\hat{\theta}_{1,n}, \hat{\theta}_{2,n})$ such that

$$\hat{\theta}_{1,n} = \bar{\xi} = \frac{1}{n} \sum_{i=1}^n \xi_i, \quad \hat{\theta}_{2,n} = \frac{1}{n-1} \sum_{i=1}^n (\xi_i - \bar{\xi})^2.$$

It is clear that $\hat{\theta}_n$ is an unbiased estimator. It follows from (3.5.33) that representation (3.5.14) does not hold for the density $f_n(x; \theta)$, since

$$\begin{aligned} f_n(x; \theta) &= (2\pi)^{-n/2} \exp \left\{ -\frac{1}{2\theta_2} \sum_{i=1}^n x_i^2 + \frac{\theta_1}{\theta_2} \sum_{i=1}^n x_i - \frac{n\theta_1^2}{2\theta_2} - \frac{n}{2} \ln \theta_2 \right\} \\ &= (2\pi)^{-n/2} \exp \left\{ \frac{n\theta_1}{\theta_2} \hat{\theta}_{1,n} - \frac{n-1}{2\theta_2} \hat{\theta}_{2,n} - \frac{n}{2\theta_2} \hat{\theta}_{1,n}^2 - \frac{n\theta_1^2}{2\theta_2} - \frac{n}{2} \ln \theta_2 \right\}. \end{aligned}$$

This means that the lower bound in the multivariate Cramér–Rao inequality is not attained and therefore $\hat{\theta}_n$ is not an efficient estimator. Nevertheless $\hat{\theta}_n$ is an asymptotically efficient estimator. Below we prove this result.

First we evaluate the matrix $I(\theta)$. We have

$$S_1(x; \theta) = \frac{x - \theta_1}{\theta_2}, \quad S_2(x; \theta) = \frac{(x - \theta_1)^2}{2\theta_2^2} - \frac{1}{2\theta_2}$$

where $S_i(x; \theta) = \partial \ln f(x; \theta) / \partial \theta_i$, $i = 1, 2$. Thus

$$\begin{aligned} I_{11}(\theta) &= E_\theta \frac{(\xi_1 - \theta_1)^2}{\theta_2^2} = \frac{1}{\theta_2}, \\ I_{12}(\theta) = I_{21}(\theta) &= E_\theta \left(\frac{(\xi_1 - \theta_1)^3}{2\theta_2^3} - \frac{\xi_1 - \theta_1}{2\theta_2^2} \right) = 0, \\ I_{22}(\theta) &= E_\theta \frac{((\xi_1 - \theta_1)^2 - \theta_2)^2}{4\theta_2^4} = \frac{1}{2\theta_2^2} \end{aligned}$$

and the lower bound in the Cramér–Rao inequality is given by

$$(3.5.34) \quad \frac{1}{n} I^{-1}(\theta) = \begin{pmatrix} \theta_2/n & 0 \\ 0 & 2\theta_2^2/n \end{pmatrix}.$$

Now we evaluate the matrix $R(\hat{\theta}_n; \theta)$. We have

$$\begin{aligned} E_\theta(\hat{\theta}_{1,n} - \theta_1)^2 &= E_\theta(\bar{\xi} - \theta_1)^2 = \frac{\theta_2}{n}, \\ E_\theta(\hat{\theta}_{2,n} - \theta_2)^2 &= \frac{2\theta_2^2}{n-1}, \quad E_\theta(\hat{\theta}_{1,n} - \theta_1)(\hat{\theta}_{2,n} - \theta_2) = 0. \end{aligned}$$

The first of the latter equalities is obvious. The last of them follows from the independence of $\hat{\theta}_{1,n}$ and $\hat{\theta}_{2,n}$ (see Theorem 1.4.2). To prove the second equality we note that

$$\mathcal{L} \left(\sum_{i=1}^n \frac{(\xi_i - \bar{\xi})^2}{\theta_2} \mid P_\theta \right) = \chi^2(n-1)$$

by Theorem 1.4.2. Since $D\chi_{n-1}^2 = 2(n-1)$, we get the desired equality. Thus the matrix $R(\hat{\theta}_n; \theta)$ is given by

$$(3.5.35) \quad R(\hat{\theta}_n; \theta) = \begin{pmatrix} \theta_2/n & 0 \\ 0 & 2\theta_2^2/(n-1) \end{pmatrix}.$$

Finally we apply (3.5.34) and (3.5.35) and obtain (3.5.32). This shows that $\hat{\theta}_n$ is an asymptotically efficient estimator.

Other results related to the Cramér–Rao inequality. Inequality (3.5.12) allows one to obtain some other results corresponding to other matrices ζ and η . We restrict our consideration to the case of a one-dimensional parameter θ . Assume that the density $f(x; \theta)$ satisfies a stronger condition as compared to the regularity condition (R), namely let

- (i₁) the density $f(x; \theta)$ be continuously differentiable $m \geq 1$ times in θ ;
- (i₂) the integrals

$$K_j(\theta) = \int_{N_\theta} \frac{|\partial^j f(x; \theta) / \partial \theta^j|^2}{f(x; \theta)} \mu(dx), \quad j = 1, 2, \dots, m,$$

converge for all $\theta \in \Theta$ and, as functions of θ , be continuous on Θ where $N_\theta = \{x: f(x; \theta) \neq 0\}$.

Using the same method as that in the proof of Lemma 3.4.4 we show that conditions (i₁)–(i₂) imply that the function $a(\theta) = E_\theta \hat{\theta}$ has m continuous derivatives for an arbitrary estimator $\hat{\theta}$ if its second moment $E_\theta \hat{\theta}^2$ is locally bounded.

Let $c = (c_1, c_2, \dots, c_m)$ be some vector of \mathbf{R}^m . Put $\zeta = \hat{\theta}(\xi) - a(\theta)$ and

$$\eta = \sum_{j=1}^m c_j \frac{\partial^j f(\xi; \theta) / \partial \theta^j}{f(\xi; \theta)} I_{N_\theta}(\xi).$$

Then it follows from (3.5.12) that

(3.5.36)

$$D_\theta \hat{\theta} \geq \sup_c \frac{\left(c_1 + \sum_{j=1}^m c_j (\partial^j b(\theta) / \partial \theta^j) \right)^2}{\sum_{i,j=1}^m c_i c_j \int_{N_\theta} (\partial^i f(x; \theta) / \partial \theta^i) (\partial^j f(x; \theta) / \partial \theta^j) f^{-1}(x; \theta) \mu(dx)}$$

where $b(\theta) = E_\theta \hat{\theta} - \theta$ is the bias of the estimator $\hat{\theta}$. Inequality (3.5.36) is called the *Bhattacharyya inequality*. More details about the Bhattacharyya inequality can be found in [36].

Assume that the set N_θ does not depend on θ and let $N_\theta = N$ for all $\theta \in \Theta$. Below we avoid the regularity conditions posed on the density $f(x; \theta)$. Denote by M the set of charges m on Θ such that

$$\int_{\Theta} (1 + f(x; u)) |m(du)| < \infty.$$

Put $\zeta = \widehat{\theta} - a(\theta)$ and

$$\eta = \frac{1}{f(\xi; \theta)} \int_{\Theta} f(\xi; u) m(du) I_N(\xi) - m(\Theta)$$

in (3.5.12). Inequality (3.5.12) becomes of the form
(3.5.37)

$$D_{\theta} \widehat{\theta} \geq \sup_{m \in M} \frac{(\int_{\Theta} (a(u) - a(\theta)) m(du))^2}{\int_{\Theta} \int_{\Theta} m(du_1) m(du_2) \int_N f(x; u_1) f(x; u_2) / f(x; \theta) \mu(dx) - m^2(\Theta)}$$

and is called the *Barankin–Kiefer inequality* in this case. If the upper bound in the Barankin–Kiefer inequality is evaluated only with respect to δ -measures on Θ instead of charges $m \in M$, then (3.5.37) becomes of the form

$$(3.5.38) \quad D_{\theta} \widehat{\theta} \geq \sup_{u \in \Theta} \frac{(a(u) - a(\theta))^2}{\int_N (f(x; u) - f(x; \theta))^2 f^{-1}(x; \theta) \mu(dx)}$$

and is called the *Chapman–Robbins inequality* (cf. inequality (3.4.43)).

3.6. Integral inequalities of Cramér–Rao type

We follow the Bayes approach and obtain lower estimates for risks of estimators in this section. The corresponding inequalities can be called integral inequalities of Cramér–Rao type, since they involve risk functions integrated with respect to the a priori measure.

Efficient and superefficient estimators. Throughout this section we assume that an observation is a sample $\xi^{(n)} = (\xi_1, \dots, \xi_n)$ from a distribution for which there exists the density $f(x; \theta)$ with respect to a measure μ where $\theta = (\theta_1, \dots, \theta_k)'$ is an unknown parameter whose value belongs to a certain set $\Theta \subset \mathbf{R}^k$, $k \geq 1$. First we consider the case $k = 1$, that is, the case of a scalar parameter θ . Let an arbitrary set of regularity conditions given in Section 3.4 hold. Then the Cramér–Rao inequality

$$(3.6.1) \quad E_{\theta}(\widehat{\theta}_n - \theta)^2 \geq \frac{1}{nI(\theta)}, \quad \theta \in \Theta,$$

holds for all unbiased estimators $\widehat{\theta}_n$ of the parameter θ where $I(\theta)$ is the Fisher information evaluated with respect to the density $f(x; \theta)$, that is,

$$I(\theta) = E_{\theta} S^2(\xi_1; \theta), \quad S(x; \theta) = \partial \ln f(x; \theta) / \partial \theta.$$

The right-hand side of inequality (3.6.1) is sometimes called the *Cramér–Rao bound*. This bound is attained if an estimator is efficient. The question is whether one can improve this result for biased estimators. In other words, the question is how precise is the Cramér–Rao bound for biased estimators.

It is quite obvious that the expectation $E_{\theta}(\widehat{\theta}_n - \theta)^2$ at a fixed point θ_0 can be smaller than the Cramér–Rao bound. Indeed, this is true for $\widehat{\theta}_n \equiv \theta_0$, for example. However the latter estimator is very bad at any other point.

Below is another example of this kind. Let $\mathcal{L}(\xi_1 | P_{\theta}) = \mathcal{N}(\theta, 1)$ where

$$\theta \in \Theta = [0, \infty).$$

It is clear that the estimator $\hat{\theta}_n = n^{-1} \sum_{i=1}^n \xi_i$ is efficient. Nevertheless the estimator $\theta_n^* = 0 \vee \hat{\theta}_n$ is even better, since it decreases the mean square deviation by substituting 0 for negative values of $\hat{\theta}_n$ that are meaningless in view of the restriction that $\theta \in [0, \infty)$. On the other hand, θ_n^* is a biased estimator, since $E_\theta \theta_n^* > E_\theta \hat{\theta}_n = \theta$. We have $I(\theta) = 1$ for all $\theta \in \Theta$. At the point $\theta = 0$ we get

$$E_0 \hat{\theta}_n^2 = \frac{1}{n} = \frac{1}{nI(0)}, \quad E_0 (\theta_n^*)^2 = \frac{1}{2n} < \frac{1}{nI(0)}.$$

The improvement of the Cramér–Rao bound is explained in this example by a restriction of the domain of the estimator $\hat{\theta}_n$ to the set Θ .

Another example is due to Hodges. The improvement of the Cramér–Rao bound in the Hodges example is not due to the restriction of the set Θ and is explained by other circumstances.

Again let $\mathcal{L}(\xi_1 | P_\theta) = \mathcal{N}(\theta, 1)$ where $\theta \in \Theta = (-\infty, \infty)$. Along with an efficient estimator $\hat{\theta}_n = n^{-1} \sum_{i=1}^n \xi_i$ we consider another estimator

$$\theta_n^* = \begin{cases} \hat{\theta}_n, & \text{if } |\hat{\theta}_n| \geq n^{-1/4}, \\ \beta \hat{\theta}_n, & \text{if } |\hat{\theta}_n| < n^{-1/4}, \end{cases}$$

where $|\beta| < 1$. It is easy to see for $\theta > 0$ that

$$P_\theta \left(|\hat{\theta}_n| < n^{-1/4} \right) \leq P_\theta \left((\hat{\theta}_n - \theta)\sqrt{n} < n^{1/4} - \theta\sqrt{n} \right) = \Phi \left(n^{1/4} - \theta\sqrt{n} \right) \rightarrow 0$$

as $n \rightarrow \infty$ where $\Phi(x)$ is the standard $\mathcal{N}(0, 1)$ distribution function. A similar result holds for the case $\theta < 0$, too. If $\theta \neq 0$, then the estimator θ_n^* coincides with $\hat{\theta}_n$ on an event whose probability approaches 1 as $n \rightarrow \infty$. Thus

$$\mathcal{L}((\theta_n^* - \theta)\sqrt{n} | P_\theta) \rightarrow \mathcal{N}(0, 1)$$

as $n \rightarrow \infty$ if $\theta \neq 0$.

If $\theta = 0$ and as $n \rightarrow \infty$

$$P_0 \left(|\hat{\theta}_n| < n^{-1/4} \right) = P_0 \left(|\hat{\theta}_n \sqrt{n}| < n^{1/4} \right) = 1 - 2\Phi \left(-n^{1/4} \right) \rightarrow 1.$$

If $\theta = 0$, then the estimator θ_n^* coincides with $\beta \hat{\theta}_n$ on an event whose probability approaches 1 as $n \rightarrow \infty$. Hence

$$\mathcal{L}(\theta_n^* \sqrt{n} | P_0) \rightarrow \mathcal{N}(0, \beta^2)$$

as $n \rightarrow \infty$.

Therefore the estimator θ_n^* is asymptotically $\mathcal{N}(\theta, \sigma^2(\theta)n^{-1})$ normal for all θ where

$$\sigma^2(\theta) = \begin{cases} 1, & \text{if } \theta \neq 0, \\ \beta^2, & \text{if } \theta = 0. \end{cases}$$

Note that $\beta^2 < 1$. Thus the asymptotic variance of the estimator θ_n^* at the point $\theta = 0$ is equal to $n^{-1}\beta^2$ which is less than the lower Cramér–Rao bound

$$(nI(0))^{-1} = n^{-1}.$$

Asymptotically normal estimators whose asymptotic variance $\sigma^2(\theta)$ is such that $\sigma^2(\theta)/n \leq 1/(nI(\theta))$ and is less than $(nI(\theta))^{-1}$ for some θ are sometimes called

superefficient. The points θ for which $\sigma^2(\theta)n^{-1} < (nI(\theta))^{-1}$ are called the *points of superefficiency*.

The examples of superefficient estimators do not change our conclusion that efficient or asymptotically efficient estimators are the best. Namely Le Cam (1953) proved that an improvement of an efficient estimator can be made only at a set of points of superefficiency whose Lebesgue measure is small.

We show in this section that $\inf_{\hat{\theta}_n} E_t(\hat{\theta}_n - t)^2 = 0$ for all t and that there is a lower bound of the integral of $E_t(\hat{\theta}_n - t)^2$ that does not depend on $\hat{\theta}_n$ but still is closely related to the function $(nI(t))^{-1}$. More precisely, we obtain a lower bound for

$$(3.6.2) \quad \inf_{\hat{\theta}_n} \int_{\Theta} E_t(\hat{\theta}_n; t)^2 q(t) dt$$

for an arbitrary nonnegative weight function $q(t)$ such that $\int_{\Theta} q(t) dt = 1$. This lower bound is close to J/n where

$$(3.6.3) \quad J = \int_{\Theta} \frac{q(t)}{I(t)} dt.$$

Note that the integral in (3.6.2) can be treated as the unconditional mathematical expectation $E(\hat{\theta}_n - \theta)^2$ for the Bayes approach where the a priori measure

$$\mathbf{Q}(A) = \int_A q(t) dt$$

is a probability distribution of the parameter θ and the density $q(t)$ of \mathbf{Q} with respect to the Lebesgue measure exists. Relation (3.6.3) in this case can be rewritten as $J = EI^{-1}(\theta)$.

Integral inequalities. Let $f_n(x; t)$ be the density of the sample $\xi^{(n)}$ with respect to the measure μ if $\theta = t$. Then $f_n(x; t)q(t) = p_n(x, t)$ is the density of the joint distribution of the vector $(\xi^{(n)}, \theta)$. Denote by $N_h \subset \Theta$ the support of a function h defined on Θ . In other words $N_h = \{t: h(t) \neq 0\}$. By N we denote the support of the function $p_n(x; t)$ in $X \times \Theta$, that is, N is the support of the distribution of the vector $(\xi^{(n)}; \theta)$.

THEOREM 3.6.1. *Let the function $f_n(x; t)$ be differentiable with respect to t , while $\sqrt{I(t)}$ is integrable on every finite interval. Then*

$$(3.6.4) \quad \begin{aligned} E(\hat{\theta}_n - \theta)^2 &\geq \frac{[E(h(\theta)/q(\theta))]^2}{nE(I(\theta)[h(\theta)/q(\theta)]^2) + E(h'(\theta)/q(\theta))^2} \\ &= \frac{(\int h(t) dt)^2}{n \int I(t)h^2(t)/q(t) dt + \int (h'(t))^2/q(t) dt} \end{aligned}$$

for all differentiable functions $h(t)$ with bounded support such that $N_h \subset N_q$ and for all estimators $\hat{\theta}_n$ of the parameter θ .

PROOF. Since the support of the function $h(t)$ is bounded, we get

$$\begin{aligned} \int (f_n(x; t)h(t))' dt &= \int d(f_n(x; t)h(t)) = 0, \\ \int t(f_n(x; t)h(t))' dt &= - \int f_n(x; t)h(t) dt. \end{aligned}$$

Thus

$$(3.6.5) \quad \int_X \int_{\Theta} (\widehat{\theta}_n(x) - t)(f_n(x; t)h(t))' dt \mu(dx) = \int_X \int_{\Theta} f_n(x; t)h(t) dt \mu(dx) \\ = \int_{\Theta} h(t) dt$$

for an arbitrary estimator $\widehat{\theta}_n$. Since $N_h \subset N_q$, equality (3.6.5) holds for integrals over the set N . Multiplying and dividing (3.6.5) by $p_n(x; t)$ we obtain

$$\mathbb{E} \left[(\widehat{\theta}_n - \theta) \frac{(f_n(\xi^{(n)}; \theta) h(\theta))'}{f_n(\xi^{(n)}; \theta) q(\theta)} \right] = \int_{N_q} h(t) dt = \mathbb{E} \frac{h(\theta)}{q(\theta)}.$$

By the Cauchy–Bunyakovskiï inequality

$$(3.6.6) \quad \mathbb{E}(\widehat{\theta}_n - \theta)^2 \geq \frac{[\mathbb{E}(h(\theta)/q(\theta))]^2}{\mathbb{E} \left[(f_n(\xi^{(n)}; \theta) h(\theta))' / (f_n(\xi^{(n)}; \theta) q(\theta)) \right]^2}.$$

It remains to rewrite the latter result in the form of (3.6.4). Note that

$$(3.6.7) \quad \mathbb{E}_t \left| S_n(\xi^{(n)}; t) \right| \leq n\sqrt{I(t)}$$

and for almost all t

$$(3.6.8) \quad \mathbb{E}_t S_n(\xi^{(n)}; t) = 0$$

where $S_n(x; t) = \partial \ln f_n(x; t) / \partial t$. Estimate (3.6.7) follows from

$$\mathbb{E}_t \left| S_n(\xi^{(n)}; t) \right| \leq n \mathbb{E}_t |S(\xi_1; t)| \leq n (\mathbb{E}_t S^2(\xi_1; t))^{1/2} = n\sqrt{I(t)},$$

since

$$S_n(\xi^{(n)}; t) = \sum_{i=1}^n S(\xi_i; t) \quad \text{and} \quad S(x; t) = \partial \ln f(x; t) / \partial t.$$

To prove equality (3.6.8) we consider an arbitrary function $g(t)$ whose support is bounded and whose derivative $g'(t)$ is continuous everywhere. Then

$$\int g(t) \frac{\partial f_n(x; t)}{\partial t} dt = - \int g'(t) f_n(x; t) dt.$$

Moreover

$$\int |g(t)| \mathbb{E}_t \left| S_n(\xi^{(n)}; t) \right| dt \leq n \int |g(t)| \sqrt{I(t)} dt < \infty.$$

This implies that one can interchange the integrals:

$$\int g(t) \mathbb{E}_t S_n(\xi^{(n)}; t) dt = \int_X \int_{\Theta} g(t) \frac{\partial f_n(x; t)}{\partial t} dt \mu(dx) \\ = - \int_X \int_{\Theta} g'(t) f_n(x; t) dt \mu(dx) \\ = - \int_{\Theta} g'(t) dt = - \int_{\Theta} dg(t) = 0.$$

Since this equality holds for all g , we prove that (3.6.8) holds for almost all t .

Now we transform the right-hand side of (3.6.6):

$$\begin{aligned}
 & \mathbb{E} \left(\frac{(f_n(\xi^{(n)}; \theta) h(\theta))'}{f_n(\xi^{(n)}; \theta) q(\theta)} \right)^2 \\
 &= \mathbb{E} \left(S_n(\xi^{(n)}; \theta) \frac{h(\theta)}{q(\theta)} + \frac{h'(\theta)}{q(\theta)} \right)^2 \\
 &= \mathbb{E} \left[\left(\frac{h(\theta)}{q(\theta)} \right)^2 \mathbb{E}_\theta S_n^2(\xi^{(n)}; \theta) \right] + 2\mathbb{E} \left[\frac{h'(\theta)h(\theta)}{q^2(\theta)} \mathbb{E}_\theta S_n(\xi^{(n)}; \theta) \right] + \mathbb{E} \left(\frac{h'(\theta)}{q(\theta)} \right)^2 \\
 &= n\mathbb{E} \left(\frac{h(\theta)}{q(\theta)} \right)^2 I(\theta) + \mathbb{E} \left(\frac{h'(\theta)}{q(\theta)} \right)^2.
 \end{aligned}$$

Here we used equalities $\mathbb{E}_t S_n^2(\xi^{(n)}; t) = nI(t)$ and

$$\mathbb{E} \left(\frac{h'(\theta)h(\theta)}{q^2(\theta)} \mathbb{E}_\theta S_n(\xi^{(n)}; \theta) \right) = \int_{N_q} \frac{h'(t)h(t)}{q(t)} \mathbb{E}_t S_n(\xi^{(n)}; t) dt = 0$$

following from (3.6.8). Thus relation (3.6.4) is proved. \square

THEOREM 3.6.2. *Let the function $f_n(x; t)$ satisfy the conditions of Theorem 3.6.1. Assume that the function $h(t) = h_0(t) \equiv q(t)/I(t)$ has finite support and is differentiable. Then for all estimators $\hat{\theta}_n$ one has*

$$(3.6.9) \quad \mathbb{E}(\hat{\theta}_n - \theta)^2 \geq \frac{J}{n} \left(1 + \frac{H}{nJ} \right)^{-1} \geq \frac{J}{n} - \frac{H}{n^2}$$

where

$$H = \int \left[\left(\frac{q(t)}{I(t)} \right)' \right]^2 \frac{dt}{q(t)}.$$

PROOF. Theorem 3.6.2 follows directly from Theorem 3.6.1 since the right-hand side of (3.6.4) is equal to $J^2/(nJ + H)$ for $h(t) = q(t)/I(t)$. \square

REMARK 3.6.1. Inequalities (3.6.4) and (3.6.9) are *integral* in the sense that they provide the lower bounds for integrals of $\mathbb{E}_t(\hat{\theta}_n - t)^2$.

We see from inequalities (3.6.4) and (3.6.9) that the lower bound of $\mathbb{E}(\hat{\theta}_n - \theta)^2$ differs slightly from

$$\frac{J}{n} = \int \frac{1}{nI(t)} q(t) dt$$

for large n . The latter integral is equal to $\mathbb{E}(\theta_n^* - \theta)^2$ if θ_n^* is an efficient estimator. This indicates that one should use efficient estimators, since $\mathbb{E}(\hat{\theta}_n - \theta)^2$ attains its minimum at efficient estimators whatever function $q(t)$ is.

The following example shows that the lower bounds (3.6.4) and (3.6.9) cannot be improved in general.

EXAMPLE 3.6.1. Let $\xi^{(n)} = (\xi_1, \dots, \xi_n)$ be a sample from the distribution $\mathcal{N}(\theta, 1)$. In this case $I(\theta) = 1$. Let the parameter θ be a random variable with a smooth density $q(t)$, $t \in (-\infty, \infty)$. Then the lower bound in (3.6.9) is of the form $(n + H)^{-1}$ where

$$H = \int \frac{q'(t)}{q(t)} dt = E[(\ln q(\theta))']^2.$$

Let θ_n^Q be the Bayes estimator of the parameter θ corresponding to the a priori probability measure \mathbf{Q} whose density is $q(t)$ and let the loss function be quadratic. The estimator θ_n^Q minimizes the risk $E(\hat{\theta}_n - \theta)^2$ and coincides with the a posteriori mean $\theta_n^Q(x) = E\{\theta / \xi^{(n)} = x\}$. Thus

$$(3.6.10) \quad \begin{aligned} \theta_n^Q(x) &= \frac{\int tq(t)f_n(x;t) dt}{\int q(t)f_n(x;t) dt} = \frac{\int tq(t) \exp(n\bar{x}t - nt^2/2) dt}{\int q(t) \exp(n\bar{x}t - nt^2/2) dt} \\ &= \frac{\int tq(t) \exp(-n(t - \bar{x})^2/2) dt}{\int q(t) \exp(-n(t - \bar{x})^2/2) dt} \end{aligned}$$

where $\bar{x} = n^{-1} \sum_{i=1}^n x_i$ and $x = (x_1, \dots, x_n)$. If the function $q(t)$ is sufficiently smooth, then (3.6.10) implies that

$$\begin{aligned} \theta_n^Q(x) &= \bar{x} + \frac{q'(\bar{x})}{nq(\bar{x})} + O\left(\frac{1}{n^2}\right), \\ E(\theta_n^Q - \theta)^2 &= \frac{1}{n} - \frac{H}{n^2} + O\left(\frac{1}{n^3}\right) \end{aligned}$$

as $n \rightarrow \infty$. In particular, let $q(t) = (2\pi)^{-1/2} \exp(-t^2/2)$. Then $H = 1$ and the lower bound in (3.6.9) is $(n+1)^{-1}$. On the other hand, we learned in Example 3.1.5 that

$$E(\theta_n^Q - \theta)^2 = \frac{1}{n+1}.$$

This proves that the lower bounds in (3.6.4) and (3.6.9) cannot be improved, indeed.

The following result follows from Theorem 3.6.1. It allows one to make some conclusions concerning the points of superefficiency.

THEOREM 3.6.3. *Let the density $f_n(x;t)$ satisfy the assumptions of Theorem 3.6.1. If the interval $(a - \varepsilon, a + \varepsilon)$ belongs to Θ , then*

$$(3.6.11) \quad \max_{t \in (a-\varepsilon, a+\varepsilon)} E_t(\hat{\theta}_n - t)^2 \geq \left(n \max_{t \in (a-\varepsilon, a+\varepsilon)} I(t) + \pi^2 \varepsilon^{-2} \right)^{-1}$$

for an arbitrary estimator $\hat{\theta}_n$.

PROOF. Let $q(t) = 0$ for $t \notin (a - \varepsilon, a + \varepsilon)$. Then

$$(3.6.12) \quad \max_{t \in (a-\varepsilon, a+\varepsilon)} E(\hat{\theta} - t)^2 \geq \int_{a-\varepsilon}^{a+\varepsilon} E_t(\hat{\theta}_n - t)^2 q(t) dt = E(\hat{\theta}_n - \theta)^2.$$

Put

$$h(t) = q(t) = \frac{1}{\varepsilon} \cos^2 \frac{\pi(t-a)}{2\varepsilon}, \quad |t-a| \leq \varepsilon,$$

in Theorem 3.6.1. Then inequality (3.6.4) implies that

$$(3.6.13) \quad E(\hat{\theta} - \theta)^2 \geq \left(n \int_{a-\varepsilon}^{a+\varepsilon} I(t)q(t) dt + \int_{a-\varepsilon}^{a+\varepsilon} (q'(t))^2/q(t) dt \right)^{-1}$$

where

$$(3.6.14) \quad \begin{aligned} \int_{a-\varepsilon}^{a+\varepsilon} \frac{(q'(t))^2}{q(t)} dt &= \int_{-\varepsilon}^{\varepsilon} \left(\frac{\pi}{2\varepsilon^2} 2 \cos \frac{\pi t}{2\varepsilon} \sin \frac{\pi t}{2\varepsilon} \right)^2 \varepsilon \cos^{-2} \frac{\pi t}{2\varepsilon} dt \\ &= \frac{1}{\varepsilon^2} \int_{-1}^1 \pi^2 \sin^2 \frac{\pi t}{2} dt = \frac{\pi^2}{\varepsilon^2}. \end{aligned}$$

Now relations (3.6.12)–(3.6.14) yield inequality (3.6.11). \square

REMARK 3.6.2. It is not hard to show that the minimum of the functional $\int_{-1}^1 (q'(t))^2 q^{-1}(t) dt$ in the class of all differentiable densities $q(t)$ whose support belongs to $[-1, 1]$ is attained for the density $q(t) = \cos^2(\pi t/2)$.

REMARK 3.6.3. Theorem 3.6.3 implies that the length of the interval of values of θ for which $\hat{\theta}_n$ is superefficient does not exceed $O(n^{-1/2})$.

Integral inequalities for nondifferentiable functions $q(t)/I(t)$. If the function $h(t) = q(t)/I(t)$ does not satisfy the assumptions of Theorem 3.6.1, then one can estimate the asymptotic behavior of $E(\hat{\theta}_n - \theta)^2$ by using the following result.

THEOREM 3.6.4. *Let the function $f_n(x; t)$ satisfy the assumptions of Theorem 3.6.1. Let the functions $h_\varepsilon(t)$ depend on a positive parameter ε , and satisfy the assumptions of Theorem 3.6.1 and*

- 1) $h_\varepsilon(t) \leq h_0(t) = q(t)/I(t)$ for all $\varepsilon > 0$,
- 2) $H(\varepsilon) = \int (h'_\varepsilon(t))^2/q(t) dt < \infty$ for all $\varepsilon > 0$.

Then for all $\varepsilon > 0$

$$E(\hat{\theta}_n - \theta)^2 \geq \frac{(\int h_\varepsilon(t) dt)^2}{nJ + H(\varepsilon)}.$$

PROOF. It is necessary to put $h(t) = h_\varepsilon(t)$ in Theorem 3.6.1. \square

Theorem 3.6.4 implies the following useful result.

THEOREM 3.6.5. *Let the function $f_n(x; t)$ satisfy the assumptions of Theorem 3.6.1. If the function $q(t)$ is Riemann integrable and $J < \infty$, then*

$$(3.6.15) \quad E(\hat{\theta}_n - \theta)^2 \geq \frac{J}{n}(1 + o(1))$$

as $n \rightarrow \infty$.

PROOF. Consider the following functions:

$$\begin{aligned} \hat{q}_\varepsilon(t) &= \min_{|u| \leq \varepsilon} q(t+u), & q_\varepsilon(t) &= \hat{q}_\varepsilon(t)I(\hat{q}_\varepsilon(t) \geq \varepsilon), \\ I_\varepsilon(t) &= \max(\varepsilon, I(t)), & h_\varepsilon(t) &= \frac{1}{2\varepsilon} \int_{t-\varepsilon}^{t+\varepsilon} \frac{q_\varepsilon(s)}{I_\varepsilon(s)} ds. \end{aligned}$$

It is clear that the support of the function $h_\varepsilon(t)$ is bounded, $h_\varepsilon(t)$ is differentiable for all $\varepsilon > 0$, and $h_\varepsilon(t) \leq h_0(t) = q(t)/I(t)$.

Since the function $q(t)$ is Riemann integrable, we obtain $q_\varepsilon(t) \nearrow q(t)$ almost everywhere as $\varepsilon \rightarrow 0$. This result follows from

$$(3.6.16) \quad \int_a^b [q(t) - q_\varepsilon(t)] dt \downarrow 0$$

as $\varepsilon \rightarrow 0$ for all $a, b \in \Theta$, $-\infty < a < b < \infty$. Moreover

$$\sum_k q_\delta(2k\delta)2\delta \uparrow \int_a^b q(t) dt \quad \text{and} \quad \sum_k q_\delta((2k+1)\delta)2\delta \uparrow \int_a^b q(t) dt$$

as $\delta \rightarrow 0$. Therefore, as $\varepsilon \rightarrow 0$

$$\begin{aligned} \int_a^b q_\varepsilon(t) dt &\geq \sum_k q_{2\varepsilon}(2k\varepsilon)2\varepsilon = \frac{1}{2} \left(\sum_k q_{2\varepsilon}(4k\varepsilon)4\varepsilon + \sum_k q_{2\varepsilon}((4k+2)\varepsilon)4\varepsilon \right) \\ &\rightarrow \int_a^b q(t) dt. \end{aligned}$$

Thus relation (3.6.16) is proved, whence the convergence $q_\varepsilon(t) \uparrow q(t)$ follows.

The convergence $q_\varepsilon(t) \uparrow q(t)$ implies that

$$\begin{aligned} \frac{q_\varepsilon(t)}{I_\varepsilon(t)} \uparrow \frac{q(t)}{I(t)} &= h_0(t), \\ \int h_\varepsilon(t) dt &= \int \frac{dt}{2\varepsilon} \int_{-e}^\varepsilon \frac{q_\varepsilon(t+s)}{I_\varepsilon(t+s)} ds = \frac{1}{2\varepsilon} \int_{-\varepsilon}^\varepsilon ds \int \frac{q_\varepsilon(t)}{I_\varepsilon(t)} dt \\ &= \int \frac{q_\varepsilon(t)}{I_\varepsilon(t)} dt \uparrow \int \frac{q(t)}{I(t)} dt = J \end{aligned}$$

as $\varepsilon \rightarrow 0$. Moreover

$$\begin{aligned} |h'_\varepsilon(t)| &= \frac{1}{2\varepsilon} \left| \frac{q_\varepsilon(t+\varepsilon)}{I_\varepsilon(t+\varepsilon)} - \frac{q_\varepsilon(t-\varepsilon)}{I_\varepsilon(t-\varepsilon)} \right| \leq \frac{q(t)}{\varepsilon^2}, \\ H(\varepsilon) &= \int \frac{(h'_\varepsilon(t))^2}{q(t)} dt \leq \int \frac{q(t)}{\varepsilon^4} dt = \frac{1}{\varepsilon^4}. \end{aligned}$$

Putting $\varepsilon = n^{-1/5}$ in Theorem 3.6.4 we obtain

$$E(\hat{\theta}_n - \theta)^2 \geq \frac{(\int h_\varepsilon(t) dt)^2}{nJ + H(\varepsilon)} \geq \frac{J^2 + o(1)}{nJ + n^{4/5}}$$

as $n \rightarrow \infty$, whence the desired estimate (3.6.15) follows. □

Asymptotically Bayes and asymptotically minimax estimators. One of the main results following from the above integral inequalities can be stated as follows. If an efficient or, at least, an asymptotically efficient estimator exists, then any other estimator is asymptotically “worse”. Below we introduce the notions of the asymptotically Bayes and the asymptotically minimax estimators. We consider the quadratic loss function and the a priori measure \mathbf{Q} and assume without loss of generality that \mathbf{Q} is a probability measure for which the density $q(t)$ exists.

An estimator θ_n^* of a parameter θ is called *asymptotically Bayes* (with respect to the quadratic loss function and the a priori measure \mathbf{Q}) if

$$(3.6.17) \quad \limsup_{n \rightarrow \infty} \left[\text{En}(\theta_n^* - \theta)^2 - \text{En}(\hat{\theta}_n - \theta)^2 \right] \leq 0$$

for an arbitrary estimator $\hat{\theta}_n$. An estimator θ_n^* is called *asymptotically R-Bayes* if

$$(3.6.18) \quad \text{En}(\theta_n^* - \theta)^2 = J + o(1)$$

as $n \rightarrow \infty$. In other words, an estimator is asymptotically R-Bayes if the lower bound for the mean square deviation given by Theorems 3.6.2 and 3.6.5 is attained at this estimator. Another name for this estimator is *asymptotically R-efficient in the mean square sense*.

The following result contains a relationship between asymptotically Bayes and asymptotically R-Bayes estimators.

THEOREM 3.6.6. *Let all the assumptions of Theorem 3.6.1 hold. If the function $q(t)$ is Riemann integrable, then every R-Bayes estimator is an asymptotically Bayes estimator.*

PROOF. Let θ_n^* be an asymptotically R-Bayes estimator, that is, (3.6.18) holds. According to Theorem 3.6.5

$$\liminf_{n \rightarrow \infty} \text{En}(\hat{\theta}_n - \theta)^2 \geq J$$

for an arbitrary estimator $\hat{\theta}_n$. This together with (3.6.18) implies (3.6.17) for the estimator θ_n^* , that is, θ_n^* is an asymptotically Bayes estimator. \square

It is clear that if an asymptotically R-Bayes estimator exists, then every asymptotically Bayes estimator is an asymptotically R-Bayes estimator.

The equivalence of all asymptotically R-Bayes estimators is established in the following result.

THEOREM 3.6.7. *Let all the assumptions of Theorem 3.6.1 hold. Assume that the function $q(t)$ is Riemann integrable. If θ_n^* and θ_n^{**} are two asymptotically R-Bayes estimators, then they are asymptotically equivalent in the following sense:*

$$(3.6.19) \quad \text{En}(\theta_n^* - \theta_n^{**})^2 \rightarrow 0, \quad (\theta_n^* - \theta_n^{**})\sqrt{n} \rightarrow 0$$

as $n \rightarrow \infty$ where the second relation means the convergence in probability with respect to the joint distribution of $\xi^{(n)}$ and θ .

PROOF. It follows from (3.6.18) that

$$(3.6.20) \quad \lim_{n \rightarrow \infty} \text{En}(\theta_n^* - \theta)^2 = \lim_{n \rightarrow \infty} \text{En}(\theta_n^{**} - \theta)^2 = J.$$

Let $\hat{\theta}_n = (\theta_n^* + \theta_n^{**})/2$. Relation (3.6.20) implies that

$$(3.6.21) \quad \lim_{n \rightarrow \infty} \text{En}(\hat{\theta}_n - \theta)^2 = J.$$

It is easy to show that

$$(\hat{\theta}_n - \theta)^2 + \left(\frac{\theta_n^* - \theta_n^{**}}{2} \right)^2 = \frac{(\theta_n^* - \theta)^2 + (\theta_n^{**} - \theta)^2}{2}.$$

This equality together with (3.6.20) and (3.6.21) yields

$$\lim_{n \rightarrow \infty} E n(\theta_n^* - \theta_n^{**})^2 = 0$$

and the first relation in (3.6.19) is proved. The second relation in (3.6.19) obviously follows from the first one. \square

Further we consider the asymptotically minimax approach. An estimator θ_n^* is called *asymptotically minimax* if

$$(3.6.22) \quad \limsup_{n \rightarrow \infty} \sup_{t \in \Theta} E_t n(\theta_n^* - t)^2 \leq \liminf_{n \rightarrow \infty} \sup_{t \in \Theta} E_t n(\hat{\theta}_n - t)^2$$

for all estimators $\hat{\theta}_n$.

The following result contains sufficient conditions that an estimator is asymptotically minimax.

THEOREM 3.6.8. *Let the Fisher information $I(\theta)$ exist and be continuous. If*

$$(3.6.23) \quad \limsup_{n \rightarrow \infty} \sup_{t \in \Theta} E_t n(\theta_n^* - t)^2 \leq \sup_{t \in \Theta} I^{-1}(t),$$

then θ_n^ is an asymptotically minimax estimator.*

PROOF. It is sufficient to show that

$$(3.6.24) \quad \liminf_{n \rightarrow \infty} \sup_{t \in \Theta} E_t n(\hat{\theta}_n - t)^2 \geq \sup_{t \in \Theta} I^{-1}(t)$$

for an arbitrary estimator $\hat{\theta}_n$. Let \mathbf{Q} be an arbitrary probability measure on Θ whose density $q(t)$ is smooth. It is obvious that

$$(3.6.25) \quad \sup_{t \in \Theta} E_t n(\hat{\theta}_n - t)^2 \geq \int E_t n(\hat{\theta}_n - t)^2 q(t) dt.$$

The right-hand side of (3.6.25) is greater than or equal to $J - H/n$ according to Theorem 3.6.2. Thus (3.6.25) implies that

$$(3.6.26) \quad \liminf_{n \rightarrow \infty} \sup_{t \in \Theta} E_t n(\hat{\theta}_n - t) \geq \int I^{-1}(t) q(t) dt.$$

Since $q(t)$ is arbitrary, it can be specified such that

$$(3.6.27) \quad \int I^{-1}(t) q(t) dt \geq \sup_{t \in \Theta} I^{-1}(t) - \varepsilon$$

for a given $\varepsilon > 0$. The number ε is arbitrary and thus (3.6.26) and (3.6.27) imply relation (3.6.24). Taking into account (3.6.22), we obtain from (3.6.24) and (3.6.23) that θ_n^* is an asymptotically minimax estimator. \square

Remarks concerning the multidimensional case. All the results of this section can be proved for the case of a multidimensional parameter $\theta \in \Theta \subset \mathbf{R}^k$, $k \geq 1$.

In particular, Theorem 3.6.5, one of the main results of this section, is of the following form in the multidimensional case:

$$E(\theta_n^* - \theta)(\theta_n^* - \theta)' \geq \frac{1 + o(1)}{n} EI^{-1}(\theta)$$

where $I(t)$ is the Fisher information matrix.

The results for asymptotic Bayes and asymptotic minimax estimators are also valid in the multidimensional case if the quality of an estimator is measured by

$$v(\theta_n^*) = E(\theta_n^* - t)'V(\theta_n^* - \theta)$$

where V is a certain nonnegative definite matrix.

Bayes and minimax (or asymptotically Bayes and asymptotically minimax) estimators can be defined in the multidimensional case as estimators whose qualities satisfy the corresponding inequalities for all nonnegative definite matrices V .

REMARK 3.6.4. Other approaches to integral inequalities of Cramér–Rao type can be found in Chapter 5 of [22] where the estimates of the Shannon information contained in an observation $\xi^{(n)}$ and in an estimator $\hat{\theta}_n$ with respect to a random parameter θ are used. The corresponding results are called *information inequalities* in [22].

Sufficient Statistics

In the preceding section we discussed the problem on how to construct different kinds of optimal estimators, namely Bayes, minimax, efficient, asymptotically Bayes, asymptotically minimax, asymptotically efficient, and others. In this section, we introduce the so-called sufficient estimators that allow one to construct optimal estimators by using a sufficient statistic instead of an observation. Sufficient statistics play an important role in mathematical statistics in general and in the theory of estimation in particular.

4.1. Sufficient statistics and a theorem on factorization

Conditional expectations, conditional probabilities, and sufficient statistics. Let (Ω, \mathcal{F}, P) be a probability space, let ξ be a nonnegative random variable, and let \mathcal{G} be a σ -algebra, $\mathcal{G} \subset \mathcal{F}$. A generalized nonnegative random variable $E(\xi/\mathcal{G})$ (the extended form of this notation is $E(\xi/\mathcal{G})(\omega)$, $\omega \in \Omega$) is called the *conditional expectation of the random variable ξ with respect to the σ -algebra \mathcal{G}* if $E(\xi/\mathcal{G})$ is \mathcal{G} -measurable and for all $A \in \mathcal{G}$

$$(4.1.1) \quad \int_A \xi(\omega) P(d\omega) = \int_A E(\xi/\mathcal{G})(\omega) P(d\omega)$$

or, equivalently,

$$(4.1.2) \quad EI_A \xi = EI_A E(\xi/\mathcal{G}) \quad \text{for all } A \in \mathcal{G}$$

where $I_A = I_A(\omega)$ is the indicator of the set A . The conditional expectation $E(\xi/\mathcal{G})$ of a random variable ξ with respect to a σ -algebra \mathcal{G} is well defined if

$$(4.1.3) \quad \min(E(\xi^+/\mathcal{G}), E(\xi^-/\mathcal{G})) < \infty \quad (\text{P-a.s.}).$$

In this case

$$(4.1.4) \quad E(\xi/\mathcal{G}) = E(\xi^+/\mathcal{G}) - E(\xi^-/\mathcal{G}) \quad (\text{P-a.s.}).$$

Here $\xi^+ = 0 \vee \xi$ and $\xi^- = -(0 \wedge \xi)$. Note that the conditional expectation $E(\xi/\mathcal{G})$ exists if $\xi \geq 0$. Indeed, let $\mathbf{Q}(A) = EI_A \xi$, $A \in \mathcal{G}$. Then \mathbf{Q} is a measure on (Ω, \mathcal{G}) and it is absolutely continuous with respect to the measure P . According to the Radon–Nikodym theorem, there exists a generalized nonnegative \mathcal{G} -measurable random variable $E(\xi/\mathcal{G})$ such that

$$\mathbf{Q}(A) = \int_A E(\xi/\mathcal{G})(\omega) P(d\omega) \quad \text{for all } A \in \mathcal{G}.$$

Thus $E(\xi/\mathcal{G})(\omega) = d\mathbf{Q}/dP(\omega)$ (P-a.s.) is the Radon–Nikodym derivative of the measure \mathbf{Q} with respect to the measure P ; both measures \mathbf{Q} and P are considered on the space (Ω, \mathcal{G}) .

Let $B \in \mathcal{F}$. Then the conditional expectation $E(I_B/\mathcal{G})$ is called the *conditional probability of an event B with respect to a σ -algebra \mathcal{G}* , $\mathcal{G} \subset \mathcal{F}$. The conditional probability is denoted by $P(B/\mathcal{G})$. Therefore the conditional probability of an event $B \in \mathcal{F}$ with respect to the σ -algebra \mathcal{G} is a \mathcal{G} -measurable random variable $P(B/\mathcal{G})$ such that

$$(4.1.5) \quad P(A \cap B) = EI_A P(B/\mathcal{G}) = \int_A P(B/\mathcal{G}) dP \quad \text{for all } A \in \mathcal{G}.$$

Let ξ be a random variable and let \mathcal{G}_η be the σ -algebra generated by some random element η . Then the conditional expectation $E(\xi/\mathcal{G}_\eta)$, if it exists, is denoted by $E(\xi/\eta)$ or by $E(\xi/\eta)(\omega)$ and is called the *conditional expectation of ξ with respect to η* . The conditional probability $P(B/\mathcal{G}_\eta)$ is denoted by $P(B/\eta)$ or by $P(B/\eta)(\omega)$ and is called the *conditional probability of an event $B \in \mathcal{F}$ with respect to η* . Let $\eta = \eta(\omega)$ be a random element assuming values in a measurable space (Y, \mathcal{S}) . Since $E(\xi/\eta)$ is a \mathcal{G}_η -measurable function, there exists a real Borel function $m = m(y)$ defined on (Y, \mathcal{S}) , assuming values in $\bar{\mathbf{R}} = [-\infty, \infty]$, and such that for all $\omega \in \Omega$

$$m(\eta(\omega)) = E(\xi/\eta)(\omega).$$

This function $m(y)$ is denoted by $E(\xi/y) = E(\xi/\eta = y)$ and is called the *conditional expectation of a random variable ξ with respect to an event $\{\eta = y\}$* or *conditional expectation of ξ given $\eta = y$* .

According to definitions (4.1.1)–(4.1.4) we have

$$(4.1.6) \quad EI_A \xi = EI_A E(\xi/\eta) = EI_A m(\eta) \quad \text{for all } A \in \mathcal{G}_\eta.$$

Changing the variables in the integral we obtain

$$(4.1.7) \quad EI_{\{\eta \in B\}} m(\eta) = \int_{\{\omega: \eta(\omega) \in B\}} m(\eta) dP = \int_B m(y) P_\eta(dy) \quad \text{for all } B \in \mathcal{S}$$

where $\{\omega: \eta(\omega) \in B\} \in \mathcal{G}_\eta$ for all $B \in \mathcal{S}$ and P_η is the probability distribution of the random element η on (Y, \mathcal{S}) . Thus equalities (4.1.6) and (4.1.7) imply that $m = m(y)$ is a Borel function defined on (Y, \mathcal{S}) and such that for all $B \in \mathcal{S}$

$$(4.1.8) \quad EI_{\{\eta \in B\}} \xi = \int_{\{\omega: \eta(\omega) \in B\}} \xi(\omega) P(d\omega) = \int_B m(y) P_\eta(dy).$$

Relation (4.1.8) can be used as an alternative definition of the conditional expectation $E(\xi/\eta = y) = E(\xi/y) = m(y)$.

The conditional expectation $E(I_A/\eta = y)$ is called the *conditional probability of an event $A \in \mathcal{F}$ given $\eta = y$* ; this expectation is denoted by

$$P(A/\eta = y) = P(A/y).$$

It is clear that the conditional probability $P(A/\eta = y)$ can be defined as a measurable function defined on (Y, \mathcal{S}) , assuming values in $([0, 1], \mathcal{B}([0, 1]))$, and such that for all $B \in \mathcal{S}$

$$(4.1.9) \quad P(A \cap \{\eta \in B\}) = \int_B P(A/\eta = y) P_\eta(dy)$$

(see (4.1.5)). Note that the conditional expectation $E(\xi/\mathcal{G})$ can be defined in a similar way for rather general random elements ξ if the expectation $E\xi$ exists. A

detailed treatment of this topic as well as a discussion of properties of conditional expectations can be found in [30].

Let ξ be a random element assuming values in a measurable space (X, \mathcal{B}) ; let the distribution of ξ be a probability measure belonging to a family $\mathcal{P} = \{P_\theta, \theta \in \Theta\}$ where $\theta = (\theta_1, \theta_2, \dots, \theta_k)'$ is an unknown parameter, $\theta \in \Theta \subset \mathbf{R}^k$, $k \geq 1$. We call ξ an *observation*. An arbitrary measurable function $T = T(x)$ mapping (X, \mathcal{B}) into some measurable space (Y, \mathcal{S}) is called a *statistic*.

For a fixed $\theta \in \Theta$ consider the probability space $(X, \mathcal{B}, P_\theta)$. Let \mathcal{B}_T be the σ -algebra in (X, \mathcal{B}) generated by the statistic $T = T(x)$. It is clear that

$$\mathcal{B}_T = T^{-1}(\mathcal{S}) \subset \mathcal{B}$$

where $T^{-1}(\mathcal{S})$ is the minimal σ -algebra generated by the family of events

$$\{x: T(x) \in B\}, \quad B \in \mathcal{S}.$$

According to definition (4.1.5), a \mathcal{B}_T -measurable function

$$P_\theta(A/\mathcal{B}_T) = P_\theta(A/\mathcal{B}_T)(x)$$

such that

$$(4.1.10) \quad P_\theta(A \cap B) = \int_B P_\theta(A/\mathcal{B}_T)(x) P_\theta(dx) \quad \text{for all } B \in \mathcal{B}_T$$

is called the conditional probability measure P_θ of the set $A \in \mathcal{B}$ with respect to the σ -algebra \mathcal{B}_T . By definition (4.1.9), the conditional probability measure P_θ of the set $A \in \mathcal{B}$ given $T = y$ is a measurable function

$$P_\theta(A/\eta = y) = P_\theta(A/y)$$

defined on (Y, \mathcal{S}) , assuming values in $([0, 1], \mathcal{B}([0, 1]))$, and such that for all $B \in \mathcal{S}$

$$(4.1.11) \quad P_\theta(A \cap \{x: T(x) \in B\}) = \int_B P_\theta(A/y) P_{\theta, T}(dy)$$

where $P_{\theta, T}$ is the distribution of the statistic T defined by

$$P_{\theta, T}(B) = P_\theta\{x: T(x) \in B\}, \quad B \in \mathcal{S}.$$

Relations (4.1.10) and (4.1.11) imply that

$$(4.1.12) \quad P_\theta(A/T(x)) = P_\theta(A/\mathcal{B}_T)(x) \quad (P_\theta\text{-a.s.}) \quad \text{for all } A \in \mathcal{B}.$$

A statistic $T = T(x)$ is called a *sufficient statistic for a family* $\mathcal{P} = \{P_\theta, \theta \in \Theta\}$ (or for a parameter θ) if for any $A \in \mathcal{B}$ there exists a measurable function

$$\psi_A = \psi_A(y)$$

defined on (Y, \mathcal{S}) , that depends on A and y and does not depend on θ , and such that

$$(4.1.13) \quad P_\theta(A/T(x)) = \psi_A(T(x)) \quad (P_\theta\text{-a.s.}).$$

This property means that the conditional distribution of the observation ξ given a fixed value of the statistic T does not depend on the parameter θ . This means that the fact that a sampling point $x \in X$ lies on the surface $T(x) = y$ gives no additional information about the parameter θ . In other words, the statistic T exhausts all the information about θ that is contained in the sample. This explains the name for

a *sufficient statistic*: knowledge of $T(x)$ is *sufficient* to construct an estimator of the parameter θ , while the other information included into the observed point x is useless.

EXAMPLE 4.1.1. Let an observation be a sample $\xi^{(n)} = (\xi_1, \dots, \xi_n)$ from a Poisson distribution with parameter θ . Consider a statistic $T_n = \sum_{i=1}^n \xi_i$. This statistic obviously has the Poisson distribution with parameter $n\theta$. Note that

$$\left\{ \xi^{(n)} = x, \sum_{i=1}^n \xi_i = y \right\} = \begin{cases} \{ \xi^{(n)} = x \}, & \text{if } \sum_{i=1}^n x_i = y, \\ \emptyset, & \text{if } \sum_{i=1}^n x_i \neq y \end{cases}$$

where $x = (x_1, \dots, x_n)$ and $x_i, y \in \{0, 1, \dots\}$ for all $i = 1, \dots, n$. Then

$$(4.1.14) \quad P_\theta \left\{ \xi^{(n)} = x / T_n = y \right\} = \begin{cases} \frac{P_\theta \{ \xi^{(n)} = x \}}{P_\theta \{ T_n = y \}}, & \text{if } \sum_{i=1}^n x_i = y, \\ 0, & \text{if } \sum_{i=1}^n x_i \neq y. \end{cases}$$

If $\sum_{i=1}^n x_i = y$, then

$$\frac{P_\theta \{ \xi^{(n)} = x \}}{P_\theta \{ T_n = y \}} = \left(e^{-n\theta} \frac{(n\theta)^y}{y!} \right)^{-1} \prod_{i=1}^n e^{-\theta} \frac{\theta^{x_i}}{x_i!} = \frac{y!}{n^y} \prod_{i=1}^n \frac{1}{x_i!}.$$

Thus relation (4.1.14) implies that the conditional probability

$$P_\theta \left(\xi^{(n)} = x / T_n = y \right)$$

does not depend on the parameter θ . This means that the statistic T_n is sufficient according to definition (4.1.13).

Relation (4.1.12) suggests the following definition of a sufficient σ -algebra. We say \mathcal{B}_T is a *sufficient σ -algebra* if the statistic T is sufficient. The notion of a sufficient statistic is important in many problems, however the notion of a σ -algebra is more convenient, at least from the point of view of the theory, than that of a sufficient statistic. Note that there are examples of sufficient σ -algebras that are not generated by any sufficient statistics assuming values in a given measurable space [3].

Dominated families of distributions. Let \mathcal{P} be a family of probability measures defined on a measurable space (X, \mathcal{B}) and let μ be some σ -finite measure on (X, \mathcal{B}) . We say that a family \mathcal{P} is *dominated by the measure μ* if every measure $P \in \mathcal{P}$ is absolutely continuous with respect to μ . A family \mathcal{P} is called *dominated* if there is a σ -finite measure dominating the family \mathcal{P} . Note that if a family \mathcal{P} is dominated, then there exists a finite dominating measure. Indeed, let a family \mathcal{P} be dominated by a σ -finite measure μ and let $X = \bigcup_{i=1}^{\infty} A_i$ where $\mu(A_i) < \infty$ for all $i = 1, 2, \dots$ and $A_i \cap A_j = \emptyset$ for $i \neq j$. Then the measure ν defined by

$$\nu(A) = \sum_{i=1}^{\infty} 2^{-i} \mu(A \cap A_i) / \mu(A_i), \quad A \in \mathcal{B},$$

also dominates the family \mathcal{P} ; moreover the measure ν is finite.

Let two families of measure \mathcal{M} and \mathcal{N} be given. We say that a *family \mathcal{M} is dominated by a family \mathcal{N}* if every measure of the family \mathcal{N} dominates the family \mathcal{M} . The families of measures \mathcal{M} and \mathcal{N} are called *equivalent* if the family \mathcal{M} is dominated by the family \mathcal{N} , and the family \mathcal{N} is dominated by the family \mathcal{M} .

THEOREM 4.1.1. *A family of probability measures \mathcal{P} is dominated a σ -finite measure if and only if the family \mathcal{P} contains a countable equivalent subfamily.*

PROOF. First we assume that the family \mathcal{P} contains a countable equivalent subfamily $\{P_1, P_2, \dots\}$. Then the family \mathcal{P} is dominated by the measure

$$\mu = \sum_{n=1}^{\infty} 2^{-n} P_n.$$

Conversely let the family \mathcal{P} be dominated by a σ -finite measure μ . Without loss of generality we assume that the measure μ is finite. Let \mathcal{K} be the class of all probability measures \mathbf{Q} of the form $\sum c_i P_i$ where $P_i \in \mathcal{P}$, all numbers c_i are positive, and $\sum c_i = 1$. The class \mathcal{K} is dominated by the measure μ . For $\mathbf{Q} \in \mathcal{K}$ we denote by $q(x) = d\mathbf{Q}/d\mu(x)$ the density of the measure \mathbf{Q} with respect to μ .

Our current goal is to prove the following assertion (which is equivalent to the statement of the theorem): *there exists a measure $\mathbf{Q}_0 \in \mathcal{K}$ such that the equality $\mathbf{Q}_0(A) = 0$ implies $\mathbf{Q}(A) = 0$ for all $\mathbf{Q} \in \mathcal{K}$.*

Consider the class S of sets $C \in \mathcal{B}$ for which there exists a measure $\mathbf{Q} \in \mathcal{K}$ such that $q(x) > 0$ almost surely with respect to the measure μ on a set C with $\mathbf{Q}(C) > 0$. Let $\mu(C_i) \rightarrow \sup_{C \in S} \mu(C)$ as $i \rightarrow \infty$ where $C_i \in S$ and $q_i(x) > 0$ almost surely with respect to the measure μ on a set C_i ($q_i(x)$ corresponds to \mathbf{Q}_i which in turn corresponds to C_i). Let $C_0 = \bigcup_{i=1}^{\infty} C_i$. Then $q_0^*(x) = \sum_{i=1}^{\infty} c_i q_i(x)$ coincides μ -almost surely with the density $d\mathbf{Q}_0/d\mu(x)$ where $\mathbf{Q}_0 = \sum_{i=1}^{\infty} c_i \mathbf{Q}_i$. It is clear that $q_0^*(x) > 0$ almost surely with respect to the measure μ on the set C_0 , whence $C_0 \in S$.

Assume that $\mathbf{Q}_0(A) = 0$. Let \mathbf{Q} be another measure of the class \mathcal{K} and let $C = \{x: q(x) > 0\}$ and $q(x) = d\mathbf{Q}/d\mu(x)$. Then $\mathbf{Q}_0(A \cap C_0) = 0$, whence $\mu(A \cap C_0) = 0$ and $\mathbf{Q}(A \cap C_0) = 0$. We also have $\mathbf{Q}(A \cap \bar{C}_0 \cap \bar{C}) = 0$ where $\bar{C} = X \setminus C$ and $\bar{C}_0 = X \setminus C_0$. Now we prove that $\mathbf{Q}(A \cap \bar{C}_0 \cap C) = 0$. Assume the converse, namely let $\mathbf{Q}(A \cap \bar{C}_0 \cap C) > 0$. This implies

$$(4.1.15) \quad \mu(C_0 \cup (A \cap \bar{C}_0 \cap C)) = \mu(C_0) + \mu(A \cap \bar{C}_0 \cap C) > \mu(C_0),$$

since $\mu(A \cap \bar{C}_0 \cap C) > 0$ in view of the inequality $\mathbf{Q}(A \cap \bar{C}_0 \cap C) > 0$ and \mathbf{Q} is absolutely continuous with respect to μ . Inequality (4.1.15) contradicts the equality $\mu(C_0) = \sup_{C \in S} \mu(C)$. Thus $\mathbf{Q}(A \cap \bar{C}_0 \cap C) = 0$ and

$$\mathbf{Q}(A) = \mathbf{Q}(A \cap C_0) + \mathbf{Q}(A \cap \bar{C}_0 \cap C) + \mathbf{Q}(A \cap \bar{C}_0 \cap \bar{C}) = 0.$$

Therefore the equality $\mathbf{Q}_0(A) = 0$ implies that $\mathbf{Q}(A) = 0$ for all $\mathbf{Q} \in \mathcal{K}$. \square

Theorem on the factorization. Let ξ be an observation that is a random element assuming values in a measurable space (X, \mathcal{B}) and whose distribution belongs to a family of probability measures $\mathcal{P} = \{P_\theta, \theta \in \Theta\}$ where $\theta = (\theta_1, \theta_2, \dots, \theta_k)'$ is an unknown parameter $\theta \in \Theta \subset \mathbf{R}^k$, $k \geq 1$.

The following result, known as the *Neyman-Fisher factorization criterion*, contains a necessary and sufficient condition for a statistic to be sufficient for a dominated family \mathcal{P} . The short name of this result is the *factorization criterion*. The first result of this type is obtained by Fisher (1920) and rediscovered by Neyman (1935). It is proved for general dominated families by Halmos and Savage (1949). Further generalization is due to Bahadur (1954). The result below is closer to the theorem of Bahadur.

THEOREM 4.1.2. *Let a family $\mathcal{P} = \{P_\theta, \theta \in \Theta\}$ be dominated by a measure μ and let $T = T(x)$ be a statistic mapping (X, \mathcal{B}) into a measurable space (Y, \mathcal{S}) . The statistic T is sufficient for a family \mathcal{P} if and only if the density*

$$f(x; \theta) = dP_\theta/d\mu(x)$$

admits the factorization

$$(4.1.16) \quad f(x; \theta) = g(T(x); \theta)r(x) \quad (\mu\text{-a.s.}) \quad \text{for all } \theta \in \Theta$$

where $g(y; \theta)$ is a nonnegative \mathcal{S} -measurable function for all $\theta \in \Theta$ and $r(x)$ is a nonnegative \mathcal{B} -measurable function.

PROOF. Since the family \mathcal{P} is dominated, Theorem 4.1.1 implies that there is a countable subfamily $\{P_{\theta_1}, P_{\theta_2}, \dots\} \subset \mathcal{P}$ that is equivalent to the family \mathcal{P} . Consider the probability measure $\lambda = \sum_i c_i P_{\theta_i}$ where $c_i > 0$ for all $i = 1, 2, \dots$ and $\sum_i c_i = 1$. It is obvious that the measure λ dominates the family \mathcal{P} . Denote by $p(x; \theta) = dP_\theta/d\lambda(x)$ the density of the measure P_θ with respect to the measure λ .

Necessity. Let T be a sufficient statistic for the family \mathcal{P} . According to definition (4.1.13), for every $A \in \mathcal{B}$ there exists an \mathcal{S} -measurable function $\psi_A(y)$ defined on (Y, \mathcal{S}) and such that $P_\theta(A/T(x)) = \psi_A(T(x))$ almost surely with respect to the measure P_θ for all $\theta \in \Theta$. Then $\lambda(A/T(x)) = \psi_A(T(x))$ almost surely with respect to the measure λ . Indeed, the definition of the conditional probability implies for all $A \in \mathcal{B}$ and $B \in \mathcal{B}_T = T^{-1}(\mathcal{S})$ that

$$\begin{aligned} \lambda(A \cap B) &= \sum_i c_i P_{\theta_i}(A \cap B) = \sum_i c_i \int_B P_{\theta_i}(A/T(x)) P_{\theta_i}(dx) \\ &= \sum_i c_i \int_B \psi_A(T(x)) P_{\theta_i}(dx) = \int_B \psi_A(T(x)) \lambda(dx) \end{aligned}$$

which together with (4.1.10)–(4.1.12) implies that $\lambda(A/T(x)) = \psi_A(T(x))$ almost surely with respect to the measure λ for all $A \in \mathcal{B}$. Further, for all $A \in \mathcal{B}$ we have

$$\begin{aligned} (4.1.17) \quad E_\lambda I_A(\xi) p(\xi; \theta) &= P_\theta(A) = E_\theta P_\theta(A/T(\xi)) \\ &= E_\theta \psi_A(T(\xi)) = E_\lambda p(\xi; \theta) \psi_A(T(\xi)) \\ &= E_\lambda p(\xi; \theta) \lambda(A/T(\xi)) = E_\lambda \lambda(A/T(\xi)) E_\lambda(p(\xi; \theta)/T(\xi)) \\ &= E_\lambda E_\lambda(I_A(\xi)/T(\xi)) E_\lambda(p(\xi; \theta)/T(\xi)) \\ &= E_\lambda E_\lambda(I_A(\xi) E_\lambda(p(\xi; \theta)/T(\xi))/T(\xi)) \\ &= E_\lambda I_A(\xi) E_\lambda(p(\xi; \theta)/T(\xi)) \end{aligned}$$

where E_λ and E_θ are the expectations with respect to the measures λ and P_θ , respectively. This implies that $p(x; \theta) = E_\lambda(p(\xi; \theta)/T(x))$ almost surely with respect to the measure λ , that is, the density $p(x; \theta)$ is \mathcal{B}_T -measurable or, in other words, $p(x; \theta) = g(T(x); \theta)$ almost surely with respect to the measure λ . Since $\lambda \ll \mu$, equality (4.1.17) implies that for all $A \in \mathcal{B}$

$$\begin{aligned} (4.1.18) \quad \int_A f(x; \theta) \mu(dx) &= P_\theta(A) = E_\lambda I_A(\xi) g(T(\xi); \theta) \\ &= \int_A g(T(x); \theta) r(x) \mu(dx) \end{aligned}$$

where $r(x) = d\lambda/d\mu(x)$. Since the set $A \in \mathcal{B}$ is arbitrary, equality (4.1.18) implies the required relation (4.1.16).

Sufficiency. Let relation (4.1.16) hold. Then

$$(4.1.19) \quad \begin{aligned} \frac{d\lambda}{d\mu}(x) &= \sum_i c_i f(x; \theta_i) = r(x) \sum_i c_i g(T(x); \theta_i) \\ &= r(x)G(T(x)) \quad (\mu\text{-a.s.}) \end{aligned}$$

where

$$G(T(x)) = \sum_i c_i g(T(x); \theta_i), \quad x \in X.$$

Consider the function

$$(4.1.20) \quad \tilde{p}(x; \theta) = \begin{cases} \frac{g(T(x); \theta)}{G(T(x))}, & \text{if } G(T(x)) > 0, \\ 0, & \text{if } G(T(x)) = 0. \end{cases}$$

It is clear that the function $\tilde{p}(x; \theta)$ is \mathcal{B}_T -measurable. Moreover

$$\frac{dP_\theta}{d\mu}(x) = \frac{dP_\theta}{d\lambda}(x) \cdot \frac{d\lambda}{d\mu}(x) \quad (\mu\text{-a.s.}),$$

since $P_\theta \ll \lambda \ll \mu$.

Therefore

$$(4.1.21) \quad \frac{dP_\theta}{d\lambda}(x) = \frac{dP_\theta}{d\mu}(x) / \frac{d\lambda}{d\mu}(x) = \frac{g(T(x); \theta)}{G(T(x))} \quad (\mu\text{-a.s.})$$

by (4.1.16) and (4.1.19). Equality (4.1.19) implies that $\lambda\{x: G(T(x)) = 0\} = 0$. Thus $p(x; \theta) = \tilde{p}(x; \theta)$ (λ -a.s.) by (4.1.20) and (4.1.21), since λ is absolutely continuous with respect to μ . The function $\tilde{p}(x; \theta)$ is \mathcal{B}_T -measurable, whence for all $A \in \mathcal{B}$ we have

$$\begin{aligned} E_\theta P_\theta(A/T(\xi)) &= P_\theta(A) = E_\lambda I_A(\xi) p(\xi; \theta) \\ &= E_\lambda E_\lambda(I_A(\xi) p(\xi; \theta) / T(\xi)) = E_\lambda p(\xi; \theta) E_\lambda(I_A(\xi) / T(\xi)) \\ &= E_\theta E_\lambda(I_A(\xi) / T(\xi)) = E_\theta \lambda(A/T(\xi)). \end{aligned}$$

If A is replaced with $A \cap B$ in this equality and if $A \in \mathcal{B}$ and $B \in \mathcal{B}_T$, then

$$E_\theta I_B(\xi) P_\theta(A/T(\xi)) = E_\theta I_B(\xi) \lambda(A/T(\xi)).$$

Since $B \in \mathcal{B}_T$ is arbitrary, it follows that for all $A \in \mathcal{B}$

$$P_\theta(A/T(x)) = \lambda(A/T(x)) \quad (P_\theta\text{-a.s.}) \quad \text{for all } \theta \in \Theta.$$

This means that the statistic $T(x)$ is sufficient. □

Below are two corollaries of Theorem 4.1.2.

COROLLARY 4.1.1. *Let a family \mathcal{P} be dominated. If a measurable function of some statistic is sufficient for the family \mathcal{P} , then the statistic itself is sufficient for this family, too.*

PROOF. Let T and \tilde{T} be two statistics such that $T = \phi(\tilde{T})$ where ϕ is a measurable function. Assume that T is a sufficient statistic for the family \mathcal{P} . According to Theorem 4.1.2 we have

$$f(x; \theta) = g(T(x); \theta)r(x) = \tilde{g}(\tilde{T}(x); \theta)r(x) \quad (\mu\text{-a.s.})$$

(see relation (4.1.16)), that is, the statistic \tilde{T} is also sufficient for the family \mathcal{P} . \square

COROLLARY 4.1.2. *Let a family \mathcal{P} be dominated. If T is a sufficient statistic for the family \mathcal{P} and a function ϕ is such that $v = \phi(y)$ is a measurable one-to-one mapping, then the statistic $\tilde{T} = \phi(T)$ is also sufficient for the family \mathcal{P} .*

PROOF. Since ϕ is a one-to-one mapping, we have $T = \phi^{-1}(\tilde{T})$. Now Corollary 4.1.2 follows from Corollary 4.1.1. \square

Applying the factorization criterion one can find sufficient statistics for dominated families. Below are some applications of the factorization criterion for obtaining sufficient statistics.

Examples of sufficient statistics. According to the factorization criterion the statistic $T(x) = x$, called the *trivial sufficient statistic*, is sufficient for every dominated family of probability measures \mathcal{P} . In the examples below we find nontrivial sufficient statistics.

EXAMPLE 4.1.2. Let $\Theta = \{\theta_1, \theta_2, \dots, \theta_s\}$. Then any family \mathcal{P} is finite, that is, $\mathcal{P} = \{P_{\theta_1}, P_{\theta_2}, \dots, P_{\theta_s}\}$. A dominating measure μ exists in this case. In particular, one can put $\mu = \sum_{i=1}^s c_i P_{\theta_i}$ for $c_i > 0$, $i = 1, 2, \dots, s$. Consider the statistic

$$T(x) = (T_1(x), T_2(x), \dots, T_{s-1}(x))$$

where $T_j(x) = \ln(f(x; \theta_{j+1})/f(x; \theta_1))$, $j = 1, 2, \dots, s-1$. Then

$$f(x; \theta_j) = \exp(T_{j-1}(x))r(x), \quad j = 2, 3, \dots, s,$$

where $r(x) = f(x; \theta_1)$. Thus

$$f(x; \theta) = g(T(x); \theta)r(x) \quad \text{for all } \theta \in \Theta$$

where

$$g(T(x); \theta) = \begin{cases} \exp(T_{j-1}(x)), & \text{if } \theta = \theta_j, \quad j = 2, 3, \dots, s, \\ 1, & \text{if } \theta = \theta_1. \end{cases}$$

Therefore $T(x) = (T_1(x), T_2(x), \dots, T_{s-1}(x))$ is a sufficient statistic.

EXAMPLE 4.1.3. Let $\xi^{(n)} = (\xi_1, \xi_2, \dots, \xi_n)$ be a sample from the normal distribution $\mathcal{N}(\theta_1, \theta_2^2)$. Then the density $f(x; \theta)$, $\theta = (\theta_1, \theta_2)$, is represented in the form

$$\begin{aligned} f(x; \theta) &= \frac{1}{(\sqrt{2\pi}\theta_2)^n} \exp \left\{ -\frac{1}{2\theta_2^2} \sum_{i=1}^n (x_i - \theta_1)^2 \right\} \\ &= \frac{1}{(\sqrt{2\pi}\theta_2)^n} \exp \left\{ -\frac{1}{2\theta_2^2} \sum_{i=1}^n x_i^2 + \frac{\theta_1}{\theta_2^2} \sum_{i=1}^n x_i - \frac{n\theta_1^2}{2\theta_2^2} \right\} \end{aligned}$$

where $x = (x_1, \dots, x_n)$. This means that the statistic $T(x) = (T_1(x), T_2(x))$, where

$$T_1(x) = \sum_{i=1}^n x_i, \quad T_2(x) = \sum_{i=1}^n x_i^2$$

is sufficient. Note that according to Corollary 4.1.1, the statistic $\tilde{T}(x) = (\tilde{T}_1(x), \tilde{T}_2(x))$, where

$$\tilde{T}_1(x) = \bar{x} = \frac{1}{n} \sum_{i=1}^n x_i, \quad \tilde{T}_2(x) = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$$

is also sufficient because $T(x) = \phi(\tilde{T}(x))$.

If $\xi^{(n)}$ is a sample from the distribution $\mathcal{N}(\theta_1, \sigma^2)$ and the variance σ^2 is known, then $T(x) = \sum_{i=1}^n x_i$ is a sufficient statistic, while if $\xi^{(n)}$ is a sample from the distribution $\mathcal{N}(\tilde{a}, \theta_2^2)$ and the expectation \tilde{a} is known, then $T(x) = \sum_{i=1}^n (x_i - \tilde{a})^2$ is a sufficient statistic. It is clear that the statistic $T(x) = (\sum_{i=1}^n x_i, \sum_{i=1}^n x_i^2)$ is also sufficient in both cases.

EXAMPLE 4.1.4. Let $\xi^{(n)} = (\xi_1, \dots, \xi_n)$ where ξ_1, \dots, ξ_n are independent random variables. Assume that ξ_i has the normal $\mathcal{N}(\theta_0, \theta_i^2)$ distribution, $i = 1, \dots, n$. Then the density $f^{(n)}$ is of the form

$$f(x; \theta) = \prod_{i=1}^n \frac{\exp\left\{-\frac{1}{2} \sum_{i=1}^n \left(\frac{x_i}{\theta_i}\right)^2 + \theta_0 \sum_{i=1}^n \frac{x_i}{\theta_i^2} - \frac{\theta_0^2}{2} \sum_{i=1}^n \frac{1}{\theta_i^2}\right\}}{\sqrt{2\pi}\theta_i}$$

where $\theta = (\theta_0, \dots, \theta_n)$ is an unknown vector parameter. It is clear that the "best" sufficient statistic in this case is $T(x) = x$. If the random variable ξ_i has the normal $\mathcal{N}(\theta, \sigma_i^2)$ distribution and the variances σ_i^2 are known, then $T(x) = \sum_{i=1}^n x_i/\sigma_i^2$ is a sufficient statistic.

EXAMPLE 4.1.5. Let $\xi^{(n)} = (\xi_1, \xi_2, \dots, \xi_n)$ be a sample from an uniform distribution on the interval $[0, \theta]$ where $\theta > 0$ is an unknown parameter. Then the density of the sample $f^{(n)}$ is of the form

$$f(x; \theta) = \theta^{-n} I_{[0, \infty)}(x_{n,1}) I_{(-\infty, \theta]}(x_{n,n})$$

where

$$x_{n,1} = \min_{1 \leq i \leq n} x_i, \quad x_{n,n} = \max_{1 \leq i \leq n} x_i, \quad (x_1, \dots, x_n) = x.$$

Thus relation (4.1.16) holds with $g(t; \theta) = \theta^{-n} I_{(-\infty, \theta]}(t)$, $r(x) = I_{[0, \infty)}(x_{n,1})$, and $T(x) = x_{n,n}$, whence it follows that $T(x) = x_{n,n}$ is a sufficient statistic.

If $\xi^{(n)}$ is a sample from the uniform distribution on the interval $[\theta, \theta + 1]$ and $\theta \in \mathbf{R}$ is an unknown parameter, then we can proceed in the same way as above to show that $T(x) = (x_{n,1}, x_{n,n})$ is a sufficient statistic of the parameter θ . The same statistic is sufficient for the two-dimensional parameter $\theta = (\theta_1, \theta_2)$, $-\infty < \theta_1 < \theta_2 < \infty$, in the case of a sample from the uniform distribution on the interval $[\theta_1, \theta_2]$.

EXAMPLE 4.1.6. Let $\xi^{(n)} = (\xi_1, \xi_2, \dots, \xi_n)$ be a sample from the Pearson type III distribution whose density is given by

$$\frac{\theta_2^{\theta_3}}{\Gamma(\theta_3)} (x - \theta_1)^{\theta_3 - 1} e^{-\theta_2(x - \theta_1)} I_{[\theta_1, \infty)}(x), \quad x \in (-\infty, \infty),$$

where $(\theta_1, \theta_2, \theta_3) = \theta$ is an unknown parameter. Then the distribution density of the sample $\xi^{(n)}$ is

$$f(x; \theta) = \frac{\theta_2^{n\theta_3}}{\Gamma^n(\theta_3)} \prod_{i=1}^n (x_i - \theta_1)^{\theta_3 - 1} \exp \left\{ -\theta_2 \sum_{i=1}^n (x_i - \theta_1) \right\} I_{[\theta_1, \infty)}(x_{n,1})$$

where $x = (x_1, \dots, x_n)$. If $\theta_3 = 1$ and an unknown parameter is $\theta = (\theta_1, \theta_2)$, then the density is of the form

$$f(x; \theta) = \theta_2^n \exp \left\{ -\theta_2 \sum_{i=1}^n (x_i - \theta_1) \right\} I_{[\theta_1, \infty)}(x_{n,1}).$$

Thus the statistic $T(x) = (\sum_{i=1}^n x_i, x_{n,1})$ is sufficient. If $\theta_1 = 0$ and $\theta = (\theta_2, \theta_3)$ is an unknown parameter, then $T(x) = (\sum_{i=1}^n x_i, \prod_{i=1}^n x_i)$ is a sufficient statistic. If $\theta_3 \neq 1$ is either a known or unknown parameter, while θ_1 is an unknown parameter, then the "simplest" sufficient statistic is $T(x) = (x_1, x_2, \dots, x_n)$.

Note that the density of the Pearson type III distribution belongs to the system of Pearson curves (its description can be found in [9], §19.4, or in [24], §5.6).

Sufficient statistics in the Bayes approach. Let ξ be an observation that is a random element assuming values in a measurable space (X, \mathcal{B}) . Let its distribution depend on an unknown parameter $\theta = (\theta_1, \theta_2, \dots, \theta_k)'$ which is random with the distribution \mathbf{Q} concentrated on a Borel set $\Theta \subset \mathbf{R}^k$, $k \geq 1$. Let P_t , $t \in \Theta$, be probability measures corresponding to the conditional distribution of the observation ξ given $\theta = t$, that is, $P_t(A) = P\{\xi \in A / \theta = t\}$ for all $A \in \mathcal{B}$ and $t \in \Theta$. Thus we are given a family of probability measures $\mathcal{P} = (P_t, t \in \Theta)$ which we assume to be dominated by some measure μ . We also assume that the measure \mathbf{Q} possesses the density $q(t)$ with respect to some measure λ .

The following result contains necessary and sufficient conditions that a statistic is sufficient for a family $\mathcal{P} = (P_t, t \in \Theta)$ expressed in terms of a posteriori density.

THEOREM 4.1.3. *A statistic $T = T(x)$ mapping (X, \mathcal{B}) into some measurable space (Y, \mathcal{S}) is sufficient for a family $\mathcal{P} = (P_t, t \in \Theta)$ if and only if for any a priori distribution \mathbf{Q} of the parameter θ , the a posteriori distribution \mathbf{Q}_x depends on x through $T(x)$. Here $\mathbf{Q}_x(A) = P\{\theta \in A / \xi = x\}$, $x \in X$, $A \in \mathcal{B}^k$.*

PROOF. Let T be a sufficient statistic for a family \mathcal{P} and let $q(t)$ be the density of the measure \mathbf{Q} with respect to the measure λ . The density $q(t/x)$ of the a posteriori measure \mathbf{Q}_x with respect to the measure λ exists for all $x \in X$. Applying relation (4.1.16) we obtain from the Bayes theorem that

$$q(t/x) = \frac{f(x; t)q(t)}{\int f(x; u)q(u) \lambda(du)} = \frac{g(T(x); t)q(t)}{\int g(T(x); u)q(u) \lambda(du)},$$

that is, $q(t/x)$ depends on x through $T(x)$.

Now we prove the converse. Consider an a priori distribution such that $q(t) > 0$ everywhere on Θ and for all t

$$f(x; t) = \frac{q(t/x)f(x)}{q(t)}, \quad f(x) = \int f(x; u)q(u) \lambda(du).$$

Put $q(t/x) = \psi(t, T(x))$. Setting $g(y; t) = \psi(t, y)/q(t)$ and $r(x) = f(x)$ we get relation (4.1.16), that is, $T(x)$ is a sufficient statistic. \square

COROLLARY 4.1.3. *If $T = T(x)$ is a sufficient statistic for a family*

$$\mathcal{P} = (\mathbf{P}_t, t \in \Theta),$$

then any Bayes estimator, as well as any minimax estimator, of the parameter θ with respect to the quadratic loss function defined as in Theorem 3.1.3 depend on the statistic T .

PROOF. Note that the Bayes estimator of the parameter θ with respect to the quadratic loss function is the a posteriori expectation

$$E(\theta/\xi = x) = \int tq(t/x) \lambda(dt).$$

Now we apply Theorem 4.1.3 to complete the proof. \square

REMARK 4.1.1. Using Theorem 4.1.3 one can provide the following equivalent definition, called the *Bayes definition of a sufficient statistic* (see [36]). A statistic $T = T(x)$ is sufficient for a family $\mathcal{P} = (\mathbf{P}_t, t \in \Theta)$ if for any a priori distribution \mathbf{Q} of the parameter θ , the a posteriori distribution \mathbf{Q}_x depends on x through $T(x)$ almost surely with respect to μ .

Fisher information and sufficient statistics. Let ξ be an observation that is a random element assuming values in a measurable space (X, \mathcal{B}) . Let its distribution belong to a family $\mathcal{P} = \{\mathbf{P}_\theta, \theta \in \Theta\}$ where θ is a nonrandom scalar unknown parameter. We assume that the family \mathcal{P} is dominated by some σ -finite measure μ whose density $f(x; \theta) = d\mathbf{P}_\theta/d\mu(x)$ satisfies an arbitrary regularity condition under which the Fisher information $I(\theta)$ is well defined. For the sake of definiteness we assume the regularity conditions (CR) (see Section 3.4). Then $I(\theta) = E_\theta S^2(\xi; \theta)$ where $S(x; \theta) = \partial \ln f(x; \theta)/\partial \theta$. Note that conditions (CR) imply that

$$E_\theta S(\xi; \theta) = 0$$

(see Lemma 3.4.1), whence it follows that $I(\theta) = D_\theta S(\xi; \theta)$ is the variance of the random variable $S(\xi; \theta)$.

Let $T = T(x)$ be some statistic mapping (X, \mathcal{B}) into some measurable space (Y, \mathcal{S}) . Denote by μ^T and \mathbf{P}_θ^T , $\theta \in \Theta$, the images of measures μ and \mathbf{P}_θ , $\theta \in \Theta$, respectively, under the mapping $T: (X, \mathcal{B}) \rightarrow (Y, \mathcal{S})$, that is,

$$\mathbf{P}_\theta^T(B) = \mathbf{P}_\theta(T^{-1}(B)) \quad \text{and} \quad \mu^T(B) = \mu(T^{-1}(B)) \quad \text{for all } B \in \mathcal{S}$$

where $T^{-1}(B) = \{x: T(x) \in B\}$ is the preimage of the set B under the mapping T . It is clear that the family of measures $\mathcal{P}^T = (\mathbf{P}_\theta^T, \theta \in \Theta)$ is dominated by the measure μ^T . Denote by $g(y; \theta) = d\mathbf{P}_\theta^T/d\mu^T(y)$, $y \in Y$, the density of the measure \mathbf{P}_θ^T with respect to the measure μ^T . If the density $g(y; \theta)$ also satisfies conditions (CR), then the Fisher information $I^T(\theta) = E_\theta(\partial \ln g(T(\xi); \theta)/\partial \theta)^2$ is well defined.

The following result provides a relation between the Fisher information $I^T(\theta)$ generated by the distribution of the statistic $T(\xi)$ and the Fisher information $I(\theta)$ generated by the distribution of the observation ξ .

THEOREM 4.1.4. *Let a family $\mathcal{P} = (P_\theta, \theta \in \Theta)$ be dominated by a measure μ and satisfy the regularity conditions (CR). Let $T = T(x)$ be a statistic mapping (X, \mathcal{B}) into a measurable space (Y, \mathcal{S}) . Assume that the family of probability measures $\mathcal{P}^T = (P_\theta^T, \theta \in \Theta)$ is generated by T on (Y, \mathcal{S}) and satisfies the regularity conditions (CR). Then*

$$(4.1.22) \quad I^T(\theta) \leq I(\theta) \quad \text{for all } \theta \in \Theta.$$

Moreover inequality (4.1.22) becomes an equality if and only if the statistic T is sufficient for the family \mathcal{P} .

PROOF. Let C be an arbitrary set of \mathcal{S} . According to the definition of the conditional expectation and in view of condition (iii) in Section 3.4 we obtain from (CR) that for all $\theta \in \Theta$

$$(4.1.23) \quad \int_{T^{-1}(C)} S(x; \theta) P_\theta(dx) = \int_C E_\theta\{S(\xi; \theta)/T(\xi) = y\} P_\theta^T(dy).$$

Since conditions (CR) holds for both families \mathcal{P} and \mathcal{P}^T , we have for any P_θ^T -nonzero set $C \in \mathcal{S}$ that

$$(4.1.24) \quad \begin{aligned} \int_{T^{-1}(C)} S(x; \theta) P_\theta(dx) &= \int_{T^{-1}(C)} \frac{\partial}{\partial \theta} f(x; \theta) \mu(dx) \\ &= \frac{\partial}{\partial \theta} \int_{T^{-1}(C)} f(x; \theta) \mu(dx) = \frac{\partial}{\partial \theta} \int_C g(y; \theta) \mu^T(dy) \\ &= \int_C \frac{\partial}{\partial \theta} g(y; \theta) \mu^T(dy) = \int_C \frac{\partial}{\partial \theta} \ln g(y; \theta) P_\theta^T(dy). \end{aligned}$$

Relations (4.1.23) and (4.1.24) yield

$$(4.1.25) \quad \frac{\partial}{\partial \theta} \ln g(y; \theta) = E_\theta\{S(\xi; \theta)/T(\xi) = y\} \quad (P_\theta^T\text{-a.s.})$$

for all $\theta \in \Theta$. By definition we have $I^T(\theta) = E_\theta(\partial \ln g(T(\xi); \theta)/\partial \theta)^2$. To prove the inequality $I^T(\theta) \leq I(\theta)$, note that

$$(4.1.26) \quad \begin{aligned} 0 &\leq E_\theta \left(\frac{\partial}{\partial \theta} \ln f(\xi; \theta) - \frac{\partial}{\partial \theta} \ln g(T(\xi); \theta) \right)^2 \\ &= E_\theta \left(\frac{\partial}{\partial \theta} \ln f(\xi; \theta) \right)^2 + E_\theta \left(\frac{\partial}{\partial \theta} \ln g(T(\xi); \theta) \right)^2 \\ &\quad - 2E_\theta \frac{\partial}{\partial \theta} \ln f(\xi; \theta) \frac{\partial}{\partial \theta} \ln g(T(\xi); \theta). \end{aligned}$$

It follows from (4.1.25) that

$$\begin{aligned}
 (4.1.27) \quad & E_{\theta} \frac{\partial}{\partial \theta} \ln f(\xi; \theta) \frac{\partial}{\partial \theta} \ln g(T(\xi); \theta) \\
 &= E_{\theta} \frac{\partial}{\partial \theta} \ln g(T(\xi); \theta) E_{\theta} \left(\frac{\partial}{\partial \theta} \ln f(\xi; \theta) / T(\xi) \right) \\
 &= E_{\theta} \left(\frac{\partial}{\partial \theta} \ln g(T(\xi); \theta) \right)^2.
 \end{aligned}$$

Substituting (4.1.27) into (4.1.26) we prove inequality (4.1.22). It remains to show that the inequality (4.1.22) becomes an equality if and only if the statistic T is sufficient for the family \mathcal{P} .

If T is a sufficient statistic for \mathcal{P} , then by the factorization criterion

$$(4.1.28) \quad f(x; \theta) = g^*(T(x); \theta) h(x) \quad (\mu\text{-a.s.})$$

where $h(x) \geq 0$, $g^*(y; \theta) \geq 0$, the function $h(x)$ is \mathcal{B} -measurable, and $g^*(y; \theta)$ is \mathcal{S} -measurable for all $\theta \in \Theta$. Thus for all $\theta \in \Theta$

$$\frac{\partial}{\partial \theta} \ln f(x; \theta) = \frac{\partial}{\partial \theta} \ln g^*(T(x); \theta) \quad (\mathbb{P}_{\theta}\text{-a.s.}).$$

Therefore

$$(4.1.29) \quad I(\theta) = E_{\theta} \left(\frac{\partial}{\partial \theta} \ln g^*(T(\xi); \theta) \right)^2 \quad \text{for all } \theta \in \Theta.$$

On the space (Y, \mathcal{S}) consider the measure

$$\lambda(B) = \int_{T^{-1}(B)} h(x) \mu(dx), \quad B \in \mathcal{S}.$$

According to (4.1.28) we have for all $B \in \mathcal{S}$

$$\begin{aligned}
 \mathbb{P}_{\theta}^T(B) &= \int_{T^{-1}(B)} f(x; \theta) \mu(dx) \\
 &= \int_{T^{-1}(B)} g^*(T(x); \theta) h(x) \mu(dx) = \int_B g^*(y; \theta) \lambda(dy).
 \end{aligned}$$

Thus the measure \mathbb{P}_{θ}^T is absolutely continuous with respect to the measure λ and the density is $g^*(y; \theta)$. Therefore (4.1.29) implies that $I(\theta) = I^T(\theta)$ for all $\theta \in \Theta$.

Now we show that if $I(\theta) = I^T(\theta)$ for all $\theta \in \Theta$, then T is a sufficient statistic for \mathcal{P} . Indeed, $I(\theta)$ is the variance of $S(\xi; \theta)$, hence

$$(4.1.30) \quad I(\theta) = D_{\theta} S(\xi; \theta) = E_{\theta} D_{\theta} (S(\xi; \theta) / T(\xi)) + D_{\theta} E(S(\xi; \theta) / T(\xi))$$

where $E_{\theta}(S(\xi; \theta) / y) = E_{\theta}\{S(\xi; \theta) / T(\xi) = y\}$ and

$$D_{\theta}(S(\xi; \theta) / y) = D_{\theta}\{S(\xi; \theta) / T(\xi) = y\}$$

are the conditional expectation and variance, respectively. Equality (4.1.25) implies

$$E_{\theta}(S(\xi; \theta) / y) = \frac{\partial}{\partial \theta} \ln g(y; \theta) \quad (\mathbb{P}_{\theta}^T\text{-a.s.}).$$

The family \mathcal{P}^T satisfies conditions (CR), thus

$$(4.1.31) \quad I^T(\theta) = D_{\theta} \ln g(T(\xi); \theta) = D_{\theta} E_{\theta}(S(\xi; \theta) / T(\xi)), \quad \theta \in \Theta.$$

Since $I(\theta) = I^T(\theta)$ for all $\theta \in \Theta$, relations (4.1.30) and (4.1.31) yield

$$E_{\theta} D_{\theta}(S(\xi; \theta)/T(\xi)) = 0 \quad \text{for all } \theta \in \Theta.$$

This implies that

$$D_{\theta}(S(\xi; \theta)/T(x)) = 0 \quad (\mathbb{P}_{\theta}\text{-a.s.}).$$

In other words, the function $S(x; \theta)$ is $T^{-1}(\mathcal{S})$ -measurable. Therefore there exists a measurable function $k(y; \theta)$ on (Y, \mathcal{S}) such that $S(x; \theta) = k(T(x); \theta)$ for all $\theta \in \Theta$, whence

$$\ln f(x; \theta) = \int_{\theta_0}^{\theta} k(T(x); t) dt + a(x) \quad \text{for all } \theta \in \Theta$$

and some $\theta_0 \in \Theta$. Finally,

$$f(x; \theta) = g^*(T(x); \theta)h(x) \quad \text{for all } \theta \in \Theta.$$

Now the factorization criterion implies that $T(x)$ is a sufficient statistic for \mathcal{P} . \square

It follows from inequality (4.1.22) that a sufficient statistic is the only statistic that compresses the sampling data without loss of information. The precise statement concerning the Fisher information is given by Theorem 4.1.4. Similar results can be given for other measures of the amount of information about a parameter contained in an observation, say for the Shannon information or Kullback information (see [22] and [20]).

REMARK 4.1.2. Theorem 4.1.4 remains true for the Fisher information matrix in the case of a multidimensional parameter θ under the multidimensional analogue of the regularity conditions (CR) (see [36], §4.3). Theorem 4.1.4 can be proved under other sets of regularity conditions, say $(CR)^*$, (R) , or $(R)^*$, in both one-dimensional and multidimensional cases.

EXAMPLE 4.1.7. Let $\xi^{(n)} = (\xi_1, \xi_2, \dots, \xi_n)$ be a sample from the Bernoulli distribution with parameter $\theta \in \Theta = (0, 1)$. The density of the distribution of ξ_i with respect to the counting measure is of the form

$$f(x; \theta) = \theta^x(1 - \theta)^{1-x} = \mathbb{P}_{\theta}\{\xi_i = x\}$$

where x is either 0 or 1. Thus the Fisher information contained in a single observation ξ_i is

$$I(\theta) = \frac{1}{\theta(1 - \theta)}, \quad \theta \in \Theta,$$

while the Fisher information contained in the whole sample $\xi^{(n)}$ is

$$I_n(\theta) = nI(\theta) = \frac{n}{\theta(1 - \theta)}, \quad \theta \in \Theta.$$

Let ν_n be the number of "successes" in the sample $\xi^{(n)}$, that is, $\nu_n = \sum_{i=1}^n \xi_i$ is a statistic assuming values in the set $Y = \{0, 1, 2, \dots, n\}$. The density of the distribution of the statistic ν_n with respect to the counting measure is given by

$$g(y; \theta) = \binom{n}{y} \theta^y (1 - \theta)^{n-y} = \mathbb{P}_{\theta}\{\nu_n = y\}, \quad y \in Y.$$

The Fisher information contained in the statistic ν_n is

$$\begin{aligned} I^{\nu_n}(\theta) &= E_{\theta} \left(\frac{\partial}{\partial \theta} \ln g(\nu_n; \theta) \right)^2 = \sum_{y=0}^n \binom{n}{y} \theta^y (1-\theta)^{n-y} \left(\frac{y}{\theta} - \frac{n-y}{1-\theta} \right)^2 \\ &= \sum_{y=0}^n \binom{n}{y} \theta^y (1-\theta)^{n-y} \frac{(y-n\theta)^2}{(\theta(1-\theta))^2} = \frac{D\nu_n}{(\theta(1-\theta))^2} = \frac{n}{\theta(1-\theta)}. \end{aligned}$$

Thus $I^{\nu_n}(\theta) = I_n(\theta)$ for all $\theta \in \Theta$. By Theorem 4.1.4 the statistic ν_n is sufficient for the parameter θ . This conclusion can also be made from the factorization criterion.

4.2. Sufficient statistics and optimal estimators

Rao–Blackwell–Kolmogorov theorem. Let ξ be an observation assuming values in a measurable space (X, \mathcal{B}) and whose distribution belongs to a family

$$\mathcal{P} = (P_{\theta}, \theta \in \Theta)$$

where $\theta = (\theta_1, \dots, \theta_k)'$ is an unknown parameter of a set $\Theta \subset \mathbf{R}^k$, $k \geq 1$. First we consider the case where θ is a scalar parameter, that is, the case $k = 1$. Let K_b be the class of estimators $\hat{\theta}$ of the parameter θ with a bias $b(\theta)$ (see Section 3.4). In other words, $a(\theta) = E_{\theta} \hat{\theta} = \theta + b(\theta)$ for all $\theta \in \Theta$.

The following result, known as the *Rao–Blackwell–Kolmogorov theorem*, highlights the role of sufficient statistics in the theory of estimation.

THEOREM 4.2.1. *Let $T = T(x)$ be a sufficient statistic for a family*

$$\mathcal{P} = (P_{\theta}, \theta \in \Theta)$$

and let $\hat{\theta} \in K_b$. Then the function $\hat{\theta}_T = E_{\theta}(\hat{\theta}/T)$ is an estimator such that

- 1) $\hat{\theta}_T \in K_b$;
- 2) the estimator $\hat{\theta}_T$ depends on x through $T(x)$;
- 3) $E_{\theta}(\hat{\theta}_T - \theta)^2 \leq E_{\theta}(\hat{\theta} - \theta)^2$ for all $\theta \in \Theta$, and moreover the inequality becomes an equality if and only if $\hat{\theta} = \hat{\theta}_T$ almost surely with respect to the measure P_{θ} .

PROOF. Let T be a sufficient statistic. The conditional probability $P_{\theta}(A/T)$ depends on A and T and does not depend on θ . Moreover $P_{\theta}(A/T)$ is a measurable function of T . Thus the conditional expectation $\hat{\theta}_T = E_{\theta}(\hat{\theta}/T)$ depends on T and does not depend on θ . Therefore the estimator $\hat{\theta}_T$ satisfies condition 2) of the theorem.

Properties of the conditional expectation imply that

$$E_{\theta} \hat{\theta}_T = E_{\theta} E_{\theta}(\hat{\theta}/T) = E_{\theta} \hat{\theta},$$

that is, $\hat{\theta}_T \in K_b$, whence condition 1) of the theorem follows.

The inequality in statement 3) of the theorem is obvious if $E_{\theta}(\hat{\theta} - \theta)^2 = \infty$. Thus we consider the case when $E_{\theta}(\hat{\theta} - \theta)^2 < \infty$. We have

$$\begin{aligned} (4.2.1) \quad E_{\theta}(\hat{\theta} - \theta)^2 &= E_{\theta}(\hat{\theta} - \hat{\theta}_T + \hat{\theta}_T - \theta)^2 \\ &= E_{\theta}(\hat{\theta} - \hat{\theta}_T)^2 + E_{\theta}(\hat{\theta}_T - \theta)^2 + 2E_{\theta}(\hat{\theta} - \hat{\theta}_T)(\hat{\theta}_T - \theta). \end{aligned}$$

By the properties of the conditional expectation

$$(4.2.2) \quad \begin{aligned} E_{\theta}(\widehat{\theta} - \widehat{\theta}_T)(\widehat{\theta}_T - \theta) &= E_{\theta}E_{\theta}((\widehat{\theta} - \widehat{\theta}_T)(\widehat{\theta}_T - \theta)/T) \\ &= E_{\theta}(\widehat{\theta}_T - \theta)E_{\theta}(\widehat{\theta} - \widehat{\theta}_T/T) = 0, \end{aligned}$$

since $E_{\theta}(\widehat{\theta} - \widehat{\theta}_T/T) = E_{\theta}(\widehat{\theta}/T) - \widehat{\theta}_T = 0$.

Equalities (4.2.1) and (4.2.2) imply

$$E_{\theta}(\widehat{\theta} - \theta)^2 = E_{\theta}(\widehat{\theta} - \widehat{\theta}_T)^2 + E_{\theta}(\widehat{\theta}_T - \theta)^2,$$

whence statement 3) of the theorem follows. \square

Theorem 4.2.1 shows that if T is a sufficient statistic, then one can improve the estimator $\widehat{\theta} \in K_b$ uniformly in $\theta \in \Theta$ by applying the operator $E_{\theta}(\cdot/T)$ to the statistic $\widehat{\theta}$.

There is another interpretation of Theorem 4.2.1. Namely let S and T be two sufficient statistics for a family \mathcal{P} . If $\widehat{\theta} = \phi(T)$ where ϕ is a measurable function and S is a measurable function of T , then

$$E_{\theta}(\widehat{\theta}_S - \theta)^2 \leq E_{\theta}(\widehat{\theta} - \theta)^2$$

where $\widehat{\theta}_S = E_{\theta}(\widehat{\theta}/S)$. In other words, one should find the so-called *minimal* sufficient statistics, that is those statistics for which any other sufficient statistic is a function of it. The procedure of the construction of an optimal estimator is as follows. One starts with a "bad" estimator $\widehat{\theta}$ and improves it with the help of sufficient statistics until the estimator becomes optimal.

Theorem 4.2.1 holds in the multidimensional case as well. In this case θ and $\widehat{\theta}$ are vectors of the space \mathbf{R}^k , $k > 1$. As in the one-dimensional case let K_b be the class of estimators $\widehat{\theta}$ of the parameter θ with the bias $b(\theta)$.

THEOREM 4.2.2. *Let T be a sufficient statistic for a family $\mathcal{P} = (P_{\theta}, \theta \in \Theta)$ and let $\widehat{\theta} \in K_b$. Then the estimator $\widehat{\theta}_T = E_{\theta}(\widehat{\theta}/T)$ is such that*

- 1) $\widehat{\theta}_T \in K_b$;
- 2) $\widehat{\theta}_T$ depends on x through $T(x)$;
- 3) for any vector $a \in \mathbf{R}^k$

$$(4.2.3) \quad E_{\theta}(a'(\widehat{\theta}_T - \theta))^2 \leq E_{\theta}(a'(\widehat{\theta} - \theta))^2.$$

This inequality becomes an equality if and only if $\widehat{\theta} = \widehat{\theta}_S$ almost surely with respect to the measure P_{θ} .

PROOF. The first two statements of the theorem are obvious. Inequality (4.2.3) follows from Theorem 4.2.1, since the proof is reduced to the one-dimensional estimators $a'\widehat{\theta}$ of the parameter $a'\theta$ and since $E_{\theta}(a'\widehat{\theta}/S) = a'\widehat{\theta}_S$. If (4.2.3) becomes an equality for all $a \in \mathbf{R}^k$, then $a'\widehat{\theta} = a'\widehat{\theta}_S$ almost surely with respect to the measure P_{θ} for all a , whence $\widehat{\theta} = \widehat{\theta}_S$ almost surely with respect to P_{θ} . \square

REMARK 4.2.1. All the vectors a , θ , $\widehat{\theta}$, and $\widehat{\theta}_T$ in Theorem 4.2.2 are column-vectors.

Sufficient statistics and efficient estimators. We proved in Theorems 4.2.1 and 4.2.2 that if an estimator is not a function of a sufficient statistic, then it can be improved by using this sufficient statistic. However we still have no tool to construct an optimal estimator by following this idea.

On the other hand, if a set of regularity conditions holds, say (CR) or (R) , and the Cramér–Rao inequality becomes an equality, then the estimator is optimal (this kind of optimality is called efficiency in Sections 3.4 and 3.5).

Below we consider the case where conditions (R) are satisfied. All other cases can be studied in a similar way. Let $\hat{\theta} = (\hat{\theta}_1, \hat{\theta}_2, \dots, \hat{\theta}_k)'$ be an estimator of the parameter $\theta = (\theta_1, \theta_2, \dots, \theta_k)'$ constructed from an observation ξ . Let K_b be the class of estimators $\hat{\theta}$ of the parameter θ with a bias $b(\theta)$, that is,

$$a(\theta) = E_{\theta}\hat{\theta} = \theta + b(\theta).$$

The following result contains a relationship between efficient estimators and sufficient statistics.

THEOREM 4.2.3. *Let conditions (R) hold. Let $\hat{\theta}$ be an estimator of the class K_b such that $\det D(\hat{\theta}; \theta) > 0$ where $D(\hat{\theta}; \theta) = E_{\theta}(\hat{\theta} - a(\theta))(\hat{\theta} - a(\theta))'$ is the covariation matrix of the estimator $\hat{\theta}$. Then $\hat{\theta}$ is an efficient estimator in the class K_b if and only if $\hat{\theta}$ is a sufficient statistic of the parameter θ ,*

$$(4.2.4) \quad f(x; \theta) = g(\hat{\theta}(x); \theta)r(x) \quad \text{for all } \theta \in \Theta,$$

and all $r(x) = h(x)$ and $g(\hat{\theta}; \theta) = \exp(A(\theta)'\hat{\theta} + C(\theta))$ and $h(x)$, $A(\theta)$, and $C(\theta)$ are the functions occurring in representation (3.5.14).

PROOF. By definition, $\hat{\theta}$ is an efficient estimator of the parameter θ if and only if the Cramér–Rao inequality (3.5.13) becomes an equality. According to Theorem 3.5.2, the Cramér–Rao inequality (3.5.13) becomes an equality if and only if relation (3.5.14) holds. Note that relation (3.5.14) coincides with (4.2.4). By the factorization criterion (Theorem 4.1.2), the estimator $\hat{\theta}$ is a sufficient statistic for the parameter θ . \square

Note that if $\hat{\theta}$ is an efficient estimator of the parameter θ in the class K_b and regularity conditions hold, then $\hat{\theta}$ is also an optimal estimator of the parameter θ in the class K_b in the sense of the definition of Section 3.1. Generally speaking, the converse is not true, namely an estimator can be optimal in a class K_b , but the lower bound in the Cramér–Rao inequality is not attained for it. Thus an important role is played by those sufficient statistics that, by the Rao–Blackwell–Kolmogorov theorem, allow one to improve estimators and, in the case where one can construct the minimal sufficient statistics, to construct the optimal estimator.

Minimal sufficient statistics. We have seen above that there exist many sufficient statistics in the general situation. In particular, there always exists the so-called trivial sufficient statistic, namely $T(x) = x$. Nevertheless we are interested in those statistics that provide the best reduction of the data. However it is not always possible to find a sufficient statistic for which the reduction of the data is essentially better than that for the trivial sufficient statistic. To make the notion of the reduction of the data precise we introduce a partial order on the set of all sufficient statistics.

We say that a statistic T_1 is *subordinated* to a statistic T_2 if T_1 is a measurable function of T_2 , that is, $T_1 = \phi(T_2)$. If a statistic T_1 is subordinated to a statistic T_2 and T_2 is subordinated to T_1 , then the statistics T_1 and T_2 are called *equivalent*. A sufficient statistic T_0 is called *minimal* if it is subordinated to any other sufficient statistic.

The reduction of the data is best for a minimal sufficient statistic. If T is a minimal sufficient statistic, then a further reduction of the information as compared to T gives no result if the statistic remains sufficient.

