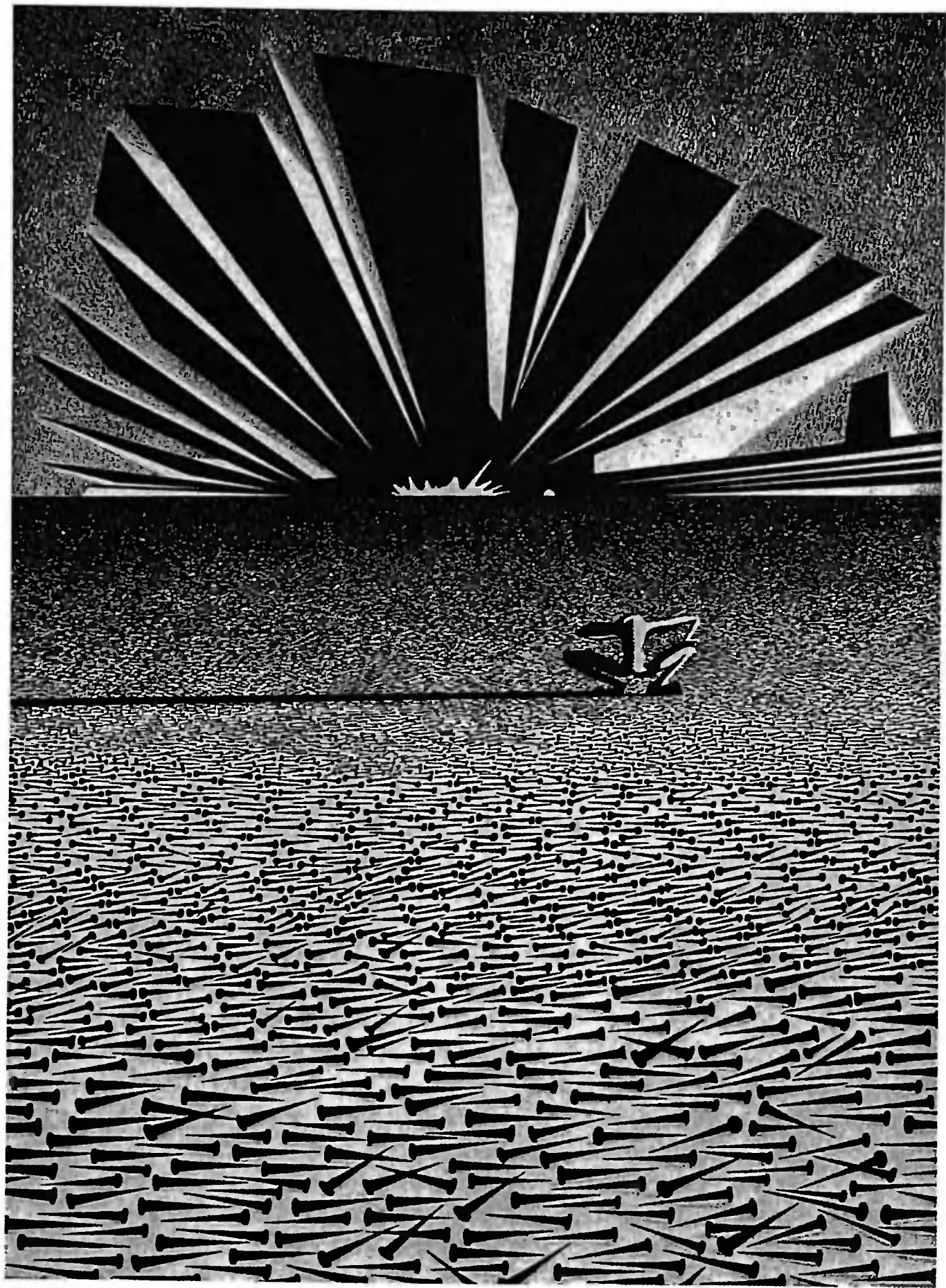


# Graduate Texts in Mathematics

**A.N. Shiryaev**

## **Probability**

**Second Edition**



**“Order out of chaos”**

*(Courtesy of Professor A. T. Fomenko of the Moscow State University)*



# Preface to the Second Edition

In the Preface to the first edition, originally published in 1980, we mentioned that this book was based on the author's lectures in the Department of Mechanics and Mathematics of the Lomonosov University in Moscow, which were issued, in part, in mimeographed form under the title "Probability, Statistics, and Stochastic Processes, I, II" and published by that University. Our original intention in writing the first edition of this book was to divide the contents into three parts: probability, mathematical statistics, and theory of stochastic processes, which corresponds to an outline of a three-semester course of lectures for university students of mathematics. However, in the course of preparing the book, it turned out to be impossible to realize this intention completely, since a full exposition would have required too much space. In this connection, we stated in the Preface to the first edition that only *probability theory* and the *theory of random processes with discrete time* were really adequately presented.

Essentially all of the first edition is reproduced in this second edition. Changes and corrections are, as a rule, editorial, taking into account comments made by both Russian and foreign readers of the Russian original and of the English and German translations [S11]. The author is grateful to all of these readers for their attention, advice, and helpful criticisms.

In this second English edition, new material also has been added, as follows: in Chapter III, §5, §§7–12; in Chapter IV, §5; in Chapter VII, §§8–10. The most important addition is the third chapter. There the reader will find expositions of a number of problems connected with a deeper study of themes such as the distance between probability measures, metrization of weak convergence, and contiguity of probability measures. In the same chapter, we have added proofs of a number of important results on the rapidity of convergence in the central limit theorem and in Poisson's theorem on the

approximation of the binomial by the Poisson distribution. These were merely stated in the first edition.

We also call attention to the new material on the probability of large deviations (Chapter IV, §5), on the central limit theorem for sums of dependent random variables (Chapter VII, §8), and on §§9 and 10 of Chapter VII.