We have seen above that the definition of sufficient statistics can be given in a more general form in terms of σ -algebras. If $\mathcal{B}' \subset \mathcal{B}$ is a σ -algebra, then \mathcal{B}' is called a *sufficient σ -algebra* for the family $\mathcal{P} = (P_\theta, \theta \in \Theta)$ (or, for the parameter θ) if there is a version of the conditional probability measure $P_\theta(A/\mathcal{B}')$, $A \in \mathcal{B}$, that does not depend on θ . Let T be a statistic mapping a measurable space (X, \mathcal{B}) into a measurable space (Y, \mathcal{S}) , and let $\mathcal{B}_T = T^{-1}(\mathcal{S})$ be the preimage of the σ -algebra \mathcal{S} under the mapping T . If the σ -algebra \mathcal{B}_T generated by the statistic T is sufficient, then T is a sufficient statistic. All the results on sufficient statistics can be stated in terms of sufficient σ -algebras. In particular, the factorization criterion remains true if the function $g(T(x); \theta)$ is substituted for a \mathcal{B}' -measurable function $g(x; \theta)$; in this case, \mathcal{B}' is a sufficient σ -algebra.

Let T_1 and T_2 be two statistics. It is clear that T_1 is subordinated to T_2 if

$$\mathcal{B}_{T_1} \subset \mathcal{B}_{T_2}.$$

Thus the statistic T_1 reduces the data in a better way than the statistic T_2 . Two statistics T_1 and T_2 are equivalent if and only if $\mathcal{B}_{T_1} = \mathcal{B}_{T_2}$.

A σ -algebra \mathcal{B}^* is called the *minimal sufficient σ -algebra* if it belongs to any other sufficient σ -algebra, that is, $\mathcal{B}^* \subset \mathcal{B}'$ for any sufficient σ -algebra \mathcal{B}' . In other words, T_0 is a minimal sufficient statistic if $\mathcal{B}_{T_0} \subset \mathcal{B}_T$ for every sufficient statistic T .

A minimal sufficient statistic always exists for dominated families

$$\mathcal{P} = (P_\theta, \theta \in \Theta)$$

(see Theorem 4.2.4). To prove this result we use Theorem 4.1.1: for a family $\mathcal{P} = (P_\theta, \theta \in \Theta)$ dominated by a σ -finite measure μ , there exists a discrete distribution \mathbf{Q} on Θ such that the family \mathcal{P} is dominated by the probability measure $P_{\mathbf{Q}} = \int P_t \mathbf{Q}(dt)$. Then the density $p(x; \theta)$ of the measure P_θ with respect to the measure $P_{\mathbf{Q}}$ can be expressed as

$$(4.2.5) \quad \frac{dP_\theta}{dP_{\mathbf{Q}}}(x) = p(x; \theta) = \frac{f(x; \theta)}{f(x; \mathbf{Q})} \quad (\mu\text{-a.s.})$$

where $f(x; \theta) = dP_\theta/d\mu(x)$ and $f(x; \mathbf{Q}) = dP_{\mathbf{Q}}/d\mu(x)$. If T is a sufficient statistic, then $p(x; \theta)$ depends on x through $T(x)$ by the factorization criterion.

THEOREM 4.2.4. *Let a family $\mathcal{P} = (P_\theta, \theta \in \Theta)$ be dominated by some σ -finite measure μ and let $\mathcal{B}^* = \sigma\{p(x; \theta); \theta \in \Theta\}$ be the σ -algebra in (X, \mathcal{B}) generated by the functions $p(x; \theta)$, $\theta \in \Theta$. Then \mathcal{B}^* is the minimal sufficient σ -algebra.*

PROOF. According to (4.2.5) we have

$$f(x; \theta) = p(x; \theta)f(x; \mathbf{Q}) \quad (\mu\text{-a.s.})$$

where the function $f(x; \mathbf{Q})$ does not depend on θ and the function $p(x; \theta)$ is \mathcal{B}^* -measurable for all $\theta \in \Theta$. From the factorization criterion for sufficient σ -algebras we obtain that \mathcal{B}^* is a sufficient σ -algebra.

Now let \mathcal{B}' be an arbitrary σ -algebra. Then $f(x; \theta) = g(x; \theta)r(x)$ (μ -a.s.) where $h(x)$ is a nonnegative \mathcal{B} -measurable function and $g(x; \theta)$ is a nonnegative \mathcal{B}' -measurable function. Consider the σ -algebra $\mathcal{B}_g = \sigma\{g(x; \theta); \theta \in \Theta\} \subset \mathcal{B}'$. It follows from (4.2.5) that

$$p(x; \theta) = \frac{g(x; \theta)}{\int g(x; t) \mathbf{Q}(dt)},$$

whence $\mathcal{B}^* \subset \mathcal{B}_g \subset \mathcal{B}'$. □

REMARK 4.2.2. Using Theorem 4.1.3 one can construct a minimal sufficient σ -algebra from the a posteriori distribution \mathbf{Q}_x . Let \mathbf{Q} be an a priori measure such that its density $q(t)$ with respect to some other measure λ is positive for all $t \in \Theta$. Then a posteriori density is given by

$$q(t/x) = \frac{f(x; t)q(t)}{f(x; \mathbf{Q})} = p(x; t)q(t).$$

Thus the minimal sufficient σ -algebra \mathcal{B}^* can be viewed as one generated by a posteriori distribution, that is, $\mathcal{B}^* = \sigma\{q(t/x); t \in \Theta\}$.

EXAMPLE 4.2.1. Let $\xi^{(n)} = (\xi_1, \xi_2, \dots, \xi_n)$ be a sample from the Poisson distribution with parameter $\theta \in \Theta = (0, \infty)$. We learned in Example 4.1.1 that

$$T_n = T_n(x) = \sum_{i=1}^n x_i$$

is a sufficient statistic. Here $x = (x_1, \dots, x_n)$ and $x_i \in \{0, 1, 2, \dots\}$ for all $i = 1, 2, \dots, n$. Then T_n is the minimal sufficient statistic by Theorem 4.2.4, since the σ -algebra \mathcal{B}_{T_n} coincides with the σ -algebra generated by the functions

$$p(x; \theta) = \frac{f(x; \theta)}{f(x; \theta_0)} = \left(\frac{\theta}{\theta_0}\right)^{T_n(x)} e^{n(\theta_0 - \theta)}$$

where $\theta_0 \in (0, \infty)$ is some fixed value of the parameter. As the distribution \mathbf{Q} on Θ we consider the distribution concentrated at the point θ_0 .

EXAMPLE 4.2.2. Let $\xi^{(n)} = (\xi_1, \xi_2, \dots, \xi_n)$ be a sample from the uniform distribution on the interval $[0, \theta]$ where $\theta > 0$ is an unknown parameter. Let the sampling space be $\mathbf{R}_+^n = \{x = (x_1, x_2, \dots, x_n): x_i > 0 \text{ for all } i = 1, 2, \dots, n\}$. As in Example 4.1.5 we prove that $T_n(x) = x_{n,n} = \max_{1 \leq i \leq n} x_i$ is a sufficient statistic. Moreover it is the minimal sufficient statistic. Indeed, let \mathbf{Q} be some distribution on $(0, \infty)$ whose density $q(t)$ is positive for all $t \in (0, \infty)$. Then

$$f(x; \theta) = \theta^{-n} I_{(-\infty, \theta]}(x_{n,n}), \quad x \in \mathbf{R}_+^n,$$

$$f(x; \mathbf{Q}) = \int_0^\infty f(x; t)q(t) dt = \int_{T_n(x)}^\infty t^{-n} q(t) dt > 0 \quad \text{for all } x \in \mathbf{R}_+^n.$$

It is also clear that $T_n(x) = \sup\{\theta: f(x; \theta)/f(x; \mathbf{Q}) = 0\}$, $x \in \mathbf{R}_+^n$. This means that the statistic T_n is measurable with respect to the minimal sufficient σ -algebra

$\mathcal{B}^* = \sigma\{f(x; \theta)/f(x; \mathbf{Q}); \theta \in \Theta\}$, $\mathcal{B}_{T_n} \subset \mathcal{B}^*$. Therefore T_n is the minimal sufficient statistic.

REMARK 4.2.3. There is another method to construct minimal sufficient statistics based on partitions of the sampling space generated by sufficient statistics. This method can be found in [36], also see [5] and [19].

Complete statistics and optimal estimators. In this section we consider complete sufficient statistics and use them to construct optimal estimators of a parameter. Let $T = T(x)$ be a statistic mapping a measurable space (X, \mathcal{B}) into a measurable space (Y, \mathcal{S}) . Assume that the dimension of the space Y is l , that is, $Y \subset \mathbf{R}^l$. A usual assumption is that $l > k$ where k is the dimension of the parameter θ .

Let $\Gamma = \{\mathbf{G}_\theta; \theta \in \Theta\}$ be some family of probability measures on $(\mathbf{R}^l, \mathcal{B}^l)$. A family Γ is called *complete* if the relation

$$(4.2.6) \quad \int \phi(x) \mathbf{G}_\theta(dx) = 0 \quad \text{for all } \theta \in \Theta$$

implies $\phi(x) = 0$ (\mathbf{G}_θ -a.s.) for all $\theta \in \Theta$. Equation (4.2.6) is considered in the class of functions $\phi: \mathbf{R}^l \rightarrow \mathbf{R}^k$ for which the integral (4.2.6) exists.

Let $\mathcal{P}^T = \{\mathbf{P}_\theta^T; \theta \in \Theta\}$ be a family of probability measures on (Y, \mathcal{S}) generated by the mapping $T: X \rightarrow Y$ where $\mathbf{P}_\theta^T(B) = \mathbf{P}_\theta(T^{-1}(B))$ and $B \in \mathcal{S}$, $\theta \in \Theta$. A statistic T is called *complete* if the family of distributions \mathcal{P}^T is complete. Equation (4.2.6) for the statistic T can be written in the following form:

$$(4.2.7) \quad \mathbf{E}_\theta \phi(T(\xi)) = 0 \quad \text{for all } \theta \in \Theta.$$

THEOREM 4.2.5. *A statistic T is complete if and only if for some $b_0(\theta)$ a \mathcal{B}_T -measurable estimator $\hat{\theta}$ is unique in the class of all \mathcal{B}_T -measurable estimators of the class K_{b_0} where $b_0(\theta)$ is the bias of the estimator.*

If a \mathcal{B}_T -measurable estimator is unique in the class K_{b_0} , then any \mathcal{B}_T -measurable estimator is unique in any other class K_b of estimators with the bias $b(\theta)$.

PROOF. Let $\hat{\theta}_1 = \phi_1(T)$ and $\hat{\theta}_2 = \phi_2(T)$ be two \mathcal{B}_T -measurable estimators of K_{b_0} . Then $\mathbf{E}_\theta \phi_i(T(\xi)) = b_0(\theta)$, $i = 1, 2$, and $\mathbf{E}_\theta(\phi_1(T(\xi)) - \phi_2(T(\xi))) = 0$ for all $\theta \in \Theta$. Since T is a complete statistic, $\phi_1(y) = \phi_2(y)$ (\mathbf{P}_θ^T -a.s.) for all $\theta \in \Theta$. Conversely, let $\mathbf{E}_\theta \phi(T(\xi)) = 0$ for all $\theta \in \Theta$ and $\hat{\theta}_1 = \phi_1(T) \in K_{b_0}$. Then $\hat{\theta}_2 = \phi_1(T) + \phi(T) \in K_{b_0}$. Since a \mathcal{B}_T -measurable estimator of K_{b_0} is unique, $\phi(T(x)) = 0$ (\mathbf{P}_θ -a.s.) for all $\theta \in \Theta$, that is, $\phi(y) = 0$ (\mathbf{P}_θ^T -a.s.) for all $\theta \in \Theta$.

The second statement of the theorem is obvious. \square

THEOREM 4.2.6. *If a sufficient statistic T is complete and $\hat{\theta} \in K_b$, then*

$$\hat{\theta}_T = \mathbf{E}_\theta(\hat{\theta}/T)$$

is a unique optimal estimator in the class K_b .

PROOF. Since T is complete, Theorem 4.2.5 implies that a \mathcal{B}_T -measurable estimator is unique in K_b .

Let $\tilde{\theta}$ be any other estimator of the class K_b . Then $\tilde{\theta}_T = \mathbf{E}_\theta(\tilde{\theta}/T) \in K_b$, thus we get by Theorem 4.2.5 that $\tilde{\theta}_T = \hat{\theta}_T$ (\mathbf{P}_θ -a.s.) for all $\theta \in \Theta$. This together with

the Rao–Blackwell–Kolmogorov theorem implies that

$$E_{\theta}(\widehat{\theta}_T - \theta)^2 = E_{\theta}(\widetilde{\theta}_T - \theta)^2 \leq E_{\theta}(\widetilde{\theta} - \theta)^2 \quad \text{for all } \theta \in \Theta.$$

Moreover the inequality becomes an equality if and only if $\widetilde{\theta} = \widehat{\theta}_T$ (P_{θ} -a.s.) for all $\theta \in \Theta$. \square

COROLLARY 4.2.1. *If T is a complete sufficient statistic and $\widehat{\theta}$ is an unbiased estimator of the parameter θ , then $\theta_T = E_{\theta}(\widehat{\theta}/T)$ is a unique optimal unbiased estimator of the parameter θ .*

PROOF. It is necessary to apply Theorem 4.2.6 to the class K_{b_0} where $b_0(\theta) = 0$ for all $\theta \in \Theta$. \square

EXAMPLE 4.2.3. Let an observation $\xi^{(n)} = (\xi_1, \xi_2, \dots, \xi_n)$ be a sample from the Poisson distribution with parameter $\theta \in \Theta = (0, \infty)$. Consider the estimator $\widehat{\theta}_n(\xi^{(n)}) = \xi_1$. It is clear that this estimator is unbiased for all $\theta \in \Theta$, that is, $E_{\theta}\widehat{\theta}_n = E_{\theta}\xi_1 = \theta$. At the same time it is not consistent, since it does not depend on n . Consider the statistic

$$T_n(x) = \sum_{i=1}^n x_i, \quad x = (x_1, x_2, \dots, x_n),$$

where $x_i \in \{0, 1, 2, \dots\}$ for all $i = 1, 2, \dots, n$. We learned in Examples 4.1.1 and 4.2.1 that T_n is a minimal sufficient statistic for the parameter θ . Moreover Example 4.1.1 implies that the conditional distribution of ξ_1 given T_n is of the form

$$P_{\theta} \left\{ \xi_1 = x/T(\xi^{(n)}) = y \right\} = \binom{y}{x} \left(\frac{1}{n} \right)^x \left(1 - \frac{1}{n} \right)^{y-x}$$

where $y \in \{0, 1, 2, \dots, n\}$ and $x \in \{0, 1, 2, \dots, y\}$. Thus

$$\begin{aligned} (4.2.8) \quad \widehat{\theta}_{T_n}(x) &= E_{\theta}(\xi_1/T_n(x)) = \sum_{k=0}^{T_n(x)} k \binom{T_n(x)}{k} \left(\frac{1}{n} \right)^k \left(1 - \frac{1}{n} \right)^{T_n(x)-k} \\ &= \frac{1}{n} T_n(x) = \bar{x}. \end{aligned}$$

Now we show that T_n is a complete statistic. Since the distribution of T_n is Poisson with parameter $n\theta$, equation (4.2.7) for T_n becomes of the form

$$\sum_{k=0}^{\infty} \phi(k) e^{-n\theta} \frac{(n\theta)^k}{k!} = 0 \quad \text{for all } \theta \in \Theta$$

or, equivalently,

$$(4.2.9) \quad v(z) = \sum_{k=0}^{\infty} \phi(k) \frac{z^k}{k!} = 0 \quad \text{for all } z \in \Theta.$$

The convergence of series (4.2.9) at $z = 1$ implies that the function $v(z)$ is analytic for $|z| < 1$, whence $\phi(k) = 0$ for all k . Taking into account equality (4.2.8), we obtain from Corollary 4.2.1 that \bar{x} is a unique optimal estimator of the parameter θ in the class of unbiased estimators.

EXAMPLE 4.2.4. Let $\xi^{(n)} = (\xi_1, \xi_2, \dots, \xi_n)$ be a sample from the uniform distribution on the interval $[0, \theta]$ where $\theta \in \Theta = (0, \infty)$ is an unknown parameter. It is natural to consider \mathbf{R}_+^n as the sample space in this case. We learned in Examples 4.1.5 and 4.2.2 that $T_n(x) = x_{n,n}$ is a sufficient and minimal statistic. Its distribution function is

$$P_\theta\{T_n < y\} = \left(\frac{y}{\theta}\right)^n, \quad 0 \leq y \leq \theta.$$

Thus equality (4.2.7) for T_n becomes of the form

$$\int_0^\theta \phi(y) \frac{ny^{n-1}}{\theta^n} dy = 0 \quad \text{for all } \theta \in \Theta,$$

whence it follows that $\phi(y)y^{n-1} = 0$ for almost all $y > 0$ and therefore $\phi(y) = 0$ for almost all $y > 0$. This means that T_n is a complete statistic.

Since θ is a scale parameter, Theorems 3.3.1 and 3.3.2 imply that

$$\hat{\theta}_n = T_n(n+1)/n$$

is an optimal estimator in the class of equivariant unbiased estimators of the scale parameter. Note that $\hat{\theta}_n$ is a \mathcal{B}_{T_n} -measurable estimator and Corollary 4.2.1 implies that $\hat{\theta}_n$ is an optimal estimator in the class of all unbiased estimators of the parameter θ .

The latter result can also be obtained directly from Corollary 4.2.1 (or from Theorem 4.2.6) by considering the estimator $\tilde{\theta} = 2x_1$ which obviously is unbiased. According to Corollary 4.2.1 $\tilde{\theta}_{T_n} = E_\theta(\tilde{\theta}/T_n)$ is a unique optimal unbiased estimator of the parameter θ . The conditional density of ξ_1 given T_n is given by

$$(4.2.10) \quad f(x/y) = \begin{cases} \left(1 - \frac{1}{n}\right) \frac{1}{y}, & 0 < x < y, \\ \frac{1}{n}, & x = y. \end{cases}$$

Indeed, since $\xi_1, \xi_2, \dots, \xi_n$ are independent identically distributed random variables and the model is symmetric, we get

$$P_\theta\{\xi_1 = T_n/T_n\} = \frac{1}{n}, \quad P_\theta\{\xi_1 < T_n/T_n\} = 1 - \frac{1}{n}.$$

Moreover given $\{T_n = y, \xi_1 < T_n\}$, the conditional distribution of ξ_1 is uniform on $(0, y)$, since the distribution of ξ_1 is uniform. This leads to equality (4.2.10). Using equality (4.2.10) we obtain

$$E_\theta\{\xi_1/\tilde{T}_n = y\} = \left(1 - \frac{1}{n}\right) \frac{y}{2} + \frac{y}{n}.$$

This implies that

$$\tilde{\theta}_{T_n} = E_\theta(\tilde{\theta}/T_n) = \left(1 - \frac{1}{n}\right) T_n + \frac{2}{n} T_n = \frac{n+1}{n} T_n.$$

Another consequence of Theorem 4.2.6 is an assertion on the optimality of a function $g(\theta)$ of the parameter θ in the class K^g of all unbiased estimators of the function $g(\theta)$.

COROLLARY 4.2.2. *If T is a complete sufficient statistic and \hat{g} is an unbiased estimator of the function $g(\theta)$, then $\hat{g}_T = E_\theta(\hat{g}/T)$ is a unique optimal unbiased estimator of the function $g(\theta)$ in the class K^g .*

PROOF. It is sufficient to apply Theorem 4.2.6 to estimators of the class K_{b_0} where $b_0(\theta) = g(\theta) - \theta$ for all $\theta \in \Theta$. \square

Corollary 4.2.2 is known as the *Lehmann–Scheffé theorem*; see [36], Theorem 3.1.2.

COROLLARY 4.2.3. *If T is a complete sufficient statistic, then any function $H(T)$ of it is a unique unbiased optimal estimator of its own expectation, that is, of the function $g(\theta) = E_\theta H(T)$.*

PROOF. To prove this result it is sufficient to consider the class of estimators K^g with $g(\theta) = E_\theta H(T)$ and to apply Corollary 4.2.2 with $\hat{g} = H(T)$. Then

$$H(T) = \hat{g}_T = E_\theta(\hat{g}/T)$$

is a unique optimal estimator of the function $g(\theta) = E_\theta H(T)$. \square

In fact we have a series of results allowing one to find optimal estimators of the function $g(\theta)$ if a complete sufficient statistic T exists, namely:

- 1) *if there exists an unbiased estimator of a function $g(\theta)$, then there exists an unbiased estimator that is a function of T ; if there is no unbiased estimator of the form $H(T)$, that is, the equation $E_\theta H(T) = g(\theta)$ has no solution, then the class K^g of unbiased estimators of the function $g(\theta)$ is empty;*
- 2) *an optimal unbiased estimator of the function $g(\theta)$ (if such an estimator exists at all) is a function of T and it is determined uniquely by the equality $E_\theta H(T) = g(\theta)$;*
- 3) *an optimal unbiased estimator g^* of the function $g(\theta)$ can be found as follows:*

$$(4.2.11) \quad g^* = \hat{g}_T = E_\theta(\hat{g}/T)$$

where \hat{g} is an arbitrary unbiased estimator of the function $g(\theta)$.

The latter method is rarely used when finding optimal estimators, since it requires the evaluation of the conditional expectation (4.2.11) which usually meets technical problems. Instead the equation

$$(4.2.12) \quad E_\theta H(T) = g(\theta), \quad \theta \in \Theta,$$

is used to determine the function H . There are several methods for solving equation (4.2.12). For example, one can expand both sides of (4.2.12) into power series of θ and equate corresponding coefficients.

The following result provides a relationship between complete and minimal statistics.

THEOREM 4.2.7. *Any complete sufficient statistic T is a minimal sufficient statistic.*

PROOF. Let \mathcal{B}^* be the minimal sufficient σ -algebra (this σ -algebra exists by Theorem 4.2.4).

First we assume that $E_\theta T$ exists and consider the function $\phi = T - E_\theta(T/\mathcal{B}^*)$. Since $\mathcal{B}^* \subset \mathcal{B}_T$ where \mathcal{B}_T is the σ -algebra generated by the statistic T , the function ϕ is \mathcal{B}_T -measurable, whence $\phi = \phi(T)$. Denote by P_θ^T the distribution of the statistic T . Then $E_\theta \phi(T) = 0$ for all $\theta \in \Theta$ or, equivalently,

$$\int \phi(y) P_\theta^T(dy) = 0 \quad \text{for all } \theta \in \Theta.$$

This implies that $\phi(y) = 0$ (P_θ^T -a.s.) for all $\theta \in \Theta$, since T is a complete statistic. This means that $T = E_\theta(T/\mathcal{B}^*)$ (P_θ^T -a.s.). Therefore T is a \mathcal{B}^* -measurable statistic and hence $\mathcal{B}_T = \mathcal{B}^*$.

If $E_\theta T$ does not exist, then one should consider the statistic $\arctan T$ instead of T . If T is either sufficient, or complete, or minimal, then so is $\arctan T$. \square

It is easy to construct an example of a minimal sufficient statistic that is not a complete statistic (see, for example, [36]).

Exponential families of distributions. Let ξ be an observation that is a random element assuming values in a measurable space (X, \mathcal{B}) whose distribution belongs to a family $\mathcal{P} = (P_\theta, \theta \in \Theta)$ where $\theta = (\theta_1, \theta_2, \dots, \theta_k)'$ is a k -dimensional parameter of the set $\Theta \subset \mathbf{R}^k$, $k \geq 1$. Let a family \mathcal{P} be dominated by some σ -finite measure μ and let $f(x; \theta) = dP_\theta/d\mu(x)$, $x \in X$, be the density of the measure P_θ with respect to the measure μ .

A family \mathcal{P} is called *exponential* if the density $f(x; \theta)$ is of the form

$$(4.2.13) \quad f(x; \theta) = h(x) \exp \left\{ \sum_{j=1}^k c_j(\theta) U_j(x) + V(\theta) \right\}$$

where all the functions on the right-hand side are finite and measurable.

Various distributions known in the literature are exponential. For example, normal, Poisson, binomial, Gamma-distributions, and others form exponential families of distributions.

If an observation is a sample $\xi^{(n)} = (\xi_1, \xi_2, \dots, \xi_n)$ from an exponential family of the form (4.2.13), then the distribution of the sample $\xi^{(n)}$ also belongs to some exponential family \mathcal{P}_n . Indeed, if $f_n(x; \theta)$ is the density of the distribution of the sample with respect to the measure μ^n , then

$$(4.2.14) \quad f_n(x; \theta) = h_n(x) \exp \{ c(\theta)' T(x) + nV(\theta) \}$$

where

$$c(\theta) = (c_1(\theta), \dots, c_k(\theta))', \quad T(x) = (T_1(x), \dots, T_k(x))',$$

$$h_n(x) = \prod_{i=1}^n h(x_i), \quad T_j(x) = \sum_{i=1}^n U_j(x_i), \quad x = (x_1, \dots, x_n) \in X^n.$$

It follows from the factorization criterion that if the family of distributions \mathcal{P} is exponential, then the statistic $U(x) = (U_1(x), \dots, U_n(x))$ is sufficient. Similarly, the statistic $T(x)$ is sufficient for a family of distributions \mathcal{P}_n whose densities are of the form (4.2.14). It turns out that these statistics are minimal and sufficient. We prove this result for families \mathcal{P}_n with densities (4.2.14).

Since the functions $c_j(\theta)$, $U_j(x)$, and $V(\theta)$ are finite, the exponent in (4.2.14) is positive everywhere. Thus as a measure \mathbf{Q} in Theorem 4.2.4 one can take a distribution concentrated at any fixed point $\theta^{(0)} \in \Theta$. All the measures \mathbf{P}_θ are absolutely continuous with respect to the measure $\mathbf{P}_\mathbf{Q} = \int \mathbf{P}_t \mathbf{Q}(dt) = \mathbf{P}_{\theta^{(0)}}$ in this case. Theorem 4.2.4 implies that the σ -algebra \mathcal{B}^* generated by the functions

$$\begin{aligned} p(x; \theta) &= f_n(x; \theta) / f_n(x; \theta^{(0)}) \\ &= \exp((c(\theta) - c(\theta^{(0)}))' T(x) + n(V(\theta) - V(\theta^{(0)}))) \end{aligned}$$

is the minimal sufficient σ -algebra for all $\theta \in \Theta$.

THEOREM 4.2.8. *Let functions $c_0(\theta) \equiv 1$, $c_1(\theta), \dots, c_k(\theta)$ be linearly independent on Θ . Then $T(x)$ in representation (4.2.14) is a minimal sufficient statistic.*

PROOF. Since the functions $1, c_1(\theta), \dots, c_k(\theta)$ are linearly independent on Θ , it follows that the functions $c_1(\theta) - c(\theta^{(0)}), \dots, c_k(\theta) - c(\theta^{(0)})$ are linearly independent. This means that there are k points $\theta^{(1)}, \dots, \theta^{(k)}$ in Θ such that the numbers

$$c_{ij} = c_i(\theta^{(j)}) - c_i(\theta^{(0)})$$

form the matrix C whose determinant is nonzero. This implies that the system of equations

$$(c(\theta^{(j)}) - c(\theta^{(0)}))' T(x) = \ln p(x; \theta^{(j)}) - n(V(\theta^{(j)}) - V(\theta^{(0)}))$$

for $j = 1, 2, \dots, k$ has a unique solution $T(x)$. Thus

$$\mathcal{B}_T \subset \sigma\{p(x; \theta^{(j)}), j = 1, 2, \dots, k\} \subset \mathcal{B}^*. \quad \square$$

Sometimes the assumptions of Theorem 4.2.8 are too restrictive if one proves only that T is a complete sufficient statistic. First we note that representation (4.2.14) yields

$$(4.2.15) \quad f_n(x; \theta) = g(T(x); \theta) r(x) \quad (\mu^n\text{-a.s.})$$

(see Theorem 4.1.2) where

$$\begin{aligned} g(y; \theta) &= \exp\{c(\theta)' y + nV(\theta)\}, \\ r(x) = h_n(x) &= \prod_{i=1}^n h(x_i), \quad x = (x_1, \dots, x_n). \end{aligned}$$

Consider the following measure on $(\mathbf{R}^k, \mathcal{B}^k)$:

$$\nu(B) = \int_{T^{-1}(B)} h_n(x) \mu^n(dx), \quad B \in \mathcal{B}^k,$$

where $T^{-1}(B) = \{x: T(x) \in B\}$. In what follows we need two auxiliary results.

LEMMA 4.2.1. *The distribution $\mathbf{P}_\theta^T(B) = \mathbf{P}_\theta\{x: T(x) \in B\}$, $B \in \mathcal{B}^k$, of the statistic T is absolutely continuous with respect to the measure ν and its density is $g(y; \theta)$.*

PROOF. It is sufficient to note that relation (4.2.15) implies

$$P_{\theta}^T(B) = \int_{T^{-1}(B)} g(T(x); \theta) r(x) \mu^n(dx) = \int_B g(y; \theta) \nu(dy)$$

where the latter equality follows from the change of variables theorem for the Lebesgue integral. \square

LEMMA 4.2.2. *Let G_1 and G_2 be two σ -finite measures on $(\mathbf{R}^k, \mathcal{B}^k)$ and for some parallelepiped $B \subset \mathbf{R}^k$ two integrals exist and are equal:*

$$\int e^{a'y} G_1(dy) = \int e^{a'y} G_2(dy)$$

for all $a \in B$. Then $G_1 = G_2$.

PROOF. We give the proof for the one-dimensional case, that is, for $k = 1$. Let $B = \{x: |x| \leq \tilde{a}\}$. Then the functions

$$\psi_j(a) = \int e^{ay} G_j(dy), \quad j = 1, 2,$$

are analytic for $|a| < \tilde{a}$. Moreover for all $b \in \mathbf{R}$ the functions

$$\psi_j(z) = \int e^{(a+ib)y} G_j(dy), \quad j = 1, 2,$$

of a complex variable $z = a + ib$ are well defined. It is clear that $\psi_j(z)$ are analytic functions in the strip $|a| < \tilde{a}$, $-\infty < b < \infty$. By assumption, $\psi_1(z) = \psi_2(z)$ on the interval $|a| < \tilde{a}$ of the line $b = 0$; thus $\psi_1(z) = \psi_2(z)$ for all z of the above strip. Thus for all $b \in (-\infty, \infty)$

$$(4.2.16) \quad \int e^{iby} G_1(dy) = \int e^{iby} G_2(dy).$$

Since $\psi_j(0) = \int G_j(dy) < \infty$, without loss of generality one can assume that G_j are probability measures. Since the correspondence between characteristic functions and distributions is one-to-one, equality (4.2.16) implies that $G_1 = G_2$.

If the parallelepiped B is of the form $\{x: |x - a_0| \leq \tilde{a}\}$, then we consider the measures $G_j^*(dy) = e^{a_0 y} G_j(dy)$ and follow the line of the above proof.

The proof for the multidimensional case $k > 1$ is the same. \square

THEOREM 4.2.9. *Let representation (4.2.14) hold for the density $f_n(x; \theta)$ of a sample $\xi^{(n)}$ and let the density belong to a family \mathcal{P}_n of distributions where the function $c(\theta)$ and the set Θ are such that the image of the set Θ under the mapping*

$$c: \Theta \rightarrow \mathbf{R}^k$$

contains some k -dimensional parallelepiped. Then the statistic T occurring in representation (4.2.14) is complete and sufficient.

PROOF. It is sufficient to show that if ϕ is a measurable function on $(\mathbf{R}^k, \mathcal{B}^k)$ and that there exists

$$(4.2.17) \quad \int \phi(y) P_{\theta}^T(dy) = 0, \quad \theta \in \Theta,$$

then $\phi(y) = 0$ (P_{θ}^T -a.s.) for all $\theta \in \Theta$ where P_{θ}^T is the distribution of the statistic T . Let $\phi = \phi^+ - \phi^-$ where $\phi^+ = 0 \vee \phi$ and $\phi^- = -(0 \wedge \phi)$. Then (4.2.17) implies that

$$\int \phi^+(y) P_{\theta}^T(dy) = \int \phi^-(y) P_{\theta}^T(dy) \quad \text{for all } \theta \in \Theta.$$

This yields by Lemma 4.2.1 that

$$\begin{aligned} \int \phi^+(y) g(y; \theta) \nu(dy) &= \int \phi^-(y) g(y; \theta) \nu(dy) \quad \text{for all } \theta \in \Theta, \\ \int \phi^+(y) e^{c(\theta)'y} \nu(dy) &= \int \phi^-(y) e^{c(\theta)'y} \nu(dy) \quad \text{for all } \theta \in \Theta. \end{aligned}$$

Consider σ -finite measures $\nu_{\pm}(dy) = \phi^{\pm}(y) \nu(dy)$. By the assumptions of the theorem we get

$$\int e^{c'y} \nu_+(dy) = \int e^{c'y} \nu_-(dy)$$

for all c of some parallelepiped in \mathbf{R}^k . Now it remains to apply Lemma 4.2.2. \square

COROLLARY 4.2.4. *Let all the assumptions of Theorem 4.2.9 hold. Let $\hat{\theta}$ be an estimator of the parameter θ of the class K_b constructed from a sample $\xi^{(n)}$. Then $\hat{\theta}_T = E_{\theta}(\hat{\theta}/T)$ is an optimal estimator of the parameter θ in the class K_b where T is the statistic occurring in representation (4.2.14).*

COROLLARY 4.2.5. *Let all the assumptions of Theorem 4.2.9 hold. Let \hat{g} be some estimator of a function $g(\theta)$ of a parameter θ of the class K^g constructed from a sample $\xi^{(n)}$. Then $\hat{g}_T = E_{\theta}(\hat{g}/T)$ is an optimal estimator of the function $g(\theta)$ in the class K^g .*

PROOF. It is sufficient to apply Corollary 4.2.4 to estimators of the class K_b where $b(\theta) = g(\theta) - \theta$ for all $\theta \in \Theta$. \square

Some applications of sufficient statistics. In a series of examples above we learned how to construct sufficient, or minimal, or complete and sufficient statistics. We consider in this section some applications of sufficient statistics for several models and construct optimal estimators of a parameter θ and a function $g(\theta)$.

EXAMPLE 4.2.5. Let an observation be a vector $\xi^{(n)} = (\xi_1, \xi_2, \dots, \xi_n)$ where

$$\xi_i = \theta + \eta_i, \quad i = 1, 2, \dots, n,$$

and random variables η_1, \dots, η_n (being dependent, generally speaking) do not depend on θ and are such that $E\eta_i = 0$ and $E\eta_i^2 < \infty$. Consider the class L of linear unbiased estimators of the parameter θ , that is, the class of functions $l = l(x) = \sum_{i=1}^n c_i x_i$ where $\sum_{i=1}^n c_i = 1$, $x = (x_1, \dots, x_n)$. Our aim is to construct an optimal estimator $l^* = \sum_{i=1}^n c_i^* x_i$ in the class L with respect to the quadratic loss function.

Assume that random variables η_1, η_2, \dots form an autoregression sequence of first order, more precisely let

$$(4.2.18) \quad \eta_1 = \varepsilon_1, \quad \eta_j = \lambda\eta_{j-1} + \varepsilon_j, \quad j = 2, 3, \dots,$$

where $\varepsilon_1, \varepsilon_2, \dots$ are independent Gaussian random variables with $E\varepsilon_j = 0$ and $E\varepsilon_j^2 = \sigma_j^2$, $0 < \sigma_j^2 < \infty$, for all j ; here $\lambda \neq 1$ is a known parameter.

Let $p(u_1, \dots, u_n)$ be the probability density of the vector (η_1, \dots, η_n) and let $\phi_j(x)$ be the probability density of the normal $\mathcal{N}(0, \sigma_j^2)$ law. Since

$$\varepsilon_1 = \eta_1, \quad \varepsilon_2 = \eta_2 - \lambda\eta_1, \quad \dots, \quad \varepsilon_n = \eta_n - \lambda\eta_{n-1}$$

by (4.2.18), the density $p(u_1, \dots, u_n)$ is of the form

$$p(u_1, \dots, u_n) = \phi_1(u_1)\phi_2(u_2 - \lambda u_1) \cdots \phi(u_n - \lambda u_{n-1}).$$

This implies that the probability density of the vector (ξ_1, \dots, ξ_n) is equal to

$$(4.2.19) \quad \begin{aligned} f(x_1, \dots, x_n; \theta) &= p(x_1 - \theta, \dots, x_n - \theta) \\ &= \phi_1(x_1 - \theta)\phi_2(x_2 - \theta - \lambda(x_1 - \theta)) \cdots \phi_n(x_n - \theta - \lambda(x_{n-1} - \theta)) \\ &= \prod_{j=1}^n \frac{1}{\sqrt{2\pi}\sigma_j} \exp \left\{ -\frac{(x_1 - \theta)^2}{2\sigma_1^2} - \sum_{j=2}^n \frac{(x_j - \theta - \lambda(x_{j-1} - \theta))^2}{2\sigma_j^2} \right\} \\ &= C(\theta)R(x_1, \dots, x_n) \exp \left\{ \theta \left[\frac{x_1}{\sigma_1^2} + (1 - \lambda) \sum_{j=2}^n \frac{x_j - \lambda x_{j-1}}{\sigma_j^2} \right] \right\}. \end{aligned}$$

In what follows we do not need the explicit expressions for $C(\theta)$ and $R(x_1, \dots, x_n)$.

It follows from (4.2.19) that the linear statistic

$$(4.2.20) \quad \begin{aligned} l(x) &= \frac{x_1}{\sigma_1^2} + (1 - \lambda) \sum_{j=2}^n \frac{x_j - \lambda x_{j-1}}{\sigma_j^2} \\ &= \left(\frac{1}{\sigma_1^2} - \frac{\lambda(1 - \lambda)}{\sigma_2^2} \right) x_1 + (1 - \lambda) \sum_{j=2}^{n-1} \left(\frac{1}{\sigma_j^2} - \frac{1}{\sigma_{j+1}^2} \right) x_j + \frac{1 - \lambda}{\sigma_n^2} x_n \end{aligned}$$

is sufficient for the parameter θ . It is also seen from (4.2.19) that the linear statistic given by (4.2.20) is complete and sufficient (according to Theorem 4.2.9). It is clear that the optimal linear estimator $l^*(x)$ is a function of the statistic $l(x)$. Thus the coefficients c_j^* are proportional to the corresponding coefficients of the statistic $l(x)$. Thus

$$\begin{aligned} c_1^* &= c \left(\frac{1}{\sigma_1^2} - \frac{\lambda(1 - \lambda)}{\sigma_2^2} \right), \\ c_j^* &= c(1 - \lambda) \left(\frac{1}{\sigma_j^2} - \frac{\lambda}{\sigma_{j+1}^2} \right), \quad 2 \leq j \leq n - 1, \\ c_n^* &= c \frac{1 - \lambda}{\sigma_n^2} \end{aligned}$$

where the constant c is defined from the condition $\sum_{j=1}^n c_j^* = 1$.

The assumption that the distribution of the random variables $\varepsilon_1, \varepsilon_2, \dots, \varepsilon_n$ is Gaussian is not crucial. The same result can be obtained for a general distribution of random variables $\varepsilon_1, \varepsilon_2, \dots, \varepsilon_n$ [15].

EXAMPLE 4.2.6. Let $\xi^{(n)} = (\xi_1, \dots, \xi_n)$ be a sample from the Bernoulli distribution, that is, $P_\theta(\xi_i = 1) = \theta$ and $P_\theta(\xi_i = 0) = 1 - \theta$, $i = 1, 2, \dots, n$, where $\theta \in (0, 1)$ is an unknown parameter. The total number of "successes" $T(x) = \sum_{i=1}^n x_i$ in n Bernoulli trials is a sufficient statistic (see Example 4.1.7). We show that T is a complete statistic.

Note that T has the binomial distribution:

$$P_\theta\{T = y\} = \binom{n}{y} \theta^y (1 - \theta)^{n-y}, \quad y = 0, 1, 2, \dots, n.$$

Let $\phi(y)$ be an arbitrary function on $\{0, 1, 2, \dots, n\}$. Then condition (4.2.7) can be rewritten as

$$\sum_{y=0}^n \phi(y) \binom{n}{y} \theta^y (1 - \theta)^{n-y} = 0 \quad \text{for all } \theta \in (0, 1)$$

or, putting $x = \theta/(1 - \theta)$,

$$\sum_{y=0}^n \phi(y) \binom{n}{y} x^y = 0 \quad \text{for all } x > 0.$$

This implies that all the coefficients of the latter polynomial are zero, that is, $\phi(y) = 0$ for all $y = 0, 1, \dots, n$. Thus T is a complete sufficient statistic and $\hat{\theta} = T/n$ is an optimal estimator of the parameter θ . Moreover, according to Corollary 4.2.3 any function of T is an optimal estimator of its own mean.

Since the moment generating function of the random variable T is

$$E_\theta z^T = \phi(z; \theta) = (1 + (z - 1)\theta)^n,$$

we put $(a)_k = a(a - 1) \cdots (a - k + 1)$, $k \geq 1$, and obtain

$$E_\theta(T)_k = \left. \frac{\partial^k \phi(z; \theta)}{\partial z^k} \right|_{z=1} = (n)_k \theta^k.$$

This implies that for any integer k , $1 \leq k \leq n$, the statistic $(T)_k/(n)_k$ is an optimal estimator of the function θ^k . At the same time, other functions θ^i with $i > n$ cannot be estimated from a sample of size n in the class of unbiased estimators. Finally, according to Corollary 4.2.3 we get that $\hat{g} = \sum_{j=1}^n c_j (T)_j / (n)_j$ is an optimal estimator for the polynomial $g(\theta) = \sum_{j=1}^k c_j \theta^j$ if $k \leq n$. Therefore

$$\hat{\tau} = \frac{T}{n} - \frac{T(T-1)}{n(n-1)} = \frac{T(n-T)}{n(n-1)}$$

is an optimal estimator of the variance $\tau(\theta) = \theta(1 - \theta)$ in the case of the binomial distribution.

EXAMPLE 4.2.7. Let a discrete random variable assume values $l, l+1, \dots$ with probabilities

$$(4.2.21) \quad f(x; \theta) = P_\theta\{\xi = x\} = \frac{a(x)\theta^x}{f(\theta)}, \quad x = l, l+1, \dots,$$

where $f(\theta) = \sum_{x=l}^{\infty} a(x)\theta^x$ is a series whose radius of convergence R is nonzero. We treat $\theta \in \Theta = (0, R)$ as an unknown parameter. Discrete distributions (4.2.21) are sometimes called *power series distributions* (see [21]).

Distributions (4.2.21) include many well-known discrete distributions with an infinite number of values, say Poisson ($f(\theta) = e^\theta$, $R = \infty$), logarithmic ($f(x; \theta) = \theta^x / \ln(1-\theta)^{-x}$, $x = 0, 1, 2, \dots$, $R = 1$), negative binomial ($f(\theta) = (1-\theta)^{-r}$, $R = 1$), and others as well as their truncated versions.

Let $\xi^{(n)} = (\xi_1, \dots, \xi_n)$ be a sample from a distribution (4.2.21). Then the distribution of the sample is of the form

$$(4.2.22) \quad f_n(x; \theta) = P_\theta\{\xi^{(n)} = x\} = \theta^{T(x)} \prod_{i=1}^n a(x_i) f^{-n}(\theta), \quad \theta \in \Theta,$$

where $T(x) = \sum_{i=1}^n x_i$, $x = (x_1, \dots, x_n)$, $x_i = l, l+1, \dots$, for all $i = 1, 2, \dots, n$. This together with the factorization criterion implies that $T(x)$ is a sufficient statistic.

Since (4.2.21) is a distribution of the exponential type, Theorem 4.2.9 implies that $T(x)$ is a complete and sufficient statistic.

Note also that (4.2.21) yields

$$E_\theta \xi = \sum_{x=l}^{\infty} x a(x) \theta^x / f(\theta) = \frac{\theta f'(\theta)}{f(\theta)} = \tau(\theta).$$

Then $\hat{\tau} = n^{-1}T(x) = n^{-1} \sum_{i=1}^n x_i$, $x = (x_1, \dots, x_n)$, is an optimal unbiased estimator of the function $\tau(\theta) = E_\theta \hat{\tau}$ by Corollary 4.2.3.

Let a function $g(\theta)$ be represented as a power series $g(\theta) = \sum_{j=r}^{\infty} a_j \theta^j$ convergent on Θ . To estimate the function $g(\theta)$ we find the distribution of the statistic $T(x)$. We have

$$(4.2.23) \quad g(t; \theta) = P_\theta\{T(\xi^{(n)}) = t\} = \sum_{\{x: T(x)=t\}} f_n(x; \theta) = \theta^t b_n(t) f^{-n}(\theta)$$

where $t = nl, nl+1, \dots$ and

$$b_n(t) = \sum_{x_1 + \dots + x_n = t} a(x_1) \cdots a(x_n).$$

Thus $b_n(t)$ is equal to the coefficient of z^t in the expansion of the function $f^n(z)$. It follows from (4.2.23) that the condition (4.2.12) can be rewritten as

$$\begin{aligned} \sum_{t=nl}^{\infty} H(t) b_n(t) \theta^t &= f^n(\theta) g(\theta) = \sum_{i=nl}^{\infty} b_n(i) \theta^i \sum_{j=r}^{\infty} a_j \theta^j \\ &= \sum_{k=nl+r}^{\infty} \theta^k \sum_{j=r}^{k-nl} a_j b_n(k-j). \end{aligned}$$

This is an identity with respect to θ . Equating the coefficients of θ^t we get

$$H(t)b_n(t) = \begin{cases} \sum_{j=r}^{t-nl} a_j b_n(k-j), & \text{if } t \geq nl+r, \\ 0, & \text{if } t < nl+r. \end{cases}$$

This implies that the optimal estimator g^* of the function $g(\theta)$ is of the form

$$g^* = H(T) = \begin{cases} b_n^{-1}(T) \sum_{j=r}^{T-nl} a_j b_n(T-j), & \text{if } T \geq nl+r, \\ 0, & \text{if } T < nl+r. \end{cases}$$

In particular, if $g(\theta) = \theta^r$ for some $r \geq 1$, then

$$g^* = \begin{cases} b_n(T-r)/b_n(T), & \text{if } T \geq nl+r, \\ 0, & \text{if } T < nl+r. \end{cases}$$

Therefore one can construct optimal estimators for an arbitrary function of the parameter represented as a power series of θ in the case of discrete distributions (4.2.21).

General Methods for Constructing Estimators

In the previous two chapters we dealt with the optimal and efficient statistical estimators of unknown parameters of distributions or of functions of parameters. We considered some methods for constructing optimal and efficient estimators based on some properties of families of distributions. In particular, we considered the method based on sufficient statistics.

In this chapter, we consider general methods of forming estimators, namely the method of moments, the maximum likelihood method, the Bayes method, and the integral estimation method.

5.1. Method of moments

The oldest general method proposed to construct estimators of unknown parameters is the *method of moments* introduced by K. Pearson (1894). This method can be described as follows.

Let $\xi^{(n)} = (\xi_1, \xi_2, \dots, \xi_n)$ be a sample from a distribution belonging to a family of distributions $\{P_\theta, \theta \in \Theta\}$ where $\theta = (\theta_1, \dots, \theta_k)$ is an unknown parameter of a set $\Theta \subset \mathbf{R}^k$, $k \geq 1$. Assume that $E_\theta |\xi_1|^k < \infty$. Then the following k functions $\alpha_j(\theta) = \alpha_j(\theta_1, \dots, \theta_k) = E_\theta \xi_1^j$, $j = 1, 2, \dots, k$, are well defined. Let a_j be the sampling moment of order j constructed from the sample $\xi^{(n)}$, that is,

$$a_j = \frac{1}{n} \sum_{i=1}^n \xi_i^j, \quad j = 1, 2, \dots, k.$$

Consider the system of equations

$$(5.1.1) \quad \alpha_j(\theta_1, \dots, \theta_k) = a_j, \quad j = 1, 2, \dots, k,$$

with unknowns $\theta_1, \dots, \theta_k$. Solutions $\widehat{\theta}_j$, $j = 1, 2, \dots, k$, of system (5.1.1), if they exist, are called *method of moments estimators of a parameter*.

Note that $E_\theta |\xi_1|^j < \infty$ for all $j = 1, 2, \dots, k$ if $E_\theta |\xi_1|^k < \infty$. Using results of Section 1.2 we obtain $a_j \rightarrow \alpha_j(\theta_1, \dots, \theta_k)$ as $n \rightarrow \infty$ in probability P_θ for all $j = 1, 2, \dots, k$. Assume that the functions $\alpha_j(\theta_1, \dots, \theta_k)$, $j = 1, 2, \dots, k$, determine a continuous one-to-one correspondence between vectors $(\theta_1, \theta_2, \dots, \theta_k)$ and (a_1, \dots, a_k) , that is, there exist continuous functions ϕ_j , $j = 1, 2, \dots, k$, such that $\alpha_j = \phi_j(a_1, \dots, a_k)$, $j = 1, 2, \dots, k$. Then solutions of system (5.1.1) can be represented as

$$(5.1.2) \quad \widehat{\theta}_j = \phi_j(a_1, \dots, a_k), \quad j = 1, 2, \dots, k.$$

Thus estimators (5.1.2) are consistent by Theorem 1.2.2.

The sampling moments a_j are asymptotically normal if $E_\theta |\xi_1|^{2j} < \infty$ (see Theorem 1.2.3). Moreover, we noticed in Section 1.2 that the asymptotic normality

can be proved for a continuous function of a finite number of sampling moments a_j provided suitable conditions on the function are posed (more details are given in [5, 9]). Therefore conditions can be posed on moments of random variables ξ_i to guarantee that estimators (5.1.2) are asymptotically normal.

On the other hand, Fisher (1921) pointed out that estimators (5.1.2) are not asymptotically efficient. Moreover, the method of moments cannot be applied in the cases where the corresponding moments do not exist (say, in the case of the Cauchy distribution). Nevertheless the method of moments has an advantage because of its practical expediency. Estimators (5.1.2) can be treated sometimes as a first approximation used for other methods to construct estimators of a higher efficiency.

REMARK 5.1.1. One can also use another form of the method of moments, namely one can use the moments $m_j(\theta) = E_{\theta}g_j(\xi_1)$, $j = 1, 2, \dots, k$, instead of moments $\alpha_j(\theta)$, $j = 1, 2, \dots, k$, where $g_j(x)$, $j = 1, \dots, k$, are some measurable functions. Then we obtain the system of equations

$$(5.1.3) \quad m_j(\theta) = \frac{1}{n} \sum_{i=1}^n g_j(\xi_i), \quad j = 1, 2, \dots, k,$$

instead of (5.1.1). Solutions of system (5.1.3) are also called *method of moments estimators* (more detail is given in [5]). Note that system (5.1.3) reduces to (5.1.1) if $g_j(x) = x^j$, $j = 1, 2, \dots, k$.

EXAMPLE 5.1.1. Let $\xi^{(n)} = (\xi_1, \dots, \xi_n)$ be a sample from the Gamma distribution, so that the density is

$$f(x; \theta) = \frac{1}{\Gamma(\theta)} x^{\theta-1} e^{-x} I_{(0, \infty)}(x)$$

where $\theta \in \Theta = (0, \infty)$ is an unknown parameter. In this case $\alpha_1(\theta) = E_{\theta}\xi_1 = \theta$ and therefore a solution of the equation

$$\alpha_1(\theta) = a_1 = \frac{1}{n} \sum_{i=1}^n \xi_i$$

is a method of moments estimator and it is given by $\hat{\theta} = n^{-1} \sum_{i=1}^n \xi_i$. This estimator is unbiased and consistent. On the other hand, this estimator is not asymptotically efficient whatever the parameter θ is (see Example 3.4.3).

EXAMPLE 5.1.2. Let $\xi^{(n)} = (\xi_1, \dots, \xi_n)$ be a sample from a distribution with the density $f(x; \theta) = \theta e^{-\theta x} I_{(0, \infty)}(x)$ where $\theta \in \Theta = (0, \infty)$ is an unknown parameter. We use two functions $g_1(x) = x$ and $g_2(x) = x^2$ to construct estimators according to the method of moments (see Remark 5.1.1). Since $m_j(\theta) = E_{\theta}g_j(\xi_1) = E_{\theta}\xi_1^j = j\theta^{-j}$ for $j = 1, 2$, equations (5.1.3) are of the form

$$m_1(\theta) = a_1 = \frac{1}{n} \sum_{i=1}^n \xi_i,$$

$$m_2(\theta) = a_2 = \frac{1}{n} \sum_{i=1}^n \xi_i^2.$$

There are two solutions of these equations with respect to θ :

$$\hat{\theta} = \left(\frac{1}{n} \sum_{i=1}^n \xi_i \right)^{-1}, \quad \tilde{\theta} = \left(\frac{1}{2n} \sum_{i=1}^n \xi_i^2 \right)^{-1/2}$$

Every solution is a method of moments estimator of the parameter θ . One can show (see [5]) that both estimators $\hat{\theta}$ and $\tilde{\theta}$ are asymptotically normal with parameters $\mathcal{N}(\theta, n^{-1}\theta^2)$ and $\mathcal{N}(\theta, (5/4)n^{-1}\theta^2)$, respectively. Thus the estimator $\hat{\theta}$ is better than $\tilde{\theta}$, since $n^{-1}\theta^2 < (5/4)n^{-1}\theta^2$. The Fisher information of $\hat{\theta}$ is $I(\theta) = \theta^{-2}$, whence it follows that $\hat{\theta}$ is asymptotically efficient in the weak sense.

EXAMPLE 5.1.3. Let $\xi^{(n)} = (\xi_1, \dots, \xi_n)$ be a sample from the Gamma distribution with the density

$$f(x; \theta) = \frac{1}{\theta_1^{\theta_2} \Gamma(\theta_2)} x^{\theta_2-1} e^{-x/\theta_1} I_{(0, \infty)}(x)$$

where $\theta = (\theta_1, \theta_2) \in \Theta = \{(\theta_1, \theta_2): \theta_1 > 0, \theta_2 > 0\}$. It is clear that

$$\alpha_j(\theta) = E_{\theta} \xi_1^j = \theta_1^j \frac{\Gamma(\theta_2 + j)}{\Gamma(\theta_2)} = \theta_1^j \theta_2 (\theta_2 + 1) \cdots (\theta_2 + j - 1).$$

In particular, $\alpha_1(\theta) = \theta_1 \theta_2$ and $\alpha_2(\theta) = \theta_1^2 \theta_2 (\theta_2 + 1)$. Thus solutions of the system of equations (5.1.1) are of the form

$$\hat{\theta}_1 = \frac{a_2 - a_1^2}{a_1}, \quad \hat{\theta}_2 = \frac{a_1^2}{a_2 - a_1^2}.$$

REMARK 5.1.2. When considering maximum likelihood estimators for samples from the normal distribution, we will show that estimators obtained by the method of moments and by the maximum likelihood method coincide and both are efficient. This is an exceptional case where the method of moments estimators are efficient.

5.2. The maximum likelihood method

From a theoretical point of view, the most important general method of estimation of parameters is the *method of moments*. In particular cases, this method is already used by F. Gauss. As a general method of estimation it was first introduced by Fisher (1912) and afterwards it was further developed by the same author. In 1925 Fisher studied asymptotic properties of maximum likelihood estimators.

Maximum likelihood estimators. Let ξ be an observation that is a random element assuming values in a measurable space (X, \mathcal{B}) and whose distribution is determined by a measure of a family of measures $\mathcal{P} = \{P_{\theta}, \theta \in \Theta\}$ where

$$\theta = (\theta_1, \dots, \theta_k)$$

is an unknown parameter belonging to a set $\Theta \subset \mathbf{R}^k$, $k \geq 1$. Let the family \mathcal{P} be dominated by a σ -finite measure μ and let $f(x; \theta)$ be the density of the measure P_{θ} with respect to the measure μ . The function $f(x; \theta)$ of an argument θ is called the *likelihood function*, while $L(x; \theta) = \ln f(x; \theta)$ is called the *logarithmic likelihood function*.

A statistic $\hat{\theta} = \hat{\theta}(x)$ such that

$$(5.2.1) \quad L(x; \hat{\theta}(x)) = \sup_{\theta \in \Theta} L(x; \theta), \quad x \in X,$$

is called the *maximum likelihood estimator of the parameter* θ if such a point

$$\hat{\theta}(x) \in \Theta$$

exists. Otherwise, if there is no $\hat{\theta}(x)$ satisfying (5.2.1), then we take an arbitrary point of Θ as $\hat{\theta}(x)$. If the function $L(x; \theta)$ is continuous with respect to θ and the set Θ is close, then the supremum on the right-hand side of (5.2.1) is attained, thus the maximum likelihood estimator is well defined.

The maximum likelihood estimator $\hat{\theta}_n$ can be defined in an equivalent way as a statistic maximizing the likelihood function $f(x; \theta)$. Below we consider the case where the likelihood functions are differentiable with respect to θ . One can substitute the closure $\bar{\Theta}$ instead of Θ in (5.2.1) in this case, the supremum on the right-hand side of (5.2.1) is attained, and the maximum likelihood estimator exists. If the supremum on the right-hand side of (5.2.1) is attained at an interior point of Θ and the function $L(x; \theta)$ is differentiable with respect to θ , then one can seek the maximum likelihood estimator $\hat{\theta}$ among solutions of the system of equations

$$(5.2.2) \quad \frac{\partial}{\partial \theta_j} L(x; \theta) = 0, \quad j = 1, 2, \dots, k.$$

The equations of system (5.2.2) are called *likelihood equations*.

The following are two properties of the maximum likelihood estimators:

- 1) if there exists an efficient unbiased estimator $T = T(x)$ of a scalar parameter θ , then the maximum likelihood estimator $\hat{\theta}$ exists and coincides with the estimator T ;
- 2) if there exists a sufficient statistic $T = T(x)$ and the maximum likelihood estimator exists and is unique, then $\hat{\theta}$ is a function of T .

It is sufficient to apply Theorem 3.4.1 (or Theorem 3.4.3) to prove property 1), and the factorization criterion to prove property 2).

Consider some examples of maximum likelihood estimators.

EXAMPLE 5.2.1. Let $\xi^{(n)} = (\xi_1, \dots, \xi_n)$ be a sample from the normal $\mathcal{N}(\theta_1, \theta_2)$ distribution where $\theta = (\theta_1, \theta_2)$ is an unknown parameter such that $\theta_1 \in (-\infty, \infty)$ and $\theta_2 > 0$. Then the logarithmic likelihood function for the distribution of the sample is given by

$$L(x; \theta) = -\frac{n}{2} \ln(2\pi\theta_2) - \frac{1}{2\theta_2} \sum_{i=1}^n (x_i - \theta_1)^2, \quad x = (x_1, \dots, x_n).$$

The system of likelihood equalities (5.2.2) is of the form in this case:

$$\begin{aligned} \frac{\partial L(x; \theta)}{\partial \theta_1} &= \frac{1}{\theta_2} \sum_{i=1}^n (x_i - \theta_1) = 0, \\ \frac{\partial L(x; \theta)}{\partial \theta_2} &= \frac{1}{2\theta_2^2} \sum_{i=1}^n (x_i - \theta_1)^2 - \frac{n}{2\theta_2} = 0. \end{aligned}$$

Solving this system of equations with respect to θ_1 and θ_2 we obtain the following maximum likelihood estimators:

$$(5.2.3) \quad \hat{\theta}_1 = \frac{1}{n} \sum_{i=1}^n x_i, \quad \hat{\theta}_2 = \frac{1}{n} \sum_{i=1}^n (x_i - \hat{\theta}_1)^2.$$

It is easy to see that these estimators coincide with the estimators of the parameters θ_1 and θ_2 obtained by the method of moments. Further, it is easy to check that the maximum of the function $L(x; \theta)$ is attained at the point $(\hat{\theta}_1, \hat{\theta}_2)$. Thus the maximum likelihood estimator exists, is unique, and is defined by (5.2.3). Note that this maximum likelihood estimator is a function of the sufficient statistic $T = (T_1, T_2)$ considered in Example 4.1.3. Moreover the maximum likelihood estimator $\hat{\theta} = (\hat{\theta}_1, \hat{\theta}_2)$ itself is a sufficient statistic (see Example 4.1.3).

EXAMPLE 5.2.2. Let $\xi^{(n)} = (\xi_1, \dots, \xi_n)$ be a sample from the uniform distribution on the interval $[0, \theta]$ where $\theta > 0$ is an unknown parameter. The likelihood function in this case is of the form

$$f(x; \theta) = \theta^{-n} I_{[0, \infty)}(x_{n,1}) I_{(-\infty, \theta]}(x_{n,n})$$

(see Example 4.1.5). The function $f(x; \theta)$ is discontinuous with respect to θ and moreover $f(x; \theta) = 0$ for $\theta < x_{n,n}$ and $f(x; \theta) = \theta^{-n}$ for $\theta \geq x_{n,n}$ if $x_{n,1} \geq 0$. Then the maximum likelihood estimator is $\hat{\theta} = x_{n,n}$. If $x_{n,1} < 0$, then $f(x; \theta) = 0$ for all θ and any number can be taken as the maximum likelihood estimator, in particular one can put $\hat{\theta} = x_{n,n}$. If the sampling space is \mathbf{R}_+^n , then the maximum likelihood estimator $\hat{\theta} = x_{n,n}$ is unique and is a complete sufficient statistic (see Example 4.2.4). We learned from Example 4.2.4 that the estimator $\tilde{\theta} = \frac{n+1}{n} x_{n,n}$ is optimal in the class of unbiased estimators of the parameter θ .

EXAMPLE 5.2.3. Let $\xi^{(n)} = (\xi_1, \dots, \xi_n)$ be a sample from the uniform distribution on the interval $[\theta, \theta + 1]$ where $\theta \in (-\infty, \infty)$ is an unknown parameter. The likelihood function in this case is given by

$$f(x; \theta) = I_{[\theta, \infty)}(x_{n,1}) I_{(-\infty, \theta+1]}(x_{n,n}), \quad x = (x_1, \dots, x_n).$$

The maximum likelihood estimator is not unique in this case. In particular, $\hat{\theta} = x_{n,1}$ is one of the maximum likelihood estimators, another one is $\tilde{\theta} = x_{n,n} - 1$. Note that $T(x) = (x_{n,1}, x_{n,n})$ is a sufficient statistic in this case.

The invariance principle for maximum likelihood estimators. The following result is known as the invariance principle for maximum likelihood estimators with respect to the change of a parameter.

THEOREM 5.2.1. *Let $\mathcal{P} = \{P_\theta, \theta \in \Theta\}$ be a family of probability measures on (X, \mathcal{B}) defining the distribution of an observation ξ , and let $g = g(\theta)$ be a one-to-one mapping of Θ into some set G . If $\hat{\theta}$ is a maximum likelihood estimator of the parameter θ constructed from an observation ξ , then $\hat{g} = g(\hat{\theta})$ is a maximum likelihood estimator of the function $g(\theta)$ constructed from the observation ξ .*

PROOF. Let $\theta(\gamma)$ be the inverse function to $g(\theta)$ and let $\mathbf{Q}_\gamma = P_{\theta(\gamma)}$ for all $\gamma \in G$. Then the logarithmic likelihood function for the family $\{\mathbf{Q}_\gamma, \gamma \in G\}$ is of the form

$$(5.2.4) \quad M(\gamma; x) = \ln(d\mathbf{Q}_\gamma/d\mu(x)) = L(\theta(\gamma); x).$$

Let $\hat{\theta}$ be the maximum likelihood estimator of the parameter θ and let $\hat{\gamma}$ be the maximum likelihood estimator of the parameter γ . Then equality (5.2.4) implies that $\hat{\theta} = \theta(\hat{\gamma})$ or, equivalently, $\hat{\gamma} = g(\hat{\theta})$. \square

EXAMPLE 5.2.4. Let $\xi^{(n)} = (\xi_1, \dots, \xi_n)$ be a sample from the logarithmic distribution, that is, $\mathcal{L}(\ln \xi_1) = \mathcal{N}(\theta_1, \theta_2)$ where $-\infty < \theta_1 < \infty$ and $0 < \theta_2 < \infty$. Let $\theta = (\theta_1, \theta_2)$. It is not hard to show that

$$\gamma_1 = E_\theta \xi_1 = \exp \left\{ \theta_1 + \frac{1}{2} \theta_2 \right\}, \quad \gamma_2 = D_\theta \xi_1 = \theta_1^2 (e^{\theta_2} - 1).$$

Consider the function $g(\theta_1, \theta_2) = (\gamma_1, \gamma_2)$ and find the maximum likelihood estimator $(\hat{\gamma}_1, \hat{\gamma}_2)$ of the function $g(\theta_1, \theta_2)$. We obtain from Example 5.2.1 that the maximum likelihood estimator $\hat{\theta} = (\hat{\theta}_1, \hat{\theta}_2)$ is of the form

$$\hat{\theta}_1 = \frac{1}{n} \sum_{i=1}^n \eta_i, \quad \hat{\theta}_2 = \frac{1}{n} \sum_{i=1}^n (\eta_i - \hat{\theta}_1)^2$$

where $\eta_i = \ln \xi_i$. We obtain from the invariance principle for maximum likelihood estimators that

$$\hat{\gamma}_1 = \exp \left\{ \hat{\theta}_1 + \frac{1}{2} \hat{\theta}_2 \right\}, \quad \hat{\gamma}_2 = \hat{\theta}_1^2 (e^{\hat{\theta}_2} - 1).$$

Asymptotic properties of maximum likelihood estimators. In this section we consider asymptotic properties of maximum likelihood estimators, namely we prove that a maximum likelihood estimator is consistent, asymptotically normal, and asymptotically efficient.

Let an observation $\xi^{(n)} = (\xi_1, \dots, \xi_n)$ be a sample from a distribution belonging to a family $\{P_\theta, \theta \in \Theta\}$ dominated by some σ -finite measure μ . Denote by $f(x; \theta)$ the density of the measure P_θ with respect to μ . Then $f_n(x; \theta) = \prod_{i=1}^n f(x_i; \theta)$, $x = (x_1, \dots, x_n)$, is the likelihood function, while $L_n(x; \theta) = \sum_{i=1}^n \ln f(x_i; \theta)$ is the logarithmic likelihood function. We denote by $\hat{\theta}_n = \hat{\theta}_n(x)$ the maximum likelihood estimator.

First we consider the case of a one-dimensional parameter θ , that is, we consider the case $k = 1$. The following result is an assertion on asymptotic properties of the maximum likelihood estimator $\hat{\theta}_n$ of a one-dimensional parameter θ .

THEOREM 5.2.2. *Let Θ be an open interval. Assume that*

- 1) *for all $\theta \in \Theta$ there exist the derivatives $\partial^j \ln f(x; \theta) / \partial \theta^j$, $j = 1, 2, 3$, for μ -almost all x ;*
- 2) *for all $\theta \in \Theta$ there exist nonnegative functions $F_1(x)$, $F_2(x)$, and $H(x)$ depending on x and such that*

$$\left| \frac{\partial^j f(x; \theta)}{\partial \theta^j} \right| \leq F_j(x), \quad j = 1, 2; \quad \left| \frac{\partial^3 \ln f(x; \theta)}{\partial \theta^3} \right| \leq M(x);$$

the functions $F_j(x)$, $j = 1, 2$, are integrable with respect to the measure μ and $E_\theta H(\xi_1) \leq M < \infty$ where the constant M does not depend on θ ;

3) $0 < I(\theta) = E_\theta(\partial \ln f(\xi_1; \theta) / \partial \theta)^2 < \infty$ for all $\theta \in \Theta$.

Then the likelihood equation (5.2.2) has a solution $\hat{\theta}_n = \hat{\theta}_n(\xi^{(n)})$ that converges in probability P_{θ_0} to the true value θ_0 of the parameter; moreover the maximum likelihood estimator $\hat{\theta}_n$ is an asymptotically normal and asymptotically efficient estimator of the parameter θ .

PROOF. Let $\theta_0 \in \Theta$ be the true value of the parameter θ . First we show that there exists a solution of the likelihood equation that converges in probability to θ_0 .

Expanding $\partial L_n(\xi^{(n)}; \theta) / \partial \theta$ into the Taylor series in a neighborhood of the point $\theta = \theta_0$ we get

$$(5.2.5) \quad \begin{aligned} \frac{\partial L_n(\xi^{(n)}; \theta)}{\partial \theta} &= \sum_{i=1}^n \frac{\partial \ln f(\xi_i; \theta)}{\partial \theta} \\ &= \sum_{i=1}^n \left[\left(\frac{\partial \ln f(\xi_i; \theta)}{\partial \theta} \right)_{\theta_0} + (\theta - \theta_0) \left(\frac{\partial^2 \ln f(\xi_i; \theta)}{\partial \theta^2} \right)_{\theta_0} \right. \\ &\quad \left. + \frac{1}{2} \lambda_i (\theta - \theta_0)^2 H(\xi_i) \right] \end{aligned}$$

where $|\lambda_i| \leq 1$ and the symbol $(\cdot)_{\theta_0}$ stands for $(\phi(\theta))_{\theta_0} = \phi(\theta_0)$. Substituting (5.2.5) into the likelihood equation (5.2.2) we get

$$(5.2.6) \quad B_0 + B_1(\theta - \theta_0) + \frac{1}{2} \lambda B_2(\theta - \theta_0)^2 = 0$$

where $|\lambda| \leq 1$ and

$$(5.2.7) \quad \begin{aligned} B_j &= \frac{1}{n} \sum_{i=1}^n \left(\frac{\partial^{j+1} \ln f(\xi_i; \theta)}{\partial \theta^{j+1}} \right)_{\theta_0}, \quad j = 0, 1, \\ B_2 &= \frac{1}{n} \sum_{i=1}^n H(\xi_i). \end{aligned}$$

Now we show that P_{θ_0} approaches 1 as $n \rightarrow \infty$ for an arbitrary $\delta > 0$ where P_{θ_0} is the probability that equation (5.2.6) has a solution on the interval $(\theta_0 - \delta, \theta_0 + \delta)$. To prove this result we study the limit behavior of B_j as $n \rightarrow \infty$.

Assumptions 1) and 2) imply that

$$\int \frac{\partial f(x; \theta)}{\partial \theta} \mu(dx) = \int \frac{\partial^2 f(x; \theta)}{\partial \theta^2} \mu(dx) = 0.$$

Thus

$$(5.2.8) \quad E_{\theta_0} \left(\frac{\partial \ln f(\xi_1; \theta)}{\partial \theta} \right)_{\theta_0} = \int \left(\frac{\partial f(x; \theta)}{\partial \theta} \right)_{\theta_0} \mu(dx) = 0,$$

$$(5.2.9) \quad \begin{aligned} E_{\theta_0} \left(\frac{\partial^2 \ln f(\xi_1; \theta)}{\partial \theta^2} \right)_{\theta_0} &= \int \left(\frac{\partial^2 f(x; \theta)}{\partial \theta^2} - \left(\frac{\partial f(x; \theta)}{\partial \theta} \right)^2 \right)_{\theta_0} f(x; \theta_0) \mu(dx) \\ &= - \int \left(\frac{\partial f(x; \theta)}{\partial \theta} \right)^2_{\theta_0} f(x; \theta_0) \mu(dx) \\ &= -E_{\theta_0} \left(\frac{\partial \ln f(\xi_1; \theta)}{\partial \theta} \right)_{\theta_0}^2 = -I(\theta_0). \end{aligned}$$

According to the law of large numbers we have $B_0 \rightarrow 0$ and $B_1 \rightarrow -I(\theta_0)$ as $n \rightarrow \infty$ in probability P_{θ_0} . Similarly $B_2 \rightarrow E_{\theta_0} H(\xi_1) \leq M < \infty$ as $n \rightarrow \infty$ in probability P_{θ_0} .

Let $\delta > 0$ and $\varepsilon > 0$ be two arbitrary numbers. Then there exists an integer number $n_0 = n_0(\delta, \varepsilon)$ such that

$$\begin{aligned} p_1 &= P_{\theta_0} \{ |B_0| \geq \delta^2 \} < \frac{\varepsilon}{3}, \\ p_2 &= P_{\theta_0} \left\{ B_1 \geq -\frac{1}{2} I(\theta_0) \right\} < \frac{\varepsilon}{3}, \\ p_3 &= P_{\theta_0} \{ B_1 \geq 2M \} \leq \frac{\varepsilon}{3} \end{aligned}$$

for $n > n_0$. Consider events $S = \{ |B_0| < \delta^2, B_1 < \frac{1}{2} I(\theta_0), B_2 < 2M \}$. Then

$$P_{\theta_0}(\Omega \setminus S) \leq p_1 + p_2 + p_3 < \varepsilon,$$

whence $P_{\theta_0}(S) > 1 - \varepsilon$ for all $n > n_0$.

The left-hand side of equality (5.2.6) is equal to $B_0 \pm B_1 \delta + 2^{-1} \lambda B_2 \delta^2$ at the point $\theta = \theta_0 \pm \delta$. If the event S occurs, then $|B_0 + 2^{-1} \lambda B_2 \delta^2| < \delta^2 + |\lambda| M \delta^2 < (M+1)\delta^2$ and $B_1 \delta < -2^{-1} I(\theta_0) \delta$. Thus the sign of $B_0 \pm B_1 \delta + 2^{-1} \lambda B_2 \delta^2$ is defined by the second term if $\delta < I(\theta_0) 2^{-1} (M+1)^{-1}$, whence $\partial L_n(\xi^{(n)}; \theta) / \partial \theta > 0$ for $\theta = \theta_0 - \delta$ and $\partial L_n(\xi^{(n)}; \theta) / \partial \theta < 0$ for $\theta = \theta_0 + \delta$. Assumption 1) implies that $\partial L_n(x; \theta) / \partial \theta$ is a continuous function of $\theta \in \Theta$ for μ -almost all x . Therefore for arbitrary $\delta > 0$ and $\varepsilon > 0$ the likelihood equation (5.2.6) has, with probability greater than $1 - \varepsilon$, a solution belonging to the interval $(\theta_0 - \delta, \theta_0 + \delta)$ if $n > n_0(\delta, \varepsilon)$. This proves the first part of the theorem.

Now we prove that the maximum likelihood estimator $\hat{\theta}_n$ is asymptotically normal and asymptotically efficient.

Let $\hat{\theta}_n = \hat{\theta}_n(\xi^{(n)})$ be a solution of the likelihood equation. It follows from (5.2.6) and (5.2.7) that

$$\hat{\theta}_n - \theta_0 = \frac{B_0}{-B_1 - 2^{-1} \lambda B_2 (\hat{\theta}_n - \theta_0)}.$$

This implies that

$$(5.2.10) \quad \sqrt{nI(\theta_0)}(\hat{\theta}_n - \theta_0) = \frac{\frac{1}{\sqrt{nI(\theta_0)}} \sum_{i=1}^n \left(\frac{\partial \ln f(\xi_i; \theta)}{\partial \theta} \right)_{\theta_0}}{-\frac{B_1}{I(\theta_0)} - \frac{1}{2} \lambda B_2 \frac{\hat{\theta}_n - \theta_0}{I(\theta_0)}}.$$

Thus the denominator on the right-hand side of (5.2.10) converges in probability P_{θ_0} to 1 as $n \rightarrow \infty$. Every term of the sum in the numerator on the right-hand side of (5.2.10) has expectation 0 and variance $I(\theta_0)$ according to equalities (5.2.8) and (5.2.9). Thus the central limit theorem implies that the numerator on the right-hand side of (5.2.10) is asymptotically $\mathcal{N}(0, 1)$ normal. Hence (5.2.10) implies that the random variable $\sqrt{nI(\theta_0)}(\hat{\theta}_n - \theta_0)$ is asymptotically $\mathcal{N}(0, 1)$ normal. Therefore the estimator $\hat{\theta}_n$ is asymptotically efficient. \square

Asymptotic properties of the maximum likelihood estimator $\hat{\theta}_n$ of a multidimensional parameter $\theta = (\theta_1, \dots, \theta_k)$ are listed in the following result.

THEOREM 5.2.3. *Let Θ be an open nondegenerate k -dimensional parallelepiped. Assume that*

- 1) *for all $\theta \in \Theta$ and for μ -almost all x there exist partial derivatives up to third order inclusive of the function $\ln f(x; \theta)$ with respect to θ ;*
- 2) *for all $\theta \in \Theta$*

$$\left| \frac{\partial f(x; \theta)}{\partial \theta_p} \right| \leq F_1(x), \quad \left| \frac{\partial^2 f(x; \theta)}{\partial \theta_p \partial \theta_q} \right| \leq F_2(x),$$

$$\left| \frac{\partial^3 \ln f(x; \theta)}{\partial \theta_p \partial \theta_q \partial \theta_s} \right| \leq H(x),$$

the functions $F_j(x)$ are integrable with respect to the measure μ , and there is a constant M such that $E_{\theta} H(\xi_1) \leq M < \infty$ for all θ ;

- 3) *for all $\theta \in \Theta$ the matrix*

$$B(\theta, \theta_0) = \left\| E_{\theta_0} \frac{\partial \ln f(\xi_1; \theta)}{\partial \theta_p} \frac{\partial \ln f(\xi_1; \theta)}{\partial \theta_q} \right\|$$

is nonsingular and $\det B(\theta, \theta_0) < \infty$.

Then the system of likelihood equations (5.2.2) has a solution that is a consistent, asymptotically $\mathcal{N}(\theta_0, n^{-1}B^{-1}(\theta_0, \theta_0))$ normal, and asymptotically efficient estimator of the parameter θ where θ_0 is the true value of the parameter.

The proof of Theorem 5.2.3 is similar to that of Theorem 5.2.2, thus we omit it.

EXAMPLE 5.2.5. Let $\xi^{(n)} = (\xi_1, \dots, \xi_n)$ be a sample from the normal $\mathcal{N}(\theta_1, \theta_2)$ distribution where $\theta = (\theta_1, \theta_2)$ is an unknown parameter, $-\infty < \theta_1 < \infty$, $\theta_2 > 0$. We learned in Example 5.2.1 that the maximum likelihood estimators of the parameter θ are of the form (5.2.3), and moreover they coincide with the estimators obtained by the method of moments. By Theorem 5.2.3 estimators (5.2.3) are consistent, asymptotically $\mathcal{N}(\theta_0, n^{-1}I^{-1}(\theta_0))$ normal, and asymptotically efficient. Here θ_0 is the true value of the parameter and the Fisher information matrix $I(\theta)$ is of the form

$$I(\theta) = \begin{pmatrix} \theta_2^{-1} & 0 \\ 0 & 2^{-1}\theta_2^{-2} \end{pmatrix}$$

(see Example 3.5.1). Note that the estimator $\hat{\theta}_1$ defined in (5.2.3) is unbiased, while the estimator $\hat{\theta}_2$ is biased. On the other hand, we learned in Example 3.5.1 that the estimator $\hat{\theta}_n = (\hat{\theta}_{1,n}, \hat{\theta}_{2,n})$ is unbiased and asymptotically efficient. Moreover it is easy to see that the estimator $\hat{\theta}_n$ is also asymptotically $\mathcal{N}(\theta_0, n^{-1}I^{-1}(\theta_0))$ normal.

EXAMPLE 5.2.6. Let $\xi^{(n)} = (\xi_1, \dots, \xi_n)$ be a sample from the Gamma distribution, so that the density is

$$f(x; \theta) = \frac{1}{\Gamma(\theta)} x^{\theta-1} e^{-x} I_{(0, \infty)}(x)$$

where $\theta \in \Theta = (0, \infty)$ is an unknown parameter. We showed in Examples 3.4.3 and 5.1.1 that the estimator $\hat{\theta}_n = n^{-1} \sum_{i=1}^n \xi_i$ of the parameter θ obtained by the method of moments is unbiased and consistent but it is not asymptotically efficient whatever the parameter θ is. On the other hand, the maximum likelihood method leads to the equation

$$(5.2.11) \quad \frac{1}{n} \sum_{i=1}^n \ln \xi_i - \frac{d \ln \Gamma(\theta)}{d\theta} = 0$$

and the maximum likelihood estimator $\tilde{\theta}_n$ is a unique positive solution of this equation. According to Theorem 5.2.2 the estimator $\tilde{\theta}_n$ is asymptotically

$$\mathcal{N} \left(\theta_0, \left(n \frac{d^2 \ln \Gamma(\theta)}{d\theta^2} \right)_{\theta_0}^{-1} \right)$$

normal and asymptotically efficient. This can easily be obtained explicitly from equation (5.2.11), since

$$E_{\theta} \ln \xi_1 = \frac{d \ln \Gamma(\theta)}{d\theta}, \quad D_{\theta} \ln \xi_1 = \frac{d^2 \ln \Gamma(\theta)}{d\theta^2}$$

and the random variable $n^{-1} \sum_{i=1}^n \ln \xi_i$ is asymptotically

$$\mathcal{N} \left(\left(\frac{d \ln \Gamma(\theta)}{d\theta} \right)_{\theta_0}, \left(\frac{1}{n} \frac{d^2 \ln \Gamma(\theta)}{d\theta^2} \right)_{\theta_0} \right)$$

normal by the central limit theorem.

Applications of regularity conditions for families of distributions for studying asymptotic properties of maximum likelihood estimators. If the derivatives up to third order of the logarithmic likelihood function $\ln f(x; \theta)$ exist, then Theorems 5.2.2 and 5.2.3 show that the maximum likelihood estimator $\hat{\theta}_n$ is consistent, asymptotically normal, and asymptotically efficient under some extra assumptions that are, in fact, not necessary for these properties.

First we consider sufficient conditions for the consistency of the maximum likelihood estimator $\hat{\theta}_n$. For all sets $A \subset \mathbf{R}^k$ such that $A \cap \Theta \neq \emptyset$ put

$$f_n^*(x; A) = \sup\{f_n(x; \theta); \theta \in \Theta \cap A\}$$

where $x = (x_1, \dots, x_n)$. For $n = 1$ we have $f^*(x; A) = f_1^*(x; A)$.

The following result contains sufficient conditions for the consistency of the maximum likelihood estimator $\hat{\theta}_n$ as $n \rightarrow \infty$.

THEOREM 5.2.4. *Let $\theta_0 \in \Theta$ be the true value of the parameter. Assume that*

- 1) *if $\theta \neq \theta_0$, then $\int |f(x; \theta) - f(x; \theta_0)| \mu(dx) > 0$;*
- 2) *for all x the density $f(x; \theta)$ is a semicontinuous function with respect to θ on the set Θ , that is, for all $\theta' \in \Theta$*

$$\limsup_{h \rightarrow 0} \{f(x; \theta); |\theta - \theta'| < h\} = f(x; \theta');$$

- 3) *for some r*

$$E_{\theta_0} \ln(f_r(\xi^{(r)}; \theta_0)/f_r^*(\xi^{(r)}; H)) > -\infty.$$

If H is a compact subset of Θ containing the point θ_0 , then there exists $\hat{\theta}_n \in H$ such that

$$f_n(x; \hat{\theta}_n) = f_n^*(x; H), \quad x = (x_1, \dots, x_n),$$

and $\hat{\theta}_n \rightarrow \theta_0$ with P_{θ_0} -probability 1 as $n \rightarrow \infty$. Moreover,

- 4) *if additionally*

$$E_{\theta_0} \ln(f_r(\xi^{(r)}; \theta_0)/f_r^*(\xi^{(r)}; \Theta \setminus H)) > 0,$$

then with P_{θ_0} -probability 1 the likelihood function has the global maximum at the point $\hat{\theta}_n$ if n is sufficiently large, that is, $f_n(x; \hat{\theta}_n) = f_n^(x; \Theta)$.*

The proof of Theorem 5.2.4 can be found in [25]. Note that Theorem 5.2.4 claims that the maximum likelihood estimator $\hat{\theta}_n$ approaches θ_0 with P_{θ_0} -probability 1 as $n \rightarrow \infty$. We say in this case that $\hat{\theta}_n$ is a *strongly consistent* estimator of the parameter.

Analyzing the assumptions of Theorem 5.2.4 one can see that the continuity of the function $f(x; \theta)$ with respect to θ is close to being a necessary condition for the consistency of the maximum likelihood estimator $\hat{\theta}_n$.

Now we discuss sufficient conditions for the asymptotic normality and asymptotic efficiency of the maximum likelihood estimator $\hat{\theta}_n$. We consider the case where θ is a one-dimensional parameter.

THEOREM 5.2.5. *Let $\theta_0 \in \Theta$ be the true value of the parameter. Assume that*

- 1) *$f(x; \theta)$ is a measurable function with respect to the pair of variables $(x; \theta)$ and $\int |f(x; \theta) - f(x; \theta')| \mu(dx) > 0$ for all $\theta \neq \theta'$;*
- 2) *for all x the function $f(x; \theta)$ is absolutely continuous with respect to θ and $\int |\partial f(x; \theta)/\partial \theta| d\theta < \infty$ for μ -almost all x ;*
- 3) *$E_{\theta_0} |\partial \ln f(\xi_1; \theta)/\partial \theta|^{2+\delta}$ for some $\delta > 0$ and all $\theta \in \Theta$; the Fisher information $I(\theta) = E_{\theta} (\partial \ln f(\xi_1; \theta)/\partial \theta)^2$ is a continuous function such that $I(\theta) \leq C(1 + |\theta|^p)$ for some $C > 0$ and $p \geq 0$;*
- 4) *$\sup_{\theta} |\theta - \theta_0|^\gamma \int (f(x; \theta)f(x; \theta_0))^{1/2} \mu(dx) < \infty$ for some $\gamma > 0$.*

Then the maximum likelihood estimator $\hat{\theta}_n$ is asymptotically $\mathcal{N}(\theta_0, (nI(\theta_0))^{-1})$ normal and asymptotically efficient as $n \rightarrow \infty$.

The proof of Theorem 5.2.5 can be found in [13].

Note that the assumptions of Theorem 5.2.5 can be weakened (see [13]). Moreover Theorem 5.2.5 can be proved for a multidimensional parameter. We also note that the asymptotic normality and asymptotic efficiency of the maximum likelihood estimator $\hat{\theta}_n$ is proved in [5] under conditions weaker than those in Theorems 5.2.2 and 5.2.3 but stronger than those in Theorem 5.2.5. In particular, it is assumed

in [5] that the function $\ln f(x; \theta)$ is twice continuously differentiable with respect to θ for μ -almost all x . It is also proved in [5] that the maximum likelihood estimator is asymptotically Bayes and asymptotically minimax. More properties of the maximum likelihood estimator can be found in [13].

5.3. Bayes and minimax methods

The Bayes approach. Let ξ be an observation that is a random element assuming values in a measurable space (X, \mathcal{B}) and whose distribution belongs to a family $\mathcal{P} = (P_\theta, \theta \in \Theta)$ where $\theta = (\theta_1, \dots, \theta_k)$ is an unknown parameter and $\Theta \subset \mathbf{R}^k$ is a Borel set. For the sake of simplicity we assume that Θ is an interval for $k = 1$ and Θ is a k -dimensional interval (parallelepiped) in \mathbf{R}^k for $k > 1$. Let $r(y, \theta)$ be a nonnegative loss function, $\theta \in \mathbf{R}^k$, $y \in \mathbf{R}^k$, and let \mathbf{Q} be a probability (a priori) measure on $(\mathbf{R}^k, \mathcal{B}^k)$ such that $\mathbf{Q}(\mathbf{R}^k \setminus \Theta) = 0$. For any estimator $T = T(\xi)$ of the parameter θ we introduce the risk function

$$(5.3.1) \quad R(T; \theta) = E_\theta r(T(\xi), \theta), \quad \theta \in \Theta.$$

Following ideas from Section 3.1 one can study estimators of a general function $g(\theta)$ of the parameter θ . Since a general function $g(\theta)$ can be studied similarly to the particular function $g(\theta) = \theta$, we restrict ourselves to the latter case.

Consider the risk of the estimator $T = T(\xi)$ defined as

$$(5.3.2) \quad R(T) = \int R(T; \theta) \mathbf{Q}(d\theta).$$

An estimator $\tilde{\theta} = \tilde{\theta}(\xi)$ is called a *Bayes estimator* of the parameter θ (with respect to the loss function $r(y, \theta)$ and the a priori measure \mathbf{Q}) if

$$(5.3.3) \quad R(\tilde{\theta}) = \inf_T R(T)$$

where the infimum is taken over all estimators T of the parameter θ .

A posteriori Bayes estimators are also considered in the literature. Moreover the same name "Bayes estimators" is used for them. Usually this does not cause any misunderstanding, since the classes of a priori Bayes estimators and a posteriori Bayes estimators coincide in most cases.

When following the Bayes approach it is natural to treat the parameter θ as a random vector with the distribution \mathbf{Q} and the measure P_y as the conditional distribution of the observation ξ given $\theta = y$, that is, $P_y(A) = P\{\xi \in A / \theta = y\}$, $A \in \mathcal{B}$. In this case the a posteriori measure $\mathbf{Q}_x(B) = P\{\theta \in B / \xi = x\}$, $B \in \mathcal{B}^k$, is well defined. We define a *posteriori risk* $R(T/x)$ of the estimator $T = T(\xi)$ by putting

$$(5.3.4) \quad R(T/x) = E\{r(T(\xi), \theta) / \xi = x\} = \int r(T(x), y) \mathbf{Q}_x(dy).$$

An estimator $\tilde{\theta} = \tilde{\theta}(x)$ is called an *a posteriori Bayes estimator* of the parameter θ (with respect to the loss function $r(y, t)$ and the a priori measure \mathbf{Q}) if the a posteriori risk (5.3.4) attains its minimum at this estimator; more precisely, if

$$(5.3.5) \quad R(\tilde{\theta}/x) = \inf_T R(T/x) \quad (\mu\text{-a.s.}).$$

There are sufficient conditions posed on the a priori measure \mathbf{Q} and loss function $r(y, t)$ under which the infimum in (5.3.5) is attained and the corresponding a posteriori Bayes estimator is unique (see [36]). Under these conditions, the a posteriori Bayes estimator coincides with the Bayes estimator defined by (5.3.3).

Equalities (5.3.3) and (5.3.5) can be used to construct estimators of parameters. It is natural to refer to the method of estimation based on (5.3.3) and (5.3.5) as the *Bayes method*. An advantage of this approach is that the corresponding estimators are optimal in the sense that they minimize the risk of the estimators. Moreover one can freely choose the a priori measure \mathbf{Q} and the loss function $r(y, t)$ to reflect the features of the case under consideration. A disadvantage of this approach is that the corresponding Bayes estimators are not easy to evaluate.

Below we give some examples of the evaluation of Bayes and a posteriori Bayes estimators. These examples also show some problems when following the Bayes approach (more examples of Bayes estimators can be found in Examples 3.1.5 and 3.1.6).

EXAMPLE 5.3.1. Let an observation be a sample $\xi^{(n)} = (\xi_1, \dots, \xi_n)$ from the normal $\mathcal{N}(\theta, 1)$ distribution where θ is an unknown parameter such that

$$-\infty < \theta < \infty.$$

Then $T_n = \sum_{i=1}^n \xi_i$ is a complete sufficient statistic for the parameter θ in this case (see Section 4.2). Let the a priori distribution \mathbf{Q} of the parameter θ be normal $\mathcal{N}(0, \tau^2)$. Our aim is to get the a posteriori Bayes estimator $\tilde{\theta}_n$ of the parameter θ with respect to the loss function

$$r(y, t) = \begin{cases} 0, & \text{if } |y - t| < \delta, \\ 1, & \text{if } |y - t| \geq \delta \end{cases}$$

where $\delta > 0$ is a certain fixed number. Since there exists a complete sufficient estimator T_n , a posteriori distribution depends on $x = (x_1, \dots, x_n)$ through $T_n(x)$ (see Remark 4.1.1 and Corollary 4.1.3). Then an a posteriori Bayes estimator should also be a function of the sufficient statistic T_n , that is, $\tilde{\theta}_n = d(T_n)$.

Given the statistic T_n , an a posteriori distribution of the parameter θ is normal

$$\mathcal{N}\left(\frac{T_n}{n + \tau^{-2}}, \left(n + \frac{1}{\tau^2}\right)^{-1}\right).$$

This result can easily be derived from properties of the normal distribution. Thus a posteriori risk of the estimator $\tilde{\theta}_n = d(T_n)$ is

$$\begin{aligned} R(\tilde{\theta}_n/x) &= 1 - \mathbb{P}\{|d(T_n) - \theta| < \delta/T_n(x)\} \\ &= 1 - \Phi\left(\frac{d(T_n(x)) + \delta - T_n(x)/(n + \tau^{-2})}{(n + \tau^{-2})^{-1/2}}\right) \\ &\quad + \Phi\left(\frac{d(T_n(x)) - \delta - T_n(x)/(n + \tau^{-2})}{(n + \tau^{-2})^{-1/2}}\right). \end{aligned}$$

To minimize $R(\tilde{\theta}_n/x)$ one should choose $d(T_n)$ to maximize

$$\Phi\left(\frac{d(T_n(x)) + \delta - T_n(x)/(n + \tau^{-2})}{(n + \tau^{-2})^{-1/2}}\right) - \Phi\left(\frac{d(T_n(x)) - \delta - T_n(x)/(n + \tau^{-2})}{(n + \tau^{-2})^{-1/2}}\right).$$

Consider the function $f(x) = \Phi(x + \varepsilon - \eta) - \Phi(x - \varepsilon - \eta)$. Differentiating it with respect to x we obtain a sufficient condition for x_0 to be a point of maximum of $f(x)$, namely

$$(5.3.6) \quad \phi(x_0 + \varepsilon - \eta) = \phi(x_0 - \varepsilon - \eta)$$

where $\phi(x)$ is the density of the normal $\mathcal{N}(0, 1)$ distribution. Since $\phi(\varepsilon) = \phi(-\varepsilon)$, we get $x_0 = \eta$. It is easy to check that $x_0 = \eta$ is a unique solution of equation (5.3.6). The second derivative of $f(x)$ at the point $x_0 = \eta$ is equal to $-2\varepsilon\phi(\varepsilon)$. Thus x_0 is the point of maximum of the function $f(x)$. Putting

$$x = \frac{d(T_n)}{(n + \tau^{-2})^{-1/2}} \quad \text{and} \quad \eta = \frac{T_n}{(n + \tau^{-2})^{1/2}}$$

we prove that a unique a posteriori Bayes estimator of the parameter θ is given by

$$\tilde{\theta}_n = \tilde{\theta}_n(x) = \frac{T_n(x)}{n + \tau^{-2}}, \quad x = (x_1, \dots, x_n).$$

Note that

$$\lim_{\tau \rightarrow \infty} \tilde{\theta}_n(x) = \frac{1}{n} \sum_{i=1}^n x_i.$$

The following example shows that there are a posteriori Bayes estimators for which the risk is infinite and thus the evaluation of a Bayes estimator does not make any sense from a practical point of view.

EXAMPLE 5.3.2. Let an observation ξ be a random variable with the distribution belonging to the family of uniform distributions on $(0, |\theta|^{-1})$ where θ is a real number such that $1 \leq |\theta| < \infty$, that is, $f(x; \theta) = |\theta|$ for $0 \leq x \leq |\theta|^{-1}$.

Let an a priori distribution \mathbf{Q} be absolutely continuous with respect to the Lebesgue measure with the density

$$q(t) = \begin{cases} \frac{1}{2}|t|^{-2}, & \text{if } 1 \leq |t| < \infty, \\ 0, & \text{if } |t| < 1. \end{cases}$$

An a posteriori distribution of the parameter θ given $\xi = x$ possesses the density

$$q(t/x) = \begin{cases} \frac{|t|^{-1}}{2 \ln x^{-1}}, & \text{if } x \leq |t|^{-1} \leq 1, \\ 0, & \text{if } |t|^{-1} \notin [x, 1]. \end{cases}$$

Let the loss function $r(y, t)$ be quadratic, that is, $r(y, t) = |y - t|^2$. Then the a posteriori Bayes estimator $\hat{\theta}(x)$ minimizes the a posteriori risk

$$R(T/x) = \int_{\{x \leq |t|^{-1} \leq 1\}} (T - t)^2 \frac{|t|^{-1}}{2 \ln x^{-1}} dt.$$

The only function $T(x)$ minimizing $R(T/x)$ is the a posteriori mean of the parameter θ given $\xi = x$, that is, $T(x) = \mathbf{E}\{\theta/\xi = x\}$. Since the a posteriori density $q(t/x)$ is symmetric with respect to $x = 0$, we get that $\hat{\theta}(x) = 0$ (μ -a.s.) is an a posteriori Bayes estimator. The a posteriori risk of this estimator is given by

$$R(\hat{\theta}/x) = \int_{\{x \leq |t|^{-1} \leq 1\}} \frac{|t|}{2 \ln x^{-1}} dt = \frac{1 - x^2}{2x^2 \ln x^{-1}} < \infty \quad (\mu\text{-a.s.}).$$

However the risk of the a posteriori Bayes estimator $\widehat{\theta}(x) = 0$ is

$$R(\widehat{\theta}) = E(\widehat{\theta} - \theta)^2 = 2 \int_1^\infty t^2 \frac{1}{2} |t|^{-2} dt = \infty.$$

Note that the a posteriori Bayes estimator $\widehat{\theta}(x) = 0$ (μ -a.s.) is of no interest at all, since it does not depend on observations. Moreover the estimator $\widehat{\theta}(x)$ assumes values outside the set of parameters $\Theta = (-\infty, 1] \cup (1, \infty)$.

Below are some concluding remarks.

REMARK 5.3.1. The problem of how to construct Bayes estimators in the case of the loss function $r(y, t) = |y - t|^m$, $m = 1, 2, \dots$, is quite well studied in the literature. It is known in the case $m = 1$ that the median of an a posteriori distribution is the Bayes estimator (and the a posteriori Bayes estimator, as well) of a parameter θ (see Remark 3.2.4 and [9], pp. 178–179). If $m = 2$, then the expectation of an a posteriori distribution is the Bayes estimator (and the a posteriori Bayes estimator, as well) of a parameter θ (see Theorem 2.2.1 concerning the general regression). In the case of a general loss function $r(y, t) = w(|y - t|)$ De Groot and Rao (1963) obtained necessary and sufficient conditions that an estimator is a posteriori Bayes (see [36], Theorem 6.2.2).

REMARK 5.3.2. If the a priori distribution \mathbf{Q} is such that a posteriori risk with respect to the loss function $r(y, t) = \lambda(t)(y - t)^2$ is finite (μ -a.s.) for all estimators $\widehat{\theta}(x)$, then an a posteriori Bayes estimator is given by

$$\widetilde{\theta}(x) = \frac{E(\theta \lambda(\theta) / \xi)}{E(\lambda(\theta) / \xi)} \quad (\mu\text{-a.s.})$$

(see [36]). Here $0 < \lambda(t) < \infty$ for all $t \in \Theta$. If $R(T/x) < \infty$ (μ -a.s.) only for $T = T_0$, then $T_0(x)$ is an a posteriori Bayes estimator. In general, an a posteriori Bayes estimator is unique (see Theorem 6.2.1 in [36]).

REMARK 5.3.3. In Section 3.1 we defined a generalized Bayes estimator as an estimator minimizing the risk if the a priori measure \mathbf{Q} is not a probability measure. Generalized Bayes estimators are sometimes defined as limits of Bayes estimators constructed with respect to a priori probability measures \mathbf{Q}_m as $m \rightarrow \infty$ (see, for example, Theorem 3.2.6). The estimator

$$\theta^*(x) = \frac{\int t f(x; t) \mathbf{Q}(dt)}{\int f(x; t) \mathbf{Q}(dt)}$$

is also called a generalized Bayes estimator where \mathbf{Q} is some σ -finite measure (see [36]). More results on Bayes estimators can be found in [36].

Asymptotic properties of Bayes estimators. Let an observation be a sample $\xi^{(n)} = (\xi_1, \dots, \xi_n)$ from a distribution belonging to a family $\{P_\theta, \theta \in \Theta\}$ dominated by some σ -finite measure μ . Let $f(x; \theta)$ be the density of the measure P_θ with respect to the measure μ . Let $\widehat{\theta}_n$ be the Bayes estimator of the parameter θ with respect to the quadratic loss function $r(y, t) = (y - t)^2$ and the a priori measure \mathbf{Q} possessing the density $q(t)$ with respect to the Lebesgue measure.

The following result describes the asymptotic behavior of $\widehat{\theta}_n$ as $n \rightarrow \infty$ in the case of a one-dimensional parameter θ .

THEOREM 5.3.1. *Let assumptions 1)–4) of Theorem 5.2.5 hold. Let additionally*

- 5) *the function $q(t)$ is continuous in a neighborhood of the point $t = t_0$, $q(t_0) \neq 0$, and $\sup_t q(t)(1 + |t|^{p_0})^{-1} < \infty$ for some $p_0 \geq 0$.*

Then the Bayes estimator $\tilde{\theta}_n$ given $\theta = t_0$ is asymptotically $\mathcal{N}(t_0, (nI(t_0))^{-1})$ normal as $n \rightarrow \infty$.

The proof of Theorem 5.3.1 can be found in [13].

The assumptions of Theorem 5.3.1 can be weakened (see [13]). Moreover this result can be generalized to the case of a multidimensional parameter θ .

The minimax approach. The minimax method allows one to construct estimators, called *minimax*, that minimize the maximum of the risk function. Some necessary and sufficient conditions for estimators to be minimax are given in Section 3.1 (see Theorems 3.1.3 and 3.1.4). In many cases the minimax estimator is a Bayes estimator with respect to the less favorable a priori distribution. In those cases the construction of a minimax estimator is reduced to the construction of an appropriate Bayes estimator.

Some examples of minimax estimators are given in Examples 3.1.5 and 3.1.6. Below we give another example related to the loss function introduced in Example 5.3.1.

EXAMPLE 5.3.3. We learned in Example 5.3.1 that

$$\tilde{\theta}_\tau(x) = \bar{x}_n \left(1 + \frac{1}{n\tau^2} \right)^{-1}$$

where $\bar{x}_n = n^{-1} \sum_{i=1}^n x_i$, $x = (x_1, \dots, x_n)$, is a Bayes estimator with respect to the a priori normal $\mathcal{N}(0, \tau^2)$ distribution and the loss function $r(y, t)$ that equals 0 for $|y - t| < \delta$ and 1 for $|y - t| \geq \delta$, $\delta > 0$. Let us show that \bar{x}_n is a minimax estimator. The risk function of the estimator \bar{x}_n is equal to

$$R(\bar{x}_n, t) = P_t\{|\bar{\xi}_n - t| > \delta\} = 2(1 - \Phi(\delta\sqrt{n})) = \rho^*$$

for all $t \in (-\infty, \infty)$ where $\bar{\xi}_n = n^{-1} \sum_{i=1}^n \xi_i$. It is easy to check that the risk of the estimator $\tilde{\theta}_\tau$ is given by

$$(5.3.7) \quad R(\tilde{\theta}_\tau, t) = 2 - \left\{ \Phi \left(\frac{\delta + t(1 + n\tau^2)^{-1}}{(1 + (n\tau^2)^{-1})^{-1}} \sqrt{n} \right) + \Phi \left(\frac{\delta - t(1 + n\tau^2)^{-1}}{(1 + (n\tau^2)^{-1})^{-1}} \sqrt{n} \right) \right\}.$$

Putting $\tau = \tau_1, \tau_2, \dots$ we get a sequence of a priori $\mathcal{N}(0, \tau_m^2)$ distributions, $m = 1, 2, \dots$. We denote the risk with respect to the a priori $\mathcal{N}(0, \tau^2)$ distribution by

$$(5.3.8) \quad R_\tau(\tilde{\theta}_\tau) = E_\tau R(\tilde{\theta}_\tau; \theta)$$

where E_τ is the expectation with respect to the $\mathcal{N}(0, \tau^2)$ distribution. For the right-hand side of (5.3.8) we apply the estimate $R(\tilde{\theta}_\tau; t) \leq 2$. Then the Lebesgue dominated convergence theorem and (5.3.7) imply

$$\lim_{\tau \rightarrow \infty} R_\tau(\tilde{\theta}_\tau) = 2(1 - \Phi(\delta\sqrt{n})) = \rho^*.$$

Now we apply Theorem 3.1.4 for the sequence of a priori $\mathcal{N}(0, \tau_m)$ distributions and Bayes estimators $\tilde{\theta}_{\tau_m}$ and obtain that \bar{x}_n is the minimax estimator of the parameter θ .

REMARK 5.3.4. Further results on minimax estimators can be found in [36].

5.4. Confidence intervals and regions

The notion of a confidence interval. In the preceding sections we considered the problem of constructing *point* estimators of an unknown parameter or a function of a parameter. Every point estimator is a statistic assuming values in the region of values of the parameter (or of the function of the parameter). It is a useful method in practice to construct an interval or a region from the observation or from the sample. The idea behind this method is that the interval or region mentioned above should contain the parameter with a probability close to 1.

Let ξ be an observation that is a random element assuming values in a measurable space (X, \mathcal{B}) . Let its distribution belong to a family $(P_\theta, \theta \in \Theta)$ where $\theta = (\theta_1, \dots, \theta_k)$ is an unknown parameter such that $\theta \in \Theta \subset \mathbf{R}^k$, $k \geq 1$. First we consider the case of a one-dimensional parameter θ , that is, we consider the case $k = 1$.

Let $T_1 = T_1(\xi)$ and $T_2 = T_2(\xi)$ be two statistics such that $T_1 < T_2$ and let

$$P_\theta\{T_1(\xi) < \theta < T_2(\xi)\} \geq \gamma \quad \text{for all } \theta \in \Theta$$

for a given $\gamma \in (0, 1)$. In this case the interval (T_1, T_2) is called a γ -*confidence interval* or a *confidence interval of level γ* for the parameter θ . The number γ is called a *confidence probability* or a *confidence level*, while T_1 and T_2 are called the *lower* and *upper confidence limits*, respectively.

Constructing a confidence interval by a given statistic. Let $\hat{\theta}$ be an estimator of a parameter θ . It is natural to seek a confidence interval of a level γ in the form of $(\hat{\theta} - \Delta(\gamma, \xi), \hat{\theta} + \Delta(\gamma, \xi))$. However the random variables $\Delta(\gamma, \xi)$ depend, generally speaking, on the unknown parameter θ , since these random variables are found from the equation

$$P_\theta\{\hat{\theta} - \Delta(\gamma, \xi) < \theta < \hat{\theta} + \Delta(\gamma, \xi)\} \geq \gamma \quad \text{for all } \theta \in \Theta.$$

Along with the estimator $\hat{\theta}$ one can use any other statistic T when constructing a confidence interval. Let $G_\theta(y) = P_\theta\{T(\xi) < y\}$ be the distribution function of the statistic T . Assume that $G_\theta(y)$ depends on the parameter θ monotonically. More precisely let

$$(5.4.1) \quad G_{\theta_1}(y) \geq G_{\theta_2}(y) \quad \text{for all } y \text{ and } \theta_1 < \theta_2.$$

If additionally the function $G_\theta(y)$ is continuous with respect to θ , then the equation

$$(5.4.2) \quad G_\theta(y) = \gamma$$

has a solution with respect to θ for every $\gamma \in (0, 1)$. We denote this solution by $b(y, \gamma)$.

THEOREM 5.4.1. *Let $\gamma = \gamma_1 + \gamma_2$. If the distribution function $G_\theta(y)$ of the statistic T is continuous with respect to θ and satisfies condition (5.4.1), then the statistics*

$$T_1 = b(T, 1 - \gamma_2), \quad T_2 = b(T, \gamma_1)$$

are lower and upper limits of a confidence interval of level $1 - \gamma$.

PROOF. The random variable $G_\theta(T(\xi))$ has the uniform distribution with respect to P_θ on the interval $[0, 1]$. Thus

$$P_\theta\{\gamma_1 < G_\theta(T(\xi)) < 1 - \gamma_2\} = 1 - \gamma$$

or, equivalently,

$$P_\theta\{b(T(\xi), 1 - \gamma_2) < \theta < b(T(\xi), \gamma_1)\} = 1 - \gamma. \quad \square$$

The inversion procedure of the function $G_\theta(T)$ used in the proof of Theorem 5.4.1 can be done in two steps. First one inverts the function $G_\theta(y)$ with respect to y , that is, one finds the quantiles of $G_\theta^{-1}(\gamma)$ that are solutions of equation (5.4.2). Then one solves the following equations with respect to θ :

$$G_\theta^{-1}(\gamma_1) = T, \quad G_\theta^{-1}(1 - \gamma_2) = T.$$

Solutions of these equations exist, since the function $G_\theta^{-1}(\gamma)$ is monotone and continuous with respect to θ for all $\gamma \in (0, 1)$.

If the function $G_\theta(y)$ is not continuous with respect to θ , then Theorem 5.4.1 still holds and the above procedure still works. The only difference is that an equality in the definition of quantiles is substituted by the inequality

$$\mathbf{G}_\theta(G_\theta^{-1}(\gamma_1), G_\theta^{-1}(1 - \gamma_2)) \geq 1 - \gamma$$

where \mathbf{G}_θ is the measure on $(-\infty, \infty)$ generated by the distribution of $G_\theta(y)$. Since we assumed continuity in Theorem 5.4.1, the quantiles were evaluated from the corresponding equalities in the proof above.

The problem of finding the most precise estimator also exists in the case of the interval setting. We will solve this problem when studying hypotheses testing.

Constructing confidence intervals for the Bayes approach. Let a parameter θ be random with a priori distribution \mathbf{Q} possessing the density $q(y)$ with respect to some σ -finite measure λ . Assume that a family $(P_\theta, \theta \in \Theta)$ of distributions of the observation ξ is dominated by some σ -finite measure μ . Thus $f(x; \theta)$ is the density of the measure P_θ with respect to the measure μ . In this case there exists an a posteriori distribution of the parameter θ given $\xi = x$. Its density with respect to λ is given by

$$q(y/x) = \frac{f(x; y)q(y)}{\int f(x; t)q(t) \lambda(dt)}.$$

As lower and upper limits for a confidence interval of level $1 - \gamma$ one can take statistics $T_1(x)$ and $T_2(x)$ such that

$$\int_{T_1(x)}^{T_2(x)} q(t/x) \lambda(dt) = 1 - \gamma$$

or

$$\int_{T_1(x)}^{T_2(x)} q(t/x) \lambda(dt) \geq 1 - \gamma$$

depending on the continuity or discontinuity of the function $\int_{-\infty}^t q(u/x) \lambda(du)$ with respect to t . In other words, as statistics T_1 and T_2 one should take a γ_1 -quantile and a $(1 - \gamma_2)$ -quantile, respectively, of an a posteriori distribution for all γ_1 and γ_2 such that $\gamma_1 + \gamma_2 = \gamma$.

In contrast to the non-Bayesian approach, in the relation $T_1 < \theta < T_2$ both T_1 and T_2 , as well as the parameter θ , are now random.

Asymptotically confidence intervals. Assume that an observation

$$\xi^{(n)} = (\xi_1, \dots, \xi_n)$$

is a sample from a distribution belonging to a family $(P_\theta, \theta \in \Theta)$. Let

$$\underline{T}_n = T_1(\gamma, \xi^{(n)}) \quad \text{and} \quad \bar{T}_n = T_2(\gamma, \xi^{(n)})$$

be two statistics such that

$$(5.4.3) \quad \liminf_{n \rightarrow \infty} P_\theta \{ \underline{T}_n < \theta < \bar{T}_n \} \geq \gamma \quad \text{for all } \theta \in \Theta.$$

Then the interval $(\underline{T}_n, \bar{T}_n)$ is called a *confidence interval of level γ* . In fact, now one should speak of a sequence of intervals $(\underline{T}_n, \bar{T}_n)$, $n = 1, 2, \dots$.

In the preceding sections we considered point estimators that in the majority of cases are asymptotically normal. Below we construct asymptotic confidence intervals from the point estimators.

Let $\hat{\theta}_n$ be an asymptotically $\mathcal{N}(\theta, \sigma^2(\theta)/n)$ normal estimator where $\sigma(\theta)$ is a continuous function. Since $\hat{\theta}_n \rightarrow \theta$ in P_θ -probability as $n \rightarrow \infty$, we also have $\sigma(\hat{\theta}_n) \rightarrow \sigma(\theta)$ in P_θ -probability as $n \rightarrow \infty$. This implies that the sequence

$$\frac{(\hat{\theta}_n - \theta)\sqrt{n}}{\sigma(\hat{\theta}_n)}, \quad n = 1, 2, \dots,$$

is asymptotically $\mathcal{N}(0, 1)$ normal. Denote by z_δ a solution of the equation

$$\Phi(z) = 1 - \delta$$

with respect to z , that is, z_δ is a $(1 - \delta)$ -quantile of the distribution $\mathcal{N}(0, 1)$. Here the symbol $\Phi(z)$ stands for the distribution function of the law $\mathcal{N}(0, 1)$. If η is a random variable distributed according to the law $\mathcal{N}(0, 1)$, then $P\{|\eta| < z_\delta\} = 1 - 2\delta$. Let $\beta = z_{\gamma/2}$ for a fixed number $\gamma > 0$. Hence

$$\lim_{n \rightarrow \infty} P_\theta \left\{ \left| \frac{(\hat{\theta}_n - \theta)\sqrt{n}}{\sigma(\hat{\theta}_n)} \right| < \beta \right\} = 1 - \gamma$$

or, in other words,

$$\lim_{n \rightarrow \infty} P_\theta \left\{ \hat{\theta}_n - \frac{\beta\sigma(\hat{\theta}_n)}{\sqrt{n}} < \theta < \hat{\theta}_n + \frac{\beta\sigma(\hat{\theta}_n)}{\sqrt{n}} \right\} = 1 - \gamma.$$

Therefore relation (5.4.3) holds for the random variables

$$(5.4.4) \quad \underline{T}_n = \hat{\theta}_n - \frac{\beta\sigma(\hat{\theta}_n)}{\sqrt{n}}, \quad \bar{T}_n = \hat{\theta}_n + \frac{\beta\sigma(\hat{\theta}_n)}{\sqrt{n}}.$$

Equalities (5.4.4) define lower and upper limits of an asymptotic confidence interval of level $1 - \gamma$.

EXAMPLE 5.4.1. Let $\xi^{(n)} = (\xi_1, \dots, \xi_n)$ be a sample from the Gamma distribution, so that the density is $f(x; \theta) = \theta e^{-\theta x}$, $x > 0$, where $\theta \in \Theta = (0, \infty)$ is an unknown parameter. The random variable $T_n = \sum_{i=1}^n \xi_i$ is a complete sufficient statistic and moreover $E_\theta T_n^{-1} = \theta/(n-1)$. Thus $\hat{\theta}_n = (n-1)T_n^{-1}$ is an unbiased optimal estimator of the parameter θ . Further $D_\theta \hat{\theta}_n = \theta^2/(n-2)$ and therefore $\sigma^2(\theta) = \theta^2$. Therefore the limits defined by (5.4.4) become of the form

$$\underline{T}_n = \hat{\theta}_n \left(1 - \frac{\beta}{\sqrt{n}} \right), \quad \bar{T}_n = \hat{\theta}_n \left(1 + \frac{\beta}{\sqrt{n}} \right).$$

The asymptotic confidence level of the interval $(\underline{T}_n, \bar{T}_n)$ is $1 - \gamma$. One can find a precise confidence level of the interval $(\underline{T}_n, \bar{T}_n)$ for a fixed n by evaluating the probability

$$P_\theta \left\{ 1 - \frac{\beta}{\sqrt{n}} < \frac{\theta T_n}{n-1} < 1 + \frac{\beta}{\sqrt{n}} \right\}$$

which is possible by taking into account that θT_n has the Gamma distribution with the density $x^{n-1} e^{-x}/(n-1)!$, $x > 0$.

The multidimensional case. If a parameter $\theta \in \Theta \subset \mathbf{R}^k$ is multidimensional, that is, $k > 1$, then we consider a confidence region instead of a confidence interval.

A random subset $\Theta^* = \Theta^*(\gamma, \xi)$ of the region of parameters Θ is called a *confidence region of level γ* if

$$P_\theta\{\theta \in \Theta^*\} \geq \gamma \quad \text{for all } \theta \in \Theta.$$

In other words, a confidence region Θ^* of level γ contains the unknown parameter θ with a probability greater than or equal to γ .

If an observation is a sample $\xi^{(n)}$, then a random set $\Theta_n^* = \Theta^*(\gamma, \xi^{(n)}) \subset \Theta$ such that

$$\liminf_{n \rightarrow \infty} P_\theta\{\theta \in \Theta_n^*\} \geq \gamma$$

is called an *asymptotic confidence region of level γ* .

The procedure for constructing confidence regions is the same as in the one-dimensional case.

Confidence intervals for normal distributions. We use the exact distributions of linear and quadratic forms of Gaussian random variables (see Section 1.4) to construct exact confidence intervals for parameters of the normal distribution.

EXAMPLE 5.4.2. Let $\xi^{(n)} = (\xi_1, \dots, \xi_n)$ be a sample from the normal $\mathcal{N}(\theta, \sigma^2)$ distribution where $\theta \in \Theta = (-\infty, \infty)$ is an unknown parameter, while the variance σ^2 is known. Thus $E_\theta \xi_1 = \theta$ and $D_\theta \xi_1 = \sigma^2$ for all $\theta \in \Theta$. Our goal is to

construct a confidence interval of level γ for the parameter θ . We use the estimator

$$\hat{\theta}_n = n^{-1} \sum_{i=1}^n \xi_i$$

of the parameter θ to find a confidence interval $(\hat{\theta}_n - \Delta_{n,\gamma}, \hat{\theta}_n + \Delta_{n,\gamma})$ where $\Delta_{n,\gamma}$ is a solution of the equation

$$P_\theta \left\{ \hat{\theta}_n - \Delta_{n,\gamma} < \theta < \hat{\theta}_n + \Delta_{n,\gamma} \right\} = \gamma.$$

Since the distribution of the estimator $\hat{\theta}_n$ is $\mathcal{N}(\theta, \sigma^2/n)$, the latter equation is equivalent to

$$P_\theta \left\{ |\hat{\theta}_n - \theta| < \Delta_{n,\gamma} \right\} = P_\theta \left\{ \left| \frac{(\hat{\theta}_n - \theta)\sqrt{n}}{\sigma} \right| < \frac{\sqrt{n}\Delta_{n,\gamma}}{\sigma} \right\} = 2\Phi \left(\frac{\sqrt{n}\Delta_{n,\gamma}}{\sigma} \right) - 1 = \gamma.$$

This implies that $\Delta_{n,\gamma} = \sigma t_{(1-\gamma)/2} / \sqrt{n}$ where t_p is a p -quantile of the law $\mathcal{N}(0, 1)$, that is, $\Phi(t_p) = p$. Thus, a confidence interval is of the form

$$\left(\hat{\theta}_n - \frac{\sigma}{\sqrt{n}} t_{(1-\gamma)/2}, \hat{\theta}_n + \frac{\sigma}{\sqrt{n}} t_{(1-\gamma)/2} \right).$$

The procedure described in Theorem 5.4.1 can also be used to construct a confidence interval.

EXAMPLE 5.4.3. Let $\xi^{(n)} = (\xi_1, \dots, \xi_n)$ be a sample from the normal $\mathcal{N}(\alpha, \theta)$ distribution where α is known and $\theta \in (0, \infty)$ is an unknown parameter. Now we use the statistic $T_n = \sum_{i=1}^n (\xi_i - \alpha)^2$ which, as we know, is a sufficient statistic for the parameter θ . It is obvious that the distribution of the random variable T_n/θ is $\chi^2(n)$. Thus there are two numbers $\underline{r}_{n,\gamma}$ and $\bar{r}_{n,\gamma}$ such that

$$P_\theta \left\{ \underline{r}_{n,\gamma} < T_n/\theta < \bar{r}_{n,\gamma} \right\} = \gamma.$$

Note that a solution of the latter equation is not unique. A confidence interval of level γ for the variance θ can be taken as follows:

$$(T_n/\bar{r}_{n,\gamma}, T_n/\underline{r}_{n,\gamma}).$$

Note that numbers $\underline{r}_{n,\gamma}$ and $\bar{r}_{n,\gamma}$ possessing this property are not unique.

EXAMPLE 5.4.4. Let $\xi^{(n)} = (\xi_1, \dots, \xi_n)$ be a sample from the normal $\mathcal{N}(\theta_1, \theta_2)$ distribution where $\theta = (\theta_1, \theta_2)$ is an unknown parameter such that $\theta_1 \in (-\infty, \infty)$ and $\theta_2 \in (0, \infty)$, that is, we assume that both the expectation θ_1 and variance θ_2 are unknown. First we construct a confidence interval of level $\gamma \in (0, 1)$ for the expectation θ_1 . We use the random variable $T_n = (a_1 - \theta_1)m_2^{-1/2}$ where

$$a_1 = \frac{1}{n} \sum_{i=1}^n \xi_i \quad \text{and} \quad m_2 = \frac{1}{n} \sum_{i=1}^n (\xi_i - a_1)^2$$

are the sampling mean and sampling variance, respectively. According to Theorem 1.4.2 the random variable T_n has the Student distribution with $n - 1$ degrees of freedom. Let $c_{n,\gamma}$ be a constant such that

$$P_\theta \left\{ |T_n| < c_{n,\gamma} \right\} = \gamma.$$

Such a constant $c_{n,\gamma}$ exists and is unique. Its approximate value can be found from tables of the Student distribution. Therefore a confidence interval of level γ is of the form $(a_1 - c_{n,\gamma}\sqrt{m_2}, a_1 + c_{n,\gamma}\sqrt{m_2})$.

Now we construct a confidence interval of level $\gamma \in (0, 1)$ for the variance θ_2 . We use the random variable $S_n = nm_2/\theta_2$. According to Theorem 1.4.2 the distribution of the random variable S_n is $\chi^2(n-1)$. Thus one can find two numbers $\underline{r}_{n-1,\gamma}$ and $\bar{r}_{n-1,\gamma}$ (see Example 5.4.3) such that

$$P_\theta \{ \underline{r}_{n-1,\gamma} < nm_2/\theta_2 < \bar{r}_{n-1,\gamma} \} = \gamma.$$

This implies that a confidence interval of level γ for the variance θ_2 can be taken as $(nm_2/\bar{r}_{n-1,\gamma}, nm_2/\underline{r}_{n-1,\gamma})$. It is obvious that the numbers $\underline{r}_{n-1,\gamma}$ and $\bar{r}_{n-1,\gamma}$ satisfying the above equality are not unique and thus confidence intervals also are not unique.

References to Part 1

1. T. W. Anderson, *An introduction to multivariate statistical analysis*, Wiley, New York–London–Sydney, 1958.
2. ———, *The statistical analysis of time series*, Wiley, New York–London–Sydney, 1971.
3. J.-R. Barra, *Notions fondamentales de statistique mathématique*, Dunod, Paris, 1971; English transl., *Mathematical basis of statistics*, Academic Press, New York–London, 1981.
4. P. J. Bickel and K. A. Doksum, *Mathematical statistics*, Holden Day, San Francisco, CA, 1977.
5. A. A. Borovkov, *Mathematical statistics. Estimation of parameters. Testing of hypotheses*, “Nauka”, Moscow, 1984; English transl., *Mathematical statistics*, Gordon & Breach, Amsterdam, 1998.
6. ———, *Mathematical statistics. Supplementary chapters*, “Nauka”, Moscow, 1984; English transl., Gordon & Breach, Amsterdam, 1998.
7. D. R. Cox and P. A. M. Lewis, *The statistical analysis of series of events*, Methuen, London, 1966.
8. D. R. Cox and D. Hinkley, *Theoretical statistics*, Chapman & Hall, London, 1974.
9. H. Cramér, *Mathematical methods of statistics*, Princeton Univ. Press, Princeton, NJ, 1999.
10. H. A. David, *Order statistics*, Wiley, New York–London–Sydney, 1970.
11. I. I. Gikhman, A. V. Skorokhod, and M. I. Yadrenko, *Theory of probability and mathematical statistics*, “Vysshcha shkola”, Kiev, 1988. (Russian)
12. B. V. Gnedenko, *The theory of probability*, “Fizmatgiz”, Moscow, 1961; English transl., Chelsea, New York, 1967.
13. I. A. Ibragimov and R. Z. Khas'minskiĭ, *Statistical estimation. Asymptotic theory*, “Nauka”, Moscow, 1979; English transl., Springer-Verlag, New York–Berlin, 1981.
14. G. I. Ivchenko and Yu. I. Medvedev, *Mathematical statistics*, “Vysshaya shkola”, Moscow, 1984. (Russian)
15. A. M. Kagan, Yu. V. Kagan, and C. R. Rao, *Characterization problems in mathematical statistics*, “Nauka”, Moscow, 1972; English transl., Wiley, New York–London–Sydney, 1973.
16. M. G. Kendall and A. Stuart, *The advanced theory of statistics. Distribution theory*, 3rd edition, vol. 1, Griffin, London, 1961.
17. ———, *The advanced theory of statistics. Inference and relationship*, vol. 2, Griffin, London, 1961.
18. ———, *The advanced theory of statistics. Design and analysis and time-series*, 2nd edition, vol. 3, Griffin, London, 1968.
19. M. V. Kozlov and A. V. Prokhorov, *Introduction to mathematical statistics*, Moscow University Press, Moscow, 1987. (Russian)
20. S. Kullback, *Information theory and statistics*, Dover, New York, 1968.
21. E. L. Lehmann, *Theory of point estimation*, Wiley, New York, 1983.
22. Yu. N. Lin'kov, *Asymptotic statistical methods for stochastic processes*, “Naukova dumka”, Kiev, 1993; English transl., Amer. Math. Soc., Providence, RI, 2001.
23. M. Loève, *Probability theory*, 4th edition, Springer-Verlag, New York–Heidelberg–Berlin, 1977.
24. A. K. Mitropol'skiĭ, *Technique of statistical computations*, “Nauka”, Moscow, 1961. (Russian)
25. E. L. G. Pitman, *Some basic theory for statistical inference*, Chapman & Hall, London, 1979.
26. C. R. Rao, *Statistical inference and its applications*, Wiley, New York–London–Sydney, 1965.
27. Yu. A. Rozanov, *Theory of probability, random processes, and mathematical statistics*, “Nauka”, Moscow, 1985. (Russian)
28. H. Scheffé, *The analysis of variance*, Wiley and Chapman & Hall, New York–London, 1959.

29. B. A. Sevast'yanov, *A course in probability theory and mathematical statistics*, "Nauka", Moscow, 1982. (Russian)
30. A. N. Shiryaev, *Probability*, "Nauka", Moscow, 1989; English transl., Springer-Verlag, New York, 1996.
31. L. Shmetterer, *Introduction to mathematical statistics*, Springer-Verlag, New York, 1973.
32. N. V. Smirnov, *Theory of probability and mathematical statistics. Selected works*, "Nauka", Moscow, 1970. (Russian)
33. N. V. Tutubalin, *Theory of probability*, Moscow University Press, Moscow, 1972. (Russian)
34. B. L. van der Waerden, *Mathematische Statistik*, Springer-Verlag, Berlin-Göttingen-Heidelberg, 1957. (German)
35. S. S. Wilks, *Mathematical statistics*, Wiley, New York-London-Sydney, 1963.
36. S. Zacks, *The theory of statistical inference*, Wiley, New York-London-Sydney, 1971.

Part 2

Preface to Part 2

Part 1 of this book dealt with the estimation of unknown parameters, while Part 2 is devoted to testing statistical hypotheses.

The theory of hypotheses testing appears, in more or less detail, in practically any textbook or monograph on mathematical statistics. We mention here the books by Lehmann [34], and Hájek and Šidák [22] that are entirely devoted to statistical tests, as well as the book by Borovkov and Mogul'skiĭ [10] that is devoted to asymptotic problems in testing statistical hypotheses.

Part 2 begins with an exposition of a general theory of testing (Chapter 1), that is, of problems related to testing statistical hypotheses in the scheme of general statistical experiments according to Ibragimov and Khas'minskiĭ [25], Barra [2], and Soler [49]. First, in Section 1.1, we deal with testing two hypotheses, we study the structure of the set formed by type I and type II error probabilities, and we introduce Neyman-Pearson tests, Bayes tests, and minimax tests. In Section 1.2, the theory of testing a finite number of simple hypotheses is presented and the most powerful tests, Bayes tests, and minimax tests are introduced. Section 1.3 deals with testing composite hypotheses and discusses different approaches to the definition of optimal tests. A relationship between tests and confidence intervals is investigated.

Chapter 2 deals with problems for asymptotically distinguishable families of simple statistical hypotheses in the scheme of general statistical experiments following the books [47] and [37]. A complete group of types of families of statistical hypotheses that can be asymptotically distinguished is introduced and characterization theorems are given, which enables one to determine the type to which a family of hypotheses belongs (Section 2.2). Complete asymptotic testing under the strong law of large numbers (Section 2.3) or under weak convergence (Section 2.4) of the logarithm of the likelihood ratio are presented. Section 2.5 deals with testing contiguous families of hypotheses.

Chapter 3 is devoted to goodness-of-fit tests for independent observations. The Kolmogorov test (Section 3.1), the Pearson test (Section 3.2), and the Smirnov test (Section 3.3) are considered in detail. Section 3.4 focuses on some other well-known goodness-of-fit tests.

Chapter 4 presents elements of sequential analysis applied to the problem of testing statistical hypotheses. Section 4.1 deals with the Bayes theory of sequential testing of, generally speaking, composite hypotheses. Sections 4.2 and 4.3 are devoted to the Wald sequential test for testing two simple hypotheses. Section 4.2 presents the basic properties of the Wald test and Section 4.3 establishes that the Wald test is optimal.

The list of references at the end of the book contains only those references that are directly related to the topics we treat in the book and is by no means a complete list of references on testing statistical hypotheses.

In Part 2 we follow the same system of notational conventions as in Part 1. We also enumerate theorems, lemmas, formulas, etc., in the same way as we did in Part 1.

General Theory of Hypotheses Testing

1.1. Testing two simple hypotheses

Statistical hypotheses and tests. Type I and type II error probabilities of a test. Let ξ be a random element assuming values in a measurable space (X, \mathcal{B}) and let $\mathcal{P} = (P, \tilde{P})$ be a pair of probability measures defined on (X, \mathcal{B}) . Assume that the distribution of the random element ξ is generated by one of the measures of the family \mathcal{P} . The random element ξ is called an *observation*. The problem is to make a decision about the distribution of the random element ξ by the observation $\xi = x$.

Any conjecture about the distribution of an observation ξ is called a *statistical hypothesis* or, simply, a *hypothesis*. If a statistical hypothesis uniquely determines the distribution of an observation, then it is called a *simple hypothesis*. Otherwise it is called a *composite hypothesis*.

Let H and \tilde{H} be two statistical hypotheses that the distribution of an observation ξ corresponds to the measure P and \tilde{P} , respectively. It is clear that the hypotheses H and \tilde{H} are simple. Therefore the problem is to decide by using the observation $\xi = x$ which of the hypotheses H or \tilde{H} is true. In other words, this is a problem of *distinguishing two simple hypotheses* H and \tilde{H} by an observation $\xi = x$.

Any measurable mapping $\delta: (X, \mathcal{B}) \rightarrow ([0, 1], \mathcal{B}([0, 1]))$ where $\mathcal{B}(A)$ is the Borel σ -algebra of the set A is called a *statistical test* for distinguishing hypotheses H and \tilde{H} . We treat $\delta(x)$ as the probability of accepting the hypothesis \tilde{H} given $\xi = x$, while $1 - \delta(x)$ is the probability of accepting the hypothesis H given $\xi = x$. The mapping δ is sometimes called a *decision rule* or a *decision function*. If the function $\delta(x)$ assumes only two values 0 and 1, then it is called a *nonrandomized test*. Otherwise δ is called a *randomized test*.

If a test δ is nonrandomized, then $X = X_0 \cup X_1$ where $X_i = \{x: \delta(x) = i\}$, $i = 0, 1$, and $X_0 \cap X_1 = \emptyset$. In this case the hypothesis H is accepted for $x \in X_0$, while the hypothesis \tilde{H} is accepted for $x \in X_1$. Thus every nonrandomized test is of the form $\delta(x) = I_{X_1}(x)$, $x \in X$, where $I_A(x)$ is the indicator of the set A , that is, $I_A(x) = 1$ for $x \in A$ and $I_A(x) = 0$ for $x \in A^c = X \setminus A$.

Throughout this chapter we write $\delta = \delta(\xi)$. To measure the quality of a test δ we consider the two numbers

$$(1.1.1) \quad \alpha(\delta) = E\delta \quad \text{and} \quad \beta(\delta) = \tilde{E}(1 - \delta)$$

where E and \tilde{E} are expectations with respect to the measures P and \tilde{P} , respectively. If $f(x)$ is a measurable function, then we write $P\{f(\xi) \in A\}$ or $P\{f \in A\}$ and $\tilde{P}\{f(\xi) \in A\}$ or $\tilde{P}\{f \in A\}$ instead of $P\{x: f(x) \in A\}$ and $\tilde{P}\{x: f(x) \in A\}$, respectively. The number $\alpha(\delta)$ is called the *type I error probability* or δ -*level* of the test δ .

Similarly, the number $\beta(\delta)$ is called the *type II error probability* of the test δ . The number $1 - \beta(\delta)$ is called the *power* of the test δ .

It is natural to say that a test δ_1 is better than a test δ_2 if $\alpha(\delta_1) \leq \alpha(\delta_2)$, $\beta(\delta_1) \leq \beta(\delta_2)$, and at least one of these two inequalities is strict. However it is not always possible to compare tests δ_1 and δ_2 in the specified way. In what follows we consider the set \mathfrak{N} of points $(\alpha(\delta), \beta(\delta))$ corresponding to all possible tests δ . It is clear that $\mathfrak{N} \subset [0, 1] \times [0, 1]$. The definition of the set \mathfrak{N} implies that $(\alpha, \beta) \in \mathfrak{N}$ if and only if there is a test δ such that $\alpha(\delta) = \alpha$ and $\beta(\delta) = \beta$.

Properties of the set \mathfrak{N} . First we consider some properties of the set \mathfrak{N} that hold for each pair of measures (P, \tilde{P}) .

LEMMA 1.1.1. *The set \mathfrak{N} is convex.*

PROOF. Let δ_1 and δ_2 be two arbitrary tests. Then $(\alpha(\delta_1), \beta(\delta_1)) \in \mathfrak{N}$ and $(\alpha(\delta_2), \beta(\delta_2)) \in \mathfrak{N}$. Let $0 \leq \lambda \leq 1$ and

$$(1.1.2) \quad \alpha = \lambda\alpha(\delta_1) + (1 - \lambda)\alpha(\delta_2), \quad \beta = \lambda\beta(\delta_1) + (1 - \lambda)\beta(\delta_2).$$

We prove that $(\alpha, \beta) \in \mathfrak{N}$ for all $\lambda \in [0, 1]$. We get from (1.1.1) and (1.1.2) that

$$\alpha = E[\lambda\delta_1 + (1 - \lambda)\delta_2], \quad \beta = \tilde{E}[1 - (\lambda\delta_1 + (1 - \lambda)\delta_2)].$$

This implies that $\alpha = \alpha(\bar{\delta})$ and $\beta = \beta(\bar{\delta})$ where $\bar{\delta} = \lambda\delta_1 + (1 - \lambda)\delta_2$. It is obvious that $\bar{\delta}$ is a test for any $\lambda \in [0, 1]$. Thus $(\alpha, \beta) \in \mathfrak{N}$ for any $\lambda \in [0, 1]$ and therefore \mathfrak{N} is a convex set. \square

LEMMA 1.1.2. *The points $(0, 1)$ and $(1, 0)$ belong to the set \mathfrak{N} .*

PROOF. Let $\delta_0(x) = 0$ for all $x \in X$. Then $\alpha(\delta_0) = 0$ and $\beta(\delta_0) = 1$, whence $(0, 1) \in \mathfrak{N}$. Further let $\delta_1(x) = 1$ for all $x \in X$. Thus $\alpha(\delta_1) = 1$ and $\beta(\delta_1) = 0$, whence $(1, 0) \in \mathfrak{N}$. \square

LEMMA 1.1.3. *The set \mathfrak{N} is symmetric about the point $(1/2, 1/2)$.*

PROOF. It is sufficient to prove that if $(\alpha, \beta) \in \mathfrak{N}$, then $(1 - \alpha, 1 - \beta) \in \mathfrak{N}$. Let δ be a test such that $\alpha(\delta) = \alpha$ and $\beta(\delta) = \beta$. It follows from (1.1.1) that

$$1 - \alpha(\delta) = E(1 - \delta), \quad 1 - \beta(\delta) = \tilde{E}\delta.$$

Putting $\tilde{\delta} = 1 - \delta$, we get

$$\alpha(\tilde{\delta}) = 1 - \alpha(\delta) = 1 - \alpha, \quad \beta(\tilde{\delta}) = 1 - \beta(\delta) = 1 - \beta,$$

that is, $(1 - \alpha, 1 - \beta) \in \mathfrak{N}$. \square

REMARK 1.1.1. Lemmas 1.1.1 and 1.1.2 imply that the diagonal of the square $[0, 1] \times [0, 1]$ joining its corners $(0, 1)$ and $(1, 0)$ belongs to the set \mathfrak{N} . Lemma 1.1.3 implies that the subset of \mathfrak{N} above this diagonal coincides with the image under the central symmetry about the point $(1/2, 1/2)$ of the subset of \mathfrak{N} below the diagonal. Therefore one can derive all the properties for the set \mathfrak{N} from their counterparts for one of the two parts of \mathfrak{N} specified above.

Now we consider other properties of the set \mathfrak{N} that depend on the measures P and \tilde{P} . We need some definitions and results from measure theory that can be found, for example, in [19, 23, 31, 32].

A measure \tilde{P} is called *absolutely continuous with respect to a measure P* if

$$\tilde{P}(A) = 0$$

for all $A \in \mathcal{B}$ such that $P(A) = 0$. We write in this case $\tilde{P} \ll P$. If $\tilde{P} \ll P$ and $P \ll \tilde{P}$, then the measures P and \tilde{P} are called *equivalent*. The equivalence of measures \tilde{P} and P is denoted by $\tilde{P} \sim P$.

LEMMA 1.1.4. *If $P \sim \tilde{P}$, then for all tests δ we have*

$$(1.1.3) \quad \alpha(\delta) = 0 \iff \beta(\delta) = 1,$$

$$(1.1.4) \quad \alpha(\delta) = 1 \iff \beta(\delta) = 0.$$

PROOF. Since $0 \leq \delta(x) \leq 1$, we get

$$\begin{aligned} \alpha(\delta) = 0 &\iff E\delta = 0 \iff P\{x: \delta(x) \neq 0\} = 0 \\ &\iff \tilde{P}\{x: \delta(x) \neq 0\} = 0 \iff \tilde{E}(1 - \delta) = 1 \end{aligned}$$

and (1.1.3) is proved. Relation (1.1.4) is proved similarly. \square

We say that a measure P is *not absolutely continuous with respect to a measure \tilde{P}* (denoted by $P \not\ll \tilde{P}$) if there is a set $C \in \mathcal{B}$ such that $\tilde{P}(C) = 0$ and $P(C) > 0$.

LEMMA 1.1.5. *If $\tilde{P} \ll P$ and $P \not\ll \tilde{P}$, then for all tests δ we have*

$$(1.1.5) \quad \alpha(\delta) = 0 \Rightarrow \beta(\delta) = 1,$$

$$(1.1.6) \quad \alpha(\delta) = 1 \Rightarrow \beta(\delta) = 0.$$

Moreover there are tests δ' and δ'' such that $\beta(\delta') = 0$, $\alpha(\delta') < 1$, $\beta(\delta'') = 1$, and $\alpha(\delta'') > 0$.

PROOF. Since $0 \leq \delta(x) \leq 1$, we have

$$\alpha(\delta) = 0 \Rightarrow P\{x: \delta(x) \neq 0\} = 0 \Rightarrow \tilde{P}\{x: \delta(x) \neq 0\} = 0 \Rightarrow \tilde{E}(1 - \delta) = 1,$$

whence (1.1.5) follows. Relation (1.1.6) is proved similarly.

Further let $C \in \mathcal{B}$ be such that $\tilde{P}(C) = 0$ and $P(C) > 0$. Putting

$$\delta'(x) = I_{X \setminus C}(x)$$

we get

$$\begin{aligned} \alpha(\delta') &= E\delta' = P(X \setminus C) = 1 - P(C) < 1, \\ \beta(\delta') &= \tilde{E}(1 - \delta') = \tilde{P}(C) = 0. \end{aligned}$$

If $\delta''(x) = I_C(x)$, then we get in a similar way that $\alpha(\delta'') = P(C) > 0$ and

$$\beta(\delta'') = \tilde{P}(X \setminus C) = 1. \quad \square$$

LEMMA 1.1.6. *If $P \ll \tilde{P}$ and $\tilde{P} \not\ll P$, then for all tests δ we have*

$$\begin{aligned}\beta(\delta) = 0 &\Rightarrow \alpha(\delta) = 1, \\ \beta(\delta) = 1 &\Rightarrow \alpha(\delta) = 0.\end{aligned}$$

Moreover there are tests δ' and δ'' such that $\alpha(\delta') = 0$, $\beta(\delta') < 1$, $\alpha(\delta'') = 1$, and $\beta(\delta'') > 0$.

The proof is similar to that of Lemma 1.1.5 and thus is omitted.

LEMMA 1.1.7. *If $\tilde{P} \not\ll P$ and $P \not\ll \tilde{P}$, then there are tests δ_1 , δ_2 , δ_3 , and δ_4 such that*

$$\begin{aligned}\alpha(\delta_1) = 0, \quad \beta(\delta_1) < 1, \quad \alpha(\delta_2) < 1, \quad \beta(\delta_2) = 0, \\ \alpha(\delta_3) = 1, \quad \beta(\delta_3) > 0, \quad \alpha(\delta_4) > 0, \quad \beta(\delta_4) = 1.\end{aligned}$$

The proof is similar to that of Lemma 1.1.5 and thus is omitted.

REMARK 1.1.2. To prove Lemmas 1.1.5 and 1.1.6 one can put $\delta'' = 1 - \delta'$ and apply Lemma 1.1.3. Similarly, to prove Lemma 1.1.7 one can put $\delta_3 = 1 - \delta_1$ and $\delta_4 = 1 - \delta_2$ and apply Lemma 1.1.3.

Measures P and \tilde{P} are called *singular* (denoted by $P \perp \tilde{P}$) if there exists $C \in \mathcal{B}$ such that $P(C) = 0$ and $\tilde{P}(X \setminus C) = 0$.

LEMMA 1.1.8. *If $P \perp \tilde{P}$, then $(0, 0) \in \mathfrak{N}$.*

PROOF. Let $C \in \mathcal{B}$ be such that $P(C) = 0$ and $\tilde{P}(C) = 1$. Putting $\delta^0(x) = I_C(x)$ we get

$$\alpha(\delta^0) = E\delta^0 = P(C) = 0, \quad \beta(\delta^0) = \tilde{E}(1 - \delta^0) = \tilde{P}(X \setminus C) = 0,$$

that is, $(0, 0) \in \mathfrak{N}$. □

COROLLARY 1.1.1. *If $P \perp \tilde{P}$, then $\mathfrak{N} = [0, 1] \times [0, 1]$.*

PROOF. According to Lemmas 1.1.8 and 1.1.3 we get $(1, 1) \in \mathfrak{N}$. Since $(0, 1) \in \mathfrak{N}$ and $(1, 0) \in \mathfrak{N}$ by Lemma 1.1.2, we apply Lemma 1.1.1 and obtain

$$\mathfrak{N} = [0, 1] \times [0, 1]. \quad \square$$

We write $P = \tilde{P}$ if $P(A) = \tilde{P}(A)$ for all $A \in \mathcal{B}$. Put

$$(1.1.7) \quad \underline{\mathfrak{N}} = \{(\alpha, \beta) : \beta = 1 - \alpha \text{ for all } \alpha \in [0, 1]\}, \quad \overline{\mathfrak{N}} = [0, 1] \times [0, 1].$$

It is clear that $\underline{\mathfrak{N}}$ is the diagonal of the square $[0, 1] \times [0, 1]$ joining its corners $(0, 1)$ and $(1, 0)$.

COROLLARY 1.1.2. *The following hold:*

$$(1.1.8) \quad P = \tilde{P} \iff \mathfrak{N} = \underline{\mathfrak{N}},$$

$$(1.1.9) \quad P \perp \tilde{P} \iff \mathfrak{N} = \overline{\mathfrak{N}}.$$

PROOF. Let $P = \tilde{P}$. Then for all tests δ

$$\beta(\delta) = \tilde{E}(1 - \delta) = E(1 - \delta) = 1 - \alpha(\delta),$$

that is, $\mathfrak{N} = \underline{\mathfrak{N}}$. Thus the implication \Rightarrow in (1.1.8) is proved.

Let $\mathfrak{N} = \underline{\mathfrak{N}}$. Then $\beta(\delta) = 1 - \alpha(\delta)$ for all tests δ . Consider $\delta(x) = I_A(x)$ for $A \in \mathcal{B}$. Then

$$(1.1.10) \quad \beta(\delta) = \tilde{E}(1 - \delta) = 1 - \tilde{P}(A), \quad \alpha(\delta) = E\delta = P(A).$$

Since $\beta(\delta) = 1 - \alpha(\delta)$ for all tests δ , we obtain from (1.1.10) that $\tilde{P}(A) = P(A)$ for all $A \in \mathcal{B}$. Thus $P = \tilde{P}$. Therefore the implication \Leftarrow in (1.1.8) is also proved.

According to Corollary 1.1.1, relation (1.1.9) follows from the implication \Leftarrow in (1.1.9).

Let $\mathfrak{N} = \overline{\mathfrak{N}}$. Then $(0, 0) \in \mathfrak{N}$ and therefore there is a test δ^0 such that $\alpha(\delta^0) = 0$ and $\beta(\delta^0) = 0$. Thus $P\{x: \delta^0(x) \neq 0\} = 0$ if $\alpha(\delta^0) = 0$. Similarly, the equality $\beta(\delta^0) = 0$ implies that $\tilde{P}\{x: \delta^0(x) \neq 1\} = 0$, that is, $\tilde{P}\{x: \delta^0(x) = 1\} = 1$. Putting $C = \{x: \delta^0(x) \neq 0\}$, we obtain from $\{x: \delta^0(x) = 1\} \subset C$ that $P(C) = 0$ and $\tilde{P}(C) = 1$, that is, $P \perp \tilde{P}$. Thus the implication \Leftarrow in (1.1.9) is also proved. \square

Likelihood ratio and Lebesgue decomposition. Let \mathbf{Q} be some σ -finite measure on (X, \mathcal{B}) dominating the family $\mathcal{P} = (P, \tilde{P})$. This means that $P \ll \mathbf{Q}$ and $\tilde{P} \ll \mathbf{Q}$. It is obvious that such a measure exists. In particular, as the measure \mathbf{Q} one can take $\mathbf{Q} = (P + \tilde{P})/2$. Let $\mathfrak{z}(x) = dP/d\mathbf{Q}(x)$ and $\tilde{\mathfrak{z}}(x) = d\tilde{P}/d\mathbf{Q}(x)$ be the Radon–Nikodym derivatives (densities) of the measures P and \tilde{P} with respect to the measure \mathbf{Q} , respectively. Note that $0 \leq \mathfrak{z}(x) < \infty$ and $0 \leq \tilde{\mathfrak{z}}(x) < \infty$ almost everywhere with respect to the measure \mathbf{Q} . Moreover, $P\{x: \mathfrak{z}(x) = 0\} = 0$ and $\tilde{P}\{x: \tilde{\mathfrak{z}}(x) = 0\} = 0$. We define the *likelihood ratios* as follows:

$$(1.1.11) \quad z(x) = \tilde{\mathfrak{z}}(x)/\mathfrak{z}(x), \quad \tilde{z}(x) = \mathfrak{z}(x)/\tilde{\mathfrak{z}}(x).$$

If we agree that $0 = 0/0$, then $z(x)$ and $\tilde{z}(x)$ in (1.1.11) are well defined. Note that

$$\begin{aligned} P\{x: \tilde{\mathfrak{z}}(x) = 0, \mathfrak{z}(x) = 0\} &= 0, \\ \tilde{P}\{x: \tilde{\mathfrak{z}}(x) = 0, \mathfrak{z}(x) = 0\} &= 0. \end{aligned}$$

The following result provides the Lebesgue decomposition of one of the measures P or \tilde{P} with respect to the other one.

LEMMA 1.1.9. For all sets $A \in \mathcal{B}$

$$(1.1.12) \quad \tilde{P}(A) = \int_A z(x) P(dx) + \tilde{P}(A \cap \{x: \mathfrak{z}(x) = 0\}),$$

$$(1.1.13) \quad P(A) = \int_A \tilde{z}(x) \tilde{P}(dx) + P(A \cap \{x: \tilde{\mathfrak{z}}(x) = 0\})$$

where $z(x)$ and $\tilde{z}(x)$ are defined in (1.1.11).

PROOF. Note that for all $A \in \mathcal{B}$ we have

$$(1.1.14) \quad \tilde{P}(A) = \tilde{P}(A \cap \{\mathfrak{z} > 0\}) + \tilde{P}(A \cap \{\mathfrak{z} = 0\}).$$

Since $P(\mathfrak{z} = 0) = 0$, we get

$$(1.1.15) \quad \begin{aligned} \tilde{P}(A \cap \{\mathfrak{z} > 0\}) &= E_{Q\mathfrak{z}} I(A \cap \{\mathfrak{z} > 0\}) \\ &= E_{Q\mathfrak{z}} \frac{\mathfrak{z}}{\mathfrak{z}} I(A \cap \{\mathfrak{z} > 0\}) = E \frac{\mathfrak{z}}{\mathfrak{z}} I(A \cap \{\mathfrak{z} > 0\}) \\ &= E z I(A) = \int_A z(x) P(dx) \end{aligned}$$

where E_Q is the integral with respect to the measure Q and $I(A)$ is the indicator of the set A (in other words, of the event $\{\xi \in A\}$). Thus $I(A) = 1$ if $x \in A$ or $\xi \in A$, while $I(A) = 0$ if $x \notin A$ or $\xi \notin A$. Decomposition (1.1.12) follows from (1.1.14) and (1.1.15). Decomposition (1.1.13) is proved similarly. \square

Put

$$(1.1.16) \quad \bar{\alpha} = P(\mathfrak{z} > 0), \quad \bar{\beta} = \tilde{P}(\mathfrak{z} > 0).$$

The Lebesgue decompositions (1.1.12) and (1.1.13) imply that

$$(1.1.17) \quad \tilde{P} \ll P \iff \bar{\beta} = 1,$$

$$(1.1.18) \quad P \ll \tilde{P} \iff \bar{\alpha} = 1.$$

It follows from (1.1.11) and (1.1.16) that

$$(1.1.19) \quad \bar{\alpha} = P(\mathfrak{z} > 0) = P(\tilde{z} < \infty),$$

$$(1.1.20) \quad \bar{\beta} = \tilde{P}(\mathfrak{z} > 0) = \tilde{P}(z < \infty).$$

The Lebesgue decompositions yield the following result.

LEMMA 1.1.10. *If η is an arbitrary nonnegative and measurable function defined on (X, \mathcal{B}) , then*

$$(1.1.21) \quad \tilde{E}\eta = E\eta z + \tilde{E}\eta I(\mathfrak{z} = 0),$$

$$(1.1.22) \quad E\eta = \tilde{E}\eta \tilde{z} + E\eta I(\mathfrak{z} = 0).$$

The following result contains more properties of the set \mathfrak{N} .

LEMMA 1.1.11. *For all tests δ*

$$(1.1.23) \quad \beta(\delta) = 0 \Rightarrow \alpha(\delta) \geq \bar{\alpha}.$$

Moreover there exists a test δ' such that $\beta(\delta') = 0$ and $\alpha(\delta') = \bar{\alpha}$. Further, for all tests δ

$$(1.1.24) \quad \alpha(\delta) = 0 \Rightarrow \beta(\delta) \geq \bar{\beta}$$

and there exists a test δ'' such that $\alpha(\delta'') = 0$ and $\beta(\delta'') = \bar{\beta}$.

PROOF. If $\beta(\delta) = 0$, then $\tilde{P}(\delta = 1) = 1$. Taking into account equalities (1.1.22) and $\tilde{P}(\tilde{z} = 0) = 0$, we obtain

$$\begin{aligned} \alpha(\delta) &= E\delta = \tilde{E}\delta\tilde{z} + E\delta I(\tilde{z} = 0) = \tilde{E}\tilde{z} + E\delta I(\tilde{z} = 0) \\ &\geq \tilde{E}\tilde{z} I(\tilde{z} > 0) = E_Q \tilde{z} I(\tilde{z} > 0) = EI(\tilde{z} > 0) = \bar{\alpha}. \end{aligned}$$

Thus the implication (1.1.23) is proved. Putting $\delta' = I(\tilde{z} > 0)$, we get

$$\alpha(\delta') = P(\tilde{z} > 0) = \bar{\alpha}, \quad \beta(\delta') = \tilde{P}(\tilde{z} = 0) = 0.$$

The proof of the implication (1.1.24) is similar and follows from (1.1.21). Putting $\delta'' = I(\tilde{z} = 0)$, we obtain $\alpha(\delta'') = 0$ and $\beta(\delta'') = \bar{\beta}$. \square

We derive the following useful relations from Lemma 1.1.11 and Corollary 1.1.2:

$$(1.1.25) \quad \tilde{P} \perp P \iff \bar{\beta} = 0,$$

$$(1.1.26) \quad \tilde{P} \perp P \iff \bar{\alpha} = 0.$$

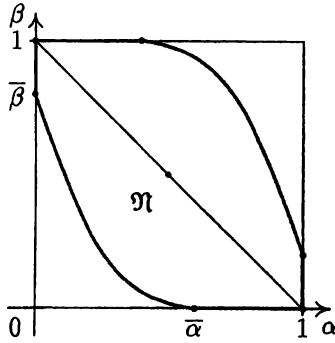


FIGURE 1.1.1

REMARK 1.1.3. The set \mathfrak{N} is shown in Figure 1.1.1. The points $(\bar{\alpha}, 0)$ and $(0, \bar{\beta})$ depicted in Figure 1.1.1 are defined by equalities (1.1.16), (1.1.19), and (1.1.20). Using (1.1.17) and (1.1.25) we get

$$(1.1.27) \quad 0 < \bar{\beta} < 1 \iff \tilde{P} \not\llcorner P, \tilde{P} \not\lrcorner P.$$

Similarly, it follows from (1.1.18) and (1.1.26) that

$$(1.1.28) \quad 0 < \bar{\alpha} < 1 \iff P \not\llcorner \tilde{P}, \tilde{P} \not\lrcorner P.$$

The properties of the set \mathfrak{N} proved above together with equivalences (1.1.27) and (1.1.28) completely describe the set \mathfrak{N} . Note that the tests corresponding to the points $(\bar{\alpha}, 0)$ and $(0, \bar{\beta})$ are defined in Lemma 1.1.11. Thus $(\bar{\alpha}, 0) \in \mathfrak{N}$ and $(0, \bar{\beta}) \in \mathfrak{N}$. Since \mathfrak{N} is convex, the two segments of the straight lines joining the points $(\bar{\alpha}, 0)$ and $(1, 0)$ and $(0, \bar{\beta})$ and $(0, 1)$, respectively, belong to the set \mathfrak{N} . Since \mathfrak{N} is symmetric about the point $(1/2, 1/2)$, two segments of the straight lines joining the points $(0, 1)$ and $(1 - \bar{\alpha}, 0)$ and $(1, 1 - \bar{\beta})$ and $(1, 0)$, respectively, also belong to the set \mathfrak{N} .

The most powerful, Bayes, and minimax tests. Consider the following two classes of tests:

$$(1.1.29) \quad \mathcal{K}_\alpha = \{\delta: \alpha(\delta) = \alpha\}, \quad \mathcal{K}^\alpha = \{\delta: \alpha(\delta) \leq \alpha\}$$

where α is some number of the interval $[0, 1]$. It is clear that $\mathcal{K}_\alpha \subset \mathcal{K}^\alpha$ for all $\alpha \in [0, 1]$. Put

$$(1.1.30) \quad \mathfrak{N}_\alpha = \{(\alpha(\delta), \beta(\delta)): \delta \in \mathcal{K}_\alpha\}, \quad \mathfrak{N}^\alpha = \{(\alpha(\delta), \beta(\delta)): \delta \in \mathcal{K}^\alpha\}.$$

Then $\mathfrak{N}_\alpha \subset \mathfrak{N}^\alpha$.

A test $\delta^{*,\alpha}$ is called the *most powerful test of level α* if

$$(1.1.31) \quad \beta(\delta^{*,\alpha}) = \min\{\beta(\delta): \delta \in \mathcal{K}_\alpha\}$$

(in what follows we show that the minimum in (1.1.31) is attained for all $\alpha \in [0, 1]$, indeed). It is clear that the test $\delta^{*,\alpha}$ has the maximal power $1 - \beta(\delta)$ among tests δ of the class \mathcal{K}_α . Also, it follows from (1.1.29), (1.1.30), and the definition of the set \mathfrak{N} that the test $\delta^{*,\alpha}$ has the maximal power among the tests of the class \mathcal{K}^α . This is an explanation of why we say that $\delta^{*,\alpha}$ is the *most powerful test in the class \mathcal{K}^α* .

The intersection of the straight line $\alpha = \alpha_0$ and the lower bound of the set \mathfrak{N} in Figure 1.1.2 determines the point A whose coordinates are $(\alpha_0, \beta(\delta^{*,\alpha_0}))$ and which corresponds to the most powerful test δ^{*,α_0} in the class \mathcal{K}^{α_0} .

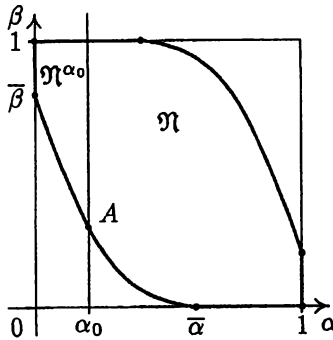


FIGURE 1.1.2

There is a different approach to compare tests. This is the *Bayes approach* based on the assumption that the tests H and \tilde{H} are random events and their probabilities $\pi = P(H)$ and $\tilde{\pi} = P(\tilde{H}) = 1 - \pi$ are known. The probabilities π and $\tilde{\pi}$ are called *a priori error probabilities of tests H and \tilde{H}* . The quality of a test δ is defined as the average of the error probabilities:

$$(1.1.32) \quad e_\pi(\delta) = \pi\alpha(\delta) + (1 - \pi)\beta(\delta).$$

A test δ_π is called a *Bayes test with respect to the a priori distribution $(\pi, 1 - \pi)$* if

$$(1.1.33) \quad e_\pi(\delta_\pi) = \min_\delta e_\pi(\delta)$$

where the minimum is considered with respect to all tests δ . In Figure 1.1.3 the straight line $\pi\alpha + (1 - \pi)\beta = c$ and set \mathfrak{N} have only one common point B and

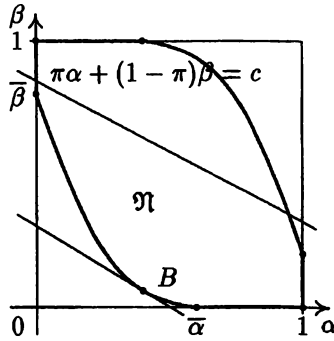


FIGURE 1.1.3

it corresponds to the Bayes test δ_π where c is some constant. It is clear that the Bayes test δ_π is the most powerful one in the class \mathcal{K}^{α_0} for $\alpha_0 = \alpha(\delta_\pi)$.

The following approach to compare tests is called *minimax* and is based on the maximal probability of errors of a test.

A test δ^* is called *minimax* if

$$(1.1.34) \quad \alpha(\delta^*) \vee \beta(\delta^*) = \min_{\delta} (\alpha(\delta) \vee \beta(\delta))$$

where $\alpha \vee \beta$ is the maximum of two numbers α and β , while the minimum in (1.1.34) is considered with respect to all tests δ .

We will discuss Bayes and minimax tests in more detail when considering the problem of testing a finite number of simple tests.

The maximum likelihood test and the Neyman–Pearson fundamental lemma. Consider the test

$$(1.1.35) \quad \delta^{c,\varepsilon} = I(z > c) + \varepsilon I(z = c)$$

where z is the likelihood ratio defined by (1.1.11), and $c \in [0, \infty]$ and $\varepsilon \in [0, 1]$ are the parameters of the test. The test $\delta^{c,\varepsilon}$ defined by (1.1.35) is called the *maximum likelihood test*.

The following result is known as the *Neyman–Pearson fundamental lemma*. It shows that every maximum likelihood test is the most powerful one and, moreover, every most powerful test coincides (in a certain sense) with some maximum likelihood test.

THEOREM 1.1.1.

- 1) For every $\alpha \in (0, \bar{\alpha})$ there exists a maximum likelihood test of level α .
- 2) The maximum likelihood test is the most powerful test of level α .
- 3) If $\delta^{*,\alpha}$ is the most powerful test of level $\alpha \in (0, \bar{\alpha})$, then there exists a constant c such that

$$P(S_c) = \tilde{P}(S_c) = 0$$

where

$$(1.1.36) \quad S_c = \{x: \delta^{*,\alpha}(x) \neq \delta^{c,\varepsilon}(x)\} \cap \{x: z(x) \neq c\}$$

and ε is an arbitrary constant of the interval $[0, 1]$.

PROOF. 1) Consider the function $F(c) = P(z < c)$. Obviously $F(c) = 0$ for $c \leq 0$ and $F(0+) = 1 - P(z > 0) = 1 - P(\mathfrak{z} > 0) = 1 - \bar{\alpha}$. Moreover,

$$F(\infty) = 1 - P(z = \infty) = 1 - P(\mathfrak{z} = 0) = 1.$$

Let $c(\alpha)$ be the minimal solution of the system of inequalities

$$(1.1.37) \quad F(c) \leq 1 - \alpha \leq F(c + 0) = F(c) + P(z = c).$$

Further let $\varepsilon(\alpha) \in [0, 1]$ be such that

$$(1.1.38) \quad 1 - \alpha = F(c(\alpha)) + (1 - \varepsilon(\alpha))P(z = c(\alpha)).$$

If $P(z = c(\alpha)) = 0$, then an arbitrary number of the interval $[0, 1]$ can be taken as $\varepsilon(\alpha)$. Otherwise, if $P(z = c(\alpha)) \neq 0$, then we get from (1.1.38) that

$$(1.1.39) \quad \varepsilon(\alpha) = \frac{F(c(\alpha) + 0) - (1 - \alpha)}{P(z = c(\alpha))}.$$

Equality (1.1.38) implies that the level of the maximum likelihood test $\delta^{c(\alpha), \varepsilon(\alpha)}$ is α .

2) Let $0 < \alpha < \bar{\alpha}$ and let $\delta^{c, \varepsilon}$ be the maximum likelihood test of level α . It is sufficient to show that $\beta(\delta) \geq \beta(\delta^{c, \varepsilon})$ for every test δ of level α . We have

$$(1.1.40) \quad \beta(\delta) - \beta(\delta^{c, \varepsilon}) = \tilde{E}(\delta^{c, \varepsilon} - \delta).$$

Since the levels of the tests $\delta^{c, \varepsilon}$ and δ are equal to α , it follows that

$$(1.1.41) \quad \alpha(\delta^{c, \varepsilon}) - \alpha(\delta) = E(\delta^{c, \varepsilon} - \delta) = 0$$

holds. When proving the first statement of the theorem we showed that

$$F(0+) = 1 - \bar{\alpha}, \quad F(\infty) = 1.$$

This implies that the minimal solution $c = c(\alpha)$ of the system of inequalities (1.1.37) for $\alpha \in (0, \bar{\alpha})$ is such that $0 < c < \infty$. Multiplying (1.1.41) by c and subtracting the result from (1.1.40) we get

$$(1.1.42) \quad \beta(\delta) - \beta(\delta^{c, \varepsilon}) = \tilde{E}(\delta^{c, \varepsilon} - \delta) - cE(\delta^{c, \varepsilon} - \delta).$$

Applying equality (1.1.21) we derive from (1.1.42) that

$$(1.1.43) \quad \beta(\delta) - \beta(\delta^{c, \varepsilon}) = E(\delta^{c, \varepsilon} - \delta)(z - c) + \tilde{E}(\delta^{c, \varepsilon} - \delta)I(\mathfrak{z} = 0).$$

Since $(\mathfrak{z} = 0) \subset (z = 0) \cup (z = \infty)$, we have $(z = c, \mathfrak{z} = 0) = \emptyset$ and

$$\tilde{P}(z < c, \mathfrak{z} = 0) = \tilde{P}(\mathfrak{z} = 0, \mathfrak{z} = 0) = 0.$$

Thus

$$(1.1.44) \quad \begin{aligned} & \tilde{E}(\delta^{c, \varepsilon} - \delta)I(\mathfrak{z} = 0) \\ &= \tilde{E}(1 - \delta)I(z > c, \mathfrak{z} = 0) \\ & \quad + \tilde{E}(\varepsilon - \delta)I(z = c, \mathfrak{z} = 0) - \tilde{E}I(z < c, \mathfrak{z} = 0) \\ &= \tilde{E}(1 - \delta)I(z > c, \mathfrak{z} = 0) \geq 0. \end{aligned}$$

Taking into account (1.1.44) we obtain from (1.1.43) that

$$(1.1.45) \quad \begin{aligned} \beta(\delta) - \beta(\delta^{c,\varepsilon}) &\geq E(\delta^{c,\varepsilon} - \delta)(z - c) \\ &= E(1 - \delta)(z - c)I(z > c) + E(-\delta)(z - c)I(z < c) \geq 0. \end{aligned}$$

Therefore $\beta(\delta) \geq \beta(\delta^{c,\varepsilon})$.

3) Let $\delta^{*,\alpha}$ be the most powerful test of level $\alpha \in (0, \bar{\alpha})$ and let $\delta^{c,\varepsilon}$ be the likelihood ratio test of level α . Since $\beta(\delta^{*,\alpha}) - \beta(\delta^{c,\varepsilon}) \leq 0$ and $\delta^{*,\alpha}$ is the most powerful test, we obtain that

$$(1.1.46) \quad E(\delta^{c,\varepsilon} - \delta^{*,\alpha})(z - c)I(S_c) + \tilde{E}(\delta^{c,\varepsilon} - \delta^{*,\alpha})I(S_c \cap (\mathfrak{J} = 0)) = 0$$

in view of relations (1.1.43) and (1.1.44) for $\delta = \delta^{*,\alpha}$ where S_c is the set defined by (1.1.36). Note that $(\delta^{c,\varepsilon} - \delta^{*,\alpha})(z - c) > 0$ on the set S_c and $\delta^{c,\varepsilon} - \delta^{*,\alpha} > 0$ on the set $S_c \cap (\mathfrak{J} = 0)$. Thus relation (1.1.46) implies that $P(S_c) = 0$ and

$$\tilde{P}(S_c \cap (\mathfrak{J} = 0)) = 0.$$

Lebesgue decomposition yields

$$\tilde{P}(S_c) = \int_{S_c} z dP + \tilde{P}(S_c \cap (\mathfrak{J} = 0)) = 0,$$

whence the third statement of Theorem 1.1.1 follows. \square

Combining Lemmas 1.1.3 and 1.1.11 with Remark 1.1.3 and Theorem 1.1.1 we obtain the following result.

LEMMA 1.1.12. *The set \mathfrak{N} is closed.*

REMARK 1.1.4. It is easy to see that the level of the likelihood ratio test $\delta^{\infty,\varepsilon}$ is $\alpha = 0$ for all $\varepsilon \in [0, 1]$. Indeed, according to relations (1.1.1) and (1.1.35) we have for all $\varepsilon \in [0, 1]$

$$\alpha(\delta^{\infty,\varepsilon}) = P(z > \infty) + \varepsilon P(z = \infty) = \varepsilon P(\mathfrak{J} = 0) = 0.$$

Thus statement 1) of Theorem 1.1.1 holds for $\alpha = 0$, too. Further, as can be seen from the proof of Lemma 1.1.11, $\delta^{*,0} = I(\mathfrak{J} = 0)$ is the most powerful test of level $\alpha = 0$ and that $\beta(\delta^{*,0}) = \bar{\beta}$. Taking into account (1.1.20) we obtain from (1.1.1) and (1.1.35) that

$$\beta(\delta^{\infty,\varepsilon}) = P(z < \infty) + (1 - \varepsilon)P(z = \infty) = \bar{\beta} + (1 - \varepsilon)(1 - \bar{\beta}) \geq \bar{\beta}.$$

This implies that the test $\delta^{\infty,\varepsilon}$ is the most powerful only for $\varepsilon = 1$. This means that statement 2) of Theorem 1.1.1 does not hold in general. As we proved above it only holds for $\varepsilon = 1$. Note that $\delta^{\infty,1} = I(z = \infty)$. Finally, the set S_c defined by (1.1.36) is of the form

$$S_{\infty} = (\delta^{*,0} \neq \delta^{\infty,\varepsilon}) \cap (z \neq \infty) = (\mathfrak{J} = 0) \cap (z \neq \infty) = (\mathfrak{J} = 0, \bar{\mathfrak{J}} = 0)$$

for all $\varepsilon \in [0, 1]$. Moreover $P(S_{\infty}) = \tilde{P}(S_{\infty}) = 0$. Therefore statement 3) of Theorem 1.1.1 holds for $\alpha = 0$, too.

REMARK 1.1.5. According to equality (1.1.19) the level of the likelihood ratio test $\delta^{0,\varepsilon}$ is

$$\alpha(\delta^{0,\varepsilon}) = P(z > 0) + \varepsilon P(z = 0) = \bar{\alpha} + \varepsilon(1 - \bar{\alpha})$$

for $\varepsilon \in [0, 1]$. This implies that if $\bar{\alpha} < 1$, then $\alpha(\delta^{0,\varepsilon}) = \alpha$ for $\alpha \in [\bar{\alpha}, 1]$ and $\varepsilon = (\alpha - \bar{\alpha}) / (1 - \bar{\alpha})$. On the other hand, if $\alpha = 1$, then $\alpha(\delta^{0,\varepsilon}) = 1$ for all $\varepsilon \in [0, 1]$. This shows that statement 1) of Theorem 1.1.1 holds for $\alpha \in [\bar{\alpha}, 1]$, too. Further, we have

$$\beta(\delta^{0,\varepsilon}) = \tilde{P}(z < 0) + (1 - \varepsilon)\tilde{P}(z = 0) = (1 - \varepsilon)\tilde{P}(\tilde{z} = 0, \tilde{z} \geq 0) = 0$$

for all $\varepsilon \in [0, 1]$, that is, the likelihood ratio test $\delta^{0,\varepsilon}$ of an arbitrary level $\alpha \in [\bar{\alpha}, 1]$ is the most powerful. Thus statement 2) of Theorem 1.1.1 holds for $\alpha \in [\bar{\alpha}, 1]$, too. The set S_c defined by (1.1.36) is of the form

$$S_0 = (\delta^{*,\alpha} \neq \delta^{0,\varepsilon}) \cap (z > 0) = (\delta^{*,\alpha} \neq 1) \cap (z > 0)$$

for $c = 0$ and arbitrary $\varepsilon \in [0, 1]$ where $\delta^{*,\alpha}$ is the most powerful test of level $\alpha \in [\bar{\alpha}, 1]$. It can be seen from the proof of Lemma 1.1.11 that the most powerful test of level $\alpha = \bar{\alpha}$ is given by $\delta^{*,\bar{\alpha}} = I(\tilde{z} > 0)$. Moreover $\delta^{*,1} \equiv 1$ (see the proof of Lemma 1.1.2). Thus the most powerful test of level α is

$$\delta^{*,\alpha} = \frac{1 - \alpha}{1 - \bar{\alpha}} I(\tilde{z} > 0) + \frac{\alpha - \bar{\alpha}}{1 - \bar{\alpha}} = I(\tilde{z} > 0) + \frac{\alpha - \bar{\alpha}}{1 - \bar{\alpha}} I(\tilde{z} = 0).$$

Note that $(\delta^{*,\alpha} \neq 1) = \emptyset$ for $\alpha = 1$ and $(\delta^{*,\alpha} \neq 1) = (\tilde{z} = 0)$ for $\alpha < 1$. Since $(z > 0) = (\tilde{z} > 0)$, it holds that $S_0 = \emptyset$, whence $P(S_0) = \tilde{P}(S_0) = 0$. Therefore statement 3) of Theorem 1.1.1 holds for $\alpha \in [\bar{\alpha}, 1]$, too, that is, all the statements of Theorem 1.1.1 hold for $\alpha \in [\bar{\alpha}, 1]$.

Neyman–Pearson test. The fundamental Neyman–Pearson lemma (Theorem 1.1.1) and Remarks 1.1.4 and 1.1.5 imply that for any $\alpha \in [0, 1]$ there exists a likelihood ratio test $\delta^{c(\alpha),\varepsilon(\alpha)}$ of level α where $(c(\alpha), \varepsilon(\alpha))$ is some solution of the equation $\alpha(\delta^{c,\varepsilon}) = \alpha$ with respect to (c, ε) . Moreover $\delta^{c(\alpha),\varepsilon(\alpha)}$ is the most powerful test in the class \mathcal{K}_α if $\varepsilon(0) = 1$ and $\alpha = 0$. The test $\delta^{c(\alpha),\varepsilon(\alpha)}$ for $\varepsilon(0) = 1$ is called the *Neyman–Pearson test of level α* for distinguishing the hypotheses H and \tilde{H} . In what follows we denote this test by $\delta^{+,\alpha}$. One can see from the proof of Theorem 1.1.1 and Remarks 1.1.4 and 1.1.5 that the functions $c(\alpha)$ and $\varepsilon(\alpha)$ can be taken of the form

$$(1.1.47) \quad c(\alpha) = \begin{cases} \infty, & \alpha = 0, \\ \bar{c}(\alpha), & 0 < \alpha < \bar{\alpha}, \\ 0, & \bar{\alpha} \leq \alpha \leq 1, \end{cases} \quad \varepsilon(\alpha) = \begin{cases} 1, & \alpha = 0, \\ \bar{\varepsilon}(\alpha), & 0 < \alpha < \bar{\alpha}, \\ \tilde{\varepsilon}(\alpha), & \bar{\alpha} \leq \alpha \leq 1, \end{cases}$$

where $\bar{c}(\alpha)$ is the minimal number c such that

$$(1.1.48) \quad P(z > c) \leq \alpha \leq P(z \geq c),$$

$$(1.1.49) \quad \bar{\varepsilon}(\alpha) = \frac{\alpha - P(z > \bar{c}(\alpha))}{P(z = \bar{c}(\alpha))}, \quad \tilde{\varepsilon}(\alpha) = \frac{\alpha - \bar{\alpha}}{1 - \bar{\alpha}}.$$

If $P(z = \bar{c}(\alpha)) = 0$, then $P(z > \bar{c}(\alpha)) = \alpha$ which leads to an expression $\bar{\varepsilon}(\alpha) = 0/0$. In this case an arbitrary number of the interval $[0, 1]$ can be taken as $\bar{\varepsilon}(\alpha)$. If $\bar{\alpha} = 1$, then $[\bar{\alpha}, 1] = \{1\}$ and this also results in an expression $\tilde{\varepsilon}(1) = 0/0$. In this case an

arbitrary number of the interval $[0, 1]$ can be taken as $\tilde{\varepsilon}(1)$. The definition of the test $\delta^{+, \alpha}$ and equalities (1.1.47) imply

$$(1.1.50) \quad \beta(\delta^{+, \alpha}) = \begin{cases} \bar{\beta}, & \alpha = 0, \\ \tilde{P}(z < \bar{c}(\alpha)) + (1 - \bar{\varepsilon}(\alpha))\tilde{P}(z = \bar{c}(\alpha)), & 0 < \alpha < \bar{\alpha}, \\ 0, & \bar{\alpha} \leq \alpha \leq 1. \end{cases}$$

It is clear that the function $\beta(\delta^{+, \alpha})$ determines the lower boundary of the set \mathfrak{N} .

EXAMPLE 1.1.1. Let an observation ξ be a Gaussian random variable with the normal $\mathcal{N}(a, 1)$ distribution under the hypothesis H and let its distribution be $\mathcal{N}(\tilde{a}, 1)$ under the hypothesis \tilde{H} . Then the measures P and \tilde{P} corresponding to the distribution of the observation ξ under the hypotheses H and \tilde{H} , respectively, are absolutely continuous and

$$z(x) = \frac{d\tilde{P}}{dP}(x) = \exp\left((\tilde{a} - a)x + \frac{a^2 - \tilde{a}^2}{2}\right).$$

It is obvious in this case that $\bar{\alpha} = 1$ and $\bar{\beta} = 1$. Moreover, the random variable $z = z(\xi)$ has a continuous distribution for both hypotheses H and \tilde{H} . Thus equalities (1.1.37) and (1.1.48) defining the constant $c(\alpha)$ for $0 < \alpha < 1$ become $P(z < c) = 1 - \alpha$. For the sake of definiteness let $\tilde{a} > a$. Then

$$P(z < c) = P\left(\xi < \frac{\ln c}{\tilde{a} - a} + \frac{\tilde{a} + a}{2}\right) = \Phi\left(\frac{\ln c}{\tilde{a} - a} + \frac{\tilde{a} - a}{2}\right) = 1 - \alpha$$

where $\Phi(x)$ is the distribution function of the normal $\mathcal{N}(0, 1)$ law. This implies that

$$\frac{\ln c(\alpha)}{\tilde{a} - a} + \frac{\tilde{a} - a}{2} = t_{1-\alpha}$$

where t_p is the p -quantile of the law $\mathcal{N}(0, 1)$, that is, $\Phi(t_p) = p$. Thus we have for all $\alpha \in (0, 1)$ that

$$c(\alpha) = \exp\left((\tilde{a} - a)t_{1-\alpha} - \frac{(\tilde{a} - a)^2}{2}\right).$$

The number $\varepsilon(\alpha)$ can be chosen arbitrarily from the interval $[0, 1]$. Taking into account (1.1.50) we get

$$\beta(\delta^{+, \alpha}) = \tilde{P}(z < c(\alpha)) = \tilde{P}(\xi < t_{1-\alpha} + a) = \Phi(t_{1-\alpha} - \tilde{a} + a).$$

Note also that the Neyman–Pearson test of level α can be represented in the form $\delta^{+, \alpha} = I(z(x) > c(\alpha)) = I(x > t_{1-\alpha} + a)$.

EXAMPLE 1.1.2. Let an observation ξ have the normal $\mathcal{N}(0, 1)$ distribution under the hypothesis H and the exponential distribution with the density

$$\tilde{f}(x) = e^{-x}I_{(0, \infty)}(x)$$

with respect to the Lebesgue measure under the hypothesis \tilde{H} . Then $\tilde{P} \ll P$, $P \not\ll \tilde{P}$, and moreover $z(x) = d\tilde{P}/dP(x)$ where $z(x) = 0$ for $x \leq 0$, while

$$z(x) = \sqrt{\frac{\pi}{2}} \exp\left(\frac{(x-1)^2}{2}\right)$$

for $x > 0$. In this case $\bar{\beta} = 1$ and $\bar{\alpha} = P(\tilde{f}(x) > 0) = P(x > 0) = 1/2$. The random variable $z = z(\xi)$ is continuous under the hypothesis H . Thus the constant $c(\alpha)$ for $\alpha \in (0, \bar{\alpha}) = (0, 1/2)$ can be determined from the equation $P(z < c) = 1 - \alpha$ or, equivalently, from

$$\Phi \left(1 + \sqrt{\ln \frac{ec^2}{2\pi}} \right) - \Phi \left(1 - \sqrt{\ln \frac{ec^2}{2\pi}} \right) = 1 - \alpha.$$

According to Remark 1.1.5 the Neyman–Pearson test of level α can be represented as

$$\delta^{+, \alpha} = I(z(x) > c(\alpha)) = I \left((x - 1)^2 > \ln \frac{ec^2(\alpha)}{2\pi} \right)$$

for $\alpha \in [0, \bar{\alpha})$ and as

$$\delta^{+, \alpha} = I(x > 0) + (2\alpha - 1)I(x \leq 0)$$

for $\alpha \in [\bar{\alpha}, 1]$, since $\bar{\alpha} = 1/2$. Thus $\varepsilon(\alpha)$ determined by equalities (1.1.47) and (1.1.49) for $\alpha \in [\bar{\alpha}, 1]$ is equal to $\tilde{\varepsilon}(\alpha) = 2\alpha - 1$. The type II error probability of the test $\delta^{+, \alpha}$ is given by relation (1.1.50). Moreover

$$\beta(\delta^{+, \alpha}) = \tilde{P}\{|\xi - 1| < C(\alpha)\} = e^{-(1-C(\alpha))\nu_0} - e^{-1-C(\alpha)}$$

for $\alpha \in (0, \bar{\alpha})$ where

$$C(\alpha) = \sqrt{\ln \frac{ec^2(\alpha)}{2\pi}}.$$

If one interchanges the hypotheses H and \tilde{H} in Example 1.1.2, then $P \ll \tilde{P}$ and $\tilde{P} \not\ll P$. In this case $\bar{\alpha} = 1$ and $\bar{\beta} < 1$. Details are left to the reader.

EXAMPLE 1.1.3. Let an observation ξ assume two values 1 and 0 with probabilities p and $q = 1 - p$ under the hypothesis H and with probabilities \tilde{p} and $\tilde{q} = 1 - \tilde{p}$ under the hypothesis \tilde{H} . Then $P \sim \tilde{P}$ and the likelihood ratio $z(x) = d\tilde{P}/dP(x)$ is given by

$$z(x) = \left(\frac{\tilde{p}}{p} \right)^x \left(\frac{\tilde{q}}{q} \right)^{1-x}, \quad x = 0, 1.$$

Thus the random variable $z = z(\xi)$ assumes two values \tilde{p}/p and \tilde{q}/q with probabilities p and q under the hypothesis H and with probabilities \tilde{p} and \tilde{q} under the hypothesis \tilde{H} . For the sake of definiteness let $\tilde{p} > p$. Then $\tilde{q}/q < 1 < \tilde{p}/p$. Since $P \sim \tilde{P}$, we get $\bar{\alpha} = \bar{\beta} = 1$. Solving equation (1.1.48) and evaluating $c(\alpha)$ and $\varepsilon(\alpha)$ we obtain

$$c(\alpha) = \begin{cases} 0, & \alpha = 1, \\ \tilde{q}/q, & p \leq \alpha < 1, \\ \tilde{p}/p, & 0 \leq \alpha < p, \end{cases}$$

and

$$\varepsilon(\alpha) = \begin{cases} 1, & \alpha = 1, \\ (\alpha - p)/q, & p \leq \alpha < 1, \\ \alpha/p, & 0 \leq \alpha < p. \end{cases}$$

This together with (1.1.50) implies

$$\beta(\delta^{+, \alpha}) = \begin{cases} \tilde{q} + \frac{\tilde{p}}{p}(p - \alpha), & 0 \leq \alpha < p, \\ \frac{\tilde{q}}{q}(1 - \alpha), & p \leq \alpha \leq 1. \end{cases}$$

The function $\beta(\delta^{+, \alpha})$ determines the lower boundary of the set \mathfrak{N} for $0 \leq \alpha \leq 1$. The set \mathfrak{N} is shown in Figure 1.1.4.

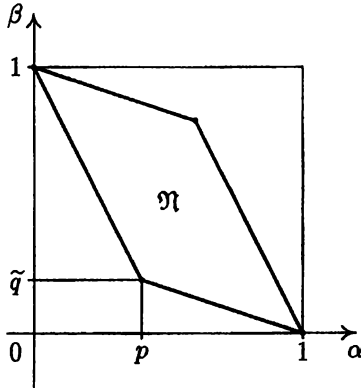


FIGURE 1.1.4

1.2. Distinguishing a finite number of simple hypotheses

Setting of the problem. The most powerful tests. Let (Ω, \mathcal{F}, P) be the main probability space, let ξ be an observation that is a measurable mapping of the space (Ω, \mathcal{F}) into some measurable space (X, \mathcal{B}) , and let $\mathcal{P} = \{P_1, P_2, \dots, P_N\}$, $1 < N < \infty$, be a family of probability measures defined on the space (X, \mathcal{B}) . We assume that the distribution of the observation ξ is generated by some measure of the family \mathcal{P} .

Let $H_j = \{\theta = j\}$ be the hypothesis that the distribution of the observation ξ is generated by the measure P_j . We write in this case

$$P_j(A) = P_j\{\xi \in A\} = P\{\xi \in A/H_j\}, \quad A \in \mathcal{B},$$

and say that $P_j(A)$ is the probability of the event $\{\xi \in A\}$ under the hypothesis $H_j = \{\theta = j\}$. The parameter θ assumes values in the set $\Theta = \{1, 2, \dots, N\}$ and is the index of the measure of the family \mathcal{P} generating the distribution of the observation ξ . Thus we deal with N simple hypotheses H_1, H_2, \dots, H_N . Given an observation ξ , the problem is to decide which hypothesis of the set of hypotheses H_1, H_2, \dots, H_N is true.

Any measurable mapping $\delta: (X, \mathcal{B}) \rightarrow \Theta$ is called a *statistical test for distinguishing N hypotheses H_1, H_2, \dots, H_N* by an observation ξ . The equality $\delta(x) = j$ means that the hypothesis H_j is accepted if $\xi = x$ (that is, $\theta = j$ in the parametric setting). A mapping δ is sometimes called a *decision rule* or a *decision function*. Every test δ uniquely determines a partition of the space (and vice versa) X into N disjoint measurable sets $X_j \in \mathcal{B}$, $j = 1, 2, \dots, N$, such that $X_j = \{x: \delta(x) = j\}$,

$j = 1, 2, \dots, N$, $\bigcup_{j=1}^N X_j = X$. The test δ corresponding to a partition is closely related to the problem of estimation of an unknown parameter θ , namely $\delta(x)$ is an estimator of an unknown parameter θ if $\xi = x$. A test δ defined in this way is a *nonrandomized test* (see also Section 1.1). Any random variable $\delta = \delta(\xi)$ is called a test.

We consider below randomized tests defined as follows. Every measurable mapping $\delta: (X \times \Omega, \mathcal{B} \times \mathcal{F}) \rightarrow \Theta$ is called a *statistical test for distinguishing N hypotheses* H_1, H_2, \dots, H_N by an observation ξ . Given $\xi = x$, the hypothesis H_j is accepted if the random variable $\delta(x, \omega)$ is equal to j . If the function $\delta(x, \omega)$ does not depend on the variable ω , then the test δ is nonrandomized. Otherwise a test δ is called *randomized*. In general, every random variable $\delta = \delta(\xi(\omega), \omega)$, $\omega \in \Omega$, is called a *statistical test*.

Consider a family of functions $q^\delta(x) = (q_1^\delta(x), q_2^\delta(x), \dots, q_N^\delta(x))$, $x \in X$, such that $q_j^\delta(x) = P\{\delta = j/\xi = x\}$ is the conditional probability that the hypothesis H_j is accepted under the test δ given $\xi = x$ (that is, q_j^δ is the conditional probability of the event $\{\delta = j\}$ given $\xi = x$). It is clear that a test δ is uniquely determined by the family of conditional probabilities

$$q^\delta(x) = (q_1^\delta(x), q_2^\delta(x), \dots, q_N^\delta(x)), \quad x \in X.$$

Sometimes this family is called a *statistical test* (see Section 1.1). Note that

$$q_1^\delta(x) + \dots + q_N^\delta(x) = 1$$

for all $x \in X$. If the functions $q_j^\delta(x)$ assume only two values 0 and 1, then the test δ is nonrandomized. Otherwise a test δ is called *randomized*. The decision domain of the hypotheses H_j , $j = 1, 2, \dots, N$, is given by $X_j = \{x: q_j^\delta(x) = 1\}$ in the case of nonrandomized tests. Note that $X_j \cap X_i = \emptyset$, $i \neq j$, and $\bigcup_{j=1}^N X_j = X$.

The definition of a statistical test given in the preceding section differs to some extent from that given in this section, since the latter definition is inconvenient in the case $N = 2$. In what follows we use the simpler definition of the preceding section if we deal with the case of only two hypotheses.

To measure the quality of a test δ we introduce, as in the preceding section, the error probabilities:

$$(1.2.1) \quad \alpha_j(\delta) = P\{\delta \neq j/H_j\} = \int_X (1 - q_j^\delta(x)) P_j(dx), \quad j = 1, 2, \dots, N.$$

The number $\alpha_j(\delta)$ is the probability to reject the hypothesis H_j by using the test δ if the hypothesis H_j is true. The number $\alpha_j(\delta)$ is called the *type j error probability of the test δ* .

It is natural to say that a test δ_1 is better than a test δ_2 if $\alpha_j(\delta_1) \leq \alpha_j(\delta_2)$ for all $j = 1, 2, \dots, N$ and at least one of these inequalities is strict. However not all tests δ_1 and δ_2 can be compared in this way. In what follows we restrict the set of tests in order to have the possibility to compare them. Let

$$(1.2.2) \quad \mathcal{K}_{\alpha_1, \alpha_2, \dots, \alpha_{N-1}} = \{\delta: \alpha_j(\delta) = \alpha_j, j = 1, 2, \dots, N-1\}$$

where $\alpha_j \in [0, 1]$, $j = 1, 2, \dots, N-1$, are some fixed numbers.

A test $\delta^* \in \mathcal{K}_{\alpha_1, \dots, \alpha_{N-1}}$ is called the *most powerful test* (MP test, for short) in the class $\mathcal{K}_{\alpha_1, \dots, \alpha_{N-1}}$ if

$$\alpha_N(\delta^*) \leq \alpha_N(\delta)$$

for all tests $\delta \in \mathcal{K}_{\alpha_1, \dots, \alpha_{N-1}}$.

Prior to constructing the most powerful tests in the class $\mathcal{K}_{\alpha_1, \dots, \alpha_{N-1}}$ we consider two other approaches for comparing the tests, namely the Bayes and the minimax approaches.

The Bayes approach. Assume that the hypotheses H_1, \dots, H_N are random events whose probabilities are known. Put

$$P(H_j) = \pi_j, \quad j = 1, 2, \dots, N.$$

The family of probabilities $\pi = (\pi_1, \pi_2, \dots, \pi_N)$ is called the *a priori distribution of the hypotheses*. This family determines a distribution on the set $\Theta = \{1, 2, \dots, N\}$. The numbers $P_j(A) = P\{\xi \in A/H_j\}$ are conditional probabilities of the event $\{\xi \in A\}$ given the event H_j occurs. Moreover, we assume that the loss is A where $A = A_{ij}$ if the hypothesis H_j is accepted, while the hypothesis H_i is true. Therefore, the loss A is a random variable whose values are uniquely determined by the test δ and the index of the true hypothesis. To compare tests in the Bayes approach we use the *risk of a test* δ defined as the expectation of the loss:

$$(1.2.3) \quad R(\delta) = EA = \sum_{i=1}^N \sum_{j=1}^N A_{ij} p_{j/i}^\delta \pi_i$$

where $p_{j/i}^\delta = P\{\delta = j/H_i\}$. Since

$$(1.2.4) \quad p_{j/i}^\delta = \int_X q_j^\delta(x) P_j(dx),$$

the risk (1.2.3) can be rewritten as

$$(1.2.5) \quad R(\delta) = \sum_{i=1}^N \sum_{j=1}^N A_{ij} \pi_i \int_X q_j^\delta(x) P_j(dx).$$

A test $\delta_{\pi, A}$ that minimizes the risk $R(\delta)$ is called the *Bayes test corresponding to a priori distribution π and loss A* .

Let μ be some σ -finite measure on (X, \mathcal{B}) dominating the family of probability measures $\mathcal{P} = \{P_1, P_2, \dots, P_N\}$, $1 < N < \infty$, and let $p_i(x)$ be the density of the measure P_i with respect to the measure μ . Note that such a measure μ always exists. In particular, the measure $\mu = \sum_{i=1}^N c_i P_i$ where $c_i > 0$ for all i possesses this property. The risk (1.2.5) can be expressed in terms of the measure μ and densities $p_i(x)$:

$$(1.2.6) \quad R(\delta) = \int_X \sum_{i=1}^N \sum_{j=1}^N A_{ij} q_j^\delta(x) p_i(x) \pi_i \mu(dx).$$

Consider the measure $\bar{P} = \sum_{i=1}^N \pi_i P_i$ which defines the unconditional distribution of the observation ξ . Note that \bar{P} is absolutely continuous with respect to μ

and $f(x) = \sum_{i=1}^N \pi_i p_i(x)$ is the density of the measure \bar{P} with respect to the measure μ . Let $C = \{x: f(x) > 0\}$. Then $\bar{P}(X \setminus C) = 0$ and $P_i(X \setminus C) = 0$ for all i . For all $x \in C$ consider the functions

$$(1.2.7) \quad \pi_i(x) = \frac{p_i(x)\pi_i}{f(x)}, \quad i = 1, 2, \dots, N.$$

Equality (1.2.7) is the well-known Bayes formula for evaluating the conditional probability $\pi_i(x)$ of a hypothesis H_i given $\xi = x$. The numbers $\pi_i(x)$ are called a priori probabilities of hypotheses H_i .

THEOREM 1.2.1. For all tests δ it holds that

$$(1.2.8) \quad R(\delta) \geq E \min_{1 \leq j \leq N} \sum_{i=1}^N A_{ij} \pi_i(\xi).$$

A test $\delta = \delta_{\pi, A}$ is a Bayes test with respect to a priori distribution π and loss A if and only if

$$(1.2.9) \quad q_k^\delta(x) = 1 \quad \text{if} \quad \sum_{i=1}^N A_{ik} \pi_i(x) = \min_{1 \leq j \leq N} \sum_{i=1}^N A_{ij} \pi_i(x)$$

almost surely with respect to \bar{P} . If $\delta = \delta_{\pi, A}$, then inequality (1.2.8) becomes an equality.

PROOF. We obtain from equalities (1.2.6) and (1.2.7) that

$$(1.2.10) \quad \begin{aligned} R(\delta) &= \int_C \sum_{j=1}^N q_j^\delta(x) \sum_{i=1}^N A_{ij} \pi_i(x) f(x) \mu(dx) \\ &\geq \int_C \min_{1 \leq j \leq N} \sum_{i=1}^N A_{ij} \pi_i(x) \bar{P}(dx) = E \min_{1 \leq j \leq N} \sum_{i=1}^N A_{ij} \pi_i(\xi) \end{aligned}$$

where $C = \{x: f(x) > 0\}$, $\bar{P} = \sum_{i=1}^N \pi_i P_i$, and $f(x)$ is the density of the measure \bar{P} with respect to the measure μ . Thus inequality (1.2.8) is proved.

The sufficiency of condition (1.2.9) follows from (1.2.10). The case of an equality in (1.2.8) also follows from (1.2.10).

Now we prove the necessity of condition (1.2.9) by contradiction. Let $\delta = \delta_{\pi, A}$ be a Bayes test such that $q_k^\delta(x) = 1$ and

$$\sum_{i=1}^N A_{ik} \pi_i(x) > \sum_{i=1}^N A_{il} \pi_i(x) = \min_{1 \leq j \leq N} \sum_{i=1}^N A_{ij} \pi_i(x)$$

for $x \in A$, where A is some event of positive probability, $\bar{P}(A) > 0$. Let δ_1 be a test that differs from δ only on the event A and such that $q_l^{\delta_1}(x) = 1$ for $x \in A$. Then

$$\begin{aligned} R(\delta_1) &= \int_A \sum_{i=1}^N A_{il} \pi_i(x) \bar{P}(dx) + \int_{X \setminus A} \sum_{j=1}^N q_j^{\delta_1}(x) \sum_{i=1}^N A_{ij} \pi_i(x) \bar{P}(dx) \\ &< \int_A \sum_{i=1}^N A_{ik} \pi_i(x) \bar{P}(dx) + \int_{X \setminus A} \sum_{j=1}^N q_j^\delta(x) \sum_{i=1}^N A_{ij} \pi_i(x) \bar{P}(dx) \\ &= R(\delta). \end{aligned}$$

This is a contradiction, since δ is a Bayes test. The necessity of condition (1.2.9) is proved. \square

If $A_{ij} = 1 - \delta_{ij}$ where δ_{ij} is the Kronecker symbol (that is, $\delta_{ij} = 0$ for $i = j$ and $\delta_{ij} = 1$ for $i \neq j$), then the Bayes test $\delta_{\pi, A}$ for the loss A is called the *maximum a posteriori probability test*. In this case

$$\sum_{i=1}^N A_{ij} \pi_i(x) = \sum_{i \neq j} \pi_i(x) = 1 - \pi_j(x)$$

and condition (1.2.9) defining the Bayes test can be rewritten as

$$(1.2.11) \quad q_k^\delta(x) = 1 \quad \text{if } \pi_k(x) = \max_{1 \leq j \leq N} \pi_j(x).$$

The risk of an arbitrary test δ with loss $A_{ij} = 1 - \delta_{ij}$ is of the form

$$(1.2.12) \quad R(\delta) = \mathbb{E}A = \sum_{i=1}^N \sum_{j \neq i} p_{j/i}^\delta \pi_i = \mathbb{P}\{\delta \neq \theta\} = e_\pi(\delta)$$

where θ is the index of a hypothesis and $\{\theta = j\} = H_j$. We see that $R(\delta)$ in this case is the unconditional probability of a wrong decision $e_\pi(\delta)$ for the test δ . Thus the maximum a posteriori probability test minimizes the error probability $e_\pi(\delta)$ of the test δ . Taking into account (1.2.1) and (1.2.4) we obtain from (1.2.12) that

$$(1.2.13) \quad e_\pi(\delta) = \sum_{i=1}^N \mathbb{P}\{\delta \neq i/H_i\} \pi_i = \sum_{i=1}^N \alpha_i(\delta) \pi_i.$$

If $N = 2$ and $A_{ij} = 1 - \delta_{ij}$, then according to (1.2.11) the maximum a posteriori probability test is of the form

$$(1.2.14) \quad q_2^\delta(x) = \begin{cases} 1, & \pi_2(x) > \pi_1(x), \\ 0, & \pi_2(x) < \pi_1(x), \end{cases} \quad q_1^\delta(x) = 1 - q_2^\delta(x).$$

Moreover, if $\pi_2(x) = \pi_1(x)$, then one can put either $q_2^\delta(x) = 1$ or $q_1^\delta(x) = 1$. Applying equality (1.2.7) one can rewrite equality (1.2.14) as

$$(1.2.15) \quad q_2^\delta(x) = \begin{cases} 1, & \pi_2 p_2(x) > \pi_1 p_1(x), \\ 0, & \pi_2 p_2(x) < \pi_1 p_1(x), \end{cases} \quad q_1^\delta(x) = 1 - q_2^\delta(x).$$

Note that condition (1.2.9) does not uniquely determine the test $\delta_{\pi, A}$. In particular, it does not uniquely determine which hypothesis should be accepted if two or more numbers among $\sum_{i=1}^N A_{ij} \pi_i(x)$ are maximal. This is a matter of definition of the probabilities $q^\delta(x) = (q_1^\delta(x), q_2^\delta(x), \dots, q_N^\delta(x))$ on the boundaries

$$\Gamma_k = \left\{ x: \sum_{i=1}^N A_{ik} \pi_i(x) = \min_{j \neq k} \sum_{i=1}^N A_{ij} \pi_i(x) \right\}$$

of the sets

$$\tilde{X}_k = \left\{ x: \sum_{i=1}^N A_{ik} \pi_i(x) < \min_{j \neq k} \sum_{i=1}^N A_{ij} \pi_i(x) \right\}.$$

Using condition (1.2.9) the hypothesis H_k is accepted on Γ_k according to the test $\delta_{\pi, A}$. Therefore the problem is to decide about the points of the boundary Γ_k to be

included into the set \tilde{X}_k where the hypothesis H_k is accepted. One of the possible approaches is to include points of Γ_k to any of the regions \tilde{X}_j adjacent to Γ_k ; in this case condition (1.2.9) holds and the risk $R(\delta_{\pi,A})$ does not change. More precisely, if $A \subset \Gamma_{k_1} \cap \dots \cap \Gamma_{k_l}$, then, according to the Bayes test, it makes no difference for $x \in A$ which hypothesis among H_{k_1}, \dots, H_{k_l} is accepted. Moreover, one can accept the hypotheses H_{k_1}, \dots, H_{k_l} randomly with probabilities $q_{k_1}^\delta(x), \dots, q_{k_l}^\delta(x)$, $\sum_{i=1}^l q_{k_i}^\delta(x) = 1$. The risk $R(\delta_{\pi,A})$ does not change in this case.

The general definition of a Bayes test $\delta_{\pi,A}$ is based on the sets

$$(1.2.16) \quad \Gamma_{k_1, \dots, k_l} = \bigcap_{i=1}^l \Gamma_{k_i} \quad \bigcap_{j \neq k_1, \dots, k_l} \bar{\Gamma}_j$$

where $\bar{\Gamma}_j = X \setminus \Gamma_j$. As $q^\delta(x)$ for $x \in \Gamma_{k_1, \dots, k_l}$ one can take an arbitrary vector of the set R_{k_1, \dots, k_l} of vectors (q_1, q_2, \dots, q_N) with $q_1 \geq 0, q_2 \geq 0, \dots, q_N \geq 0$ and $\sum_{i=1}^N q_i = 1$ and whose coordinates with indices different from k_1, \dots, k_l are zero. It is clear that the set R_k includes only one vector e_k whose k -th coordinate is 1, while all others are zero. Thus one should put $q^\delta(x) = e_k$ for $x \in \tilde{X}_k$. This implies the following improvement of Theorem 1.2.1.

THEOREM 1.2.2. *A test δ is Bayes if and only if*

$$q^\delta(x) = \begin{cases} e_k, & \text{for } x \in \tilde{X}_k, \\ R_{k_1, \dots, k_l}, & \text{for } x \in \Gamma_{k_1, \dots, k_l} \end{cases}$$

for \bar{P} -almost all x where Γ_{k_1, \dots, k_l} are the sets defined by (1.2.16).

Theorems 1.2.1 and 1.2.2 show that randomized tests do not decrease the risk $R(\delta)$, however they enlarge the set of different Bayes tests $\delta_{\pi,A}$. Moreover Theorems 1.2.1 and 1.2.2 imply that among Bayes tests $\delta_{\pi,A}$ there is at least one nonrandomized test.

Let $N = 2$. Then

$$\begin{aligned} \tilde{X}_1 &= \left\{ x: \sum_{i=1}^2 A_{i1} \pi_i(x) < \sum_{i=1}^2 A_{i2} \pi_i(x) \right\}, \\ \tilde{X}_2 &= \left\{ x: \sum_{i=1}^2 A_{i2} \pi_i(x) < \sum_{i=1}^2 A_{i1} \pi_i(x) \right\}, \\ \Gamma_1 = \Gamma_2 &= \left\{ x: \sum_{i=1}^2 A_{i1} \pi_i(x) = \sum_{i=1}^2 A_{i2} \pi_i(x) \right\}. \end{aligned}$$

Hence we deal with a single set $\Gamma_{1,2} = \Gamma_1 = \Gamma_2$ instead of sets (1.2.16); moreover $R_{1,2} = \{(q_1, q_2): q_1 \geq 0, q_2 \geq 0, q_1 + q_2 = 1\}$ in this case. According to Theorem 1.2.2 we obtain for the Bayes test $\delta_{\pi,A}$ that $q_1^\delta(x) = 1$ for $x \in \tilde{X}_1$ and $q_2^\delta(x) = 1$ for $x \in \tilde{X}_2$. As $q^\delta(x) = (q_1^\delta(x), q_2^\delta(x))$ for $x \in \Gamma_{1,2}$ one can take an arbitrary function with values in $R_{1,2}$.

In the case of $N = 2$, the Bayes test $\delta = \delta_{\pi,A}$ equals the maximum a posteriori probability for $A_{ij} = 1 - \delta_{ij}$. Applying (1.2.14) and (1.2.15) we represent the test

$\delta = \delta_{\pi, A}$ in the form

$$(1.2.17) \quad q_2^\delta(x) = \begin{cases} 1, & z_{2,1}(x) > c, \\ q(x), & z_{2,1}(x) = c, \\ 0, & z_{2,1}(x) < c, \end{cases} \quad q_1^\delta(x) = 1 - q_2^\delta(x)$$

where $c = \pi_1/\pi_2$, $z_{2,1}(x) = p_2(x)/p_1(x)$ is the likelihood ratio (we assume that $0/0 = 0$), and $q(x)$ is an arbitrary measurable function with values in $[0, 1]$. A test of the form (1.2.17) for an arbitrary function $q(x)$ is called the *likelihood ratio test*. Note that we considered in the preceding section a likelihood ratio test of the form (1.2.17) for a specific function $q(x)$ being constant on the set $\{z_{2,1}(x) = c\}$ (see (1.1.35)). Like the preceding section, we denote by $\delta^{c,q}$ the likelihood ratio test δ defined by relation (1.2.17).

The minimax approach. The quality of a test δ in the minimax approach is measured by

$$(1.2.18) \quad e(\delta) = \max_{1 \leq j \leq N} \alpha_j(\delta) = \max_{\pi} e_{\pi}(\delta)$$

where $e_{\pi}(\delta)$ is the unconditional error probability of the test δ defined in (1.2.13). Recall that $e_{\pi}(\delta)$ is the risk of the test δ if a priori distribution is determined by the vector $\pi = (\pi_1, \pi_2, \dots, \pi_N)$ and loss is $A_{ij} = 1 - \delta_{ij}$.

A test δ^* such that

$$(1.2.19) \quad e(\delta^*) = \min_{\delta} e(\delta)$$

is called *minimax*, where $e(\delta)$ is the maximal error probability of the test δ (see relation (1.2.18)).

The following result contains a sufficient condition that a test is minimax.

THEOREM 1.2.3. *Let there exist a Bayes test $\bar{\delta}$ (with respect to some a priori distribution $\bar{\pi} = (\bar{\pi}_1, \dots, \bar{\pi}_N)$ and loss $A_{ij} = 1 - \delta_{ij}$) such that*

$$(1.2.20) \quad \alpha_1(\bar{\delta}) = \dots = \alpha_N(\bar{\delta}).$$

Then the test $\bar{\delta}$ is minimax.

PROOF. For all tests δ it holds that

$$e(\delta) \geq \sum_{j=1}^N \bar{\pi}_j \alpha_j(\delta) \geq \sum_{j=1}^N \bar{\pi}_j \alpha_j(\bar{\delta}) = \max_{1 \leq j \leq N} \alpha_j(\bar{\delta}) = e(\bar{\delta}),$$

whence it follows that the test $\bar{\delta}$ is minimax. □

Let $\bar{\pi}$ be a test satisfying (1.2.20). The a priori distribution $\bar{\pi}$ corresponding to the test $\bar{\pi}$ is called the *worst* or the *least favorable*. This notion is explained by saying that the maximum

$$(1.2.21) \quad \max_{\pi} e_{\pi}(\delta_{\pi}) = \max_{\pi} \min_{\delta} e_{\pi}(\delta)$$

is attained at $\pi = \bar{\pi}$ where δ_{π} is the *Bayes test with respect to a priori distribution π and loss $A_{ij} = 1 - \delta_{ij}$* . Therefore the minimax test satisfying (1.2.20) is the Bayes test with the maximal error probability. The proof of equality (1.2.21) and of the

existence of the worse a priori distribution and the minimax test can be found in [8, 10]. More detail on the minimax approach can be found in [4, 52].

In the case $N = 2$, the minimax test δ^* can be found by applying the fundamental Neyman–Pearson lemma and the set \mathfrak{N} described in Section 1.1. By Theorem 1.2.3 the minimax test δ^* is a Bayes test $\bar{\delta}$ with respect to some a priori distribution $\bar{\pi} = (\bar{\pi}_1, \bar{\pi}_2)$ and loss $A_{ij} = 1 - \delta_{ij}$ such that $\alpha_1(\bar{\delta}) = \alpha_2(\bar{\delta})$ provided such a test exists. Such a Bayes test $\bar{\delta}$ exists, indeed, and it corresponds to the point $A = (\alpha_1(\bar{\delta}), \alpha_2(\bar{\delta})) \in \mathfrak{N}$ of the intersection of the lower boundary of the set \mathfrak{N} and the diagonal of the square $[0, 1] \times [0, 1]$ joining its corners $(0, 0)$ and $(1, 1)$ (see Figure 1.2.1).

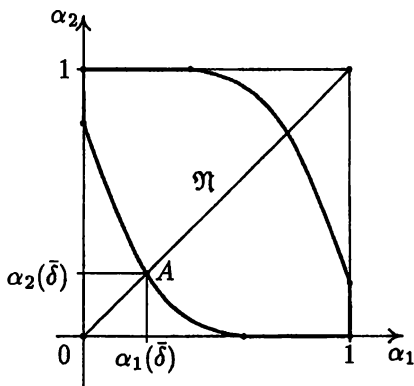


FIGURE 1.2.1

Moreover the following result gives an explicit form of the test $\bar{\delta}$.

THEOREM 1.2.4. *There is a likelihood test $\delta^{c,q}$ that is a minimax test. The parameters c and $q(x) \equiv q = \text{const}$ of the likelihood test $\delta^{c,q}$ are determined by the equation $\alpha_1(\delta^{c,q}) = \alpha_2(\delta^{c,q})$.*

PROOF. According to Theorem 1.2.3 it is sufficient to find a Bayes test $\bar{\delta}$ corresponding to some a priori distribution $\bar{\pi} = (\bar{\pi}_1, \bar{\pi}_2)$ and loss $A_{ij} = 1 - \delta_{ij}$ such that $\alpha_1(\bar{\delta}) = \alpha_2(\bar{\delta})$ and $\bar{\delta}$ coincides with a likelihood ratio test $\delta^{\bar{c},\bar{q}}$ for some parameters \bar{c} and $\bar{q}(x) \equiv \bar{q} = \text{const}$. We have seen above that such a Bayes test exists and it corresponds to the point $A = (\alpha_1(\bar{\delta}), \alpha_2(\bar{\delta})) \in \mathfrak{N}$ of the intersection of the lower boundary of the set \mathfrak{N} and the straight line joining the points $(0, 0)$ and $(1, 1)$ (see Figure 1.2.1). According to Theorems 1.2.1 and 1.2.2 this Bayes test $\bar{\delta}$ coincides with the likelihood ratio test $\delta^{\bar{c},\bar{q}}$ for some parameters \bar{c} and \bar{q} where $\bar{c} = \bar{\pi}_1/\bar{\pi}_2$ and $\bar{\pi} = (\bar{\pi}_1, \bar{\pi}_2)$ is the a priori distribution corresponding to the Bayes test $\bar{\delta}$. \square

REMARK 1.2.1. Since an arbitrary Bayes test coincides with a likelihood ratio test $\delta^{c,q}$ for some constants c and q , the proof of Theorem 1.2.4 follows from the existence of a solution of the equation $\alpha_1(\delta^{c,q}) = \alpha_2(\delta^{c,q})$ with respect to (c, q) . Note that this equation is of the form

$$P_1(z_{2,1} > c) + P_2(z_{2,1} > c) + q[P_1(z_{2,1} = c) + P_2(z_{2,1} = c)] = 1.$$

The proof that this equation has a solution with respect to (c, q) is the same as that of the existence of a solution of the equation $\alpha_1(\delta^{c,q}) = \alpha \in [0, 1]$ used in Theorem 1.1.1 and in Remarks 1.1.4 and 1.1.5.

The most powerful tests. We turn to the construction of the most powerful test in the class $\mathcal{K}_{\alpha_1, \dots, \alpha_{N-1}}$ defined by relation (1.2.2). Like the case of the minimax test, one can apply Bayes tests for this purpose, too.

THEOREM 1.2.5. *Let there exist an a priori distribution $\pi = (\pi_1, \pi_2, \dots, \pi_N)$ such that*

$$(1.2.22) \quad \alpha_j(\delta_\pi) = \alpha_j, \quad j = 1, 2, \dots, N-1,$$

where δ_π is the Bayes test corresponding to a priori distribution π and loss $A_{ij} = 1 - \delta_{ij}$. Then δ_π is the most powerful test in the class $\mathcal{K}_{\alpha_1, \dots, \alpha_{N-1}}$.

PROOF. The definition of a Bayes test implies that

$$e_\pi(\delta_\pi) \leq e_\pi(\delta)$$

for all tests δ , whence

$$(1.2.23) \quad \sum_{j=1}^N \pi_j \alpha_j(\delta_\pi) \leq \sum_{j=1}^{N-1} \pi_j \alpha_j + \pi_N \alpha_N(\delta)$$

for all tests $\delta \in \mathcal{K}_{\alpha_1, \dots, \alpha_{N-1}}$. Condition (1.2.22) implies that $\alpha_j(\delta_\pi) = \alpha_j$ for $j \leq N-1$. Then it follows from (1.2.23) that $\alpha_N(\delta_\pi) \leq \alpha_N(\delta)$, that is, δ_π is the most powerful test in the class $\mathcal{K}_{\alpha_1, \dots, \alpha_{N-1}}$. \square

REMARK 1.2.2. One can treat equalities (1.2.22) (and equalities (1.2.20), too) as the system of $N-1$ equations and use it to evaluate the a priori distribution π and the corresponding Bayes test δ_π . Generally speaking, this Bayes test δ_π is randomized. See Section 1.1 for another method of finding the most powerful test for $N=2$ that does not use the Bayes tests.

EXAMPLE 1.2.1 (Change point problem). Let an observation be a sample $\xi^{(n)} = (\xi_1, \dots, \xi_n)$ where ξ_1, \dots, ξ_n are independent random variables. The first $\theta-1$ of them have a distribution G_1 , while all other random variables have a distribution G_2 and $G_2 \neq G_1$. The number θ is called the *change point* (of the distribution). Possible values of θ are $1, 2, \dots, n$. Let $H_j = \{\theta = j\}$, $j = 1, 2, \dots, n$, be the statistical hypotheses about the parameter θ . Without loss of generality we assume that the measures G_1 and G_2 have densities $g_1(x)$ and $g_2(x)$, respectively, with respect to some σ -finite measure μ . Then the measure $P_j^{(n)}$ generating the distribution of $\xi^{(n)}$ under the hypothesis H_j is absolutely continuous with respect to the measure μ^n and its density is given by

$$p_j(x) = \prod_{i=1}^{j-1} g_1(x_i) \prod_{i=j}^n g_2(x_i), \quad x = (x_1, x_2, \dots, x_n),$$

where we put $\prod_{i=1}^0 = 1$.

Let δ^* be the maximum a posteriori probability test and let the a priori distribution be uniform, that is, δ^* is the Bayes test corresponding to the loss $A_{ij} = 1 - \delta_{ij}$ where δ_{ij} is the Kronecker symbol and the a priori distribution is $\pi = \pi_n^* = (1/n, 1/n, \dots, 1/n)$. According to (1.2.11) the test δ^* is such that

$$(1.2.24) \quad p_{\delta(x)^*}(x) = \max_{1 \leq j \leq n} p_j(x), \quad x \in \mathbf{R}^n,$$

that is, the test δ^* maximizes the density $p_j(x)$. A test δ^* satisfying equality (1.2.24) is called the *maximum likelihood test* (see [9]). Condition (1.2.24) means that for all $j = 1, 2, \dots, n$

$$(1.2.25) \quad \prod_{i=1}^{\delta^*(x)-1} g_1(x_i) \prod_{i=\delta^*(x)}^n g_2(x_i) \geq \prod_{i=1}^{j-1} g_1(x_i) \prod_{i=j}^n g_2(x_i).$$

Dividing inequality (1.2.25) by $p_1(x)$ one can prove that the test $\delta^*(x)$ is such that

$$\prod_{i=1}^{\delta^*(x)-1} \frac{g_1(x_i)}{g_2(x_i)} \geq \prod_{i=1}^{j-1} \frac{g_1(x_i)}{g_2(x_i)}$$

for all $j = 1, 2, \dots, n$. Properties of the test $\delta^*(x)$ are studied in detail in [9], §72.

To complete a brief discussion of the change point problem we mention [9], a survey paper [30], and a monograph [46] where more detail is given about the change point analysis (however there is an extensive literature devoted to this topic).

1.3. Distinguishing composite hypotheses

The setting of the problem and main definitions. Let ξ be an observation that is a random element assuming values in a measurable space (X, \mathcal{B}) and let $\mathcal{P} = \{P_\theta; \theta \in \Theta\}$ be a parametric family of probability measures defined on (X, \mathcal{B}) where Θ is some set containing more than two points. Let the distribution of the random element ξ be generated by a measure of the family \mathcal{P} .

Let $\mathcal{P}_i = \{P_\theta; \theta \in \Theta_i\}$, $i = 1, 2$, where $\Theta_1 \cap \Theta_2 = \emptyset$ and $\Theta_1 \cup \Theta_2 = \Theta$, so that $\mathcal{P}_1 \cup \mathcal{P}_2 = \mathcal{P}$. Consider the hypotheses H_1 and H_2 that the distribution of the element ξ belongs to the sets \mathcal{P}_1 and \mathcal{P}_2 , respectively. For the sake of brevity we write $H_i: \theta \in \Theta_i$, $i = 1, 2$. It is clear that at least one of the hypotheses H_1 and H_2 is *composite*, since at least one of the sets Θ_1 and Θ_2 contains at least two points. Consider the problem of distinguishing two hypotheses H_1 and H_2 by the observation $\xi = x$, $x \in X$.

As in Section 1.1 we consider a statistical test δ for distinguishing hypotheses H_1 and H_2 by the observation $\xi = x$. The test is a measurable mapping

$$\delta: (X, \mathcal{B}) \rightarrow ([0, 1], \mathcal{B}([0, 1])).$$

We treat $\delta(x)$ as the probability that the hypothesis H_2 is accepted if $\xi = x$, while $1 - \delta(x)$ is the probability that the hypothesis H_1 is accepted if $\xi = x$. We also put $\delta = \delta(\xi)$.

To measure the quality of a test δ we consider the function

$$(1.3.1) \quad \beta(\delta; \theta) = E_\theta \delta, \quad \theta \in \Theta,$$

where E_θ is the expectation with respect to the distribution P_θ . The function $\beta(\delta; \theta)$ is called the *power function of the test* δ . It is clear that $\beta(\delta; \theta)$ for $\theta \in \Theta_1$ is the probability of a wrong decision, while for $\theta \in \Theta_2$ it is the probability of a correct decision.

DEFINITION 1.3.1. We say that a test δ_1 is *uniformly more powerful* than a test δ_2 if

$$(1.3.2) \quad \beta(\delta_1; \theta) \leq \beta(\delta_2, \theta) \quad \text{for all } \theta \in \Theta_1,$$

$$(1.3.3) \quad \beta(\delta_1; \theta) \geq \beta(\delta_2, \theta) \quad \text{for all } \theta \in \Theta_2,$$

and at least one of the inequalities for at least one θ is strict.

DEFINITION 1.3.2. A test that is uniformly more powerful than any other test is called the *uniformly most powerful test* (UMP test).

There is a different approach for measuring the quality of tests, namely the Bayes approach. When following this approach we consider a random variable A treated as a loss: it assumes the value $A_i(\theta)$ if the hypothesis H_i is accepted and the true parameter is θ . The mean loss of the test δ is

$$(1.3.4) \quad E_\theta^\delta A = A_1(\theta)E_\theta(1 - \delta) + A_2(\theta)E_\theta\delta = (A_2(\theta) - A_1(\theta))\beta(\delta; \theta) + A_1(\theta)$$

if the parameter is θ . Here E_θ^δ is the expectation with respect to the distribution generated by the test δ if the true parameter is θ and $\beta(\delta; \theta)$ is the power function of the test δ defined by (1.3.1).

Let a σ -algebra $\mathcal{B}(\Theta)$ of measurable subsets of Θ be given and let a probability measure \mathbf{Q} be defined on the measurable space $(\Theta, \mathcal{B}(\Theta))$, that is, θ is a random parameter and $\mathbf{Q}(B) = P\{\theta \in B\}$, $B \in \mathcal{B}(\Theta)$. Applying (1.3.4) one can evaluate the mean loss:

$$(1.3.5) \quad E^\delta A = \int_{\Theta} E_i^\delta A \mathbf{Q}(dt) = \int_{\Theta} A_1(t) \mathbf{Q}(dt) + \int_{\Theta} (A_2(t) - A_1(t))\beta(\delta_i; t) \mathbf{Q}(dt)$$

where E^δ is the expectation with respect to the distribution generated by the test δ . The measure \mathbf{Q} is sometimes called a *a priori measure* or a *a priori distribution*.

DEFINITION 1.3.3. A test $\delta_{A, \mathbf{Q}}$ is called *Bayes* with respect to the loss A and a priori distribution \mathbf{Q} if

$$(1.3.6) \quad E^{\delta_{A, \mathbf{Q}}} A \leq E^\delta A$$

for any test δ where $E^\delta A$ is the mean loss defined by (1.3.5).

The following result allows one to compare the quality of tests in the Bayes approach and in an approach based on the power function.

THEOREM 1.3.1. *Let \mathbf{Q} be an arbitrary a priori measure. Assume that the loss function $A_i(t)$ is such that*

$$(1.3.7) \quad A_1(t) \leq A_2(t) \quad \text{for all } t \in \Theta_1,$$

$$(1.3.8) \quad A_1(t) \geq A_2(t) \quad \text{for all } t \in \Theta_2.$$

If a test δ_1 is uniformly more powerful than a test δ_2 , then their mean losses are such that

$$(1.3.9) \quad E^{\delta_1} A \leq E^{\delta_2} A.$$

PROOF. Let a test δ_1 be uniformly more powerful than a test δ_2 . Then inequalities (1.3.2) and (1.3.3) hold. Taking into account inequalities (1.3.2) and (1.3.3) and conditions (1.3.7) and (1.3.8) we obtain from (1.3.5) that

$$\begin{aligned} E^{\delta_1} A - E^{\delta_2} A &= \int_{\Theta_1} (A_2(t) - A_1(t))(\beta(\delta_1; t) - \beta(\delta_2; t)) \mathbf{Q}(dt) \\ &\quad + \int_{\Theta_2} (A_2(t) - A_1(t))(\beta(\delta_1; t) - \beta(\delta_2; t)) \mathbf{Q}(dt) \leq 0, \end{aligned}$$

that is, inequality (1.3.9) is proved. \square

COROLLARY 1.3.1. *If δ^* is a UMP test, then it also is a Bayes test with respect to an arbitrary a priori distribution \mathbf{Q} and any loss function $A_i(t)$ satisfying conditions (1.3.7) and (1.3.8).*

REMARK 1.3.1. Conditions (1.3.7) and (1.3.8) posed on the loss functions $A_1(t)$ and $A_2(t)$ are natural in the sense that if $t \in \Theta_1$, then $A_1(t)$ is the loss due to the acceptance of the hypothesis H_1 if it is true, while $A_2(t)$ is the loss due to the acceptance of the alternative hypothesis H_2 if the hypothesis H_1 is true. Thus it is reasonable to assume that $A_1(t) \leq A_2(t)$ for $t \in \Theta_1$. The same remark can be made regarding inequality (1.3.8).

REMARK 1.3.2. The method of comparing the quality of tests based on Definition 1.3.3 is sometimes called the *complete Bayes approach* (see [7]). Following this approach, one treats the numbers $\pi_i = \mathbf{Q}(\Theta_i) = \mathbf{P}\{\theta \in \Theta_i\}$, $i = 1, 2$, as a priori probabilities of the hypotheses H_1 and H_2 , respectively. Considered in [7] the so-called *partial Bayes approach* does not require that the probabilities π_1 and π_2 are known. Instead, the distributions of the parameter θ are known on both sets Θ_1 and Θ_2 .

Along with the Bayes approach we consider the minimax approach under which one seeks a test minimizing the maximum of the conditional mean loss $E_t^\delta A$.

DEFINITION 1.3.4. A test $\bar{\delta}$ is called *minimax* for the loss A if

$$(1.3.10) \quad \sup_{t \in \bar{\Theta}} E_t^{\bar{\delta}} A \leq \sup_{t \in \bar{\Theta}} E_t^\delta A$$

for all tests δ where $E_t^\delta A$ is the conditional mean loss defined by (1.3.4).

Generally speaking, the UMP tests do not exist in the class of all possible tests. Thus we consider proper subsets of tests and look for the UMP tests there. The following example exhibits this idea. The method below is suitable for finding both Bayes and minimax tests.

EXAMPLE 1.3.1. Let an observation ξ be a vector $\xi = (\xi_1, \xi_2, \dots, \xi_n)$ whose components are independent and have normal $\mathcal{N}(\theta, \sigma_i^2)$, $i = 1, 2, \dots, n$, distributions. Put $\Theta_1 = (-\infty, 0)$ and $\Theta_2 = [0, \infty)$. To distinguish the hypotheses $H_1: \theta \in \Theta_1$ and $H_2: \theta \in \Theta_2$ we consider the class of linear nonrandomized tests $\delta(r)$ of the form $\delta(x; r) = I((x, r) \geq 0)$ where $x = (x_1, x_2, \dots, x_n)$, $r = (r_1, r_2, \dots, r_n)$, and $(x, r) = \sum_{i=1}^n r_i x_i$; the vector r is such that $\sum_{i=1}^n r_i = 1$. Put $G(r; \theta) = \beta(\delta(r); \theta)$. It is clear that

$$G(r; \theta) = \mathbf{P}_\theta\{(\xi, r) \geq 0\} = \Phi \left(\theta \left(\sum_{i=1}^n \sigma_i^2 r_i^2 \right)^{-1/2} \right)$$

where $\Phi(x)$ is the distribution function of the law $\mathcal{N}(0, 1)$. Taking into account inequalities (1.3.2) and (1.3.3) we prove that $\delta(\tilde{r})$ is a UMP test if

$$\begin{aligned} \Phi \left(\theta \left(\sum_{i=1}^n \sigma_i^2 \tilde{r}_i^2 \right)^{-1/2} \right) &\leq \Phi \left(\theta \left(\sum_{i=1}^n \sigma_i^2 r_i^2 \right)^{-1/2} \right) \quad \text{for all } \theta < 0, \\ \Phi \left(\theta \left(\sum_{i=1}^n \sigma_i^2 \tilde{r}_i^2 \right)^{-1/2} \right) &\geq \Phi \left(\theta \left(\sum_{i=1}^n \sigma_i^2 r_i^2 \right)^{-1/2} \right) \quad \text{for all } \theta \geq 0 \end{aligned}$$

for all vectors $r = (r_1, r_2, \dots, r_n)$ where $\tilde{r} = (\tilde{r}_1, \tilde{r}_2, \dots, \tilde{r}_n)$. The latter inequalities are equivalent to

$$\sum_{i=1}^n \sigma_i^2 \tilde{r}_i^2 \leq \sum_{i=1}^n \sigma_i^2 r_i^2 \quad \text{for all } r = (r_1, r_2, \dots, r_n).$$

Thus in order to find a UMP test it is necessary to find a vector \tilde{r} for which the function $\sum_{i=1}^n \sigma_i^2 r_i^2$ assumes its minimal value on the set of vectors r such that $\sum_{i=1}^n r_i = 1$. It is clear that the components of the vector \tilde{r} with this property are such that

$$\tilde{r}_i = \sigma_i^{-2} \left(\sum_{j=1}^n \sigma_j^{-2} \right)^{-1}, \quad i = 1, 2, \dots, n.$$

The following result provides necessary and sufficient conditions for the existence of a UMP test.

THEOREM 1.3.2. *In order that a test is UMP for distinguishing the hypotheses $H_1: \theta \in \Theta_1$ and $H_2: \theta \in \Theta_2$ it is necessary and sufficient that it is an MP test for distinguishing two arbitrary simple hypotheses $H'_1: \theta = \theta_1$ and $H'_2: \theta = \theta_2$ where $\theta_1 \in \Theta_1$ and $\theta_2 \in \Theta_2$.*

PROOF. *Necessity.* Let δ^* be a UMP test and let δ be an arbitrary test. Then $\beta(\delta^*; \theta) \leq \beta(\delta; \theta)$ for all $\theta \in \Theta_1$ and $\beta(\delta^*; \theta) \geq \beta(\delta; \theta)$ for all $\theta \in \Theta_2$. Let $\theta_1 \in \Theta_1$ and $\theta_2 \in \Theta_2$. Consider two simple hypotheses $H'_1: \theta = \theta_1$ and $H'_2: \theta = \theta_2$. Then type I and type II error probabilities of the test are such that

$$(1.3.11) \quad \begin{aligned} \alpha(\delta^*) &= E_{\theta_1} \delta^* = \beta(\delta^*; \theta_1) \leq \beta(\delta; \theta_1) = \alpha(\delta), \\ \beta(\delta^*) &= E_{\theta_2} (1 - \delta^*) = 1 - \beta(\delta^*; \theta_2) \leq 1 - \beta(\delta; \theta_2) = \beta(\delta) \end{aligned}$$

(see Section 1.1), that is, δ^* is an MP test.

Sufficiency. Let δ^* be an MP test for distinguishing the hypotheses $H'_1: \theta = \theta_1$ and $H'_2: \theta = \theta_2$ for all $\theta_i \in \Theta_i$, $i = 1, 2$. The Neyman–Pearson fundamental lemma (see Theorem 1.1.1) implies that $\beta(\delta^*) \leq \beta(\delta)$ for all tests δ of level $E_{\theta_1} \delta^* = \beta(\delta^*; \theta_1)$. It follows from (1.3.11) that $\beta(\delta; \theta_2) \leq \beta(\delta^*; \theta_2)$ for all $\theta_2 \in \Theta_2$ and all tests δ such that $\beta(\delta; \theta_1) = \beta(\delta^*; \theta_1)$. Interchanging the hypotheses H'_1 and H'_2 and applying the Neyman–Pearson fundamental lemma once more we prove that $\beta(\delta^*; \theta_1) \leq \beta(\delta; \theta_1)$ for all tests δ such that $\beta(\delta; \theta_2) = \beta(\delta^*; \theta_2)$ and all $\theta_1 \in \Theta_1$. \square

The following example is a continuation of Example 1.3.1. It exhibits an application of Theorem 1.3.2 for finding a UMP test for another restriction of the class of tests as compared to that studied in Example 1.3.1.

EXAMPLE 1.3.2. Let $\xi = (\xi_1, \xi_2, \dots, \xi_n)$ and random variables ξ_i be independent and distributed according to the $\mathcal{N}(\theta, \sigma_i^2)$, $i = 1, 2, \dots, n$, laws. Consider the hypotheses $H_1: \theta \leq 0$ and $H_2: \theta > 0$. In order to construct an UMP test we use Theorem 1.3.2 and find an MP test for distinguishing the hypotheses $H'_1: \theta = \theta_1$ and $H'_2: \theta = \theta_2$ for all $\theta_1 \leq 0$ and $\theta_2 > 0$. The Neyman–Pearson fundamental lemma implies that this can be done by constructing a likelihood ratio test for distinguishing the hypotheses H'_1 and H'_2 . Let $p(x; \theta_i)$ be the density of the distribution of the vector ξ in the case of $\theta = \theta_i$. Then

$$(1.3.12) \quad \Lambda(x) = \ln \frac{p(x; \theta_2)}{p(x; \theta_1)} = (\theta_2 - \theta_1) \sum_{i=1}^n \frac{x_i}{\sigma_i^2} + \frac{\theta_2^2 - \theta_1^2}{2} \sum_{i=1}^n \frac{1}{\sigma_i^2}$$

where $x = (x_1, x_2, \dots, x_n)$. This means that $\delta(x) = I(\Lambda(x) > c)$ is the desired likelihood ratio test where c is some constant. Since $\theta_2 > \theta_1$, we obtain from (1.3.12) that this test is of the form

$$(1.3.13) \quad \delta(x) = I\left(\sum_{i=1}^n \sigma_i^{-2} x_i > k\right)$$

where k is some constant. If the constant k is such that $E_0 \delta(\xi) = \alpha$ where $\alpha \in (0, 1)$ is a certain number, then the test δ given by (1.3.13) for some constant k does not depend on θ_1 and θ_2 . Now Theorem 1.3.2 implies that the test defined by (1.3.13) is UMP in the class of all tests such that $\beta(\delta; 0) = \alpha$ where the constant k is specified above.

The latter example suggests a general idea on how to restrict the class of tests under consideration.

DEFINITION 1.3.5. The number

$$\alpha_1(\delta) = \sup_{\theta \in \Theta_1} \beta(\delta; \theta)$$

is called a *level* or *type I error probability* of the test δ .

The number $\alpha_1(\delta)$ is sometimes called the *size of the test* δ . This is the maximal probability of rejecting the hypothesis H_1 if it is true.

Consider the class of tests

$$K_\alpha = \{\delta: \alpha_1(\delta) \leq \alpha\}$$

where α is some number of the interval $[0, 1]$.

DEFINITION 1.3.6. A test δ^* is called a *uniformly most powerful (UMP) test* in the class K_α if

$$\beta(\delta^*; \theta) \geq \beta(\delta; \theta), \quad \theta \in \Theta_2,$$

for an arbitrary test $\delta \in K_\alpha$.

Similar definitions in the class K_α can be introduced for Bayes and minimax tests [7].

UMP test for distributions with a monotone likelihood ratio. Let

$$\Theta = (-\infty, \infty), \quad \Theta_1 = (-\infty, \theta_0], \quad \text{and} \quad \Theta_2 = (\theta_0, \infty)$$

where θ_0 is a fixed point of Θ . The hypotheses $H_1: \theta \in \Theta_1$ and $H_2: \theta \in \Theta_2$ are called *one-sided hypotheses* in contrast to the case of $H_2: \theta \neq \theta_0$ and $H_1: \theta = \theta_0$ where they are called *two-sided hypotheses*, since the sign of $\theta - \theta_0$ can be arbitrary in the latter case.

In what follows we assume that the measure P_θ for all $\theta \in \Theta$ is absolutely continuous with respect to some σ -finite measure μ and the density is

$$p(x; \theta) = dP_\theta/d\mu(x), \quad x \in X.$$

Moreover we assume that for all $\theta_1 < \theta_2$ the likelihood ratio

$$(1.3.14) \quad z(x; \theta_2, \theta_1) = \frac{p(x; \theta_2)}{p(x; \theta_1)}, \quad x \in X,$$

is a monotone (either nondecreasing or nonincreasing) function of some statistic $T(x)$. We say in this case that the family $\mathcal{P} = \{P_\theta; \theta \in \Theta\}$ has a *monotone likelihood ratio*. For the sake of definiteness we assume that the likelihood ratio (1.3.14) is a nondecreasing function of $T(x)$. The case where $z(x; \theta_2, \theta_1)$ is a nonincreasing function of $T(x)$ can be considered analogously.

THEOREM 1.3.3. *Let $\theta \in \Theta = (-\infty, \infty)$ be a one-dimensional parameter. Assume that a family $\mathcal{P} = \{P_\theta; \theta \in \Theta\}$ has a monotone likelihood ratio $z(x; \theta_2, \theta_1)$. Then*

1. *In the class of tests K_α , $\alpha \in (0, 1)$, there exists a UMP test for distinguishing the hypotheses $H_1: \theta \in \Theta_1 = (-\infty, \theta_0]$ and $H_2: \theta \in \Theta_2 = (\theta_0, \infty)$. The test is given by*

$$(1.3.15) \quad \delta^*(x) = I(T(x) > c) + qI(T(x) = c)$$

where $c \in (-\infty; \infty)$ and $q \in [0, 1]$ are the parameters defined by

$$(1.3.16) \quad E_{\theta_0} \delta^* = P_{\theta_0}(T(\xi) > c) + qP_{\theta_0}(T(\xi) = c) = \alpha.$$

2. *The power function $\beta(\delta^*; \theta)$ of the test δ^* defined by equalities (1.3.15) and (1.3.16) is a nondecreasing function of $\theta \in \Theta$.*
3. *For all $\theta' \in \Theta$ the test (1.3.15) is a UMP test in the class $K_{\beta(\delta^*; \theta')}$ for distinguishing the hypotheses $H_1': \theta \leq \theta'$ and $H_2': \theta > \theta'$.*
4. *For every $\theta < \theta_0$ the test δ^* defined by equalities (1.3.15) and (1.3.16) minimizes the function $\beta(\delta; \theta)$ in the class K_α .*

PROOF. Consider the two simple hypotheses $\tilde{H}_1: \theta = \theta_0$ and $\tilde{H}_2: \theta = \theta_2$ where $\theta_2 > \theta_0$. According to the Neyman–Pearson fundamental lemma (Theorem 1.1.1), a most powerful test for distinguishing the hypotheses \tilde{H}_1 and \tilde{H}_2 in the class of tests δ such that $E_{\theta_0} \delta = \alpha$ is of the form (1.3.15), since the inequality $z(x; \theta_2, \theta_0) > c$ is equivalent to the inequality $T(x) > c$ in view of the monotonicity of the likelihood ratio where the constants c and q are defined by (1.3.16). Since the parameters c and q do not depend on θ_2 , the test δ^* is the most powerful for distinguishing the hypotheses $\tilde{H}_1: \theta = \theta_0$ and $\tilde{H}_2: \theta = \theta_2$ for all $\theta_2 \in \Theta_2$. Thus Theorem 1.3.2 implies that the test δ^* maximizes $\beta(\delta; \theta)$ for all $\theta \in \Theta_2$ in the class of tests δ such that $\beta(\delta; \theta_0) = \alpha$.

Now let θ' and θ'' be two arbitrary points such that $\theta' < \theta''$. Again by the Neyman–Pearson fundamental lemma the test δ^* is the most powerful for distinguishing the simple hypotheses $\tilde{H}_1': \theta = \theta'_1$ and $\tilde{H}_2': \theta = \theta''$ in the class of tests of level $\alpha' = \beta(\delta^*; \theta')$. By the definition of the set \mathfrak{N} (see Section 1.1) $\alpha(\delta^*) \leq 1 - \beta(\delta^*)$ for the most powerful test δ^* . Thus

$$\alpha' = \beta(\delta^*; \theta') = \alpha(\delta^*) \leq 1 - \beta(\delta^*) = \beta(\delta^*; \theta'').$$

Therefore $\beta(\delta^*; \theta') \leq \beta(\delta^*; \theta'')$ for all $\theta' < \theta''$ and statement 2 of the theorem is proved.

Since the function $\beta(\delta^*; \theta)$ is nondecreasing, the test δ^* is such that $\beta(\delta^*; \theta) \leq \alpha$ for all $\theta \leq \theta_0$, that is, the test δ^* belongs to the class K_α . In its turn K_α belongs to a wider class $\{\delta: E_{\theta_0} \delta = \alpha\}$. Since δ^* is a UMP test in the class $\{\delta: E_{\theta_0} \delta = \alpha\}$, it also is a UMP test in the class K_α .

Statement 3 of the theorem can be proved in the same manner.

Statement 4 follows from statements 1–3 applied to the problem of distinguishing the hypotheses $H_1^0: \theta \geq \theta_0$ and $H_1^0: \theta < \theta_0$. A UMP test in the class $\{\delta: \sup_{\theta \geq \theta_0} E_\theta \delta \leq 1 - \alpha\}$ is $\delta^0(x) = 1 - \delta(x)$ for this problem and the power function $1 - \beta(\delta^*; \theta) = E_\theta \delta^0$ is maximal for $\theta < \theta_0$. Therefore the test δ^* minimizes $\beta(\delta; \theta)$ for $\theta < \theta_0$ and for tests δ of the class K_α . \square

REMARK 1.3.3. Equality (1.3.12) shows that the likelihood ratio in Example 1.3.2 is a monotone function of the statistic

$$T(x) = \sum_{i=1}^n \sigma_i^{-2} x_i, \quad x = (x_1, x_2, \dots, x_n).$$

Thus Theorem 1.3.3 is applicable in this case and a UMP test exists for distinguishing the one-sided hypotheses.

EXAMPLE 1.3.3. Let $\xi = (\xi_1, \xi_2, \dots, \xi_n)$ where $\xi_1, \xi_2, \dots, \xi_n$ are independent identically distributed random variables whose distribution depends on a parameter $\theta \in (0, 1)$ such that $P_\theta\{\xi_i = 1\} = \theta$ and $P_\theta\{\xi_i = 0\} = 1 - \theta$. The space X of possible values of the random vector ξ consists of the vectors $x = (x_1, x_2, \dots, x_n)$ whose coordinates x_i are either 0 or 1. The distribution of the vector ξ is given by

$$P_\theta(x) = P_\theta\{\xi = x\} = \theta^{\sum_{i=1}^n x_i} (1 - \theta)^{n - \sum_{i=1}^n x_i}, \quad x = (x_1, \dots, x_n),$$

whence we obtain the likelihood ratio:

$$z(x; \theta_2, \theta_1) = \left(\frac{\theta_2}{\theta_1}\right)^{\sum_{i=1}^n x_i} \left(\frac{1 - \theta_2}{1 - \theta_1}\right)^{n - \sum_{i=1}^n x_i} = \left(\frac{\theta_2(1 - \theta_1)}{\theta_1(1 - \theta_2)}\right)^{\sum_{i=1}^n x_i} \left(\frac{1 - \theta_2}{1 - \theta_1}\right)^n.$$

We see that the likelihood ratio $z(x; \theta_2, \theta_1)$ for $\theta_1 < \theta_2$ is an increasing function of the statistic $T(x) = \sum_{i=1}^n x_i$. According to Theorem 1.3.3 there exists a UMP test for distinguishing one-sided hypotheses in the class K_α .

An important class of distributions for which the likelihood ratio is monotone is presented by the *one parameter exponential family*. The density $p(x; \theta)$ in this case is given by

$$(1.3.17) \quad p(x; \theta) = h(x) \exp\{a(\theta)T(x) + V(\theta)\}, \quad x \in X,$$

where all the functions on the right-hand side are finite and measurable (see [38], Section 1.2). In view of the factorization criterion (Theorem 4.1.2 in [38]), the statistic $T(x)$ is sufficient. The likelihood ratio in this case is given by

$$z(x; \theta_2, \theta_1) = \exp\{(a(\theta_2) - a(\theta_1))T(x) + (V(\theta_2) - V(\theta_1))\}.$$

This implies that if $a(\theta_2) - a(\theta_1)$ does not change its sign for all $\theta_1 < \theta_2$, then the likelihood ratio $z(x; \theta_2, \theta_1)$ is a monotone function of the statistic $T(x)$.

Theorem 1.3.3 implies the following assertion.

COROLLARY 1.3.2. *Let the density $p(x, \theta)$ be of the form (1.3.17) where $a(\theta)$ is a monotone function. Then there exists a UMP test δ^* in the class K_α for distinguishing the hypotheses $H_1: \theta \leq \theta_0$ and $H_2: \theta > \theta_0$. If the function $a(\theta)$ increases, then the test δ^* is defined by (1.3.15) and (1.3.16). If the function $a(\theta)$ decreases, then the test δ^* is defined by (1.3.15) and (1.3.16) where $T(x) \leq c$ and $T(\xi) \leq c$ substitute $T(x) > c$ and $T(\xi) > c$, respectively.*

REMARK 1.3.4. If we distinguish the hypothesis $H_1: \theta = \theta_0$ and its two-sided alternative $H_2: \theta \neq \theta_0$, then a UMP test does not exist in the case of exponential distributions (1.3.17). Indeed, for simplicity let the function $a(\theta)$ increase and let the P_θ -distribution of $T(\xi)$ for all θ possess the density. Then by the Neyman–Pearson fundamental lemma a most powerful test for distinguishing the hypotheses

$$H_1: \theta = \theta_0 \quad \text{and} \quad H_2': \theta = \theta_2$$

with $\theta_2 > \theta_0$ is nonrandomized and moreover $\delta^*(x) = I(T(x) \geq c)$. On the other hand, if $\theta_2 < \theta_0$, then the most powerful test is $\delta'(x) = I(T(x) < c)$. Thus there is no unique UMP test for all $\theta_2 > \theta_0$ and $\theta_2 < \theta_0$. In a similar manner we get that there is no UMP test for distinguishing the hypotheses $H_1: \theta \in (\theta_1, \theta_2)$ and $H_2: \theta \notin (\theta_1, \theta_2)$ where $\theta_1 < \theta_2$. However if $H_1: \theta \notin (\theta_1, \theta_2)$ is the null hypothesis and $H_2: \theta \in (\theta_1, \theta_2)$ is its alternative, then a UMP test exists. This case is studied in the next section.

Two-sided null hypotheses. Exponential families of distributions. Let a distribution P_θ be absolutely continuous with respect to some σ -finite measure μ and let the density $p(x; \theta) = dP_\theta/d\mu(x)$ be of the form (1.3.17).

THEOREM 1.3.4. *Let $p(x; \theta)$ be of the form (1.3.17) where the function $a(\theta)$ is monotone. Let $H_1: \theta \notin (\theta_1, \theta_2)$ be the null hypothesis and let $H_2: \theta \in (\theta_1, \theta_2)$ be its alternative where $\theta_1 < \theta_2$ are two fixed numbers. Then*

1) *in the class*

$$K_\alpha = \left\{ \delta: \sup_{\theta \notin (\theta_1, \theta_2)} E_\theta \delta(\xi) \leq \alpha \right\}$$

there exists a UMP test δ^ such that*

$$(1.3.18) \quad \delta^*(x) = I(c_1 < T(x) < c_2) + q_1 I(T(x) = c_1) + q_2 I(T(x) = c_2)$$

where $c_1, c_2, q_1,$ and q_2 are the constants defined by

$$(1.3.19) \quad E_{\theta_1} \delta^*(\xi) = E_{\theta_2} \delta^*(\xi) = \alpha;$$

2) *the test δ^* defined by equalities (1.3.18) and (1.3.19) maximizes the power function $\beta(\delta; \theta)$ inside the interval (θ_1, θ_2) and minimizes it outside this interval;*

- 3) for $0 < \alpha < 1$ the power function $\beta(\delta^*; \theta)$ attains its maximum at some point $\theta_0 \in (\theta_1, \theta_2)$; moreover it strictly decreases in both cases if the argument goes to the left of θ_0 or if it goes to the right of θ_0 . Note that the case where the distribution of $T(\xi)$ is concentrated at two points is not excluded, that is, the case where there are t_1 and t_2 such that for all θ

$$P_\theta\{T(\xi) = t_1\} + P_\theta\{T(\xi) = t_2\} = 1.$$

We omit the proof of Theorem 1.3.4 that can be found in [7, 9], or [34].

The generalized Neyman–Pearson fundamental lemma. The construction of MP and UMP tests requires, in fact, the solution of a variational problem and finding a maximum of a certain functional of the test satisfying some restrictions. In particular, we deal with the test δ^* in Theorem 1.3.4 for which we maximize the functional

$$\int_X \delta(x) p(x; \theta) \mu(dx)$$

in the class of tests δ such that

$$\int_X \delta(x) p(x; \theta_i) \mu(dx) = \alpha, \quad i = 1, 2.$$

The following result is sometimes called the *generalized Neyman–Pearson fundamental lemma*.

THEOREM 1.3.5. Let f_1, f_2, \dots, f_{m+1} be real Borel functions defined on (X, \mathcal{B}) that are integrable with respect to a measure μ . Consider the tests δ such that

$$(1.3.20) \quad \int_X \delta(x) f_i(x) \mu(dx) = \alpha_i, \quad i = 1, 2, \dots, m,$$

where $\alpha_1, \alpha_2, \dots, \alpha_m$ are some numbers. Then the test $\delta^*(x)$ that maximizes the functional $\int_X \delta(x) f_{m+1}(x) \mu(dx)$ is of the form

$$\delta^*(x) = \begin{cases} 1, & \text{if } f_{m+1}(x) > \sum_{i=1}^m k_i f_i(x), \\ 0, & \text{if } f_{m+1}(x) < \sum_{i=1}^m k_i f_i(x), \end{cases}$$

where the constants k_1, \dots, k_m are defined by conditions (1.3.20).

PROOF. Put

$$F_i(\delta) = \int_X \delta(x) f_i(x) \mu(dx), \quad i = 1, 2, \dots, m + 1.$$

A test δ such that

$$F_i(\delta) = \alpha_i, \quad i = 1, 2, \dots, m,$$

maximizes $F_{m+1}(\delta)$ if and only if it maximizes $F_{m+1}(\delta) - \sum_{i=1}^m k_i F_i(\delta)$ for some constants k_1, k_2, \dots, k_m ($\sum_{i=1}^m k_i F_i(\delta)$ is fixed in this expression). This is the case if the test $\delta(x)$ maximizes the functional

$$\int_X \left(f_{m+1}(x) - \sum_{i=1}^m k_i f_i(x) \right) \delta(x) \mu(dx).$$

The latter expression is maximal for the test δ such that $\delta(x) = 1$ if

$$f_{m+1}(x) - \sum_{i=1}^m k_i f_i(x) > 0$$

and $\delta(x) = 0$ otherwise. The constants k_1, k_2, \dots, k_m occurring in the definition of the test δ , as well as the values of $\delta(x)$ on the set

$$\left\{ f_{m+1}(x) = \sum_{i=1}^m k_i f_i(x) \right\},$$

should be chosen to satisfy conditions (1.3.20). □

Unbiased tests. Another restricted class we use to construct UMP tests consists of the so-called unbiased tests.

Consider the general problem of distinguishing the hypotheses $H_1: \theta \in \Theta_1$ and $H_2: \theta \in \Theta_2$ where $\Theta_1 \cap \Theta_2 = \emptyset$ and $\Theta_1 \cup \Theta_2 = \Theta$. Let δ be a test of the class

$$K_\alpha = \left\{ \delta: \alpha_1(\delta) = \sup_{\theta \in \Theta_1} E_\theta \delta \leq \alpha \right\}.$$

If Θ_1 contains only a single point θ_1 and $E_{\theta_1} \delta = \alpha$, then α is the probability of rejecting the hypothesis H_1 if it is true. It is natural to require that a test δ is such that the probability of rejecting the hypothesis H_1 , if it is wrong, is bigger than α , that is, $\beta(\delta; \theta) \geq \alpha$ for all $\theta \in \Theta_2$. If this is not the case, then there are alternative hypotheses $\theta \in \Theta_2$ such that the probability of accepting the hypothesis H_1 is bigger than $1 - \alpha = 1 - E_{\theta_1} \delta$ and the latter is the probability of accepting the hypothesis H_1 if it is true. It is reasonable to exclude such cases from our consideration.

DEFINITION 1.3.7. A test δ is called *unbiased* if

$$(1.3.21) \quad \inf_{\theta \in \Theta_2} \beta(\delta; \theta) \geq \sup_{\theta \in \Theta_1} \beta(\delta; \theta).$$

Condition (1.3.21) implies that a test $\delta \in K_\alpha$ of level $\alpha_1(\delta) = \alpha$ is unbiased if $\beta(\delta; \theta) \geq \alpha$ for all $\theta \in \Theta_2$. The class of unbiased tests of level α is denoted by \widehat{K}_α . By $\partial\Theta_i$ we denote the boundary of the set Θ_i , that is, all the limit points of the set Θ_i .

LEMMA 1.3.1. Let $\Gamma = \partial\Theta_1 \cap \partial\Theta_2 \neq \emptyset$. Assume that the density $p(x; \theta)$ is continuous in θ for μ -almost all $x \in X$. Then

$$(1.3.22) \quad \beta(\delta; \theta) = \alpha \quad \text{for all } \theta \in \Gamma$$

for any test $\delta \in \widehat{K}_\alpha$.

PROOF. Since

$$\beta(\delta; \theta) = \int_X \delta(x) p(x; \theta) \mu(dx), \quad 0 \leq \delta(x) \leq 1,$$

and the function $p(x; \theta)$ is continuous in θ by Corollary 3.4.1 in [38], the power function $\beta(\delta; \theta)$ is also continuous for any test δ . This implies equality (1.3.22) for any test $\delta \in \widehat{K}_\alpha$. □

We denote by \overline{K}_α the class of tests δ satisfying condition (1.3.22).

LEMMA 1.3.2. Let $\Gamma = \partial\Theta_1 \cap \partial\Theta_2 \neq \emptyset$ and $\widehat{K}_\alpha \subset \overline{K}_\alpha$. Then any UMP test in the class $\overline{K}_\alpha \cap K_\alpha$ is UMP in the class \widehat{K}_α , too.

PROOF. Let δ' be a UMP test in the class $\overline{K}_\alpha \cap K_\alpha$. It is sufficient to prove that $\delta' \in \widehat{K}_\alpha$ and $\widehat{K}_\alpha \subset \overline{K}_\alpha \cap K_\alpha$. The inclusion $\widehat{K}_\alpha \subset \overline{K}_\alpha \cap K_\alpha$ follows from $\widehat{K}_\alpha \subset \overline{K}_\alpha$ and $\widehat{K}_\alpha \subset K_\alpha$. Since δ' is a UMP test in the class $\overline{K}_\alpha \cap K_\alpha$ and the test $\delta \equiv \alpha$ belongs to $\overline{K}_\alpha \cap K_\alpha$, we have $\beta(\delta'; \theta) \geq \beta(\delta; \theta) = \alpha$ for all $\theta \in \Theta_2$, whence

$$\inf_{\theta \in \Theta_2} \beta(\delta'; \theta) \geq \inf_{\theta \in \Theta_2} \beta(\delta; \theta) = \alpha.$$

Therefore $\delta' \in \widehat{K}_\alpha$. □

Lemma 1.3.2 implies that the problem of finding an unbiased UMP test can be reduced to the problem of finding a usual UMP test satisfying condition (1.3.22). If the number of points of the boundary Γ is finite, then the conditions of Theorem 1.3.5 hold and one needs to prove that an optimal test δ^* does not depend on the point $\theta \in \Theta_2$ at which the test maximizes the functional $\beta(\delta; \theta) = E_\theta \delta$. This means that δ^* is a UMP test.

Consider condition (1.3.22) for the following case. Let

$$\Theta = (-\infty, \infty), \quad \Theta_1 = [\theta_1, \theta_2], \quad \text{and} \quad \Theta_2 = [\theta_1, \theta_2]^c = (-\infty, \infty) \setminus [\theta_1, \theta_2].$$

If $\theta_1 < \theta_2$, then the common boundary Γ of the sets Θ_1 and Θ_2 contains only two points θ_1 and θ_2 . Therefore condition (1.3.22) becomes of the form $\beta(\delta; \theta_i) = \alpha$, $i = 1, 2$. If $\Theta_1 = \{\theta_1\}$, then condition (1.3.22) is equivalent to $\beta(\delta; \theta_1) = \alpha$.

The following result gives a necessary condition for a test of level α to be unbiased.

LEMMA 1.3.3. Let δ be a test of level α for distinguishing the hypotheses

$$H_1: \theta = \theta_1 \quad \text{and} \quad H_2: \theta \neq \theta_1.$$

Assume that the regularity conditions (R) hold for the density $p(x; \theta)$ of the measure P_θ with respect to a σ -finite measure μ . If the test δ is unbiased, then

$$(1.3.23) \quad E_{\theta_1} \delta(\xi) S(\xi; \theta_1) = 0$$

where $S(x; \theta) = \frac{\partial}{\partial \theta} \ln p(x; \theta)$.

PROOF. Since the test δ is unbiased, the power function $\beta(\delta; \theta)$ attains its minimum at the point θ_1 . Since the function $p(x; \theta)$ satisfies the regularity conditions (R), the power function $\beta(\delta; \theta)$ is differentiable by Lemma 3.4.4 of [38]. Therefore the equality $\beta'(\delta; \theta_1) = 0$ is satisfied. Applying again Lemma 3.4.4 of [38], we obtain

$$\beta'(\delta; \theta_1) = \int_X \delta(x) p'_\theta(x; \theta_1) \mu(dx) = \int_X \delta(x) S(x; \theta_1) p(x; \theta_1) \mu(dx).$$

This together with the equality $\beta'(\delta; \theta_1) = 0$ implies (1.3.23). □

REMARK 1.3.5. The regularity conditions (R) can be found in [38] (also see [7] or [9]). Lemma 1.3.3 implies that an unbiased test of level α for distinguishing the hypotheses $H_1: \theta = \theta_1$ and $H_2: \theta \neq \theta_1$ is a solution of the following two equations:

$$(1.3.24) \quad E_{\theta_1} \delta(\xi) = \alpha, \quad E_{\theta_1} \delta(\xi) S(\xi; \theta_1) = 0.$$

EXAMPLE 1.3.4. Let the distribution of ξ be exponential with a density of the form (1.3.17) where the functions $a(\theta)$ and $V(\theta)$ are differentiable. Then

$$S(x; \theta) = a'(\theta)T(x) + V'(\theta).$$

Since $E_\theta S(\xi; \theta) = 0$, we get $V'(\theta) = -a'(\theta)E_\theta T(\xi)$. Thus

$$E_\theta \delta(\xi) S(\xi; \theta) = a'(\theta)E_\theta \delta(\xi) T(\xi) - a'(\theta)E_\theta \delta(\xi) E_\theta T(\xi).$$

This implies that equations (1.3.24) become of the form

$$E_{\theta_1}(\delta(\xi) - \alpha) = 0, \quad E_{\theta_1}(\delta(\xi) - \alpha)T(\xi) = 0.$$

The following result describes the form of an unbiased UMP test for distinguishing a null hypothesis and its two-sided alternative hypothesis for exponential families.

THEOREM 1.3.6. *Let $p(x; \theta)$ be of the form (1.3.17) where the function $a(\theta)$ is monotone. Assume that the problem is to distinguish the hypotheses $H_1: \theta \in [\theta_1, \theta_2]$ and $H_2: \theta \notin [\theta_1, \theta_2]$ where $\theta_1 \leq \theta_2$. Then*

1. *In the class of tests \hat{K}_α there exists a UMP test $\hat{\delta}$ such that*

$$(1.3.25) \quad \hat{\delta}(x) = I(T(x) \notin [c_1, c_2]) + q_1 I(T(x) = c_1) + q_2 I(T(x) = c_2)$$

where the constants c_i and q_i , $i = 1, 2$, are defined by

$$(1.3.26) \quad E_{\theta_i} \hat{\delta}(\xi) = \alpha, \quad i = 1, 2,$$

if $\theta_1 < \theta_2$, and by

$$(1.3.27) \quad E_{\theta_1} \hat{\delta}(\xi) = \alpha, \quad E_{\theta_1}(\hat{\delta}(\xi) - \alpha)T(\xi) = 0$$

if $\theta_1 = \theta_2$.

2. *The test $\hat{\delta}$ defined by (1.3.25)–(1.3.27) minimizes the power function $\beta(\delta; \theta)$ inside the interval $[\theta_1, \theta_2]$ if (1.3.26) holds and maximizes it outside the interval $[\theta_1, \theta_2]$ if (1.3.26) holds and $\theta_1 < \theta_2$ or (1.3.27) holds and $\theta_1 = \theta_2$.*
3. *If $0 < \alpha < 1$ and $\theta_1 < \theta_2$, then the function $\beta(\hat{\delta}; \theta)$ attains its minimum at some point $\theta_0 \in (\theta_1, \theta_2)$; moreover it strictly decreases in both cases: either the argument goes to the left of θ_0 or it goes to the right of θ_0 . Note that the case where the distribution of $T(\xi)$ is concentrated at two points is not excluded, that is, we do not exclude the case where there are t_1 and t_2 such that for all θ*

$$P_\theta\{T(\xi) = t_1\} + P_\theta\{T(\xi) = t_2\} = 1.$$

Theorem 1.3.6 is similar to Theorem 1.3.4; however the words “minimizes” and “maximizes” are interchanged and the case of $\theta_1 = \theta_2$ is not excluded in Theorem 1.3.6. The proof of Theorem 1.3.6 is omitted (it can be found in [7, 9], or [34]).

EXAMPLE 1.3.5. Let $\xi = (\xi_1, \xi_2, \dots, \xi_n)$ where $\xi_1, \xi_2, \dots, \xi_n$ are independent identically distributed random variables with the $\mathcal{N}(0, \sigma^2)$ distribution. Consider the hypotheses $H_1: \sigma = \sigma_0$ and $H_2: \sigma \neq \sigma_0$ where $\sigma_0 \in (0, \infty)$ is a fixed number. The density of the distribution belongs to the family of densities

$$(2\pi\sigma)^{-n/2} \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^n x_i^2\right), \quad \sigma \in (0, \infty), \quad x = (x_1, \dots, x_n).$$

The statistic $T(x) = \sum_{i=1}^n x_i^2$ is essential for this family and the density of $T(\xi)$ is $\sigma^{-2} f_n(y/\sigma^2)$ where

$$f_n(y) = \frac{1}{2^{n/2} \Gamma(n/2)} y^{(n/2)-1} e^{-y/2}, \quad y > 0,$$

is the $\chi^2(n)$ density with n degrees of freedom.

An unbiased UMP test $\delta^*(x)$ of level α is of the form (1.3.25) where the constants c_1 and c_2 satisfy condition (1.3.27), while the constants q_1 and q_2 are arbitrary, since the distribution of the statistic $T(\xi)$ is continuous. Putting $q_1 = q_2 = 1$ one can represent the test $\delta^*(x)$ as

$$\delta^*(x) = I\left(\frac{1}{\sigma_0^2} \sum_{i=1}^n x_i^2 \notin (C_1, C_2)\right)$$

where $C_i = c_i/\sigma_0^2$, $i = 1, 2$. Then condition (1.3.27) becomes of the form

$$(1.3.28) \quad \int_{C_1}^{C_2} f_n(y) dy = 1 - \alpha, \quad \int_{C_1}^{C_2} y f_n(y) dy = n(1 - \alpha).$$

To determine the constants c_1 and c_2 one can use, for example, the tables of the $\chi^2(n)$ distribution. Using the equality $y f_n(y) = n f_{n+2}(y)$ the second equation in (1.3.28) can be rewritten as

$$\int_{C_1}^{C_2} f_{n+2}(y) dy = 1 - \alpha.$$

Another way to determine C_1 and C_2 is to integrate by parts the second equation in (1.3.28) and obtain

$$C_1^{n/2} e^{-C_1/2} = C_2^{n/2} e^{-C_2/2}.$$

A relationship between tests and confidence sets. Let ξ be an observation whose distribution belongs to a family $\{P_\theta; \theta \in \Theta\}$.

DEFINITION 1.3.8. A random set $\Theta^*(\xi, \gamma)$ is called a *confidence set of level γ* if $\Theta^*(\xi, \gamma) \subseteq \Theta$ and

$$(1.3.29) \quad P_\theta\{\theta \in \Theta^*(\xi, \gamma)\} \geq \gamma$$

for all $\theta \in \Theta$.

Put

$$(1.3.30) \quad X(\theta, \gamma) = \{x \in X: \theta \in \Theta^*(x, \gamma)\}.$$

Then the inclusions

$$(1.3.31) \quad \theta \in \Theta^*(x, \gamma) \quad \text{and} \quad x \in X(\theta, \gamma)$$

are equivalent.

In the definition of a confidence set we assume that the set $X(\theta, \gamma)$ in (1.3.30) is measurable, thus the probability in (1.3.29) is well defined. In view of the equivalence of inclusions (1.3.31), the latter probability is equal to

$$(1.3.32) \quad P_\theta\{\theta \in \Theta^*(\xi, \gamma)\} = P_\theta\{\xi \in X(\theta, \gamma)\}.$$

The following result describes a relationship between confidence sets and statistical tests for distinguishing the hypotheses $H(\theta_0): \theta = \theta_0$ and

$$K(\theta_0): \theta \in \Theta \setminus \{\theta_0\} = \Theta(\theta_0).$$

THEOREM 1.3.7. 1) For every θ_0 let a nonrandomized test $\delta(\theta_0)$ of level $1 - \gamma$ be given for distinguishing the hypotheses $H(\theta_0)$ and $K(\theta_0)$. Let $X(\theta_0, \gamma)$ be the acceptance set for the hypothesis $H(\theta_0)$ defined by (1.3.30). Then

$$\Theta^*(\xi; \gamma) = \{\theta \in \Theta: \xi \in X(\theta_0, \gamma)\}$$

is a confidence set of level γ .

Conversely, if $\Theta^*(\xi; \gamma)$ is a confidence set of level γ and $\theta_0 \in \Theta^*(\xi; \gamma)$, then the acceptance set $X(\theta_0, \gamma)$ defined by (1.3.30) for the hypothesis $H(\theta_0)$ determines a test for distinguishing the hypotheses $H(\theta_0)$ and $K(\theta_0)$.

2) Let $X(\theta_0, \gamma)$ be the set defined by (1.3.30) for the hypothesis $H(\theta_0)$. If $\delta(\theta_0)$ is a UMP test of level $1 - \gamma$ for all θ_0 , then the corresponding set $\Theta^*(\xi, \gamma)$ minimizes the probability

$$(1.3.33) \quad P_\theta\{\theta' \in \Theta^*(\xi, \gamma)\} \quad \text{for all } \theta \text{ and } \theta' \text{ such that } \theta \in \Theta(\theta')$$

in the class of all confidence sets of level γ .

Conversely, the minimal probability in (1.3.33) corresponds to a set $X(\theta, \gamma)$ that generates a UMP test.

PROOF. Equality (1.3.32) yields

$$P_\theta\{\theta \in \Theta^*(\xi, \gamma)\} = P_\theta\{\xi \in X(\theta, \gamma)\} \geq \gamma,$$

whence the first statement of the theorem follows. To prove the second statement of the theorem we consider another confidence set $\tilde{\Theta}^*(\xi, \gamma)$ and the corresponding subset $\tilde{X}(\theta, \gamma)$ in X . Then

$$P_\theta\{\xi \in \tilde{X}(\theta, \gamma)\} = P_\theta\{\theta \in \tilde{\Theta}^*(\theta, \gamma)\} \geq \gamma.$$

Since $X(\theta_0, \gamma)$ is the acceptance set for a UMP test,

$$P_\theta\{\xi \in \tilde{X}(\theta_0, \gamma)\} \geq P_\theta\{\xi \in X(\theta_0, \gamma)\}$$

for all $\theta \in \Theta(\theta_0)$. Thus

$$P_\theta\{\theta_0 \in \tilde{\Theta}^*(\xi, \gamma)\} \geq P_\theta\{\theta_0 \in \Theta^*(\xi, \gamma)\}$$

for all $\theta \in \Theta(\theta_0)$. □

DEFINITION 1.3.9. Confidence sets for which the probability in (1.3.33) is minimal given (1.3.29) are called *uniformly most precise confidence sets* of level γ with respect to alternatives θ' such that $\theta \in \Theta(\theta')$.

Consider in more detail the notion introduced above for a particular case of a one-dimensional parameter. The following result holds in this case.

THEOREM 1.3.8. *Let the set $X(\theta, \gamma)$ for a UMP test described in Theorem 1.3.7 be of the form*

$$c_1(\theta, \gamma) \leq T(x) \leq c_2(\theta, \gamma)$$

where $c_i(\theta, y)$ are monotone and continuous in θ . If the functions $c_i(\theta, y)$ increase in θ , then a uniformly most precise confidence set of level γ with respect to alternatives θ' such that $\theta \in \Theta(\theta')$ is the interval

$$c_2^{-1}(T(x), \gamma) \leq \theta \leq c_1^{-1}(T(x), \gamma)$$

where $c_i^{-1}(t, y)$ is a solution in θ of the equation $c_i(\theta, y) = t$.

The proof of Theorem 1.3.8 is obvious and omitted.

Consider in more detail *one-sided confidence intervals for a one-dimensional parameter* θ , namely we consider the intervals $(\underline{\theta}(\xi, \gamma), \infty)$ and $(-\infty, \bar{\theta}(\xi, \gamma))$. We restrict our consideration to the case of a *lower confidence bound* $\underline{\theta}(\xi, \gamma)$ for which

$$(1.3.34) \quad P_\theta\{\underline{\theta}(\xi, \gamma) \leq \theta\} \geq \gamma,$$

since an *upper confidence bound* $\bar{\theta}(\xi, \gamma)$ is considered similarly.

DEFINITION 1.3.10. A bound $\underline{\theta} = \underline{\theta}(\xi, \gamma)$ for which the probability $P_\theta\{\underline{\theta} \leq \theta'\}$ is minimal for all $\theta' < \theta$ is called a *uniformly most precise lower confidence bound of level γ* .

Below we consider another definition of an optimal confidence interval. Let $L(\theta, \underline{\theta})$ be a loss arising if θ is underestimated, so that for all fixed θ the function $L(\theta, \underline{\theta})$ is defined as follows: $L(\theta, \underline{\theta}) = 0$ for $\underline{\theta} \geq \theta$ and $L(\theta, \underline{\theta}) \geq 0$ for $\underline{\theta} < \theta$. We also assume that $L(\theta, \underline{\theta})$ is continuously increasing if $\underline{\theta}$ goes away from θ and that $E_\theta L(\theta, \underline{\theta}) < \infty$ for all θ . Given (1.3.34) our goal is to minimize $E_\theta L(\theta, \underline{\theta})$.

The following auxiliary result establishes a relationship between two notions of optimality.

LEMMA 1.3.4. *Given (1.3.34) a uniformly most precise lower bound $\underline{\theta}$ minimizes $E_\theta L(\theta, \underline{\theta})$ for any loss function $L(\theta, \underline{\theta})$ satisfying the above conditions.*

PROOF. Let $\underline{\theta}'$ be an arbitrary lower bound. Since the increments of the loss function $L(\theta, u)$ in u in the domain $u < \theta$ are negative, we get

$$\begin{aligned} E_\theta L(\theta, \underline{\theta}) &= \int_{-\infty}^{\theta} L(\theta, u) d_u P_\theta(\underline{\theta} < u) = - \int_{-\infty}^{\theta} P_\theta(\underline{\theta} < u) d_u L(\theta, u) \\ &\leq - \int_{-\infty}^{\theta} P_\theta(\underline{\theta}' < u) d_u L(\theta, u) = E_\theta L(\theta, \underline{\theta}') \end{aligned}$$

where d_u is the differential with respect to the variable u . □

It is natural to call the number $E_\theta L(\theta, \underline{\theta})$ a *risk of the underestimation of the parameter θ* . Therefore Lemma 1.3.4 implies that a uniformly most precise lower

bound θ minimizes the risk of the underestimation of the parameter θ . This together with Theorems 1.3.7 and 1.3.8 allows one to construct a uniformly most precise one-sided interval in an explicit form for families with monotone likelihood ratio.

THEOREM 1.3.9. *Let a family $\{P_\theta, \theta \in \Theta\}$ have the monotone likelihood ratio with respect to a statistic $T(x)$ whose distribution function $G_\theta(t) = P_\theta\{T(\xi) < t\}$ is continuous in θ and t . Then the distribution of the statistic $T(x)$ monotonically and continuously depends on the parameter θ , that is, $G_\theta(t)$ continuously decreases if θ increases (see relation (5.4.1) in [38]). If $b(t, \gamma)$ is a solution in θ of the equation $G_\theta(t) = \gamma$, then a uniformly most precise lower bound $\underline{\theta}(\xi, \gamma)$ of level γ is*

$$\underline{\theta}(\xi, \gamma) = b(T(\xi), \gamma).$$

PROOF. We put $\Theta(\theta) = \{t: t > \theta\}$ in Theorems 1.3.7 and 1.3.8. According to Theorem 1.3.3 there exists a nonrandomized UMP test for distinguishing the hypotheses $H_1: \theta = \theta_0$ and $H_2: \theta > \theta_0$. Moreover the acceptance set for the hypothesis H_1 is $X(\theta_0, \gamma) = \{x: T(x) < c\}$ where the constant $c = c(\theta_0, \gamma) = G_{\theta_0}^{-1}(\gamma)$ is such that

$$P_{\theta_0}\{T(\xi) < c(\theta_0, \gamma)\} = \gamma.$$

By the assumptions of the theorem we have

$$P_\theta\{T(\xi) \geq c\} > 1 - \gamma = P_{\theta_0}\{T(\xi) \geq c\}$$

for all $\theta > \theta_0$. The latter relation means that $c(\theta_0, \gamma) < c(\theta, \gamma)$ for $\theta_0 < \theta$, that is, the function $c(\theta, \gamma)$ increases in θ . The continuity of $c(\theta, \gamma) = G_\theta^{-1}(\gamma)$ in θ follows from that of $G_\theta(t)$.

Thus the conditions of Theorems 1.3.7 and 1.3.8 hold for $c_2(\theta, \gamma) = c(\theta, \gamma)$ and therefore a uniformly most precise confidence interval is $(c^{-1}(T(\xi); \gamma), \infty)$ where obviously $c^{-1}(t, \gamma) = b(t, \gamma)$. \square

A uniformly most precise upper bound $\bar{\theta}(\xi, \gamma)$ can be constructed by the same method if the assumptions of Theorem 1.3.9 hold.

Now let $\underline{\theta}(\xi, \gamma_1) < \bar{\theta}(\xi, \gamma_2)$ where $\underline{\theta}(\xi, \gamma_1)$ and $\bar{\theta}(\xi, \gamma_2)$ are lower and upper bounds of levels γ_1 and γ_2 , respectively. Let γ_1 and γ_2 be such that the events $\{\underline{\theta}(\xi, \gamma_1) > \theta\}$ and $\{\bar{\theta}(\xi, \gamma_2) < \theta\}$ are disjoint. Then

$$P_\theta\{\underline{\theta}(\xi, \gamma_1) < \theta < \bar{\theta}(\xi, \gamma_2)\} \geq \gamma_1 + \gamma_2 - 1,$$

that is, $(\underline{\theta}(\xi, \gamma_1), \bar{\theta}(\xi, \gamma_2))$ is a confidence interval of level $\gamma_1 + \gamma_2 - 1$.

Let $L_1(\theta, \underline{\theta})$ and $L_2(\theta, \bar{\theta})$ be the loss functions due to the underestimation of the parameter θ for bounds $\underline{\theta}$ and $\bar{\theta}$, respectively. Assume that $L_1(\theta, \underline{\theta})$ and $L_2(\theta, \bar{\theta})$ satisfy the conditions indicated above.

LEMMA 1.3.5. *Let $L(\theta, \underline{\theta}, \bar{\theta}) = L(\theta, \underline{\theta}) + L(\theta, \bar{\theta})$. Then the confidence interval $(\underline{\theta}, \bar{\theta})$ whose end points are uniformly most precise lower and upper bounds minimizes $E_\theta L(\theta, \underline{\theta}, \bar{\theta})$ under the conditions*

$$P_\theta\{\underline{\theta} > \theta\} \leq 1 - \gamma_1, \quad P_\theta\{\bar{\theta} < \theta\} \leq 1 - \gamma_2.$$

This result is an obvious corollary of Lemma 1.3.4. Using Theorem 1.3.9 and Lemma 1.3.5 one can construct optimal intervals in an explicit form for families with the monotone likelihood ratio.

To conclude this section we show how to construct confidence sets and intervals with the help of unbiased tests.

As before let a set $\Theta(\theta)$ correspond to every θ such that $\theta \notin \Theta(\theta)$.

DEFINITION 1.3.11. A confidence set $\Theta^*(\xi, \gamma)$ of level γ for a parameter θ is called *unbiased with respect to the alternative hypothesis* θ' such that $\theta \in \Theta(\theta')$ if

$$(1.3.35) \quad P_{\theta}\{\theta' \in \Theta^*(\xi, \gamma)\} \leq \gamma \quad \text{for all } \theta \text{ and } \theta' \text{ such that } \theta \in \Theta(\theta').$$

The set $\Theta^*(\xi, \gamma)$ is called *unbiased* if (1.3.35) holds for all $\theta' \neq \theta$.

If a confidence set is unbiased, then the *probability that it contains a wrong value θ' of the parameter is less than or equal to the probability that it contains the true value of the parameter.*

DEFINITION 1.3.12. If conditions (1.3.29) and (1.3.35) hold, then a confidence set for which the probability (1.3.33) is minimal is called a *uniformly most precise unbiased confidence set* of level γ against the alternatives θ' such that $\theta \in \Theta(\theta')$.

THEOREM 1.3.10. 1) *Since inclusions (1.3.31) are equivalent, unbiased non-randomized tests generate unbiased confidence sets, and vice versa.*

2) *If $X(\theta_0, \gamma)$ for any $\theta_0 \in \Theta$ is the acceptance set of the null hypothesis*

$$H_1: \theta = \theta_0$$

for a nonrandomized UMP test against the alternative $H_2: \theta \in \Theta(\theta_0)$, then the corresponding set $\Theta^(\xi, \gamma)$ is a uniformly most precise unbiased confidence set, and vice versa.*

PROOF. The method of proof is the same as that of Theorem 1.3.7. An additional reasoning is that if a test is unbiased, then so is the corresponding confidence set, and vice versa. Indeed relations (1.3.29) and (1.3.35) are equivalent to inequalities

$$\sup_{\theta \in \Theta(\theta_0)} P_{\theta}\{\xi \in X(\theta_0, \gamma)\} \leq \gamma \leq P_{\theta_0}\{\xi \in X(\theta_0, \gamma)\}.$$

If δ is a nonrandomized test (that is, $\delta(x) = 0$ for $x \in X(\theta_0, \gamma)$), then

$$E_{\theta}\delta(\xi) = 1 - P_{\theta}\{\xi \in X(\theta_0, \gamma)\}, \quad \inf_{\theta \in \Theta(\theta_0)} E_{\theta}\delta(\xi) \geq 1 - \gamma \geq E_{\theta_0}\delta(\xi).$$

These conditions obviously mean that the test is unbiased and this property is equivalent to (1.3.35). \square

Applying Theorem 1.3.10 one can construct a uniformly most precise unbiased confidence interval for a parameter of the exponential family. The method of this construction is the same as above.

EXAMPLE 1.3.6. Let $\xi = (\xi_1, \xi_2, \dots, \xi_n)$ where $\xi_1, \xi_2, \dots, \xi_n$ are independent identically distributed random variables with the density

$$\theta^{-1}I(x > 0) \exp(-x/\theta), \quad 0 < \theta < \infty.$$

The statistic $T_n(x) = \sum_{i=1}^n x_i$, $x = (x_1, \dots, x_n)$, is a minimal essential statistic and moreover the distribution of $T_n(\xi)/(2\theta)$ is $\chi^2(2n)$. Denote by $\chi_{\gamma}^2(2n)$ the γ -quantile of the $\chi^2(2n)$ law. Then

$$P_{\theta}\{\theta \geq T_n(\xi)/(2\chi_{\gamma}^2(2n))\} = \gamma,$$

that is, $\underline{\theta} = T_n(\xi)/(2\chi_\gamma^2(2n))$ is a lower confidence bound of level γ for the parameter θ . Further let $\zeta_{n,1} = \min_{1 \leq i \leq n} \xi_i$ be the minimal order statistic. Then the distribution of $n\zeta_{n,1}/(2\theta)$ is $\chi^2(2)$. This implies that $\tilde{\theta} = n\zeta_{n,1}/(2\chi_\gamma^2(2))$ also is a lower bound of level γ for the parameter θ .

Which of these two bounds is better? To answer this question note that the distribution of ξ has the monotone likelihood ratio with respect to the statistic $T_n(x)$. Thus the UMP test of level α for distinguishing the hypotheses $H_1: \theta = \theta_0$ and $H_2: \theta > \theta_0$ is determined by the acceptance set of the hypothesis H_1 which is of the form $X(\theta_0, 1 - \alpha) = \{x: T_n(x) \leq 2\theta_0\chi_{1-\alpha}^2(2n)\}$. Thus a uniformly most precise lower bound of level γ for the parameter θ is $\underline{\theta} = T_n(\xi)/(2\chi_\gamma^2(2n))$. Therefore

$$P_\theta \{\theta' \geq \underline{\theta}\} < P_{\theta'} \{\theta' \geq \tilde{\theta}\} \quad \text{for all } \theta \text{ and } \theta' \text{ such that } \theta > \theta'.$$

This implies that the bound $\underline{\theta}$ is better, since it is a uniformly most precise lower bound.

More details on confidence sets and intervals and on an application of statistical tests to construct confidence sets as well as various examples can be found in [54].

Bayes tests for distinguishing a finite number of composite hypotheses. The last topic of this section is the problem of distinguishing N composite hypotheses for $N \geq 2$. As before let the distribution of an observation ξ belong to a family $\mathcal{P} = \{P_\theta, \theta \in \Theta\}$ and let $\mathcal{P} = \bigcup_{i=1}^N \mathcal{P}_i$ where $\mathcal{P}_i = \{P_\theta: \theta \in \Theta_i\}$, $i = 1, 2, \dots, N$; $\Theta_i \cap \Theta_j = \emptyset$, $i \neq j$; $\bigcup_{i=1}^N \Theta_i = \Theta$. Assume that at least one of the sets $\Theta_1, \Theta_2, \dots, \Theta_N$ contains at least two points. Let the hypothesis H_i be that the distribution of ξ belongs to the set \mathcal{P}_i . We write in this case $H_i: \theta \in \Theta_i$. Consider a randomized test δ whose values are $1, 2, \dots, N$ and the corresponding probabilities are $q_i^\delta(x) = P\{\delta = i/\xi = x\}$. The event $\{\delta = i\}$ means that the hypothesis H_i is accepted, and $q_i^\delta(x)$ is the conditional probability of accepting the hypothesis H_i according to the test δ given $\xi = x$.

Let $A_i(t)$, $i = 1, 2, \dots, N$, be nonnegative functions defined on Θ that are measurable with respect to some σ -algebra of subsets of Θ . Let A be a random variable treated as the loss; it is equal to $A_i(t)$ if the hypothesis H_i is accepted and the parameter θ is t . Further we assume that a probability measure \mathbf{Q} (a priori measure) is given on Θ , so that the parameter can be treated to be random with the distribution \mathbf{Q} . Then we define the mean loss or *risk of the test* δ (see (1.3.4) and (1.3.5)):

$$(1.3.36) \quad R(\delta) = E^\delta A = \int_{\Theta} E_t^\delta A \mathbf{Q}(dt)$$

where

$$(1.3.37) \quad E_t^\delta A = \sum_{i=1}^N A_i(t) P_t\{\delta(\xi) = i\}.$$

It is obvious that

$$(1.3.38) \quad P_t\{\delta(\xi) = i\} = \int_{\mathcal{X}} q_i^\delta(x) P_t(dx).$$

Combining (1.3.36)–(1.3.38) we obtain

$$(1.3.39) \quad R(\delta) = \int_{\Theta} \sum_{i=1}^N A_i(t) \int_X q_i^\delta(x) P_t(dx) Q(dt).$$

According to Definition 1.3.3, a test $\delta_{A,Q}$ is Bayes with respect to a loss A and a priori measure Q if it minimizes the risk (1.3.39). Let us find the test $\delta_{A,Q}$ in an explicit form.

Let the family \mathcal{P} be dominated by some σ -finite measure μ and let

$$p(x; t) = dP_t/d\mu(x)$$

be the density of the measure P_t with respect to μ . Then the risk $R(\delta)$ can be written as

$$(1.3.40) \quad R(\delta) = \int_X \sum_{i=1}^N q_i^\delta(x) \int_{\Theta} A_i(t) p(x; t) Q(dt) \mu(dx).$$

Put

$$(1.3.41) \quad R_i(x) = \int_{\Theta} A_i(t) p(x; t) Q(dt), \quad i = 1, 2, \dots, N, \quad x \in X.$$

It follows from (1.3.40) that the test δ minimizes the risk if

$$(1.3.42) \quad q_i^\delta(x) = \begin{cases} 1, & \text{if } R_i(x) = \min_{1 \leq j \leq N} R_j(x), \\ 0, & \text{otherwise,} \end{cases}$$

where $R_i(x)$ are defined by (1.3.41). If $\min_{1 \leq j \leq N} R_j(x)$ is attained for several indices i_1, \dots, i_k , then one can proceed in the same way as in the case of simple hypotheses (see Section 1.2). Thus equality (1.3.42) defines a Bayes test $\delta_{A,Q}$.

More details about Bayes tests for distinguishing composite hypotheses can be found in [9, 54].

EXAMPLE 1.3.7. Let the distribution of ξ be $\mathcal{N}(\theta, 1)$ and $\theta \in \Theta = (-\infty, \infty)$. Let $\Theta = \Theta_1 \cup \Theta_2 \cup \Theta_3$ where

$$\Theta_1 = (-\infty, -1), \quad \Theta_2 = [-1, 1], \quad \Theta_3 = (1, \infty).$$

Consider the loss functions

$$A_1(t) = I(t \geq -1), \quad A_2(t) = I(|t| \geq 1), \quad A_3(t) = I(t \leq 1).$$

Let the a priori distribution Q of the parameter θ be $\mathcal{N}(0, \tau^2)$. Then the distribution of ξ is $\mathcal{N}(0, 1 + \tau^2)$, whence

$$p(x; t) Q(dt) = n(x; 0, 1 + \tau^2) n\left(t; x \frac{\tau^2}{1 + \tau^2}, \frac{\tau^4}{1 + \tau^2}\right) dt$$

where $n(x; a, b^2)$ is the density of the $\mathcal{N}(a, b^2)$ law. Thus in view of (1.3.41) we get

$$\begin{aligned} R_1(x) &= n(x; 0, 1 + \tau^2) \int_{-1}^{\infty} n\left(t; x \frac{\tau^2}{1 + \tau^2}, \frac{\tau^4}{1 + \tau^2}\right) dt \\ &= n(x; 0, 1 + \tau^2) \Phi\left(\frac{1 + x^*}{\sigma}\right) \end{aligned}$$

where

$$x^* = x \frac{\tau^2}{1 + \tau^2}, \quad \sigma^2 = \frac{\tau^4}{1 + \tau^2}.$$

Similarly we get

$$\begin{aligned} R_2(x) &= n(x; 0, 1 + \tau^2) \int_{|x| > 1} n(t; x^*, \sigma^2) dt \\ &= n(x; 0, 1 + \tau^2) \left[2 - \Phi\left(\frac{1 + x^*}{\sigma}\right) - \Phi\left(\frac{1 - x^*}{\sigma}\right) \right], \\ R_3(x) &= n(x; 0, 1 + \tau^2) \Phi\left(\frac{1 - x^*}{\sigma}\right). \end{aligned}$$

Let

$$R_j^*(x) = R_j(x) / n(x; 0, 1 + \tau^2).$$

The function $R_2^*(x)$ is symmetric with respect to x^* ; its minimum is 0 and it is attained at the point $x = 0$. Further

$$(1.3.43) \quad \frac{d}{dx^*} R_2^*(x) = -\frac{1}{\sigma} \left[\varphi\left(\frac{1 + x^*}{\sigma}\right) - \varphi\left(\frac{1 - x^*}{\sigma}\right) \right]$$

where $\varphi(x)$ is the density of the $\mathcal{N}(0, 1)$ law.

For all $x^* < 0$ it holds that $|1 + x^*| < 1 - x^*$. Hence $\varphi((1 + x^*)/\sigma) > \varphi((1 - x^*)/\sigma)$ for all $x^* < 0$. Equality (1.3.43) implies that $R_2^*(x)$ is decreasing in $(-\infty, 0)$ and $\lim_{x \rightarrow -\infty} R_2^*(x) = 1$. The function $R_1^*(x)$ is increasing and

$$\lim_{x \rightarrow -\infty} R_1^*(x) = 0, \quad \lim_{x \rightarrow \infty} R_1^*(x) = 1.$$

Thus there exists a unique point x' in the interval $(-\infty, 0)$ such that

$$R_1^*(x') = R_2^*(x').$$

Since the problem is symmetric, there exists a unique point $x'' = -x'$ in the interval $(0, \infty)$ such that $R_2^*(x'') = R_3^*(x'')$. Thus the partition $X = X_1 \cup X_2 \cup X_3$ where $X_1 = (-\infty, x')$, $X_2 = [x', x'']$, and $X_3 = (x'', \infty)$ determines a Bayes test. According to this test the hypothesis H_i , $i = 1, 2, 3$, is accepted if $x \in X_i$.

REMARK 1.3.6. More details about distinguishing composite hypotheses can be found in [9, 34, 54]. Distinguishing composite hypotheses can be viewed as a problem of the general theory of statistical decisions; see [4, 9, 15, 52, 54]. Asymptotic problems of distinguishing composite hypotheses for independent observations are considered in [10]. Tests of significance are studied in Chapter 3 below.

Asymptotic Distinguishability of Simple Hypotheses

2.1. Statistical hypotheses and tests

Let ξ^t , $t \in \mathbf{R}_+$, be a family of observations assuming values in a measurable space (X^t, \mathcal{B}^t) and let $\mathcal{P}^t = \{\mathbf{P}^t, \tilde{\mathbf{P}}^t\}$ be a family of two probability measures defined on (X^t, \mathcal{B}^t) . Let H^t and \tilde{H}^t be two statistical hypotheses that the distribution of the observation ξ^t is generated by the measures \mathbf{P}^t and $\tilde{\mathbf{P}}^t$, respectively. Denote by δ_t a measurable mapping of the space (X^t, \mathcal{B}^t) into the space $([0, 1], \mathcal{B}([0, 1]))$. The mapping δ_t (as well as the random variable $\delta(\xi^t)$ denoted by the same symbol δ_t) is called a *test for distinguishing the hypotheses H^t and \tilde{H}^t by the observation ξ^t* ; here we treat $\delta_t(x)$ as the conditional probability of rejecting the hypothesis H^t (or, equivalently, of accepting the hypothesis \tilde{H}^t) given $\xi^t = x$. Denote by Σ^t the set of all tests δ_t for distinguishing the hypotheses H^t and \tilde{H}^t ; for all $\delta_t \in \Sigma^t$ we introduce the *type I and type II error probabilities*

$$(2.1.1) \quad \alpha(\delta_t) = \mathbf{E}^t \delta_t \quad \text{and} \quad \beta(\delta_t) = \tilde{\mathbf{E}}^t(1 - \delta_t),$$

respectively, where \mathbf{E}^t and $\tilde{\mathbf{E}}^t$ are expectations with respect to the distributions \mathbf{P}^t and $\tilde{\mathbf{P}}^t$, respectively. For all $\alpha \in [0, 1]$ we denote by Σ_α^t the set of all tests δ_t of Σ^t such that $\alpha(\delta_t) \leq \alpha$.

Let $\mathbf{Q}^t = (\mathbf{P}^t + \tilde{\mathbf{P}}^t)/2$ be another probability measure on the measurable space (X^t, \mathcal{B}^t) and let $\mathfrak{z}_t = d\mathbf{P}^t/d\mathbf{Q}^t$ and $\tilde{\mathfrak{z}}_t = d\tilde{\mathbf{P}}^t/d\mathbf{Q}^t$ be two finite versions of the Radon–Nikodym derivatives of the measures \mathbf{P}^t and $\tilde{\mathbf{P}}^t$, respectively, with respect to the measure \mathbf{Q}^t . Consider the likelihood ratios

$$(2.1.2) \quad z_t = \tilde{\mathfrak{z}}_t/\mathfrak{z}_t, \quad \tilde{z}_t = \mathfrak{z}_t/\tilde{\mathfrak{z}}_t.$$

Here we agree that $0/0 = 0$ and the likelihood ratios are well defined for all t . Note that $\mathfrak{z}_t + \tilde{\mathfrak{z}}_t = 2$.

Put

$$(2.1.3) \quad \bar{\alpha}_t = \mathbf{P}^t(\tilde{\mathfrak{z}}_t > 0), \quad \bar{\beta}_t = \tilde{\mathbf{P}}^t(\mathfrak{z}_t > 0).$$

Then

$$(2.1.4) \quad \bar{\alpha}_t = \mathbf{P}^t(z_t > 0) = \mathbf{P}^t(\tilde{z}_t < \infty),$$

$$(2.1.5) \quad \bar{\beta}_t = \tilde{\mathbf{P}}^t(\tilde{z}_t > 0) = \tilde{\mathbf{P}}^t(z_t < \infty).$$

Equalities (2.1.3)–(2.1.5) and Lemma 1.1.9 imply the following Lebesgue decomposition of any of the measures \mathbf{P}^t or $\tilde{\mathbf{P}}^t$ with respect to the other one.

LEMMA 2.1.1. For all $A \in \mathcal{B}^t$ it holds that

$$(2.1.6) \quad \tilde{P}^t(A) = \int_A z_t dP^t + \tilde{P}^t(A \cap \{z_t = \infty\}),$$

$$(2.1.7) \quad P^t(A) = \int_A \tilde{z}_t d\tilde{P}^t + P^t(A \cap \{\tilde{z}_t = \infty\})$$

where z_t and \tilde{z}_t are the likelihood ratios defined by equalities (2.1.2).

Lemma 1.1.10 can also be rewritten in the following form.

LEMMA 2.1.2. For all nonnegative measurable functions η defined on (X^t, \mathcal{B}^t) it holds that

$$(2.1.8) \quad \tilde{E}^t \eta = E^t \eta z_t + \tilde{E}^t \eta I(z_t = \infty),$$

$$(2.1.9) \quad E^t \eta = \tilde{E}^t \eta \tilde{z}_t + E^t \eta I(\tilde{z}_t = \infty).$$

Consider the set

$$(2.1.10) \quad \mathfrak{N}^t = \{(\alpha(\delta_t), \beta(\delta_t)) : \delta_t \in \Sigma^t\}$$

where $\alpha(\delta_t)$ and $\beta(\delta_t)$ are type I and type II error probabilities of the test $\delta_t \in \Sigma^t$ defined by (2.1.1).

The properties of the set \mathfrak{N}^t are studied in Section 1.1. In particular, the set \mathfrak{N}^t is convex, closed, and symmetric about the point $(1/2, 1/2)$, \mathfrak{N}^t contains the points $(0, 1)$ and $(1, 0)$, and $\mathfrak{N}^t \subseteq [0, 1] \times [0, 1]$. An example of the set \mathfrak{N}^t is shown in Figure 2.1.1.

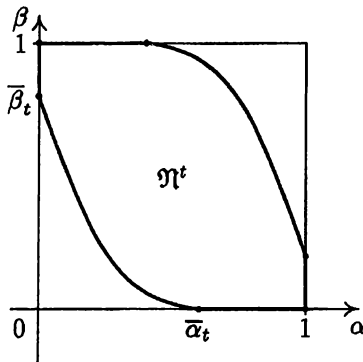


FIGURE 2.1.1

Now we introduce the *likelihood ratio test* by

$$(2.1.11) \quad \delta_t^{c, \varepsilon} = I(z_t > c) + \varepsilon I(z_t = c)$$

where $c \in [0, \infty]$ and $\varepsilon \in [0, 1]$ are parameters of the test. Let $(c_t(\alpha), \varepsilon_t(\alpha))$ be some solution of the equation $\alpha(\delta_t^{c, \varepsilon}) = \alpha$ with respect to (c, ε) .

A likelihood ratio test $\delta_t^{c_t(\alpha), \varepsilon_t(\alpha)}$ with $\varepsilon_t(0) = 1$ is called the *Neyman-Pearson test of level α for distinguishing the hypotheses H^t and \tilde{H}^t* . In what follows we denote this test by $\delta_t^{+, \alpha}$. The functions $c_t(\alpha)$ and $\varepsilon_t(\alpha)$ as well as the type II error probability $\beta(\delta_t^{+, \alpha})$ can be obtained by applying results of Section 1.1.

According to the Neyman–Pearson fundamental lemma, the point $(\alpha, \beta(\delta_t^{+, \alpha}))$ belongs to the boundary of the set \mathfrak{N}^t for all $\alpha \in [0, 1]$. In various cases where the level α depends on t , one can study the behavior of the set \mathfrak{N}^t as $t \rightarrow \infty$ instead of that of $\beta(\delta_t^{+, \alpha})$ as $t \rightarrow \infty$.

Below we consider the Neyman–Pearson test δ_t^{+, α_t} of level α_t that depends, generally speaking, on t . We will use the notation $c_t = c_t(\alpha_t)$ and $\varepsilon_t = \varepsilon_t(\alpha_t)$. Moreover we put $\Lambda_t = \ln z_t$ and $d_t = \ln c_t$ assuming that $\ln 0 = -\infty$. Then the test δ_t^{+, α_t} can be rewritten as

$$(2.1.12) \quad \delta_t^{+, \alpha_t} = I(\Lambda_t > d_t) + \varepsilon_t I(\Lambda_t = d_t).$$

Consider the Bayes test δ_t^π with respect to the a priori distribution $(\pi, 1 - \pi)$ for distinguishing the hypotheses H^t and \tilde{H}^t that minimizes the error probability $e_\pi(\delta_t)$ (see (1.1.32)) in the class Σ^t of all tests δ_t .

Denote by δ_t^* the minimax test for distinguishing the hypotheses H^t and \tilde{H}^t that minimizes $\alpha(\delta_t) \vee \beta(\delta_t)$ in the class Σ^t of all tests δ_t . It follows from Theorem 1.2.4 that the likelihood ratio test $\delta_t^{c, \varepsilon}$ is minimax if $\alpha(\delta_t^{c, \varepsilon}) = \beta(\delta_t^{c, \varepsilon})$. According to (2.1.11) and (2.1.12) the Neyman–Pearson test is minimax if $\beta(\delta_t^{+, \alpha_t}) = \alpha_t$. Moreover we learned in Section 1.2 that the Neyman–Pearson test δ_t^{+, α_t} coincides with some Bayes test $\delta_t^{\pi_t}$, and the probability π_t depends, generally speaking, on t .

Below we study the asymptotic behavior (as $t \rightarrow \infty$) of the Neyman–Pearson test δ_t^{+, α_t} , the minimax δ_t^* test, and the Bayes $\delta_t^{\pi_t}$ test. The asymptotic behavior of these tests depends on the behavior of the set \mathfrak{N}^t which in turn is determined by the behavior of the measures P^t and \tilde{P}^t as $t \rightarrow \infty$ (see Section 1.1). In the next section we consider all the possible types of the asymptotic behavior of the set \mathfrak{N}^t as $t \rightarrow \infty$.

2.2. Types of the asymptotic distinguishability of families of hypotheses. The characterization of types

Consider two families of statistical hypotheses (H^t) and (\tilde{H}^t) and the two corresponding families of probability measures (P^t) and (\tilde{P}^t) . Here and throughout this section the symbols (H^t) and (\tilde{H}^t) stand for $(H^t)_{t \in \mathbf{R}_+}$ and $(\tilde{H}^t)_{t \in \mathbf{R}_+}$, respectively. For other families of measures, random variables, etc. we follow the same notation to make them shorter.

Below we define the types of the asymptotic distinguishability of families of hypotheses (H^t) and (\tilde{H}^t) as $t \rightarrow \infty$. Our approach is based on the asymptotic behavior of sets \mathfrak{N}^t as $t \rightarrow \infty$.

The distance in variance between measures. The Kakutani–Hellinger distance and Hellinger integrals. To state the main results on the characterization of types we need the following notions.

DEFINITION 2.2.1. Let P^t and \tilde{P}^t be two measures. The full variation of the charge $\tilde{P}^t - P^t$ is called *the distance in variance between measures P^t and \tilde{P}^t* and is denoted by $\|\tilde{P}^t - P^t\|$, namely

$$(2.2.1) \quad \|\tilde{P}^t - P^t\| = E_Q^t |\tilde{z}_t - z_t|$$

where E_Q^t stands for the expectation with respect to the measure Q^t .

DEFINITION 2.2.2. Let

$$(2.2.2) \quad \rho^2(P^t, \tilde{P}^t) = \frac{1}{2} E_Q^t \left(\tilde{\mathfrak{z}}_t^{1/2} - \mathfrak{z}_t^{1/2} \right)^2.$$

The nonnegative number $\rho(P^t, \tilde{P}^t)$ is called the *Kakutani–Hellinger distance between the measures P^t and \tilde{P}^t* .

DEFINITION 2.2.3. Let

$$(2.2.3) \quad H(\varepsilon; \tilde{P}^t, P^t) = E_Q^t \tilde{\mathfrak{z}}_t^\varepsilon \mathfrak{z}_t^{1-\varepsilon}.$$

Here we put

$$\tilde{\mathfrak{z}}_t^\varepsilon \mathfrak{z}_t^{1-\varepsilon} = \begin{cases} 0, & \varepsilon < 0, \tilde{\mathfrak{z}}_t = 0, \text{ and } \mathfrak{z}_t = 0, \\ \infty, & \varepsilon < 0, \tilde{\mathfrak{z}}_t = 0, \text{ and } \mathfrak{z}_t > 0, \\ \mathfrak{z}_t I(\tilde{\mathfrak{z}}_t > 0), & \varepsilon = 0, \\ \tilde{\mathfrak{z}}_t I(\mathfrak{z}_t > 0), & \varepsilon = 1, \\ 0, & \varepsilon > 1, \mathfrak{z}_t = 0, \text{ and } \tilde{\mathfrak{z}}_t = 0, \\ \infty, & \varepsilon > 1, \mathfrak{z}_t = 0, \text{ and } \tilde{\mathfrak{z}}_t > 0. \end{cases}$$

Then $H(\varepsilon; \tilde{P}^t, P^t)$ is well defined for all ε and t . The number $H(\varepsilon; \tilde{P}^t, P^t)$ is called the *Hellinger integral of order $\varepsilon \in \mathbf{R} = (-\infty, \infty)$ for the measures \tilde{P}^t and P^t* . The number $H(1/2; \tilde{P}^t, P^t)$ is called the *Hellinger integral for the measures \tilde{P}^t and P^t* and is denoted by $H(\tilde{P}^t, P^t)$.

Properties of $\|\tilde{P}^t - P^t\|$, $\rho(\tilde{P}^t, P^t)$, and $H(\varepsilon; \tilde{P}^t, P^t)$ defined by (2.2.1), (2.2.2), and (2.2.3), respectively, can be found in [28, 33, 35, 47]. In particular, neither $\|\tilde{P}^t - P^t\|$ nor $\rho(\tilde{P}^t, P^t)$ nor $H(\varepsilon; \tilde{P}^t, P^t)$ depend on the dominating measure Q^t .

Below we give an auxiliary result on some relationships between these notions.

LEMMA 2.2.1. *We have*

$$(2.2.4) \quad 2 \left(1 - H(\tilde{P}^t, P^t) \right) \leq \|\tilde{P}^t - P^t\| \leq \sqrt{8 \left(1 - H(\tilde{P}^t, P^t) \right)},$$

$$(2.2.5) \quad \|\tilde{P}^t - P^t\| \leq 2 \sqrt{1 - H^2(\tilde{P}^t, P^t)}.$$

In particular

$$(2.2.6) \quad 2\rho^2(P^t, \tilde{P}^t) \leq \|\tilde{P}^t - P^t\| \leq \sqrt{8} \rho(P^t, \tilde{P}^t).$$

PROOF. By the definitions of the Kakutani–Hellinger distance and Hellinger integral

$$(2.2.7) \quad \rho^2(P^t, \tilde{P}^t) = 1 - H(\tilde{P}^t, P^t).$$

Thus inequalities (2.2.6) follow from (2.2.4). Let us prove inequalities (2.2.4) and (2.2.5).

Since $\mathfrak{z}_t + \tilde{\mathfrak{z}}_t = 2$, it follows from the Jensen inequality that

$$\begin{aligned} \frac{1}{2} \|\tilde{P}^t - P^t\| &= \frac{1}{2} E_Q^t |\tilde{\mathfrak{z}}_t - \mathfrak{z}_t| = E_Q^t |1 - \mathfrak{z}_t| \leq \sqrt{E_Q^t (1 - \mathfrak{z}_t)^2} \\ &= \sqrt{1 - E_Q^t \mathfrak{z}_t (2 - \mathfrak{z}_t)} = \sqrt{1 - E_Q^t \mathfrak{z}_t \tilde{\mathfrak{z}}_t}. \end{aligned}$$

It follows from the Cauchy–Bunyakovskii inequality that

$$E_Q^t \sqrt{\delta_t \tilde{\delta}_t} \leq \sqrt{E_Q^t \delta_t \tilde{\delta}_t},$$

and thus

$$\frac{1}{2} \|\tilde{P}^t - P^t\| \leq \sqrt{1 - E_Q^t \delta_t \tilde{\delta}_t} \leq \sqrt{1 - (E_Q^t \sqrt{\delta_t \tilde{\delta}_t})^2} = \sqrt{1 - H^2(\tilde{P}^t, P^t)},$$

that is, (2.2.5) is proved.

Since $H(\tilde{P}^t, P^t) \leq 1$, we have

$$1 - H^2(\tilde{P}^t, P^t) = (1 - H(\tilde{P}^t, P^t)) (1 + H(\tilde{P}^t, P^t)) \leq 2 (1 - H(\tilde{P}^t, P^t)).$$

This estimate together with (2.2.5) implies the second inequality in (2.2.4).

Since $(a - b)^2 \leq |a^2 - b^2|$ for $a > 0$ and $b > 0$, we get

$$(2.2.8) \quad \frac{1}{2} (\sqrt{z} - \sqrt{2-z})^2 \leq |z - 1|, \quad 0 \leq z \leq 2.$$

Using (2.2.7), equality $\delta_t + \tilde{\delta}_t = 2$, and inequality (2.2.8) we obtain

$$\begin{aligned} 1 - H(\tilde{P}^t, P^t) &= \rho^2(P^t, \tilde{P}^t) = \frac{1}{2} E_Q^t (\sqrt{\delta_t} - \sqrt{1 - \delta_t})^2 \\ &\leq E_Q^t |\delta_t - 1| = \frac{1}{2} E_Q^t |\delta_t - \tilde{\delta}_t| = \frac{1}{2} \|P^t - \tilde{P}^t\|, \end{aligned}$$

that is, the first inequality in (2.2.4) is also proved. \square

Let

$$(2.2.9) \quad \|P^t \wedge \tilde{P}^t\| = \inf \{ \alpha(\delta_t) + \beta(\tilde{\delta}_t) : \delta_t \in \Sigma^t \}.$$

It is clear that $\|P^t \wedge \tilde{P}^t\|$ is the doubled error of the Bayes test with respect to the a priori distribution $(1/2, 1/2)$.

LEMMA 2.2.2. *We have*

$$(2.2.10) \quad \|P^t \wedge \tilde{P}^t\| = 1 - \frac{1}{2} \|\tilde{P}^t - P^t\|.$$

PROOF. Since the Bayes test with respect to the a priori distribution $(1/2, 1/2)$ can be represented as $\delta_t^{1,1} = I(z_t \geq 1)$, relation (2.2.9) implies

$$\begin{aligned} (2.2.11) \quad \|P^t \wedge \tilde{P}^t\| &= E^t I(z_t \geq 1) + \tilde{E}^t (1 - I(z_t \geq 1)) \\ &= 1 + E_Q^t \delta_t I(z_t \geq 1) - E_Q^t \tilde{\delta}_t I(z_t \geq 1) \\ &= 1 - E_Q^t (\tilde{\delta}_t - \delta_t) I(z_t \geq 1). \end{aligned}$$

Since $E_Q^t (\tilde{\delta}_t - \delta_t) = 0$, we obtain from (2.2.1) that

$$\begin{aligned} (2.2.12) \quad \|\tilde{P}^t - P^t\| &= E_Q^t (\tilde{\delta}_t - \delta_t) I(z_t \geq 1) + E_Q^t (\delta_t - \tilde{\delta}_t) I(z_t < 1) \\ &= 2E_Q^t (\tilde{\delta}_t - \delta_t) I(z_t \geq 1). \end{aligned}$$

Now (2.2.11) and (2.2.12) imply equality (2.2.10). \square

REMARK 2.2.1. Similarly to (2.2.11) we prove that

$$\begin{aligned} \|\mathbf{P}^t \wedge \tilde{\mathbf{P}}^t\| &= \mathbf{E}^t I(z_t \geq 1) + \tilde{\mathbf{E}}^t I(z_t < 1) = \mathbf{E}_Q^t \mathfrak{z}_t I(z_t \geq 1) + \mathbf{E}_Q^t \tilde{\mathfrak{z}}_t I(z_t < 1) \\ &= \mathbf{E}_Q^t (\mathfrak{z}_t \wedge \tilde{\mathfrak{z}}_t) \end{aligned}$$

where $a \wedge b$ stands for the minimum of two numbers a and b . The latter equality makes the notation $\|\mathbf{P}^t \wedge \tilde{\mathbf{P}}^t\|$ clear.

The complete asymptotic distinguishability. First we give some necessary definitions.

DEFINITION 2.2.4. Families of hypotheses (H^t) and (\tilde{H}^t) are called *completely asymptotically distinguishable* (denoted by $(H^t) \Delta (\tilde{H}^t)$) if there exist a sequence $t_n \uparrow \infty$, $n \rightarrow \infty$, and tests $\delta_{t_n} \in \Sigma^{t_n}$ such that

$$(2.2.13) \quad \alpha(\delta_{t_n}) \rightarrow 0 \quad \text{and} \quad \beta(\delta_{t_n}) \rightarrow 0 \quad \text{as } n \rightarrow \infty.$$

DEFINITION 2.2.5. Families of measures (\mathbf{P}^t) and $(\tilde{\mathbf{P}}^t)$ are called *completely asymptotically separable* (denoted by $(\mathbf{P}^t) \Delta (\tilde{\mathbf{P}}^t)$) if there exist a sequence $t_n \uparrow \infty$, $n \rightarrow \infty$, and sets $A_n \in \mathcal{B}^{t_n}$ such that

$$(2.2.14) \quad \mathbf{P}^{t_n}(A_n) \rightarrow 0 \quad \text{and} \quad \tilde{\mathbf{P}}^{t_n}(A_n) \rightarrow 1 \quad \text{as } n \rightarrow \infty.$$

The following result contains a characterization of the complete asymptotic distinguishability $(H^t) \Delta (\tilde{H}^t)$.

THEOREM 2.2.1. *The following statements are equivalent:*

- a) $(H^t) \Delta (\tilde{H}^t)$;
- b) $(\mathbf{P}^t) \Delta (\tilde{\mathbf{P}}^t)$;
- c) $\limsup_{t \rightarrow \infty} \tilde{\mathbf{P}}^t(z_t > N) = 1$ for all $N < \infty$;
- d) $\limsup_{t \rightarrow \infty} \mathbf{P}^t(z_t < N) = 1$ for all $N > 0$;
- e) $\liminf_{t \rightarrow \infty} \|\mathbf{P}^t \wedge \tilde{\mathbf{P}}^t\| = 0$;
- f) $\liminf_{t \rightarrow \infty} H(\varepsilon; \tilde{\mathbf{P}}^t, \mathbf{P}^t) = 0$ for all $\varepsilon \in (0, 1)$;
- g) $\limsup_{t \rightarrow \infty} \|\tilde{\mathbf{P}}^t - \mathbf{P}^t\| = 2$;
- h) $\limsup_{t \rightarrow \infty} \rho(\mathbf{P}^t, \tilde{\mathbf{P}}^t) = 1$.

PROOF. a) \Rightarrow b). Let $t_n \uparrow \infty$ as $n \rightarrow \infty$ and let tests $\delta_{t_n} \in \Sigma^{t_n}$ be such that $\alpha(\delta_{t_n}) \rightarrow 0$ and $\beta(\delta_{t_n}) \rightarrow 0$. Put $A_n = I(\delta_{t_n} > \gamma)$ for $0 < \gamma < 1$. By the Chebyshev inequality we have

$$\mathbf{P}^{t_n}(A_n) = \mathbf{P}^{t_n}(\delta_{t_n} > \gamma) \leq \gamma^{-1} \alpha(\delta_{t_n}) \rightarrow 0.$$

Similarly we obtain

$$\tilde{\mathbf{P}}^{t_n}(\delta_{t_n} \leq \gamma) \leq \tilde{\mathbf{P}}^{t_n}(1 - \delta_{t_n} \geq 1 - \gamma) \leq (1 - \gamma)^{-1} \beta(\delta_{t_n}) \rightarrow 0,$$

that is, $\tilde{\mathbf{P}}^{t_n}(A_n) \rightarrow 1$. Thus the implication a) \Rightarrow b) is proved.

b) \Rightarrow d). Let $t_n \uparrow \infty$ as $n \rightarrow \infty$ and let sets $A_n \in \mathcal{B}^{t_n}$ be such that

$$(2.2.15) \quad \mathbf{P}^{t_n}(A_n) \rightarrow 0, \quad \tilde{\mathbf{P}}^{t_n}(A_n) \rightarrow 1.$$

By the Lebesgue decomposition (2.1.7) we have for all $A \in \mathcal{B}^t$ that

$$(2.2.16) \quad \begin{aligned} P^t(A) &= P^t(A \cap \{\tilde{z}_t \leq N\}) + P^t(A \cap \{\tilde{z}_t > N\}) \\ &= \int_{A \cap \{\tilde{z}_t \leq N\}} \tilde{z}_t d\tilde{P}^t + P^t(A \cap \{\tilde{z}_t > N\}) \leq N\tilde{P}^t(A) + P^t(\tilde{z}_t > N) \end{aligned}$$

for $0 < N < \infty$. Relations (2.2.15) and estimate (2.2.16) for $A = A_n^c = X^t \setminus A_n$ and $t = t_n$ imply that $P^{t_n}(\tilde{z}_{t_n} > N) \rightarrow 1$ for all $N \in (0, \infty)$, whence we obtain d), since

$$(2.2.17) \quad P^t(\tilde{z}_t > N) = P^t(z_t < N^{-1}).$$

d) \Rightarrow f). For all $\varepsilon, \lambda \in (0, 1)$ and $\gamma > 0$ we have

$$(2.2.18) \quad \begin{aligned} H(\varepsilon; \tilde{P}^t, P^t) &= E_Q^t \tilde{\delta}_t^\varepsilon \delta_t^{1-\varepsilon} I(\tilde{\beta}_t < \gamma) + E_Q^t \tilde{\delta}_t^\varepsilon \delta_t^{1-\varepsilon} I(\tilde{\beta}_t \geq \gamma, \beta_t \leq \lambda) \\ &\quad + E_Q^t \tilde{\delta}_t^\varepsilon \delta_t^{1-\varepsilon} I(\tilde{\beta}_t \geq \gamma, \beta_t > \lambda) \\ &\leq \gamma^\varepsilon E_Q^t \delta_t^{1-\varepsilon} + \lambda^{1-\varepsilon} E_Q^t \tilde{\delta}_t^\varepsilon + E_Q^t \delta_t^\varepsilon I(\tilde{\beta}_t \geq \gamma, \beta_t > \lambda) \\ &\leq \gamma^\varepsilon + \lambda^{1-\varepsilon} + E^t z_t^\varepsilon I(\tilde{\beta}_t \geq \gamma) \\ &\leq \gamma^\varepsilon + \lambda^{1-\varepsilon} + (E^t z_t)^\varepsilon (E^t I(\tilde{\beta}_t \geq \gamma))^{1-\varepsilon} \\ &\leq \gamma^\varepsilon + \lambda^{1-\varepsilon} + [P^t(\tilde{\beta}_t \geq \gamma)]^{1-\varepsilon}, \end{aligned}$$

since $E_Q^t \delta_t^{1-\varepsilon} \leq 1$ and $E_Q^t \tilde{\delta}_t^\varepsilon \leq 1$ by the Hölder inequality, and $E^t z_t \leq 1$ by (2.1.8) for $\eta = 1$.

Equality (2.2.17) and relation d) imply that $\limsup_{t \rightarrow \infty} P^t(\tilde{z}_t > N) = 1$. Since $\beta_t + \tilde{\beta}_t = 2$, we have $\tilde{z}_t = 2/\tilde{\beta}_t - 1$, whence

$$(2.2.19) \quad \liminf_{t \rightarrow \infty} P^t(\tilde{\beta}_t \geq \gamma) = 0$$

for all $\gamma > 0$. It follows from (2.2.18) and (2.2.19) that

$$\liminf_{t \rightarrow \infty} H(\varepsilon; \tilde{P}^t, P^t) \leq \gamma^\varepsilon + \lambda^{1-\varepsilon}.$$

Since γ and λ are arbitrary, the latter inequality proves f).

f) \Rightarrow h). Follows from (2.2.7).

h) \Rightarrow g). This implication follows from the first inequality in (2.2.6) and

$$\|\tilde{P}^t - P^t\| \leq 2$$

(the latter estimate holds in view of (2.2.1)).

g) \Rightarrow e). Follows from inequality (2.2.10).

e) \Rightarrow a). Relation (2.2.9) and Theorem 1.2.2 for $N = 2$, $A_{ij} = 1 - \delta_{ij}$, and $\pi_1 = \pi_2 = 1/2$ yield

$$(2.2.20) \quad \|P^t \wedge \tilde{P}^t\| = \alpha(\delta_t^{1,1}) + \beta(\delta_t^{1,1})$$

where $\delta_t^{1,1} = I(z_t \geq 1)$. Condition e) and equality (2.2.20) imply that there exists a sequence $t_n \uparrow \infty$, $n \rightarrow \infty$, such that $\alpha(\delta_{t_n}^{1,1}) \rightarrow 0$ and $\beta(\delta_{t_n}^{1,1}) \rightarrow 0$ as $n \rightarrow \infty$, that is, condition a) holds.

Therefore conditions a), b), d), e), f), g), and h) are equivalent. It remains to prove that conditions b) and c) are equivalent.

b) \Rightarrow c). Let $t_n \rightarrow \infty$ as $n \rightarrow \infty$ and let sets $A_n \in \mathcal{B}^{t_n}$ be such that relations (2.2.14) hold. The Lebesgue decomposition (2.1.6) for all $A \in \mathcal{B}^t$ and $0 < N < \infty$ imply similarly to (2.2.16) that

$$\tilde{P}^t(A) \leq NP^t(A) + \tilde{P}^t(z_t > N).$$

This together with (2.2.14) yields c).

c) \Rightarrow b). Condition c) implies that there exists a sequence $t_n \uparrow \infty$, $n \rightarrow \infty$, such that

$$\tilde{P}^{t_n}(z_{t_n} > n) \geq 1 - n^{-1}.$$

Thus $\tilde{P}^{t_n}(z_{t_n} > n) \rightarrow 1$ as $n \rightarrow \infty$. On the other hand, the Chebyshev inequality implies that

$$P^{t_n}(z_{t_n} > n) \leq n^{-1}E^{t_n}z_{t_n} \leq n^{-1},$$

since $E^t z_t \leq 1$ in view of (2.1.8). This implies that $P^{t_n}(z_{t_n} > n) \rightarrow 0$ as $n \rightarrow \infty$ and condition b) is proved. \square

REMARK 2.2.2. If $(H^t) \Delta (\tilde{H}^t)$, then there exist a sequence $t_n \uparrow \infty$, $n \rightarrow \infty$, and tests $\delta_{t_n} \in \Sigma^{t_n}$ such that conditions (2.2.13) hold. Then, obviously,

$$\mathfrak{N}^{t_n} \rightarrow [0, 1] \times [0, 1]$$

as $t \rightarrow \infty$ where the set \mathfrak{N}^t is defined by (2.1.10). On the other hand, if

$$(P^t) \Delta (\tilde{P}^t),$$

then there exist a sequence $t_n \uparrow \infty$, $n \rightarrow \infty$, and sets $A_n \in \mathcal{B}^{t_n}$ such that relations (2.2.14) hold. These relations mean that the sequences of measures P^{t_n} and \tilde{P}^{t_n} , $n = 1, 2, \dots$, are *asymptotically singular* (cf. (1.1.9)).

Now we define the counterparts of the notions of the complete asymptotic distinguishability $(H^t) \Delta (\tilde{H}^t)$ and complete asymptotic separability $(P^t) \Delta (\tilde{P}^t)$.

DEFINITION 2.2.6. We say that *families of hypotheses* (H^t) and (\tilde{H}^t) are *not completely asymptotically distinguishable* (denoted by $(H^t) \bar{\Delta} (\tilde{H}^t)$) if there is no sequence $t_n \uparrow \infty$, $n \rightarrow \infty$, and tests $\delta_{t_n} \in \Sigma^{t_n}$ such that relations (2.2.13) hold.

DEFINITION 2.2.7. We say that *families of measures* (P^t) and (\tilde{P}^t) are *not completely asymptotically separable* (denoted by $(P^t) \bar{\Delta} (\tilde{P}^t)$) if there is no sequence $t_n \uparrow \infty$, $n \rightarrow \infty$, and sets $A_n \in \mathcal{B}^{t_n}$ such that relations (2.2.14) hold.

REMARK 2.2.3. If $(H^t) \bar{\Delta} (\tilde{H}^t)$, then according to Definition 2.2.6

$$\liminf_{n \rightarrow \infty} \beta(\delta_{t_n}) > 0$$

for all sequences $t_n \uparrow \infty$, $n \rightarrow \infty$, and all tests $\delta_{t_n} \in \Sigma^{t_n}$ such that $\alpha(\delta_{t_n}) \rightarrow 0$ as $n \rightarrow \infty$. A similar remark regarding the notion $(P^t) \bar{\Delta} (\tilde{P}^t)$ also holds. Therefore we have the following dichotomy: either $(H^t) \Delta (\tilde{H}^t)$ or $(H^t) \bar{\Delta} (\tilde{H}^t)$ (respectively, either $(P^t) \Delta (\tilde{P}^t)$ or $(P^t) \bar{\Delta} (\tilde{P}^t)$). Since Theorem 2.2.1 provides the necessary and sufficient conditions for $(H^t) \Delta (\tilde{H}^t)$, it can be used to characterize the notion $(H^t) \bar{\Delta} (\tilde{H}^t)$. For example, $(H^t) \bar{\Delta} (\tilde{H}^t) \iff \liminf_{t \rightarrow \infty} \|P^t \wedge \tilde{P}^t\| > 0$.

The complete asymptotic indistinguishability. Now we consider an asymptotic analog of the indistinguishability of hypotheses (cf. (1.1.8)).

DEFINITION 2.2.8. We say that families of hypotheses (H^t) and (\tilde{H}^t) are *completely asymptotically indistinguishable* (denoted by $(H^t) \cong (\tilde{H}^t)$) if

$$(2.2.21) \quad \lim_{t \rightarrow \infty} \alpha(\delta_t) = \alpha \Rightarrow \lim_{t \rightarrow \infty} \beta(\delta_t) = 1 - \alpha$$

for all $\alpha \in [0, 1]$ and all families (δ_t) of tests $\delta_t \in \Sigma^t$ such that the limit $\lim_{t \rightarrow \infty} \alpha(\delta_t)$ exists.

DEFINITION 2.2.9. We say that families of measures (P^t) and (\tilde{P}^t) are *completely asymptotically inseparable* (denoted by $(P^t) \cong (\tilde{P}^t)$) if

$$(2.2.22) \quad \lim_{t \rightarrow \infty} P^t(A_t) = \alpha \Rightarrow \lim_{t \rightarrow \infty} \tilde{P}^t(A_t) = \alpha$$

for all $\alpha \in [0, 1]$ and families (A_t) of sets $A_t \in \mathcal{B}^t$ such that the limit $\lim_{t \rightarrow \infty} P^t(A_t)$ exists.

The following result contains a characterization of the complete asymptotic indistinguishability of families of hypotheses.

THEOREM 2.2.2. *The following statements are equivalent:*

- a) $(H^t) \cong (\tilde{H}^t)$;
- b) $(P^t) \cong (\tilde{P}^t)$;
- c) $\lim_{t \rightarrow \infty} P^t(|\Lambda_t| > \gamma) = 0$ for all $\gamma > 0$;
- d) $\lim_{t \rightarrow \infty} \tilde{P}^t(|\Lambda_t| > \gamma) = 0$ for all $\gamma > 0$;
- e) $\lim_{t \rightarrow \infty} H(\varepsilon; \tilde{P}^t, P^t) = 1$ for all $\varepsilon \in (0, 1)$;
- f) $\lim_{t \rightarrow \infty} \rho(P^t, \tilde{P}^t) = 0$;
- g) $\lim_{t \rightarrow \infty} \|\tilde{P}^t - P^t\| = 0$;
- h) $\lim_{t \rightarrow \infty} \|P^t \wedge \tilde{P}^t\| = 1$.

PROOF. a) \Rightarrow b). Let α be an arbitrary number of the interval $[0, 1]$ and let (A_t) be an arbitrary family of sets $A_t \in \mathcal{B}^t$ such that $P^t(A_t) \rightarrow \alpha$ as $t \rightarrow \infty$. The test $\delta_t = I(A_t)$ is such that $\alpha(\delta_t) = P^t(A_t) \rightarrow \alpha$. Then condition a) implies that $\tilde{P}^t(A_t) = 1 - \beta(\delta_t) \rightarrow \alpha$ as $t \rightarrow \infty$, that is, the implication (2.2.22) is proved.

b) \Rightarrow c). We prove this implication by contradiction. Assume that condition b) holds and condition c) does not hold. Then there exist a number $\gamma_0 > 0$ and a sequence (t_n) such that $t_n \rightarrow \infty$ and $P^{t_n}(|\Lambda_{t_n}| > \gamma_0) \rightarrow a > 0$ as $n \rightarrow \infty$. Thus there are sequences $(m) \subset (n)$ and $(k) \subset (n)$ such that either $P^{t_m}(\Lambda_{t_m} > \gamma_0) \rightarrow b$ as $m \rightarrow \infty$ for some constant $b > 0$, or $P^{t_k}(\Lambda_{t_k} < -\gamma_0) \rightarrow c$ as $k \rightarrow \infty$ for some constant $c > 0$.

Let $P^{t_m}(\Lambda_{t_m} > \gamma_0) \rightarrow b > 0$ as $m \rightarrow \infty$. Then condition b) implies that $\tilde{P}^{t_m}(\Lambda_{t_m} > \gamma_0) \rightarrow b$ as $m \rightarrow \infty$. From the Lebesgue decomposition (2.1.6) we obtain

$$\tilde{P}^t(\Lambda_t > \gamma_0) = \int_{(\Lambda_t > \gamma_0)} \exp(\Lambda_t) dP^t + \tilde{P}^t(\Lambda_t = \infty) \geq e^{\gamma_0} P^t(\Lambda_t > \gamma_0).$$

Passing to the limit in this inequality along the sequence (t_m) we get $b \geq e^{\gamma_0} b$. Since $b > 0$ and $\gamma_0 > 0$, we obtain the contradiction $b > b$. Now we let

$$P^{t_k}(\Lambda_{t_k} < -\gamma_0) \rightarrow c > 0$$

as $k \rightarrow \infty$ and arrive at a similar contradiction. These contradictions prove the implication b) \Rightarrow c).

c) \Rightarrow d). Using the Lebesgue decomposition (2.1.6) we obtain

$$\tilde{P}^t(|\Lambda_t| \leq \gamma) \geq \tilde{P}^t(|\Lambda_t| \leq \gamma') = \int_{(|\Lambda_t| \leq \gamma')} \exp(\Lambda_t) dP^t \geq e^{-\gamma'} P^t(|\Lambda_t| \leq \gamma')$$

for all $\gamma > 0$ and $\gamma' \in (0, \gamma)$, whence $\liminf_{t \rightarrow \infty} \tilde{P}^t(|\Lambda_t| \leq \gamma) \geq e^{-\gamma'}$ by condition c). Passing to the limit as $\gamma' \rightarrow 0$ we prove that $\tilde{P}^t(|\Lambda_t| \leq \gamma) \rightarrow 1$ as $t \rightarrow \infty$ for all $\gamma > 0$, that is, condition d) holds.

d) \Rightarrow e). Let $\varepsilon \in (0, 1)$ and $\gamma > 0$. Since $(|\Lambda_t| \leq \gamma) \subset (\tilde{\mathfrak{z}}_t > 0, \mathfrak{z}_t > 0)$, it holds that

$$\begin{aligned} H(\varepsilon; \tilde{P}^t, P^t) &\geq E_Q^t \tilde{\mathfrak{z}}_t^{\varepsilon} \mathfrak{z}_t^{1-\varepsilon} I(|\Lambda_t| \leq \gamma) \\ &= E_Q^t \tilde{\mathfrak{z}}_t^{\varepsilon} \mathfrak{z}_t^{\varepsilon-1} I(|\Lambda_t| \leq \gamma) \geq e^{(\varepsilon-1)\gamma} \tilde{P}^t(|\Lambda_t| \leq \gamma). \end{aligned}$$

In view of d) this implies that

$$(2.2.23) \quad \liminf_{t \rightarrow \infty} H(\varepsilon; \tilde{P}^t, P^t) \geq e^{(\varepsilon-1)\gamma}$$

for all $\varepsilon \in (0, 1)$ and $\gamma > 0$. Passing to the limit in (2.2.23) as $\gamma \rightarrow 0$ and taking into account the inequality $H(\varepsilon; \tilde{P}^t, P^t) \leq 1$ we obtain condition e).

e) \Rightarrow f). Follows from equality (2.2.7).

f) \Rightarrow g). Follows from inequalities (2.2.6).

g) \Rightarrow h). Follows from equality (2.2.10).

h) \Rightarrow a). Let α be an arbitrary number of the interval $[0, 1]$ and let $(\bar{\delta}_t)$ be an arbitrary family of tests $\bar{\delta}_t \in \Sigma^t$ such that $\alpha(\bar{\delta}_t) \rightarrow \alpha$ as $t \rightarrow \infty$. By Lemma 1.1.3 the set \mathfrak{N}^t is symmetric about the point $(1/2, 1/2)$, so that it contains both points (α, β) and $(1 - \alpha, 1 - \beta)$. Thus

$$(2.2.24) \quad \inf \{ \alpha(\delta_t) + \beta(\delta_t) : \delta_t \in \Sigma^t \} = 2 - \sup \{ \alpha(\delta_t) + \beta(\delta_t) : \delta_t \in \Sigma^t \}.$$

Now we obtain from h) that

$$(2.2.25) \quad \lim_{t \rightarrow \infty} \sup \{ \alpha(\delta_t) + \beta(\delta_t) : \delta_t \in \Sigma^t \} = 1.$$

Applying the inequality

$$(2.2.26) \quad \begin{aligned} \inf \{ \alpha(\delta_t) + \beta(\delta_t) : \delta_t \in \Sigma^t \} &\leq \alpha(\bar{\delta}_t) + \beta(\bar{\delta}_t) \\ &\leq \sup \{ \alpha(\delta_t) + \beta(\delta_t) : \delta_t \in \Sigma^t \}, \end{aligned}$$

condition h), and relation (2.2.25) we prove that $\beta(\bar{\delta}_t) \rightarrow 1 - \alpha$ as $t \rightarrow \infty$, that is, the implication (2.2.21) is proved. \square

Definitions 2.2.8 and 2.2.9 are nonsymmetric with respect to $\alpha(\delta_t)$ and $\beta(\delta_t)$. Nevertheless using (2.2.24)–(2.2.26) one can prove the following result showing that (2.2.21) and (2.2.22) are, in fact, equivalent.

LEMMA 2.2.3. *If $(H^t) \cong (\tilde{H}^t)$, then*

$$\lim_{t \rightarrow \infty} \beta(\delta_t) = \beta \Rightarrow \lim_{t \rightarrow \infty} \alpha(\delta_t) = 1 - \beta$$

for all numbers $\beta \in [0, 1]$ and all families (δ_t) of tests $\delta_t \in \Sigma^t$ such that the limit $\lim_{t \rightarrow \infty} \beta(\delta_t)$ exists.

If $(P^t) \cong (\tilde{P}^t)$, then

$$\lim_{t \rightarrow \infty} \tilde{P}^t(A_t) = \alpha \Rightarrow \lim_{t \rightarrow \infty} P^t(A_t) = \alpha$$

for all numbers $\alpha \in [0, 1]$ and for all families (A_t) of sets $A_t \in \mathcal{B}^t$ such that the limit $\lim_{t \rightarrow \infty} \tilde{P}^t(A_t)$ exists.

REMARK 2.2.4. Definition 2.2.8 in the case of $(H^t) \cong (\tilde{H}^t)$ implies that

$$\mathfrak{N}^t \rightarrow \mathfrak{N} \quad \text{as } t \rightarrow \infty,$$

that is, the set \mathfrak{N}^t “approaches”, as $t \rightarrow \infty$, the diagonal of the square

$$[0, 1] \times [0, 1]$$

joining its corners $(1, 0)$ and $(0, 1)$. The measures P^t and \tilde{P}^t corresponding to the hypotheses H^t and \tilde{H}^t , respectively, asymptotically coincide in this case (cf. (1.1.8)).

Contiguous families of hypotheses. Now we consider families of hypotheses (H^t) and (\tilde{H}^t) whose asymptotic behavior differs from the complete asymptotic indistinguishability $(H^t) \cong (\tilde{H}^t)$ and complete asymptotic distinguishability $(H^t) \Delta (\tilde{H}^t)$.

DEFINITION 2.2.10. We say that a family of hypotheses (\tilde{H}^t) is contiguous to a family of hypotheses (H^t) (denoted by $(\tilde{H}^t) \triangleleft (H^t)$) if $\beta(\delta_t) \rightarrow 1$ as $t \rightarrow \infty$ for all tests $\delta_t \in \Sigma^t$ such that $\alpha(\delta_t) \rightarrow 0$ as $t \rightarrow \infty$. Otherwise, that is, if there exists a family (δ_t) of tests $\delta_t \in \Sigma^t$ such that

$$\lim_{t \rightarrow \infty} \alpha(\delta_t) = 0, \quad \liminf_{t \rightarrow \infty} \beta(\delta_t) < 1,$$

we say that a family of hypotheses (\tilde{H}^t) is noncontiguous to a family of hypotheses (H_t) (denoted by $(\tilde{H}^t) \ntriangleleft (H_t)$).

DEFINITION 2.2.11. We say that a family of measures (\tilde{P}^t) is contiguous to a family (P^t) (denoted by $(\tilde{P}^t) \triangleleft (P^t)$) if $\tilde{P}^t(A_t) \rightarrow 0$ as $t \rightarrow \infty$ for all sets $A_t \in \mathcal{B}^t$ such that $P^t(A_t) \rightarrow 0$ as $t \rightarrow \infty$. Otherwise, that is, if there exists a family (A_t) of sets $A_t \in \mathcal{B}^t$ such that

$$\lim_{t \rightarrow \infty} P^t(A_t) = 0, \quad \limsup_{t \rightarrow \infty} \tilde{P}^t(A_t) > 0,$$

we say that a family of measures (\tilde{P}^t) is noncontiguous to a family of measures (P_t) (denoted by $(\tilde{P}^t) \ntriangleleft (P_t)$).

Let (X^t, \mathcal{B}^t) , $t \in \mathbf{R}_+$, be a family of measurable spaces, let \mathbf{S}^t be a probability measure on (X^t, \mathcal{B}^t) , and let ζ_t , $t \in \mathbf{R}_+$, be a measurable function defined on (X^t, \mathcal{B}^t) and assuming values in $\overline{\mathbf{R}} = [-\infty, \infty]$.

DEFINITION 2.2.12. We say that a family (ζ_t) is dense with respect to a family of measures (\mathbf{S}^t) (denoted by $(\zeta_t | \mathbf{S}^t)$) if

$$\lim_{N \rightarrow \infty} \limsup_{t \rightarrow \infty} \mathbf{S}^t(|\zeta_t| > N) = 0.$$

DEFINITION 2.2.13. We say that a family (ζ_t) is uniformly integrable with respect to a family of measures (\mathbf{S}^t) if

$$\lim_{N \rightarrow \infty} \sup_{t \in \mathbf{R}_+} \int I(|\zeta_t| > N) |\zeta_t| d\mathbf{S}^t = 0.$$

The characterization of $(\tilde{H}^t) \triangleleft (H^t)$ is given in the following result.

THEOREM 2.2.3. The following statements are equivalent:

- a) $(\tilde{H}^t) \triangleleft (H^t)$;
- b) $(\tilde{P}^t) \triangleleft (P^t)$;
- c) $\lim_{t \rightarrow \infty} \tilde{P}^t(z_t = \infty) = 0$ and the family (z_t) is uniformly integrable with respect to the family of measures (P^t) ;
- d) $(z_t | \tilde{P}^t)$;
- e) $(1/\delta_t | \tilde{P}^t)$;
- f) $\lim_{\varepsilon \uparrow 1} \liminf_{t \rightarrow \infty} H(\varepsilon; \tilde{P}^t, P^t) = 1$.

PROOF. a) \Rightarrow b). Let $A_t \in \mathcal{B}^t$ and $P^t(A_t) \rightarrow 0$ as $t \rightarrow \infty$. Then $\alpha(\bar{\delta}_t) \rightarrow 0$ as $t \rightarrow \infty$ for the test $\bar{\delta}_t = I(A_t)$. Condition a) implies that $\tilde{P}^t(A_t) = 1 - \beta(\bar{\delta}_t) \rightarrow 0$ as $t \rightarrow \infty$.

b) \Rightarrow a). Let (δ_t) be an arbitrary family of tests such that $\alpha(\delta_t) \rightarrow 0$ as $t \rightarrow \infty$, and let ε be an arbitrary positive number. Put $A_t^\varepsilon = I(\delta_t \geq \varepsilon)$. Then

$$P^t(A_t^\varepsilon) \leq \varepsilon^{-1} \int_{A_t^\varepsilon} \delta_t dP^t \leq \varepsilon^{-1} \alpha(\delta_t),$$

whence it follows that $P^t(A_t^\varepsilon) \rightarrow 0$ as $t \rightarrow \infty$ for all $\varepsilon > 0$. This together with condition b) implies that $\tilde{P}^t(A_t^\varepsilon) \rightarrow 0$ as $t \rightarrow \infty$ for all $\varepsilon > 0$. Since

$$1 - \beta(\delta_t) = \int_{A_t^\varepsilon} \delta_t d\tilde{P}^t + \int_{(A_t^\varepsilon)^c} \delta_t d\tilde{P}^t \leq \tilde{P}^t(A_t^\varepsilon) + \varepsilon$$

and ε is arbitrary, we deduce that $\beta(\delta_t) \rightarrow 1$ as $t \rightarrow \infty$.

b) \Rightarrow c). Since $P^t(z_t = \infty) = 0$, condition b) implies that $\tilde{P}^t(z_t = \infty) \rightarrow 0$ as $t \rightarrow \infty$. A family (z_t) is uniformly integrable with respect to (P^t) if and only if

$$(2.2.27) \quad \sup_{t \in \mathbf{R}_+} \int z_t dP^t < \infty,$$

$$(2.2.28) \quad \text{if } P^t(A_t) \rightarrow 0 \text{ for } A_t \in \mathcal{B}^t, \text{ then } \int_{A_t} z_t dP^t \rightarrow 0$$

(see Lemma 2.6.2 in [47]). It follows from the Lebesgue decomposition that

$$\int_{A_t} z_t dP^t \leq \tilde{P}^t(A_t) \leq 1.$$

This implies (2.2.27), while condition b) implies (2.2.28).

c) \Rightarrow d). According to the Lebesgue decomposition

$$\tilde{P}^t(z_t > N) = \int_{(z_t > N)} z_t dP^t + \tilde{P}^t(z_t = \infty).$$

This together with condition c) implies d).

d) \Rightarrow b). Let $A_t \in \mathcal{B}^t$ be sets such that $P^t(A_t) \rightarrow 0$ as $t \rightarrow \infty$. According to the Lebesgue decomposition

$$\begin{aligned} \tilde{P}^t(A_t) &= \tilde{P}^t(A_t \cap (z_t \leq N)) + \tilde{P}^t(A_t \cap (z_t > N)) \\ &\leq \int_{A_t \cap (z_t \leq N)} z_t dP^t + \tilde{P}^t(z_t > N) \leq NP^t(A_t) + \tilde{P}^t(z_t > N). \end{aligned}$$

Using the latter result and condition d) we prove that $\tilde{P}^t(A_t) \rightarrow 0$ as $t \rightarrow \infty$.

d) \Leftrightarrow e). Follows from equality $z_t = 2/\beta_t - 1$.

e) \Rightarrow f). Let $\gamma > 0$. Then for $\varepsilon \in (0, 1)$

$$\begin{aligned} H(\varepsilon; \tilde{P}^t, P^t) &\geq E_Q^t \tilde{\beta}_t \left(\frac{\beta_t}{\tilde{\beta}_t} \right)^{1-\varepsilon} I(\beta_t \geq \gamma, \tilde{\beta}_t > 0) = \tilde{E}^t(\beta_t/\tilde{\beta}_t)^{1-\varepsilon} I(\beta_t \geq \gamma) \\ &\geq (\gamma/2)^{1-\varepsilon} \tilde{P}^t(\beta_t \geq \gamma), \end{aligned}$$

since $\beta_t + \tilde{\beta}_t = 2$. Thus for all $\gamma > 0$

$$\lim_{\varepsilon \uparrow 1} \liminf_{t \rightarrow \infty} H(\varepsilon; \tilde{P}^t, P^t) \geq \liminf_{\varepsilon \uparrow 1} \left(\frac{\gamma}{2} \right)^{1-\varepsilon} \liminf_{t \rightarrow \infty} \tilde{P}^t(\beta_t \geq \gamma) \geq \liminf_{t \rightarrow \infty} \tilde{P}^t(\beta_t \geq \gamma).$$

Condition e) implies that $\lim_{\gamma \downarrow 0} \lim_{t \rightarrow \infty} \tilde{P}^t(\beta_t \geq \gamma) = 1$, therefore the latter result and inequality $H(\varepsilon; \tilde{P}^t, P^t) \leq 1$ prove condition f).

f) \Rightarrow e). For all $\varepsilon, \lambda \in (0, 1)$ and $\gamma > 0$ we obtain similarly to (2.2.18) that

$$\begin{aligned} H(\varepsilon; \tilde{P}^t, P^t) &\leq \gamma^{1-\varepsilon} + \lambda^\varepsilon + \tilde{E}^t(\beta_t/\tilde{\beta}_t)^{1-\varepsilon} I(\beta_t \geq \gamma, \tilde{\beta}_t \geq \lambda) \\ &\leq \gamma^{1-\varepsilon} + \lambda^\varepsilon + (2/\lambda)^{1-\varepsilon} \tilde{P}^t(\beta_t \geq \gamma). \end{aligned}$$

Then for all $\varepsilon, \lambda \in (0, 1)$ we get

$$\liminf_{\gamma \downarrow 0} \liminf_{t \rightarrow \infty} \tilde{P}^t(\beta_t \geq \gamma) \geq \left(\frac{\lambda}{2} \right)^{1-\varepsilon} \liminf_{t \rightarrow \infty} H(\varepsilon; \tilde{P}^t, P^t) - \frac{\lambda}{2^{1-\varepsilon}}.$$

Passing to the limit in this inequality as $\varepsilon \uparrow 1$ and using condition f), then passing to the limit as $\lambda \downarrow 0$ we prove that

$$\liminf_{\gamma \downarrow 0} \liminf_{t \rightarrow \infty} \tilde{P}^t(\beta_t \geq \gamma) \geq 1,$$

whence condition e) follows. \square

REMARK 2.2.5. If $(\tilde{H}^t) \triangleleft (H^t)$ and $\mathfrak{N}^t \rightarrow \mathfrak{N}^\infty$ as $t \rightarrow \infty$, then according to Definition 2.2.10 the limit set \mathfrak{N}^∞ does not contain any point of the interval of the straight line joining the points $(0, 0)$ and $(0, 1)$, except for the point $(0, 1)$.

DEFINITION 2.2.14. If $(H^t) \triangleleft (\tilde{H}_t)$ and $(\tilde{H}^t) \triangleleft (H_t)$, then the families of hypotheses (H^t) and (\tilde{H}^t) are called *mutually contiguous* (denoted by $(H^t) \triangleleft \triangleright (\tilde{H}_t)$). If $(\tilde{H}^t) \triangleleft (H_t)$ and $(H^t) \triangleleft (\tilde{H}_t)$, then they are called *mutually noncontiguous* (denoted by $(H^t) \triangleleft \triangleright (\tilde{H}_t)$). If either $(\tilde{H}^t) \triangleleft (H_t)$ or $(H^t) \triangleleft (\tilde{H}_t)$, then we say that the families of hypotheses (H^t) and (\tilde{H}^t) are not mutually contiguous (denoted by $(H^t) \triangleleft \triangleright (\tilde{H}_t)$).

DEFINITION 2.2.15. If $(\tilde{P}^t) \triangleleft (P_t)$ and $(P^t) \triangleleft (\tilde{P}_t)$, then the families of measures (P^t) and (\tilde{P}^t) are called *mutually contiguous* (denoted by $(P^t) \triangleleft \triangleright (\tilde{P}^t)$). If $(\tilde{P}^t) \triangleleft (P_t)$ and $(P^t) \triangleleft (\tilde{P}_t)$, then they are called *mutually noncontiguous* (denoted by $(P^t) \triangleleft \triangleright (\tilde{P}^t)$). If either $(\tilde{P}^t) \triangleleft (P_t)$ or $(P^t) \triangleleft (\tilde{P}_t)$, then we say that the *families of measures* (P^t) and (\tilde{P}^t) are *not mutually contiguous* (denoted by $(P^t) \triangleleft \triangleright (\tilde{P}^t)$).

REMARK 2.2.6. Theorem 2.2.3 implies the characterization of all types

$$(\tilde{H}^t) \triangleleft (H^t), \quad (H^t) \triangleleft \triangleright (\tilde{H}^t), \quad (H^t) \triangleleft \triangleright (\tilde{H}^t), \quad \text{and} \quad (H^t) \triangleleft \triangleright (\tilde{H}^t).$$

For example,

$$(\tilde{H}^t) \triangleleft (H^t) \iff \lim_{\varepsilon \uparrow 1} \liminf_{t \rightarrow \infty} H(\varepsilon; \tilde{P}^t, P^t) < 1.$$

Further results on the contiguous families can be found in [21, 22, 37, 45].

The whole range of types of the asymptotic distinguishability. Using the notions of the asymptotic distinguishability of families of hypotheses we obtain the whole range of types of the asymptotic distinguishability of families of hypotheses (H^t) and (\tilde{H}^t) . We will use the following conditions:

- a₀) $(H^t) \cong (\tilde{H}^t)$;
- a) $(H^t) \triangleleft \triangleright (\tilde{H}^t)$;
- b) $(H^t) \triangleleft (\tilde{H}^t)$, $(\tilde{H}^t) \triangleleft (H^t)$;
- c) $(H^t) \triangleleft (\tilde{H}^t)$, $(\tilde{H}^t) \triangleleft (H^t)$;
- d) $(H^t) \triangleleft \triangleright (\tilde{H}^t)$, $(H^t) \triangleleft (\tilde{H}^t)$;
- e) $(H^t) \triangleleft (\tilde{H}^t)$.

DEFINITION 2.2.16. We say that the *asymptotic distinguishability of families of hypotheses* (H^t) and (\tilde{H}^t) is of type **a₀** (respectively, of type **a**, **b**, **c**, **d**, or **e**), if condition **a₀** (respectively **a**), **b**), **c**), **d**), or **e**) holds.

Note that the types **a**, **b**, **c**, **d**, and **e** are disjoint and form the *whole range of types of the asymptotic distinguishability of families of hypotheses*. Since **a₀** \Rightarrow **a**), type **a₀** is a subtype of type **a**. If the type **a₁** is defined as a subtype of the type **a** for which condition **a₀** does not hold, then the types **a₀**, **a₁**, **b**, **c**, **d**, and **e** still form the whole range of disjoint types of the asymptotic distinguishability of families of hypotheses (H^t) and (\tilde{H}^t) .

A characterization of types **e** and **a₀** is given in Theorems 2.2.1 and 2.2.2, respectively. A characterization of other types can easily be obtained by combining Theorems 2.2.1–2.2.3 and taking into account Remarks 2.2.3 and 2.2.6. We do not give this characterization and leave it to the reader.

EXAMPLE 2.2.1. Let an observation be the vector $\xi^{(n)} = (\xi_{n1}, \xi_{n2}, \dots, \xi_{nn})$ where $\xi_{n1}, \xi_{n2}, \dots, \xi_{nn}$ are independent random variables such that the distribution of ξ_{ni} is $\mathcal{N}(a_{ni}, 1)$ under the hypothesis H^n or $\mathcal{N}(\tilde{a}_{ni}, 1)$ under the hypothesis \tilde{H}^n . Then the likelihood ratio $z_n(x)$, $x \in \mathbf{R}^n$, is the density of the measure \tilde{P}^n corresponding to the hypothesis \tilde{H}^n with respect to the measure P^t corresponding to the hypothesis H^n . The likelihood ratio is given by

$$z_n(x) = \exp \left(\sum_{i=1}^n (\tilde{a}_{ni} - a_{ni}) x_i - \frac{1}{2} \sum_{i=1}^n (\tilde{a}_{ni}^2 - a_{ni}^2) \right)$$

where $x = (x_1, x_2, \dots, x_n)$. Put $\Lambda_n = \Lambda_n(\xi^{(n)})$ for $\Lambda_n(x) = \ln z_n(x)$. It is clear that

$$(2.2.29) \quad \mathcal{L}(\Lambda_n | H^n) = \mathcal{N}\left(-\frac{1}{2}v_n^2, v_n^2\right),$$

$$(2.2.30) \quad \mathcal{L}(\Lambda_n | \tilde{H}^n) = \mathcal{N}\left(\frac{1}{2}v_n^2, v_n^2\right)$$

where

$$(2.2.31) \quad v_n^2 = \sum_{i=1}^n (\tilde{a}_{ni} - a_{ni})^2$$

and $\mathcal{L}(\Lambda_n | H^n)$ is the distribution of Λ_n under the hypothesis H^n . From equality (2.2.29) we derive that

$$(2.2.32) \quad P^n(\Lambda_n < N) = \Phi\left(\frac{v_n}{2} + \frac{N}{v_n}\right),$$

whence it follows by Theorem 2.2.1 that

$$(2.2.33) \quad (H^n) \Delta (\tilde{H}^n) \iff \limsup_{n \rightarrow \infty} v_n = \infty.$$

Using equality (2.2.30) we obtain

$$(2.2.34) \quad \tilde{P}^n(\Lambda_n > N) = \Phi\left(\frac{v_n}{2} - \frac{N}{v_n}\right),$$

whence

$$(2.2.35) \quad \limsup_{n \rightarrow \infty} v_n < \infty \implies (\tilde{H}^n) \triangleleft (H^n)$$

in view of Theorem 2.2.3. By contradiction, we derive from (2.2.33) that

$$(2.2.36) \quad (\tilde{H}^n) \triangleleft (H^n) \implies \limsup_{n \rightarrow \infty} v_n < \infty.$$

Analogously, using (2.2.32) and Theorem 2.2.3 with the hypotheses H^n and \tilde{H}^n interchanged we obtain

$$(2.2.37) \quad (H^n) \triangleleft (\tilde{H}^n) \iff \limsup_{n \rightarrow \infty} v_n < \infty.$$

Combining (2.2.35)–(2.2.37) we prove that

$$(2.2.38) \quad (H^n) \triangleleft \triangleright (\tilde{H}^n) \iff \limsup_{n \rightarrow \infty} v_n < \infty.$$

It follows from (2.2.33) and (2.2.38) that either $(H^n) \Delta (\tilde{H}^n)$ or $(H^n) \triangleleft \triangleright (\tilde{H}^n)$. In other words, either the distinguishability is of type a or of type e. Moreover using Theorem 2.2.2 and relations (2.2.29) and (2.2.30) one can show that

$$(H^n) \cong (\tilde{H}^n) \iff \limsup_{n \rightarrow \infty} v_n = 0.$$

Therefore either the distinguishability is of type \mathbf{a}_0 or of type \mathbf{a}_1 or of type \mathbf{e} . Namely

$$\begin{aligned} \text{type } \mathbf{a}_0 &\iff \limsup_{n \rightarrow \infty} v_n = 0, \\ \text{type } \mathbf{a}_1 &\iff 0 < \limsup_{n \rightarrow \infty} v_n < \infty, \\ \text{type } \mathbf{e} &\iff \limsup_{n \rightarrow \infty} v_n = \infty. \end{aligned}$$

Further results on types of the asymptotic distinguishability of families of hypotheses and various examples can be found in [37].

2.3. Complete asymptotic distinguishability under the strong law of large numbers

Consider the Neyman–Pearson test δ_t^{+, α_t} of level $\alpha_t \in [0, 1]$ defined by (2.1.12) and assume that the complete asymptotic distinguishability $(H^t) \Delta (\tilde{H}^t)$ holds. By Theorem 2.2.1, we have

$$(2.3.1) \quad (H^t) \Delta (\tilde{H}^t) \iff \liminf_{t \rightarrow \infty} [\alpha(\delta_t^{c,1}) + \beta(\delta_t^{c,1})] = 0 \quad \text{for all } c \in (0, \infty)$$

where $\delta_t^{c,1} = I(z_t \geq c)$ is the likelihood ratio test defined by (2.1.11).

This implies that if $(H^t) \Delta (\tilde{H}^t)$, for any $c \in (0, \infty)$ there exists a sequence (t_n) for the test δ_t^{+, α_t} with $\alpha_t = \alpha(\delta_t^{c,1})$ such that $t_n \rightarrow \infty$, $\alpha_{t_n} \rightarrow 0$, and $\beta(\delta_{t_n}^{+, \alpha_{t_n}}) \rightarrow 0$ as $n \rightarrow \infty$. By equivalence (2.3.1), obtaining more refined properties of the function $\beta(\delta_{t_n}^{+, \alpha_{t_n}})$ requires that we consider a more stringent constraint than $(H^t) \Delta (\tilde{H}^t)$.

Relative entropy and the law of large numbers. Consider the following condition:

$\Lambda 1$. $\lim_{t \rightarrow \infty} \mathbf{P}^t(|\chi_t^{-1} \Lambda_t + 1| > \gamma) = 0$ for any $\gamma > 0$ where χ_t is a nonrandom positive function such that $\chi_t \rightarrow \infty$ as $t \rightarrow \infty$.

It is easy to see that, in view of Theorem 2.2.1,

$$\Lambda 1 \Rightarrow (H^t) \Delta (\tilde{H}^t),$$

that is, the complete asymptotic distinguishability holds under condition $\Lambda 1$.

DEFINITION 2.3.1. The number

$$(2.3.2) \quad I(\mathbf{P}^t | \tilde{\mathbf{P}}^t) = \mathbf{E}_{\tilde{Q}}^t \mathfrak{z}_t \ln(\mathfrak{z}_t / \tilde{\mathfrak{z}}_t)$$

is called *the relative entropy of a measure \mathbf{P}^t with respect to a measure $\tilde{\mathbf{P}}^t$* . We agree that $\mathfrak{z}_t \ln(\mathfrak{z}_t / \tilde{\mathfrak{z}}_t)$ equals 0 if $\mathfrak{z}_t = 0$ and equals ∞ if $\tilde{\mathfrak{z}}_t = 0$. Then the relative entropy is well defined for all t .

The relative entropy $I(\mathbf{P}^t | \tilde{\mathbf{P}}^t)$ is often called the *Kullback–Leibler divergence*, or *distance* or *deviation* [7, 9, 11, 33], or the *Kullback–Leibler information* defined for the measures \mathbf{P}^t and $\tilde{\mathbf{P}}^t$ [1, 33].

LEMMA 2.3.1. *We have*

$$(2.3.3) \quad I(\mathbf{P}^t | \tilde{\mathbf{P}}^t) = \mathbf{E}^t \ln \tilde{z}_t = -\mathbf{E}^t \ln z_t$$

where z_t and \tilde{z}_t are the likelihood ratios defined by (2.1.2). Moreover

$$(2.3.4) \quad I(\mathbf{P}^t | \tilde{\mathbf{P}}^t) = \tilde{\mathbf{E}}^t \tilde{z}_t \ln \tilde{z}_t$$

if $\mathbf{P}^t \ll \tilde{\mathbf{P}}^t$ and $I(\mathbf{P}^t | \tilde{\mathbf{P}}^t) = \infty$ if $\mathbf{P}^t \not\ll \tilde{\mathbf{P}}^t$.

PROOF. The first equality in (2.3.3) follows immediately from definitions (2.3.2) and (2.1.2). The second equality in (2.3.3) holds, since $\tilde{z}_t = z_t^{-1}$ (\mathbf{P}^t -a.s.).

If $\mathbf{P}^t \ll \tilde{\mathbf{P}}^t$, then $\tilde{z}_t = d\mathbf{P}^t/d\tilde{\mathbf{P}}^t$ ($\tilde{\mathbf{P}}^t$ -a.s.). Thus equality (2.3.4) follows from the first equality in (2.3.3). If $\mathbf{P}^t \not\ll \tilde{\mathbf{P}}^t$, then $\mathbf{P}^t(\tilde{\mathfrak{J}}_t = 0) > 0$. Therefore

$$\mathbf{P}^t(\tilde{\mathfrak{J}}_t > 0, \tilde{\mathfrak{J}}_t = 0) > 0$$

implying that $\mathbf{Q}^t(\tilde{\mathfrak{J}}_t > 0, \tilde{\mathfrak{J}}_t = 0) > 0$. Then definition (2.3.2) gives

$$I(\mathbf{P}^t | \tilde{\mathbf{P}}^t) = \infty. \quad \square$$

LEMMA 2.3.2. *The relative entropy $I(\mathbf{P}^t | \tilde{\mathbf{P}}^t)$ of a measure \mathbf{P}^t with respect to a measure $\tilde{\mathbf{P}}^t$ is nonnegative. Moreover $I(\mathbf{P}^t | \tilde{\mathbf{P}}^t) = 0$ if and only if $\mathbf{P}^t = \tilde{\mathbf{P}}^t$.*

PROOF. By Lemma 2.3.1, it is sufficient to consider the case $\mathbf{P}^t \ll \tilde{\mathbf{P}}^t$ in view of equality (2.3.4). Put $\varphi(t) = t \ln t$. Equality (2.3.4) and the Jensen inequality imply that

$$I(\mathbf{P}^t | \tilde{\mathbf{P}}^t) = \tilde{\mathbf{E}}^t \varphi(\tilde{z}_t) \geq \varphi(\tilde{\mathbf{E}}^t \tilde{z}_t) = \varphi(1) = 0$$

where the inequality becomes an equality if and only if $\tilde{z}_t = \text{const} = a$ ($\tilde{\mathbf{P}}^t$ -a.s.). Note that $a = \tilde{\mathbf{E}}^t \tilde{z}_t = 1$ in this case. This means that if the equality holds, we have for any $A \in \mathcal{B}^t$

$$\mathbf{P}^t(A) = \int_A \tilde{z}_t d\tilde{\mathbf{P}}^t = \int_A d\tilde{\mathbf{P}}^t = \tilde{\mathbf{P}}^t(A),$$

that is, $\mathbf{P}^t = \tilde{\mathbf{P}}^t$. □

REMARK 2.3.1. Condition $\Lambda 1$ is known as *relative stability* of Λ_t as $t \rightarrow \infty$ and is the most natural and general form of the law of large numbers [20]. If the relative entropy $I(\mathbf{P}^t | \tilde{\mathbf{P}}^t)$ is finite for any $t \in \mathbf{R}_+$ and $I(\mathbf{P}^t | \tilde{\mathbf{P}}^t) \rightarrow \infty$ as $t \rightarrow \infty$, then, by putting $\chi_t = I(\mathbf{P}^t | \tilde{\mathbf{P}}^t)$, condition $\Lambda 1$ can be transformed into the following form resembling the law of large numbers:

$$\lim_{t \rightarrow \infty} \mathbf{P}^t \left(\left| \frac{\Lambda_t - \mathbf{E}^t \Lambda_t}{\chi_t} \right| > \varepsilon \right) = 0 \quad \text{for all } \varepsilon > 0$$

where, by (2.3.3), we have $\mathbf{E}^t \Lambda_t = -I(\mathbf{P}^t | \tilde{\mathbf{P}}^t)$.

The following result describes the behavior of the points $(\bar{\alpha}_t, 0)$ and $(0, \bar{\beta}_t)$ of the set \mathfrak{N}^t under the law of large numbers.

LEMMA 2.3.3. *The following relations hold:*

$$(2.3.5) \quad \Lambda 1 \Rightarrow \lim_{t \rightarrow \infty} \bar{\alpha}_t = 1,$$

$$(2.3.6) \quad \Lambda 1 \Rightarrow \liminf_{t \rightarrow \infty} \chi_t^{-1} \ln \bar{\beta}_t \geq -1.$$

PROOF. Implication (2.3.5) follows from (2.1.4) and the estimate

$$P^t(\Lambda_t < -a\chi_t) \geq P^t(\Lambda_t = -\infty) = P^t(z_t = 0)$$

where $1 < a < \infty$. By the Lebesgue decomposition (2.1.7) and in view of equality (2.1.5), we obtain for any $a \in (1, \infty)$ that

$$P^t(\Lambda_t \geq -a\chi_t) \leq e^{a\chi_t} E^t e^{\Lambda_t} \chi(\Lambda_t \geq -a\chi_t) \leq e^{a\chi_t} E^t z_t = e^{a\chi_t} \tilde{P}^t(z_t < \infty) = e^{a\chi_t} \bar{\beta}_t$$

which together with condition $\Lambda 1$ yields

$$\liminf_{t \rightarrow \infty} \chi_t^{-1} \ln \bar{\beta}_t \geq -a.$$

Approaching the limit as $a \rightarrow 1$, we obtain implication (2.3.6). \square

REMARK 2.3.2. By equality (1.1.50), implication (2.3.6) is equivalent to

$$\Lambda 1 \Rightarrow \liminf_{t \rightarrow \infty} \chi_t^{-1} \ln \beta(\delta_t^{+,0}) \geq -1.$$

Behavior of the Neyman–Pearson test under the law of large numbers. Introduce the following conditions:

$$\begin{array}{ll} \alpha 1) \liminf_{t \rightarrow \infty} \alpha_t > 0; & \alpha 2) \limsup_{t \rightarrow \infty} \alpha_t < 1; \\ d 1) \limsup_{t \rightarrow \infty} \chi_t^{-1} d_t \leq -1; & d 2) \liminf_{t \rightarrow \infty} \chi_t^{-1} d_t \geq -1; \\ \beta 1) \limsup_{t \rightarrow \infty} \chi_t^{-1} \ln \beta(\delta_t^{+, \alpha_t}) \leq -1; & \beta 2) \liminf_{t \rightarrow \infty} \chi_t^{-1} \ln \beta(\delta_t^{+, \alpha_t}) \geq -1, \end{array}$$

where χ_t is the normalizing term appearing in the law of large numbers $\Lambda 1$.

The following result establishes relationships between the behavior of α_t , d_t , and $\beta(\delta_t^{+, \alpha_t})$ under condition $\Lambda 1$.

THEOREM 2.3.1. *The following implications hold:*

$$(2.3.7) \quad \Lambda 1, \alpha 1 \Rightarrow d 1 \Rightarrow \beta 1,$$

$$(2.3.8) \quad \Lambda 1, \alpha 2 \Rightarrow \beta 2 \Rightarrow d 2.$$

PROOF. The implication $\Lambda 1, \alpha 1 \Rightarrow d 1$ is proved by contradiction. Representation (2.1.12) implies

$$(2.3.9) \quad \alpha_t = P^t(Y_t > y_t) + \varepsilon_t P^t(Y_t = y_t)$$

where $Y_t = \chi_t^{-1} \Lambda_t$ and $y_t = \chi_t^{-1} d_t$.

Let (t_n) be a sequence such that $t_n \rightarrow \infty$ and $y_{t_n} \rightarrow \bar{y} = \limsup_{t \rightarrow \infty} y_t$ as $n \rightarrow \infty$. Assume that $\underline{\alpha} = \liminf_{t \rightarrow \infty} \alpha_t > 0$ and condition $\Lambda 1$ holds, but $\bar{y} > -1$. Then, by condition $\Lambda 1$, we obtain that $P^t(Y_{t_n} \geq y_{t_n}) \rightarrow 0$ as $n \rightarrow \infty$. On the other hand, the inequality $\underline{\alpha} > 0$ and equality (2.3.9) imply that

$$\lim_{n \rightarrow \infty} P^{t_n}(Y_{t_n} \geq y_{t_n}) \geq \liminf_{n \rightarrow \infty} \alpha_{t_n} \geq \liminf_{t \rightarrow \infty} \alpha_t = \underline{\alpha} > 0.$$

The contradiction we have obtained proves the implication $\Lambda 1, \alpha 1 \Rightarrow d 1$.

Using equality (2.1.8) we obtain that

$$(2.3.10) \quad \begin{aligned} \beta(\delta_t^{+, \alpha_t}) &= \tilde{E}^t(1 - \delta_t^{+, \alpha_t}) = E^t(1 - \delta_t^{+, \alpha_t})z_t + \tilde{E}^t(1 - \delta_t^{+, \alpha_t})I(z_t = \infty) \\ &= E^t(1 - \delta_t^{+, \alpha_t})z_t \leq e^{d_t}, \end{aligned}$$

proving the implications $d1 \Rightarrow \beta1$ and $\beta2 \Rightarrow d2$.

Now we prove the implication $\Lambda1, \alpha2 \Rightarrow \beta2$. Let $a > 1$ and $\gamma > 0$ be arbitrary numbers and let $S_t = \{\chi_t^{-1}\Lambda_t \geq -a\}$. In view of condition $\Lambda1$, there exists a number $t_0 = t_0(a, \gamma)$ depending on a and γ such that $P^t(S_t^c) < \gamma$ for all $t > t_0$. Therefore we obtain from (2.1.8) for $t > t_0$ that

$$(2.3.11) \quad \begin{aligned} \beta(\delta_t^{+, \alpha_t}) &= E^t z_t(1 - \delta_t^{+, \alpha_t}) \geq E^t I(S_t) z_t(1 - \delta_t^{+, \alpha_t}) \\ &\geq e^{-a\chi_t} E^t I(S_t)(1 - \delta_t^{+, \alpha_t}) \geq e^{-a\chi_t}(1 - \alpha_t - \gamma). \end{aligned}$$

In view of condition $\alpha2$ and since $a > 1$ and $\gamma > 0$ are arbitrary, relation (2.3.11) proves the implication $\Lambda1, \alpha2 \Rightarrow \beta2$. \square

Implications (2.3.7) and (2.3.8) give the following result.

COROLLARY 2.3.1. *If the law of large numbers $\Lambda1$ holds and the level α_t satisfies conditions $\alpha1$ and $\alpha2$, then*

$$(2.3.12) \quad \lim_{t \rightarrow \infty} \chi_t^{-1} \ln \beta(\delta_t^{+, \alpha_t}) = -1.$$

Corollary 2.3.1 implies that if the law of large numbers $\Lambda1$ holds for the Neyman–Pearson test δ_t^{+, α_t} whose level α_t tends to a limit $\alpha \in (0, 1)$ as $t \rightarrow \infty$, then the rate of decay of the type II error probability $\beta(\delta_t^{+, \alpha_t})$ does not depend on α . More specifically, (2.3.12) means that

$$(2.3.13) \quad \beta(\delta_t^{+, \alpha_t}) = \exp(-b\chi_t(1 + o(1))), \quad t \rightarrow \infty,$$

where $b = 1$ for all $\alpha \in (0, 1)$.

Independent observations. Stein's lemma. Assume that an observation is the vector $\xi^{(n)} = (\xi_{n1}, \xi_{n2}, \dots, \xi_{nn})$, $n = 1, 2, \dots$, where $\xi_{n1}, \xi_{n2}, \dots, \xi_{nn}$ are independent random variables. The hypothesis H^n is that ξ_{ni} has a distribution P_{ni} with the density $p_{ni}(x)$ with respect to the Lebesgue measure, while the hypothesis \tilde{H}^n is that the distribution of ξ_{ni} is \tilde{P}_{ni} with the density $\tilde{p}_{ni}(x)$. Then the distribution of the vector $\xi^{(n)}$ under the hypothesis H^n is $P^n = P_{n1} \times P_{n2} \times \dots \times P_{nn}$, and the density of this distribution with respect to the Lebesgue measure is

$$p_n(x) = \prod_{i=1}^n p_{ni}(x_i), \quad x = (x_1, \dots, x_n).$$

Similarly, the distribution of the vector $\xi^{(n)}$ under the hypothesis \tilde{H}^n is

$$\tilde{P}^n = \tilde{P}_{n1} \times \tilde{P}_{n2} \times \dots \times \tilde{P}_{nn}$$

and the density with respect to the Lebesgue measure equals

$$\tilde{p}_n(x) = \prod_{i=1}^n \tilde{p}_{ni}(x_i), \quad x = (x_1, \dots, x_n).$$

The likelihood ratio $z_n(x)$ for the measures \tilde{P}^n and P^n is

$$(2.3.14) \quad z_n(x) = \prod_{i=1}^n z_{ni}(x_i), \quad z_{ni}(x_i) = \frac{\tilde{p}_{ni}(x_i)}{p_{ni}(x_i)}, \quad x = (x_1, \dots, x_n),$$

where we agree that $0/0 = 0$. It is clear that the relative entropies $I(P^n|\tilde{P}^n)$ and $I(P_{ni}|\tilde{P}_{ni})$, $i = 1, 2, \dots, n$, are related as follows:

$$(2.3.15) \quad I(P^n|\tilde{P}^n) = \sum_{i=1}^n I(P_{ni}|\tilde{P}_{ni}).$$

Let

$$(2.3.16) \quad \Lambda_n(x) = \ln z_n(x), \quad \lambda_{ni}(x_i) = \ln z_{ni}(x_i), \quad i = 1, 2, \dots, n,$$

where we agree that $\ln 0 = -\infty$. If

$$\Lambda_n = \Lambda_n(\xi^{(n)}) \quad \text{and} \quad \lambda_{ni} = \lambda_{ni}(\xi_{ni}), \quad i = 1, 2, \dots, n,$$

then $\Lambda_n = \sum_{i=1}^n \lambda_{ni}$. It follows from (2.3.3) that

$$(2.3.17) \quad I(P^n|\tilde{P}^n) = -E^n \Lambda_n, \quad I(P_{ni}|\tilde{P}_{ni}) = -E_{ni} \lambda_{ni}$$

where E^n and E_{ni} are mathematical expectations with respect to the distributions P^n and P_{ni} , respectively.

Corollary 2.3.1 can be stated in the following form.

COROLLARY 2.3.2. *Let $I(P^n|\tilde{P}^n) < \infty$ for all $n = 1, 2, \dots$, and let*

$$I(P^n|\tilde{P}^n) \rightarrow \infty \quad \text{as } n \rightarrow \infty.$$

If the law of large numbers $\Lambda 1$ holds for Λ_n , $n = 1, 2, \dots$, with $\chi_n = I(P^n|\tilde{P}^n)$ and the level α_n satisfies conditions $\alpha 1$ and $\alpha 2$, then

$$(2.3.18) \quad \lim_{n \rightarrow \infty} \chi_n^{-1} \ln \beta(\delta_n^{+, \alpha_n}) = -1.$$

Now let $\xi^{(n)} = (\xi_1, \xi_2, \dots, \xi_n)$ where $\xi_1, \xi_2, \dots, \xi_n$ are independent identically distributed random variables such that the distribution of ξ_i under the hypothesis H^n is P with the density $p(x)$ with respect to the Lebesgue measure, while the distribution of ξ_i under the hypothesis \tilde{H}^n is \tilde{P} with the density $\tilde{p}(x)$. We assume that the distribution of ξ_i is independent of n both under H^n and under \tilde{H}^n . Then (2.3.14) and (2.3.16) can be rewritten as

$$(2.3.19) \quad z_n(x) = \prod_{i=1}^n z(x_i), \quad z(x_i) = \frac{\tilde{p}(x_i)}{p(x_i)},$$

$$(2.3.20) \quad \Lambda_n(x) = \ln z_n(x), \quad \lambda(x_i) = \ln z(x_i)$$

where $x = (x_1, \dots, x_n)$. Put $\Lambda_n = \Lambda_n(\xi^{(n)})$ and $\lambda_i = \lambda(\xi_i)$, $i = 1, 2, \dots, n$. Then $\Lambda_n = \sum_{i=1}^n \lambda_i$ and (2.3.15) and (2.3.17) imply

$$(2.3.21) \quad I(P^n|\tilde{P}^n) = nI(P|\tilde{P}), \quad I(P|\tilde{P}) = -E\lambda_1.$$

Therefore the following classical result follows from Corollary 2.3.2.

COROLLARY 2.3.3. Let $\xi^{(n)} = (\xi_1, \dots, \xi_n)$ where ξ_1, \dots, ξ_n are independent identically distributed random variables both under the hypothesis H^n and under the hypothesis \tilde{H}^n . Assume that the distributions of ξ_n are independent of n and the relative entropy $I(P|\tilde{P})$ is positive and finite. If $\alpha_n \rightarrow \alpha \in (0, 1)$ as $n \rightarrow \infty$, then

$$(2.3.22) \quad \lim_{n \rightarrow \infty} \frac{1}{n} \ln \beta(\delta_n^{+, \alpha_n}) = -I(P|\tilde{P}).$$

PROOF. It follows from (2.3.19) and (2.3.20) that $\Lambda_n = \sum_{i=1}^n \lambda_i$ where the random variables $\lambda_1, \lambda_2, \dots, \lambda_n$ are independent and identically distributed under H^n with mean $E\lambda_i = -I(P|\tilde{P})$. By the Khintchine law of large numbers for sums of independent identically distributed random variables [20], we obtain that condition $\Lambda 1$ holds with $\chi_n = nI(P|\tilde{P})$. Since $\alpha_n \rightarrow \alpha \in (0, 1)$ as $n \rightarrow \infty$, the assumptions of Corollary 2.3.2 hold. Therefore (2.3.22) follows from (2.3.18). \square

Corollary 2.3.3 is proved by Rao [43]. Corollary 2.3.3 with $\alpha_n = \alpha \in (0, 1)$ for all n is called the *Stein lemma* (see [1, 33]).

EXAMPLE 2.3.1. Let an observation be $\xi^{(n)} = (\xi_{n1}, \xi_{n2}, \dots, \xi_{nn})$ where $\xi_{n1}, \xi_{n2}, \dots, \xi_{nn}$ are independent random variables such that the distribution of ξ_{ni} is $\mathcal{N}(a_{ni}, 1)$ under the hypothesis H^n and $\mathcal{N}(\tilde{a}_{ni}, 1)$ under the hypothesis \tilde{H}^n (see Example 2.2.1). Using the notation of Example 2.2.1 we obtain from (2.2.29)–(2.2.31) that

$$(2.3.23) \quad \Lambda_n = v_n \eta - \frac{1}{2} v_n^2,$$

$$(2.3.24) \quad \Lambda_n = v_n \tilde{\eta} + \frac{1}{2} v_n^2$$

where

$$(2.3.25) \quad \mathcal{L}(\eta|H^n) = \mathcal{N}(0, 1), \quad \mathcal{L}(\tilde{\eta}|\tilde{H}^n) = \mathcal{N}(0, 1).$$

Assume that the complete asymptotic distinguishability $(H^n) \Delta (\tilde{H}^n)$ holds. Considering (2.2.33) we assume without loss of generality that $v_n \rightarrow \infty$ as $n \rightarrow \infty$. It follows from (2.3.23)–(2.3.25) that

$$(2.3.26) \quad I(P^n|\tilde{P}^n) = I(\tilde{P}^n|P^n) = \frac{1}{2} v_n^2.$$

Since $v_n \rightarrow \infty$ as $n \rightarrow \infty$, relation (2.3.23) implies that the law of large numbers holds with $\chi_n = v_n^2/2$. Therefore Corollary 2.3.2 holds if conditions $\alpha 1$ and $\alpha 2$ are satisfied. Note that (2.2.32) yields

$$(2.3.27) \quad \alpha(\delta_n^{+, \alpha_n}) = P^n(\Lambda_n > d_n) = 1 - \Phi\left(\frac{v_n}{2} + \frac{d_n}{v_n}\right),$$

whence

$$(2.3.28) \quad d_n = -\frac{1}{2} v_n^2 + v_n t_{1-\alpha_n}$$

where t_p is the p -quantile of the distribution $\mathcal{N}(0, 1)$. This implies that conditions $\alpha 1$ and $\alpha 2$ hold if the parameter d_n of the test δ_n^{+, α_n} is such that

$$(2.3.29) \quad \liminf_{n \rightarrow \infty} \left(\frac{d_n}{v_n} + \frac{v_n}{2} \right) > -\infty, \quad \limsup_{n \rightarrow \infty} \left(\frac{d_n}{v_n} + \frac{v_n}{2} \right) < \infty.$$

Therefore Corollary 2.3.2 implies that

$$(2.3.30) \quad \beta(\delta_n^{+, \alpha_n}) = \exp \left\{ -\frac{1}{2} v_n^2 (1 + o(1)) \right\}$$

as $n \rightarrow \infty$ if conditions (2.3.29) hold.

Moreover $\alpha_n \rightarrow \alpha \in (0, 1)$ as $n \rightarrow \infty$ if conditions (2.3.27) and (2.3.28) hold and if

$$(2.3.31) \quad d_n = -\frac{1}{2} v_n^2 + v_n t_{1-\alpha} + o(v_n).$$

It follows from (2.2.34) and (2.3.28) that

$$(2.3.32) \quad \beta(\delta_n^{+, \alpha_n}) = \tilde{\mathbb{P}}^n(\Lambda_n < d_n) = 1 - \Phi(v_n - t_{1-\alpha_n}).$$

By the well-known asymptotic expansion

$$(2.3.33) \quad 1 - \Phi(x) = \frac{1}{\sqrt{2\pi}} \int_x^\infty e^{-z^2/2} dz = \frac{1}{x\sqrt{2\pi}} e^{-x^2/2} (1 + o(1))$$

as $x \rightarrow \infty$ (see [34]) we obtain (2.3.30) from equality (2.3.32) if condition (2.3.31) holds. If we put $\alpha_n = \alpha \in (0, 1)$ for all n (as in the Stein lemma), then (2.3.32) and (2.3.33) imply a stronger result:

$$(2.3.34) \quad \beta(\delta_n^{+, \alpha}) = \frac{1}{\sqrt{2\pi} v_n} \exp \left(-\frac{1}{2} v_n^2 + t_{1-\alpha} v_n - \frac{1}{2} t_{1-\alpha}^2 + o(1) \right)$$

or, in other words,

$$(2.3.35) \quad \ln \beta(\delta_n^{+, \alpha}) = -\frac{1}{2} v_n^2 + t_{1-\alpha} v_n - \ln v_n - \frac{1}{2} t_{1-\alpha}^2 - \ln \sqrt{2\pi} + o(1).$$

It is clear from asymptotic expansion (2.3.35) that the dependence on α in the behavior of $\beta(\delta_t^{+, \alpha})$ shows up in the second term of the asymptotic expansion of $\ln \beta(\delta_n^{+, \alpha})$ only, while the first term is independent of α (cf. (2.3.13)).

Large deviations. Implications (2.3.7) and (2.3.8) show that under condition $\Lambda 1$, relations $\beta 1$ and $\beta 2$ require conditions $\alpha 1$ and $\alpha 2$ to be satisfied; the latter conditions prohibit the level α_t to approach 0 and 1, respectively. However, if we impose a more restrictive condition on the likelihood ratio z_t , relations $\beta 1$ and $\beta 2$ can also be obtained for levels α_t that tend to 0 or 1 as $t \rightarrow \infty$, but rather slowly. To be more specific, let us introduce the following conditions:

$$\alpha 1') \quad \lim_{t \rightarrow \infty} \chi_t^{-1} \ln \alpha_t = 0;$$

$$\alpha 2') \quad \lim_{t \rightarrow \infty} \chi_t^{-1} \ln(1 - \alpha_t) = 0;$$

$$\Lambda 2) \quad \limsup_{\varepsilon \downarrow 0} \limsup_{t \rightarrow \infty} \varepsilon^{-1} \chi_t^{-1} \ln H_t(\varepsilon) \leq -1;$$

$$\Lambda 3) \quad \liminf_{\varepsilon \uparrow 0} \liminf_{t \rightarrow \infty} \varepsilon^{-1} \chi_t^{-1} \ln H_t(\varepsilon) \geq -1$$

where χ_t is the normalization occurring in condition $\Lambda 1$ and $H_t(\varepsilon) = H(\varepsilon; \tilde{\mathbf{P}}^t, \mathbf{P}^t)$ is the Hellinger integral of order ε for the measures $\tilde{\mathbf{P}}^t$ and \mathbf{P}^t .

Observe that the definition of $H_t(\varepsilon)$ gives

$$(2.3.36) \quad H_t(\varepsilon) = H(\varepsilon; \tilde{\mathbf{P}}^t, \mathbf{P}^t) = \begin{cases} \mathbf{P}^t(\tilde{\beta}_t > 0), & \varepsilon = 0, \\ \tilde{\mathbf{P}}^t(\beta_t > 0), & \varepsilon = 1, \\ \mathbf{E}_Q^t \tilde{\beta}_t^\varepsilon \beta_t^{1-\varepsilon}, & \varepsilon \neq 0, \varepsilon \neq 1. \end{cases}$$

Introduce the following notation:

$$(2.3.37) \quad \varepsilon_-^t = \inf\{\varepsilon: H_t(\varepsilon) > -\infty\}, \quad \varepsilon_+^t = \sup\{\varepsilon: H_t(\varepsilon) < \infty\}.$$

It is clear that $\varepsilon_-^t \leq 0$ and $\varepsilon_+^t \geq 1$, since $H_t(\varepsilon) \leq 1$ for $\varepsilon \in [0, 1]$. The following result gives a useful representation for the Hellinger integral $H_t(\varepsilon)$ in terms of the likelihood ratios z_t and \tilde{z}_t .

LEMMA 2.3.4. *For any $\varepsilon \in (\varepsilon_-^t, \varepsilon_+^t)$ different from 0 and 1,*

$$(2.3.38) \quad H_t(\varepsilon) = \mathbf{E}^t z_t^\varepsilon = \tilde{\mathbf{E}}^t \tilde{z}_t^{1-\varepsilon}.$$

PROOF. If $0 < \varepsilon < 1$, then

$$\begin{aligned} H_t(\varepsilon) &= \mathbf{E}_Q^t \tilde{\beta}_t^\varepsilon \beta_t^{1-\varepsilon} I(\beta_t > 0) = \mathbf{E}^t z_t^\varepsilon I(\beta_t > 0) = \mathbf{E}^t z_t^\varepsilon, \\ H_t(\varepsilon) &= \mathbf{E}_Q^t \tilde{\beta}_t^\varepsilon \beta_t^{1-\varepsilon} I(\tilde{\beta}_t > 0) = \tilde{\mathbf{E}}^t \tilde{z}_t^{1-\varepsilon} I(\tilde{\beta}_t > 0) = \tilde{\mathbf{E}}^t \tilde{z}_t^{1-\varepsilon} \end{aligned}$$

proving equalities (2.3.38) for $0 < \varepsilon < 1$.

Now let $\varepsilon_-^t < 0$ and $\varepsilon \in (\varepsilon_-^t, 0)$. Since $H_t(\varepsilon) < \infty$ by (2.3.37), we have $\mathbf{E}_Q^t \tilde{\beta}_t^\varepsilon \beta_t^{1-\varepsilon} I(\tilde{\beta}_t = 0, \beta_t > 0) < \infty$. This implies that $\mathbf{Q}(\tilde{\beta}_t = 0, \beta_t > 0) = 0$ and therefore $\mathbf{P}^t(\tilde{\beta}_t = 0)$. Thus

$$\begin{aligned} H_t(\varepsilon) &= \mathbf{E}_Q^t \tilde{\beta}_t^\varepsilon \beta_t^{1-\varepsilon} I(\tilde{\beta}_t > 0) = \mathbf{E}^t (\beta_t / \tilde{\beta}_t)^{-\varepsilon} = \mathbf{E}^t z_t^\varepsilon, \\ H_t(\varepsilon) &= \mathbf{E}_Q^t \tilde{\beta}_t^\varepsilon \beta_t^{1-\varepsilon} I(\tilde{\beta}_t > 0) = \tilde{\mathbf{E}}^t \tilde{z}_t^{1-\varepsilon} I(\tilde{\beta}_t > 0) = \tilde{\mathbf{E}}^t \tilde{z}_t^{1-\varepsilon}. \end{aligned}$$

Therefore equalities (2.3.38) are proved for $\varepsilon \in (\varepsilon_-^t, 0)$ and $\varepsilon_-^t < 0$.

Finally let $\varepsilon_+^t > 1$ and $\varepsilon \in (1, \varepsilon_+^t)$. Since $H_t(\varepsilon) < \infty$, we also have

$$\mathbf{E}_Q^t \tilde{\beta}_t^\varepsilon \beta_t^{1-\varepsilon} I(\tilde{\beta}_t > 0, \beta_t = 0) < \infty.$$

Then $\mathbf{Q}(\tilde{\beta}_t > 0, \beta_t = 0) = 0$ and hence $\tilde{\mathbf{P}}^t(\beta_t = 0)$. Therefore

$$\begin{aligned} H_t(\varepsilon) &= \tilde{\mathbf{E}}_Q^t \tilde{\beta}_t^{\varepsilon-1} \beta_t^{1-\varepsilon} I(\beta_t > 0) = \tilde{\mathbf{E}}^t (\tilde{\beta}_t / \beta_t)^{\varepsilon-1} = \tilde{\mathbf{E}}^t \tilde{z}_t^{1-\varepsilon}, \\ H_t(\varepsilon) &= \mathbf{E}_Q^t \tilde{\beta}_t^\varepsilon \beta_t^{1-\varepsilon} I(\beta_t > 0) = \mathbf{E}^t z_t^\varepsilon I(\beta_t > 0) = \mathbf{E}^t z_t^\varepsilon. \end{aligned}$$

This proves equalities (2.3.38) for $\varepsilon \in (1, \varepsilon_+^t)$ and $\varepsilon_+^t > 1$. □

REMARK 2.3.3. If $H_t(\varepsilon_-^t) < \infty$, then the same argument as that used to prove Lemma 2.3.4 shows that equality (2.3.38) holds for $\varepsilon = \varepsilon_-^t$, too. If $H_t(\varepsilon_+^t) < \infty$, then equality (2.3.38) holds for $\varepsilon = \varepsilon_+^t$, too. Put $E^t z_t^\varepsilon = P^t(\tilde{\mathfrak{J}}_t > 0)$ in the case $\varepsilon = 0$ and $P^t(\tilde{\mathfrak{J}}_t = 0) > 0$ and put $\tilde{E}^t \tilde{z}_t^\varepsilon = \tilde{P}^t(\tilde{\mathfrak{J}}_t > 0)$ in the case $\varepsilon = 0$ and $\tilde{P}^t(\tilde{\mathfrak{J}}_t = 0) > 0$. Then equality (2.3.38) holds for any $\varepsilon \in (-\infty, \infty)$. Moreover we get for any $\varepsilon \in (-\infty, \infty)$

$$(2.3.39) \quad H_t(\varepsilon) = E^t e^{\varepsilon \Lambda_t},$$

that is, $H_t(\varepsilon)$ is the moment generating function of the random variable Λ_t which is, generally speaking, an extended random variable since $P^t(\Lambda_t = -\infty)$ can be positive.

REMARK 2.3.4. If $\varepsilon_-^t < 0$, then the proof of Lemma 2.3.4 shows that

$$P^t(\tilde{\mathfrak{J}}_t = 0) = 0.$$

By (1.1.18), this means that $P^t \ll \tilde{P}^t$. Moreover $H_t(0) = 1$ by (2.3.36) in this case. If $\varepsilon_+^t > 1$, then it has also been shown in the proof of Lemma 2.3.4 that $\tilde{P}^t(\tilde{\mathfrak{J}}_t = 0) = 0$. By (1.1.17), this means that $\tilde{P}^t \ll P^t$. Relation (2.3.36) implies in this case that $H_t(1) = 1$.

REMARK 2.3.5. Condition $\Lambda 3$ implies that there exist numbers $\varepsilon_0 < 0$ and $t_0 < \infty$ such that $H_t(\varepsilon) < \infty$ for all $\varepsilon \in (\varepsilon_0, 0)$ and $t > t_0$. Then we have by Remark 2.3.4 that

$$\Lambda 3 \Rightarrow P^t \ll \tilde{P}^t \quad \text{for all } t > t_0.$$

In view of (1.1.18), this implies

$$\Lambda 3 \Rightarrow \bar{\alpha}_t = 1 \quad \text{for all } t > t_0.$$

The following result establishes a relationship between conditions $\Lambda 2$, $\alpha 1$, $\alpha 2$ and conditions $\Lambda 2$, $\Lambda 3$, $\alpha 1'$, $\alpha 2'$.

LEMMA 2.3.5. *We have*

$$(2.3.40) \quad \alpha 1 \Rightarrow \alpha 1'; \quad \alpha 2 \Rightarrow \alpha 2';$$

$$(2.3.41) \quad \Lambda 2, \Lambda 3 \Rightarrow \Lambda 1.$$

PROOF. Implications (2.3.40) are obvious. To prove implication (2.3.41) let $\gamma > 0$ be an arbitrary number. Then, by condition $\Lambda 2$, there exists a positive number $\varepsilon_0 = \varepsilon_0(\gamma)$ such that

$$(2.3.42) \quad \limsup_{t \rightarrow \infty} \chi_t^{-1} \ln H_t(\varepsilon) \leq -\varepsilon + \frac{1}{2} \varepsilon \gamma$$

for all $\varepsilon \in (0, \varepsilon_0)$. Fix some $\varepsilon \in (0, \varepsilon_0)$. Then, in view of (2.3.42), there exists a number $t_0 = t_0(\varepsilon, \gamma)$ such that for all $t > t_0$

$$\chi_t^{-1} \ln H_t(\varepsilon) \leq -\varepsilon + \frac{1}{2} \varepsilon \gamma + \frac{1}{4} \varepsilon \gamma,$$

whence we obtain for $t > t_0$ that

$$(2.3.43) \quad \begin{aligned} \mathbf{P}^t (\chi_t^{-1} \Lambda_t > -1 + \gamma) &= \mathbf{P}^t (z_t^\varepsilon > e^{-\varepsilon(1-\gamma)\chi_t}) \leq e^{\varepsilon(1-\gamma)\chi_t} \mathbf{E}^t z_t^\varepsilon \\ &= e^{\varepsilon(1-\gamma)\chi_t} H_t(\varepsilon) \leq \exp\left(-\frac{1}{4}\varepsilon\gamma\chi_t\right) \end{aligned}$$

by the Chebyshev inequality and equality (2.3.38) for $\varepsilon \in (0, 1)$. Following a similar argument, we obtain from condition $\Lambda 3$ that for any $\gamma > 0$ there exists $\varepsilon_1 < 0$ such that

$$(2.3.44) \quad \mathbf{P}^t (\chi_t^{-1} \Lambda_t < -1 - \gamma) \leq \exp\left(\frac{1}{4}\varepsilon\gamma\chi_t\right)$$

for all $\varepsilon \in (\varepsilon_1, 0)$ and all $t > t_1(\varepsilon, \gamma)$. Bounds (2.3.43) and (2.3.44) imply that condition $\Lambda 1$ holds. Thus implication (2.3.41) is proved. \square

Conditions $\Lambda 2$ and $\Lambda 3$ are related to a theorem on large deviations for Λ_t as $t \rightarrow \infty$. To state it we introduce the following condition.

Λ^* . For any $\varepsilon \in (-\infty, \infty)$, the limit

$$(2.3.45) \quad \lim_{t \rightarrow \infty} \varphi_t^{-1} \ln H_t(\varepsilon) = \varkappa(\varepsilon)$$

exists where $\varphi_t \rightarrow \infty$ as $t \rightarrow \infty$ and $\varkappa(\varepsilon)$ is a proper convex function differentiable in the interval $(\varepsilon_-, \varepsilon_+)$ where

$$(2.3.46) \quad \varepsilon_- = \inf\{\varepsilon: \varkappa(\varepsilon) < \infty\}, \quad \varepsilon_+ = \sup\{\varepsilon: \varkappa(\varepsilon) < \infty\}.$$

It is clear that $\varepsilon_- \leq 0$ and $\varepsilon_+ \geq 1$. Let

$$(2.3.47) \quad \begin{aligned} \gamma_0 &= \varkappa'(0), & \gamma_1 &= \varkappa'(1), \\ \gamma_- &= \lim_{\varepsilon \downarrow \varepsilon_-} \varkappa'(\varepsilon), & \gamma_+ &= \lim_{\varepsilon \uparrow \varepsilon_+} \varkappa'(\varepsilon). \end{aligned}$$

Note that γ_0 is defined for $\varepsilon_- < 0$ only and γ_1 is defined for $\varepsilon_+ > 1$. If condition Λ^* holds with $\varepsilon_- < 0$, then it is clear that conditions $\Lambda 2$ and $\Lambda 3$ hold with $\chi_t = -\gamma_0\varphi_t$. It follows from the properties of the function $H_t(\varepsilon)$ that $\gamma_0 < 0$. In what follows, we will reveal a tighter relationship between conditions $\Lambda 2$, $\Lambda 3$ and condition Λ^* providing a theorem on large deviations for Λ_t .

The following result gives upper and lower bounds for $\beta(\delta_t^{+, \alpha})$ for all $\alpha \in (0, 1)$ enabling us to obtain $\beta 1$ and $\beta 2$ if conditions $\alpha 1'$ and $\alpha 2'$ hold.

THEOREM 2.3.2. *For all $\alpha \in (0, 1)$ and all $t \in \mathbf{R}_+$*

$$(2.3.48) \quad \beta(\delta_t^{+, \alpha}) \geq (1 - \alpha)^{\varepsilon/(\varepsilon-1)} (H_t(1 - \varepsilon))^{1/(1-\varepsilon)}, \quad \varepsilon > 1,$$

$$(2.3.49) \quad \beta(\delta_t^{+, \alpha}) \leq (1 - \varepsilon)(\varepsilon/\alpha)^{\varepsilon/(1-\varepsilon)} (H_t(1 - \varepsilon))^{1/(1-\varepsilon)}, \quad 0 < \varepsilon < 1.$$

PROOF. If $H_t(1 - \varepsilon) = \infty$ for $\varepsilon > 1$, then estimate (2.3.48) is trivial. Therefore we assume that $H_t(1 - \varepsilon) < \infty$ for $\varepsilon > 1$. By the definition of a Bayes test, we have for any $c \in (0, \infty)$

$$(2.3.50) \quad \begin{aligned} c\alpha + \beta(\delta_t^{+, \alpha}) &\geq \inf\{c\alpha(\delta_t) + \beta(\delta_t): \delta \in \Sigma^t\} = c\alpha \left(\delta_t^{c, 1}\right) + \beta \left(\delta_t^{c, 1}\right) \\ &= 1 - \mathbf{E}_Q^t (\tilde{\beta}_t - c\beta_t)^+ \end{aligned}$$

where $a^+ = a \vee 0$. By the fundamental Neyman–Pearson lemma and by the definition of a Bayes test, there exists a constant $c \in (0, \infty)$ such that inequality (2.3.50) becomes an equality. Therefore

$$(2.3.51) \quad \beta(\delta_t^{+, \alpha}) = 1 - \inf \{ c\alpha + E_Q^t(\tilde{\mathfrak{J}}_t - c\mathfrak{J}_t)^+; c \geq 0 \}.$$

The condition $H_t(1 - \varepsilon) < \infty$ for $\varepsilon > 1$ implies that $\mathbf{Q}(\tilde{\mathfrak{J}}_t = 0) = 0$. Hence

$$(2.3.52) \quad E_Q^t(\tilde{\mathfrak{J}}_t - c\mathfrak{J}_t)^+ = \tilde{E}^t(1 \vee cz_t) - c.$$

Since $1 \vee z \leq az^\varepsilon + 1$ for $\varepsilon > 1$, $z \geq 0$, and $a = \varepsilon^{-\varepsilon}(\varepsilon - 1)^{\varepsilon-1}$, relations (2.3.51) and (2.3.52) imply that

$$(2.3.53) \quad \beta(\delta_t^{+, \alpha}) \geq \sup \{ (1 - \alpha)c - ac^\varepsilon H_t(1 - \varepsilon); c \geq 0 \}$$

by (2.3.38) for $\varepsilon > 1$. The upper bound (2.3.53) is attained at

$$(2.3.54) \quad c = c^* = (1 - \alpha)^{1/(\varepsilon-1)}(a\varepsilon H_t(1 - \varepsilon))^{1/(1-\varepsilon)}.$$

Therefore (2.3.53) and (2.3.54) imply estimate (2.3.48).

Now we prove estimate (2.3.49). As in the proof of equality (2.3.52) we get

$$(2.3.55) \quad E_Q^t(\tilde{\mathfrak{J}}_t - c\mathfrak{J}_t)^+ = 1 - \tilde{E}^t(1 \wedge c\tilde{z}_t).$$

Since

$$(2.3.56) \quad z \wedge 1 \leq z^\varepsilon, \quad z \geq 0, \quad 0 < \varepsilon < 1,$$

we obtain from (2.3.51), (2.3.55), and (2.3.38) for $0 < \varepsilon < 1$ that

$$(2.3.57) \quad \beta(\delta_t^{+, \alpha}) \leq \sup \{ c^\varepsilon H_t(1 - \varepsilon) - c\alpha; c \geq 0 \}.$$

The upper bound in (2.3.57) is attained at

$$(2.3.58) \quad c = c^* = \alpha^{1/(\varepsilon-1)}(\varepsilon H_t(1 - \varepsilon))^{1/(1-\varepsilon)}.$$

Therefore (2.3.57) and (2.3.58) imply estimate (2.3.49). \square

COROLLARY 2.3.4. *For all $\varepsilon, \alpha \in (0, 1)$ and $t \in \mathbf{R}_+$*

$$H_t(\varepsilon) \geq \beta^\varepsilon(\delta_t^{+, \alpha})\alpha^{1-\varepsilon}e^{h(\varepsilon)}$$

where $h(\varepsilon) = -\varepsilon \ln \varepsilon - (1 - \varepsilon) \ln(1 - \varepsilon)$ is the Shannon entropy of the distribution of the random variable taking two values with probabilities ε and $1 - \varepsilon$.

REMARK 2.3.6. If $\bar{\alpha}_t < 1$, then $H_t(\varepsilon) = \infty$ for all $\varepsilon < 0$. Therefore, in the case of $\bar{\alpha}_t < 1$, estimate (2.3.48) acquires the trivial form $\beta(\delta_t^{+, \alpha}) \geq 0$ for all $\alpha \in (0, 1)$. We also note that (1.1.50) implies $\beta(\delta_t^{+, \alpha}) = 0$ for $\alpha \in [\bar{\alpha}_t, 1]$. Estimate (2.3.49) is rather rough for $\alpha \in [\bar{\alpha}_t, 1]$, however $\bar{\alpha}_t \rightarrow 1$ as $t \rightarrow \infty$ by (2.3.5) and if condition $\Lambda 1$ holds. Moreover if condition $\Lambda 3$ holds, then Remark 2.3.5 yields that there exists $t_0 < \infty$ such that $\bar{\alpha}_t = 1$ for all $t > t_0$ (see Theorem 2.7.2 in [37] for improved estimates (2.3.48) and (2.3.49)).

THEOREM 2.3.3. *The following implications hold:*

$$(2.3.59) \quad \Lambda 2, \alpha 1' \Rightarrow d1 \Rightarrow \beta 1,$$

$$(2.3.60) \quad \Lambda 3, \alpha 2' \Rightarrow \beta 2 \Rightarrow d2.$$

PROOF. According to Theorem 2.3.1, it is sufficient to prove the first implications in (2.3.59) and (2.3.60). The first equality in (2.3.50) implies that for all $c > 0$

$$c\alpha\left(\delta_t^{c,1}\right) + \beta\left(\delta_t^{c,1}\right) \leq c^\varepsilon \tilde{E}^t \tilde{z}_t^\varepsilon, \quad 0 < \varepsilon < 1,$$

in view of (2.3.55) and (2.3.56). Thus

$$e^{d_t} \alpha_t + \beta\left(\delta_t^{+, \alpha_t}\right) \leq e^{\varepsilon d_t} H_t(1 - \varepsilon)$$

for $0 < \varepsilon < 1$ by the first equality in (2.3.50) and by equality (2.3.38). Therefore we have for $0 < \varepsilon < 1$

$$d_t \leq (1 - \varepsilon)^{-1} \ln H_t(1 - \varepsilon) - (1 - \varepsilon)^{-1} \ln \alpha_t,$$

whence the implication $\Lambda 2, \alpha 1' \Rightarrow \beta 1$ follows.

The implication $\Lambda 3, \alpha 2' \Rightarrow \beta 2$ follows from estimate (2.3.48). \square

REMARK 2.3.7. The implication $\Lambda 2, \alpha 1' \Rightarrow \beta 1$ can be obtained directly from estimate (2.3.49).

Theorem 2.3.3 implies the following well-known result of Krafft and Plachky (see [37]).

COROLLARY 2.3.5. *Let the assumptions of Corollary 2.3.3 hold. If conditions $\alpha 1'$ and $\alpha 2'$ are satisfied with $\chi_n = n$, then (2.3.22) holds.*

EXAMPLE 2.3.2. This is a continuation of Example 2.3.1. Relations (2.3.23)–(2.3.25) imply for all $t \in \mathbf{R}_+$ and $\varepsilon \in (-\infty, \infty)$ that

$$(2.3.61) \quad H_n(\varepsilon) = \exp \left\{ -\frac{\varepsilon(1 - \varepsilon)}{2} v_n^2 \right\},$$

whence

$$\lim_{\varepsilon \rightarrow 0} \lim_{n \rightarrow \infty} (\varepsilon \chi_n)^{-1} \ln H_n(\varepsilon) = -1$$

for $\chi_n = 2^{-1} v_n^2$, that is, conditions $\Lambda 2$ and $\Lambda 3$ hold. Theorem 2.3.3 implies (2.3.30) for the test δ_n^{+, α_n} of level α_n satisfying conditions $\alpha 1'$ and $\alpha 2'$.

The following result of independent interest holds under the assumptions of the latter example.

THEOREM 2.3.4. *If $v_n \rightarrow \infty$ as $n \rightarrow \infty$, then*

1) *if $\alpha_n \rightarrow 0$ as $n \rightarrow \infty$, then*

$$(2.3.62) \quad \alpha 1' \iff d 1 \iff \beta 1 \iff z_{1 - \alpha_n} = o(v_n);$$

2) *if $\alpha_n \rightarrow 1$ as $n \rightarrow \infty$, then*

$$(2.3.63) \quad \alpha 2' \iff d 2 \iff \beta 2 \iff z_{1 - \alpha_n} = o(v_n)$$

where z_p is a p -quantile of the distribution $\mathcal{N}(0, 1)$ and $\chi_n = 2^{-1} v_n^2$.

PROOF. Let $\alpha_n \rightarrow 0$ as $n \rightarrow \infty$. Then it follows from (2.3.33) that

$$(2.3.64) \quad \ln \alpha_n = -2^{-1} z_{1-\alpha_n}^2 (1 + o(1)), \quad n \rightarrow \infty.$$

Therefore

$$(2.3.65) \quad \alpha 1' \iff z_{1-\alpha_n} = o(v_n).$$

Since $v_n \rightarrow \infty$ as $n \rightarrow \infty$, conditions $\Lambda 1$, $\Lambda 2$, and $\Lambda 3$ hold for $\chi_n = 2^{-1} v_n^2$. Therefore, in view of Theorems 2.3.1 and 2.3.3, relation (2.3.62) follows from the chain of implications

$$(2.3.66) \quad \beta 1 \Rightarrow d 1 \Rightarrow \alpha 1'.$$

Since $z_{1-\alpha_n} \rightarrow \infty$ as $n \rightarrow \infty$, it follows from (2.3.28) that $d 1 \Rightarrow z_{1-\alpha_n} = o(v_n)$, where the implication $d 1 \Rightarrow \alpha 1'$ follows in view of (2.3.65).

If $\beta 1$ holds, we have by (2.3.32) that $v_n - z_{1-\alpha_n} \rightarrow \infty$ as $n \rightarrow \infty$. Applying (2.3.33) we get

$$\lim_{n \rightarrow \infty} (v_n - z_{1-\alpha_n})^{-2} \ln \beta(\delta_n^{+, \alpha_n}) = -2^{-1}.$$

Thus $z_{1-\alpha_n} = o(v_n)$ by $\beta 1$. Hence the implication $\beta 1 \Rightarrow d 1$ follows from (2.3.65) and the equivalence $d 1 \iff \alpha 1'$ proved above. Therefore implication (2.3.66) is proved.

The proof of (2.3.63) is similar to that of (2.3.62) and thus is omitted. \square

REMARK 2.3.8. Since condition $\Lambda 2$ holds, we see from statement 1) in Theorem 2.3.4 that conditions $\alpha 1'$ and $d 1$ cannot be weakened to prove the chain of implications $\Lambda 2$, $\alpha 1' \Rightarrow d 1 \Rightarrow \beta 1$ in Theorem 2.3.3. Since condition $\Lambda 3$ also holds, we obtain from statement 2) in Theorem 2.3.4 that conditions $\alpha 2'$ and $\beta 2$ cannot be weakened in the chain of implications $\Lambda 3$, $\alpha 2' \Rightarrow \beta 2 \Rightarrow d 2$ in Theorem 2.3.3. Moreover, it follows from Theorem 2.3.4 that

$$(\alpha 1', \alpha 2') \iff (d 1, d 2) \iff (\beta 1, \beta 2) \iff z_{1-\alpha_n} = o(v_n),$$

that is, relation (2.3.30) is equivalent to conditions $\alpha 1'$ and $\alpha 2'$.

Rates of decay of probabilities of error of the Neyman–Pearson, Bayes, and minimax tests under condition Λ^* . Throughout this section we assume that condition Λ^* holds. In this case, the Neyman–Pearson test $\delta_t^{+, \alpha}$ can be used even if the level α_t tends to zero as $t \rightarrow \infty$ faster than is allowed by condition $\alpha 1'$. In particular, if $\varphi_t^{-1} \ln \alpha_t \rightarrow -a$ for some positive number a , then one can prove that $\varphi_t^{-1} \ln \beta(\delta_t^{+, \alpha}) \rightarrow -b(a)$ where $b(a)$ is a positive function of a . In order to provide an exact statement, we need the following result on large deviations of Λ_t . Below we use the following notation (see also (2.3.45)–(2.3.47)):

$$(2.3.67) \quad \Gamma_0 = \gamma_0 I(\varepsilon_- < 0) + \gamma_- I(\varepsilon_- = 0),$$

$$(2.3.68) \quad \Gamma_1 = \gamma_1 I(\varepsilon_+ > 1) + \gamma_+ I(\varepsilon_+ = 1).$$

THEOREM 2.3.5. *Let condition Λ^* be satisfied. If*

$$\Gamma_0 < \gamma_+,$$

then for all $\gamma \in (\Gamma_0, \gamma_+)$

$$(2.3.69) \quad \lim_{t \rightarrow \infty} \varphi_t^{-1} \ln P^t(\varphi_t^{-1} \Lambda_t > \gamma) = \lim_{t \rightarrow \infty} \varphi_t^{-1} \ln P^t(\varphi_t^{-1} \Lambda_t \geq \gamma) = -I(\gamma).$$

Further if $\gamma_- < \Gamma_1$, then for all $\gamma \in (\gamma_-, \Gamma_1)$

$$(2.3.70) \quad \lim_{t \rightarrow \infty} \varphi_t^{-1} \ln \tilde{P}^t(\varphi_t^{-1} \Lambda_t < \gamma) = \lim_{t \rightarrow \infty} \varphi_t^{-1} \ln \tilde{P}^t(\varphi_t^{-1} \Lambda_t \leq \gamma) = -I(\gamma) + \gamma$$

where $I(\gamma) = \gamma \varepsilon(\gamma) - \varkappa(\varepsilon(\gamma))$ and $\varepsilon(\gamma)$ is an arbitrary solution of the equation $\varkappa'(\varepsilon) = \gamma$.

Theorem 2.3.5 can be deduced from Theorem 2.6.3 of [17] in view of equality (2.3.39). This proof can be found in [39].

REMARK 2.3.9. Applying the methods of convex analysis [44], we can readily obtain that the function $I(\gamma)$ is strictly convex in the interval (γ_-, γ_+) and has the unique minimum at $\gamma = \gamma_0$ if $\varepsilon_- < 0$. Moreover $I(\gamma_0) = 0$ in this case, while the minimum is attained at $\gamma = \gamma_-$ if $\varepsilon_- = 0$.

First we consider the Bayes test δ_t^π with respect to the a priori distribution $(\pi, \tilde{\pi})$, $\pi + \tilde{\pi} = 1$, and the loss $A_{ij} = 1 - \delta_{ij}$ (see Section 2.1). It follows from Section 1.2 that we can put $\delta_t^\pi = \delta_t^{c, \varepsilon}$ where $c = \pi/\tilde{\pi}$ and $\varepsilon \in [0, 1]$ is an arbitrary number.

THEOREM 2.3.6. *Let condition Λ^* hold with $\Gamma_0 < 0 < \Gamma_1$. Then*

$$(2.3.71) \quad \lim_{t \rightarrow \infty} \varphi_t^{-1} \ln \alpha(\delta_t^\pi) = \lim_{t \rightarrow \infty} \varphi_t^{-1} \ln \beta(\delta_t^\pi) = -I(0),$$

$$(2.3.72) \quad \lim_{t \rightarrow \infty} \varphi_t^{-1} \ln e_\pi(\delta_t^\pi) = -I(0)$$

where e_π is the probability of error of the test δ_t^π (see (1.1.32)).

PROOF. It is clear that for any Bayes test δ_t^π

$$(2.3.73) \quad P^t(\Lambda_t > \ln(\pi/\tilde{\pi})) \leq \alpha(\delta_t^\pi) \leq P^t(\Lambda_t \geq \ln(\pi/\tilde{\pi})),$$

$$(2.3.74) \quad \tilde{P}^t(\Lambda_t < \ln(\pi/\tilde{\pi})) \leq \beta(\delta_t^\pi) \leq \tilde{P}^t(\Lambda_t \leq \ln(\pi/\tilde{\pi})).$$

Then (2.3.69), (2.3.70), (2.3.73), and (2.3.74) imply (2.3.71). Now the equality

$$e_\pi(\delta_t^\pi) = \pi \alpha(\delta_t^\pi) + \tilde{\pi} \beta(\delta_t^\pi)$$

and (2.3.71) imply (2.3.72). □

REMARK 2.3.10. By (2.3.67) and (2.3.68), the condition $\Gamma_0 < 0 < \Gamma_1$ implies that

$$\Gamma_0 < \gamma_+ \quad \text{and} \quad \gamma_- < \Gamma_1,$$

which enables us to apply relation (2.3.69) for $\gamma \in (\Gamma_0, \gamma_+)$ and relation (2.3.70) for $\gamma \in (\gamma_-, \Gamma_1)$. On the other hand,

$$\varphi_t^{-1} \ln(\pi/\tilde{\pi}) \rightarrow \gamma = 0 \in (\Gamma_0, \Gamma_1) \quad \text{as } t \rightarrow \infty.$$

Now let δ_t^* be the minimax test for distinguishing the hypotheses H^t and \tilde{H}^t (see Section 2.1). According to the results of Section 1.2, $e(\delta_t^*)$ is the probability of error for the test δ_t^* (see (1.2.19)).

THEOREM 2.3.7. *Let condition Λ^* hold with $\Gamma_0 < 0 < \Gamma_1$. Then*

$$(2.3.75) \quad \lim_{t \rightarrow \infty} \varphi_t^{-1} \ln \alpha(\delta_t^*) = \lim_{t \rightarrow \infty} \varphi_t^{-1} \ln \beta(\delta_t^*) = -I(0),$$

$$(2.3.76) \quad \lim_{t \rightarrow \infty} \varphi_t^{-1} \ln e_\pi(\delta_t^*) = -I(0).$$

PROOF. Let $(\pi, \tilde{\pi})$ be the a priori distribution of the hypotheses, $\pi + \tilde{\pi} = 1$. Then the definitions of Bayes tests and minimax tests imply that

$$\begin{aligned} e_\pi(\delta_t^\pi) &= \min_{\delta_t} e_\pi(\delta_t) \leq \min_{\delta_t} (\alpha(\delta_t) \vee \beta(\delta_t)) = e(\delta_t^*), \\ e_\pi(\delta_t^\pi) &\geq (\pi \wedge \tilde{\pi})(\alpha(\delta_t^\pi) \vee \beta(\delta_t^\pi)) \geq (\pi \wedge \tilde{\pi})e(\delta_t^*), \end{aligned}$$

that is,

$$(2.3.77) \quad (\pi \wedge \tilde{\pi})e(\delta_t^*) \leq e_\pi(\delta_t^*) \leq e(\delta_t^*).$$

Combining (2.3.72) and (2.3.77), we obtain (2.3.76). Relation (2.3.75) follows from (2.3.76), since $\alpha(\delta_t^*) = \beta(\delta_t^*) = e(\delta_t^*)$ (see Theorem 1.2.4). \square

The following result describes a relationship between the rates of decay of the level α_t and the type II error probability $\beta(\delta_t^{+, \alpha_t})$ for the Neyman–Pearson test under condition Λ^* .

THEOREM 2.3.8. *Let condition Λ^* hold with $\Gamma_0 < \Gamma_1$. Then*

$$(2.3.78) \quad \lim_{t \rightarrow \infty} \varphi_t^{-1} \ln \alpha_t = -a \iff \lim_{t \rightarrow \infty} \varphi_t^{-1} \ln \beta_t(\delta_t^{+, \alpha_t}) = -b(a)$$

for any $a \in (I(\Gamma_0), I(\Gamma_1))$ where

$$b(a) = a - \gamma(a) \in (I(\Gamma_1) - \Gamma_1, I(\Gamma_0) - \Gamma_0)$$

and $\gamma(a)$ is a unique solution of the equation $I(\gamma) = a$.

PROOF. We prove the implication \Rightarrow in (2.3.78). Assume that $\varphi_t^{-1} \ln \alpha_t \rightarrow -a$ as $t \rightarrow \infty$ for $a \in (I(\Gamma_0), I(\Gamma_1))$. We have

$$\alpha_t = \mathbf{P}^t(Y_t > y_t) + \varepsilon \mathbf{P}^t(Y_t = y_t)$$

where $Y_t = \varphi_t^{-1} \Lambda_t$ and $y_t = \varphi_t^{-1} d_t$. Our current goal is to show that $y_t \rightarrow \gamma(a)$ as $t \rightarrow \infty$. Put

$$\underline{y} = \liminf_{t \rightarrow \infty} y_t, \quad \bar{y} = \limsup_{t \rightarrow \infty} y_t.$$

Then it is sufficient to prove that $\underline{y} = \gamma(a)$ and $\bar{y} = \gamma(a)$.

First we show that $\underline{y} = \gamma(a)$. Assume for contradiction that $\underline{y} \neq \gamma(a)$. Below we use the obvious estimates

$$(2.3.79) \quad \mathbf{P}^t(Y_t > y_t) \leq \alpha_t \leq \mathbf{P}^t(Y_t \geq y_t).$$

By the definition of \underline{y} , there exists a sequence (t_n) such that $t_n \rightarrow \infty$ and $y_{t_n} \rightarrow \underline{y}$ as $n \rightarrow \infty$. Then, if \underline{y} is finite, for any $y' < \underline{y}$ and $y'' > \underline{y}$ there exists $n_0 = n_0(y', y'')$ such that $y' < y_{t_n} < y''$ for all $n > n_0$. If $\underline{y} = +\infty$, then for any $y' < \underline{y} = \infty$ there exists $n_0 = n_0(y')$ such that $y_{t_n} > y'$ for all $n > n_0$, while if $\underline{y} = -\infty$, then for any $y'' > \underline{y} = -\infty$ there exists $n_0 = n_0(y'')$ such that $y_{t_n} < y''$ for all $n > n_0$.

First let $\underline{y} \leq \Gamma_0$ and $y'' \in (\Gamma_0, \gamma(a))$. Since $y_{t_n} < y''$ for all $n > n_0$, we have by Theorem 2.3.5

$$\liminf_{n \rightarrow \infty} \varphi_{t_n}^{-1} \ln P^{t_n}(Y_{t_n} > y_{t_n}) \geq \lim_{n \rightarrow \infty} \varphi_{t_n}^{-1} \ln P^{t_n}(Y_{t_n} > y'') = -I(y'').$$

Therefore (2.3.79) implies

$$\lim_{n \rightarrow \infty} \varphi_{t_n}^{-1} \ln \alpha^{t_n} \geq -I(y'') > -I(\gamma(a)) = -a,$$

since $\Gamma_0 < y'' < \gamma(a)$. By Remark 2.3.9, the function $I(\gamma)$ is strictly increasing on the interval (Γ_0, γ_+) , giving a contradiction.

Now assume that $\underline{y} \in (\Gamma_0, \gamma(a)) \cup (\gamma(a), \gamma_+)$ and that y' and y'' are such that $(y', y'') \subset (\Gamma_0, \gamma(a)) \cup (\gamma(a), \gamma_+)$. Then by Theorem 2.3.5 we have

$$\liminf_{n \rightarrow \infty} \varphi_{t_n}^{-1} \ln P^{t_n}(Y_{t_n} > y_{t_n}) \geq \lim_{n \rightarrow \infty} \varphi_{t_n}^{-1} \ln P^{t_n}(Y_{t_n} > y'') = -I(y''),$$

$$\limsup_{n \rightarrow \infty} \varphi_{t_n}^{-1} \ln P^{t_n}(Y_{t_n} \geq y_{t_n}) \leq \lim_{n \rightarrow \infty} \varphi_{t_n}^{-1} \ln P^{t_n}(Y_{t_n} \geq y') = -I(y').$$

Since y' and y'' are arbitrary and since the function $I(\gamma)$ is continuous on the interval (γ_-, γ_+) , we obtain in view of inequalities (2.3.79) that

$$\lim_{n \rightarrow \infty} \varphi_{t_n}^{-1} \ln \alpha_{t_n} = -I(\underline{y}).$$

Since $\underline{y} \neq \gamma(a)$ by assumption and since $I(\gamma)$ is strictly increasing on (Γ_0, γ_+) , this again gives a contradiction.

Finally, let $\underline{y} > \gamma_+$ and $y' \in (\gamma(a), \gamma_+)$. Since $y_{t_n} > y'$ for all $n > n_0$, we apply Theorem 2.3.5 once more to obtain

$$\limsup_{n \rightarrow \infty} \varphi_{t_n}^{-1} \ln P^{t_n}(Y_{t_n} \geq y_{t_n}) \leq \lim_{n \rightarrow \infty} \varphi_{t_n}^{-1} \ln P^{t_n}(Y_{t_n} \geq y') = -I(y').$$

Since $\gamma(a) < y' < \gamma_+$, we have $I(y') > I(\gamma(a)) = a$, giving a contradiction by (2.3.79).

The above contradictions show that $\underline{y} = \gamma(a)$. By a similar argument, we can prove that $\bar{y} = \gamma(a)$. Therefore $y_t \rightarrow \gamma(a)$ as $t \rightarrow \infty$.

Let $\varepsilon > 0$ be arbitrary and let $t = t_0(\varepsilon)$ be such that $|y_t - \gamma(a)| < \varepsilon$ for all $t > t_0$. Assume that $\varepsilon > 0$ is small enough in order that $(\gamma(a) - \varepsilon, \gamma(a) + \varepsilon) \subset (\Gamma_0, \Gamma_1)$. It is clear that the following inequalities hold for all $t > t_0$:

$$\tilde{P}^t(\varphi_t^{-1} \Lambda_t < \gamma(a) - \varepsilon) \leq \beta(\delta_t^{+, \alpha_t}) \leq \tilde{P}^t(\varphi_t^{-1} \Lambda_t \leq \gamma(a) + \varepsilon),$$

whence

$$\liminf_{t \rightarrow \infty} \varphi_t^{-1} \ln \beta(\delta_t^{+, \alpha_t}) \geq -I(\gamma(a) - \varepsilon) + \gamma(a) - \varepsilon,$$

$$\limsup_{t \rightarrow \infty} \varphi_t^{-1} \ln \beta(\delta_t^{+, \alpha_t}) \leq -I(\gamma(a) + \varepsilon) + \gamma(a) + \varepsilon$$

by Theorem 2.3.5. Approaching the limit in these inequalities as $\varepsilon \rightarrow 0$ we obtain that the upper and lower bounds in these inequalities coincide and

$$\lim_{t \rightarrow \infty} \varphi_t^{-1} \ln \beta(\delta_t^{+, \alpha_t}) = -a + \gamma(a) = -b(a),$$

since the function $I(\gamma)$ is continuous in the interval (γ_-, γ_+) and $I(\gamma(a)) = a$.

Therefore the implication \Rightarrow in (2.3.78) is proved. The inverse implication \Leftarrow in (2.3.78) can be proved along similar lines. \square

REMARK 2.3.11. The behavior of $\beta(\delta_t^{+, \alpha_t})$ in the case where $\varphi_t^{-1} \ln \alpha_t \rightarrow -a$ for $a \notin (I(\Gamma_0), I(\Gamma_1))$ is studied in [39]. In particular, the results in [39] show that for any $a \in [0, I(\Gamma_0)]$ we have

$$(2.3.80) \quad \lim_{t \rightarrow \infty} \varphi_t^{-1} \ln \alpha_t = -a \Rightarrow \limsup_{t \rightarrow \infty} \varphi_t^{-1} \ln \beta(\delta_t^{+, \alpha_t}) \leq \Gamma_0 - I(\Gamma_0).$$

Let $\varepsilon_- < 0$. Then $\Gamma_0 = \gamma_0$ and $I(\Gamma_0) = 0$. Therefore (2.3.80) becomes of the form

$$(2.3.81) \quad \lim_{t \rightarrow \infty} \varphi_t^{-1} \ln \alpha_t = a = 0 \Rightarrow \limsup_{t \rightarrow \infty} \varphi_t^{-1} \ln \beta(\delta_t^{+, \alpha_t}) \leq \gamma_0.$$

As we have already observed, in this case conditions $\Lambda 2$ and $\Lambda 3$ hold with

$$\chi_t = -\gamma_0 \varphi_t.$$

Therefore implication (2.3.81) is equivalent to the implication $\alpha 1' \Rightarrow \beta 1$. If condition $\alpha 2'$ holds, then we obtain from the implication $\alpha 2' \Rightarrow \beta 2$ that

$$\alpha 1', \alpha 2' \Rightarrow \lim_{t \rightarrow \infty} \varphi_t^{-1} \ln \beta(\delta_t^{+, \alpha_t}) = -\gamma_0$$

if condition Λ^* holds with $\varepsilon_- < 0$ and $\Gamma_0 < \Gamma_1$.

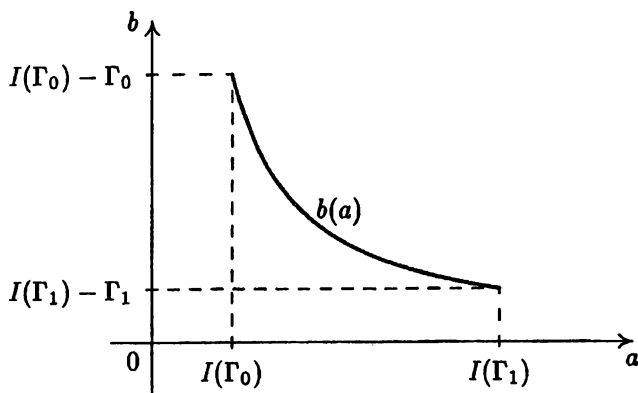


FIGURE 2.3.1

REMARK 2.3.12. The function $b(a)$ is shown in Figure 2.3.1. According to Remark 2.3.11 we have $I(\Gamma_0) = 0$ and $I(\Gamma_0) - \Gamma_0 = -\gamma_0$ if $\varepsilon_- < 0$. If $\varepsilon_- = 0$, then $\Gamma_0 = \gamma_- = \varkappa'(0+)$ and $I(\Gamma_0) - \Gamma_0 = -\varkappa(0+) - \varkappa'(0+)$. If $\varepsilon_+ > 1$, then $\Gamma_1 = \gamma_1$, $I(\Gamma_1) = \gamma_1$, and $I(\Gamma_1) - \Gamma_1 = 0$. If $\varepsilon_+ = 1$, then $\Gamma_1 = \gamma_+ = \varkappa'(1-)$ and $I(\Gamma_1) - \Gamma_1 = -\varkappa(1-)$.

EXAMPLE 2.3.3. This is a continuation of Example 2.3.2. Relation (2.3.61) implies condition Λ^* with $\varphi_n = v_n^2$, $\varkappa(\varepsilon) = -2^{-1}\varepsilon(1 - \varepsilon)$, $\varepsilon_- = -\infty$, and $\varepsilon_+ = \infty$. Hence $\gamma_- = -\infty$, $\gamma_+ = \infty$, $\gamma_0 = -2^{-1}$, and $\gamma_1 = 2^{-1}$. It is easy to show that $I(\gamma) = 2^{-1}(\gamma + 2^{-1})^2$. It is clear that $\Gamma_0 < 0 < \Gamma_1$. Therefore Theorems 2.3.6–2.3.8 hold. In particular, Theorems 2.3.6 and 2.3.7 imply that

$$\varepsilon_\pi(\delta_n^\pi) = \exp \left\{ -\frac{1}{8} v_n^2 (1 + o(1)) \right\}, \quad \varepsilon_\pi(\delta_n^*) = \exp \left\{ -\frac{1}{8} v_n^2 (1 + o(1)) \right\}.$$

Since $\gamma(a) = \sqrt{2a} - 2^{-1}$ is a solution of the equation $I(\gamma) = a$, we obtain

$$b(a) = 2^{-1}(\sqrt{2a} - 1)^2$$

for all

$$a \in (0, 2^{-1}) = (I(\Gamma_0), I(\Gamma_1)).$$

Therefore, in view of Theorem 2.3.8 we have that

$$\lim_{n \rightarrow \infty} v_n^{-2} \ln \alpha_n = -a \iff \lim_{n \rightarrow \infty} v_n^{-2} \ln \beta(\delta_n^{+, \alpha_n}) = -2^{-1} (\sqrt{2a} - 1)^2$$

for all $a \in (0, 2^{-1})$

EXAMPLE 2.3.4. Let an observation be

$$\xi^{(n)} = (\xi_{n1}, \xi_{n2}, \dots, \xi_{nn})$$

where $\xi_{n1}, \xi_{n2}, \dots, \xi_{nn}$ are independent both under H^n and \tilde{H}^n . Assume that the random variable ξ_{ni} has the exponential distribution with the density

$$\lambda_{ni} \exp(-\lambda_{ni}(x - b_{ni}))I(x \geq b_{ni})$$

under the hypothesis H^n , and the exponential distribution with the density

$$\tilde{\lambda}_{ni} \exp(-\tilde{\lambda}_{ni}(x - \tilde{b}_{ni}))I(x \geq \tilde{b}_{ni})$$

under the hypothesis \tilde{H}^n . Here $\lambda_{ni} \in (0, \infty)$, $\tilde{\lambda}_{ni} \in (0, \infty)$, $b_{ni} \in (-\infty, \infty)$, and $\tilde{b}_{ni} \in (-\infty, \infty)$. It is easy to show that

$$(2.3.82) \quad \varepsilon_-^n = - \min_{1 \leq i \leq n} \left[\left(\left(\frac{\tilde{\lambda}_{ni}}{\lambda_{ni}} \vee 1 \right) - 1 \right)^{-1} I(\tilde{b}_{ni} \leq b_{ni}) \right],$$

$$(2.3.83) \quad \varepsilon_+^n = \min_{1 \leq i \leq n} \left[1 - \left(\frac{\tilde{\lambda}_{ni}}{\lambda_{ni}} \wedge 1 \right) I(\tilde{b}_{ni} \geq b_{ni}) \right]^{-1},$$

and that for all $\varepsilon \in (\varepsilon_-^n, \varepsilon_+^n)$

$$(2.3.84) \quad \ln H_n(\varepsilon) = \sum_{i=1}^n \left\{ \varepsilon \ln \frac{\tilde{\lambda}_{ni}}{\lambda_{ni}} - \ln \left[\varepsilon \left(\frac{\tilde{\lambda}_{ni}}{\lambda_{ni}} - 1 \right) + 1 \right] \right\} \\ + \varepsilon \ln H_n(1) + (1 - \varepsilon) \ln H_n(0)$$

where

$$(2.3.85) \quad \ln H_n(1) = - \sum_{i=1}^n \tilde{\lambda}_{ni}(b_{ni} - \tilde{b}_{ni})I(b_{ni} > \tilde{b}_{ni}),$$

$$(2.3.86) \quad \ln H_n(0) = - \sum_{i=1}^n \lambda_{ni}(\tilde{b}_{ni} - b_{ni})I(\tilde{b}_{ni} > b_{ni}).$$

Relations (2.3.82)–(2.3.86) allow us to completely study the problem of the asymptotic distinguishability of the hypotheses H^n and \tilde{H}^n , although we consider some particular cases only.

EXAMPLE 2.3.5. This is a continuation of Example 2.3.4. Let $\tilde{b}_{ni} = b_{ni}$ for all $i = 1, 2, \dots, n$ and all $n = 1, 2, \dots$. Then it is clear that the measures \mathbf{P}^n and $\tilde{\mathbf{P}}^n$ are equivalent and therefore $\bar{\alpha}_n = \bar{\beta}_n = 1$ for all n . Assume that $\tilde{\lambda}_{ni} = \lambda_{ni} = \rho_n$ for all $i = 1, 2, \dots, n$ and all $n = 1, 2, \dots$, and let $\rho_n \rightarrow \rho$ as $n \rightarrow \infty$. Relations (2.3.82)–(2.3.86) imply that

$$\begin{aligned} \ln H_n(\varepsilon) &= n(\varepsilon \ln \rho_n - \ln[\varepsilon(\rho_n - 1) + 1]), & \varepsilon &\in (\varepsilon_-^n, \varepsilon_+^n), \\ \varepsilon_-^n &= -((\rho_n \vee 1) - 1)^{-1}, & \varepsilon_+^n &= (1 - (\rho_n \wedge 1))^{-1}. \end{aligned}$$

Now we find ρ for which condition Λ^* holds and evaluate the function $b(a)$ defined in Theorem 2.3.8.

If $\rho \in (0, \infty) \setminus \{1\}$, then condition Λ^* holds and

$$\begin{aligned} \varphi_n &= n, & \varepsilon_- &= -((\rho \vee 1) - 1)^{-1}, & \varepsilon_+ &= (1 - (\rho \wedge 1))^{-1}, \\ \varkappa(\varepsilon) &= \varepsilon \ln \rho - \ln(\varepsilon(\rho - 1) + 1), & \varepsilon &\in (\varepsilon_-, \varepsilon_+). \end{aligned}$$

In this case, $\gamma_- = \ln \rho$ and $\gamma_+ = \infty$ if $\rho \in (0, 1)$, while $\gamma_- = -\infty$ and $\gamma_+ = \ln \rho$ if $\rho \in (1, \infty)$. Moreover, we have for any ρ

$$\gamma_0 = -(\rho - 1 - \ln \rho), \quad \gamma_1 = \rho^{-1} - 1 - \ln \rho^{-1}.$$

It is easy to show that for all $\gamma \in (\gamma_-, \gamma_+)$ we have

$$I(\gamma) = z(\gamma) - 1 - \ln z(\gamma)$$

where

$$z(\gamma) = (\ln \rho - \gamma)/(\rho - 1).$$

Observe that the condition $\Gamma_0 < 0 < \Gamma_1$ (and therefore the condition $\Gamma_0 < \Gamma_1$) holds, since $\Gamma_0 = \gamma_0 < 0$ and $\Gamma_1 = \gamma_1 > 0$. Therefore Theorems 2.3.6–2.3.8 hold and in Theorems 2.3.6 and 2.3.7 we have that

$$I(0) = \frac{\ln \rho}{\rho - 1} - 1 - \ln \frac{\ln \rho}{\rho - 1}.$$

If $\rho \in (0, 1)$, then $\gamma > \gamma_0$ for $z(\gamma) > 1$. Let z_a be a solution of the equation $z - 1 - \ln z = a$ for $z > 1$ and $a \in (0, \gamma_1)$. Then $\gamma(a) = \ln \rho + (1 - \rho)z_a$ and therefore

$$b(a) = a - \ln \rho - (1 - \rho)z_a.$$

If $\rho \in (1, \infty)$, then $\gamma > \gamma_0$ for $z(\gamma) \in (0, 1)$. Let \tilde{z}_a be a solution of the equation $z - 1 - \ln z = a$ for $z \in (0, 1)$ and $a \in (0, \gamma_1)$. Then $\gamma(a) = \ln \rho - (\rho - 1)\tilde{z}_a$ and therefore

$$b(a) = a - \ln \rho + (\rho - 1)\tilde{z}_a.$$

Now let $\rho = 1$. Assume that $n(\rho_n - 1)^2 \rightarrow \infty$ as $n \rightarrow \infty$. Then condition Λ^* holds with

$$\varphi_n = n(\rho_n - 1)^2, \quad \varkappa(\varepsilon) = -\varepsilon(1 - \varepsilon)/2, \quad \varepsilon_- = -\infty, \quad \varepsilon_+ = \infty.$$

This implies that

$$\gamma_- = -\infty, \quad \gamma_0 = -\frac{1}{2}, \quad \gamma_1 = \frac{1}{2}, \quad \gamma_+ = \infty, \quad I(\gamma) = \frac{1}{2} \left(\gamma + \frac{1}{2} \right)^2.$$

Therefore we have $I(0) = 1/8$ in Theorems 2.3.6 and 2.3.7 and $b(a) = (1 - \sqrt{2a})^2/2$ in Theorem 2.3.8 (cf. Example 2.3.3).

If $\rho = 0$, then condition Λ^* holds with $\varphi_n = n \ln \rho_n^{-1}$, $\varepsilon_- = -\infty$, $\varepsilon_+ = 1$, and $\varkappa(\varepsilon) = -\varepsilon$. Since $\varkappa'(\varepsilon) = -1$ for all ε , we have $\gamma_- = \gamma_+ = -1$. Therefore the conditions $\Gamma_0 < \Gamma_1$ and $\Gamma_0 < 0 < \Gamma_1$ fail to hold and Theorems 2.3.6–2.3.8 do not apply. Since $\varepsilon_- < 0$, conditions Λ_1 , Λ_2 , and Λ_3 hold and therefore Theorems 2.3.1 and 2.3.3 apply.

If $\rho = \infty$, then condition Λ^* holds with $\varphi_n = n \ln \rho_n^{-1}$, $\varepsilon_- = 0$, $\varepsilon_+ = \infty$, and $\varkappa(\varepsilon) = \varepsilon - 1$. In this case, $\varkappa'(\varepsilon) = 1$ for all ε . Therefore $\gamma_- = \gamma_+ = 1$ and Theorems 2.3.6–2.3.8 do not apply. Observe that conditions Λ_1 , Λ_2 , and Λ_3 fail to hold in this case.

EXAMPLE 2.3.6. This is a continuation of Example 2.3.4. Let $\lambda_{ni} = \lambda$, $\tilde{\lambda}_{ni} = \tilde{\lambda}$, $b_{ni} = b$, and $\tilde{b}_{ni} = \tilde{b}$ for all $i = 1, 2, \dots, n$ and all $n = 1, 2, \dots$. Then condition Λ^* holds with $\varphi_n = n$,

$$\varepsilon_- = -((\rho \vee 1) - 1)^{-1} I(\tilde{b} \leq b), \quad \varepsilon_+ = -[1 - (\rho \wedge 1) I(\tilde{b} \geq b)]^{-1},$$

and

$$\varkappa(\varepsilon) = \varepsilon \ln \rho - \ln[\varepsilon(\rho - 1) + 1] + \varepsilon \varkappa(1) + (1 - \varepsilon) \varkappa(0)$$

for all $\varepsilon \in (\varepsilon_-, \varepsilon_+)$ where

$$\varkappa(1) = -\tilde{\lambda}(b - \tilde{b}) I(b > \tilde{b}), \quad \varkappa(0) = -\lambda(\tilde{b} - b) I(\tilde{b} > b), \quad \rho = \frac{\tilde{\lambda}}{\lambda}.$$

If $\rho = 1$ and $\tilde{b} \neq b$, then $\gamma_- = \gamma_+$ and therefore Theorems 2.3.6–2.3.8 do not apply. In this case conditions Λ_1 , Λ_2 , and Λ_3 hold and therefore Theorems 2.3.1 and 2.3.3 apply if $\tilde{b} < b$.

Let $\rho \neq 1$ and $\tilde{b} > b$. Then

$$\varepsilon_- = 0, \quad \varepsilon_+ = (1 - (\rho \wedge 1))^{-1}, \quad \varkappa(\varepsilon) = \varepsilon \ln \rho - \ln[\varepsilon(\rho - 1) + 1] - (1 - \varepsilon)(\tilde{b} - b)\lambda$$

in condition Λ^* . In this case, $\gamma_- = -(\rho - 1 - \ln \rho) + (\tilde{b} - b)\lambda$ for all $\rho \in (0, \infty) \setminus \{1\}$, while $\gamma_+ = \infty$ if $\rho \in (0, 1)$, and $\gamma_+ = \ln \rho + (\tilde{b} - b)\lambda$ if $\rho \in (1, \infty)$. Observe that $\Gamma_0 = \gamma_-$ and $\Gamma_1 = \gamma_1 = (\rho^{-1} - 1 - \ln \rho^{-1}) + (\tilde{b} - b)\lambda$. The condition $\Gamma_0 < \Gamma_1$ clearly holds and therefore Theorem 2.3.8 applies. The condition $\Gamma_0 < 0 < \Gamma_1$ holds if $(\tilde{b} - b)\lambda < \rho - 1 - \ln \rho$ and therefore Theorems 2.3.6 and 2.3.7 apply only in this case. It is easy to show that for all $\gamma \in (\gamma_-, \gamma_+)$

$$I(\gamma) = z(\gamma) - 1 - \ln z(\gamma) + (\tilde{b} - b)\lambda$$

where $z(\gamma) = (\rho - 1)^{-1}(\ln \rho + (\tilde{b} - b)\lambda - \gamma)$.

If $\rho \in (0, 1)$, then $\gamma > \gamma_-$ for $z(\gamma) > 1$. Let z_a be a solution of the equation $z - 1 - \ln z + (\tilde{b} - b)\lambda = a$ with respect to $z \in (1, \infty)$ where

$$a \in (I(\Gamma_0), I(\Gamma_1)) = ((\tilde{b} - b)\lambda, \gamma_1).$$

Then $\gamma(a) = \ln \rho + (1 - \rho)z_a + (\tilde{b} - b)\lambda$ and therefore

$$b(a) = a - \ln \rho - (1 - \rho)z_a - (\tilde{b} - b)\lambda.$$

If $\rho \in (1, \infty)$, then $\gamma > \gamma_-$ for $z(\gamma) \in (0, 1)$. Let \tilde{z}_a be a solution of the equation $z - 1 - \ln z + (\tilde{b} - b)\lambda = a$ with respect to $z \in (0, 1)$ where $a \in ((\tilde{b} - b)\lambda, \gamma_1)$. Then

$\gamma(a) = \ln \rho - (\rho - 1)\tilde{z}_a + (\tilde{b} - b)\lambda$ and therefore

$$b(a) = a - \ln \rho + (\rho - 1)\tilde{z}_a - (\tilde{b} - b)\lambda.$$

If $\rho \neq 1$ and $\tilde{b} < b$, then $\varepsilon_- = -((\rho \wedge 1) - 1)^{-1}$, $\varepsilon_+ = 1$, and

$$\varkappa(\varepsilon) = \varepsilon \ln \rho - \ln[\varepsilon(\rho - 1) + 1] - \varepsilon(b - \tilde{b})\lambda.$$

This case can be considered similarly to the case $\tilde{b} > b$.

2.4. Complete asymptotic distinguishability under the weak convergence

Consider completely asymptotically distinguishable families of hypotheses (H^t) and (\tilde{H}^t) in two cases:

- 1) the law of large numbers (LLN) for Λ_t does not hold;
- 2) the law of large numbers for Λ_t holds and, moreover, $(\Lambda_t - \varphi_t)/\psi_t$ converges weakly with some φ_t and ψ_t .

The case where the LLN does not hold. Let $\mathcal{L}(\eta_t | P^t)$ be the distribution of η_t with respect to the measure P^t and let the symbol \xrightarrow{w} stand for the weak convergence of probability distributions. Introduce the following condition:

- $\Lambda 4.$ $\mathcal{L}(\psi_t^{-1}\Lambda_t | P^t) \xrightarrow{w} L$ as $t \rightarrow \infty$ where ψ_t is a positive function such that $\psi_t \rightarrow \infty$ as $t \rightarrow \infty$ and L is a probability distribution on \mathbf{R} whose distribution function is $L(x)$.

LEMMA 2.4.1. *If condition $\Lambda 4$ holds, then $\bar{\alpha}_t \rightarrow 1$ as $t \rightarrow \infty$ and $L(x) = 1$ for any $x > 0$. Moreover,*

$$(2.4.1) \quad (H^t) \Delta (\tilde{H}^t) \iff L(0) = 1.$$

PROOF. For any $a \in (-\infty, 0)$, we have

$$P^t(\psi_t^{-1}\Lambda_t < a) \geq P^t(\Lambda_t = -\infty) = 1 - \bar{\alpha}_t.$$

If a is a point of continuity of the function $L(x)$, then $L(a) \geq \limsup_{t \rightarrow \infty} (1 - \bar{\alpha}_t)$ by condition $\Lambda 4$. Passing to the limit over the points of continuity $a \rightarrow -\infty$, we obtain $\bar{\alpha}_t \rightarrow 1$ as $t \rightarrow \infty$.

Choose $a > 0$ and apply the Chebyshev inequality to obtain

$$P^t(\psi_t^{-1}\Lambda_t > a) = P^t(z_t > e^{a\psi_t}) \leq e^{-a\psi_t} E^t z_t \leq e^{-a\psi_t}.$$

Passing to the limit as $t \rightarrow \infty$ in this inequality, we have $1 - L(a) = 0$ and therefore $L(a) = 1$ for any $a > 0$.

Further, in view of condition $\Lambda 4$, we have as $t \rightarrow \infty$

$$P^t(\Lambda_t > N) = P^t\left(\frac{\Lambda_t}{\psi_t} > \frac{N}{\psi_t}\right) \rightarrow \begin{cases} 1 - L(0), & N < 0, \\ 1 - L(0+), & N > 0. \end{cases}$$

Since $L(0+) = 1$, Theorem 2.2.1 gives the required equivalence (2.4.1). □

REMARK 2.4.1. Equivalence (2.4.1) means the following: If condition $\Lambda 4$ holds, then the complete asymptotic distinguishability $(H^t) \Delta (\tilde{H}^t)$ holds if and only if the distribution function $L(x)$ is continuous at the point $x = 0$.

The next result, which is an analog of Theorem 2.3.1 (see also Corollary 2.3.1), describes the behavior of the Neyman–Pearson test δ_t^{+, α_t} as $t \rightarrow \infty$ under condition $\Lambda 4$.

THEOREM 2.4.1. *If condition $\Lambda 4$ holds, then for any $\alpha \in (0, 1)$*

$$(2.4.2) \quad \lim_{t \rightarrow \infty} \alpha_t = \alpha \Rightarrow \limsup_{t \rightarrow \infty} \frac{d_t}{\psi_t} \leq \bar{l}_{1-\alpha} \Rightarrow \limsup_{t \rightarrow \infty} \frac{\ln \beta(\delta_t^{+, \alpha_t})}{\psi_t} \leq \bar{l}_{1-\alpha},$$

$$(2.4.3) \quad \lim_{t \rightarrow \infty} \alpha_t = \alpha \Rightarrow \liminf_{t \rightarrow \infty} \frac{\ln \beta(\delta_t^{+, \alpha_t})}{\psi_t} \geq \underline{l}_{1-\alpha} \Rightarrow \liminf_{t \rightarrow \infty} \frac{d_t}{\psi_t} \geq \underline{l}_{1-\alpha}$$

where

$$(2.4.4) \quad \bar{l}_p = \inf\{u: L(u) > p\}, \quad \underline{l}_p = \sup\{u: L(u) < p\}$$

for $p \in (0, 1)$.

PROOF. We have

$$(2.4.5) \quad \alpha_t = P^t(Y_t > y_t) + \varepsilon_t P^t(Y_t = y_t)$$

where $Y_t = \psi_t^{-1} \Lambda_t$ and $y_t = \psi_t^{-1} d_t$.

First we prove the first implication in (2.4.2). Suppose that $\alpha_t \rightarrow \alpha$, but

$$\limsup_{t \rightarrow \infty} y_t = \bar{y} > u$$

where u is a point of continuity of $L(x)$ such that $L(u) > 1 - \alpha$. Let t_n be a sequence such that $y_{t_n} \rightarrow \bar{y}$ as $n \rightarrow \infty$. Then

$$\alpha = \lim_{t \rightarrow \infty} \alpha_t \leq \limsup_{n \rightarrow \infty} P^{t_n}(Y_{t_n} \geq y_{t_n}) \leq \limsup_{n \rightarrow \infty} P^{t_n}(Y_{t_n} \geq u) = 1 - L(u) < \alpha.$$

This contradiction shows that $\limsup_{t \rightarrow \infty} y_t \leq u$ for any point u of continuity of $L(x)$. This proves the first implication in (2.4.2).

The second implications in (2.4.2) and (2.4.3) follow from estimate (2.3.10).

Assume again that $\alpha_t \rightarrow \alpha$ as $t \rightarrow \infty$ and let $u < 0$ be a point of continuity of the function $L(x)$ such that $L(u) < 1 - \alpha$. By equality (2.1.8), we obtain

$$\begin{aligned} \beta(\delta_t^{+, \alpha_t}) &\geq \tilde{E}^t I(Y_t \geq u) (1 - \delta_t^{+, \alpha_t}) \geq E^t I(Y_t \geq u) (1 - \delta_t^{+, \alpha_t}) z_t \\ &\geq e^{u\psi_t} E^t I(Y_t \geq u) (1 - \delta_t^{+, \alpha_t}) \geq e^{u\psi_t} (P^t(Y_t \geq u) - \alpha_t). \end{aligned}$$

This yields in view of condition $\Lambda 4$ that

$$\lim_{t \rightarrow \infty} \psi_t^{-1} \ln \beta(\delta_t^{+, \alpha_t}) \geq u,$$

whence the first implication in (2.4.3) follows, since u is arbitrary. \square

If the distribution function $L(x)$ in condition $\Lambda 4$ is continuous, then Theorem 2.4.1 can be made more precise. First we give a necessary definition.

DEFINITION 2.4.1. Let $p \in [0, 1]$. Let S be a probability distribution and let $S(x)$ be its distribution function. Any number $x_p \in \mathbf{R}^1$ such that

$$S(x_p) \leq p \leq S(x_p + 0)$$

is called a *quantile of order $p \in [0, 1]$* or, simply, a *p -quantile* of a probability distribution S . Such a number x_p is also called a *p -quantile* of the distribution function $S(x)$. If for some $p \in [0, 1]$ we have $S(x) > p$ for all $x \in \mathbf{R}^1$, then we put $x_p = -\infty$. If $S(x) < p$ for all $x \in \mathbf{R}^1$, then we put $x_p = +\infty$.

REMARK 2.4.2. If L is the probability distribution appearing in condition $\Lambda 4$, then it is clear that for $p \in (0, 1)$, the p -quantile l_p of the distribution L and the quantities \underline{l}_p and \bar{l}_p defined by relations (2.4.4) are related in the following way:

$$\underline{l}_p \leq l_p \leq \bar{l}_p.$$

THEOREM 2.4.2. Assume that condition $\Lambda 4$ holds and the function $L(x)$ is continuous and strictly increasing in the interval (\underline{l}, \bar{l}) where

$$(2.4.6) \quad \underline{l} = \sup\{x: L(x) = 0\}, \quad \bar{l} = \inf\{x: L(x) = 1\}.$$

(We agree that $\sup(\emptyset) = -\infty$ and $\inf(\emptyset) = \infty$.) Then for any $\alpha \in (0, 1)$

$$(2.4.7) \quad \lim_{t \rightarrow \infty} \alpha_t = \alpha \iff \lim_{t \rightarrow \infty} \frac{d_t}{\psi_t} = l_{1-\alpha} \iff \lim_{t \rightarrow \infty} \frac{\ln \beta(\delta_t^{+, \alpha_t})}{\psi_t} = l_{1-\alpha}$$

where $p \in (0, 1)$, $l_p = \underline{l}_p = \bar{l}_p$ is a p -quantile of the distribution L , and where the quantities \underline{l}_p and \bar{l}_p are defined by (2.4.4).

PROOF. Since the function $L(x)$ is continuous and strictly increasing in the interval $(\underline{l}_p, \bar{l}_p)$, we have $l_p = \underline{l}_p = \bar{l}_p$ for $p \in (0, 1)$. Put $L_t(x) = \mathbf{P}^t(Y_t < x)$ where $Y_t = \psi_t^{-1} \Lambda_t$. Then for any $\varepsilon > 0$ and $y_t \in \mathbf{R}^1$

$$(2.4.8) \quad \begin{aligned} \mathbf{P}^t(Y_t = y_t) &\leq L_t(y_t + \varepsilon) - L_t(y_t) \\ &= [L_t(y_t + \varepsilon) - L(y_t + \varepsilon)] \\ &\quad - [L_t(y_t) - L(y_t)] + [L(y_t + \varepsilon) - L(y_t)]. \end{aligned}$$

By the Pólya theorem (see, for example, [16]) we have

$$(2.4.9) \quad \lim_{t \rightarrow \infty} \sup_y |L_t(y) - L(y)| = 0.$$

Since the continuity of the function $L(x)$ implies its uniform continuity, we obtain from (2.4.8) and (2.4.9)

$$(2.4.10) \quad \lim_{t \rightarrow \infty} \mathbf{P}^t(Y_t = y_t) = 0.$$

Therefore it follows from (2.4.5) and (2.4.10) that

$$(2.4.11) \quad \lim_{t \rightarrow \infty} \alpha_t = \alpha \iff \lim_{t \rightarrow \infty} L_t(y_t) = 1 - \alpha,$$

whence

$$(2.4.12) \quad \lim_{t \rightarrow \infty} \alpha_t = \alpha \iff \lim_{t \rightarrow \infty} y_t = l_{1-\alpha},$$

since $L(x)$ is continuous and strictly increasing in the interval (\underline{l}, \bar{l}) .

By (2.1.9), we have for any $\varepsilon > 0$

$$(2.4.13) \quad \begin{aligned} \beta(\delta_t^{+, \alpha_t}) &\geq E^t z_t (1 - \delta_t^{+, \alpha_t}) \geq E^t I(y_t - \varepsilon \leq Y_t) z_t (1 - \delta_t^{+, \alpha_t}) \\ &\geq P^t(y_t - \varepsilon \leq Y_t < y_t) \exp((y_t - \varepsilon)\psi_t). \end{aligned}$$

Using (2.4.13) we get

$$(2.4.14) \quad \lim_{t \rightarrow \infty} y_t = l_{1-\alpha} \Rightarrow \lim_{t \rightarrow \infty} \psi_t^{-1} \ln \beta(\delta_t^{+, \alpha_t}) = l_{1-\alpha}$$

in view of the second implication in (2.4.2) and since ε is arbitrary. The relation

$$\lim_{t \rightarrow \infty} y_t = l_{1-\alpha} \Leftarrow \lim_{t \rightarrow \infty} \psi_t^{-1} \ln \beta(\delta_t^{+, \alpha_t}) = l_{1-\alpha}$$

follows from

$$(2.4.15) \quad \lim_{t \rightarrow \infty} \psi_t^{-1} \ln \beta(\delta_t^{+, \alpha_t}) = l_{1-\alpha} \Rightarrow \limsup_{t \rightarrow \infty} y_t \leq l_{1-\alpha}$$

in view of the second implication in (2.4.3).

We prove (2.4.15) by contradiction. Let

$$(2.4.16) \quad \lim_{t \rightarrow \infty} \psi_t^{-1} \ln \beta(\delta_t^{+, \alpha_t}) = l_{1-\alpha}, \quad \limsup_{t \rightarrow \infty} y_t > l_{1-\alpha},$$

and let (t_n) be a sequence such that $t_n \rightarrow \infty$ and $y_{t_n} \rightarrow \limsup_{t \rightarrow \infty} y_t$ as $n \rightarrow \infty$. Then estimate (2.4.13) implies that

$$\limsup_{t \rightarrow \infty} \psi_t^{-1} \ln \beta(\delta_t^{+, \alpha_t}) \geq \limsup_{n \rightarrow \infty} \psi_{t_n}^{-1} \ln \beta(\delta_{t_n}^{+, \alpha_{t_n}}) > l_{1-\alpha},$$

contradicting the equality in (2.4.16). This proves implication (2.4.15). By combining (2.4.12), (2.4.14), (2.4.15), and (2.4.3), we obtain (2.4.7). \square

REMARK 2.4.3. If the limit distribution L in condition $\Lambda 4$ is singular and concentrated at a point $x_0 \in (-\infty, 0)$, then the law of large numbers $\Lambda 1$ holds, and the statements of the previous section apply. If the function $L(x)$ in condition $\Lambda 4$ is discontinuous at a point $x_0 \in (-\infty, 0)$ such that $0 < L(x_0) < L(x_0 + 0) < 1$, then Theorem 2.4.1 implies the following result: *If $\alpha_t \rightarrow \alpha$ as $t \rightarrow \infty$, then we have for any $\alpha \in (1 - L(x_0 + 0), 1 - L(x_0))$*

$$(2.4.17) \quad \lim_{t \rightarrow \infty} \frac{d_t}{\psi_t} = l_{1-\alpha}, \quad \lim_{t \rightarrow \infty} \frac{\ln \beta(\delta_t^{+, \alpha_t})}{\psi_t} = l_{1-\alpha}$$

where $l_{1-\alpha} = x_0$. Moreover, if condition $\Lambda 4$ holds and $L^{-1}(\alpha) = \inf\{u: L(u) > \alpha\}$, $\alpha \in [0, 1]$, is the inverse function of $L(x)$, then Theorems 2.4.1 and 2.4.2 imply the following result: *If $\alpha_t \rightarrow \alpha$ as $t \rightarrow \infty$, then for any $\alpha \in (0, 1)$ that is a point of continuity of $L^{-1}(\alpha)$, relations (2.4.17) hold with $l_{1-\alpha} = L^{-1}(\alpha)$.*

REMARK 2.4.4. If the assumptions of Theorem 2.4.2 hold and δ_t^{+, α_t} is the Neyman–Pearson test whose level α_t has a limit value $\alpha \in (0, 1)$, then Theorem 2.4.2 implies that

$$d_t = l_{1-\alpha} \psi_t + o(\psi_t), \quad \ln \beta(\delta_t^{+, \alpha_t}) = l_{1-\alpha} \psi_t + o(\psi_t).$$

Here, $l_{1-\alpha}$ is strictly decreasing from 0 to $-\infty$ as α increases from 0 to 1. Therefore the rate of decay of $\beta(\delta_t^{+, \alpha_t})$ depends on α for each particular value of α . This differs

crucially from the case where the law of large numbers $\Lambda 1$ holds; in the latter case, the rate of decay of $\beta(\delta_t^{+, \alpha_t})$ is independent of α (see the preceding section).

Second-order behavior of the Neyman–Pearson test. Now we consider the case where the law of large numbers holds for Λ_t . In this case, the rate of decay of the type II error probability of the Neyman–Pearson test δ_t^{+, α_t} is independent of $\alpha = \lim_{t \rightarrow \infty} \alpha_t$ if $\alpha \in (0, 1)$. Below is a condition enabling us to evaluate the second term in the asymptotic expansion of $\ln \beta(\delta_t^{+, \alpha_t})$ as $t \rightarrow \infty$.

$\Lambda 5$. $\mathcal{L}(\varphi_t^{-1}(\Lambda_t + \psi_t) | P^t) \xrightarrow{w} L$ as $t \rightarrow \infty$ where φ_t and ψ_t are positive functions such that $\varphi_t \rightarrow \infty$ and $\psi_t \rightarrow \infty$ as $t \rightarrow \infty$, $\varphi_t = o(\psi_t)$, and L is a probability distribution on \mathbf{R}^1 whose distribution function is $L(x)$.

It is easy to see that $\Lambda 5 \Rightarrow \Lambda 1$ for $\chi_t = \psi_t$.

THEOREM 2.4.3. *If condition $\Lambda 5$ holds, then for any $\alpha \in (0, 1)$*

$$(2.4.18) \quad \begin{aligned} \lim_{t \rightarrow \infty} \alpha_t = \alpha &\Rightarrow \limsup_{t \rightarrow \infty} \frac{d_t + \psi_t}{\varphi_t} \leq \bar{l}_{1-\alpha} \\ &\Rightarrow \limsup_{t \rightarrow \infty} \frac{\ln \beta(\delta_t^{+, \alpha_t}) + \psi_t}{\varphi_t} \leq \bar{l}_{1-\alpha}, \end{aligned}$$

$$(2.4.19) \quad \begin{aligned} \lim_{t \rightarrow \infty} \alpha_t = \alpha &\Rightarrow \liminf_{t \rightarrow \infty} \frac{\ln \beta(\delta_t^{+, \alpha_t}) + \psi_t}{\varphi_t} \geq l_{1-\alpha} \\ &\Rightarrow \liminf_{t \rightarrow \infty} \frac{d_t + \psi_t}{\varphi_t} \geq l_{1-\alpha} \end{aligned}$$

where l_p and \bar{l}_p are defined by (2.4.4).

PROOF. It is sufficient to follow the proof of Theorem 2.4.1 with appropriate modifications. First of all we note that equality (2.4.5) holds with

$$(2.4.20) \quad Y_t = \varphi_t^{-1}(\Lambda_t + \psi_t), \quad y_t = \varphi_t^{-1}(d_t + \psi_t).$$

The proof of the first implication in (2.4.18) is the same as that of the first implication in (2.4.2). The second implications in (2.4.18) and (2.4.19) follow from the estimate (2.3.10).

Now let $\alpha_t \rightarrow \alpha$ as $t \rightarrow \infty$ and let $u \in \mathbf{R}^1$ be a point of continuity of the function $L(x)$ such that $L(u) < 1 - \alpha$. We have

$$\beta(\delta_t^{+, \alpha_t}) \geq E^t I(Y_t \geq u) (1 - \delta_t^{+, \alpha_t}) z_t \geq \exp(u\varphi_t - \psi_t) [P^t(Y_t \geq u) - \alpha_t]$$

by equality (2.1.8), whence the first implication in (2.4.19) follows. \square

The following result improves Theorem 2.4.3 for the case where the function $L(x)$ in condition $\Lambda 5$ is continuous.

THEOREM 2.4.4. *Assume that condition $\Lambda 5$ holds and the function $L(x)$ is continuous and strictly increasing on the interval (\underline{l}, \bar{l}) where \underline{l} and \bar{l} are defined by (2.4.6). Then for any $\alpha \in (0, 1)$*

$$(2.4.21) \quad \lim_{t \rightarrow \infty} \alpha_t = \alpha \iff \lim_{t \rightarrow \infty} \frac{d_t + \psi_t}{\varphi_t} = l_{1-\alpha} \iff \lim_{t \rightarrow \infty} \frac{\ln \beta(\delta_t^{+, \alpha_t}) + \psi_t}{\varphi_t} = l_{1-\alpha}$$

where l_p is a p -quantile of the distribution L .

PROOF. It suffices to repeat the proof of Theorem 2.4.2 with appropriate modifications. First, we should use equality (2.4.5) where Y_t and y_t are given by (2.4.20). By literally repeating the proof of equivalence (2.4.12), the first equivalence in (2.4.21) is proved.

Next, we have for any $\varepsilon > 0$

$$(2.4.22) \quad \begin{aligned} \beta(\delta_t^{+, \alpha_t}) &\geq E^t I(y_t - \varepsilon \leq Y_t) z_t (1 - \delta_t^{+, \alpha_t}) \\ &\geq P^t(y_t - \varepsilon \leq Y_t < y_t) \exp((y_t - \varepsilon)\varphi_t - \psi_t). \end{aligned}$$

Then the second implication in (2.4.18) implies

$$(2.4.23) \quad \lim_{t \rightarrow \infty} y_t = l_{1-\alpha} \Rightarrow \lim_{t \rightarrow \infty} \varphi_t^{-1} (\ln \beta(\delta_t^{+, \alpha_t}) + \psi_t) = l_{1-\alpha},$$

since ε is arbitrary in (2.4.22) and $\bar{l}_p = l_p$ for $p \in (0, 1)$ in view of the continuity and strict monotonicity of the function $L(x)$ in the interval (l, \bar{l}) . Further, the relation

$$\lim_{t \rightarrow \infty} y_t = l_{1-\alpha} \Rightarrow \lim_{t \rightarrow \infty} \varphi_t^{-1} (\ln \beta(\delta_t^{+, \alpha_t}) + \psi_t) = l_{1-\alpha}$$

follows from the second implication in (2.4.19) if

$$(2.4.24) \quad \lim_{t \rightarrow \infty} \varphi_t^{-1} (\ln \beta(\delta_t^{+, \alpha_t}) + \psi_t) = l_{1-\alpha} \Rightarrow \limsup_{t \rightarrow \infty} y_t \leq l_{1-\alpha}.$$

Relation (2.4.24) is readily proved by contradiction. Indeed, assume that

$$(2.4.25) \quad \lim_{t \rightarrow \infty} \varphi_t^{-1} (\ln \beta(\delta_t^{+, \alpha_t}) + \psi_t) = l_{1-\alpha}, \quad \limsup_{t \rightarrow \infty} y_t > l_{1-\alpha}$$

and let (t_n) be a sequence such that $t_n \rightarrow \infty$ and $y_{t_n} \rightarrow \limsup_{t \rightarrow \infty} y_t$ as $n \rightarrow \infty$. In view of estimate (2.4.22)

$$\limsup_{t \rightarrow \infty} \varphi_t^{-1} (\ln \beta(\delta_t^{+, \alpha_t}) + \psi_t) \geq \overline{\lim}_{n \rightarrow \infty} \varphi_{t_n}^{-1} (\ln \beta(\delta_{t_n}^{+, \alpha_{t_n}}) + \psi_{t_n}) > l_{1-\alpha},$$

contradicting the equality in (2.4.25). This proves implication (2.4.24), completing the proof of the second equivalence in (2.4.21). \square

REMARK 2.4.5. If the assumptions of Theorem 2.4.4 hold and if $\alpha_t \rightarrow \alpha \in (0, 1)$ as $t \rightarrow \infty$, then Theorem 2.4.4 gives the following asymptotic expansions for the Neyman–Pearson test δ_t^{+, α_t} :

$$(2.4.26) \quad d_t = -\psi_t + l_{1-\alpha}\varphi_t + o(\varphi_t),$$

$$(2.4.27) \quad \ln \beta(\delta_t^{+, \alpha_t}) = -\psi_t + l_{1-\alpha}\varphi_t + o(\varphi_t).$$

Expansions (2.4.26) and (2.4.27) show that, under these assumptions, the asymptotic behavior of the test δ_t^{+, α_t} depends, in the second term of the asymptotic expansion, on the limit value α of the level α_t (cf. Remark 2.4.4 and relation (2.3.12) in the case where the law of large numbers $\Lambda 1$ holds).

EXAMPLE 2.4.1. Let an observation ξ be the vector

$$\xi^n = (\xi_1, \xi_2, \dots, \xi_n), \quad n = 1, 2, \dots,$$

where the random variables $\xi_1, \xi_2, \dots, \xi_n$ form a first order autoregressive process

$$(2.4.28) \quad \xi_i = \theta \xi_{i-1} + w_i, \quad i = 1, 2, \dots, \xi_0 = 0,$$

where $\theta \in \mathbf{R}$ is an unknown parameter and w_1, w_2, \dots are independent Gaussian random variables with the $\mathcal{N}(0, 1)$ distribution (which are independent of θ). Denote by P_θ^n the probability measure generating the distribution of the observation ξ^n . Let the measures P_θ^n and $P_{\tilde{\theta}}^n$ correspond to the hypotheses H^n and \tilde{H}^n , respectively, where θ and $\tilde{\theta}$ are some points in \mathbf{R} such that $|\theta| > 1$ and $\tilde{\theta} \neq \theta$. We assume that θ is independent of n , while $\tilde{\theta}$ depends, generally speaking, on n . We write $\tilde{\theta} = \theta_n$ if $\tilde{\theta}$ depends on n . It is clear that the measures P_θ^n and $P_{\tilde{\theta}}^n$ generate Gaussian distributions. Moreover, $P_{\tilde{\theta}}^n \sim P_\theta^n$ and, in view of (2.4.28), the logarithm of the density of the measure $P_{\tilde{\theta}}^n$ with respect to the measure P_θ^n can be represented as follows (P_θ^n -a.s.):

$$(2.4.29) \quad \Lambda_n = (\tilde{\theta} - \theta) \sum_{i=1}^n \xi_{i-1} w_i - \frac{1}{2} (\tilde{\theta} - \theta)^2 \sum_{i=1}^n \xi_{i-1}^2.$$

By (2.4.28), we have

$$E_\theta^n (\theta^{-n-1} \xi_{n+1} - \theta^{-n} \xi_n)^2 = E_\theta^n (\theta^{-n-1} w_{n+1})^2 = \theta^{-2(n+1)}.$$

Using Proposition II.4.2 in [42] we obtain that $\theta^{-n} \xi_n$ is a Cauchy sequence with probability 1 and therefore the limit $\lim_{n \rightarrow \infty} \theta^{-n} \xi_n$ exists almost surely with respect to the probability P_θ^n . On the other hand (2.4.28) implies that the random variable ξ_i is normally distributed with mean 0 and variance

$$(2.4.30) \quad E_\theta^n \xi_i^2 = \theta^{2i} (1 + \theta^{-2} + \theta^{-4} + \dots + \theta^{-2i}).$$

Thus

$$(2.4.31) \quad \sqrt{\theta^2 - 1} \theta^{-n} \xi_n \rightarrow \eta, \quad n \rightarrow \infty,$$

almost surely with respect to the probability P_θ^n where the random variable η is normally $\mathcal{N}(0, 1)$ distributed. Further, we obtain by (2.4.31)

$$(2.4.32) \quad \begin{aligned} \theta^{-2n} (\theta^2 - 1)^2 \sum_{i=1}^n \xi_{i-1}^2 &= (\theta^2 - 1)^2 \sum_{i=1}^n (\theta^{i-n} \xi_{n-i})^2 \theta^{-2i} \\ &= (\theta^2 - 1)^2 \sum_{i=1}^{\infty} (\theta^{i-n} \xi_{n-i})^2 I_{[1, n]}(i) \theta^{-2i} \rightarrow \eta \end{aligned}$$

as $n \rightarrow \infty$ with probability 1. Here, we have used the dominated convergence, since for any ω where convergence (2.4.31) holds, there exists a constant $C(\omega)$ such that $|\theta^{-j} \xi_j(\omega)| \leq C(\omega)$ for all $j = 1, 2, \dots$. It is easy to show, in view of (2.4.28) and (2.4.30), that

$$(2.4.33) \quad E_\theta^n \left(\sum_{i=1}^n \xi_{i-1} w_i \right)^2 = \sum_{i=1}^n E_\theta^n \xi_{i-1}^2 \leq \sum_{i=1}^{2n} \theta^{2(i-1)} \frac{\theta^2}{\theta^2 - 1} \leq \frac{\theta^{2n}}{(\theta^2 - 1)^2}.$$

Now assume that $\tilde{\theta} = \theta_n$ depends, generally speaking, on n and $\theta^{2n}\Delta_n^2 \rightarrow \infty$ as $n \rightarrow \infty$ where $\Delta_n = \theta_n - \theta$. Then (2.4.29) implies the following representation (P_θ^n -a.s.):

$$(2.4.34) \quad \psi_n^{-1}\Lambda_n = 2(\theta^2 - 1)^2 \Delta_n^{-1} \theta^{-2n} \sum_{i=1}^n \xi_{i-1} w_i - (\theta^2 - 1)^2 \theta^{-2n} \sum_{i=1}^n \xi_{i-1}^2$$

where $\psi_n = 2^{-1}(\theta^2 - 1)^{-2} \theta^{2n} \Delta_n^2$. But, in view of estimate (2.4.33), we have

$$(2.4.35) \quad E_\theta^n \left(\theta^{-2n} \Delta_n^{-1} \sum_{i=1}^n \xi_{i-1} w_i \right)^2 \leq \Delta_n^{-2} \theta^{-2n} (\theta^2 - 1)^{-2} \rightarrow 0, \quad n \rightarrow \infty.$$

Relations (2.4.32), (2.4.34), and (2.4.35) imply that condition $\Lambda 4$ holds where $\psi_n = 2^{-1}(\theta^2 - 1)^{-2} \theta^{2n} \Delta_n^2$, L is the distribution of the random variable $-\eta^2$ whose distribution function is

$$L(x) = P(-\eta^2 < x) = 2(1 - \Phi(\sqrt{-x})), \quad x \leq 0,$$

$L(x) = 1$ for $x > 0$, and $\Phi(x)$ is the distribution function of the normal $\mathcal{N}(0, 1)$ distribution.

Since the function $L(x)$ is continuous, Lemma 2.4.1 implies that the complete asymptotic distinguishability $(H^n) \Delta (\tilde{H}^n)$ holds and, moreover, Theorem 2.4.2 applies. Observe that $\underline{l} = -\infty$, $\bar{l} = 0$, and for any $p \in (0, 1)$ a p -quantile of the distribution L can be represented in the form $l_p = -z_{1-p/2}^2$ where z_p is a p -quantile of the distribution $\mathcal{N}(0, 1)$. Therefore relations (2.4.7) become of the following form for any $\alpha \in (0, 1)$:

$$\lim_{n \rightarrow \infty} \alpha_n = \alpha \iff \lim_{n \rightarrow \infty} \frac{d_n}{\psi_n} = -z_{1+\frac{\alpha}{2}}^2 \iff \lim_{n \rightarrow \infty} \frac{\ln \beta(\delta_n^{+, \alpha_n})}{\psi_n} = -z_{1+\frac{\alpha}{2}}^2.$$

EXAMPLE 2.4.2. Let an observation be a sample $\xi^n = (\xi_1, \xi_2, \dots, \xi_n)$ where $\xi_1, \xi_2, \dots, \xi_n$ are independent identically distributed random variables. Assume that the distribution of ξ_i under the hypothesis H^n is generated by a measure P whose density with respect to some σ -finite measure μ is $p(x)$, while under the hypothesis \tilde{H}^n the distribution is given by a measure \tilde{P} having density $\tilde{p}(x)$ with respect to μ . Consider the likelihood ratios

$$z_n(x) = \prod_{i=1}^n z(x_i), \quad z(x_i) = \frac{p(x_i)}{\tilde{p}(x_i)}, \quad x = (x_1, x_2, \dots, x_n).$$

Put

$$\Lambda_n(x) = \ln z_n(x), \quad \lambda(x_i) = \ln z(x_i), \quad x = (x_1, x_2, \dots, x_n), \\ \Lambda_n = \Lambda_n(\xi^n), \quad \lambda_i = \ln z(\xi_i), \quad i = 1, 2, \dots, n.$$

Then

$$\Lambda_n = \sum_{i=1}^n \lambda_i, \quad n = 1, 2, \dots, \\ I(P^n | \tilde{P}^n) = nI(P | \tilde{P}), \quad I(P | \tilde{P}) = -E\lambda_1,$$

where P^n and \tilde{P}^n are the measures generating the distribution of ξ^n under the hypothesis H^n and \tilde{H}^n , respectively (see Section 2.3).

Assume that

$$0 < I(P|\tilde{P}) < \infty, \quad 0 < \sigma^2(P|\tilde{P}) = D\lambda_1 < \infty$$

where the symbol D stands for the variance under the hypothesis H^n . Then the central limit theorem [47] implies

$$\mathcal{L}\left(\sqrt{n}\sigma(P|\tilde{P})\left(\Lambda_n + nI(P|\tilde{P})\right) | P^n\right) \xrightarrow{w} \mathcal{N}(0, 1), \quad n \rightarrow \infty.$$

Therefore condition $\Lambda 5$ holds with

$$\varphi_n = \sqrt{n}\sigma(P|\tilde{P}), \quad \psi_n = nI(P|\tilde{P}), \quad L = \mathcal{N}(0, 1).$$

Thus the assumptions of Theorem 2.4.4 hold with $\bar{l} = -\infty$, $\bar{l} = \infty$, and $L(x) = \Phi(x)$. Hence relation (2.4.21) holds for any $\alpha \in (0, 1)$ where $l_p = z_p$ is a p -quantile of the distribution $\mathcal{N}(0, 1)$. Therefore the following asymptotic expansion holds for the Neyman–Pearson test δ_n^{+, α_n} if $\alpha_n \rightarrow \alpha \in (0, 1)$, $n \rightarrow \infty$:

$$\ln \beta(\delta_n^{+, \alpha_n}) = -nI(P|\tilde{P}) - \sqrt{n}\sigma(P|\tilde{P})z_{1-\alpha} + o(\sqrt{n})$$

(cf. relation (2.3.22)).

EXAMPLE 2.4.3. Let an observation be a column-vector $\xi^n = (\xi_1, \xi_2, \dots, \xi_n)'$ having Gaussian distributions P^n and \tilde{P}^n under the hypotheses H^n and \tilde{H}^n , respectively, namely

$$\mathcal{L}(\xi_n | H^n) = \mathcal{N}(\mu, \sigma^2 R_n), \quad \mathcal{L}(\xi_n | \tilde{H}^n) = \mathcal{N}(\tilde{\mu}, \tilde{\sigma}^2 R_n)$$

where $\mathcal{N}(a, B)$ is the Gaussian distribution with the vector of means a and matrix of the second order mixed moments B , $\mu = (\mu_1, \mu_2, \dots, \mu_n)'$, $\tilde{\mu} = (\tilde{\mu}_1, \tilde{\mu}_2, \dots, \tilde{\mu}_n)'$, and $R_n = (R_{ij})$ is an $n \times n$ -matrix. This model can be written as

$$\xi^n = \theta + v\zeta^n$$

where $\zeta^n = (\zeta_1, \zeta_2, \dots, \zeta_n)'$ is a Gaussian vector with distribution

$$\mathcal{L}(\zeta | H^n) = \mathcal{L}(\zeta | \tilde{H}^n) = \mathcal{N}(0, R_n)$$

and where $\theta = (\theta_1, \theta_2, \dots, \theta_n)'$ and $v \in (0, \infty)$ are unknown parameters such that $\theta = \mu$ and $v = \sigma$ under the hypothesis H^n , while $\theta = \tilde{\mu}$ and $v = \tilde{\sigma}$ under the hypothesis \tilde{H}^n .

Assume that the matrix R_n is nondegenerate for all $n = 1, 2, \dots$. Then the measures \tilde{P}^n and P^n are mutually absolutely continuous. Therefore the likelihood ratio z_n has the form $z_n(x) = d\tilde{P}^n/dP^n(x)$, $x = (x_1, x_2, \dots, x_n)$, where

$$\begin{aligned} \Lambda_n(x) &= \ln z_n(x) = n \ln \frac{\sigma}{\tilde{\sigma}} - \frac{1}{2\tilde{\sigma}^2}(x - \tilde{\mu})' R_n^{-1}(x - \tilde{\mu}) + \frac{1}{2\sigma^2}(x - \mu)' R_n^{-1}(x - \mu) \\ &= n \ln \frac{\sigma}{\tilde{\sigma}} + \frac{1}{2} \left(\frac{1}{\sigma^2} - \frac{1}{\tilde{\sigma}^2} \right) y' R_n^{-1} y - \frac{1}{\tilde{\sigma}^2} m' R_n^{-1} y - \frac{1}{2\tilde{\sigma}^2} m' R_n^{-1} m \end{aligned}$$

and $m = \tilde{\mu} - \mu$ and $y = x - \mu$. This implies that

$$(2.4.36) \quad \begin{aligned} \Lambda_n &= \Lambda_n(\xi^n) \\ &= n \ln \frac{\sigma}{\tilde{\sigma}} + \frac{1}{2} \left(\frac{1}{\sigma^2} - \frac{1}{\tilde{\sigma}^2} \right) \eta' R_n^{-1} \eta - \frac{1}{\tilde{\sigma}^2} m' R_n^{-1} \eta - \frac{1}{2\tilde{\sigma}^2} m' R_n^{-1} m \end{aligned}$$

where $\eta = \xi^n - \mu$. Note that Theorem 1.4.1 in [38] implies $E^n \eta' R_n^{-1} \eta = n\sigma^2$, whence we obtain

$$(2.4.37) \quad I(\mathbb{P}^n | \tilde{\mathbb{P}}^n) = -E^n \Lambda_n = \frac{n}{2} \left(\frac{\sigma^2}{\tilde{\sigma}^2} - 1 - \ln \frac{\sigma^2}{\tilde{\sigma}^2} \right) + \frac{1}{2\tilde{\sigma}^2} m' R_n^{-1} m$$

by equality (2.4.36).

Put $c_n = m' R_n^{-1} m$ and let $I(\mathbb{P}^n | \tilde{\mathbb{P}}^n) \rightarrow \infty$ as $n \rightarrow \infty$. Considering (2.4.37) we distinguish the following three cases.

a) Let $\tilde{\sigma} = \sigma$ and $c_n \rightarrow \infty$ as $n \rightarrow \infty$. Then (2.4.36) and (2.4.37) imply for any $n = 1, 2, \dots$ that

$$\mathcal{L} \left(c_n^{-1/2} \tilde{\sigma} \left[\Lambda_n + \frac{c_n}{2\tilde{\sigma}^2} \right] \mid \mathbb{P}^n \right) = \mathcal{N}(0, 1)$$

and thus condition A5 holds with

$$\varphi_n = \frac{c_n^{1/2}}{\tilde{\sigma}}, \quad \psi_n = \frac{c_n}{2\tilde{\sigma}^2}, \quad L = \mathcal{N}(0, 1), \quad L(x) = \Phi(x).$$

Therefore the assumptions of Theorem 2.4.4 are satisfied and thus

$$(2.4.38) \quad \ln \beta(\delta_n^{+, \alpha_n}) = -\frac{c_n}{2\tilde{\sigma}^2} - \frac{c_n^{1/2}}{\tilde{\sigma}} z_{1-\alpha} + o(c_n^{1/2})$$

if $\alpha_n \rightarrow \alpha \in (0, 1)$ as $n \rightarrow \infty$.

b) Let $\tilde{\sigma} \neq \sigma$ and $n/c_n \rightarrow \infty$ as $n \rightarrow \infty$. Equalities (2.4.36) and (2.4.37) imply that

$$(2.4.39) \quad \frac{1}{\sqrt{n}} \left[\Lambda_n + I(\mathbb{P}^n | \tilde{\mathbb{P}}^n) \right] = \frac{1}{\sqrt{2}} \left(\frac{\sigma^2}{\tilde{\sigma}^2} - 1 \right) G_n - H_n$$

where

$$(2.4.40) \quad H_n = \frac{1}{\sqrt{n}\tilde{\sigma}^2} m' R_n^{-1} \eta, \quad G_n = \frac{1}{\sqrt{2n}} \left(\frac{1}{\tilde{\sigma}^2} \eta' R_n^{-1} \eta - n \right).$$

It is clear that $H_n \rightarrow 0$ in probability as $n \rightarrow \infty$ under the hypothesis H^n . By the central limit theorem and in view of Theorem 1.4.1 in [38], we have

$$\mathcal{L}(G_n | \mathbb{P}^n) \xrightarrow{w} \mathcal{N}(0, 1) \quad \text{as } n \rightarrow \infty.$$

Therefore condition A5 holds with

$$\varphi_n = \sqrt{\frac{n}{2}} \left| \frac{\sigma^2}{\tilde{\sigma}^2} - 1 \right|, \quad \psi_n = I(\mathbb{P}^n | \tilde{\mathbb{P}}^n), \quad L = \mathcal{N}(0, 1).$$

Then Theorem 2.4.4 implies that

$$\ln \beta(\delta_n^{+, \alpha_n}) = -I(\mathbb{P}^n | \tilde{\mathbb{P}}^n) - \sqrt{\frac{n}{2}} \left| \frac{\sigma^2}{\tilde{\sigma}^2} - 1 \right| z_{1-\alpha} + o(\sqrt{n})$$

if $\alpha_n \rightarrow \alpha \in (0, 1)$ as $n \rightarrow \infty$. Observe that this expansion takes the following form in the case of $c_n = o(\sqrt{n})$:

$$\ln \beta(\delta_n^{+, \alpha_n}) = -\frac{n}{2} \left(\frac{\sigma^2}{\tilde{\sigma}^2} - 1 - \ln \frac{\sigma^2}{\tilde{\sigma}^2} \right) - \sqrt{\frac{n}{2}} \left| \frac{\sigma^2}{\tilde{\sigma}^2} - 1 \right| z_{1-\alpha} + o(\sqrt{n}).$$

c) Let $\tilde{\sigma} \neq \sigma$ and $n/c_n \rightarrow 0$ as $n \rightarrow \infty$. Then relations (2.4.39) and (2.4.40) imply that

$$\mathcal{L} \left(\frac{\tilde{\sigma}^2}{\sigma \sqrt{c_n}} \left[\Lambda_n + I(\mathbf{P}^n | \tilde{\mathbf{P}}^n) \right] \mid \mathbf{P}^n \right) \xrightarrow{w} \mathcal{N}(0, 1)$$

as $n \rightarrow \infty$. Therefore condition $\Lambda 5$ holds with

$$\varphi_n = \frac{\sigma}{\tilde{\sigma}^2} c_n^{1/2}, \quad \psi_n = I(\mathbf{P}^n | \tilde{\mathbf{P}}^n), \quad L = \mathcal{N}(0, 1).$$

Hence Theorem 2.4.4 implies that

$$\ln \beta(\delta_n^{+, \alpha_n}) = -I(\mathbf{P}^n | \tilde{\mathbf{P}}^n) - \frac{\sigma}{\tilde{\sigma}^2} c_n^{1/2} z_{1-\alpha} + o(c_n^{1/2}) \quad \text{if } \alpha_n \rightarrow \alpha \in (0, 1).$$

The latter expansion becomes of the following form if $n = o(c_n^{1/2})$:

$$\ln \beta(\delta_n^{+, \alpha_n}) = -\frac{1}{2\tilde{\sigma}^2} c_n - \frac{\sigma}{\tilde{\sigma}^2} c_n^{1/2} z_{1-\alpha} + o(c_n^{1/2})$$

(cf. (2.4.38)).

Examples where conditions $\Lambda 4$ and $\Lambda 5$ hold and the observation ξ^t is a stochastic process on the interval $[0, t]$ can be found in the monograph [37].

2.5. Contiguous families of hypotheses

Relative compactness and tightness of a family of probability measures. The concepts of relative compactness and tightness of families of probability measures play a fundamental role in studying contiguous families of statistical hypotheses. Assume that all underlying measures are defined on a metric space (E, \mathcal{E}, ρ) equipped with a distance ρ where the σ -algebra \mathcal{E} is generated by the metric ρ . In what follows, we often consider the case $(E, \mathcal{E}) = (\mathbf{R}^m, \mathcal{B}^m)$ with the Euclidean metric ρ .

DEFINITION 2.5.1. A family of probability measures $\mathcal{Q} = (\mathbf{Q}_u; u \in \mathfrak{A})$ is called *relatively compact* if any sequence of measures belonging to \mathcal{Q} contains a subsequence that converges weakly to a probability measure.

Note that this definition does not assume that the limit probability measure belongs to the family \mathcal{Q} .

DEFINITION 2.5.2. A family of probability measures $\mathcal{Q} = (\mathbf{Q}_u; u \in \mathfrak{A})$ is called *tight* if for any $\varepsilon > 0$ there exists a compact set $K_\varepsilon \subset E$ such that

$$\sup \{ \mathbf{Q}_u(E \setminus K_\varepsilon); u \in \mathfrak{A} \} \leq \varepsilon.$$

The following result is fundamental for the theory of weak convergence of probability measures.

THEOREM 2.5.1 (Prokhorov's theorem). *Assume that $\mathcal{Q} = (\mathbf{Q}_u; u \in \mathfrak{A})$ is a family of probability measures defined on a complete separable metric space (E, \mathcal{E}, ρ) . The family \mathcal{Q} is relatively compact if and only if \mathcal{Q} is tight.*

The proof of this theorem can be found in various textbooks on probability theory (see, for example, [3, 47]).

Consider a family of probability measures $\mathcal{Q} = (\mathbf{Q}_t; t \in \mathbf{R}_+)$ on the space $(\mathbf{R}, \mathcal{B})$ where \mathbf{Q}_t is the probability measure defining the likelihood ratio z_t with respect to the measure $\tilde{\mathbf{P}}^t$. If a family of hypotheses (\tilde{H}^t) is contiguous to a family (H^t) , then Theorem 2.2.3 implies that the family (z_t) is tight with respect to $(\tilde{\mathbf{P}}^t)$. Therefore the family of measures \mathcal{Q} is tight and thus, by the Prokhorov theorem, relatively compact. In a similar way, if a family (H^t) is contiguous to the family (\tilde{H}^t) , then the family of distributions of \tilde{z}_t with respect to the measure \mathbf{P}^t is relatively compact. It is clear that the family (z_t) is tight with respect to (\mathbf{P}^t) . Therefore the family of distributions of $\Lambda_t = \ln z_t$ is relatively compact with respect to \mathbf{P}^t , since $(H^t) \triangleleft (\tilde{H}^t)$. Therefore if $(H^t) \triangleleft (\tilde{H}^t)$, then every sequence of distributions of Λ_t with respect to \mathbf{P}^t has a weakly convergent subsequence. For the sake of brevity we assume throughout this section that the distribution of Λ_t with respect to \mathbf{P}^t is weakly convergent.

Weak convergence of the logarithm of the likelihood ratio. We introduce the following condition:

$\Lambda 6$. $\mathcal{L}(\Lambda_t | \mathbf{P}^t) \xrightarrow{w} L$ where L is a probability distribution on \mathbf{R} whose distribution function is denoted by $L(x)$.

THEOREM 2.5.2. *If condition $\Lambda 6$ holds, then*

$$(2.5.1) \quad \mathcal{L}(\Lambda_t | \tilde{\mathbf{P}}^t) \xrightarrow{w} \tilde{L}$$

where \tilde{L} is a probability distribution on the extended real line $\bar{\mathbf{R}} = [-\infty, \infty]$ whose distribution function is

$$(2.5.2) \quad \tilde{L}(x) = \int_{-\infty}^x e^y dL(y).$$

In this case, $\tilde{L}(\infty) \leq 1$ and

$$(2.5.3) \quad \lim_{N \rightarrow \infty} \lim_{t \rightarrow \infty} \tilde{\mathbf{P}}^t(\Lambda_t \geq N) = 1 - \tilde{L}(\infty).$$

PROOF. By the Lebesgue decomposition we have for any $x \in \mathbf{R}$

$$(2.5.4) \quad \tilde{\mathbf{P}}^t(\Lambda_t < x) = \mathbf{E}^t I(\ln z_t < x) z_t = \int_0^{\exp(x)} z dG_t(z)$$

where $G_t(z) = \mathbf{P}^t(z_t < z)$, $z \in \mathbf{R}_+$. By the Helly theorem [47], condition $\Lambda 6$ implies that

$$(2.5.5) \quad \lim_{t \rightarrow \infty} \int_0^x z dG_t(z) = \int_0^x z dG(z)$$

for any $x \in (0, \infty)$ such that $\ln x$ is a point of continuity of $L(x)$, where

$$G(z) = L(\ln z), \quad z \in (0, \infty).$$

It is clear that

$$(2.5.6) \quad \int_0^{\exp(x)} z dG(z) = \int_{-\infty}^y e^x dL(x).$$

It follows from (2.5.4)–(2.5.6) that for any $y \in \mathbf{R}$ which is a point of continuity of $L(x)$,

$$(2.5.7) \quad \lim_{t \rightarrow \infty} \tilde{P}^t(\Lambda_t < y) = \int_{-\infty}^y e^x dL(x),$$

whence relations (2.5.1) and (2.5.2) follow. The inequality $\tilde{L}(\infty) \leq 1$ and relation (2.5.3) also follow from (2.5.7). \square

COROLLARY 2.5.1. *We have*

$$(2.5.8) \quad \Lambda 6 \Rightarrow \lim_{t \rightarrow \infty} \bar{\alpha}_t = 1, \quad \liminf_{t \rightarrow \infty} \bar{\beta}_t \geq \tilde{L}(\infty).$$

PROOF. It is clear that for any $x \in \mathbf{R}$

$$1 - \bar{\alpha}_t = P^t(\Lambda_t = -\infty) \leq P^t(\Lambda_t < x), \quad \bar{\beta}_t = \tilde{P}^t(\Lambda_t < \infty) \geq \tilde{P}^t(\Lambda_t < x).$$

Since the number $x \in \mathbf{R}$ is arbitrary, we obtain (2.5.8) from condition $\Lambda 6$ and relation (2.5.3). \square

COROLLARY 2.5.2. *If condition $\Lambda 6$ holds for $L = \mathcal{N}(a, \sigma^2)$ where $a \in (-\infty, \infty)$ and $\sigma \in (0, \infty)$, then $a \leq -\sigma^2/2$ and (2.5.1) holds with the distribution function*

$$(2.5.9) \quad \tilde{L}(x) = h\Phi\left(\frac{x - a - \sigma^2}{\sigma}\right)$$

where $h = \exp(a + \sigma^2/2)$.

PROOF. It is clear that

$$\int_{-\infty}^x e^y dL(y) = \int_{-\infty}^x e^y d\Phi\left(\frac{y - a}{\sigma}\right) dy = h\Phi\left(\frac{x - a - \sigma^2}{\sigma}\right).$$

By relation (2.5.7) and by equality (2.5.2), this implies (2.5.1) and equality (2.5.9). Since $\tilde{\Lambda}(\infty) \leq 1$ by Theorem 2.5.1, equality (2.5.9) implies that $a \leq -\sigma^2/2$ in view of $\Phi(\infty) = 1$. \square

REMARK 2.5.1. If condition $\Lambda 6$ holds, then it follows from Theorem 2.5.2 that the limit distribution \tilde{L} is, in general, a mixture of two probability distributions \tilde{L}^+ and $\varepsilon_{\{\infty\}}$ weighted by $\tilde{L}(\infty)$ and $1 - \tilde{L}(\infty)$, respectively, that is,

$$\tilde{L} = \tilde{L}(\infty)\tilde{L}^+ + (1 - \tilde{L}(\infty))\varepsilon_{\{\infty\}}$$

where $\varepsilon_{\{\infty\}}$ is the Dirac measure concentrated at ∞ and \tilde{L}^+ is a distribution on \mathbf{R} determined by its distribution function $\tilde{L}^+(x) = \tilde{L}(x)/\tilde{L}(\infty)$, $x \in \mathbf{R}$. In particular, by Corollary 2.5.2

$$\tilde{L}(\infty) = h = \exp(a + \sigma^2/2), \quad \tilde{L}^+ = \mathcal{N}(a + \sigma^2, \sigma^2)$$

for $L = \mathcal{N}(a, \sigma^2)$.

REMARK 2.5.2. Assume that condition $\Lambda 6$ holds where the distribution L is a mixture of the normal distributions $\mathcal{N}(-\sigma^2/2, \sigma^2)$ with respect to the parameter σ with a probability distribution K on $(0, \infty)$. Then the distribution function $L(x)$ is continuous and increasing in $(-\infty, \infty)$. It is clear that the distribution \tilde{L} is a mixture of the normal distributions $\mathcal{N}(\sigma^2/2, \sigma^2)$ with respect to the parameter σ distributed according to the same distribution K . Therefore $\tilde{L}(\infty) = 1$.

THEOREM 2.5.3. *We have*

$$(2.5.10) \quad \Lambda 6 \Rightarrow (H^t) \triangleleft (\tilde{H}^t).$$

In particular, if condition $\Lambda 6$ holds, then

$$(2.5.11) \quad (\tilde{H}^t) \triangleleft (H^t) \iff \tilde{L}(\infty) = 1,$$

$$(2.5.12) \quad (\tilde{H}^t) \triangleleft\!\!\!\!\!\! \triangleright (H^t) \iff \tilde{L}(\infty) < 1.$$

PROOF. Condition $\Lambda 6$ implies that

$$(2.5.13) \quad \lim_{N \rightarrow -\infty} \lim_{t \rightarrow \infty} P^t(\Lambda_t < N) = 0.$$

By Theorem 2.2.3, with the hypotheses H^t and \tilde{H}^t interchanged, we obtain (2.5.10) in view of the equality

$$P^t(z_t < N) = P^t(\tilde{z}_t > N^{-1})$$

and relation (2.5.13).

Now we assume that condition $\Lambda 6$ is satisfied. Then $\tilde{L}(\infty) \leq 1$ by Theorem 2.2.3 and thus (2.5.3) holds. Hence Theorem 2.2.3 implies (2.5.11). Relation (2.5.12) follows from the inequality $\tilde{L}(\infty) \leq 1$, equivalence (2.5.11), and the following property: either $(\tilde{H}^t) \triangleleft (H^t)$ or $(\tilde{H}^t) \triangleleft\!\!\!\!\!\! \triangleright (H^t)$. \square

REMARK 2.5.3. The implication \Leftarrow in (2.5.12) is known as the *first Le Cam theorem* (see [22]).

REMARK 2.5.4. The following property follows from Theorem 2.5.3 under condition $\Lambda 6$: either $(\tilde{H}^t) \triangleleft\!\!\!\!\!\! \triangleright (H^t)$ (type a) or $(\tilde{H}^t) \triangleleft (H^t)$ (type b). Moreover

$$(\tilde{H}^t) \triangleleft\!\!\!\!\!\! \triangleright (H^t) \iff \tilde{L}(\infty) = 1,$$

$$(\tilde{H}^t) \triangleleft (H^t) \iff \tilde{L}(\infty) < 1.$$

If $L = \mathcal{N}(a, \sigma^2)$ in condition $\Lambda 6$, then $a \leq -\sigma^2/2$ by Corollary 2.5.2. Moreover,

$$\tilde{L}(\infty) = 1 \iff a = -\sigma^2/2,$$

$$\tilde{L}(\infty) < 1 \iff a < -\sigma^2/2.$$

Behavior of the Neyman–Pearson tests. The following result establishes a relationship between the behavior of the level α_t and that of the type II error probability $\beta(\delta_t^+, \alpha_t)$ for the Neyman–Pearson test under condition $\Lambda 6$.

THEOREM 2.5.4. *If condition $\Lambda 6$ holds, then for any $\alpha \in (0, 1)$:*

$$(2.5.14) \quad \lim_{t \rightarrow \infty} \alpha_t = \alpha \Rightarrow \limsup_{t \rightarrow \infty} \beta(\delta_t^{+, \alpha_t}) \leq \tilde{L}(\bar{l}_{1-\alpha} + 0),$$

$$(2.5.15) \quad \lim_{t \rightarrow \infty} \alpha_t = \alpha \Rightarrow \liminf_{t \rightarrow \infty} \beta(\delta_t^{+, \alpha_t}) \geq \tilde{L}(l_{1-\alpha})$$

where l_p and \bar{l}_p are defined by equalities (2.4.4).

PROOF. Assume that $\alpha_t \rightarrow \alpha \in (0, 1)$ as $t \rightarrow \infty$. First suppose that u is a point of continuity of the function $L(x)$ satisfying $L(u) > 1 - \alpha$. As in the proof of Theorem 2.4.1, we obtain that $\limsup_{t \rightarrow \infty} d_t \leq u$. Therefore

$$(2.5.16) \quad \limsup_{t \rightarrow \infty} \beta(\delta_t^{+, \alpha_t}) \leq \limsup_{t \rightarrow \infty} \tilde{P}^t(\Lambda_t \leq d_t) \leq \lim_{t \rightarrow \infty} \tilde{P}^t(\Lambda_t \leq u) = \tilde{L}(u),$$

since, by (2.5.2), the point u is also a point of continuity of $\tilde{L}(x)$.

Inequality (2.5.16) yields

$$\limsup_{t \rightarrow \infty} \beta(\delta_t^{+, \alpha_t}) \leq \inf\{\tilde{L}(u) : L(u) > 1 - \alpha\} = \tilde{L}(\bar{l}_{1-\alpha} + 0)$$

where the infimum is taken over all points u of continuity of $L(x)$ such that

$$L(u) > 1 - \alpha.$$

Therefore implication (2.5.14) is proved.

Now we suppose that u is a point of continuity of $L(x)$ such that $L(u) < 1 - \alpha$. Then $\liminf_{t \rightarrow \infty} d_t \geq u$, whence

$$\liminf_{t \rightarrow \infty} \beta(\delta_t^{+, \alpha_t}) \geq \liminf_{t \rightarrow \infty} \tilde{P}^t(\Lambda_t < d_t) \geq \lim_{t \rightarrow \infty} \tilde{P}^t(\Lambda_t < u) = \tilde{L}(u).$$

Therefore

$$\liminf_{t \rightarrow \infty} \beta(\delta_t^{+, \alpha_t}) \geq \sup\{\tilde{L}(u) : L(u) < 1 - \alpha\} = \tilde{L}(l_{1-\alpha})$$

where the supremum is taken over all points u of continuity of $L(x)$ such that $L(u) < 1 - \alpha$. Implication (2.5.15) is also proved. \square

If condition $\Lambda 6$ holds and the function $L(x)$ is continuous, then Theorem 2.5.4 can be sharpened. First we prove an auxiliary result which is also of interest on its own.

LEMMA 2.5.1. *Let $(Z^t, \mathcal{A}^t, \mathbf{S}^t)$, $t \in \mathbf{R}_+$, be a family of probability measures and let Y_t be a measurable mapping of the space (Z^t, \mathcal{A}^t) into the space $(\mathbf{R}, \mathcal{B})$ such that*

$$(2.5.17) \quad \mathcal{L}(Y_t | \mathbf{S}^t) \xrightarrow{w} S, \quad t \rightarrow \infty,$$

where S is a probability distribution on $\bar{\mathbf{R}}$ whose distribution function $S(x)$ is continuous for $x \in (-\infty, \infty)$ and such that $S(-\infty) = 0$ and $S(\infty) \leq 1$. Then

$$(2.5.18) \quad \lim_{t \rightarrow \infty} \mathbf{S}^t(Y_t = y_t) = 0$$

for any family (y_t) of numbers such that $y_t \in \mathbf{R}$ and $\limsup_{t \rightarrow \infty} y_t < \infty$ if $S(\infty) < 1$. Further assume that the function $S(x)$ is strictly increasing in the interval (x, \bar{x}) where

$$\underline{x} = \sup\{x : S(x) = 0\}, \quad \bar{x} = \inf\{x : S(x) = S(\infty)\}.$$

Let (y_t) and (ε_t) be arbitrary families of numbers such that $y_t \in \mathbf{R}$, $\varepsilon_t \in [0, 1]$, and the limit

$$(2.5.19) \quad \lim_{t \rightarrow \infty} [\mathbf{S}^t(Y_t > y_t) + \varepsilon_t \mathbf{S}^t(Y_t = y_t)] = \beta$$

exists. Then the limit $\lim_{t \rightarrow \infty} y_t$ exists for any $\beta \in (1 - S(\infty), 1)$ and

$$(2.5.20) \quad \lim_{t \rightarrow \infty} y_t = s_{1-\beta}$$

where s_p is a p -quantile of the distribution S . Moreover, if $\beta = 1$, then

$$(2.5.21) \quad \limsup_{t \rightarrow \infty} y_t \leq \underline{x},$$

while if $\beta = 1 - S(\infty)$, then

$$(2.5.22) \quad \liminf_{t \rightarrow \infty} y_t \geq \bar{x}.$$

PROOF. First assume that $S(\infty) = 1$ and prove relations (2.5.18) and (2.5.20). Put $S_t(y) = \mathbf{S}^t(Y_t < y)$. Then

$$\begin{aligned} \mathbf{S}^t(Y_t = y_t) &\leq S_t(y_t + \varepsilon) - S_t(y_t) \\ &= [S_t(y_t + \varepsilon) - S(y_t + \varepsilon)] - [S_t(y_t) - S(y_t)] + [S(y_t + \varepsilon) - S(y_t)] \end{aligned}$$

for any $\varepsilon > 0$. By the Pólya theorem (see [16]) and since the function $S(x)$ is uniformly continuous, the latter inequality and condition (2.5.17) imply (2.5.18).

It is clear that relation (2.5.18) implies that

$$(2.5.19) \iff \lim_{t \rightarrow \infty} S_t(y_t) = 1 - \beta.$$

Therefore, again by the Pólya theorem, we obtain

$$(2.5.23) \quad (2.5.19) \iff \lim_{t \rightarrow \infty} S(y_t) = 1 - \beta.$$

The properties of the function $S(x)$ and relation (2.5.23) imply that the limit $\lim_{t \rightarrow \infty} y_t$ exists and (2.5.20) holds.

Now we assume that $S(\infty) < 1$. It follows from condition (2.5.17) and the equality $S(-\infty) = 0$ that

$$(2.5.24) \quad \lim_{t \rightarrow \infty} S_t(x) = S(x)$$

uniformly in $x \leq N$ for any $N < \infty$. Now the proof of the required relation (2.5.18) follows the lines of the proof of the same relation in the case where $S(\infty) = 1$. The only difference is that the reference to the Pólya theorem is replaced with the uniform convergence (2.5.24) in the interval $(-\infty, \limsup_{t \rightarrow \infty} y_t + \varepsilon)$ where $\varepsilon > 0$.

It is clear that $\beta \in [1 - S(\infty), 1]$. First assume that $\beta \in (1 - S(\infty), 1)$. In order to prove that the limit $\lim_{t \rightarrow \infty} y_t$ exists and relation (2.5.20) holds in this case, note that

$$\lim_{N \rightarrow -\infty} \lim_{t \rightarrow \infty} S_t(N) = 0$$

by (2.5.17) and in view of $S(-\infty) = 0$. Therefore conditions (2.5.17) and (2.5.19) yield $\limsup_{t \rightarrow \infty} y_t < \infty$. Let us prove this fact by contradiction. Assume that $\limsup_{t \rightarrow \infty} y_t = \infty$. Then there exists a sequence (t_n) such that $t_n \rightarrow \infty$ and $y_{t_n} \rightarrow \infty$ as $n \rightarrow \infty$. Let $\varepsilon > 0$ be arbitrary and let $N_\varepsilon \in \mathbf{R}$ be such that $S(N_\varepsilon) > S(\infty) - \varepsilon$. By (2.5.24), there exists an integer $n_0 = n_0(\varepsilon)$ such that

$|S_{t_n}(N_\varepsilon) - S(N_\varepsilon)| < \varepsilon$ for all $n > n_0$. Now suppose that $n_1 = n_1(\varepsilon)$ is an integer and $y_{t_n} > N_\varepsilon$ for all $n > n_1$. Therefore

$$S_{t_n}(y_{t_n}) \geq S_{t_n}(N_\varepsilon) \geq S(N_\varepsilon) - \varepsilon \geq S(\infty) - 2\varepsilon$$

for all $n > n_0 \vee n_1$. Hence

$$\liminf_{t \rightarrow \infty} \mathbf{S}^t(Y_t \geq y_t) \leq \liminf_{n \rightarrow \infty} \mathbf{S}^{t_n}(Y_{t_n} \geq y_{t_n}) \leq 1 - S(\infty),$$

since ε is arbitrary. On the other hand,

$$\liminf_{t \rightarrow \infty} \mathbf{S}^t(Y_t \geq y_t) \geq \lim_{t \rightarrow \infty} [\mathbf{S}^t(Y_t > y_t) + \varepsilon_t \mathbf{S}^t(Y_t = y_t)] = \beta > 1 - S(\infty)$$

by condition (2.5.19). This contradiction proves that $\limsup_{t \rightarrow \infty} y_t < \infty$. By (2.5.18), we obtain $\mathbf{S}^t(Y_t = y_t) \rightarrow 0$ as $t \rightarrow \infty$. Now the proof of equality (2.5.20) follows the lines of the same proof in the case $S(\infty) = 1$.

Now we prove (2.5.21) and (2.5.22) assuming that $S(\infty) \leq 1$. Both proofs are carried out by contradiction.

Suppose that $\beta = 1$ but $\limsup_{t \rightarrow \infty} y_t > \underline{x}$. Then for any $N \in (\underline{x}, \bar{y})$, where $\bar{y} = \limsup_{t \rightarrow \infty} y_t$, there exists $t_0 = t_0(N) \in \mathbf{R}_+$ such that $y_t > N$ for all $t > t_0$. Therefore the inequality

$$\mathbf{S}^t(Y_t < y_t) + (1 - \varepsilon_t) \mathbf{S}^t(Y_t = y_t) \geq \mathbf{S}^t(Y_t < y_t) \geq \mathbf{S}^t(Y_t < N)$$

holds for $t > t_0$. By inequalities (2.5.17) and (2.5.19) and by the equality $S(\underline{x}) = 0$ we obtain

$$\lim_{t \rightarrow \infty} [\mathbf{S}^t(Y_t < y_t) + (1 - \varepsilon_t) \mathbf{S}^t(Y_t = y_t)] > 0,$$

since $S(x)$ is strictly monotone in the interval (\underline{x}, \bar{x}) . The latter inequality contradicts (2.5.19) for $\beta = 1$. Therefore inequality (2.5.21) is true.

Now assume that $\beta = 1 - S(\infty)$, but $\underline{y} = \liminf_{t \rightarrow \infty} y_t < \bar{x}$. Then there exists a sequence (t_n) such that $t_n \rightarrow \infty$ and $y_{t_n} \rightarrow \underline{y}$ as $n \rightarrow \infty$. By conditions (2.5.17) and (2.5.19) and by relation (2.5.18) we obtain

$$\lim_{n \rightarrow \infty} [\mathbf{S}^{t_n}(Y_{t_n} > y_{t_n}) + \varepsilon_{t_n} \mathbf{S}^{t_n}(Y_{t_n} = y_{t_n})] = 1 - S(\underline{y}) > 1 - S(\bar{x}) = 1 - S(\infty),$$

since $S(x)$ is strictly monotone in the interval (\underline{x}, \bar{x}) . This contradicts (2.5.19) for $\beta = 1 - S(\infty)$. Therefore (2.5.22) is true. \square

REMARK 2.5.5. If $\underline{x} = -\infty$ under the assumptions of Theorem 2.5.1, then the limit $\lim_{t \rightarrow \infty} y_t$ exists for $\beta = 1$ and equality (2.5.20) holds with $s_0 = \underline{x} = -\infty$. In a similar way, if $\bar{x} = \infty$, then the limit $\lim_{t \rightarrow \infty} y_t$ exists for $\beta = 1 - S(\infty)$ and equality (2.5.20) holds with $s_{1-\beta} = s_{S(\infty)} = \bar{x} = \infty$. In the case $S = \mathcal{N}(a, \sigma^2)$, see [45] for additional information if $a = -\sigma^2/2$, and [36] if $a < -\sigma^2/2$.

THEOREM 2.5.5. Assume that condition $\Lambda 6$ holds and the function $L(x)$ is continuous and strictly increasing in the interval (\underline{l}, \bar{l}) where

$$\underline{l} = \sup\{x: L(x) = 0\}, \quad \bar{l} = \inf\{x: L(x) = 1\}.$$

Then for any $\alpha \in [0, 1]$

$$(2.5.25) \quad \lim_{t \rightarrow \infty} \alpha_t = \alpha \iff \lim_{t \rightarrow \infty} \beta(\delta_t^{\dagger, \alpha_t}) = \tilde{L}(l_{1-\alpha})$$

where l_p is a p -quantile of the distribution L and $\tilde{L}(x)$ is the distribution function of the distribution \tilde{L} defined by equality (2.5.2).

PROOF. Condition Λ_6 and Theorem 2.5.1 imply weak convergence (2.5.1), where \tilde{L} is the probability distribution on $\bar{\mathbf{R}}$ whose distribution function $\tilde{L}(x)$ is given by (2.5.2) and $\tilde{L}(\infty) \leq 1$. Equality (2.5.2) shows that

$$\sup\{x: \tilde{L}(x) = 0\} = \underline{l}, \quad \inf\{x: \tilde{L}(x) = \tilde{L}(\infty)\} = \bar{l},$$

and that the function $\tilde{L}(x)$ is continuous on \mathbf{R} and strictly increasing on the interval (\underline{l}, \bar{l}) .

Assume that $\alpha_t \rightarrow \alpha$ as $t \rightarrow \infty$. If $0 < \alpha < 1$, then the limit $\lim_{t \rightarrow \infty} d_t = l_{1-\alpha}$ exists by Theorem 2.5.1. Therefore $\limsup_{t \rightarrow \infty} d_t < \infty$. Then, again by Theorem 2.5.1, we have $\tilde{P}^t(\Lambda_t = d_t) \rightarrow 0$ as $t \rightarrow \infty$. By (2.5.1) and (2.5.2), by the inequality $\limsup_{t \rightarrow \infty} d_t < \infty$, and since $\tilde{P}^t(\Lambda_t < y)$ converges to $\tilde{L}(y)$ uniformly in $y \leq N$ for any $N < \infty$, we obtain

$$(2.5.26) \quad \beta(\delta_t^{+, \alpha_t}) \rightarrow \tilde{L}(l_{1-\alpha}), \quad t \rightarrow \infty.$$

If $\alpha = 0$, then Lemma 2.5.1 implies that $\liminf_{t \rightarrow \infty} d_t \geq \bar{l}$. Therefore, for any $N \in (-\infty, \bar{l})$, there exists $t' = t'(N)$ such that $d_t > N$ for all $t > t'$. Hence

$$\beta(\delta_t^{+, \alpha_t}) \geq \tilde{P}^t(\Lambda_t < d_t) \geq \tilde{P}^t(\Lambda_t < N)$$

for all $t > t'$. Since $l_{1-\alpha} \geq \bar{l}$ for $\alpha = 0$ and N is arbitrary, convergence (2.5.1) implies convergence (2.5.26) for $\alpha = 0$.

If $\alpha = 1$, Lemma 2.5.1 implies $\limsup_{t \rightarrow \infty} d_t \leq \underline{l}$. Then for any $N \in (\underline{l}, \infty)$ there exists $t'' = t''(N)$ such that $d_t < N$ for any $t > t''$. Therefore

$$(2.5.27) \quad \beta(\delta_t^{+, \alpha_t}) \leq \tilde{P}^t(\Lambda_t < N) + \tilde{P}^t(\Lambda_t = d_t)$$

for all $t > t''$.

It is clear that $\limsup_{t \rightarrow \infty} d_t < \infty$ and $\tilde{P}^t(\Lambda_t = d_t) \rightarrow 0$ as $t \rightarrow \infty$ by Lemma 2.5.1. Therefore we obtain (2.5.26) for $\alpha = 1$ from estimate (2.5.27) and convergence (2.5.1). This completes the proof of the implication \Rightarrow in (2.5.25).

Now assume that $\beta(\delta_t^{+, \alpha_t}) \rightarrow \tilde{L}(l_{1-\alpha})$ as $t \rightarrow \infty$. Then $1 - \beta(\delta_t^{+, \alpha_t}) \rightarrow \beta$ where $\beta = 1 - \tilde{L}(l_{1-\alpha}) \in [1 - \tilde{L}(\infty), 1]$. For $\alpha \in (0, 1)$, we have $\beta \in (1 - \tilde{L}(\infty), 1)$ and, by Lemma 2.5.1, $d_t \rightarrow \tilde{l}_{1-\beta} = \tilde{l}_{\tilde{L}(l_{1-\alpha})} = l_{1-\alpha}$ where \tilde{l}_p is a p -quantile of the distribution \tilde{L} . By condition Λ_6 and by Lemma 2.5.1, we obtain $\alpha_t \rightarrow 1 - L(l_{1-\alpha}) = \alpha$ as $t \rightarrow \infty$.

If $\alpha = 0$, then $\beta = 1 - \tilde{L}(\infty)$. By Lemma 2.5.1, we have $\liminf_{t \rightarrow \infty} d_t \geq \bar{l}$. Then for any $N \in (-\infty, \bar{l})$ there exists $t_0 = t_0(N)$ such that $d_t > N$ for all $t > t_0$. Hence

$$\alpha_t \leq P^t(\Lambda_t > N) + \varepsilon_t P^t(\Lambda_t = d_t)$$

for all $t > t_0$. By condition Λ_6 and Lemma 2.5.1, we obtain $\alpha_t \rightarrow 0$ as $t \rightarrow \infty$, since N is arbitrary.

The proof of the relation $\alpha_t \rightarrow \alpha$ as $t \rightarrow \infty$ for $\alpha = 1$ is similar. Therefore the implication \Leftarrow in (2.5.25) is also proved. \square

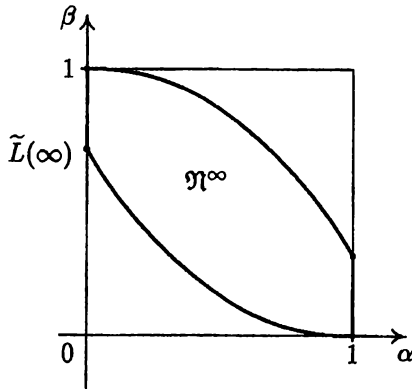


FIGURE 2.5.1

REMARK 2.5.6. Under the assumptions of Theorem 2.5.5, the function $\tilde{L}(1-\alpha)$ determines the equation for the lower bound of the limit \mathfrak{N}^∞ of the set \mathfrak{N}^t as $t \rightarrow \infty$. It is clear that this limit exists in our case. An example of the set \mathfrak{N}^∞ is shown in Figure 2.5.1.

Behavior of Bayes tests and minimax tests. Theorem 2.5.2 combined with Lemma 2.5.1 enables us to obtain results on the asymptotic behavior of the probabilities of error for Bayes tests and minimax tests.

Let δ_t^π be the Bayes test with respect to the a priori distribution $(\pi, \tilde{\pi})$, $\pi + \tilde{\pi} = 1$, and the loss function $A_{ij} = 1 - \delta_{ij}$.

THEOREM 2.5.6. *Assume that condition $\Lambda 6$ holds and $x = \ln(\pi/\tilde{\pi})$ is a point of continuity of the function $L(x)$. Then*

$$(2.5.28) \quad \lim_{t \rightarrow \infty} \alpha(\delta_t^\pi) = 1 - L\left(\ln \frac{\pi}{\tilde{\pi}}\right),$$

$$(2.5.29) \quad \lim_{t \rightarrow \infty} \beta(\delta_t^\pi) = \tilde{L}\left(\ln \frac{\pi}{\tilde{\pi}}\right),$$

$$(2.5.30) \quad \lim_{t \rightarrow \infty} e_\pi(\delta_t^\pi) = \pi \left(1 - L\left(\ln \frac{\pi}{\tilde{\pi}}\right)\right) + \tilde{\pi} \tilde{L}\left(\ln \frac{\pi}{\tilde{\pi}}\right)$$

where $\tilde{L}(x)$ is defined by (2.5.2).

PROOF. Since the Bayes test δ_t^π can be represented as $\delta_t^\pi = I(\Lambda_t \geq \ln(\pi/\tilde{\pi}))$, condition $\Lambda 6$ implies (2.5.28), while relation (2.5.29) follows from Theorem 2.5.2 and definition (2.5.2) of the function $\tilde{L}(x)$. Relation (2.5.30) follows from the equality

$$e_\pi(\delta_t^\pi) = \pi \alpha(\delta_t^\pi) + \tilde{\pi} \beta(\delta_t^\pi)$$

in view of (2.5.28) and (2.5.29). \square

Now assume that δ_t^* is the minimax test for testing the hypotheses H^t and \tilde{H}^t .

THEOREM 2.5.7. *Let condition $\Lambda 6$ hold where the function $L(x)$ is continuous and strictly increasing in the interval (l, \bar{l}) . Then*

$$(2.5.31) \quad \lim_{t \rightarrow \infty} \alpha(\delta_t^*) = \lim_{t \rightarrow \infty} \beta(\delta_t^*) = \lim_{t \rightarrow \infty} e(\delta_t^*) = \alpha^*$$

where α^* is a unique solution of the equation $\tilde{L}(l_{1-\alpha}) = \alpha$, $\tilde{L}(x)$ is defined by (2.5.2), and l_p is a p -quantile of the distribution L in condition $\Lambda 6$.

PROOF. Observe that Theorem 1.2.4 yields

$$(2.5.32) \quad \alpha(\delta_t^*) = \beta(\delta_t^*) = e(\delta_t^*).$$

Now, by Remark 2.5.6 and by Theorem 1.2.4 we obtain (2.5.31) from (2.5.32). The existence and uniqueness of a solution of the equation $\tilde{L}(l_{1-\alpha}) = \alpha$ follows, since $L(x)$ is continuous and strictly monotone in the interval (\underline{l}, \bar{l}) . \square

REMARK 2.5.7. It is clear that it is sufficient to find solution d^* of the equation $1 - L(d) = \tilde{L}(d)$ in order to find solution α^* of the equation $\tilde{L}(l_{1-\alpha}) = \alpha$. Then $\alpha^* = 1 - L(d^*) = \tilde{L}(d^*)$, that is, $d^* = l_{1-\alpha^*} = \tilde{l}_{\alpha^*}$. Further, by Theorem 1.2.4, there exists a Bayes test $\delta_t^{\pi_t^*}$ with respect to the a prior distribution $(\pi_t^*, 1 - \pi_t^*)$ and a loss function A_{ij} such that $\alpha(\delta_t^{\pi_t^*}) = \beta(\delta_t^{\pi_t^*})$. This implies that $\alpha(\delta_t^{\pi_t^*}) \rightarrow \alpha^*$. Therefore, by Lemma 2.5.1, $\pi_t^* \rightarrow \pi^*$ as $t \rightarrow \infty$ where $\pi^* \in (0, 1)$ is such that $\ln(\pi^*/(1 - \pi^*)) = d^*$.

Independent observations. Let an observation be the vector

$$\xi^{(n)} = (\xi_1, \xi_2, \dots, \xi_n)$$

where $\xi_1, \xi_2, \dots, \xi_n$ are independent identically distributed random variables with the distribution P_θ having density $p(x; \theta)$ with respect to some σ -finite measure μ . Here, θ is an unknown parameter taking values in an open set $\Theta \subset \mathbf{R}$. Then the distribution P_θ^n of the vector $\xi^{(n)}$ has density with respect to the measure μ^n and this density is of the form $p_n(x; \theta) = \prod_{i=1}^n p(x_i; \theta)$, $x = (x_1, x_2, \dots, x_n)$. For a fixed point $t \in \Theta$, introduce the following regularity conditions (R_t) on the family of probability distributions $\{P_\theta, \theta \in \Theta\}$:

- 1) the function $p(x; \theta)$ is absolutely continuous with respect to θ in some neighborhood of the point $\theta = t$ for all $x \in \mathbf{R}$;
- 2) the derivative $p'_\theta(x; \theta) = \partial p(x; \theta) / \partial \theta$ exists for any θ belonging to a neighborhood of the point $\theta = t$ for μ -almost all $x \in \mathbf{R}$;
- 3) the function $I(\theta) = E_\theta(\partial \ln p(\xi_1; \theta) / \partial \theta)^2$ is continuous and positive for $\theta = t$.

The function $I(\theta)$ is the *Fisher information* (see [25, 38]) and E_θ is expectation with respect to the measure P_θ .

Suppose that the hypothesis H^n is that the distribution of $\xi^{(n)}$ is determined by the measure P_t^n , while the hypothesis \tilde{H}^n is that the distribution of $\xi^{(n)}$ is determined by the measure $P_{t+u/\sqrt{n}}^n$ where u is a fixed number such that

$$t + u/\sqrt{n} \in \Theta.$$

Then the logarithm of the likelihood ratio is given by

$$(2.5.33) \quad \Lambda_n = \sum_{i=1}^n \ln (p(\xi_i; t + u/\sqrt{n}) / p(\xi_i; t)).$$

Consider the random variables

$$(2.5.34) \quad \eta_{in} = (p(\xi_i; t + u/\sqrt{n}) / p(\xi_i; t))^{1/2} - 1, \quad i = 1, 2, \dots, n,$$

and events $A_n = \{\max_{1 \leq i \leq n} |\eta_{in}| < \varepsilon\}$ where $\varepsilon > 0$ is a fixed (small) number. If the event A_n occurs, then we use (2.5.33) and (2.5.34) to expand the logarithm into the Taylor series and obtain

$$(2.5.35) \quad \Lambda_n = 2 \sum_{i=1}^n \ln(1 + \eta_{in}) = 2 \sum_{i=1}^n \eta_{in} - \sum_{i=1}^n \eta_{in}^2 + \sum_{i=1}^n \alpha_{in} |\eta_{in}|^3$$

where α_{in} are some numbers such that $|\alpha_{in}| < 1$.

Before studying the asymptotic behavior of Λ_n as $n \rightarrow \infty$, consider three auxiliary results.

LEMMA 2.5.2. *Let a nonnegative function $g(y)$ be absolutely continuous in the interval $[a, b]$ and, moreover,*

$$\int_a^b \frac{|g'(y)|}{\sqrt{g(y)}} dy < \infty.$$

Then the function $\sqrt{g(y)}$ is also absolutely continuous in the interval $[a, b]$.

PROOF. Assume that $g(y) > 0$ in an interval $(\alpha, \beta) \subset [a, b]$. Then it is clear that the function $\sqrt{g(y)}$ is absolutely continuous in the interval $[\alpha, \beta]$ and

$$\sqrt{g(\beta)} - \sqrt{g(\alpha)} = \int_{\alpha}^{\beta} \frac{g'(y)}{\sqrt{g(y)}} dy.$$

Given $c \in [a, b]$, consider the open set $\{y \in (a, c): g(y) > 0\}$. As is well known [41], this set can be represented as the union of an at most countable number of disjoint intervals (α_i, β_i) such that $g(\alpha_i) = g(\beta_i) = 0$ if $\alpha_i \neq a$ and $\beta_i \neq c$. Therefore

$$\int_a^c \frac{g'(y)}{\sqrt{g(y)}} dy = \sum_{i=1}^{\infty} \int_{\alpha_i}^{\beta_i} \frac{g'(y)}{\sqrt{g(y)}} dy = \sqrt{g(c)} - \sqrt{g(a)},$$

that is, the function $\sqrt{g(y)}$ is also absolutely continuous. □

Consider the random variables

$$\zeta(u) = \sqrt{\frac{p(\xi_1, t+u)}{p(\xi_1, t)}} - 1, \quad \varphi(\xi_1) = \frac{1}{2} \frac{\partial \ln(\xi_1; t)}{\partial t}$$

where u is a number such that $t+u \in \Theta$.

LEMMA 2.5.3. *If regularity conditions (R_t) hold, then*

$$(2.5.36) \quad E_t \left(\frac{\zeta(u)}{u} - \varphi(\xi_1) \right)^2 \rightarrow 0$$

as $u \rightarrow 0$.

PROOF. Since the function $I(\theta)$ is continuous in a neighborhood of the point $\theta = t$, we have by the Fubini theorem

$$\int_{-\infty}^{\infty} \left(\int_{t-\varepsilon}^{t+\varepsilon} \frac{(p'_\theta(x; \theta))^2}{p(x; \theta)} d\theta \right) d\mu = \int_{t-\varepsilon}^{t+\varepsilon} I(\theta) d\theta < \infty$$

where the internal integral on the left-hand side is μ -finite, whence it follows that the integral

$$\int_{t-\varepsilon}^{t+\varepsilon} \frac{|p'_\theta(x; \theta)|}{\sqrt{p(x; \theta)}} d\theta$$

also is μ -finite. Then we obtain in view of Lemma 2.5.2 that the function $\sqrt{p(x; \theta)}$ is absolutely continuous in a neighborhood of the point $\theta = t$ for μ -almost all x . Then we apply the Cauchy–Bunyakovskiĭ inequality to obtain

$$\begin{aligned} \text{E}_t \left(\frac{\zeta(u)}{u} \right)^2 &\leq \frac{1}{u^2} \int \left(\sqrt{p(x; t+u)} - \sqrt{p(x; t)} \right)^2 d\mu \\ (2.5.37) \quad &= \frac{1}{u^2} \int \left(\int_t^{t+u} \frac{p'_\theta(x; \theta)}{2\sqrt{p(x; \theta)}} d\theta \right)^2 d\mu \leq \frac{1}{4u} \int_t^{t+u} I(\theta) d\theta. \end{aligned}$$

Since the function $I(\theta)$ is continuous at $\theta = t$, we deduce from (2.5.37) that

$$(2.5.38) \quad \limsup_{u \rightarrow 0} \text{E}_t \left(\frac{\zeta|u|}{u} \right)^2 \leq \frac{I(t)}{4}.$$

It is clear that

$$(2.5.39) \quad \text{E}_t \varphi^2(\xi_1) = \frac{1}{4} I(t),$$

$$(2.5.40) \quad \frac{\zeta(u)}{u} \rightarrow \varphi(\xi_1), \quad u \rightarrow 0, \quad (\text{P}_t\text{-a.s.}).$$

Therefore the required relation (2.5.36) follows from Theorem 1.A.4 in [25]. □

LEMMA 2.5.4. *If regularity conditions (R_t) hold, then*

$$(2.5.41) \quad \text{E}_t \zeta^2(u) - \frac{1}{4} I(t) u^2 = o(u^2),$$

$$(2.5.42) \quad \text{E}_t \left(\zeta^2(u) - \left(\frac{\partial \sqrt{p(\xi_1; t)}}{\partial t} \right)^2 u^2 \right) = o(u^2),$$

$$(2.5.43) \quad \text{P}_t \{ |\zeta(u)| > \varepsilon \} = o(u^2),$$

$$(2.5.44) \quad \text{E}_t \zeta(u) + \frac{1}{8} I(t) u^2 = o(u^2)$$

as $u \rightarrow 0$.

PROOF. By Lemma 2.5.3, relation (2.5.36) holds and moreover it can be rewritten as

$$(2.5.45) \quad \text{E}_t (\zeta(u) - \varphi(\xi_1) u)^2 = o(u^2), \quad u \rightarrow 0,$$

whence (2.5.42) follows. Further, equality (2.5.41) follows from (2.5.39) and (2.5.40) in view of Theorem 5.3 in [2] and inequality (2.5.38). Using (2.5.41), we obtain from (2.5.37)

$$(2.5.46) \quad \lim_{u \rightarrow 0} \frac{1}{u^2} \int \left(\sqrt{p(x; t+u)} - \sqrt{p(x; t)} \right)^2 d\mu = \frac{1}{4} I(t).$$

We split the domain of integration in (2.5.46) into the sets $\{x: p(x; t) = 0\}$ and $\{x: p(x; t) \neq 0\}$ obtaining

$$(2.5.47) \quad \int_{\{x:p(x;t)=0\}} p(x; t+u) d\mu = o(u^2), \quad u \rightarrow 0,$$

in view of (2.5.41). Further we use (2.5.47) and obtain

$$\begin{aligned} E_t \zeta^2(u) &= \int_{\{x:p(x;t) \neq 0\}} \left(\sqrt{p(x; t+u)} - \sqrt{p(x; t)} \right)^2 d\mu \\ &= 2 + o(u^2) - 2E_t(p(\xi_1; t+u)/p(\xi_1; t))^{1/2} = -2E_t \zeta(u) + o(u^2). \end{aligned}$$

This implies (2.5.44) in view of (2.5.41).

Now we prove (2.5.43). It is clear that

$$\begin{aligned} P_t \{ |\zeta(u)| > \varepsilon \} &\leq P_t \left\{ \left| \zeta(u) - \frac{1}{2} \frac{\partial \ln p(\xi_1; t)}{\partial t} u \right| > \frac{\varepsilon}{2} \right\} + P_t \left\{ \left(\frac{\partial \ln p(\xi_1; t)}{\partial t} u \right)^2 > \varepsilon^2 \right\} \\ &\leq \frac{4}{\varepsilon^2} E_t \left(\zeta(u) - \frac{u}{2} \frac{\partial \ln p(\xi_1; t)}{\partial t} \right)^2 \\ &\quad + \frac{u^2}{\varepsilon^2} \int I \left(x: \left| \frac{\partial \ln p(x; t)}{\partial t} \right| > \frac{\varepsilon}{|u|} \right) \left(\frac{\partial \ln p(x; t)}{\partial t} \right)^2 d\mu. \end{aligned}$$

The first term on the right-hand side is $o(u^2)$ by (2.5.45) and the second term is also $o(u^2)$, since $I(t)$ is finite. □

REMARK 2.5.8. It is clear that (2.5.45) implies

$$E_t \zeta(u) - \frac{u}{2} E_t \frac{\partial \ln p(\xi_1; t)}{\partial t} = o(u), \quad u \rightarrow 0.$$

By (2.5.44), we obtain then

$$(2.5.48) \quad E_t \frac{\partial}{\partial t} \ln p(\xi_1; t) = \int_{\{x:p(x;t) \neq 0\}} \frac{\partial p(x; t)}{\partial t} \mu(dx) = 0.$$

The following result gives an asymptotic expansion of Λ_n as $n \rightarrow \infty$ under regularity conditions (R_t) .

THEOREM 2.5.8. *If regularity conditions (R_t) hold, then*

$$(2.5.49) \quad \Lambda_n = \frac{u}{\sqrt{n}} \sum_{i=1}^n \frac{\partial \ln p(\xi_i; t)}{\partial t} - \frac{u^2}{2} I(t) + \psi_n(u, t)$$

where

$$(2.5.50) \quad P_t^n \{ |\psi_n(u, t)| > \varepsilon \} \rightarrow 0, \quad n \rightarrow \infty,$$

for any $\varepsilon > 0$ and all $u \in \mathbf{R}$. We also have

$$(2.5.51) \quad \mathcal{L} \left(\frac{1}{\sqrt{n}} \sum_{i=1}^n \frac{\partial \ln p(\xi_i; t)}{\partial t} \mid P_t^n \right) \xrightarrow{w} \mathcal{N}(0, I(t)), \quad n \rightarrow \infty.$$

PROOF. The Taylor expansion (2.5.35) is valid if the event A_n occurs. In this case (2.5.35) implies representation (2.5.49)–(2.5.51) if for all $\varepsilon > 0$

$$(2.5.52) \quad P_t^n(A_n^C) \rightarrow 0,$$

$$(2.5.53) \quad P_t^n \left\{ \left| 2 \sum_{i=1}^n \eta_{in} - \frac{u}{\sqrt{n}} \sum_{i=1}^n \frac{\partial \ln p(\xi_i; t)}{\partial t} + \frac{u^2}{4} I(t) \right| > \varepsilon \right\} \rightarrow 0,$$

$$(2.5.54) \quad P_t^n \left\{ \left| \sum_{i=1}^n \eta_{in}^2 - \frac{u^2}{4} I(t) \right| > \varepsilon \right\} \rightarrow 0,$$

$$(2.5.55) \quad P_t^n \left\{ \sum_{i=1}^n |\eta_{in}|^3 > \varepsilon \right\} \rightarrow 0$$

as $n \rightarrow \infty$.

First we prove (2.5.52). Since the random variables $\eta_{1n}, \dots, \eta_{nn}$ are identically distributed and $\eta_{1n} = \zeta(u/\sqrt{n})$, relation (2.5.43) implies that

$$P_t^n(A_n^C) \leq \sum_{i=1}^n P_t\{|\eta_{in}| > \varepsilon\} = n P_t \left\{ \left| \zeta \left(\frac{u}{\sqrt{n}} \right) \right| > \varepsilon \right\} = o(1)$$

as $n \rightarrow \infty$. Further,

$$\begin{aligned} P_t^n \left\{ \left| \sum_{i=1}^n \eta_{in}^2 - \frac{u^2}{n} \sum_{i=1}^n \left(\frac{\partial \sqrt{p(\xi_i; t)}}{\partial t} \right)^2 \right| > \varepsilon \right\} &\leq \frac{1}{\varepsilon} E_t^n \sum_{i=1}^n \left| \eta_{in}^2 - \frac{u^2}{n} \left(\frac{\partial \sqrt{p(\xi_i; t)}}{\partial t} \right)^2 \right| \\ &\leq \frac{n}{\varepsilon} E_t \left| \zeta^2 \left(\frac{u}{\sqrt{n}} \right) - \frac{u^2}{n} \left(\frac{\partial \sqrt{p(\xi_1; t)}}{\partial t} \right)^2 \right| \\ &\rightarrow 0 \end{aligned}$$

by (2.5.42) where E_t^n is mathematical expectation with respect to the measure P_t^n . Moreover, since the sum

$$\sum_{i=1}^n \left(\partial \sqrt{p(\xi_i; t)} / \partial t \right)^2$$

converges in probability P_t^n to $I(t)/4$ by the law of large numbers, relation (2.5.54) is also proved.

Relation (2.5.55) follows from (2.5.52) and (2.5.54), since

$$\begin{aligned} P_t^n \left\{ \sum_{i=1}^n |\eta_{in}|^3 > \varepsilon \right\} &\leq P_t^n \left\{ \max_{1 \leq i \leq n} |\eta_{in}| > \frac{\varepsilon}{1 + I(t)u^2} \right\} \\ &\quad + P_t^n \left\{ \sum_{i=1}^n \eta_{in}^2 > 1 + I(t)u^2 \right\}. \end{aligned}$$

It remains to prove (2.5.53). Applying (2.5.44), we get

$$E_t \eta_{mi} = E_t \zeta \left(\frac{u}{\sqrt{n}} \right) = -\frac{I(t)u^2}{8n} + o\left(\frac{1}{n}\right)$$

as $n \rightarrow \infty$. Therefore (2.5.53) is equivalent to

$$J_n = P_t^n \left\{ 2 \left| \sum_{i=1}^n \left(\eta_{in} - E_t \eta_{in} - \frac{u}{2\sqrt{n}} \frac{\partial \ln p(\xi_i; t)}{\partial t} \right) \right| > \varepsilon \right\} \rightarrow 0.$$

Since ξ_1, \dots, ξ_n are independent, we obtain from (2.5.48) that

$$\begin{aligned} J_n &\leq \frac{4}{\varepsilon^2} E_t^n \left[\sum_{i=1}^n \left(\eta_{in} - E_t \eta_{in} - \frac{u}{2\sqrt{n}} \frac{\partial \ln p(\xi_i; t)}{\partial t} \right) \right]^2 \\ &= \frac{4n}{\varepsilon^2} E_t \left(\zeta \left(\frac{u}{\sqrt{n}} \right) - E_t \zeta \left(\frac{u}{\sqrt{n}} \right) - \frac{u}{2\sqrt{n}} \frac{\partial \ln p(\xi_1; t)}{\partial t} \right)^2 \\ &= \frac{4n}{\varepsilon^2} \left[E_t \left(\zeta \left(\frac{u}{\sqrt{n}} \right) - \frac{u}{2\sqrt{n}} \frac{\partial \ln p(\xi_1; t)}{\partial t} \right)^2 - \left(E_t \zeta \left(\frac{u}{\sqrt{n}} \right) \right)^2 \right]. \end{aligned}$$

Since the right-hand side of the latter inequality tends to zero as $n \rightarrow \infty$ by (2.5.36) and (2.5.44), relation (2.5.53) as well as representation (2.5.49)–(2.5.50) is proved.

Relation (2.5.51) follows from equality (2.5.48), since the Fisher information $I(t)$ is finite in view of the central limit theorem for sums of independent identically distributed random variables. \square

COROLLARY 2.5.3. *If regularity conditions (R_t) hold, then*

$$(2.5.56) \quad \mathcal{L}(\Lambda_n | P_t^n) \xrightarrow{w} \mathcal{N} \left(-\frac{1}{2} I(t) u^2, I(t) u^2 \right), \quad n \rightarrow \infty,$$

that is, condition $\Lambda 6$ is satisfied with $L = \mathcal{N} \left(-\frac{1}{2} I(t) u^2, I(t) u^2 \right)$.

REMARK 2.5.9. Asymptotic representations like (2.5.49)–(2.5.51) are known as the *local asymptotic normality* (LAN) property of a family of probability measures $\{P_\theta^n, \theta \in \Theta\}$ at the point $\theta = t$ as $n \rightarrow \infty$. Representation (2.5.49)–(2.5.51) for the observation $\xi^{(n)} = (\xi_1, \xi_2, \dots, \xi_n)$ can be found in the monographs [25, 37] for the case of independent but not necessarily identically distributed random variables $\xi_1, \xi_2, \dots, \xi_n$.

REMARK 2.5.10. Relation (2.5.56) and Corollary 2.5.2 imply that convergence (2.5.1) holds where $\tilde{L}(x)$ is the distribution function of the law $\mathcal{N} \left(\frac{1}{2} I(t) u^2, I(t) u^2 \right)$. According to Remark 2.5.4, the mutual contiguity $(H^n) \triangleleft \triangleright (\tilde{H}^n)$ holds in this case.

Goodness-of-Fit Tests

3.1. The setting of the problem. Kolmogorov test

Main definitions. Throughout this chapter we assume that the observation is a sample $\xi^{(n)} = (\xi_1, \xi_2, \dots, \xi_n)$ where $\xi_1, \xi_2, \dots, \xi_n$ are independent identically distributed random variables. Thus the random variables $\xi_1, \xi_2, \dots, \xi_n$ are independent observations of a random variable ξ . Consider a hypothesis H about the distribution of the random variable ξ . We call H the *main* or *null hypothesis*. Our aim is to test the hypothesis H , that is, either accept H or reject it. The decision is to be made on the basis of the information contained in the sample $\xi^{(n)}$. The alternative hypothesis is that the null hypothesis H is false. The alternative hypothesis is denoted by K . The null hypothesis H can be either simple or composite. Therefore our aim is to decide whether results of the observation $\xi^{(n)}$ are in agreement with the hypothesis H . The tests described above are called *goodness-of-fit tests*.

We follow the general procedure to construct a goodness-of-fit test. Namely we introduce a statistic $T = T(\xi^{(n)})$ treated as the measure of disagreement between the data $\xi^{(n)}$ and the hypothesis H . We require that, if H is true, the distribution of this statistic is known exactly or at least approximately. In particular, if the hypothesis H is composite, then the distribution of the statistic $T(\xi^{(n)})$ should be the same for all simple hypotheses forming H . If we treat $T(x)$ as a mapping of $(\mathbf{R}^n, \mathcal{B}^n)$ into a measurable space (Y, \mathcal{S}) , then the probability of the event $\{T(\xi^{(n)}) \in B\}$, $B \in \mathcal{S}$, is well defined if H is true. We write in this case $P\{T(\xi^{(n)}) \in B/H\}$.

Consider a set $Y_\alpha \in \mathcal{S}$ of large deviations of the hypothesis H from the data $\xi^{(n)}$ such that $P\{T(\xi^{(n)}) \in Y_\alpha/H\} \leq \alpha$ where $\alpha > 0$ is a sufficiently small number. Then the goodness-of-fit test can be described as follows. If $T(\xi^{(n)}) \in Y_\alpha$, then, under the assumption that the hypothesis H is true, an event occurs whose probability is small and thus the hypothesis H should be rejected, since it contradicts the observation $\xi^{(n)}$. Otherwise, that is, if $T(\xi^{(n)}) \notin Y_\alpha$, then there is no reason to reject the hypothesis H , since the observation does not contradict the hypothesis.

The goodness-of-fit test $\delta(x)$, $x \in \mathbf{R}^n$, for the hypothesis H is then such that $\delta(x) = 1$ for all $x \in \{x: T(x) \in Y_\alpha\}$ and $\delta(x) = 0$ for all $x \in \{x: T(x) \in Y \setminus Y_\alpha\}$, that is, $\delta(x)$ is a nonrandomized test for distinguishing two composite, generally speaking, hypotheses (see Section 1.3). The statistic T is called the *statistic of the test* δ , while the set Y_α (or the set $\{x: T(x) \in Y_\alpha\}$) is called the *critical set* for the hypothesis H . As usual, the number α is called the *significance level* or *type I error probability* for the test δ . Below we consider some examples of goodness-of-fit tests.

Kolmogorov goodness-of-fit test. Let the simple hypothesis H be that the distribution function of a random variable ξ equals $F(x)$. As the measure of

disagreement between the data $\xi^{(n)}$ and hypothesis H we consider the statistic

$$(3.1.1) \quad D_n = D_n(\xi^{(n)}) = \sup_{-\infty < x < \infty} |F_n(x) - F(x)|$$

where $F_n(x)$ is the empirical distribution function constructed from the sample $\xi^{(n)}$, that is,

$$(3.1.2) \quad F_n(x) = \nu_n(x)/n, \quad x \in \mathbf{R},$$

$$(3.1.3) \quad \nu_n(x) = \sum_{i=1}^n I_{(-\infty, x)}(\xi_i)$$

where $I_A(x)$ is the indicator of the set $A \subset \mathbf{R}$. It follows from the Glivenko theorem that

$$(3.1.4) \quad \mathbf{P} \left\{ \lim_{n \rightarrow \infty} D_n = 0/H \right\} = 1$$

(see [38], Theorem 1.1.1).

The latter relation allows one to use the statistic D_n to construct a goodness-of-fit test for the hypothesis H by treating small values of the statistic D_n in favor of the hypothesis H . Large values of D_n suggest to a statistician that the hypothesis H is false and it should be rejected. The following result, known as the Kolmogorov theorem, allows one to construct the test for the hypothesis H if $F(x)$ is a continuous function.

THEOREM 3.1.1 (Kolmogorov). *If the function $F(x)$ is continuous, then for all $z > 0$*

$$(3.1.5) \quad \lim_{n \rightarrow \infty} \mathbf{P} \left\{ \sqrt{n}D_n < z/H \right\} = K(z) = \sum_{j=-\infty}^{\infty} (-1)^j e^{-2j^2 z^2}.$$

Note that $K(z)$ is called the *Kolmogorov distribution function* (obviously, we have $K(z) = 0$ for $z \leq 0$).

The proof of Theorem 3.1.1 is quite complicated and we omit it. We only mention that the method of proof of the Kolmogorov theorem and many other limit results for functionals of empirical distribution functions is to show first that

$$(3.1.6) \quad \lim_{n \rightarrow \infty} \mathbf{P} \left\{ \sqrt{n}D_n < z/H \right\} = \mathbf{P} \left\{ \sup_{0 \leq t \leq 1} |w^0(t)| < z \right\}$$

for all $z \in \mathbf{R}$ where $w^0(t)$, $0 \leq t \leq 1$, is the Brownian bridge (see [48], §8). Using then the exact distribution of $\sup_{0 \leq t \leq 1} |w^0(t)|$ (see [48]) we obtain (3.1.5) from (3.1.6). Also see [16, 24] for the proof of Theorem 3.1.1.

Now we construct the goodness-of-fit test for the hypothesis H based on the statistic D_n defined by (3.1.1). Let $\alpha > 0$ be the significance level of the test and let $z(\alpha)$ be a solution of the equation $K(z) = 1 - \alpha$ with respect to z . Then (3.1.5) implies for large n that

$$(3.1.7) \quad \mathbf{P} \left\{ \sqrt{n}D_n \geq z(\alpha)/H \right\} \cong 1 - K(z(\alpha)) = \alpha.$$

Put $\delta_n = I(\sqrt{n}D_n \geq z(\alpha))$. It follows from (3.1.7) that the level of the test δ_n for large n is approximately equal to α . The test δ_n is called the *Kolmogorov goodness-of-fit test*. This test rejects the hypothesis H if $\sqrt{n}D_n \geq z(\alpha)$, that is,

if $D_n \geq z(\alpha)/\sqrt{n}$. Otherwise $\sqrt{n}D_n < z(\alpha)$ and there is no reason to reject the hypothesis H and it is accepted.

Let K_G be the simple hypothesis that the distribution function of a random variable ξ is G . If $\sup_x |G(x) - F(x)| \neq 0$, that is, if the distribution function G differs from F , then the hypothesis H is false given the alternative hypothesis K_G is true. The behavior of the statistic D_n if the hypothesis K_G is true is described in the following result.

THEOREM 3.1.2. *Let the distribution function $F(x)$ be continuous, and let $G(x)$ be another distribution function such that $\sup_x |G(x) - F(x)| \neq 0$. Then for all $z > 0$*

$$(3.1.8) \quad \lim_{n \rightarrow \infty} \mathbf{P} \left\{ \sqrt{n}D_n < z/K_G \right\} = 0.$$

PROOF. According to the Glivenko theorem, for all $\varepsilon > 0$ and $\delta > 0$ there exists $n_0 = n_0(\varepsilon, \delta)$ such that for all $n > n_0$

$$\mathbf{P} \left\{ \sup_x |F_n(x) - G(x)| > \varepsilon/K_G \right\} < \delta.$$

Let $\varepsilon < \sup_x |G(x) - F(x)|$. Then

$$(3.1.9) \quad \begin{aligned} \mathbf{P} \left\{ \sqrt{n}D_n < z/K_G \right\} &= \mathbf{P} \left\{ \sqrt{n}D_n < z, \sup_x |F_n(x) - G(x)| \leq \varepsilon/K_G \right\} \\ &\quad + \mathbf{P} \left\{ \sqrt{n}D_n < z, \sup_x |F_n(x) - G(x)| > \varepsilon/K_G \right\} \\ &\leq \mathbf{P} \left\{ \sqrt{n}(\sup_x |G(x) - F(x)| - \varepsilon) < z/K_G \right\} \\ &\quad + \mathbf{P} \left\{ \sup_x |F_n(x) - G(x)| > \varepsilon/K_G \right\}. \end{aligned}$$

Since $\{\sqrt{n}(\sup_x |G(x) - F(x)| - \varepsilon) < z\}$ is a null event for $n > n_1 = n_1(\varepsilon, z)$ where

$$n_1(\varepsilon, z) = z^2 \left(\sup_x |G(x) - F(x)| - \varepsilon \right)^2,$$

the first term on the right-hand side of inequality (3.1.9) is zero for $n > n_1$. Taking into account (3.1.9) we obtain for $n > n_0 \vee n_1$ that $\mathbf{P}\{\sqrt{n}D_n < z/K_G\} \leq \delta$. Thus relation (3.1.8) is proved. \square

If we put $z = z(\alpha)$ in relation (3.1.8), then

$$(3.1.10) \quad \lim_{n \rightarrow \infty} \mathbf{P} \left\{ \sqrt{n}D_n < z(\alpha)/K_G \right\} = 0$$

which means that the probability to accept the hypothesis H for the Kolmogorov test of level α given the hypothesis K_G is true and if $\sup_x |G(x) - F(x)| \neq 0$ tends to zero as $n \rightarrow \infty$. The tests satisfying condition (3.1.10) are called *consistent*. Thus the Kolmogorov test of level α is consistent. The behavior of the probability $\mathbf{P}\{\sqrt{n}D_n < z(\alpha)/K_G\}$ is studied in [7], §3.12, for distribution functions $G(x)$ that are close to $F(x)$ in a certain sense.

3.2. The Pearson test

The hypothesis. Measure of disagreement between a sample and the hypothesis. Let ξ be a random variable assuming values in a measurable space (X, \mathcal{B}) . Consider a partition of the set X into r , $r \geq 2$, domains:

$$(3.2.1) \quad X = \bigcup_{i=1}^r S_i$$

where $S_i \cap S_j = \emptyset$, $i \neq j$, and $S_i \in \mathcal{B}$ for all i . Consider the following hypothesis H about the distribution of ξ :

$$P\{\xi \in S_i/H\} = p_i, \quad i = 1, 2, \dots, r,$$

where p_1, p_2, \dots, p_r are given positive numbers such that $p_1 + p_2 + \dots + p_r = 1$. Let $\xi^{(n)} = (\xi_1, \xi_2, \dots, \xi_n)$ be a sample from the distribution of ξ and denote by ν_i the total number of members of the sample $\xi^{(n)}$ belonging to the domain S_i , that is,

$$(3.2.2) \quad \nu_i = \sum_{j=1}^n I_{S_i}(\xi_j), \quad i = 1, 2, \dots, r.$$

It is clear that $\nu_1 + \nu_2 + \dots + \nu_r = n$. We consider the following measure of disagreement between the sample $\xi^{(n)}$ and hypothesis H :

$$(3.2.3) \quad \zeta_n = \sum_{i=1}^r c_i \left(\frac{\nu_i}{n} - p_i \right)^2$$

where c_i , $i = 1, 2, \dots, r$, are some positive constants. By the Borel strong law of large numbers

$$P \left\{ \lim_{n \rightarrow \infty} \frac{\nu_i}{n} = p_i/H \right\} = 1, \quad i = 1, 2, \dots, r.$$

Thus $P \{ \lim_{n \rightarrow \infty} \zeta_n = 0/H \} = 1$. Moreover if $K_{\tilde{p}}$ is a hypothesis of the form

$$P\{\xi \in S_i/K_{\tilde{p}}\} = \tilde{p}_i, \quad i = 1, 2, \dots, r,$$

where $\tilde{p} = (\tilde{p}_1, \tilde{p}_2, \dots, \tilde{p}_r) \neq (p_1, p_2, \dots, p_r)$, then again by the Borel strong law of large numbers

$$P \left\{ \lim_{n \rightarrow \infty} \zeta_n = \sum_{i=1}^r c_i (\tilde{p}_i - p_i)^2 > 0/K_{\tilde{p}} \right\} = 1.$$

Thus the random variable ζ_n defined by (3.2.3) can be used as a measure of disagreement between the data and hypothesis H .

Pearson studied the behavior of ζ_n in the case of $c_i = n/p_i$, $i = 1, 2, \dots, r$. The random variable ζ_n can be rewritten in this case as

$$(3.2.4) \quad \zeta_n = \sum_{i=1}^r \frac{(\nu_i - np_i)^2}{np_i} = \sum_{i=1}^r \frac{\nu_i^2}{np_i} - n.$$

Consider the distribution of the random variable ζ_n defined by (3.2.4) under the condition that the hypothesis H is true. Then

$$(3.2.5) \quad E\{\zeta_n/H\} = r - 1.$$

Indeed $E\{\nu_i/H\} = np_i$ and $D\{\nu_i/H\} = np_i q_i$ where $q_i = 1 - p_i$. This implies that

$$E\{\nu_i^2/H\} = np_i q_i + (np_i)^2 = np_i(1 + (n-1)p_i).$$

Thus

$$E\{\zeta_n/H\} = \sum_{i=1}^r (1 + (n-1)p_i) - n = r - 1$$

and equality (3.2.5) is proved. Similarly we obtain that

$$(3.2.6) \quad D\{\zeta_n/H\} = 2(r-1) + \frac{1}{n} \left(\sum_{i=1}^r \frac{1}{p_i} - r^2 - 2r + 2 \right).$$

Relations (3.2.5) and (3.2.6) imply that the first two moments of the random variable ζ_n under the condition that the hypothesis H is true converge to the corresponding moments of the $\chi^2(r-1)$ distribution as $n \rightarrow \infty$. Here $\chi^2(r-1)$ stands for the chi-square distribution with $r-1$ degrees of freedom.

Pearson theorem. The following result, known as the *Pearson theorem*, provides the limit distribution of random variables ζ_n as $n \rightarrow \infty$ under the condition that the hypothesis H is true.

THEOREM 3.2.1. For all $x > 0$

$$(3.2.7) \quad \lim_{n \rightarrow \infty} P\{\zeta_n < x/H\} = K_{r-1}(x)$$

where $K_{r-1}(x)$ is the chi-square distribution function with $r-1$ degrees of freedom.

PROOF. Consider the random variables

$$\mu_k^{(i)} = I_{S_i}(\xi_k), \quad i = 1, 2, \dots, r, \quad k = 1, 2, \dots, n.$$

Then

$$\nu_i = \mu_1^{(i)} + \mu_2^{(i)} + \dots + \mu_n^{(i)}, \quad i = 1, 2, \dots, r.$$

The random variables $\mu_1^{(i)}, \mu_2^{(i)}, \dots, \mu_n^{(i)}$ are independent, identically distributed, and such that

$$P\{\mu_k^{(i)} = 1/H\} = P\{\xi_k \in S_i/H\} = p_i, \quad i = 1, 2, \dots, r.$$

Put $\nu = (\nu_1, \nu_2, \dots, \nu_r)'$. Then $\nu = \sum_{j=1}^n \mu_j$ where $\mu_j = (\mu_j^{(1)}, \mu_j^{(2)}, \dots, \mu_j^{(r)})'$, $j = 1, 2, \dots, n$, are independent identically distributed random variables such that

$$P\{\mu_j = e_k/H\} = P\{\mu_j^{(k)} = 1/H\} = p_k, \quad k = 1, 2, \dots, r.$$

Here $e_k = (0, \dots, 0, 1, 0, \dots, 0)'$ is the r -dimensional vector whose k -th coordinate equals 1 and all other coordinates are zero. If $\varphi_{\mu_j}(t)$ denotes the characteristic function of the vector μ_j , then

$$\varphi_{\mu_j}(t) = E\left\{e^{it' \mu_j/H}\right\} = \sum_{k=1}^r p_k e^{it' e_k} = \sum_{k=1}^r p_k e^{it_k}$$

where $t = (t_1, \dots, t_r)'$. This determines the characteristic function of the vector ν :

$$(3.2.8) \quad \varphi_\nu(t) = \prod_{j=1}^n \varphi_{\mu_j}(t) = \left(\sum_{k=1}^r p_k e^{it_k} \right)^n.$$

Now we introduce the random variables

$$(3.2.9) \quad \eta_i = \frac{\nu_i - np_i}{\sqrt{np_i}}, \quad i = 1, 2, \dots, r.$$

It is clear that

$$(3.2.10) \quad \sum_{i=1}^r \eta_i \sqrt{p_i} = 0.$$

It follows from (3.2.4) and (3.2.9) that

$$(3.2.11) \quad \zeta_n = \sum_{i=1}^r \eta_i^2.$$

According to (3.2.8), the characteristic function of the vector $\eta = (\eta_1, \eta_2, \dots, \eta_r)'$ can be rewritten as

$$(3.2.12) \quad \begin{aligned} \varphi_\eta(t) &= E\{e^{it'\eta}/H\} \\ &= E\left\{\exp\left(i\left[\nu_1 \frac{t_1}{\sqrt{np_1}} + \dots + \nu_r \frac{t_r}{\sqrt{np_r}}\right] - i[t_1\sqrt{np_1} + \dots + t_r\sqrt{np_r}]\right)/H\right\} \\ &= \exp\left(-i \sum_{j=1}^r t_j \sqrt{np_j}\right) \varphi_\nu\left(\frac{t_1}{\sqrt{np_1}}, \dots, \frac{t_r}{\sqrt{np_r}}\right) \\ &= \exp\left(-i\sqrt{n} \sum_{j=1}^r t_j \sqrt{p_j}\right) \left(\sum_{k=1}^r p_k \exp\left(i \frac{t_k}{\sqrt{np_k}}\right)\right)^n. \end{aligned}$$

Taking the logarithm of both sides of equality (3.2.12) and then expanding the exponent $\exp(it_k/\sqrt{np_k})$ and logarithm $\ln(1+x)$ into the Taylor series we obtain

$$(3.2.13) \quad \begin{aligned} \ln \varphi_\eta(t) &= n \ln \sum_{k=1}^r p_k e^{it_k/\sqrt{np_k}} - i\sqrt{n} \sum_{j=1}^r t_j \sqrt{p_j} \\ &= n \ln \left[1 + \frac{i}{\sqrt{n}} \sum_{k=1}^r t_k \sqrt{p_k} - \frac{1}{2n} \sum_{k=1}^r t_k^2 + O(n^{-3/2})\right] \\ &\quad - i\sqrt{n} \sum_{j=1}^r t_j \sqrt{p_j} \\ &= -\frac{1}{2} \sum_{k=1}^r t_k^2 + \frac{1}{2} \left(\sum_{k=1}^r t_k \sqrt{p_k}\right)^2 + O(n^{-1/2}). \end{aligned}$$

Passing to the limit in (3.2.13) as $n \rightarrow \infty$ we prove that

$$(3.2.14) \quad \lim_{n \rightarrow \infty} \varphi_\eta(t) = \exp\left(-\frac{1}{2}Q(t)\right)$$

uniformly with respect to t in any bounded domain where

$$(3.2.15) \quad Q(t) = \sum_{k=1}^r t_k^2 - \left(\sum_{k=1}^r t_k \sqrt{p_k}\right)^2, \quad t = (t_1, t_2, \dots, t_r)'$$

It is clear that the quadratic form $Q(t)$ defined by (3.2.15) can be represented as $Q(t) = t'\Lambda t$ where Λ is a matrix such that $\Lambda = I - pp'$. The symbol I stands for the $r \times r$ unit matrix, while $p = (\sqrt{p_1}, \dots, \sqrt{p_r})'$. Thus the right-hand side of (3.2.14)

is $\varphi_{\eta^0}(t) = \exp\{-\frac{1}{2}t'\Lambda t\}$ which is the characteristic function of the $\mathcal{N}(0, \Lambda)$ normal vector $\eta^0 = (\eta_1^0, \eta_2^0, \dots, \eta_r^0)'$. Therefore convergence (3.2.14) implies that

$$\lim_{n \rightarrow \infty} \varphi_{\eta}(t) = \varphi_{\eta^0}(t) = \exp\left(-\frac{1}{2}t'\Lambda t\right)$$

uniformly with respect to t in any bounded domain. The continuity theorem for characteristic functions implies then that as $n \rightarrow \infty$

$$\mathcal{L}(\eta|H) \xrightarrow{w} \mathcal{L}(\eta^0|H) = \mathcal{N}(0, \Lambda),$$

whence

$$(3.2.16) \quad \mathcal{L}(\zeta_n|H) \xrightarrow{w} \mathcal{L}(\zeta^0|H), \quad n \rightarrow \infty,$$

in view of (3.2.11) where

$$(3.2.17) \quad \zeta^0 = \sum_{i=1}^r (\eta_i^0)^2.$$

Equality (3.2.10) for the limit vector η^0 implies that

$$(3.2.18) \quad \sum_{i=1}^r \eta_i^0 \sqrt{p_i} = 0.$$

Let A be an orthogonal $r \times r$ matrix whose bottom row is $(\sqrt{p_1}, \sqrt{p_2}, \dots, \sqrt{p_r})$. Then by (3.2.18)

$$(3.2.19) \quad \varkappa_r = \sum_{i=1}^r \eta_i^0 \sqrt{p_i} = 0$$

where $\varkappa = (\varkappa_1, \dots, \varkappa_r)' = A\eta^0$. On the other hand, it is known that orthogonal transformations do not change the canonical representation of quadratic forms. Thus the quadratic form (3.2.17) is given by

$$(3.2.20) \quad \zeta^0 = \sum_{i=1}^r (\eta_i^0)^2 = \sum_{i=1}^r \varkappa_i^2 = \sum_{i=1}^{r-1} \varkappa_i^2$$

in view of (3.2.19).

The quadratic form $Q(t)$ can be rewritten for new coordinates $u = At$ as

$$Q(t) = \sum_{k=1}^r t_k^2 - \left(\sum_{k=1}^r t_k \sqrt{p_k} \right)^2 = \sum_{k=1}^r u_k^2 - u_r^2 = \sum_{k=1}^{r-1} u_k^2 = Q(A^{-1}u).$$

Therefore the characteristic function of the vector \varkappa can be represented for new coordinates as

$$\begin{aligned} \varphi_{\varkappa}(u) &= \mathbf{E}\left(e^{iu'\varkappa}/H\right) = \mathbf{E}\left(e^{iu'A\eta^0}/H\right) = \mathbf{E}\left(e^{i(A'u)'\eta^0}/H\right) \\ &= \exp\left(-\frac{1}{2}Q(A^{-1}u)\right) = \exp\left(-\frac{1}{2}\sum_{k=1}^{r-1} u_k^2\right), \end{aligned}$$

that is, the coordinates $\varkappa_1, \dots, \varkappa_{r-1}$ of the vector $\varkappa = (\varkappa_1, \dots, \varkappa_r)$ are independent $\mathcal{N}(0, 1)$ identically distributed random variables and $\varkappa_r = 0$. Now it follows from (3.2.20) that $\mathcal{L}(\zeta^0|H) = \chi^2(r-1)$. Therefore (3.2.7) follows from (3.2.16). \square

The Pearson goodness-of-fit test. Applying Theorem 3.2.1 one can construct a goodness-of-fit test for the hypothesis H in a way similar to that used to construct the Kolmogorov goodness-of-fit test based on the Kolmogorov limit theorem.

Let $\alpha > 0$ be a significance level and let $z_{r-1}(\alpha)$ be a solution of the equation $K_{r-1}(z) = 1 - \alpha$ with respect to z where $K_{r-1}(z)$ is the chi-square distribution function with $r - 1$ degrees of freedom. Then relation (3.2.7) implies for large n that

$$(3.2.21) \quad P\{\zeta_n \geq z_{r-1}(\alpha)/H\} \approx 1 - K_{r-1}(z_{r-1}(\alpha)) = \alpha.$$

Now the test for the hypothesis H is $\delta_n = I(\zeta_n \geq z_{r-1}(\alpha))$. Relation (3.2.21) implies that the level of this test is approximately equal to α . The test is called the *Pearson goodness-of-fit test*. Sometimes it is called the *chi-square test*. The Pearson test rejects the hypothesis H if $\zeta_n \geq z_{r-1}(\alpha)$, and it accepts the hypothesis H if $\zeta_n < z_{r-1}(\alpha)$.

Let $K_{\tilde{p}}$ be the simple hypothesis of the form $P\{\xi \in S_i/K_{\tilde{p}}\} = \tilde{p}_i$, $i = 1, 2, \dots, r$, where $\tilde{p} = (\tilde{p}_1, \tilde{p}_2, \dots, \tilde{p}_r) \neq p = (p_1, p_2, \dots, p_r)$. It turns out that the Pearson test of level $\alpha \in (0, 1)$ is consistent.

THEOREM 3.2.2. For all vectors $\tilde{p} \neq p$

$$\lim_{n \rightarrow \infty} P\{\zeta_n < z_{r-1}(\alpha)/K_{\tilde{p}}\} = 0.$$

The proof of Theorem 3.2.2 can be found in [37], Theorem 3.2.

Examples of the Pearson goodness-of-fit tests for special models of observations can be found in [7, 9, 14, 26, 34].

REMARK 3.2.1. The hypothesis H tested with the help of the Pearson test is, generally speaking, composite. This hypothesis is simple only in the case where the vector $\mu = (\mu^{(1)}, \mu^{(2)}, \dots, \mu^{(r)})$ assumes values e_k , $k = 1, 2, \dots, r$, and the hypothesis H is that $P\{\mu = e_k/H\} = p_k$, $k = 1, 2, \dots, r$. One can show in this case that the Pearson test coincides with the likelihood ratio test (see [7, 9]).

REMARK 3.2.2. Let H_0 be the hypothesis that the distribution function of ξ is $F(x)$, $x \in \mathbf{R} = (-\infty, \infty)$, and $\mathbf{R} = \bigcup_{i=1}^r X_i$ where $X_i \cap X_j = \emptyset$, $i \neq j$. Moreover let $P\{\xi \in X_i/H_0\} = p_i$, $i = 1, 2, \dots, r$, and $p_1 + p_2 + \dots + p_r = 1$. As above let H be the hypothesis that $P\{\xi \in X_i/H\} = p_i$, $i = 1, 2, \dots, r$. The Pearson test constructed above for the hypothesis H is sometimes used to test the hypothesis H_0 . However there are distributions $G(x)$ such that the hypothesis H is true but $\sup_x |G(x) - F(x)| > 0$, that is, the null hypothesis H_0 is false. The Pearson test does not detect the difference between such functions $G(x)$ and $F(x)$ and therefore is not consistent for testing the hypothesis H_0 .

Quantile test. Sign test. Let a random variable ξ be real-valued, that is, $X = \mathbf{R} = (-\infty, \infty)$. Assume that the hypothesis H is that

$$F(y_i) = p_i, \quad i = 1, 2, \dots, r - 1,$$

where $F(x)$ is the distribution function of the random variable ξ ,

$$0 < p_1 < \dots < p_{r-1} < 1$$

are given numbers,

$$-\infty < y_1 < \cdots < y_{r-1} < \infty$$

are quantiles of levels p_1, p_2, \dots, p_{r-1} , respectively, and $r \geq 2$. Thus the hypothesis H is composite and deals with all distributions with fixed quantiles and their levels.

Put $S_i = [y_{i-1}, y_i)$, $i = 1, 2, \dots, r$, $y_0 = -\infty$, and $y_r = \infty$. Let

$$\xi^{(n)} = (\xi_1, \xi_2, \dots, \xi_n)$$

be a sample of size n and let ν_i be the number of its members ξ_j belonging to the interval S_i . One can apply the Pearson test for the hypothesis H based on the statistic ζ_n defined by (3.2.4). The test in this case is called the *quantile test*. If $r = 2$ and $p_1 = 0,5$, then the corresponding test is called the *sign test*. The null hypothesis H in the latter case is that the median of the distribution of a random variable ξ is y_1 . The statistic (3.2.4) in this case is $(4/n)(\nu_1 - n/2)^2$ where ν_1 is the total number of members of the sample $\xi^{(n)}$ belonging to the interval $(-\infty, y_1)$. In other words, ν_1 is the number of negative signs in the sequence $\xi_i - y_1$, $i = 1, 2, \dots, n$.

The sign test is used under the following assumptions. Let $(\xi_1, \eta_1), \dots, (\xi_n, \eta_n)$ be a sample of size n whose members are independent observations of the vector (ξ, η) . One needs to test the hypothesis H_0 that the coordinates ξ and η are independent and identically distributed, that is, $F(x, y) = F(x)F(y)$ where $F(x, y)$ is the distribution function of the vector (ξ, η) and $F(x)$ is the distribution function of the random variable ξ (and, of course, of η if the hypothesis H_0 is true). Let $\zeta_i = \xi_i - \eta_i$, $i = 1, 2, \dots, n$. Then $P\{\zeta_i < 0/H_0\} = P\{\zeta_i > 0/H_0\} = 1/2$ and the null hypothesis H_0 is that the data $(\zeta_1, \zeta_2, \dots, \zeta_n)$ is sampled from a distribution whose median is 0. The statistic ν_1 in this case is the total number of negative members in the sequence $\zeta_1, \zeta_2, \dots, \zeta_n$. According to Theorem 3.2.1

$$\mathcal{L}(4(\nu_1 - n/2)^2/n|H_0) \xrightarrow{w} \chi^2(1), \quad n \rightarrow \infty.$$

This relation allows one to construct the sign test for the null hypothesis H_0 that random variables ξ and η are independent and identically distributed for a given level α . The procedure is the same as in the case of the Pearson test. More details on the sign tests are given in [5].

The Pearson goodness-of-fit test for distributions with unknown parameters. Let a random variable ξ assume values in a measurable space (X, \mathcal{B}) and let its distribution depend on an unknown s -dimensional parameter

$$\theta = (\theta_1, \theta_2, \dots, \theta_s)'$$

where $\theta_1, \theta_2, \dots, \theta_s$ are real numbers. As above we introduce a partition (3.2.1) of the set X consisting of r domains.

Let the hypothesis H about the distribution of ξ be that

$$(3.2.22) \quad P\{\xi \in S_i/H\} = p_i(\theta), \quad i = 1, 2, \dots, r,$$

where $p_1(\theta), p_2(\theta), \dots, p_r(\theta)$ are known functions of the parameter θ such that $p_1(\theta) + p_2(\theta) + \cdots + p_r(\theta) = 1$. The parameter θ is unknown and our aim is to estimate it by the observation ξ . Let $\xi^{(n)} = (\xi_1, \xi_2, \dots, \xi_n)$ be a sample of size n and let $\nu_1, \nu_2, \dots, \nu_r$ be the numbers of members of the sample $\xi_1, \xi_2, \dots, \xi_n$ belonging to the domains S_1, S_2, \dots, S_r , respectively (see relation (3.2.2)). The measure

of disagreement (3.2.4) between the sample and hypothesis (3.2.22) is in this case given by

$$(3.2.23) \quad \zeta_n = \zeta_n(\theta) = \sum_{i=1}^r \frac{(\nu_i - np_i(\theta))^2}{np_i(\theta)}.$$

Note that ζ_n depends on the unknown parameter θ and this does not allow one to use ζ_n to construct a test for the hypothesis H . Thus the first step is to exclude the unknown parameter from (3.2.23). To do so we substitute into (3.2.23) an estimator $\hat{\theta}_n = \hat{\theta}_n(\xi^{(n)})$ for θ and obtain the statistic

$$(3.2.24) \quad \hat{\zeta}_n = \zeta_n(\hat{\theta}_n) = \sum_{i=1}^r \frac{(\nu_i - np_i(\hat{\theta}_n))^2}{np_i(\hat{\theta}_n)}.$$

The statistic $\hat{\zeta}_n$ defined by (3.2.24) depends on the sample $\xi^{(n)}$ and does not depend on θ and thus it can be used to test the hypothesis H . When constructing a test for the hypothesis H based on the statistic $\xi^{(n)}$ one needs to know its distribution or, at least, its limit distribution (the latter depends on the estimator $\hat{\theta}_n$). Below we consider the most famous method of estimation of the parameter θ which leads to a simple limit distribution of the statistic $\hat{\zeta}_n$. This method was successfully used by R. Fisher, J. Neyman, and K. Pearson in the early twentieth century.

As an estimator of θ it is natural to take a value of the parameter for which $\zeta_n(\theta)$ defined by (3.2.23) attains its minimum. This is the so-called *minimum χ^2 method*. If the derivatives exist, then the problem of finding such a value is reduced to solving the following system of equations with respect to θ :

$$(3.2.25) \quad \frac{\partial \zeta_n(\theta)}{\partial \theta_j} = -2 \sum_{i=1}^r \left(\frac{\nu_i - np_i(\theta)}{p_i(\theta)} + \frac{(\nu_i - np_i(\theta))^2}{2np_i^2(\theta)} \right) \frac{\partial p_i(\theta)}{\partial \theta_j} = 0, \\ j = 1, 2, \dots, s.$$

Note however that this system is not easy to solve even in the simplest cases. On the other hand, one can show that the influence of the second term in parentheses is negligible for large n . Omitting this term, system (3.2.25) becomes of the form

$$(3.2.26) \quad \sum_{i=1}^r \frac{\nu_i - np_i(\theta)}{p_i(\theta)} \frac{\partial p_i(\theta)}{\partial \theta_j} = 0, \quad j = 1, 2, \dots, s.$$

The method of estimation based on solving the system (3.2.26) is called the *modified minimum χ^2 method*. Under rather general assumptions both methods have the same limit distribution $\zeta_n(\hat{\theta}_n)$ as $n \rightarrow \infty$. Below we consider a simpler method based on solving the system (3.2.26).

Since $p_1(\theta) + p_2(\theta) + \dots + p_r(\theta) = 1$ for all θ , system (3.2.26) becomes of the form

$$(3.2.27) \quad \sum_{i=1}^r \frac{\nu_i}{p_i(\theta)} \frac{\partial p_i(\theta)}{\partial \theta_j} = 0, \quad j = 1, 2, \dots, s.$$

The system (3.2.27) can be rewritten as

$$(3.2.28) \quad \frac{\partial \ln L_n(\theta)}{\partial \theta_j} = 0, \quad j = 1, 2, \dots, s,$$

where $L_n(\theta) = p_1^{\nu_1}(\theta) \dots p_r^{\nu_r}(\theta)$. The method of estimation based on solving the system (3.2.28) is nothing else but the maximum likelihood method for the polynomial distribution. Thus the estimator $\hat{\theta}_n$ obtained as a solution of system (3.2.27) (as well as that of (3.2.26)) is called the *maximum likelihood polynomial estimator*.

The limit distribution of the random variable $\hat{\zeta}_n = \zeta_n(\hat{\theta}_n)$ as $n \rightarrow \infty$ is described in the following result where $\hat{\theta}_n$ is the maximum likelihood polynomial estimator.

THEOREM 3.2.3. *Let the functions $p_j(\theta)$, $j = 1, 2, \dots, r$, $\theta = (\theta_1, \theta_2, \dots, \theta_s)'$, $s < r$, be such that:*

- 1) $p_1(\theta) + p_2(\theta) + \dots + p_r(\theta) = 1$ for all θ ;
- 2) $p_i(\theta) \geq c > 0$ for all $i = 1, 2, \dots, r$ and the derivatives $\frac{\partial p_i(\theta)}{\partial \theta_j}$ and $\frac{\partial^2 p_i(\theta)}{\partial \theta_j \partial \theta_k}$, $j, k = 1, 2, \dots, s$, $i = 1, 2, \dots, r$, are continuous;
- 3) the rank of the $r \times s$ matrix $\left\| \frac{\partial p_i(\theta)}{\partial \theta_j} \right\|$ is equal to s for all θ .

Then

$$(3.2.29) \quad \lim_{n \rightarrow \infty} P \left\{ \hat{\zeta}_n < z/H \right\} = K_{r-s-1}(z) \quad \text{for all } z > 0$$

where $\hat{\zeta}_n = \zeta_n(\hat{\theta}_n)$ and $\hat{\theta}_n$ is the maximum likelihood polynomial estimator.

The proof of Theorem 3.2.3 can be found in [14].

Based on Theorem 3.2.3 the goodness-of-fit test for the hypothesis (3.2.22) is constructed in the same way as in the case of the Pearson test for distributions whose parameters are known. The constructed test is also called the *Pearson test*.

REMARK 3.2.3. An estimator of θ can be evaluated without ranking the data. This can be done, for example, by maximizing the likelihood function

$$f(x_1; \theta) f(x_2; \theta) \dots f(x_n; \theta)$$

with respect to θ where $f(x; \theta)$ is the density of the distribution of ξ with respect to some measure. The estimator of the parameter θ is not based in this case on frequencies $\nu_1, \nu_2, \dots, \nu_r$ for domains S_1, S_2, \dots, S_r but uses the observations $\xi_1, \xi_2, \dots, \xi_n$ instead. Chernoff and Lehmann (see [12]) showed however that the limit relation (3.2.29) does not hold for this method of estimation.

EXAMPLE 3.2.1. Let a random variable ξ assume values $0, 1, 2, \dots$. Set

$$S_j = \{j - 1\}, \quad j = 1, 2, \dots, r - 1,$$

and $S_r = \{r - 1, r, \dots\}$. Let the hypothesis H be defined by (3.2.22) where

$$p_j(\theta) = \frac{\theta^{j-1}}{(j-1)!} e^{-\theta}, \quad j = 1, 2, \dots, r - 1,$$

$$p_r(\theta) = \sum_{i=r-1}^{\infty} \frac{\theta^i}{i!} e^{-\theta}, \quad \theta > 0.$$

In this case $s = 1$ and thus the system (3.2.27) is reduced to the equation

$$\sum_{j=\nu}^{r-2} \binom{j}{\theta - 1} \nu_{j+1} + \nu_r \sum_{i=r-1}^{\infty} \binom{i}{\theta - 1} \frac{\theta^i}{i!} \left(\sum_{i=r-1}^{\infty} \frac{\theta^i}{i!} \right)^{-1} = 0,$$

whence

$$\theta = \frac{1}{n} \left[\sum_{j=\nu}^{r-2} j\nu_{j+1} + \nu_r \sum_{i=r-1}^{\infty} i \frac{\theta^i}{i!} \left(\sum_{i=r-1}^{\infty} \frac{\theta^i}{i!} \right)^{-1} \right]$$

where $\nu_1, \nu_2, \dots, \nu_r$ are the numbers of members of the sample $\xi^{(n)} = (\xi_1, \xi_2, \dots, \xi_n)$ belonging to S_1, S_2, \dots, S_r , respectively. The first term in the square brackets is equal to the sum of all members of the sample $\xi_1, \xi_2, \dots, \xi_n$ such that $\xi_i \leq r-2$, while the second term is approximately equal to the sum of all members of the sample $\xi_1, \xi_2, \dots, \xi_n$ that are greater than or equal to $r-1$. Thus as an estimator $\hat{\theta}_n$ of the parameter θ one can take the sampling mean $\hat{\theta}_n = \bar{\xi}$. Note that the maximum likelihood estimator of the parameter θ is equal to $\bar{\xi}$ in the case of the Poisson distribution. Note also that the limit distribution as $n \rightarrow \infty$ of the statistic $\hat{\zeta}_n = \zeta_n(\hat{\theta}_n)$ is the chi-square distribution with $r-2$ degrees of freedom (this result follows from Theorem 3.2.3).

EXAMPLE 3.2.2. Let

$$\xi^{(n)} = (\xi_1, \xi_2, \dots, \xi_n)$$

be a sample of size n . Let $\nu_1, \nu_2, \dots, \nu_r$ be the number of members of the sample $\xi^{(n)}$ belonging to S_1, S_2, \dots, S_r , respectively (see relation (3.2.2)), where $S_1 = (-\infty, x_1 + \frac{1}{2}h)$, $S_i = (x_i - \frac{1}{2}h, x_i + \frac{1}{2}h)$ for $i = 2, 3, \dots, r-1$, $S_r = (x_r - \frac{1}{2}h, \infty)$, $x_i = x_1 + (i-1)h$ for $i = 2, 3, \dots, r$, and x_1 is some number of $(-\infty, \infty)$. Let the hypothesis H be defined by (3.2.22) where

$$p_i(\theta) = \int_{S_i} \varphi(x; \theta) dx, \quad i = 1, 2, \dots, r,$$

$\varphi(x; \theta)$ is the density of the normal $\mathcal{N}(m, \sigma^2)$ law, and $\theta = (m, \sigma)$. Then the system of equations (3.2.27) becomes of the form

$$m = \frac{1}{n} \sum_{i=1}^r \nu_i \int_{S_i} x \varphi(x; \theta) dx \left(\int_{S_i} \varphi(x; \theta) dx \right)^{-1},$$

$$\sigma^2 = \frac{1}{n} \sum_{i=1}^r \nu_i \int_{S_i} (x - m)^2 \varphi(x; \theta) dx \left(\int_{S_i} \varphi(x; \theta) dx \right)^{-1}.$$

First we assume that x_1 and r are such that $\nu_1 = \nu_r = 0$. If h is small, then an approximate solution can be obtained by substituting the values of integrands at the middle points x_i of the corresponding intervals S_i instead of the integrands. Then we get the estimators \hat{m}_n and $\hat{\sigma}_n$ defined by

$$\hat{m}_n = \frac{1}{n} \sum_i \nu_i x_i, \quad \hat{\sigma}_n^2 = \frac{1}{n} \sum_i \nu_i (x_i - \hat{m}_n)^2.$$

To improve an approximation of the solution one can expand the integrands into the Taylor series in the neighborhoods of points x_i . Then the expressions for \hat{m}_n and $\hat{\sigma}_n$ for small h become of the form

$$\hat{m}_n = \frac{1}{n} \sum_i \nu_i x_i + O(h^4), \quad \hat{\sigma}_n^2 = \frac{1}{n} \sum_i \nu_i (x_i - \hat{m}_n)^2 - \frac{h^2}{12} + O(h^4).$$

Thus omitting the terms of order h^4 we obtain estimators for m and σ^2 . These estimators are the mean value and variance of the ranked sample \bar{x} corrected with

the help of the *Sheppard correction* $\frac{h^2}{12}$ [14]. The *ranked sample* is obtained from the original sample if one substitutes the middle point of an underlying interval instead of any member of the sample belonging to this interval.

This procedure gives a better approximation even if h is not small but the end intervals are not empty and contain a small proportion of sampling values. Often it is convenient to merge the end intervals such that the union contains at least 10 sampling values. As estimators for m and σ^2 one can take \bar{x} and s^2 evaluated by the original ranked sample with the Sheppard correction applied to s^2 . If r' is the number of groups in the merged sample used for the evaluation of $\hat{\zeta}_n$, then the limit distribution of $\hat{\zeta}_n$ has $r' - 3$ degrees of freedom, since two parameters are already estimated from the sample.

More details on the Sheppard correction are given in [14] (our Examples 3.2.1 and 3.2.2 are considered there for particular numerical values).

3.3. Smirnov test

The hypothesis on the homogeneity. Measure of disagreement between a sample and the hypothesis. It is an important applied problem to check whether the data is homogeneous. More precisely let $\xi^{(n)} = (\xi_1, \xi_2, \dots, \xi_n)$ be a sample of size n consisting of n independent observations of a random variable ξ and let $\eta^{(m)} = (\eta_1, \eta_2, \dots, \eta_m)$ be a sample of size m , independent of $\xi^{(n)}$ and consisting of m independent observations of a random variable η . The *hypothesis H on the homogeneity* is that the distributions of the random variables ξ and η coincide, that is, $P\{\xi < x/H\} = P\{\eta < x/H\}$ for all $x \in (-\infty, \infty)$. In other words, the hypothesis H means that the samples $\xi^{(n)}$ and $\eta^{(m)}$ are, in fact, the observations of the same random variable.

Let $S_n(x)$ be the empirical distribution function constructed from the sample $\xi^{(n)}$ and let $T_m(x)$ be the empirical distribution function constructed from the sample $\eta^{(m)}$. Consider the following measures of disagreement between the sample and hypothesis H :

$$(3.3.1) \quad D_{n,m}^+ = \sup_x [S_n(x) - T_m(x)],$$

$$(3.3.2) \quad D_{n,m} = \sup_x |S_n(x) - T_m(x)|.$$

Since

$$S_n(x) - T_m(x) = (S_n(x) - F(x)) - (T_m(x) - F(x)),$$

where $F(x) = P\{\xi < x/H\} = P\{\eta < x/H\}$ is the common distribution function of ξ and η if the hypothesis H is true, the Glivenko theorem implies

$$P \left\{ \lim_{n,m \rightarrow \infty} D_{n,m}^+ = 0/H \right\} = P \left\{ \lim_{n,m \rightarrow \infty} D_{n,m} = 0/H \right\} = 1.$$

On the other hand, if $H_{F,G}$ is the hypothesis that $P\{\xi < x\} = F(x)$ and $P\{\eta < x\} = G(x)$ where $F(x)$ and $G(x)$ are distribution functions such that $\sup_x |F(x) - G(x)| \neq$

0, then again by the Glivenko theorem

$$P \left\{ \lim_{n,m \rightarrow \infty} D_{n,m}^+ = \sup_x [F(x) - G(x)]/H_{F,G} \right\} = 1,$$

$$P \left\{ \lim_{n,m \rightarrow \infty} D_{n,m} = \sup_x |F(x) - G(x)|/H_{F,G} \right\} = 1.$$

The above argument shows that the statistics defined by (3.3.1) and (3.3.2) can be used as measures of disagreement between the data and hypothesis on the homogeneity.

Similarly to the Kolmogorov and Pearson goodness-of-fit tests we apply limit results on the behavior of measures of disagreement between the sample and the hypothesis. In the case of the hypothesis H we deal with $D_{n,m}^+$ and $D_{n,m}$ as $n, m \rightarrow \infty$. The method described below is due to Gnedenko. It allows one to find the distributions of statistics $D_{n,m}^+$ and $D_{n,m}$. We restrict the discussion to the case $m = n$. Other methods for studying the limit behavior of $D_{n,m}^+$ and $D_{n,m}$ can be found in [24].

The distributions of statistics $D_{n,m}^+$ and $D_{n,m}$. The following result contains the explicit expressions for the distribution functions of $D_{n,m}^+$ and $D_{n,m}$ if the hypothesis H is true.

THEOREM 3.3.1. *If the distribution function $P\{\xi < x/H\} = P\{\eta < x/H\}$ is continuous, then*

$$(3.3.3) \quad P \left\{ \sqrt{\frac{n}{2}} D_{n,m}^+ < z/H \right\} = \begin{cases} 0, & z \leq 0, \\ 1 - \binom{2n}{n-c} / \binom{2n}{n}, & 0 < z \leq \sqrt{n/2}, \\ 1, & z > \sqrt{n/2}, \end{cases}$$

$$(3.3.4) \quad P \left\{ \sqrt{\frac{n}{2}} D_{n,m} < z/H \right\} = \begin{cases} 0, & z \leq 1/\sqrt{2n}, \\ \sum_{|k| \leq [n/c]} (-1)^k \binom{2n}{n-kc} / \binom{2n}{n}, & 1/\sqrt{2n} < z \leq \sqrt{n/2}, \\ 1, & z > \sqrt{n/2}, \end{cases}$$

where $c =]z\sqrt{2n}[$ is the minimal integer number which is greater than or equal to $z\sqrt{2n}$.

PROOF. We assume below that the hypothesis H is true. Then the random variables ξ_1, \dots, ξ_n and η_1, \dots, η_n are independent and identically distributed and their common distribution function is $P\{\xi < x/H\} = P\{\eta < x/H\}$. We rearrange the random variables $\xi_1, \dots, \xi_n, \eta_1, \dots, \eta_n$ in ascending order:

$$\zeta_1 < \zeta_2 < \dots < \zeta_{2n}.$$

The equalities in this sequence may occur with probability 0, since the distribution function is continuous. Define the random variables $\chi_1, \chi_2, \dots, \chi_{2n}$ as follows: $\chi_k = 1$ if ζ_k is a member of the sample $\xi_1, \xi_2, \dots, \xi_n$, while $\chi_k = -1$ if ζ_k is a member of the sample $\eta_1, \eta_2, \dots, \eta_n$.

Put

$$S_0 = 0, \quad S_k = \sum_{i=1}^k \chi_i, \quad k = 1, 2, \dots, 2n.$$

It is clear that

$$(3.3.5) \quad nD_{n,m}^+ = \sup_{0 \leq k \leq 2n} S_k, \quad nD_{n,m} = \sup_{0 \leq k \leq 2n} |S_k|.$$

Consider the points (k, S_k) on the plane (t, x) for $k = 0, 1, \dots, 2n$ and join them by segments of straight lines. As a result we get a polygonal line for which there are n subintervals $[k-1, k]$ of $[0, 2n]$ where the line goes up and n subintervals $[l-1, l]$ of $[0, 2n]$ where it goes down. This line starts at the point $(0, 0)$ and ends at the point $(2n, 0)$. Polygonal lines with these properties are called *trajectories*. Since the number of intervals of any trajectory where it rises is equal to n and the number of intervals where it descends is equal to n , the total number of trajectories is $\binom{2n}{n}$. All these trajectories are equiprobable. Indeed, any trajectory corresponds to the event $\{\xi_{i_1} < \xi_{i_2} < \dots < \xi_{i_{2n}}\}$ where i_1, i_2, \dots, i_{2n} is a permutation of the numbers $1, 2, \dots, 2n$ and $\xi_{i_1} < \xi_{i_2} < \dots < \xi_{i_{2n}}$ are random variables $\xi_1, \xi_2, \dots, \xi_n, \xi_{n+1}, \dots, \xi_{2n}$ arranged in ascending order, $\xi_{n+k} = \eta_k$, $k = 1, 2, \dots, n$. Since the random variables $\xi_1, \xi_2, \dots, \xi_{2n}$ are identically distributed, we have

$$P\{\xi_{i_1} < \xi_{i_2} < \dots < \xi_{i_{2n}}/H\} = P\{\xi_{j_1} < \xi_{j_2} < \dots < \xi_{j_{2n}}/H\}$$

for all permutations i_1, i_2, \dots, i_{2n} and j_1, j_2, \dots, j_{2n} of the numbers $1, 2, \dots, 2n$. Thus the probability of any of the trajectories is $1/\binom{2n}{n}$.

According to (3.3.5) the random variables $nD_{n,n}^+$ and $nD_{n,n}$ assume integer values and thus

$$P\left\{\sqrt{\frac{n}{2}}D_{n,n}^+ < z/H\right\} = P\{nD_{n,n}^+ < c/H\},$$

$$P\left\{\sqrt{\frac{n}{2}}D_{n,n} < z/H\right\} = P\{nD_{n,n} < c/H\}$$

where $c =]z\sqrt{2n}[$. It remains to evaluate the probabilities

$$P\left\{\sup_{0 \leq k \leq 2n} S_k < c/H\right\} \quad \text{and} \quad P\left\{\sup_{0 \leq k \leq 2n} |S_k| < c/H\right\}$$

to complete the proof of (3.3.3) and (3.3.4). Since all the trajectories are equiprobable, we need to determine the total numbers of trajectories favorable to the events

$$\left\{\sup_{0 \leq k \leq 2n} S_k < c\right\} \quad \text{and} \quad \left\{\sup_{0 \leq k \leq 2n} |S_k| < c\right\}.$$

First we prove equality (3.3.3). We determine the total number of trajectories favorable to the event

$$\left\{\sup_{0 \leq k \leq 2n} S_k < c\right\} = \{nD_{n,n}^+ < c\},$$

that is, the total number of trajectories below the straight line $x = c$ (line α). We obtain this number by evaluating the total number of trajectories favorable to the converse event

$$\left\{\sup_{0 \leq k \leq 2n} S_k \geq c\right\} = \{nD_{n,n}^+ \geq c\},$$

that is, the total number of trajectories that have common points with the line α (we say in this case that a trajectory meets α). Every trajectory meeting the line α

(and called an old trajectory in this case) corresponds to another (new) trajectory defined as follows: the new trajectory coincides with the old one from the point $(0, 0)$ until it meets the line α for the first time; after this point the new trajectory is the mirror reflection of the old one. Thus the new trajectory starts at the point $(0, 0)$ and ends at the point $(2n, 2c)$. The total number of different new trajectories (hence the trajectories meeting the line α) is equal to $\binom{2n}{n-c}$, since the number of intervals where a new trajectory rises is equal to $n+c$, while the number of intervals where it descends is equal to $n-c$. Thus the total number of trajectories that do not meet the line α is equal to $\binom{2n}{n} - \binom{2n}{n-c}$ and equality (3.3.3) is proved.

Now we turn to equality (3.3.4). We split the set \mathfrak{M} of all trajectories into disjoint subsets \mathfrak{A}_i , $i \geq 0$, and \mathfrak{B}_i , $i \geq 1$: the set \mathfrak{A}_0 consists of trajectories that do not meet both lines $x = c$ (line α) and $x = -c$ (line β); the set \mathfrak{A}_1 consists of trajectories meeting the line α but not the line β ; the set \mathfrak{B}_1 consists of trajectories meeting β but not α ; the set \mathfrak{A}_2 consists of trajectories meeting first α , then β , and then not meeting α anymore; the set \mathfrak{B}_2 consists of trajectories meeting first β , then α , and then not meeting β anymore; the set \mathfrak{A}_3 consists of trajectories meeting first α , then β , then again α , and then not meeting β anymore, and so on. Obviously these sets are eventually empty. Moreover

$$(3.3.6) \quad \mathfrak{M} = \mathfrak{A}_0 \cup \left(\bigcup_{i \geq 1} (\mathfrak{A}_i \cup \mathfrak{B}_i) \right)$$

and the sets $\mathfrak{A}_1, \mathfrak{A}_2, \dots$ and $\mathfrak{B}_1, \mathfrak{B}_2, \dots$ are disjoint.

Along with the sets defined above we introduce the following sequence: the set A_1 consists of trajectories meeting α at least one time; the set B_1 consists of trajectories meeting β at least one time; the set A_2 consists of trajectories meeting α at least one time and then meeting β ; the set B_2 consists of trajectories meeting β at least one time and then meeting α ; the set A_3 consists of trajectories meeting α at least two times and β at least one time each in the following order: first α , then β , then α , and so on. It is clear that

$$\begin{aligned} A_1 &= \mathfrak{A}_1 \cup \left(\bigcup_{i \geq 2} (\mathfrak{A}_i \cup \mathfrak{B}_i) \right), & B_1 &= \mathfrak{B}_1 \cup \left(\bigcup_{i \geq 2} (\mathfrak{A}_i \cup \mathfrak{B}_i) \right), \\ A_2 &= \mathfrak{A}_2 \cup \left(\bigcup_{i \geq 3} (\mathfrak{A}_i \cup \mathfrak{B}_i) \right), & B_2 &= \mathfrak{B}_2 \cup \left(\bigcup_{i \geq 3} (\mathfrak{A}_i \cup \mathfrak{B}_i) \right), \end{aligned}$$

and so on. This implies for all $i \geq 1$ that

$$(A_{2i-1} \setminus A_{2i}) \cup (B_{2i-1} \setminus B_{2i}) = \mathfrak{A}_{2i-1} \cup \mathfrak{A}_{2i} \cup \mathfrak{B}_{2i-1} \cup \mathfrak{B}_{2i}.$$

The latter equality together with (3.3.6) implies

$$(3.3.7) \quad \mathfrak{A}_0 = \mathfrak{M} \setminus \left(\bigcup_{i \geq 1} [(A_{2i-1} \setminus A_{2i}) \cup (B_{2i-1} \setminus B_{2i})] \right).$$

To complete the proof we determine the total number of trajectories in the sets A_{2i-1} , A_{2i} , B_{2i-1} , and B_{2i} for $i = 1, 2, \dots$. We demonstrate the method for the case of sets A_1 and A_2 . Every trajectory starting from the point $(0, 0)$

and meeting the line α corresponds to a new trajectory starting from the point $(0, 0)$ and coinciding with the original trajectory until the point where it meets the line α ; then the new trajectory is the mirror reflection of the old trajectory about the line α . The new trajectory ends at the point $(2n, 2c)$. The number of such trajectories is already determined above and it is equal to $\binom{2n}{n-c}$. Note that the cardinality of the set A_1 also is $\binom{2n}{n-c}$. If the original trajectory meets the line α first and then meets the line β , then the new trajectory meets the line $x = 3c$ (this is the mirror reflection of the line β).

To determine the cardinality of the set A_2 we introduce new trajectories as follows: the new trajectory coincides with the original one from the point $(0, 0)$ until it meets α , then it coincides with the first mirror reflection about the line α until meeting the line $x = 3c$, and, finally, the last part of the new trajectory is the second reflection about the line $x = 3c$ of the trajectory reflected first. Such new trajectories end at the point $(2n, 4c)$. The total number of such trajectories (thus, the cardinality of the set A_2) is equal to $\binom{2n}{n-2c}$.

A similar reasoning (using an appropriate number of reflections) proves that the cardinality of the set A_i is $\binom{2n}{n-ic}$. In the same way we find that the number of trajectories in the set B_i also is $\binom{2n}{n-ic}$. Since the sets $A_{2i-1} \setminus A_{2i}$ and $B_{2i-1} \setminus B_{2i}$ are disjoint and the terms in (3.3.7) also are disjoint, we obtain from (3.3.7) that the total number of trajectories in the set \mathfrak{A}_0 is

$$\binom{2n}{n} - 2 \sum_{i \geq 1} \left(\binom{2n}{n-(2i-1)c} - \binom{2n}{n-2ic} \right).$$

Thus equality (3.3.4) is proved. \square

REMARK 3.3.1. The process S_k , $k = 0, 1, \dots, 2n$, used in the proof of Theorem 3.3.1 is called a *random walk* on the axis. The method of evaluation of the number of trajectories in the sets A_i and B_i applied in the proof of Theorem 3.3.1 is well known in the theory of random walks as the *reflection method* [24, 47].

The Smirnov limit theorem. The following result describes the asymptotic behavior of the statistics $D_{n,n}^+$ and $D_{n,n}$ as $n \rightarrow \infty$.

THEOREM 3.3.2 (Smirnov). *If the distribution function*

$$P\{\xi < x/H\} = P\{\eta < x/H\}$$

is continuous, then

$$(3.3.8) \quad \lim_{n \rightarrow \infty} P \left\{ \sqrt{\frac{n}{2}} D_{n,n}^+ < z/H \right\} = \begin{cases} 0, & z \leq 0, \\ 1 - e^{-2z^2}, & z > 0, \end{cases}$$

$$(3.3.9) \quad \lim_{n \rightarrow \infty} P \left\{ \sqrt{\frac{n}{2}} D_{n,n} < z/H \right\} = K(z)$$

where $K(z)$ is the Kolmogorov distribution function defined by (3.1.5).

PROOF. Let $z > 0$ be a fixed number. Consider the ratio

$$I_k = \frac{\binom{2n}{n-kc}}{\binom{2n}{n}} = \frac{(n!)^2}{(n-kc)!(n+kc)!}$$

where k is a fixed constant independent of n . Using the Stirling formula

$$m! = \sqrt{2\pi m} m^m e^{-m} (1 + o(1)), \quad m \rightarrow \infty,$$

we obtain

$$(3.3.10) \quad I_k = \left(1 - \frac{kc}{n}\right)^{-n+kc} \left(1 + \frac{kc}{n}\right)^{-n-kc} (1 + o(1)).$$

Since $c =]z\sqrt{2n}[= z\sqrt{2n}(1 + o(1))$ as $n \rightarrow \infty$, we have $kc/n = az\sqrt{2/n}(1 + o(1))$. Then taking the logarithm of both sides of (3.3.10) and expanding the result into the Taylor series we get that

$$\ln I_k = -\frac{k^2 c^2}{n} + o(1) = -2k^2 z^2 + o(1), \quad n \rightarrow \infty.$$

Thus

$$(3.3.11) \quad I_k = e^{-2k^2 z^2} (1 + o(1)).$$

This equality for $k = 1$ proves relation (3.3.8).

Now we turn to relation (3.3.9). For given $z > 0$ and $\varepsilon > 0$ there is a number $N(\varepsilon, z) > 0$ such that

$$(3.3.12) \quad e^{-2N^2 z^2} < \frac{\varepsilon}{16}, \quad \left| \sum_{|k| > N} (-1)^k e^{-2k^2 z^2} \right| < \frac{\varepsilon}{4}$$

for $N > N(\varepsilon, z)$. Since

$$\binom{2n}{n - kc} > \binom{2n}{n - (k+1)c}, \quad k > 0,$$

we have

$$\left| \sum_{N < |k| \leq [n/c]} (-1)^k I_k \right| < 2I_N.$$

Taking into account (3.3.11) and the first inequality in (3.3.12) we obtain for sufficiently large n that

$$(3.3.13) \quad \left| \sum_{N < |k| \leq [n/c]} (-1)^k I_k \right| \leq 4e^{-2N^2 z^2} < \frac{\varepsilon}{4}.$$

Now we apply the second inequality in (3.3.12) and inequality (3.3.13) to prove that for large n

$$(3.3.14) \quad \left| \sum_{N < |k| \leq [n/c]} (-1)^k I_k - \sum_{|k| > N} (-1)^k e^{-2k^2 z^2} \right| < \frac{\varepsilon}{4} + \frac{\varepsilon}{4} = \frac{\varepsilon}{2}.$$

Further we deduce from (3.3.11) that for a fixed N and sufficiently large n

$$(3.3.15) \quad \left| \sum_{|k| \leq N} (-1)^k I_k - \sum_{|k| \leq N} (-1)^k e^{-2k^2 z^2} \right| < \frac{\varepsilon}{2}.$$

Thus (3.3.14) and (3.3.15) imply for fixed $z > 0$ and $\varepsilon > 0$ and sufficiently large n that

$$\left| \mathbb{P} \left\{ \sqrt{\frac{n}{2}} D_{n,n}^+ < z/H \right\} - K(z) \right| < \varepsilon.$$

Thus relation (3.3.9) is proved. \square

The Smirnov homogeneity test. The hypothesis H on the homogeneity, that is, $\mathbb{P}\{\xi < x\} = \mathbb{P}\{\eta < x\}$, can be tested in the case of continuous distribution functions by using either the statistic $D_{n,m}^+$ or the statistic $D_{n,m}$. As above we consider the case $m = n$. To construct a test one can use either Theorem 3.3.1 or Theorem 3.3.2. For the sake of simplicity we construct a test by using Theorem 3.3.2.

We consider a goodness-of-fit test for the hypothesis H on the homogeneity based on the statistic $D_{n,n}$. Let $\alpha > 0$ be a significance level and let $z(\alpha)$ be a solution of the equation $K(z) = 1 - \alpha$ with respect to z . Then relation (3.3.9) implies that for sufficiently large n

$$(3.3.16) \quad \mathbb{P} \left\{ \sqrt{\frac{n}{2}} D_{n,n} \geq z(\alpha)/H \right\} \approx 1 - K(z(\alpha)) = \alpha.$$

The test is defined by

$$\delta_n = I \left\{ \sqrt{\frac{n}{2}} D_{n,n} \geq z(\alpha) \right\}.$$

It follows from relation (3.3.16) that the level of the test δ_n is approximately equal to α for large n . The test δ_n is called the *Smirnov test* or *Smirnov homogeneity test* or *Smirnov goodness-of-fit test*. The Smirnov test rejects the hypothesis H if $\sqrt{n/2} D_{n,n} \geq z(\alpha)$, that is, if $D_{n,n} \geq z(\alpha)\sqrt{2/n}$. Otherwise there is no reason to reject the hypothesis H and it is accepted.

It is not hard to show that the test δ_n is consistent. Indeed, let $K_{F,G}$ be the hypothesis that $\mathbb{P}\{\xi < x/K_{F,G}\} = F(x)$ and $\mathbb{P}\{\eta < x/K_{F,G}\} = G(x)$ where $\sup_x |F(x) - G(x)| \neq 0$. Similarly to the proof of Theorem 3.1.2 we obtain from the Glivenko theorem that

$$\lim_{n \rightarrow \infty} \mathbb{P} \left\{ \sqrt{\frac{n}{2}} D_{n,n} < z(\alpha)/K_{F,G} \right\} = 0,$$

that is, the test δ_n is consistent.

In the same way one can construct a goodness-of-fit test for the hypothesis H based on the statistic $D_{n,n}^+$. Namely, we introduce a solution $z^+(\alpha)$ of the equation $e^{-2z^2} = \alpha$ with respect to z where $\alpha > 0$ is a given significance level. It follows from relation (3.3.8) that

$$\mathbb{P} \left\{ \sqrt{\frac{n}{2}} D_{n,n}^+ \geq z^+(\alpha)/H \right\} \approx e^{-2(z^+(\alpha))^2} = \alpha$$

for sufficiently large n . Thus the level of the test

$$\delta_n^+ = I \left(\sqrt{\frac{n}{2}} D_{n,n}^+ \geq z^+(\alpha) \right)$$

is approximately equal to α for sufficiently large n . The test δ_n^+ can also be used for the hypothesis H on the homogeneity and it is called the *Smirnov test*, too.

3.4. Other goodness-of-fit tests

We studied Kolmogorov, Smirnov, and Pearson tests in the preceding sections. Many other goodness-of-fit tests are well known in mathematical statistics. Some of them are considered in this section.

Symmetric tests. Let a random variable ξ be real valued and let

$$\xi^{(n)} = (\xi_1, \xi_2, \dots, \xi_n)$$

be a sample of size n . We split the real axis $\mathbf{R} = (-\infty, \infty)$ into r intervals

$$(3.4.1) \quad S_i = [y_{i-1}, y_i), \quad i = 1, 2, \dots, r, \quad y_0 = -\infty, \quad y_r = \infty.$$

Let the hypothesis H be such that $P\{\xi \in S_i/H\} = p_i = 1/r$ for all $i = 1, 2, \dots, r$.

Let $\nu = (\nu_1, \nu_2, \dots, \nu_r)$ be the vector whose coordinates equal the number of members of the sample $\xi^{(n)}$ belonging to intervals (3.4.1), that is, $\nu_i, i = 1, 2, \dots, r$, are evaluated by (3.2.2). Consider the class of statistics

$$(3.4.2) \quad \zeta_{n,r}(g) = \sum_{i=1}^n g(\nu_i)$$

where $g(x)$ is some real function defined for all nonnegative integer arguments $x = 0, 1, 2, \dots$. Since the random variables $\nu_1, \nu_2, \dots, \nu_n$ form a symmetric expression in (3.4.2), the statistic $\zeta_{n,r}(g)$ is called *symmetric* and tests based on statistic (3.4.2) are called *symmetric*.

If $g(x) = I_{\{k\}}(x)$ in (3.4.2) for some $k \in \{0, 1, 2, \dots, n\}$, then the statistic $\zeta_{n,r}(g)$ is equal to the number of intervals among S_1, S_2, \dots, S_r that contain exactly k members of the sample $\xi_1, \xi_2, \dots, \xi_n$. We denote this statistic by

$$(3.4.3) \quad \mu_k = \mu_k(n, r) = \sum_{i=1}^n I_{\{k\}}(\nu_i).$$

Thus we have a collection of symmetric statistics $\mu_0, \mu_1, \dots, \mu_n$ defined by (3.4.3). It is obvious that these statistics are such that

$$(3.4.4) \quad \sum_{k=0}^n \mu_k(n, r) = r, \quad \sum_{k=0}^n k \mu_k(n, r) = n.$$

Relation (3.4.2) can be rewritten in terms of the statistics $\mu_0, \mu_1, \dots, \mu_n$:

$$(3.4.5) \quad \zeta_{n,r}(g) = \sum_{k=0}^n g(k) \mu_k,$$

whence it follows that an arbitrary symmetric statistic is a linear combination of $\mu_0, \mu_1, \dots, \mu_n$. The converse is also true, namely any linear combination of $\mu_0, \mu_1, \dots, \mu_n$ is a symmetric statistic, since

$$\sum_{k=0}^n c_k \mu_k = \sum_{i=1}^r g(\nu_i)$$

where $g(k) = c_k, k = 0, 1, \dots, n$. Thus the *class of symmetric statistics coincides with the class of all linear combinations of $\mu_0, \mu_1, \dots, \mu_n$.*

Note that the statistic ζ_n defined by (3.2.4) and used in Section 3.2 to construct the Pearson test coincides with the symmetric statistic (3.4.2) for $g(x) = rx^2/n - x$ and if intervals (3.4.1) are equiprobable.

Empty boxes test. The symmetric test based on the statistic μ_0 is called the *empty boxes test*. Note that μ_0 is the number of intervals S_1, S_2, \dots, S_r that do not contain any member of the sample $\xi_1, \xi_2, \dots, \xi_n$.

Below we evaluate the first two moments of the statistic μ_0 if the hypothesis H_p is true where H_p is such that $P\{\xi \in S_i/H_p\} = p_i, i = 1, 2, \dots, r$. Here the numbers p_i are arbitrary and $p = (p_1, p_2, \dots, p_r)$. Consider the random variables $\eta_1, \eta_2, \dots, \eta_r$ such that $\eta_i = 1$ if the interval S_i does not contain any member of the sample ξ_1, \dots, ξ_n , and $\eta_i = 0$ otherwise. It is clear that $\mu_0 = \eta_1 + \dots + \eta_n$, whence

$$\begin{aligned} E\{\mu_0/H_p\} &= \sum_{i=1}^r E\{\eta_i/H_p\} = \sum_{i=1}^r P\{\eta_i = 1/H_p\}, \\ D\{\mu_0/H_p\} &= \sum_{i=1}^r D\{\eta_i/H_p\} + 2 \sum_{i < j} E\{(\eta_i - E\{\eta_i/H_p\})(\eta_j - E\{\eta_j/H_p\})/H_p\} \\ &= \sum_{i=1}^r P\{\eta_i = 1/H_p\}[1 - P\{\eta_i = 1/H_p\}] \\ &\quad + 2 \sum_{i < j} [P\{\eta_i = 1, \eta_j = 1/H_p\} - P\{\eta_i = 1/H_p\}P\{\eta_j = 1/H_p\}]. \end{aligned}$$

Since the random variables $\xi_1, \xi_2, \dots, \xi_n$ are independent and identically distributed and

$$\{\eta_i = 1\} = \bigcap_{j=1}^n \{\xi_j \notin S_i\}, \quad \{\eta_i = 1, \eta_j = 1\} = \bigcap_{k=1}^n \{\xi_k \notin S_i, \xi_k \notin S_j\},$$

we have

$$P\{\eta_i = 1/H_p\} = (1 - p_i)^n, \quad P\{\eta_i = 1, \eta_j = 1/H_p\} = (1 - p_i - p_j)^n.$$

Thus

$$(3.4.6) \quad E\{\mu_0/H_p\} = \sum_{i=1}^r (1 - p_i)^n,$$

$$(3.4.7) \quad D\{\mu_0/H_p\} = 2 \sum_{i < j} (1 - p_i - p_j)^n + E\{\mu_0/H_p\} - (E\{\mu_0/H_p\})^2.$$

It is easy to show that $E\{\mu_0/H_p\}$ as a function of $p = (p_1, p_2, \dots, p_r)$ attains its minimum at $p = p^0 = (p_1^0, p_2^0, \dots, p_r^0)$ where $p_1^0 = p_2^0 = \dots = p_r^0 = 1/r$. Equalities (3.4.6) and (3.4.7) for $p = p_0$ become of the form

$$\begin{aligned} E\{\mu_0/H\} &= r(1 - 1/r)^n, \\ D\{\mu_0/H\} &= r(r - 1)(1 - 2/r)^n + r(1 - 1/r)^n - r^2(1 - 1/r)^{2n} \end{aligned}$$

(note that $H = H_{p^0}$ in this case).

This shows that if the hypothesis H is not true, that is, if not all probabilities for intervals S_1, S_2, \dots, S_r are equal to $1/r$, the statistic μ_0 tends to increase, since $E\{\mu_0/H_p\} > E\{\mu_0/H_{p^0}\}$ for $p \neq p_0$. Thus large values of μ_0 lead to the rejection

of the hypothesis H . The empty boxes test for testing the hypothesis H is then as follows: the hypothesis H is rejected if $\mu_0 \geq t_\alpha(n, r)$, while it is accepted otherwise. The number $t_\alpha(n, r)$ can be evaluated by using the distribution of the statistic μ_0 given the hypothesis H is true. However this distribution is complicated and the limit results for μ_0 are often used instead (see [26, 27]).

The empty boxes test has the same disadvantage as the Pearson test (see Remark 3.2.2). To avoid this disadvantage one should assume that r is large enough, that is, $r \rightarrow \infty$ as $n \rightarrow \infty$. This assumption allows one to apply the empty boxes test for testing the simple null hypothesis H_0 that $P\{\xi < x/H_0\} = F(x)$ and $P\{\xi \in S_i/H_0\} = p_i = 1/2$, $i = 1, 2, \dots, r$, where $F(x)$ is a given distribution function. More details are given in [26, 27].

General symmetric tests. The empty boxes test is based on the statistic μ_0 that does not contain all the information available from the sample $\xi^{(n)}$. It is clear that the statistic

$$\sum_{k=0}^n c_k \mu_k$$

contains more statistical information than μ_0 where $c_1 > 0, \dots, c_n > 0$ are some weights. It is natural to choose the weights c_1, c_2, \dots, c_n such that the test is the most powerful in the class of all symmetric tests. The theory of symmetric tests is considered in [27]. Note also that the Pearson test has the maximal asymptotic power among all symmetric tests under appropriate conditions (see [26, 27]).

Tests of the homogeneity. We considered the Smirnov test in Section 3.3, however there exist many other tests of the homogeneity. We briefly discuss some of them below.

Chi-square test of the homogeneity. Consider s independent samples of independent observations. Assume that samples $1, 2, \dots, s$ contain n_1, n_2, \dots, n_s members, respectively. We assume that an attribute ξ_i is checked in the sample i and denote by $\xi_{i1}, \xi_{i2}, \dots, \xi_{in_i}$ the observations of the attribute, so that we deal with independent random variables ξ_{ij} , $j = 1, 2, \dots, n_i$, $i = 1, 2, \dots, s$. Further let the results of every observation fall into r groups according to their values. Namely the domain X of the attribute ξ_i fall into parts S_1, S_2, \dots, S_r such that $\bigcup_{k=1}^r S_k = X$ and $S_k \cap S_l = \emptyset$, $k \neq l$. Let ν_{ij} be the number of members of the sample $\xi_{i1}, \xi_{i2}, \dots, \xi_{in_i}$ belonging to S_j , $j = 1, 2, \dots, r$. Put $p_{ij} = P\{\xi_i \in S_j\}$, $i = 1, 2, \dots, s$, $j = 1, 2, \dots, r$. Consider the hypothesis $H: (p_{i1}, \dots, p_{ir}) = (p_{11}, \dots, p_{1r})$ for all $i = 1, 2, \dots, s$. We form the vector $p = (p_1, p_2, \dots, p_r)$ and note that $p_1 + \dots + p_r = 1$.

The vector p is, generally speaking, unknown. However if the vector p is known, then $E\{\nu_{ij}/H\} = n_i p_j$ and

$$(3.4.8) \quad \zeta_n(p) = \sum_{j=1}^r \sum_{i=1}^s \frac{(\nu_{ij} - n_i p_j)^2}{n_i p_j}$$

can be viewed as the measure of disagreement between the data and the null hypothesis (cf. relation (3.2.4)). Since p_1, p_2, \dots, p_r are unknown in general, we modify the χ^2 method and substitute the estimators

$$(3.4.9) \quad \hat{p}_j = \frac{\nu_{.j}}{n}, \quad j = 1, 2, \dots, r,$$

for p_1, \dots, p_r where

$$(3.4.10) \quad \nu_{\cdot j} = \sum_{i=1}^s \nu_{ij}, \quad n = \sum_{i=1}^s n_i = \sum_{i=1}^s \sum_{j=1}^r \nu_{ij}.$$

Substituting (3.4.9) into (3.4.8) and taking into account (3.4.10) we get the statistic

$$(3.4.11) \quad \hat{\zeta}_n = \zeta_n(\hat{p}) = n \sum_{j=1}^r \sum_{i=1}^s \frac{(\nu_{ij} - n_i \nu_{\cdot j} / n)^2}{n_i \nu_{\cdot j}} = n \left(\sum_{j=1}^r \sum_{i=1}^s \frac{\nu_{ij}^2}{n_i \nu_{\cdot j}} - 1 \right)$$

where $\hat{p} = (\hat{p}_1, \dots, \hat{p}_r)$. Note that

$$(3.4.12) \quad \mathcal{L} \left(\hat{\zeta}_n | H \right) \xrightarrow{w} \chi^2((r-1)(s-1)), \quad n \rightarrow \infty$$

(see [14], §30.6).

Based on the statistic $\hat{\zeta}_n$ and using the limit relation (3.4.12) we construct the goodness-of-fit test for the hypothesis H . Some examples of applications of this test can be found in [14, 26].

Empty blocks test. Let $\xi^{(n)} = (\xi_1, \dots, \xi_n)$ and $\eta^{(m)} = (\eta_1, \dots, \eta_m)$ be two independent samples of independent observations of random variables ξ and η , respectively. Assume that both ξ and η have continuous distribution functions. The hypothesis H is that $P\{\xi < x/H\} = P\{\eta < x/H\}$ for all $x \in (-\infty, \infty)$. Consider the order statistics $\zeta_{n,1} \leq \zeta_{n,2} \leq \dots \leq \zeta_{n,n}$ related to the sample $\xi^{(n)}$. These statistics split the set $(-\infty, \infty)$ into the intervals

$$S_i = (\zeta_{n,i-1}, \zeta_{n,i}), \quad i = 1, 2, \dots, n+1,$$

where we put $\zeta_{n,0} = -\infty$ and $\zeta_{n,n+1} = \infty$. These intervals are called *sampling blocks*. Consider those sampling blocks that contain exactly r random variables of $\eta_1, \eta_2, \dots, \eta_m$ and let $s_r = s_r(n, m)$ be the number of such blocks, $r = 0, 1, 2, \dots, m$. Every linear combination

$$S_l(n, m) = \sum_{r=0}^l c_r s_r(n, m), \quad l = 0, 1, 2, \dots, m,$$

can be viewed as a test for the goodness-of-fit test of the hypothesis H where c_0, c_1, \dots, c_l are some positive numbers. The test corresponding to the case $l = 0$ is called the *empty blocks test*. The number of blocks that do not contain any observation of the second sample is the test statistic (denoted by s_0) in this case. The following result on the asymptotic distribution of the statistic s_0 is crucial for constructing this test: if $n, m \rightarrow \infty$ such that $m/n \rightarrow \rho \in (0, \infty)$, then

$$\mathcal{L} \left(\left(\frac{(1+\rho)^3}{n\rho^2} \right)^{1/2} \left(s_0(n, m) - \frac{n}{1+\rho} \right) \middle| H \right) \xrightarrow{w} \mathcal{N}(0, 1)$$

(see [53]). It is proved in [53] that the empty blocks test against the alternative $P\{\xi < x\} = F_1(x) \neq F_2(x) = P\{\eta < x\}$ is consistent if the derivative $g(u)$ of the function $F_2(F_1^{-1}(u))$, $u \in [0, 1]$, differs from 1 on a nonzero Lebesgue set. If the hypothesis H is true, then $g(u) \equiv 1$, $u \in [0, 1]$. The class of such alternative hypotheses is denoted by H^* .

Test of series. As in the case of empty blocks test we deal with two independent samples $\xi^{(n)} = (\xi_1, \dots, \xi_n)$ and $\eta^{(m)} = (\eta_1, \dots, \eta_m)$ that are independent observations of random variables ξ and η , respectively. We also assume that ξ and η have continuous distribution functions. Consider the hypothesis H that

$$P\{\xi < x\} = P\{\eta < x\} \quad \text{for all } x \in (-\infty, \infty).$$

It is a quite interesting case when the alternative hypothesis H_{F_1, F_2} is such that $P\{\xi < x/H_{F_1, F_2}\} = F_1$, $P\{\eta < x/H_{F_1, F_2}\} = F_2$, and $F_1(x) > F_2(x)$ for all x . The random variable η in this case is stochastically bigger than the random variable ξ , since for all x the random variable η exceeds x with a larger probability than ξ does. The test detecting a disagreement between the data and the hypothesis H can be constructed as follows.

First we merge the samples $\xi^{(n)}$ and $\eta^{(m)}$ and obtain the sample

$$\zeta^{n+m} = (\xi_1, \dots, \xi_n, \eta_1, \dots, \eta_m)$$

of size $n + m$. Then we construct the order statistics for the sample ζ^{n+m} . Finally, in the sequence of order statistics, we substitute the symbol C for all members of the sample $\xi^{(n)}$ and the symbol \bar{C} for all members of the sample $\eta^{(m)}$. As a result we get a sequence of n symbols C and m symbols \bar{C} . The total number of such sequences is $\binom{n+m}{n}$. It is clear that if H is true, then all such sequences are equiprobable (the proof is the same as that in the case of Theorem 3.3.1; also see the proof of Theorem 14.3.1 in [53]). If the alternative hypothesis is H_{F_1, F_2} where $F_1 > F_2$, then it is more likely that symbols \bar{C} appear far away from the origin of the sequence. The measure of the displacement of the symbols \bar{C} to the right can be characterized by the statistic $W(n, m)$ which is the number of series of symbols C and \bar{C} . Any sequence of symbols C or \bar{C} is called a *series*. The number of series is small if the symbols C (or \bar{C}) are grouped in a specified place of the sequence. Thus the critical set for testing the hypothesis H can be taken in the form $\{W(n, m) \leq t_\alpha(n, m)\}$ where $t_\alpha(n, m)$ is a certain number defined by the level α . The test related to such a critical set was proposed by Wald and Wolfowitz in 1940 and is called the *test of series*.

The following result on the limit behavior of the statistic $W(n, m)$ is useful for the evaluation of $t_\alpha(n, m)$: if $n, m \rightarrow \infty$ such that $m/n \rightarrow \rho \in (0, \infty)$, then

$$\mathcal{L} \left(\left(\frac{(1+\rho)^3}{4n\rho^2} \right)^{1/2} \left(W(n, m) - \frac{2n\rho}{1+\rho} \right) | H \right) \xrightarrow{w} \mathcal{N}(0, 1)$$

(see [53]).

Wald and Wolfowitz proved that the test of series is consistent if the alternative hypothesis belongs to the class H^* . The test of series is discussed in detail in [53].

Rank tests. Most nonparametric methods use observations ranked in order of their magnitude. Statistics constructed from ranks of observations are called *rank statistics*. Tests based on rank statistics are called *rank tests*.

Let $\xi^{(n)} = (\xi_1, \dots, \xi_n)$ and $\eta^{(m)} = (\eta_1, \dots, \eta_m)$ be two independent samples of independent observations of random variables ξ and η whose distribution functions are continuous. Using the samples $\xi^{(n)}$ and $\eta^{(m)}$ we want to test the hypothesis H that $P\{\xi < x\} = P\{\eta < x\}$ for all $x \in (-\infty, \infty)$. We merge the samples $\xi^{(n)}$ and $\eta^{(m)}$ and obtain the sample $\zeta^{n+m} = (\xi_1, \dots, \xi_n, \eta_1, \dots, \eta_m)$. Then we construct the sequence of order statistics $\zeta_{N,1} \leq \zeta_{N,2} \leq \dots \leq \zeta_{N,N}$, $N = n + m$,

from the sample ζ^{n+m} . Let R_i be the index in this sequence corresponding to the member ξ_i of the sample $\xi^{(n)}$, that is, $\xi_i = \zeta_{N, R_i}$, $i = 1, 2, \dots, n$. Thus we deal with ranks R_1, R_2, \dots, R_n of the observations $\xi_1, \xi_2, \dots, \xi_n$. Consider the statistic $T = R_1 + \dots + R_n$, that is, T is the sum of indices of members of the first sample ξ_1, \dots, ξ_n in the sequence of order statistics constructed from the merged sample. The test based on the rank statistic T was proposed by Wilcoxon in 1945 for the case of identical sizes of samples ($n = m$) and is called the *Wilcoxon rank sum test*.

Consider the random variable Z_{ij} that equals 1 if $\xi_i < \eta_j$ and 0 otherwise. Put

$$(3.4.13) \quad U = U(n, m) = \sum_{i=1}^n \sum_{j=1}^m Z_{ij}.$$

It is clear that U is the total number of the cases in the merged sample where members of the sample $\xi^{(n)}$ precede members of the sample $\eta^{(m)}$. The test for testing the hypothesis H on the homogeneity based on the statistic U was studied by Mann and Whitney in 1947 and is called the *Mann-Whitney U -test*.

One can show that $T + U = nm + n(n+1)/2$. Thus the Wilcoxon rank sum test and Mann-Whitney U -test are equivalent.

Let H_{F_1, F_2} be the hypothesis that $P\{\xi < x\} = F_1(x)$ and $P\{\eta < x\} = F_2(x)$ where $F_1(x)$ and $F_2(x)$ are continuous distribution functions. It follows from (3.4.13) that

$$E\{U/H_{F_1, F_2}\} = nmE\{Z_{11}/H_{F_1, F_2}\} = nma$$

where

$$a = E\{\xi_1 < \eta_1/H_{F_1, F_2}\} = \int_{-\infty}^{\infty} F_1(x) dF_2(x).$$

If the hypothesis H is true, that is, if $F_1(x) \equiv F_2(x)$, then $a = 1/2$. Similarly

$$D\{U/H_{F_1, F_2}\} = nm[a + (n-1)b + (m-1)c - (n+m-1)a^2]$$

where

$$b = \int_{-\infty}^{\infty} F_1^2(x) dF_2(x), \quad c = \int_{-\infty}^{\infty} (1 - F_2(x))^2 dF_1(x).$$

If the hypothesis H is true, then $b = c = 1/3$ and thus $D\{U/H\} = nm(n+m+1)/12$. It is known that

$$(3.4.14) \quad \mathcal{L} \left(\left(\frac{nm(n+m+1)}{12} \right)^{-1/2} \left(U(n, m) - \frac{nm}{2} \right) \middle| H \right) \xrightarrow{w} \mathcal{N}(0, 1)$$

as $n, m \rightarrow \infty$

Relation (3.4.14) is useful for the Wilcoxon test of homogeneity (and for the Mann-Whitney test, too). The critical set for this test depends on the alternative hypothesis H_{F_1, F_2} and especially on the value of a : either $a < 1/2$ or $a > 1/2$ or $a = 1/2$. More details on the Wilcoxon and Mann-Whitney tests are given in [53]. The general theory of rank tests as well as examples of various rank tests is presented in [22].

Tests of independence. Below we consider a couple of tests for testing the hypothesis that two random variables ξ and η are independent. The statistical inference is based on independent observations $(\xi_1, \eta_1), (\xi_2, \eta_2), \dots, (\xi_n, \eta_n)$ of the vector (ξ, η) . If $F_{(\xi, \eta)}(x, y)$ is the distribution function of the vector (ξ, η) , then the hypothesis H is that $F_{(\xi, \eta)}(x, y) = F_\xi(x)F_\eta(y)$ where $F_\xi(x)$ and $F_\eta(y)$ are

distribution functions of the random variables ξ and η , respectively. Note that the distribution functions $F_{(\xi,\eta)}(x, y)$, $F_{\xi}(x)$, and $F_{\eta}(y)$ are unknown.

Chi-square test of independence. Let X and Y be the sets of values of random variables ξ and η , respectively. Consider the following partitions of the sets X and Y :

$$X = \bigcup_{i=1}^r S_i^{(1)} \quad \text{and} \quad Y = \bigcup_{i=1}^s S_i^{(2)}$$

where

$$S_i^{(1)} \cap S_j^{(1)} = \emptyset, \quad i \neq j, \quad \text{and} \quad S_k^{(2)} \cap S_l^{(2)} = \emptyset, \quad k \neq l.$$

These partitions generate the partition of the set $X \times Y$:

$$X \times Y = \bigcup_{i=1}^r \bigcup_{j=1}^s (S_i^{(1)} \times S_j^{(2)})$$

where

$$(S_i^{(1)} \times S_j^{(2)}) \cap (S_k^{(1)} \times S_l^{(2)}) = \emptyset, \quad (i, j) \neq (k, l).$$

Here $S_i^{(1)} \times S_j^{(2)} = \{(x, y): x \in S_i^{(1)}, y \in S_j^{(2)}\}$ are rectangles in the set $X \times Y$.

The hypothesis of the independence H means in this case that $p_{ij} = p_i q_j$ for all $i = 1, 2, \dots, r$ and $j = 1, 2, \dots, s$ where

$$p_{ij} = P \left\{ \xi \in S_i^{(1)}, \eta \in S_j^{(2)} / H \right\}, \quad p_i = P \left\{ \xi \in S_i^{(1)} / H \right\},$$

and

$$q_j = P \left\{ \eta \in S_j^{(2)} / H \right\}.$$

Denote by ν_{ij} the total number of observations (ξ_k, η_k) , $k = 1, \dots, n$, belonging to the set $(S_i^{(1)} \times S_j^{(2)})$, so that $\sum_{i=1}^r \sum_{j=1}^s \nu_{ij} = n$. If the probabilities p_i and q_j are known, then one can use the statistic

$$(3.4.15) \quad \zeta_n(p, q) = \sum_{i=1}^r \sum_{j=1}^s \frac{(\nu_{ij} - np_i q_j)^2}{np_i q_j}$$

to test the hypothesis H . However these probabilities are usually unknown and thus we use their estimators

$$(3.4.16) \quad \hat{p}_i = \frac{\nu_{i.}}{n}, \quad i = 1, \dots, r, \quad \text{and} \quad \hat{q}_j = \frac{\nu_{.j}}{n}, \quad j = 1, \dots, s,$$

where

$$\nu_{i.} = \sum_{j=1}^s \nu_{ij}, \quad \nu_{.j} = \sum_{i=1}^r \nu_{ij}.$$

By substituting (3.4.16) into (3.4.15) instead of the probabilities p_i and q_j we obtain the statistic

$$(3.4.17) \quad \hat{\zeta}_n = \zeta_n(\hat{p}, \hat{q}) = n \sum_{i=1}^r \sum_{j=1}^s \frac{(\nu_{ij} - \nu_{i.} \nu_{.j} / n)^2}{\nu_{i.} \nu_{.j}} = n \left(\sum_{i,j} \frac{\nu_{ij}^2}{\nu_{i.} \nu_{.j}} - 1 \right).$$

The Pearson theorem implies that

$$(3.4.18) \quad \mathcal{L} \left(\hat{\zeta}_n | H \right) \xrightarrow{w} \chi^2((r-1)(s-1)), \quad n \rightarrow \infty.$$

Now we use statistic (3.4.17) and construct the goodness-of-fit test for testing the hypothesis H of the independence of two random variables. This is a standard procedure using the limit relation (3.4.18) and leading to the so-called χ^2 test of independence. More details on the χ^2 test of independence are given in [14, 26].

Spearman test. Let R_i be the rank of the member ξ_i in the sequence of order statistics $\xi_{n1} \leq \xi_{n2} \leq \dots \leq \xi_{nn}$ constructed from the sample $(\xi_1, \xi_2, \dots, \xi_n)$. Note that the sample (ξ_1, \dots, ξ_n) consists of the first components of the members of the sample $(\xi_1, \eta_1), (\xi_2, \eta_2), \dots, (\xi_n, \eta_n)$. Similarly, let S_i be the rank of the member η_i in the sequence of order statistics $\eta_{n1} \leq \eta_{n2} \leq \dots \leq \eta_{nn}$ constructed from the sample (η_1, \dots, η_n) . The sample $(\xi_1, \eta_1), \dots, (\xi_n, \eta_n)$ thus generates the set of pairs of ranks $(R_1, S_1), \dots, (R_n, S_n)$. We rearrange these pairs in ascending order with respect to their first component and denote the rearranged set of pairs by $(1, T_1), (2, T_2), \dots, (n, T_n)$.

Consider the rank statistic

$$(3.4.19) \quad \rho = \frac{\sum_{i=1}^n (R_i - \bar{R})(S_i - \bar{S})}{\left(\sum_{i=1}^n (R_i - \bar{R})^2 \sum_{i=1}^n (S_i - \bar{S})^2 \right)^{-1/2}}$$

which is the coefficient of correlation between two sets of ranks (R_1, \dots, R_n) and (S_1, \dots, S_n) where

$$\bar{R} = \frac{1}{n} \sum_{i=1}^n R_i, \quad \bar{S} = \frac{1}{n} \sum_{i=1}^n S_i.$$

The statistic ρ defined by (3.4.19) is called the *Spearman rank correlation coefficient*. The test based on the statistic ρ is called the *Spearman test*.

Since (R_1, \dots, R_n) and (S_1, \dots, S_n) are certain permutations of the numbers $(1, 2, \dots, n)$, we have

$$(3.4.20) \quad \bar{R} = \bar{S} = \frac{1}{n} \sum_{i=1}^n i = \frac{n+1}{2},$$

$$(3.4.21) \quad \sum_{i=1}^n (R_i - \bar{R})^2 = \sum_{i=1}^n (S_i - \bar{S})^2 = \sum_{i=1}^n i^2 - n \left(\frac{n+1}{2} \right)^2 = \frac{n(n^2-1)}{12}.$$

Combining (3.4.19)–(3.4.21) we get

$$(3.4.22) \quad \begin{aligned} \rho &= \frac{12}{n(n^2-1)} \sum_{i=1}^n \left(R_i - \frac{n+1}{2} \right) \left(S_i - \frac{n+1}{2} \right) \\ &= \frac{12}{n(n^2-1)} \sum_{i=1}^n \left(i - \frac{n+1}{2} \right) \left(T_i - \frac{n+1}{2} \right). \end{aligned}$$

There is another useful formula for the Spearman coefficient, namely

$$(3.4.23) \quad \rho = 1 - \frac{6}{n(n^2-1)} \sum_{i=1}^n (R_i - S_i)^2 = 1 - \frac{6}{n(n^2-1)} \sum_{i=1}^n (i - T_i)^2.$$

It is straightforward to check that (3.4.22) and (3.4.23) are equivalent.

Let two distribution functions $P\{\xi < x/H\}$ and $P\{\eta < x/H\}$ be continuous. Since all $n!$ permutations (T_1, T_2, \dots, T_n) of the numbers $(1, 2, \dots, n)$ are equiprobable, we have

$$E\{T_i/H\} = \sum_{j=1}^n j \frac{(n-1)!}{n!} = \frac{n+1}{2},$$

whence we obtain by (3.4.23) that

$$E\{\rho/H\} = 1 - \frac{6}{n(n^2-1)} \left[2 \sum_{i=1}^n i^2 - 2 \sum_{i=1}^n i E\{T_i/H\} \right] = 0.$$

Similarly

$$D\{\rho/H\} = \frac{1}{n-1}.$$

If the ranks coincide, that is, if $R_i = S_i$, $i = 1, 2, \dots, n$, then $\rho = 1$, while if the ranks are opposite, that is, if $T_i = n - i + 1$, $i = 1, 2, \dots, n$, then $\rho = -1$. In general, $-1 \leq \rho \leq 1$. If ρ is close to either -1 or 1 , then the hypothesis H is false. Thus the critical set of the Spearman test is $\{|\rho| \geq t_\alpha(n)\}$ where $t_\alpha(n)$ is defined for the level α by using the distribution of the statistic ρ . One approach to evaluate $t_\alpha(n)$ is to use for $n = 2, 3, \dots, 30$ the tables of the distribution of the statistic ρ (see references in [22]). Another approach is based on the limit relation

$$\mathcal{L}(\sqrt{n}\rho|H) \xrightarrow{w} \mathcal{N}(0, 1), \quad n \rightarrow \infty,$$

whence one can also find $t_\alpha(n)$ (see [29], §37.28).

Kendall test. Another rank test of independence, called the *Kendall test*, is based on the statistic

$$\tau = \frac{1}{c_n^2} \sum_{i < j} \text{sign}(T_j - T_i)$$

where $\text{sign}(a) = 1$ for $a > 0$ and $\text{sign}(a) = -1$ for $a < 0$; c_n is a certain constant. The statistic τ is called the *Kendall statistic*. It is known that

$$(3.4.24) \quad E\{\tau/H\} = 0, \quad D\{\tau/H\} = \frac{2(2n+5)}{9n(n-1)},$$

$$\mathcal{L}\left(\frac{3}{2}\sqrt{n}\tau|H\right) \xrightarrow{w} \mathcal{N}(0, 1), \quad n \rightarrow \infty$$

(see [29]). The critical set for the Kendall test is $\{|\tau| \geq t_\alpha(n)\}$ where the constant $t_\alpha(n)$ can be found for the level α from relation (3.4.24). It is shown in [29] that the Spearman and Kendall tests are asymptotically equivalent as $n \rightarrow \infty$, since the coefficient of the correlation between the statistics ρ and τ approaches 1 as $n \rightarrow \infty$.

Other rank tests of independence can be found in [22].

The von Mises–Smirnov test. Let $\xi^{(n)} = (\xi_1, \dots, \xi_n)$ be a sample of size n . We treat ξ_1, \dots, ξ_n as independent observations of a random variable ξ . Let the hypothesis H be such that $P\{\xi < x/H\} = F(x)$ for all $x \in (-\infty, \infty)$. One possible approach to test the hypothesis H is to use the Kolmogorov test. Consider another goodness-of-fit test for testing the hypothesis H based on the statistic

$$(3.4.25) \quad \omega_n^2 = \int_{-\infty}^{\infty} (F_n(x) - F(x))^2 dF(x)$$

where $F_n(x)$ is the empirical distribution function constructed from the sample $\xi^{(n)}$. It is proved in [9] that if $F(x)$ is a continuous function, then

$$(3.4.26) \quad \lim_{n \rightarrow \infty} \mathbb{P}\{n\omega_n^2 < x/H\} = \Omega(x) = \mathbb{P}\left\{\int_0^1 (w^0(t))^2 dt < x\right\}$$

where $w^0(t)$, $0 \leq t \leq 1$, is the Brownian bridge. We follow standard procedure to construct the goodness-of-fit test for testing the hypothesis H based on statistic (3.4.25) and on the limit relation (3.4.26). This procedure leads to the so-called *von Mises-Smirnov test*. Sometimes it also is called the ω^2 test.

The distribution function of $\Omega(x)$ is complicated. However one can use the tables of values of the function $\Omega(x)$ (see [6]). Note that the distribution of the statistic ω_n^2 does not depend on the function $F(x)$ and moreover

$$(3.4.27) \quad \mathbb{E}\{\omega_n^2/H\} = \frac{1}{6n}, \quad \mathbb{D}\{\omega_n^2/H\} = \frac{4n-3}{180n^3}$$

(see [14]). More details about the ω^2 test can be found in [40].

Moran test. As in the preceding section let ξ^n be a sample and let the hypothesis H be such that $\mathbb{P}\{\xi < x/H\} = F(x)$ for all x where $F(x)$ is a continuous distribution function. Consider the statistic

$$(3.4.28) \quad M_n = \sum_{k=0}^n [F(\zeta_{n,k+1}) - F(\zeta_{n,k})]^2$$

where $\zeta_{n,k}$, $k = 1, 2, \dots, n$, are order statistics constructed from the sample $\xi^{(n)}$,

$$\zeta_{n,1} \leq \zeta_{n,2} \leq \dots \leq \zeta_{n,n},$$

and $F(\zeta_{n,0}) = 0$ and $F(\zeta_{n,n+1}) = 1$. The test based on the statistic (3.4.28) is called the *Moran test*. It rejects the hypothesis H if $M_n > c_n(\alpha)$ where $c_n(\alpha)$ is a constant determined by a level α and the distribution of the statistic M_n .

Since the random variable $F(\zeta_{n,k})$ is uniformly distributed on the interval $[0, 1]$, the distribution of the statistic M_n does not depend on the function $F(x)$. Thus one can consider the test based on the statistic

$$(3.4.29) \quad M_n = \sum_{k=0}^n (\zeta_{n,k+1} - \zeta_{n,k})^2.$$

The aim of the Moran test is to test whether the distribution of the random variable ξ is uniform on the interval $[0, 1]$. The number $c_n(\alpha)$ for the Moran test can be evaluated by applying the following assertion: *if the distribution function $F(x)$ is continuous, then*

$$\mathcal{L}\left(\sqrt{n}\left(\frac{nM_n}{2} - 1\right) \mid H\right) \xrightarrow{w} \mathcal{N}(0, 1)$$

as $n \rightarrow \infty$ (see [9]).

It is also proved in [9] that the Moran test is consistent. On the other hand, the Moran test does not distinguish close hypotheses (see [9]).

To conclude this section we note that many other goodness-of-fit tests are known and an extensive literature is devoted to them. Some references are given above. We also mention the book [50] on the nonparametric statistics where special attention is given to the goodness-of-fit tests.

Sequential Tests

4.1. Bayes sequential tests of hypotheses

Setting of the problem. Let (Θ, \mathcal{E}) be a measurable space and let $\xi = (\xi_1, \xi_2, \dots)$ be a sequence of independent identically distributed random variables whose distribution P_θ depends on a parameter $\theta \in \Theta$. Let $\{\Theta_1, \Theta_2, \dots, \Theta_m\}$ be a partition of the space Θ , that is, $\Theta = \bigcup_{i=1}^m \Theta_i$ and $\Theta_i \cap \Theta_j = \emptyset$, $i \neq j$. Assume that m loss functions $A_i(\theta)$, $i = 1, 2, \dots, m$, are defined on Θ . The parameter θ determining the distribution P_θ is chosen in Θ according to the a priori distribution \mathbf{Q} on (Θ, \mathcal{E}) .

Consider the problem of testing m hypotheses $H_i: \theta \in \Theta_i$, $i = 1, 2, \dots, m$, by the sequence of random variables ξ_1, ξ_2, \dots . The difference between *sequential tests* and tests with a fixed size of the sample is as follows. In the case of a sequential test, a statistician is free to decide at any time whether it is necessary to terminate the sampling. More precisely, a statistician may terminate the sampling at any time n and decide to accept a certain hypothesis H_i on the basis of the sample $\xi_1, \xi_2, \dots, \xi_n$. When deciding to terminate the sampling, the statistician takes into account the cost per observation on the one hand and the amount of information about the parameter θ available in the next observation on the other hand. Let $K_n(\xi^{(n)})$ be the *expenses caused by the observation* $\xi^{(n)} = (\xi_1, \xi_2, \dots, \xi_n)$. Assume that the expenses are such that

- 1) $K_{n+1}(\xi^{(n+1)}) \geq K_n(\xi^{(n)})$ a.s. for all $n \geq 0$;
- 2) $\lim_{n \rightarrow \infty} K_n(\xi^{(n)}) = \infty$ a.s.

Every sequential test is determined by two components: a stopping rule and a decision rule. Denote by $s(\mathbf{Q}, \xi)$ a *stopping rule* and let $\nu(s)$ be a random variable denoting the size of the sample if the stopping rule $s(\mathbf{Q}, \xi)$ is applied. A *decision rule* is denoted by $\delta(\mathbf{Q}, \xi)$. Our assumption is that the decision depends only on $\xi^{(n)} = (\xi_1, \xi_2, \dots, \xi_n)$ if $\nu(s) = n$. Consider the following set of sequential tests.

1. After a statistician terminates the sampling he applies a decision rule $\delta(\mathbf{Q}, \xi)$ that is assumed to be Bayes under the a priori distribution \mathbf{Q} . This means that if \mathbf{Q} and $\xi^{(n)} = (\xi_1, \xi_2, \dots, \xi_n)$ are given and $\nu(s) = n$, then

$$\delta(\mathbf{Q}, \xi) = \left(q_1^\delta(\xi^{(n)}), q_2^\delta(\xi^{(n)}), \dots, q_n^\delta(\xi^{(n)}) \right)$$

where the functions $q_i^\delta(\xi^{(n)})$ are defined by (1.3.42) with $\xi^{(n)}$ instead of x . We also assume that the measure P_θ is absolutely continuous with respect to some σ -finite measure μ and that its density is $p(x; \theta)$. Then the distribution of the sample $\xi^{(n)}$ is absolutely continuous with respect to the measure $\mu^n = \mu \times \mu \times \dots \times \mu$ (n times) and its density is $p_n(x^{(n)}; \theta) = \prod_{i=1}^n p(x_i; \theta)$, $x^{(n)} = (x_1, x_2, \dots, x_n)$. Thus $p_n(x^{(n)}, t)$ should be used in (1.3.40) and (1.3.41) instead of $p(x; t)$.

2. The class of stopping rules $s(\mathbf{Q}, \xi)$ is the collection of all randomized rules that can be described as follows. For any $n \geq 1$, denote by \mathcal{F}_n the σ -algebra generated by the vector $(\theta, \xi_1, \xi_2, \dots, \xi_n)$ and let $\mathcal{F}_0 = \mathcal{C}$. It is clear that

$$\mathcal{F}_0 \subset \mathcal{F}_1 \subset \mathcal{F}_2 \subset \dots$$

Let \mathcal{S} denote the class of all sequences $s = (s_1, s_2, \dots)$ such that $s_1 = s_1(\mathbf{Q})$ and $s_{n+1} = s_{n+1}(\mathbf{Q}, \xi_1, \xi_2, \dots, \xi_n)$ is an \mathcal{F}_n -measurable random variable,

$$0 \leq s_{n+1}(\mathbf{Q}, \xi_1, \xi_2, \dots, \xi_n) \leq 1, \quad n \geq 1.$$

We treat s_{n+1} as the probability that the random variable ξ_{n+1} occurs in the sampling. Set

$$J_{n+1} = J_{n+1}(\mathbf{Q}, \xi^{(n)}) = \begin{cases} 1, & \text{with probability } s_{n+1}(\mathbf{Q}, \xi^{(n)}), \\ 0, & \text{with probability } 1 - s_{n+1}(\mathbf{Q}, \xi^{(n)}). \end{cases}$$

Following the stopping rule generated by the sequence $s \in \mathcal{S}$ a statistician decides to terminate the sampling at the minimal $n \geq 0$ such that $J_{n+1} = 0$. The corresponding size of the sample is a random variable given by

$$\nu(s) = \min \left\{ n \geq 0: J_{n+1}(\mathbf{Q}, \xi^{(n)}) = 0 \right\}.$$

It is clear that $\{\nu(s) = n\} \in \mathcal{F}_n, n \geq 0$, that is, $\nu(s)$ is a *stopping rule* with respect to the family of σ -algebras $\mathcal{F}_0, \mathcal{F}_1, \mathcal{F}_2, \dots$.

Given $(\mathbf{Q}, \xi^{(n)})$ denote by $\mathbf{Q}(\cdot/\xi^{(n)})$ the a posteriori distribution of the parameter θ . Let $P_{\mathbf{Q}}^n$ be the joint unconditional (weighted) distribution of the vector $\xi^{(n)}$ given \mathbf{Q} . Denote by $\rho_0(\mathbf{Q}, \xi^{(n)})$ the Bayes risk corresponding to the Bayes decision for given $(\mathbf{Q}, \xi^{(n)})$ and $\nu(s) = n$, that is,

$$\rho_0(\mathbf{Q}, \xi^{(n)}) = \min_{1 \leq i \leq m} \int_{\Theta} A_i(\theta) \mathbf{Q}(d\theta/\xi^{(n)}).$$

Then *a priori risk* corresponding to the stopping rule s and the Bayes decision rule is defined by

$$\begin{aligned} R(\mathbf{Q}, s) &= \sum_{n=1}^{\infty} \int \left[K_n(x^{(n)}) + \rho_0(\mathbf{Q}, x^{(n)}) \right] \\ &\quad \times \prod_{k=0}^{n-1} s_{k+1}(\mathbf{Q}, x^{(k)}) (1 - s_{n+1}(\mathbf{Q}, x^{(n)})) P_{\mathbf{Q}}^n(dx^{(n)}) \\ &\quad + (1 - s_1(\mathbf{Q}))\rho_0(\mathbf{Q}) \end{aligned}$$

where $s_1(\mathbf{Q}, x^{(0)}) = s_1(\mathbf{Q})$ and $\rho_0(\mathbf{Q})$ is the Bayes risk to make a decision without sampling, that is,

$$\rho_0(\mathbf{Q}) = \min_{1 \leq i \leq m} \int_{\Theta} A_i(\theta) \mathbf{Q}(d\theta).$$

A stopping rule s^* is called *Bayes* (or *optimal*) for the a priori distribution \mathbf{Q} if

$$R(\mathbf{Q}, s^*) = \inf_{s \in \mathcal{S}} R(\mathbf{Q}, s).$$

In what follows we assume that $R(\mathbf{Q}, s^*) < \infty$.

The problem of evaluating a Bayes stopping rule s^* for a given a priori distribution is not easy. Below we consider some general properties of the rule s^* .

Properties of a Bayes stopping rule. Let $\mathcal{S}^{(N)}$ be a subclass of stopping rules $s \in \mathcal{S}$ truncated at the moment $\nu = N$, that is, $s_{j+1}(\mathbf{Q}, \xi^{(j)}) = 0$ for all $j \geq N$ and all $s \in \mathcal{S}^{(N)}$. Below we obtain a Bayes truncated rule $s^{(N)} \in \mathcal{S}^{(N)}$ such that

$$R(\mathbf{Q}, s^{(N)}) = \inf_{s \in \mathcal{S}^{(N)}} R(\mathbf{Q}, s).$$

According to the dynamic programming method (see [15]), we define the rule $s^{(N)}$ by constructing the $N + 1$ functions

$$\begin{aligned} \rho_0^{(N)}(\mathbf{Q}, \xi^{(N)}) &= \rho_0(\mathbf{Q}, \xi^{(N)}), \\ \rho_j^{(N)}(\mathbf{Q}, \xi^{(N-j)}) &= \min \left\{ \rho_0(\mathbf{Q}, \xi^{(N-j)}), E \left\{ \Delta(\eta; \xi^{(N-j)}) + \rho_{j-1}^{(N)}(\mathbf{Q}, (\xi^{(N-j)}, \eta)) / \mathcal{F}_{N-j} \right\} \right\}, \end{aligned}$$

$j = 1, 2, \dots, N$, where $E\{\cdot / \mathcal{F}_{N-j}\}$ is the conditional expectation with respect to the distribution

$$P(\cdot / \mathcal{F}_{N-j}) = \int_{\Theta} P_{\theta}(\cdot) \mathbf{Q}(d\theta / \xi^{(N-j)}), \quad j = 0, 1, \dots, N,$$

and

$$\Delta(\eta; \xi^{(N-j)}) = K_{N-j+1}(\xi^{(N-j)}, \eta) - K_{N-j}(\xi^{(N-j)}, \eta).$$

Let

$$\nu^{(N)} = \min \left\{ n: 0 \leq n \leq N, \rho_{N-n}^{(N)}(\mathbf{Q}, \xi^{(n)}) = \rho_0(\mathbf{Q}, \xi^{(n)}) \right\}.$$

The number $\nu^{(N)}$ is the size of the sample corresponding to the Bayes truncated stopping rule $s^{(N)}$. Since $\{\nu^{(N)} = n\} \in \mathcal{F}_n$ for all $n \geq 0$, $\nu^{(N)}$ is a stopping rule. The sequence $s^{(N)} = (s_1^{(N)}, s_2^{(N)}, \dots)$ is a Bayes stopping rule, and moreover

$$s_j^{(N)} = s_j^{(N)}(\mathbf{Q}, \xi^{(j-1)}) = \begin{cases} 1, & \text{if } j \leq \nu^{(N)}, \\ 0, & \text{if } j > \nu^{(N)}. \end{cases}$$

For any $N \geq 1$, $\rho_N^{(N)}(\mathbf{Q})$ is the a priori risk corresponding to \mathbf{Q} and $\mathcal{S}^{(N)}$. Denote by $\mathcal{S}_1^{(N)}$, $N \geq 1$, the subclass of stopping rules of $\mathcal{S}^{(N)}$ for which at least one observation is taken. Then

$$(4.1.1) \quad \rho_N^{(N)}(\mathbf{Q}) = \min \left\{ \rho_0(\mathbf{Q}), \inf_{s \in \mathcal{S}_1^{(N)}} R(\mathbf{Q}, s) \right\}.$$

Since $\mathcal{S}_1^{(N)} \subset \mathcal{S}_1^{(N+1)}$ for all $N \geq 1$, we have

$$(4.1.2) \quad \inf_{s \in \mathcal{S}_1^{(N+1)}} R(\mathbf{Q}, s) \leq \inf_{s \in \mathcal{S}_1^{(N)}} R(\mathbf{Q}, s).$$

It follows from (4.1.1) and (4.1.2) that $\rho_N^{(N)}(\mathbf{Q}) \geq \rho_{N+1}^{(N+1)}(\mathbf{Q})$ for all $N \geq 1$. Since $\rho_n^{(N)}(\mathbf{Q}) \geq 0$ for all $N \geq 1$, the limit

$$\rho^{(\infty)}(\mathbf{Q}) = \lim_{N \rightarrow \infty} \rho_N^{(N)}(\mathbf{Q})$$

exists.

By $\rho(\mathbf{Q})$ we denote the a priori Bayes risk related to the class \mathcal{S} of all stopping rules, that is, $\rho(\mathbf{Q}) = R(\mathbf{Q}, s^*)$. The risk $\rho(\mathbf{Q})$ satisfies the equation

$$(4.1.3) \quad \rho(\mathbf{Q}) = \min\{\rho_0(\mathbf{Q}), E_Q\{K_1(\eta) + \rho(\mathbf{Q}_\eta)\}\}$$

where $\mathbf{Q}_\eta(\cdot) = \mathbf{Q}(\cdot/\eta)$ is the a posteriori distribution of the parameter θ given (\mathbf{Q}, η) .

Indeed, if \mathcal{S}_1 is the class of all stopping rules for which at least one observation is taken, then

$$(4.1.4) \quad \rho(\mathbf{Q}) = \min\left\{\rho_0(\mathbf{Q}), \inf_{s \in \mathcal{S}_1} R(\mathbf{Q}, s)\right\}.$$

For all $s \in \mathcal{S}_1$, $s = (s_1(\mathbf{Q}), s_2(\mathbf{Q}, \xi_1), s_3(\mathbf{Q}, \xi_1, \xi_2), \dots)$, we define the *reduced* stopping rule $s^+ = (s_2(\mathbf{Q}, \xi_1), s_3(\mathbf{Q}, \xi_1, \xi_2), \dots)$. Thus we have for all $s \in \mathcal{S}_1$ that

$$(4.1.5) \quad R(\mathbf{Q}, s) = E_Q\{K_1(\eta) + R(\mathbf{Q}_\eta, s^+)\}.$$

Equality (4.1.5) and the Fatou lemma imply that

$$(4.1.6) \quad \begin{aligned} \inf_{s \in \mathcal{S}_1} R(\mathbf{Q}, s) &\geq E_Q\left\{K_1(\eta) + \inf_{s \in \mathcal{S}_1} R(\mathbf{Q}_\eta, s^+)\right\} \\ &= E_Q\left\{K_1(\eta) + \inf_{s \in \mathcal{S}_1} R(\mathbf{Q}_\eta, s)\right\} = E_Q\{K_1(\eta) + \rho(\mathbf{Q}_\eta)\}. \end{aligned}$$

We obtain from (4.1.4) and (4.1.6) that

$$\rho(\mathbf{Q}) \geq \min\{\rho_0(\mathbf{Q}), E_Q\{K_1(\eta) + \rho(\mathbf{Q}_\eta)\}\} \geq \inf_{s \in \mathcal{S}} R(\mathbf{Q}, s) = \rho(\mathbf{Q}),$$

that is, $\rho(\mathbf{Q})$ satisfies equation (4.1.3).

Similarly we show that for all $N \geq 1$

$$(4.1.7) \quad \rho_{N+1}^{(N+1)}(\mathbf{Q}) = \min\left\{\rho_0(\mathbf{Q}), E_Q\{K_1(\eta) + \rho_N^{(N)}(\mathbf{Q}_\eta)\}\right\}.$$

If $E_Q \rho_1^{(1)}(\mathbf{Q}_\eta) < \infty$, then the Lebesgue dominated convergence theorem implies that

$$\begin{aligned} \rho^{(\infty)} &= \lim_{N \rightarrow \infty} \rho_{N+1}^{(N+1)}(\mathbf{Q}) = \min\left\{\rho_0(\mathbf{Q}), \lim_{N \rightarrow \infty} \{K_1(\eta) + \rho_N^{(N)}(\mathbf{Q}_\eta)\}\right\} \\ &= \min\left\{\rho_0(\mathbf{Q}), E_Q\{K_1(\eta) + \rho^{(\infty)}(\mathbf{Q}_\eta)\}\right\}, \end{aligned}$$

since the risk $\rho_N^{(N)}$, $N \geq 1$, is monotone. Thus the limit $\rho^{(\infty)}(\mathbf{Q})$ also satisfies equation (4.1.3). Therefore we have proved the following result.

THEOREM 4.1.1. *If $E_Q K_1(\eta) < \infty$ and $E_Q \rho_1^{(1)}(\eta_\eta) < \infty$ for a given a priori distribution \mathbf{Q} , then*

$$\rho(\mathbf{Q}) = \lim_{N \rightarrow \infty} \rho_N^{(N)}(\mathbf{Q}).$$

Theorem 4.1.1 implies for all $\varepsilon > 0$ that there exists an integer number $N(\varepsilon)$ such that the Bayes truncated stopping rule $s^{(N)}$ for all $N \geq N(\varepsilon)$ is an ε -Bayes stopping rule for the nontruncated problem, that is,

$$0 \leq \rho_N^{(N)}(\mathbf{Q}) - \rho(\mathbf{Q}) < \varepsilon.$$

Below are two results about the existence of a Bayes (optimal) stopping rule.

THEOREM 4.1.2. *If $\rho(\mathbf{Q}) = \lim_{N \rightarrow \infty} \rho_N^{(N)}(\mathbf{Q})$ and for all $n \geq N_0$*

$$\rho_0(\mathbf{Q}, \xi^{(n)}) - E \left\{ \rho_0(\mathbf{Q}, (\xi^{(n)}, \eta)) / \mathcal{F}_n \right\} \leq E \left\{ \Delta(\eta; \xi^{(n)}) / \mathcal{F}_n \right\}$$

almost surely, then $\rho_{N_0}^{(N_0)} = \rho(\mathbf{Q})$ and the truncated Bayes sequential rule $s^{(N_0)}$ is optimal for the nontruncated problem.

THEOREM 4.1.3. *If for all $n \geq 1$ the a posteriori risk $\rho_0(\mathbf{Q}, \xi)$ does not depend on $\xi^{(n)}$, that is, $\rho_0(\mathbf{Q}, \xi^{(n)}) = \rho_0(\mathbf{Q}, n)$ almost surely for all $n \geq 1$, then the Bayes sequential stopping rule is the one with a fixed size n_0 of the sample where n_0 is the minimal positive integer number n that minimizes $\rho_0(\mathbf{Q}, n) + E_{\mathbf{Q}} K_n(\xi^{(n)})$.*

The proofs of Theorems 4.1.2 and 4.1.3 can be found in [54].

More details on the general theory of Bayes sequential rules are given in [4, 15, 54].

The evaluation of Bayes stopping rules is a complicated problem, especially in the case of composite hypotheses. The following example shows that this problem is complicated even in the case of two simple hypotheses.

EXAMPLE 4.1.1. Let ξ_1, ξ_2, \dots be a sequence of independent identically distributed random variables assuming two values 1 and 0 with probabilities θ and $1 - \theta$, respectively, where θ is an unknown parameter. Let θ assume only two values $1/3$ and $2/3$, that is, $\Theta = \{1/3, 2/3\}$. Let the a priori measure \mathbf{Q} be determined by the number $q = P\{\theta = 1/3\} = 1 - P\{\theta = 2/3\}$. Thus we deal with two simple hypotheses $H_1: \theta = 1/3$ and $H_2: \theta = 2/3$. Let the loss functions $A_i(\theta)$ be defined by the numbers $A_1(1/3) = A_2(2/3) = 0$ and $A_1(2/3) = A_2(1/3) = 20$, while the cost per observation is determined by the equality $K_n(\xi^{(n)}) = n$. This means that a single observation costs 1 dollar. Put $\rho_N(q) = \rho_N^{(N)}(\mathbf{Q})$, $N = 0, 1, 2, \dots$. Now we evaluate $\rho_0(q)$, $\rho_1(q)$, and $\rho_2(q)$. The definition of $\rho_0(q)$ implies that without sampling one accepts the hypothesis H_2 if $0 \leq q \leq 1/2$, and the hypothesis H_1 if $1/2 \leq q \leq 1$. Moreover

$$(4.1.8) \quad \rho_0(q) = \begin{cases} 20q, & \text{if } 0 \leq q \leq 1/2, \\ 20(1 - q), & \text{if } 1/2 \leq q \leq 1. \end{cases}$$

This implies that $\rho_0(q) = \rho_0(1 - q)$ for $0 \leq q \leq 1$. Since the problem is symmetric, $\rho_j(q) = \rho_j(1 - q)$ for $0 \leq q \leq 1$ and for all $j = 1, 2, \dots$. Thus we need to evaluate $\rho_1(q)$ and $\rho_2(q)$ only for $q \in [0, 1/2]$.

Let $q(x)$ stand for the a priori probability of the event $\{\theta = 1/3\}$ given the observation η is equal to x where either $x = 0$ or $x = 1$. Applying the Bayes formula we obtain

$$(4.1.9) \quad q(1) = \frac{q}{q + 2(1 - q)}, \quad q(0) = \frac{2q}{2q + (1 - q)}.$$

It follows from (4.1.8) and (4.1.9) that

$$(4.1.10) \quad \begin{aligned} &\rho_0(q(1)) = 20q(1) \quad \text{for } 0 \leq q \leq 1/2, \\ \rho_0(q(0)) &= \begin{cases} 20q(0), & \text{if } 0 \leq q \leq 1/3, \\ 20(1 - q(0)), & \text{if } 1/3 \leq q \leq 1/2. \end{cases} \end{aligned}$$

The unconditional distribution of the observation η is given by

$$(4.1.11) \quad P\{\eta = 1\} = \frac{1}{3}q + \frac{2}{3}(1 - q) = 1 - P\{\eta = 0\}.$$

It is clear that

$$(4.1.12) \quad \begin{aligned} E\rho_0(q(\eta)) &= \rho_0(q(1))P\{\eta = 1\} + \rho_0(q(0))P\{\eta = 0\} \\ &= \begin{cases} 20q, & \text{if } 0 \leq q \leq 1/3, \\ 20/3, & \text{if } 1/3 \leq q \leq 1/2. \end{cases} \end{aligned}$$

Since $K_1(\xi^{(1)}) = 1$, we derive from (4.1.7), (4.1.8), and (4.1.12) that

$$(4.1.13) \quad \rho_1(q) = \min\{\rho_0(q), E\rho_0(q(\eta) + 1)\} = \begin{cases} 20q, & \text{if } 0 \leq q \leq 23/60, \\ 23/3, & \text{if } 23/60 \leq q \leq 1/2. \end{cases}$$

First we consider the case $N \leq 1$. Relations (4.1.8) and (4.1.13) imply that if $0 \leq q \leq 23/60$ or $37/60 \leq q \leq 1$, then one should make the final decision without sampling, while if $23/60 < q < 37/60$, then it is necessary to take the first observation and one should make the final decision based on a sample consisting of a single observation.

Now let $N \leq 2$. To evaluate $p_2(q)$ we note that (4.1.9) implies that there are three pairs of equivalent inequalities, namely

$$\begin{aligned} q(1) &\leq 23/60 \text{ and } q \leq 46/83 \iff \\ q(0) &\leq 23/60 \text{ and } q \leq 23/97 \iff \\ q(0) &\geq 37/60 \text{ and } q \geq 37/83. \end{aligned}$$

Using (4.1.11) and (4.1.13) and taking into account the symmetry of the function $\rho_1(q)$ we get

$$(4.1.14) \quad E\rho_1(q(\eta)) = \begin{cases} 20q, & \text{if } 0 \leq q \leq \frac{23}{97}, \\ \frac{83q+23}{9}, & \text{if } \frac{23}{97} < q \leq \frac{37}{83}, \\ \frac{20}{3}, & \text{if } \frac{37}{83} < q \leq \frac{1}{2}. \end{cases}$$

Thus (4.1.7) yields

$$(4.1.15) \quad \rho_2(q) = \min\{\rho_0(q), E\rho_1(q(\eta) + 1)\} = \begin{cases} 20q, & \text{if } 0 \leq q \leq \frac{32}{97}, \\ \frac{83q+32}{9}, & \text{if } \frac{32}{97} < q \leq \frac{37}{83}, \\ \frac{23}{3}, & \text{if } \frac{37}{83} < q \leq \frac{1}{2}. \end{cases}$$

Relations (4.1.8), (4.1.13), and (4.1.15) show in the case $N \leq 2$ that if

$$0 \leq q \leq 32/97 \quad \text{or} \quad 65/97 \leq q \leq 1,$$

then a final decision should be made without sampling, while if $37/83 \leq q \leq 46/83$, then it is necessary to take the first observation and a final decision should be made based on a sample consisting of a single observation. Finally if

$$32/97 < q < 37/83 \quad \text{or} \quad 46/83 < q < 65/97,$$

then it is necessary to take the second observation and a final decision should be made based on a sample consisting of two observations.

The graphs of the functions $\rho_0(q)$, $\rho_1(q)$, and $\rho_2(q)$ are shown in Figure 4.1.1. The evaluation of the functions $\rho_N(q)$ becomes a time consuming procedure for

large N , since the number of intervals where they change their slopes increases with N . If the size of the sample is not bounded in advance, then the optimal procedure is described in Theorems 4.1.1–4.1.3 (see also Chapter 4 in [15]).

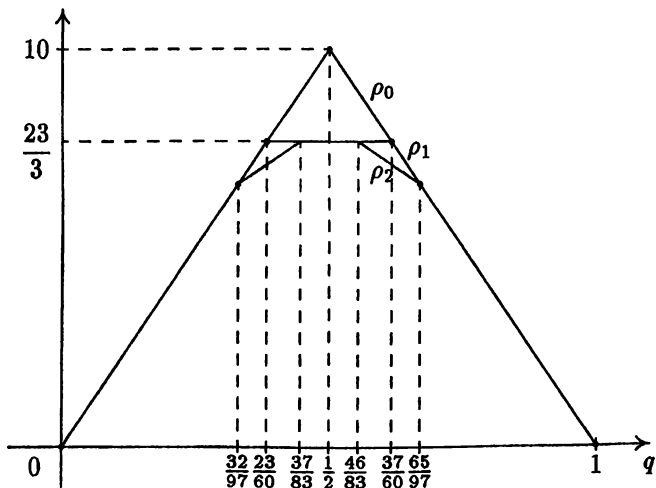


FIGURE 4.1.1. Risk functions $\rho_0(q)$, $\rho_1(q)$, and $\rho_2(q)$

EXAMPLE 4.1.2. This is a continuation of Example 4.1.1. Let the loss functions $A_i(\theta)$ be determined by the numbers $A_1(1/3) = A_2(2/3) = 0$ and

$$A_1(2/3) = A_2(1/3) = 10.$$

As in Example 4.1.1 we have $\rho_j(q) = \rho_j(1-q)$ for $0 \leq q \leq 1$ and $j = 0, 1, 2, \dots$. Thus $\rho_0(q) = 10q$ for $0 \leq q \leq 1/2$ and $\rho_0(q) = 10(1-q)$ for $1/2 \leq q \leq 1$. Similarly to Example 4.1.1 we get

$$E\rho_0(q(\eta)) = \begin{cases} 10q, & \text{if } 0 \leq q \leq 1/3, \\ 10/3, & \text{if } 1/3 < q \leq 1/2. \end{cases}$$

Hence (4.1.7) implies that

$$(4.1.16) \quad \rho_1(q) = \begin{cases} 10q, & \text{if } 0 \leq q \leq 13/30, \\ 13/3, & \text{if } 13/30 < q \leq 1/2. \end{cases}$$

Relations (4.1.9) imply that $q(1) \leq 13/30 \iff q \leq 26/43$, $q(0) \leq 13/30 \iff q \leq 13/47$ and $q(0) \geq 17/30 \iff q \geq 17/43$. Thus

$$E\rho_1(q(\eta)) = \begin{cases} 10q, & \text{if } 0 \leq q \leq \frac{13}{47}, \\ \frac{43q+13}{9}, & \text{if } \frac{13}{47} < q \leq \frac{17}{43}, \\ \frac{10}{3}, & \text{if } \frac{17}{43} < q \leq \frac{1}{2}. \end{cases}$$

Using (4.1.7) we obtain for $N = 2$ that

$$(4.1.17) \quad \rho_2(q) = \begin{cases} 10q, & \text{if } 0 \leq q \leq 13/30, \\ 13/3, & \text{if } 13/30 < q \leq 1/2. \end{cases}$$

Relations (4.1.16) and (4.1.17) imply that $\rho_1(q) = \rho_2(q)$ for all $q \in [0, 1]$. Moreover we derive from (4.1.7) that

$$\rho_3(q) = \min\{\rho_0(q), E\rho_2(q(\eta)) + 1\} = \min\{\rho_0(q), E\rho_1(q(\eta)) + 1\} = \rho_2(q)$$

if $0 \leq q \leq 1$. We conclude by induction that $\rho_N(q) = \rho_1(q)$ for $N = 2, 3, \dots$, that is, the conditions of Theorem 4.1.3 hold.

The above discussion shows that if $13/30 < q < 17/30$, then it is necessary to take the first observation and a decision should be made based on a sample consisting of a single observation, while otherwise a decision can be made without any observation.

In the next section we consider the Wald sequential test for distinguishing two simple hypotheses in the case of general distributions of observations ξ_1, ξ_2, \dots .

4.2. Wald sequential tests

Main definitions and notation. Let $\xi = (\xi_1, \xi_2, \dots)$ be a sequence of independent identically distributed random variables with a distribution P_θ where θ is an unknown parameter assuming only two values θ_1 and θ_2 . Thus we deal with the case $\Theta = \{\theta_1, \theta_2\}$ and each of the sets $\Theta_1 = \{\theta_1\}$ and $\Theta_2 = \{\theta_2\}$ contains only a single point. Therefore there are two simple hypotheses about the distribution of an observation, namely $H_1: \theta = \theta_1$ and $H_2: \theta = \theta_2$. Let the distribution P_θ be absolutely continuous with respect to some σ -finite measure μ and denote the density by $p(x; \theta)$. Then the measure $P_\theta^{(n)}$ determining the distribution of the sample $\xi^{(n)} = (\xi_1, \xi_2, \dots, \xi_n)$ is absolutely continuous with respect to the measure $\mu^n = \mu \times \mu \times \dots \times \mu$ and its density is $p_n(x^{(n)}; \theta) = \prod_{k=1}^n p(x_k; \theta)$, $x^{(n)} = (x_1, x_2, \dots, x_n)$. Let $z_n(x^{(n)})$ be the likelihood ratio

$$(4.2.1) \quad z_n(x^{(n)}) = \frac{p_n(x^{(n)}; \theta_2)}{p_n(x^{(n)}; \theta_1)} = \prod_{k=1}^n \frac{p(x_k; \theta_2)}{p(x_k; \theta_1)}$$

where $x^{(n)} = (x_1, x_2, \dots, x_n)$ (we agree that $0/0 = 0$). Put

$$(4.2.2) \quad \Lambda_n(x^{(n)}) = \ln z_n(x^{(n)}), \quad \Lambda_n = \Lambda_n(\xi^{(n)}), \quad z_n = z_n(\xi^{(n)}).$$

Therefore

$$(4.2.3) \quad \Lambda_n = \sum_{k=1}^n \lambda_k, \quad \lambda_k = \lambda(\xi_k) = \ln \frac{p(\xi_k; \theta_2)}{p(\xi_k; \theta_1)}, \quad k = 1, 2, \dots, n.$$

Throughout this section we consider sequential tests whose stopping times belong to the class \mathcal{S}_1 (in other words, the tests for which it is necessary to take at least one observation) and which depend on two constants a and b such that $0 < a$ and $b < \infty$. Such a test is called a *Wald sequential test* for distinguishing two simple hypotheses H_1 and H_2 if its stopping rule says that a statistician continues the sampling until $-b < \Lambda_n < a$, that is, the Wald stopping time ν is

$$(4.2.4) \quad \nu = \inf\{n \geq 0: \Lambda_n \notin (-b, a)\}.$$

If $\nu = n$, then the decision is as follows: the hypothesis H_2 is accepted if $\Lambda_n \geq a$, while the alternative H_1 is accepted if $\Lambda_n \leq -b$. The test described above is also called the *Wald sequential test with limit points* $(-b, a)$.

The above test is often called a *Wald sequential likelihood ratio test* or *sequential likelihood ratio test* or *Wald test*. For a Wald test we denote by α and β the *type I* and *type II error probabilities*, respectively, that is,

$$(4.2.5) \quad \alpha = P_{\theta_1}\{\Lambda_\nu \geq a\}, \quad \beta = P_{\theta_2}\{\Lambda_\nu \leq -b\}$$

where $P_{\theta_i}\{\cdot\}$ are the probabilities of events if the distribution of ξ_1, ξ_2, \dots is determined by the measure P_{θ_i} . The pair of numbers (α, β) is called the *power of a test*. We say that a *Wald sequential test is terminated with probability one during a finite time* if $P_{\theta_i}\{\nu < \infty\} = 1$ for all $i = 1, 2$.

Inequalities for error probabilities of a Wald test. The following two results establish relationships between the power (α, β) of a Wald test and its limit points $(-b, a)$.

LEMMA 4.2.1. *If a Wald sequential test of power (α, β) and with limit points $(-b, a)$ is terminated with probability one during a finite time, then*

$$(4.2.6) \quad A \leq \frac{1-\beta}{\alpha}, \quad B \geq \frac{\beta}{1-\alpha}$$

where $0 < A < 1 < B < \infty$ and $\ln B = -b$, $\ln A = a$.

PROOF. Consider the sets

$$(4.2.7) \quad W_n = \left\{ x^{(n)} : B < z_k(x^{(k)}) < A \text{ for all } k = 1, 2, \dots, n-1, \right. \\ \left. z_n(x^{(n)}) \geq A \right\},$$

$$(4.2.8) \quad V_n = \left\{ x^{(n)} : B < z_k(x^{(k)}) < A \text{ for all } k = 1, 2, \dots, n-1, \right. \\ \left. z_n(x^{(n)}) \leq B \right\}$$

where $n = 1, 2, \dots$, $x^{(n)} = (x_1, x_2, \dots, x_n)$, and $z_n(x^{(n)})$ is the likelihood ratio (4.2.1). Using (4.2.2)–(4.2.4) and (4.2.7)–(4.2.8) we get

$$(4.2.9) \quad \{\Lambda_\nu \geq a\} = \bigcup_{n=1}^{\infty} \{\xi^{(n)} \in W_n\}, \quad \{\Lambda_\nu \leq b\} = \bigcup_{n=1}^{\infty} \{\xi^{(n)} \in V_n\}.$$

Relations (4.2.5) and (4.2.9) imply that

$$\alpha = \sum_{n=1}^{\infty} \int_{W_n} P_{\theta_1}^{(n)}(dx) = \sum_{n=1}^{\infty} \int_{W_n} \frac{1}{z_n(x)} P_{\theta_2}^{(n)}(dx) \leq \frac{1}{A} \sum_{n=1}^{\infty} \int_{W_n} P_{\theta_2}^{(n)}(dx) = \frac{1-\beta}{A},$$

that is, the first inequality in (4.2.6) is proved. Here we used for $i = 1$ and $i = 2$ that

$$(4.2.10) \quad \sum_{n=1}^{\infty} \int_{W_n} P_{\theta_1}^{(n)}(dx) + \sum_{n=1}^{\infty} \int_{V_n} P_{\theta_1}^{(n)}(dx) \\ = \sum_{n=1}^{\infty} P_{\theta_1}\{\nu = n\} = P_{\theta_1}\{\nu < \infty\} = 1$$

by the assumptions of the lemma. Similarly, relations (4.2.5) and (4.2.9) imply

$$1 - \alpha = \sum_{n=1}^{\infty} \int_{V_n} P_{\theta_1}^{(n)}(dx) = \sum_{n=1}^{\infty} \int_{V_n} \frac{1}{z_n(x)} P_{\theta_2}^{(n)}(dx) \geq \frac{1}{B} \sum_{n=1}^{\infty} \int_{V_n} P_{\theta_2}^{(n)}(dx) = \frac{\beta}{B},$$

that is, the second inequality in (4.2.6) is also proved. \square

LEMMA 4.2.2. *Let $A = (1 - \beta)/\alpha$, $B = \beta/(1 - \alpha)$ and $-b = \ln B$, $a = \ln A$. If a Wald sequential test with limit points $(-b, a)$ is terminated with probability one during a finite time and its power is (α', β') , then*

$$(4.2.11) \quad \alpha' \leq \frac{\alpha}{1 - \beta}, \quad \beta' \leq \frac{\beta}{1 - \alpha},$$

$$(4.2.12) \quad \alpha' + \beta' \leq \alpha + \beta.$$

PROOF. According to Lemma 4.2.1,

$$(4.2.13) \quad A = \frac{1 - \beta}{\alpha} \leq \frac{1 - \beta'}{\alpha'}, \quad B = \frac{\beta}{1 - \alpha} \geq \frac{\beta'}{1 - \alpha'}.$$

This implies

$$\alpha' \leq (1 - \beta') \frac{\alpha}{1 - \beta} \leq \frac{\alpha}{1 - \beta}, \quad \beta' \leq (1 - \alpha') \frac{\beta}{1 - \alpha} \leq \frac{\beta}{1 - \alpha},$$

that is, inequalities (4.2.11) hold. Moreover (4.2.13) yields

$$\alpha'(1 - \beta) \leq \alpha(1 - \beta'), \quad \beta'(1 - \alpha) \leq \beta(1 - \alpha').$$

Combining these inequalities we obtain

$$\alpha' - \alpha'\beta + \beta' - \alpha\beta' \leq \alpha - \alpha\beta' + \beta - \alpha'\beta,$$

whence inequality (4.2.12) follows. \square

REMARK 4.2.1. Lemma 4.2.2 for small α and β implies that the power of a Wald test with limit points $(\ln(\beta/(1 - \alpha)), \ln((1 - \beta)/\alpha))$ is (α', β') where α' is close to α , β' is close to β , and always $\alpha' + \beta' \leq \alpha + \beta$, that is, $\beta' < \beta$ if $\alpha' > \alpha$ and $\alpha' < \alpha$ if $\beta' > \beta$.

Properties of the stopping time of a Wald sequential test. The following result contains sufficient conditions for the finiteness of the Wald stopping time ν and its moment generating function.

LEMMA 4.2.3. *Let ν be the stopping time of a Wald sequential test with limit points $(-b, a)$. Let $P_{\theta}\{|\lambda_1| > 0\} > 0$ where either $\theta = \theta_1$ or $\theta = \theta_2$, and $\lambda_1 = \lambda_1(\xi_1)$. Then*

a) $P_{\theta}\{\nu < \infty\} = 1;$

b) $E_{\theta}e^{t\nu} < \infty$ for all $t < t_0$ where t_0 is a positive number and E_{θ} is the expectation with respect to the distribution P_{θ} .

PROOF. Let m and k be fixed integer numbers such that $m > k$ and $r = [m/k]$ where $[c]$ is the integer part of a number c . Consider the random variables

$$T_1 = \Lambda_k, \quad T_i = \Lambda_{ik} - \Lambda_{(i-1)k}, \quad i = 2, 3, \dots, r.$$

If $\nu > m$, then $\Lambda_i \in (-b, a)$ for all $i = 1, 2, \dots, m$. In particular, this inclusion holds for $i = k, 2k, \dots, rk$. Thus

$$|T_i| < b \vee a = c, \quad i = 1, 2, \dots, r.$$

Since the random variables T_1, T_2, \dots, T_r are independent and identically distributed, we have

$$(4.2.14) \quad P_\theta\{\nu > m\} \leq P_\theta\{|T_i| < c \text{ for all } i = 1, 2, \dots, r\} = (P_\theta\{|T_1| < c\})^r.$$

Note that $P_\theta\{|\lambda_1| > 0\} > 0$, whence it follows that there exists a positive number h such that either $P_\theta\{\lambda_1 > h\} > 0$ or $P_\theta\{\lambda_1 < -h\} > 0$. If k is greater than c/h , then

$$\begin{aligned} P_\theta\{|T_1| \geq c\} &= P_\theta\{|\lambda_1 + \dots + \lambda_k| \geq c\} \\ &\geq P_\theta\left\{\lambda_i \geq \frac{c}{k} \text{ for all } i = 1, 2, \dots, k\right\} \\ &\quad + P_\theta\left\{\lambda_i \leq -\frac{c}{k} \text{ for all } i = 1, 2, \dots, k\right\} \\ &\geq (P_\theta\{\lambda_1 > h\})^k + (P_\theta\{\lambda_1 < -h\})^k > 0. \end{aligned}$$

This yields $P_\theta\{|T_1| < c\} < 1$. In view of (4.2.14) we therefore get

$$\lim_{m \rightarrow \infty} P_\theta\{\nu > m\} = P_\theta\{\nu = \infty\} = 0,$$

that is, $P_\theta\{\nu < \infty\} = 1$ and statement a) is proved.

Relation (4.2.14) for $t \geq 0$ implies that

$$\begin{aligned} E_\theta e^{t\nu} &= \sum_{n=1}^{\infty} e^{tn} P_\theta\{\nu = n\} = \sum_{r=1}^{\infty} \sum_{n=(r-1)k+1}^{kr} e^{tn} P_\theta\{\nu = n\} \\ (4.2.15) \quad &\leq \sum_{r=1}^{\infty} e^{tkr} P_\theta\{(r-1)k < \nu \leq kr\} \leq \sum_{r=1}^{\infty} e^{tkr} P_\theta\{\nu > (r-1)k\} \\ &\leq \sum_{r=1}^{\infty} e^{tkr} \gamma^{r-1} \end{aligned}$$

where $\gamma = P_\theta\{|T_1| < c\}$ and k is such that $\gamma < 1$. It follows from (4.2.15) that $E_\theta e^{t\nu} < \infty$ for $\gamma e^{tk} < 1$, thus for $t \in [0, t_0)$ where $t_0 = k^{-1} \ln \gamma^{-1}$.

If $t < 0$, then the equality $P_\theta\{\nu \geq 0\} = 1$ implies that $E_\theta e^{t\nu} \leq e^{t \cdot 0} = 1$ and statement b) is also proved. \square

REMARK 4.2.2. Lemma 4.2.3 holds in the case where θ is different from both θ_1 and θ_2 but if the random variables ξ_1, ξ_2, \dots are independent and identically distributed with respect to the measure P_θ . If the conditions of Lemma 4.2.3 are satisfied for both $\theta = \theta_1$ and $\theta = \theta_2$, then statement a) implies that the Wald sequential test with limit points $(-b, a)$ is terminated during a finite time with probability one.

REMARK 4.2.3. The assumption that a Wald sequential test is terminated during a finite time with probability one ($P_{\theta_i}\{\lambda_1 \neq 0\} > 0$, $i = 1, 2$) is not too restrictive. Moreover this assumption is quite natural, since otherwise $P_{\theta_i}\{\lambda_1 \neq 0\} = 0$ or $P_{\theta_i}\{\lambda_1 = 0\} = 1$ for $i = 1, 2$ and the measures P_{θ_1} and P_{θ_2} coincide. The hypotheses H_1 and H_2 are indistinguishable in this case (see Section 1.1), that is, the problem of distinguishing the hypotheses H_1 and H_2 makes no sense.

In what follows we need the following auxiliary result.

LEMMA 4.2.4. *Let ν be the stopping time of a Wald sequential test with limit points $(-b, a)$. If $E_{\theta}\nu < \infty$ and $E_{\theta}|\lambda_1| < \infty$ for $\theta = \theta_1$ or $\theta = \theta_2$, then*

$$(4.2.16) \quad E_{\theta}\Lambda_{\nu} = E_{\theta}\nu E_{\theta}\lambda_1.$$

PROOF. Consider random variables $\eta_i = I\{\nu > i - 1\}$ where $i = 1, 2, \dots$. Here $I(A)$ is the indicator of an event A , that is, $\eta_i = 1$ if a decision is not made by observations $\xi_1, \xi_2, \dots, \xi_{i-1}$. The random variable η_i is a function of $\xi_1, \xi_2, \dots, \xi_{i-1}$ and does not depend on ξ_i , thus it does not depend on $\lambda_i = \lambda(\xi_i)$. It is clear that

$$\Lambda_{\nu} = \lambda_1\eta_1 + \lambda_2\eta_2 + \dots,$$

whence

$$\begin{aligned} E_{\theta}\Lambda_{\nu} &= E_{\theta} \sum_{i=1}^{\infty} \lambda_i \eta_i = \sum_{i=1}^{\infty} E_{\theta} \lambda_i \eta_i = \sum_{i=1}^{\infty} E_{\theta} \lambda_i E_{\theta} \eta_i \\ &= E_{\theta} \lambda_1 \sum_{i=1}^{\infty} E_{\theta} \eta_i = E_{\theta} \lambda_1 \sum_{i=1}^{\infty} P_{\theta}\{\nu \geq i\} = E_{\theta} \lambda_1 E_{\theta} \nu. \end{aligned}$$

We interchanged the summation and expectation in the preceding relation, since

$$\sum_{i=1}^{\infty} E_{\theta} |\lambda_i \eta_i| = E_{\theta} |\lambda_1| \sum_{i=1}^{\infty} P_{\theta}\{\nu \geq i\} = E_{\theta} |\lambda_1| E_{\theta} \nu < \infty.$$

We also used the obvious equality

$$E_{\theta} \nu = \sum_{i=1}^{\infty} P_{\theta}\{\nu \geq i\}.$$

Thus equality (4.2.16) is proved. \square

Equality (4.2.16) is called the *Wald identity*. It also holds in the case where the assumptions of Lemma 4.2.4 are satisfied for some θ different from both θ_1 and θ_2 (see Remark 4.2.2).

The expectation of the stopping time of a Wald sequential test. The following result contains a lower bound of expectations $E_{\theta_i}\nu$ for $i = 1, 2$.

LEMMA 4.2.5. *Let ν be the stopping time of a Wald sequential test of power (α, β) and with limit points $(-b, a)$. If $P_{\theta_i}\{\lambda_1 \neq 0\} > 0$ and $E_{\theta_i}|\lambda_1| < \infty$ for $i = 1, 2$, then*

$$(4.2.17) \quad E_{\theta_1}\nu \geq -H(\alpha|1 - \beta)/E_{\theta_1}\lambda_1,$$

$$(4.2.18) \quad E_{\theta_2}\nu \geq H(\beta|1 - \alpha)/E_{\theta_2}\lambda_1$$

where

$$(4.2.19) \quad H(x|y) = x \ln \frac{x}{y} + (1-x) \ln \frac{1-x}{1-y}, \quad 0 \leq x, y \leq 1,$$

is a relative entropy of the distribution $(x, 1-x)$ with respect to the distribution $(y, 1-y)$.

PROOF. By the assumptions of the lemma $P_{\theta_i}(\lambda_1 \neq 0) > 0$ for $i = 1, 2$. Thus $E_{\theta_i} \nu < \infty$ for $i = 1, 2$ according to statement b) of Lemma 4.2.3, whence

$$(4.2.20) \quad E_{\theta_i} \Lambda_\nu = E_{\theta_i} \nu E_{\theta_i} \lambda_1, \quad i = 1, 2,$$

in view of Lemma 4.2.4. Now we conclude that the Wald test is terminated during a finite time with probability one by Lemma 4.2.3 and Remark 4.2.2, since

$$P_{\theta_i} \{\lambda_1 \neq 0\} > 0, \quad i = 1, 2.$$

This together with (4.2.9) and (4.2.10) implies

$$(4.2.21) \quad P_{\theta_i} \{\Lambda_\nu \geq a\} + P_{\theta_i} \{\Lambda_\nu \leq -b\} = P_{\theta_i} \{\nu < \infty\} = 1, \quad i = 1, 2.$$

Taking into account (4.2.21) for $i = 1$ we derive from (4.2.5) that

$$(4.2.22) \quad \begin{aligned} E_{\theta_1} \Lambda_\nu &= P_{\theta_1} \{\Lambda_\nu \leq -b\} E_{\theta_1} \{\Lambda_\nu / \Lambda_\nu \leq -b\} \\ &\quad + P_{\theta_1} \{\Lambda_\nu \geq a\} E_{\theta_1} \{\Lambda_\nu / \Lambda_\nu \geq a\} \\ &= (1-\alpha) E_{\theta_1} \{\Lambda_\nu / \Lambda_\nu \leq -b\} + \alpha E_{\theta_1} \{\Lambda_\nu / \Lambda_\nu \geq a\}. \end{aligned}$$

The Jensen inequality and definition (4.2.8) of the set V_n imply that

$$(4.2.23) \quad \begin{aligned} E_{\theta_1} \{\Lambda_\nu / \Lambda_\nu \leq -b\} &\leq \ln E_{\theta_1} \{e^{\Lambda_\nu / \Lambda_\nu} \leq -b\} \\ &= \ln \frac{1}{1-\alpha} \sum_{n=1}^{\infty} \int_{V_n} e^{\Lambda_n(x)} P_{\theta_1}^{(n)}(dx) \\ &= \ln \frac{1}{1-\alpha} \sum_{n=1}^{\infty} \int_{V_n} P_{\theta_2}^{(n)}(dx) = \ln \frac{\beta}{1-\alpha}. \end{aligned}$$

Similarly we get

$$(4.2.24) \quad \begin{aligned} E_{\theta_1} \{\Lambda_\nu / \Lambda_\nu \geq a\} &\leq \ln E_{\theta_1} \{e^{\Lambda_\nu / \Lambda_\nu} \geq a\} = \ln \frac{1}{\alpha} \sum_{n=1}^{\infty} \int_{W_n} e^{\Lambda_n(x)} P_{\theta_1}^{(n)}(dx) \\ &= \ln \frac{1}{\alpha} \sum_{n=1}^{\infty} \int_{W_n} P_{\theta_2}^{(n)}(dx) = \ln \frac{1-\beta}{\alpha}. \end{aligned}$$

It follows from equality (4.2.22) and inequalities (4.2.23) and (4.2.24) that

$$(4.2.25) \quad E_{\theta_1} \nu \leq -H(\alpha|1-\beta)$$

where $H(x|y)$ is the function defined by (4.2.19).

Since $P_{\theta_1}(\lambda_1 \neq 0) > 0$, we have $E_{\theta_1} \lambda_1 < 0$. Indeed, applying the elementary inequality $\ln x \leq x - 1$, $x > 0$, we get

$$(4.2.26) \quad \begin{aligned} E_{\theta_1} \lambda_1 &= \int \ln \frac{p(x; \theta_2)}{p(x; \theta_1)} P_{\theta_1}(dx) < \int \left(\frac{p(x; \theta_2)}{p(x; \theta_1)} - 1 \right) P_{\theta_1}(dx) \\ &= \int p(x; \theta_2) \mu(dx) - \int p(x; \theta_1) \mu(dx) = 0 \end{aligned}$$

in view of

$$P_{\theta_1}\{\lambda_1 = 0\} = P_{\theta_1}\left\{x: \ln \frac{p(x; \theta_2)}{p(x; \theta_1)} = 0\right\} = P_{\theta_1}\left\{x: \frac{p(x; \theta_2)}{p(x; \theta_1)} = 1\right\} < 1.$$

The lower bound (4.2.17) for $E_{\theta_1}\nu$ follows from (4.2.20) for $i = 1$ and inequality (4.2.25) by considering the sign of $E_{\theta_1}\lambda_1$.

Similarly to inequality (4.2.26) we obtain from $P_{\theta_2}\{\lambda_1 \neq 0\} > 0$ that $E_{\theta_2}\lambda_1 > 0$. Further we use the Jensen inequality and definition (4.2.8) of the set V_n and obtain

$$\begin{aligned} (4.2.27) \quad E_{\theta_2}\{\Lambda_\nu/\Lambda_\nu \leq -b\} &= -E_{\theta_2}\{-\Lambda_\nu/\Lambda_\nu \leq -b\} \geq -\ln E_{\theta_2}\{e^{-\Lambda_\nu}/\Lambda_\nu \leq -b\} \\ &= -\ln \frac{1}{\beta} \sum_{n=1}^{\infty} \int_{V_n} e^{-\Lambda_\nu(x)} P_{\theta_2}^{(n)}(dx) \\ &= -\ln \frac{1}{\beta} \sum_{n=1}^{\infty} \int_{V_n} P_{\theta_1}^{(n)}(dx) = -\ln \frac{1-\alpha}{\beta} = \ln \frac{\beta}{1-\alpha}. \end{aligned}$$

Following the same reasoning we get

$$(4.2.28) \quad E_{\theta_2}\{\Lambda_\nu/\Lambda_\nu \geq a\} \geq \ln \frac{1-\beta}{\alpha}.$$

Now (4.2.22), (4.2.27), and (4.2.28) imply

$$(4.2.29) \quad E_{\theta_2}\Lambda_\nu \geq H(\beta|1-\alpha).$$

Since $E_{\theta_2}\lambda_1$ is positive, equality (4.2.20) for $i = 2$ and inequality (4.2.29) imply the lower bound (4.2.18). \square

REMARK 4.2.4. The relative entropy of a measure P with respect to a measure \tilde{P} is defined in Section 2.3 and is denoted there by $I(P|\tilde{P})$ for arbitrary probability measures P and \tilde{P} . If P is a measure concentrated at two points with probabilities x and $1-x$, while \tilde{P} is a measure concentrated at the same points with probabilities y and $1-y$, then Definition 2.3.1 and equality (4.2.19) imply that $I(P|\tilde{P}) = H(x|y)$.

REMARK 4.2.5. If the assumptions of Lemma 4.2.5 are satisfied, then considering the sign of $E_{\theta_i}\lambda_1$, $i = 1, 2$, one can rewrite inequalities (4.2.17) and (4.2.18) in the form

$$(4.2.30) \quad E_{\theta_1}\nu \geq H(\alpha|1-\beta)|E_{\theta_1}\lambda_1|^{-1},$$

$$(4.2.31) \quad E_{\theta_2}\nu \geq H(\beta|1-\alpha)|E_{\theta_2}\lambda_1|^{-1},$$

that is, inequality (4.2.31) follows from (4.2.30) if we interchange the hypotheses $H_1: \theta = \theta_1$ and $H_2: \theta = \theta_2$. The latter means that θ_2 substitutes θ_1 , β substitutes α , and α substitutes β . In a similar manner, inequality (4.2.31) follows from (4.2.30).

REMARK 4.2.6. One can evaluate approximate values of the expectations $E_{\theta_i}\nu$, $i = 1, 2$, as follows. Let $E_{\theta_i}|\lambda_1| < \infty$ and $P_{\theta_i}\{\lambda_1 \neq 0\} > 0$ for $i = 1, 2$. Then $E_{\theta_i}\nu < \infty$ for $i = 1, 2$ by Lemma 4.2.3 and

$$(4.2.32) \quad E_{\theta_i}\Lambda_\nu = E_{\theta_i}\nu E_{\theta_i}\lambda_1, \quad i = 1, 2,$$

by Lemma 4.2.4. On the other hand

$$\begin{aligned} E_{\theta_1} \Lambda_\nu &= E_{\theta_1} \{ \Lambda_\nu / \Lambda_\nu \geq a \} P_{\theta_1} \{ \Lambda_\nu \geq a \} + E_{\theta_1} \{ \Lambda_\nu / \Lambda_\nu \leq -b \} P_{\theta_1} \{ \Lambda_\nu \leq -b \} \\ &\approx a\alpha - b(1 - \alpha) \end{aligned}$$

where the approximation appears, since we neglect the exits of Λ_ν from the interval $(-b, a)$. This approximation and equality (4.2.32) yield

$$(4.2.33) \quad E_{\theta_1} \nu \approx (a\alpha - b(1 - \alpha)) / E_{\theta_1} \lambda_1.$$

Similarly we obtain

$$(4.2.34) \quad E_{\theta_2} \nu \approx (a(1 - \beta) - b\beta) / E_{\theta_2} \lambda_1.$$

If

$$a = \ln \frac{1 - \beta}{\alpha} \quad \text{and} \quad b = \ln \frac{1 - \alpha}{\beta},$$

then approximations (4.2.33) and (4.2.34) for $E_{\theta_1} \nu$ and $E_{\theta_2} \nu$ coincide with the lower bounds in (4.2.17) and (4.2.18), respectively, in the case of a Wald sequential test with limit points $(-b, a)$.

EXAMPLE 4.2.1. Let ξ_1, ξ_2, \dots be independent identically distributed random variables whose distribution under the hypothesis H_i is $\mathcal{N}(\theta_i, \sigma^2)$ normal, $i = 1, 2$, where $\theta_1 < \theta_2$ and the variance σ^2 is known. Then

$$\Lambda_n = \frac{\theta_2 - \theta_1}{2\sigma^2} \left(2 \sum_{i=1}^n \xi_i - n(\theta_2 + \theta_1) \right).$$

It is clear that $P_{\theta_i}(\lambda_1 = 0) = 0$ for $i = 1, 2$. Let ν be the stopping time of a Wald sequential test of power (α, β) with limit points $(-b, a)$. Since

$$E_{\theta_1} \lambda_1 = -\frac{(\theta_2 - \theta_1)^2}{2\sigma^2} \quad \text{and} \quad E_{\theta_2} \lambda_1 = \frac{(\theta_2 - \theta_1)^2}{2\sigma^2},$$

approximations (4.2.33) and (4.2.34) become of the form

$$(4.2.35) \quad E_{\theta_1} \nu \approx \frac{2(b(1 - \alpha) - a\alpha)\sigma^2}{(\theta_2 - \theta_1)^2},$$

$$(4.2.36) \quad E_{\theta_2} \nu \approx \frac{2(a(1 - \beta) - b\beta)\sigma^2}{(\theta_2 - \theta_1)^2}.$$

Consider a test with a fixed nonrandom size n of the sample and whose power is (α, β) . For example, the Neyman-Pearson test $\delta_n^{+, \alpha}$ of level α satisfies these conditions. Then

$$\alpha = P_{\theta_1} \{ \Lambda_n > c_\alpha \} = P_{\theta_1} \left\{ \sum_{i=1}^n (\xi_i - \theta_1) > z_{1-\alpha} \sigma \sqrt{n} \right\}$$

where z_p is a p -quantile of the $\mathcal{N}(0, 1)$ law, that is, $\Phi(z_p) = p$. The test $\delta_n^{+, \alpha}$ has the type II error probability β if

$$\begin{aligned}\beta &= P_{\theta_2} \left\{ \sum_{i=1}^n (\xi_i - \theta_1) < z_{1-\alpha} \sigma \sqrt{n} \right\} \\ &= P_{\theta_2} \left\{ \sum_{i=1}^n (\xi_i - \theta_2) < z_{1-\alpha} \sigma \sqrt{n} - n(\theta_2 - \theta_1) \right\}\end{aligned}$$

or, equivalently, if

$$z_{1-\alpha} \sigma - \sqrt{n}(\theta_2 - \theta_1) = z_\beta \sigma.$$

Since $z_{1-\alpha} = -z_\alpha$, we have

$$(4.2.37) \quad n = \frac{\sigma^2(z_\alpha + z_\beta)^2}{(\theta_2 - \theta_1)^2}.$$

Relations (4.2.35)–(4.2.37) provide the following approximations:

$$\frac{E_{\theta_1} \nu}{n} \approx \frac{2(b(1-\alpha) - a\alpha)}{(z_\alpha + z_\beta)^2}, \quad \frac{E_{\theta_2} \nu}{n} \approx \frac{2(a(1-\beta) - b\beta)}{(z_\alpha + z_\beta)^2}.$$

If $\alpha = \beta = 0.05$, then $z_\alpha = z_\beta \approx 1.6449$ and $a = -b \approx 2.9444$, whence

$$(4.2.38) \quad \frac{E_{\theta_i} \nu}{n} \approx 0.4897, \quad i = 1, 2.$$

It is seen from equalities (4.2.38) that the above Wald sequential test of power (0.05; 0.05) requires two times less observations than the Neyman–Pearson test of the same power (0.05; 0.05) and with a nonrandom size of the sample.

The fundamental identity of sequential analysis. First we prove two auxiliary results.

LEMMA 4.2.6. *Let ζ be a random variable defined on the main probability space (Ω, \mathcal{F}, P) and such that:*

- $P\{\zeta > 0\} > 0$ and $P\{\zeta < 0\} > 0$;
- $\varphi(t) = Ee^{t\zeta}$ exists for all $t \in (-\infty, \infty)$;
- $E\zeta \neq 0$.

Then there is a unique number $\tau \neq 0$ such that $\varphi(\tau) = 1$, and moreover $\tau < 0$ if $E\zeta > 0$ and $\tau > 0$ if $E\zeta < 0$.

PROOF. The condition $P\{\zeta > 0\} > 0$ implies that there is a constant $c > 0$ such that $P\{\zeta > c\} > 0$. Thus for all $t > 0$

$$\varphi(t) = Ee^{t\zeta} \geq EI(\zeta > c)e^{t\zeta} > e^{tc}P\{\zeta > c\},$$

whence $\varphi(t) \rightarrow \infty$ as $t \rightarrow \infty$. Similarly, the condition $P\{\zeta < 0\} > 0$ implies that $\varphi(t) \rightarrow \infty$ as $t \rightarrow -\infty$. Moreover, $\varphi(0) = 1$ and $\varphi'(0) = E\zeta \neq 0$. If $E\zeta > 0$, then $\varphi'(0) > 0$ and thus there is $\tau < 0$ such that $\varphi(\tau) = 1$. Similarly, if $E\zeta < 0$, then $\varphi'(0) < 0$ and thus there is $\tau > 0$ such that $\varphi(\tau) = 1$. It is easy to show that $\varphi''(t) = E\zeta^2 e^{t\zeta} > 0$, hence the function $\varphi(t)$ is strictly convex. The latter property means, in particular, that φ has a unique minimum. Therefore the solution τ is unique. \square

LEMMA 4.2.7. Let ν be the stopping time of a Wald sequential test with limit points $(-b, a)$. Assume that $P_\theta\{\nu < \infty\} = 1$ where either $\theta = \theta_1$ or $\theta = \theta_2$. Then

$$(4.2.39) \quad E_\theta e^{t\Lambda_\nu} (\varphi(t))^{-\nu} = P\{\nu < \infty / H_t\}$$

for t such that $\varphi(t) = E_\theta e^{t\lambda_1} < \infty$ where $P\{\cdot / H_t\}$ is the conditional probability given H_t where the hypothesis H_t is that the random variables ξ_1, ξ_2, \dots are independent, identically distributed, and whose density is

$$(4.2.40) \quad p(x/H_t) = \frac{e^{t\lambda(x)}}{\varphi(t)} p_\theta(x).$$

PROOF. Taking into account $P_\theta\{\nu < \infty\} = 1$ we conclude that

$$\begin{aligned} E_\theta e^{t\Lambda_\nu} (\varphi(t))^{-\nu} &= \sum_{n=1}^{\infty} \int_{V_n \cup W_n} (\varphi(t))^{-n} \exp\{t\Lambda_n(x^{(n)})\} \prod_{i=1}^n p_\theta(x_i) \mu^n(dx^{(n)}) \\ &= \sum_{n=1}^{\infty} \int_{V_n \cup W_n} \prod_{i=1}^n p_\theta(x_i/H_t) \mu^n(dx^{(n)}) = P\{\nu < \infty / H_t\} \end{aligned}$$

where V_n and W_n are the sets defined by (4.2.7) and (4.2.8), respectively, while $p(x/H_t)$ is the density of the distribution given by (4.2.40). Thus equality (4.2.39) is proved. \square

LEMMA 4.2.8. Let ν be the stopping time of a Wald sequential test with limit points $(-b, a)$. Assume that $P_\theta\{\lambda_1 \neq 0\} > 0$ where either $\theta = \theta_1$ or $\theta = \theta_2$. Then

$$(4.2.41) \quad E_\theta e^{t\Lambda_\nu} (\varphi(t))^{-\nu} = 1$$

for t such that $\varphi(t) = E_\theta e^{t\lambda_1} < \infty$.

PROOF. The condition $P_\theta\{\lambda_1 \neq 0\} > 0$ implies that $P_\theta\{\nu < \infty\} = 1$ by Lemma 4.2.3 and that $P_\theta\{\lambda(\xi_1) \neq 0 / H_t\} > 0$, since

$$P_\theta\{\lambda(\xi_1) \neq 0 / H_t\} = \frac{1}{\varphi(t)} \int_{\{\lambda(x) \neq 0\}} e^{t\lambda(x)} P_\theta(dx) \neq 0.$$

Applying Lemma 4.2.3 once more we get

$$P_\theta\{\nu < \infty / H_t\} = 1.$$

Now equality (4.2.41) follows from (4.2.39). \square

Equality (4.2.41) is called the *fundamental identity of sequential analysis*.

REMARK 4.2.7. Equalities (4.2.39) and (4.2.41) hold for t such that $\varphi(t) < \infty$. If $P_\theta\{\lambda_1 \neq 0\} > 0$, then it follows from Lemma 4.2.3 that there is $t_0 > 0$ such that $\varphi(t) < \infty$ for all $t < t_0$.

REMARK 4.2.8. Lemmas 4.2.7 and 4.2.8 hold for measures P_θ if θ is different from both θ_1 and θ_2 .

The fundamental identity of sequential analysis can be applied to the problem of distinguishing composite hypotheses.

Let ξ_1, ξ_2, \dots be independent identically distributed random variables with a distribution P_θ depending on an unknown parameter $\theta \in \Theta$. Assume that P_θ is absolutely continuous with respect to some σ -finite measure μ and the density is $p(x; \theta)$. Let $\Theta = \Theta_1 \cup \Theta_2$ and $\Theta_1 \cap \Theta_2 = \emptyset$. Consider the problem of distinguishing the hypotheses $H_1: \theta \in \Theta_1$ and $H_2: \theta \in \Theta_2$ by observations ξ_1, ξ_2, \dots with the help of a Wald sequential test that distinguishes two simple hypotheses $H'_1: \theta = \theta_1$ and $H'_2: \theta = \theta_2$ where $\theta_1 \in \Theta_1$ and $\theta_2 \in \Theta_2$ are some fixed points. Let the random variable $\zeta = \lambda_1$ satisfy the assumptions of Lemma 4.2.6 concerning the measure P_θ and let $\theta \in \Theta$ be an arbitrary fixed point. According to Lemma 4.2.6, there is a number $\tau(\theta) \neq 0$ such that $\varphi_\theta(\tau(\theta)) = 1$ where $\varphi_\theta(t) = E_\theta e^{t\lambda_1}$. Further

$$(4.2.42) \quad E_\theta e^{\tau(\theta)\Lambda_\nu} = 1$$

by Lemma 4.2.8 and Remark 4.2.8.

The function $\beta(\theta) = P_\theta\{\Lambda_\nu \leq -b\}$ is the probability to accept the hypothesis H_1 if the parameter is θ , that is, $\beta(\theta)$ is the power function of the Wald test for distinguishing the hypotheses H_1 and H_2 . The function $\beta(\theta)$ is also called the *operating characteristic of the test* in sequential analysis.

Put

$$(4.2.43) \quad E_\theta^* = E_\theta\{e^{\tau(\theta)\Lambda_\nu} / \Lambda_\nu \leq -b\}, \quad E_\theta^{**} = E_\theta\{e^{\tau(\theta)\Lambda_\nu} / \Lambda_\nu \geq a\}.$$

It follows from equalities (4.2.42), (4.2.43) and (4.2.21) for $\theta_i = \theta$ that

$$(4.2.44) \quad 1 = E_\theta e^{\tau(\theta)\Lambda_\nu} = \beta(\theta)E_\theta^* + (1 - \beta(\theta))E_\theta^{**}.$$

Using the approximations

$$E_\theta^* \approx e^{-b\tau(\theta)}, \quad E_\theta^{**} \approx e^{a\tau(\theta)}$$

we derive from (4.2.44) an approximation for the operating characteristic:

$$\beta(\theta) \approx \frac{1 - e^{a\tau(\theta)}}{e^{-b\tau(\theta)} - e^{a\tau(\theta)}}, \quad \theta \in \Theta.$$

Similarly one can obtain approximations for $E_\theta\nu$, $\theta \in \Theta$.

More details on sequential Wald tests and their properties as well as on the other sequential tests can be found in [15, 51, 54]. The sequential analysis is described in [13, 46].

4.3. The optimality of a sequential Wald test

The main theorem. As in the preceding section we consider the problem of distinguishing two simple hypotheses $H_1: \theta = \theta_1$ and $H_2: \theta = \theta_2$ by observations ξ_1, ξ_2, \dots where ξ_1, ξ_2, \dots are independent identically distributed random variables whose distribution P_θ depends on an unknown parameter θ . Moreover we assume that their distribution possesses the density $p(x; \theta)$ with respect to a σ -finite measure μ . Throughout this section we also assume that $P_{\theta_i}\{\lambda_1 \neq 0\} > 0$ for $i = 1, 2$.

Generally speaking, we consider a sequential test for distinguishing the hypotheses H_1 and H_2 by observations $\xi^{(\nu)} = (\xi_1, \xi_2, \dots, \xi_\nu)$ where ν is a *stopping*

time. A stopping time can be either random or deterministic. Sequential tests are also called *sequential decision procedures*. *Decision functions* $d_\nu = d_\nu(\xi^{(\nu)})$ are defined as follows: if $\nu = n$, then $d_n(x^{(n)})$ assumes only two values d_1 and d_2 . If $d_n(x^{(n)}) = d_1$, then the hypothesis H_1 is accepted, while H_2 is accepted if $d_n(x^{(n)}) = d_2$. We define the type I and type II *error probabilities* for a sequential test with a decision function d_ν by

$$(4.3.1) \quad \alpha_1(d_\nu) = P_{\theta_1}\{d_\nu = d_2\}, \quad \alpha_2(d_\nu) = P_{\theta_2}\{d_\nu = d_1\}.$$

Note that

$$(4.3.2) \quad \alpha_1 = \alpha_1(d_\nu^*) = P_{\theta_1}\{d_\nu^* = d_2\} = P_{\theta_1}\{\Lambda_\nu \geq a\},$$

$$(4.3.3) \quad \alpha_2 = \alpha_2(d_\nu^*) = P_{\theta_2}\{d_\nu^* = d_1\} = P_{\theta_2}\{\Lambda_\nu \leq -b\}$$

for a Wald test of power (α_1, α_2) with limit points $(-b, a)$ and decision function d_ν^* .

The following result asserts that a sequential Wald test is optimal in the sense that $E_{\theta_1}\nu$ and $E_{\theta_2}\nu$ are minimal for it.

THEOREM 4.3.1. *The sequential Wald test of power (α_1, α_2) minimizes both expectations $E_{\theta_1}\nu$ and $E_{\theta_2}\nu$ in the set of all tests (including nonsequential tests) such that $E_{\theta_1}\nu$ and $E_{\theta_2}\nu$ are finite and*

$$(4.3.4) \quad P_{\theta_1}\{d_\nu = d_2\} \leq \alpha_1, \quad P_{\theta_2}\{d_\nu = d_1\} \leq \alpha_2.$$

To prove Theorem 4.3.1 we consider an auxiliary Bayes problem and use it to show that the Wald test is optimal.

An auxiliary problem. Consider the following sequential Bayes problem for distinguishing the hypotheses $H_1: \theta = \theta_1$ and $H_2: \theta = \theta_2$. Let $w_i > 0$ be the loss caused by a wrong decision given the hypothesis H_i is true and let the loss caused by a correct decision be zero. Assume that the cost of every observation is $c > 0$. The risk of the sequential test δ when making a decision given the hypothesis H_i is true equals

$$\alpha_i w_i + c E_{\theta_i} \nu, \quad i = 1, 2,$$

where ν is the stopping time of the sequential test δ and α_1 and α_2 are the type I and type II error probabilities, respectively. The risk of a test includes the mean loss caused by making a decision and the mean cost per observation. Let $q = P\{\theta = \theta_1\}$ and $1 - q = P\{\theta = \theta_2\}$ be a priori probabilities of the hypotheses H_1 and H_2 . Then the (unconditional) risk of the test δ is

$$(4.3.5) \quad r(q, \delta) = q(\alpha_1 w_1 + c E_{\theta_1} \nu) + (1 - q)(\alpha_2 w_2 + c E_{\theta_2} \nu).$$

DEFINITION 4.3.1. A sequential test δ^* is called *q-Bayes* if $r(q, \delta^*) \leq r(q, \delta)$ for all tests δ where $q \in [0, 1]$ is given and $r(q, \delta)$ is the risk of a test δ defined by (4.3.5).

DEFINITION 4.3.2. A sequential test δ^* is called *Bayes* if $r(q, \delta^*) \leq r(q, \delta)$ for all tests δ and all $q \in [0, 1]$.

The *q-Bayes* test for the above auxiliary problem is described in the following result.

LEMMA 4.3.1. Let $q' \leq q''$ be solutions of the equations

$$(4.3.6) \quad r(q', \delta_1) = \rho(q'), \quad r(q'', \delta_2) = \rho(q''),$$

respectively (provided the solutions exist), where

$$(4.3.7) \quad \rho(q) = \inf_{\delta \in \mathcal{S}_1} r(q, \delta),$$

\mathcal{S}_1 is the class of tests for which it is necessary to take at least one observation, and δ_i is the test rejecting the hypothesis H_i without sampling. Put

$$(4.3.8) \quad q' = q'' = \frac{w_2}{w_1 + w_2}$$

for the case where equations (4.3.6) have no solutions. If $0 < q' < q'' < 1$, then for all $q \in (q', q'')$ the Wald sequential test with limit points

$$(4.3.9) \quad -b = \ln \left(\frac{q}{1-q} \frac{1-q''}{q''} \right), \quad a = \ln \left(\frac{q}{1-q} \frac{1-q'}{q'} \right)$$

is q -Bayes.

PROOF. Step I. First we find q for which the better decision is the one made without sampling. We get from (4.3.5) that

$$r(q, \delta_1) = qw_1, \quad r(q, \delta_2) = (1-q)w_2.$$

Further, (4.3.5) and (4.3.7) for all $\lambda \in (0, 1)$ and all $q_1, q_2 \in [0, 1]$ imply that

$$\rho(\lambda q_1 + (1-\lambda)q_2) = \inf_{\delta \in \mathcal{S}_1} [\lambda r(q_1, \delta) + (1-\lambda)r(q_2, \delta)] \geq \lambda \rho(q_1) + (1-\lambda)\rho(q_2),$$

that is, $\rho(q)$ is a convex function. Since $\rho(q) \geq 0$, the function $\rho(q)$ is continuous in the interval $(0, 1)$.

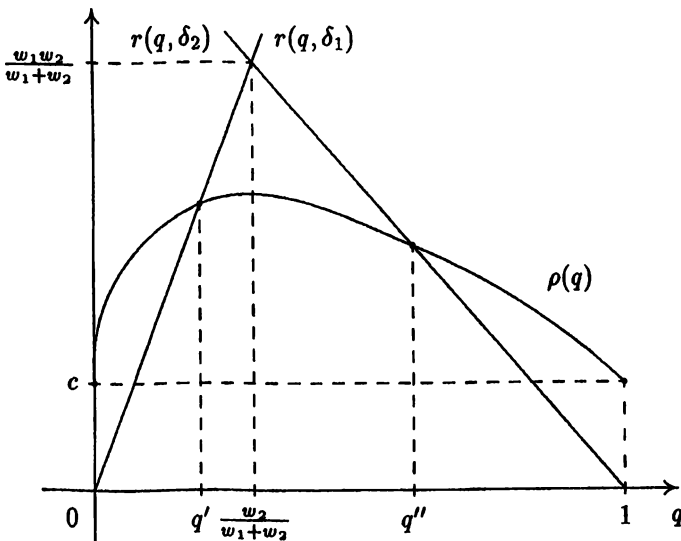


FIGURE 4.3.1. Graphs of the functions $\rho(q)$, $r(q, \delta_1)$, and $r(q, \delta_2)$

The graphs of the functions $\rho(q)$, $r(q, \delta_1)$, and $r(q, \delta_2)$ are shown in Figure 4.3.1. If

$$(4.3.10) \quad \rho\left(\frac{w_1 w_2}{w_1 + w_2}\right) < \frac{w_1 w_2}{w_1 + w_2} = r\left(\frac{w_2}{w_1 + w_2}, \delta_1\right) = r\left(\frac{w_2}{w_1 + w_2}, \delta_2\right),$$

then solutions q' and q'' of equations (4.3.6) exist. Otherwise we accept convention (4.3.8). Assume that relation (4.3.10) holds. Then $0 < q' < q'' < 1$ and the test δ_1 minimizes $r(q, \delta)$ if and only if $q \leq q'$, while δ_2 minimizes $r(q, \delta)$ if and only if $q' \geq q''$. This implies that the unique optimal decision on the first step is as follows: if $q \leq q'$, then the hypothesis H_1 is rejected and the hypothesis H_2 is accepted without sampling; if $q \geq q''$, then the hypothesis H_2 is rejected and the hypothesis H_1 is accepted without sampling; if $q' < q < q''$, then it is necessary to take the first observation ξ_1 .

Step II. We use induction to complete the proof. Let $q' < q < q''$ and n observations $\xi_1 = x_1, \xi_2 = x_2, \dots, \xi_n = x_n$ be given. Then the procedure is the same as that described in Step I above. The loss nc caused by making n observations does not change the problem, since further observations cannot reimburse this loss. If the probability that the hypothesis H_1 is true does not exceed q' or is not less than q'' , then we terminate the sampling; otherwise it is necessary to take one more observation ξ_{n+1} . According to the Bayes formula the probability that the hypothesis H_1 is true given $\xi_1 = x_1, \xi_2 = x_2, \dots, \xi_n = x_n$ equals

$$q(x^{(n)}) = \frac{qp_n(x^{(n)}; \theta_1)}{qp_n(x^{(n)}; \theta_1) + (1 - q)p_n(x^{(n)}; \theta_2)}$$

where $x^{(n)} = (x_1, x_2, \dots, x_n)$ and $p_n(x^{(n)}; \theta)$ is the density of the vector

$$\xi^{(n)} = (\xi_1, \xi_2, \dots, \xi_n)$$

with respect to the measure μ^n . Thus we keep sampling if $q' < q(x^{(n)}) < q''$, that is, if

$$e^{-b} < z_n(x^{(n)}) = \frac{p_n(x^{(n)}; \theta_2)}{p_n(x^{(n)}; \theta_1)} < e^a$$

where b and a are the constants defined by (4.3.9). If $z_n(x^{(n)}) < e^{-b}$, then the hypothesis H_1 is accepted, while if $z_n(x^{(n)}) > e^a$, then the hypothesis H_2 is accepted. Thus we proved for $q' < q < q''$ that the q -Bayes test coincides with the Wald sequential test with limit points (4.3.9). \square

REMARK 4.3.1. In Step I of the proof of Lemma 4.3.1 we determined the q -Bayes procedure (now we denote it by δ^*) as follows: $\delta^* = \delta_1$ if $q < q'$, $\delta^* = \delta_2$ if $q > q''$, and δ^* requires the first observation if $q' < q < q''$. The test δ_1 minimizes risk (4.3.5) if $q = q'$. However δ_1 is not a unique optimal test, since there exists $\delta \in \mathcal{S}_1$ such that $r(q', \delta) = \rho(q')$. If $q = q'$ and it is necessary to take an observation ξ_1 , then we showed in Step II of the proof of Lemma 4.3.1 that there is a test in \mathcal{S}_1 that minimizes the risk. This means that it makes no difference how one constructs the test in the case of $q = q'$. The same is true for the case of $q = q''$. Moreover this also is true for the next steps. This therefore proves that if $q' \leq q \leq q''$, then the test coinciding with the Wald sequential test with limit points (4.3.9) is q -Bayes.

A relationship between the auxiliary and main problems is established in the following result.

LEMMA 4.3.2. *For all $0 < q'_0 < q''_0 < 1$ there are numbers $w \in (0, 1)$ and $c > 0$ such that the Bayes solution of the auxiliary problem with $w_1 = 1 - w$, $w_2 = w$, and with a priori probability $q \in (q'_0, q''_0)$ is the Wald sequential test with limit points $(-b_0, a_0)$ such that*

$$-b_0 = \ln \left(\frac{q}{1-q} \frac{1-q''_0}{q'_0} \right), \quad a_0 = \ln \left(\frac{q}{1-q} \frac{1-q'_0}{q''_0} \right).$$

PROOF. *Step I.* Let $q'(w, c)$ and $q''(w, c)$ be solutions of equations (4.3.6) where $r(q, \delta)$ is defined by (4.3.5) for $w_1 = 1 - w$ and $w_2 = w$. Thus we need to find w and c such that $q'(w, c) = q'_0$ and $q''(w, c) = q''_0$. Given a fixed w let

$$q'(c) = q'(w, c), \quad q''(c) = q''(w, c).$$

Let $c_0 = c_0(w)$ be the minimal number c such that $q'(c) = q''(c)$. Then $q'(c)$ and $q''(c)$ for $0 < c < c_0$ are defined from the equations

$$(1-w)q' = \rho(q', c), \quad w(1-q'') = \rho(q'', c)$$

where $\rho(q, c)$ stands for $\rho(q)$ defined in (4.3.7).

Given a fixed q the function $\rho(q, c)$ of the argument c is such that

- 1) $\rho(q, c)$ is continuous with respect to c ;
- 2) $\rho(q, c)$ increases with respect to c , since for any $\delta \in \mathcal{S}_1$ the risk increases with respect to c and the minimal risk $\rho(q, c)$ is attained for the test $\delta \in \mathcal{S}_1$;
- 3) $\rho(q, c) \rightarrow 0$ as $c \rightarrow 0$.

The latter property holds, since the type I and type II error probabilities for samples with fixed size n can be arbitrarily small if n is sufficiently large.

The above properties of the function ρ imply that for $0 < c < c_0$ the function $q'(c)$ (respectively, $q''(c)$) is continuous, increasing (respectively, decreasing), and $q'(c) \rightarrow 0$ (respectively, $q''(c) \rightarrow 1$) as $c \rightarrow 0$. On the other hand,

$$q''(c) - q'(c) \rightarrow 0 \quad \text{as } c \rightarrow c_0,$$

so that both functions $q'(c)$ and $q''(c)$ approach the solution $q'(c) = q''(c) = w$ of the equation $q'(1-w) = (1-q')w$. The above properties also imply that for fixed w the function

$$\lambda(c) = \frac{q'(c)}{1-q'(c)} \cdot \frac{1-q''(c)}{q''(c)}$$

is continuous, increasing, and varying from 0 to 1 as c is varying from 0 to

$$c_0 = c_0(w).$$

Step II. Put

$$\lambda(w, c) = \frac{q'(w, c)}{1-q'(w, c)} \cdot \frac{1-q''(w, c)}{q''(w, c)}, \quad \gamma(w, c) = \frac{q''_0(w, c)}{1-q''_0(w, c)}.$$

We prove that there are w and c such that

$$\lambda(w, c) = \frac{q'_0}{1-q'_0} \cdot \frac{1-q''_0}{q''_0} = \lambda_0, \quad \gamma(w, c) = \frac{q''_0}{1-q''_0} = \gamma_0.$$

We proved in Step I that for all fixed w there exists a unique $c = c(w)$ such that $\lambda(w, c) = \lambda_0$. In Step III below we prove that the function $\gamma(w) = \gamma(w, c(w))$ is a one-to-one correspondence between $w \in (0, 1)$ and $\gamma \in (0, \infty)$. Therefore there exists a unique number $w \in (0, 1)$ such that $\gamma(w) = \gamma_0$. This will complete the proof of the lemma.

Step III. According to Lemma 4.3.1 for the auxiliary problem with $w_1 = 1 - w$ and $w_2 = w$, the cost per observation $c = c(w)$, and the a priori probability $q = q'(w, c(w))$, there exists a sequential q -Bayes test δ' which is a Wald sequential test with limit points $(-b', 0)$ where

$$-b' = \ln \left(\frac{q'(w, c(w))}{1 - q'(w, c(w))} \cdot \frac{1 - q''(w, c(w))}{q''(w, c(w))} \right) = \ln \lambda(w, c(w)) = \ln \lambda_0.$$

Further let δ'' be the Wald sequential test for the auxiliary problem with constants $w_1 = 1 - w$, $w_2 = w$, $c = c(w)$, and $q = q''(w, c(w))$ that is a Wald sequential test with limit points $(0, a'')$ where

$$a'' = \ln \left(\frac{q''(w, c(w))}{1 - q''(w, c(w))} \cdot \frac{1 - q'(w, c(w))}{q'(w, c(w))} \right) = \ln \frac{1}{\lambda_0}.$$

Then the error probabilities α'_1 and α'_2 and the expectations $E_{\theta_1} \nu'$ and $E_{\theta_2} \nu'$ of the test δ' as well as error probabilities α''_1 and α''_2 and expectations $E_{\theta_1} \nu''$ and $E_{\theta_2} \nu''$ of the test δ'' depend on w and c through λ_0 but not through γ . Thus they are fixed numbers for a fixed λ_0 . The Bayes risks for $q' = q'(w, c(w))$ and $q'' = q''(w, c(w))$ are equal to

$$\rho(q') = r(q', \delta'), \quad \rho = (q'')r(q'', \delta''),$$

respectively. Relations (4.3.6) imply that

$$r(q', \delta_1) = r(q', \delta'), \quad r(q'', \delta_2) = r(q'', \delta'').$$

The latter equalities can be rewritten as

$$\begin{aligned} q'(1 - w) &= q'[\alpha'_1(1 - w) + cE_{\theta_1} \nu'] + (1 - q')[\alpha'_2 w + cE_{\theta_2} \nu'], \\ (1 - q'')w &= q''[\alpha''_1(1 - w) + cE_{\theta_1} \nu''] + (1 - q'')[\alpha''_2 w + cE_{\theta_2} \nu'']. \end{aligned}$$

Using

$$\frac{q'}{1 - q'} = \lambda_0 \gamma, \quad \frac{q''}{1 - q''} = \gamma$$

in the latter equalities and excluding c we obtain

$$\begin{aligned} \{ \lambda_0 \gamma (1 - \alpha'_1) - w [\lambda_0 \gamma (1 - \alpha'_1) + \alpha'_2] \} (\gamma E_{\theta_1} \nu'' + E_{\theta_2} \nu'') \\ = \{ -\gamma \alpha''_1 + w [(1 - \alpha''_2) + \gamma \alpha''_1] \} (\lambda_0 \gamma E_{\theta_1} \nu' + E_{\theta_2} \nu'). \end{aligned}$$

This equation is linear with respect to w , thus it has a solution $w \in (0, 1)$ for all $\gamma > 0$. Collecting all the terms on one side of this equality we obtain a polynomial of the second degree with respect to γ such that the coefficient of γ^2 and constant term have different signs if $w \in (0, 1)$. Thus there exists a unique positive solution γ which is the desired one-to-one correspondence between γ and w . \square

REMARK 4.3.2. Property 3) of the function $\rho(q, c)$ mentioned in Step I of the proof of Lemma 4.3.2 follows from the following reasoning. By assumption

$$P_{\theta_i}\{\lambda_1 \neq 0\} > 0$$

for $i = 1$ and $i = 2$. Then $E_{\theta_1}\lambda_1 < 0$ and $E_{\theta_2}\lambda_1 > 0$ (see (4.2.26)). If $E_{\theta_1}\lambda_1 > -\infty$, then the Khinchine law of large numbers implies that $n^{-1}\Lambda_n \rightarrow E_{\theta_1}\lambda_1$ as $n \rightarrow \infty$ in probability P_{θ_1} . Further if $E_{\theta_1}\lambda_1 = -\infty$, then one can prove that $n^{-1}\Lambda_n \rightarrow -\infty$ as $n \rightarrow \infty$ in probability P_{θ_1} . Thus $\Lambda_n \rightarrow -\infty$ as $n \rightarrow \infty$ in probability P_{θ_1} . Then the type I error probabilities $\alpha_1(\delta_n^q)$ approach zero as $n \rightarrow \infty$ for the Bayes test δ_n^q constructed from a sample of size n by inequalities (2.3.73). Similarly we obtain that $\alpha_2(\delta_n^q) \rightarrow 0$ as $n \rightarrow \infty$. Thus $\rho(q, c) \rightarrow 0$ as $c \rightarrow 0$.

Proof of the main theorem. Now we use Lemmas 4.3.1 and 4.3.2 to prove the main result that the Wald test is optimal.

PROOF OF THEOREM 4.3.1. Consider the Wald sequential test of power (α_1, α_2) with limit points $(-b, a)$ where $a > 0$ and $b > 0$. Let ν be the stopping time of this test. Consider an arbitrary number q of the interval $(0, 1)$ and put

$$q' = \frac{q}{e^a(1-q) + q}, \quad q'' = \frac{q}{e^{-b}(1-q) + q}.$$

The numbers q' and q'' satisfy (4.3.9) and moreover $0 < q' < q < q'' < 1$. According to Lemma 4.3.2 there are numbers $w \in (0, 1)$ and $c > 0$ such that this test is a Bayes solution of the auxiliary problem for which a priori probabilities of the hypotheses H_1 and H_2 are q and $1 - q$, the loss due to a wrong decision is $w_1 = 1 - w$ and $w_2 = w$, respectively, and the cost per observation is c . Consider an arbitrary test δ^* (not necessarily sequential) with error probabilities α_1^* and α_2^* and the stopping time ν^* where $\alpha_i^* \leq \alpha_i$ and $E_{\theta_i}\nu^* < \infty$ for $i = 1, 2$. Again by Lemma 4.3.2

$$(4.3.11) \quad \begin{aligned} & q[(1-w)\alpha_1 + cE_{\theta_1}\nu] + (1-q)[w\alpha_2 + cE_{\theta_2}\nu] \\ & \leq q[(1-w)\alpha_1^* + cE_{\theta_1}\nu^*] + (1-q)[w\alpha_2^* + cE_{\theta_2}\nu^*] \\ & \leq q[(1-w)\alpha_1 + cE_{\theta_1}\nu^*] + (1-q)[w\alpha_2 + cE_{\theta_2}\nu^*] \end{aligned}$$

where the latter inequality holds, since $\alpha_1^* \leq \alpha_1$ and $\alpha_2^* \leq \alpha_2$ by condition. Then inequalities (4.3.11) imply

$$(4.3.12) \quad qE_{\theta_1}\nu + (1-q)E_{\theta_2}\nu \leq qE_{\theta_1}\nu^* + (1-q)E_{\theta_2}\nu^*.$$

Since (4.3.12) holds for all $q \in (0, 1)$, we pass to the limit as $q \rightarrow 0$ and obtain from (4.3.12) that $E_{\theta_2}\nu \leq E_{\theta_2}\nu^*$. Similarly we pass to the limit as $q \rightarrow 1$ and obtain from (4.3.12) that $E_{\theta_1}\nu \leq E_{\theta_1}\nu^*$. \square

REMARK 4.3.3. In the proof above we constructed a q -Bayes sequential test for distinguishing the hypotheses H_1 and H_2 with a priori distribution $(q, 1 - q)$ of the hypotheses and for the loss matrix

$$A = \begin{pmatrix} 0 & w_1 \\ w_2 & 0 \end{pmatrix}, \quad w_1 > 0, w_2 > 0.$$

The general case of the problem of constructing the Wald sequential tests is reduced to the solution of the Bellman equation (also known as the optimality equation in dynamic programming) (see Section 4.1 and [13, 15, 46]).

References to Part 2

1. R. R. Bahadur, *Some limit theorems in statistics*, SIAM, Philadelphia, PA, 1971.
2. J.-R. Barra, *Notions fondamentales de statistique mathématique*, Dunod, Paris, 1971; English transl., *Mathematical basis of statistics*, Academic Press, New York–London, 1981.
3. P. Billingsley, *Convergence of probability measures*, Wiley, New York–London–Sydney, 1968.
4. D. Blackwell and M. A. Girshick, *Theory of games and statistical decisions*, Wiley and Chapman and Hall, New York and London, 1954.
5. M. V. Boldin, G. I. Simonova, and Yu. N. Tyurin, *Sign-based methods in linear statistical models*, “Nauka”, Moscow, 1997; English transl., Amer. Math. Soc., Providence, RI, 1997.
6. L. N. Bol’shev and N. V. Smirnov, *Tables of mathematical statistics*, “Nauka”, Moscow, 1965. (Russian)
7. A. A. Borovkov, *Mathematical statistics. Estimation of parameters. Testing of hypotheses*, “Nauka”, Moscow, 1984; English transl., *Mathematical statistics*, Gordon & Breach, Amsterdam, 1998.
8. ———, *Mathematical statistics. Supplementary chapters*, “Nauka”, Moscow, 1984; English transl., Gordon & Breach, Amsterdam, 1998.
9. ———, *Mathematical statistics*, “Nauka”, Sibirskoe otdelenie RAN, Novosibirsk, 1997. (Russian)
10. A. A. Borovkov and A. A. Mogul’skiĭ, *Large deviations and the testing of statistical hypotheses*, Proceedings of the Institute of Mathematics, 19, “Nauka”, Sibirskoe otdelenie RAN, Novosibirsk, 1992. (Russian)
11. N. N. Čencov, *Statistical decision rules and optimal inference*, “Nauka”, Moscow, 1972; English transl., Amer. Math. Soc., Providence, RI, 1982.
12. D. M. Chibisov, *Certain tests of the chi-square type for continuous distributions*, Teor. Veroyatnost. i Primenen. **16** (1971), no. 1, 3–20; English transl. in Theor. Probability Appl. **16** (1971), no. 1, 1–22.
13. Y. S. Chow, H. Robbins, and D. Siegmund, *The theory of optimal stopping*, Corrected reprint of the 1971 original, Dover, New York, 1991.
14. H. Cramér, *Mathematical methods of statistics*, reprint of the 1946 original, Princeton Univ. Press, Princeton, NJ, 1999.
15. M. H. DeGroot, *Optimal statistical decisions*, McGraw-Hill, New York–London–Sydney, 1970.
16. D. Dugué, *Traité statistique théorique et appliquée: analyse aléatoire, algèbre aléatoire*, Masson et Cie, Paris, 1958. (French)
17. R. S. Ellis, *Entropy, large deviations, and statistical mechanics*, Springer-Verlag, Berlin, 1985.
18. W. Feller, *An introduction to probability theory and its applications*, Third edition, vol. 1, Wiley, New York–London–Sydney, 1968; vol. 2, 1971.
19. I. I. Gikhman, *An introduction to the general theory of measure and integral*, Donetsk University Press, Donetsk, 1971. (Russian)
20. B. V. Gnedenko and A. N. Kolmogorov, *Limit distributions for sums of independent random variables*, Gostekhizdat, Leningrad–Moscow, 1949; English transl., Addison-Wesley, Reading, MA, 1968.
21. P. E. Greenwood and A. N. Shiryaev, *Contiguity and the statistical invariance principle*, Gordon & Breach, New York, 1985.
22. J. Hájek and Z. Šidák, *Theory of rank tests*, Academic Press and Academia Publishing House of the Czechoslovak Academy of Sciences, New York–London and Prague, 1967.
23. P. R. Halmos, *Measure theory*, Van Nostrand, New York, 1950.
24. P.-L. Hennequin and A. Tortrat, *Théorie des probabilités et quelques applications*, Masson et Cie, Éditeurs, Paris, 1965. (French)

25. I. A. Ibragimov and R. Z. Khas'minskiĭ, *Statistical estimation. Asymptotic theory*, "Nauka", Moscow, 1979; English transl., Springer-Verlag, New York–Berlin, 1981.
26. G. I. Ivchenko and Yu. I. Medvedev, *Mathematical statistics*, "Vysshaya shkola", Moscow, 1984. (Russian)
27. ———, *Decomposable statistics and hypotheses testing for grouped data*, Teor. Veroyatnost. i Primenen. **25** (1980), no. 3, 549–560; English transl. in *Theory Probab. Appl.* **25** (1981), no. 3, 540–551.
28. L. Jacod and A. N. Shiryaev, *Limit theorems for stochastic processes*, 2nd edition, Berlin, Springer-Verlag, 2003.
29. M. G. Kendall and A. Stuart, *The advanced theory of statistics. Inference and relationship*, vol. 2, Griffin, London, 1961.
30. N. Kligene and L. Telksnis, *Methods of detecting instants of change of random processes properties*, Avtomat. i Telemekh. (1983), no. 10, 5–56; English transl. in *Automat. Remote Control* **44** (1984), no. 10, part 1, 1241–1283.
31. A. N. Kolmogorov and S. V. Fomin, *Introductory real analysis*, "Nauka", Moscow, 1968; English transl., Prentice-Hall, New York, 1970.
32. A. M. Kolodĭĭ, *Basics of the general theory of measure and integral*, Volgograd University Press, Volgograd, 1999. (Russian)
33. S. Kullback, *Information theory and statistics*, Dover, New York, 1968.
34. E. L. Lehmann, *Theory of point estimation*, Wiley, New York, 1983.
35. F. Liese and I. Vajda, *Convex statistical distances*, Teubner, Leipzig, 1987.
36. Yu. N. Lin'kov, *The asymptotic distinguishability of two simple statistical hypotheses*, Preprint 86.45, Institute of Mathematics of Ukrainian Academy of Sciences, Kiev, 1986.
37. ———, *Asymptotic statistical methods for stochastic processes*, "Naukova dumka", Kiev, 1993; English transl., Amer. Math. Soc., Providence, RI, 2001.
38. ———, *Lectures on mathematical statistics*, vol. 1, "Istoki", Donetsk, 1999; English transl., Part 1 of this book.
39. ———, *Large deviation theorems for extended random variables and some applications*, J. Math. Sci. **93** (1999), no. 4, 563–573.
40. G. V. Martynov, *Omega-square tests*, "Nauka", Moscow, 1978. (Russian)
41. I. P. Natanson, *Theory of functions of real variable*, GITTL, Moscow, 1957; English transl., vol. 1, Ungar, New York, 1955; vol. 2, 1961.
42. J. Neveu, *Mathematical foundations of the calculus of probability*, Masson et Cie, Éditeurs, Paris, 1964; English transl., Holden-Day, San Francisco–London–Amsterdam, 1965.
43. C. R. Rao, *Statistical inference and its applications*, Wiley, New York–London–Sydney, 1965.
44. R. T. Rockafellar, *Convex analysis*, Princeton Univ. Press, Princeton, NJ, 1970.
45. G. G. Roussas, *Contiguity of probability measures: some applications in statistics*, Cambridge Univ. Press, London–New York, 1972.
46. A. N. Shiryaev, *Optimal stopping rules*, "Nauka", Moscow, 1976; English transl., Springer-Verlag, New York–Heidelberg, 1978.
47. ———, *Probability*, "Nauka", Moscow, 1989; English transl., Springer-Verlag, New York, 1996.
48. A. V. Skorokhod, *Random processes with independent increments*, "Nauka", Moscow, 1964; English transl., Kluwer, Dordrecht, 1991.
49. J.-L. Soler, *Notion de liberté en statistique mathématique*, Thèse de Docteur de Troisième Cycle, Université de Grenoble, 1970. (French)
50. F. P. Tarasenko, *Nonparametric statistics*, Tomsk University Press, Tomsk, 1976. (Russian)
51. A. Wald, *Sequential analysis*, Wiley and Chapman and Hall, New York and London, 1947.
52. A. Wald, *Statistical decision functions*, Wiley and Chapman and Hall, New York and London, 1950.
53. S. S. Wilks, *Mathematical statistics*, Wiley, New York–London, 1967.
54. S. Zacks, *The theory of statistical inference*, Wiley, New York–London–Sydney, 1971.

Index

- σ -algebra
 - sufficient, 102, 116
 - minimal, 116
- a priori probability, 176
- Bayes approach
 - complete, 184
 - partial, 184
- Bayes estimation method, 143
- bias (of the estimator), 40
- canonical equation, 33
- conditional expectation, 99
- conditional probability, 100
- confidence bounds, 40
- confidence interval, 39, 147
- confidence level, 147
- confidence limits, 147
- confidence probability, 40, 147
- confidence region, 150
 - asymptotic, 150
- confidence set, 194
 - uniformly most precise, 196
 - unbiased, 198
- convergence
 - weak, 12
- correlation coefficient, 27
 - sampling, 29
- Cramér–Rao bound, 88
- critical set, 263
- decision function, 159, 311
- distance
 - in variance, 205
 - Kakutani–Hellinger, 206
- distribution
 - chi-square, 21
 - Fisher, 24
 - least favorable, 179
 - Snedekor, 24
 - standard normal, 20
 - Student, 24
- distribution function
 - empirical, 5, 25
 - Kolmogorov, 8, 264
- entropy
 - relative, 218
- error probability, 203
 - of type I, 159, 186, 263
 - of type II, 160
- estimator
 - absolutely admissible, 44
 - admissible, 44
 - asymptotically Bayes, 96
 - asymptotically efficient, 77, 85
 - in the strong (weak) sense, 77
 - asymptotically minimax, 97
 - asymptotically R-Bayes, 96
 - asymptotically unbiased, 40
 - Bayes, 45
 - a posteriori, 142
 - generalized, 45
 - consistent, 42
 - efficient, 76, 85
 - equivariant, 49, 57
 - likelihood, 134
 - maximum likelihood, 134
 - polynomial, 273
 - minimax, 146
 - optimal, 44
 - Pitman, 50, 57
 - point, 39
 - statistical, 39
 - strongly consistent, 43
 - superefficient, 90
 - unbiased, 40
- excess, 12
- families of hypotheses
 - completely asymptotically distinguishable, 208
 - completely asymptotically indistinguishable, 211
 - mutually contiguous, 215
 - mutually noncontiguous, 215
- family of functions
 - dense, 213
 - uniformly integrable, 214

- family of hypotheses
 - contiguous, 213
- family of measures
 - complete, 118
 - dominated by a measure, 102
 - exponential, 122
 - relatively compact, 248
 - tight, 248
- Fisher information, 62, 257
 - matrix, 80
- Gamma distribution, 16
- Hellinger integral, 206
- hypothesis, 159
 - composite, 159
 - main, 263
 - null, 263
 - one-sided, 187
 - simple, 159
 - two-sided, 187
- inequality
 - Barankin–Kiefer, 88
 - Bhattacharyya, 87
 - Chapman–Robbins, 75, 88
 - Cramér–Rao, 68
 - matrix analog, 81
- Kullback–Leibler divergence, 218
- least variance, 33
- lemma
 - Neyman–Pearson, 167
 - Stein, 223
- likelihood function, 133
 - logarithmic, 133
- likelihood ratio, 163
- location parameter, 49
- mean square approximation, 32
- measure
 - absolutely continuous, 161
- measures
 - equivalent, 161
 - singular, 162
- method of moments, 131
- minimax, 179
- mixed moment, 26
 - central, 26
 - sampling, 28
 - sampling, 28
- moment, 8
 - central, 8, 29
 - sampling, 9, 30
 - central, 9, 30
- Neyman–Fisher factorization criterion, 103
- observation, 101, 159
- operating characteristic (of a test), 310
- order (of the moment), 29
- order statistic, 5
 - central, 16
- power function, 182
- quantile, 16, 240
 - sampling, 16
- random variable
 - uncorrelated, 26
- random vector
 - Gaussian, 20
 - normal, 20
- random walk, 279
- rank statistic, 286
- reflection method, 279
- regression, 31
 - linear, 32
 - coefficient of, 32
 - sampling, 35
 - sampling coefficient of, 35
 - parabolic, 34
 - sampling, 37
- regularity conditions
 - Cramér–Rao (CR), 61
 - Cramér–Rao (CR)*, 69
- relative stability, 219
- risk
 - a posteriori, 142
 - of the estimator, 45
 - of the test, 199
- risk function, 44
- sample, 25
- sampling space, 39
- scale parameter, 56
- sequence
 - asymptotically normal, 7
- Sheppard correction, 275
- skewness, 12
- Spearman rank correlation coefficient, 289
- statistic, 39, 101
 - complete, 118
 - minimal, 116
 - of the test, 263
 - subordinated, 116
 - sufficient, 101
- statistics
 - equivalent, 116
- stopping rule, 293
 - Bayes, 294
 - truncated, 295

test

- Bayes, 166, 175, 183, 311
- chi-square, 270
- empty blocks, 285
- empty boxes, 283
- for independence, 288
- goodness-of-fit, 263
 - Kolmogorov, 264
 - Pearson, 270
 - Smirnov, 281
 - symmetric, 282
- Kendall, 290
- likelihood ratio, 204
- Mann–Whitney, 287
- maximum likelihood, 167, 182
- minimax, 167, 184
- Moran, 291
- Neyman–Pearson, 170, 204
- nonrandomized, 159, 174
- of series, 286
- Pearson, 273
- q -Bayes, 311
- quantile, 271
- randomized, 159, 174
- rank, 286
- sequential, 293
 - Wald, 300
- sign, 271
- Spearman, 289
- statistical, 174
- unbiased, 191
- uniformly more powerful, 183
- uniformly most powerful (UMP), 183
- von Mises–Smirnov, 291
- Wilcoxon, 287

theorem

- Glivenko, 6
 - Kolmogorov, 8, 264
 - Le Cam, first, 251
 - Lehmann–Scheffé, 121
 - Pearson, 267
 - Rao–Blackwell–Kolmogorov, 113
- trajectory, 277
-
- Wald identity, 304

Titles in This Series

- 229 Yu. N. Lin'kov, Lectures in mathematical statistics, 2005
- 228 D. Zhelobenko, Principal structures and methods of representation theory, 2005
- 227 Takahiro Kawai and Yoshitsugu Takei, Algebraic analysis of singular perturbation theory, 2005
- 226 V. M. Manuilov and E. V. Troitsky, Hilbert C^* -modules, 2005
- 225 S. M. Natanzon, Moduli of Riemann surfaces, real algebraic curves, and their superanalogues, 2004
- 224 Ichiro Shigekawa, Stochastic analysis, 2004
- 223 Masatoshi Noumi, Painlevé equations through symmetry, 2004
- 222 G. G. Magaril-Il'yaev and V. M. Tikhomirov, Convex analysis: Theory and applications, 2003
- 221 Katsuei Kenmotsu, Surfaces with constant mean curvature, 2003
- 220 I. M. Gelfand, S. G. Gindikin, and M. I. Graev, Selected topics in integral geometry, 2003
- 219 S. V. Kerov, Asymptotic representation theory of the symmetric group and its applications to analysis, 2003
- 218 Kenji Ueno, Algebraic geometry 3: Further study of schemes, 2003
- 217 Masaki Kashiwara, D -modules and microlocal calculus, 2003
- 216 G. V. Badalyan, Quasipower series and quasianalytic classes of functions, 2002
- 215 Tatsuo Kimura, Introduction to prehomogeneous vector spaces, 2003
- 214 L. Š. Grinblat, Algebras of sets and combinatorics, 2002
- 213 V. N. Sachkov and V. E. Tarakanov, Combinatorics of nonnegative matrices, 2002
- 212 A. V. Mel'nikov, S. N. Volkov, and M. L. Nechaev, Mathematics of financial obligations, 2002
- 211 Takeo Ohsawa, Analysis of several complex variables, 2002
- 210 Toshitake Kohno, Conformal field theory and topology, 2002
- 209 Yasumasa Nishiura, Far-from-equilibrium dynamics, 2002
- 208 Yukio Matsumoto, An introduction to Morse theory, 2002
- 207 Ken'ichi Ohshika, Discrete groups, 2002
- 206 Yuji Shimizu and Kenji Ueno, Advances in moduli theory, 2002
- 205 Seiki Nishikawa, Variational problems in geometry, 2001
- 204 A. M. Vinogradov, Cohomological analysis of partial differential equations and Secondary Calculus, 2001
- 203 Te Sun Han and Kingo Kobayashi, Mathematics of information and coding, 2002
- 202 V. P. Maslov and G. A. Omel'yanov, Geometric asymptotics for nonlinear PDE. I, 2001
- 201 Shigeyuki Morita, Geometry of differential forms, 2001
- 200 V. V. Prasolov and V. M. Tikhomirov, Geometry, 2001
- 199 Shigeyuki Morita, Geometry of characteristic classes, 2001
- 198 V. A. Smirnov, Simplicial and operad methods in algebraic topology, 2001
- 197 Kenji Ueno, Algebraic geometry 2: Sheaves and cohomology, 2001
- 196 Yu. N. Lin'kov, Asymptotic statistical methods for stochastic processes, 2001
- 195 Minoru Wakimoto, Infinite-dimensional Lie algebras, 2001
- 194 Valery B. Nevzorov, Records: Mathematical theory, 2001
- 193 Toshio Nishino, Function theory in several complex variables, 2001
- 192 Yu. P. Solov'yov and E. V. Troitsky, C^* -algebras and elliptic operators in differential topology, 2001

TITLES IN THIS SERIES

- 191 **Shun-ichi Amari and Hiroshi Nagaoka**, *Methods of information geometry*, 2000
- 190 **Alexander N. Starkov**, *Dynamical systems on homogeneous spaces*, 2000
- 189 **Mitsuru Ikawa**, *Hyperbolic partial differential equations and wave phenomena*, 2000
- 188 **V. V. Buldygin and Yu. V. Kozachenko**, *Metric characterization of random variables and random processes*, 2000
- 187 **A. V. Fursikov**, *Optimal control of distributed systems. Theory and applications*, 2000
- 186 **Kazuya Kato, Nobushige Kurokawa, and Takeshi Saito**, *Number theory 1: Fermat's dream*, 2000
- 185 **Kenji Ueno**, *Algebraic Geometry 1: From algebraic varieties to schemes*, 1999
- 184 **A. V. Mel'nikov**, *Financial markets*, 1999
- 183 **Hajime Sato**, *Algebraic topology: an intuitive approach*, 1999
- 182 **I. S. Krasil'shchik and A. M. Vinogradov, Editors**, *Symmetries and conservation laws for differential equations of mathematical physics*, 1999
- 181 **Ya. G. Berkovich and E. M. Zhmud'**, *Characters of finite groups. Part 2*, 1999
- 180 **A. A. Milyutin and N. P. Osmolovskii**, *Calculus of variations and optimal control*, 1998
- 179 **V. E. Voskresenskiĭ**, *Algebraic groups and their birational invariants*, 1998
- 178 **Mitsuo Morimoto**, *Analytic functionals on the sphere*, 1998
- 177 **Satoru Igari**, *Real analysis—with an introduction to wavelet theory*, 1998
- 176 **L. M. Lerman and Ya. L. Umanskiy**, *Four-dimensional integrable Hamiltonian systems with simple singular points (topological aspects)*, 1998
- 175 **S. K. Godunov**, *Modern aspects of linear algebra*, 1998
- 174 **Ya-Zhe Chen and Lan-Cheng Wu**, *Second order elliptic equations and elliptic systems*, 1998
- 173 **Yu. A. Davydov, M. A. Lifshits, and N. V. Smorodina**, *Local properties of distributions of stochastic functionals*, 1998
- 172 **Ya. G. Berkovich and E. M. Zhmud'**, *Characters of finite groups. Part 1*, 1998
- 171 **E. M. Landis**, *Second order equations of elliptic and parabolic type*, 1998
- 170 **Viktor Prasolov and Yuri Solovyev**, *Elliptic functions and elliptic integrals*, 1997
- 169 **S. K. Godunov**, *Ordinary differential equations with constant coefficient*, 1997
- 168 **Junjiro Noguchi**, *Introduction to complex analysis*, 1998
- 167 **Masaya Yamaguti, Masayoshi Hata, and Jun Kigami**, *Mathematics of fractals*, 1997
- 166 **Kenji Ueno**, *An introduction to algebraic geometry*, 1997
- 165 **V. V. Ishkhanov, B. B. Lur'e, and D. K. Faddeev**, *The embedding problem in Galois theory*, 1997
- 164 **E. I. Gordon**, *Nonstandard methods in commutative harmonic analysis*, 1997
- 163 **A. Ya. Dorogovtsev, D. S. Silvestrov, A. V. Skorokhod, and M. I. Yadrenko**, *Probability theory: Collection of problems*, 1997
- 162 **M. V. Boldin, G. I. Simonova, and Yu. N. Tyurin**, *Sign-based methods in linear statistical models*, 1997
- 161 **Michael Blank**, *Discreteness and continuity in problems of chaotic dynamics*, 1997

For a complete list of titles in this series, visit the
AMS Bookstore at www.ams.org/bookstore/.

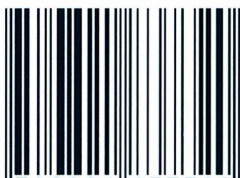
This volume is intended for the advanced study of topics in mathematical statistics. The first part of the book is devoted to sampling theory (from one-dimensional and multidimensional distributions), asymptotic properties of sampling, parameter estimation, sufficient statistics, and statistical estimates. The second part is devoted to hypothesis testing and includes the discussion of families of statistical hypotheses that can be asymptotically distinguished. In particular, the author describes goodness-of-fit and sequential statistical criteria (Kolmogorov, Pearson, Smirnov, and Wald) and studies their main properties.



For additional information
and updates on this book, visit

www.ams.org/bookpages/mmono-229

ISBN 0-8218-3732-X



9 780821 837320

MMONO/229

AMS *on the Web*
www.ams.org