During the last few years, the literature on probability published in Russia by Nauka has been extended by Sevastyanov [S10], 1982; Rozanov [R6], 1985; Borovkov [B4], 1986; and Gnedenko [G4], 1988. It appears that these publications, together with the present volume, being quite different and complementing each other, cover an extensive amount of material that is essentially broad enough to satisfy contemporary demands by students in various branches of mathematics and physics for instruction in topics in probability theory.

Gnedenko's textbook [G4] contains many well-chosen examples, including applications, together with pedagogical material and extensive surveys of the history of probability theory. Borovkov's textbook [B4] is perhaps the most like the present book in the style of exposition. Chapters 9 (Elements of Renewal Theory), 11 (Factorization of the Identity) and 17 (Functional Limit Theorems), which distinguish [B4] from this book and from [G4] and [R6], deserve special mention. Rozanov's textbook contains a great deal of material on a variety of mathematical models which the theory of probability and mathematical statistics provides for describing random phenomena and their evolution. The textbook by Sevastyanov is based on his two-semester course at the Moscow State University. The material in its last four chapters covers the minimum amount of probability and mathematical statistics required in a one-year university program. In our text, perhaps to a greater extent than in those mentioned above, a significant amount of space is given to set-theoretic aspects and mathematical foundations of probability theory.

Exercises and problems are given in the books by Gnedenko and Sevastyanov at the ends of chapters, and in the present textbook at the end of each section. These, together with, for example, the problem sets by A. V. Prokhorov and V. G. and N. G. Ushakov (*Problems in Probability Theory*, Nauka, Moscow, 1986) and by Zubkov, Sevastyanov, and Chistyakov (*Collected Problems in Probability Theory*, Nauka, Moscow, 1988) can be used by readers for independent study, and by teachers as a basis for seminars for students.

Special thanks to Harold Boas, who kindly translated the revisions from Russian to English for this new edition.

# Preface to the First Edition

This textbook is based on a three-semester course of lectures given by the author in recent years in the Mechanics–Mathematics Faculty of Moscow State University and issued, in part, in mimeographed form under the title *Probability, Statistics, Stochastic Processes, I, II* by the Moscow State University Press.

We follow tradition by devoting the first part of the course (roughly one semester) to the elementary theory of probability (Chapter I). This begins with the construction of probabilistic models with finitely many outcomes and introduces such fundamental probabilistic concepts as sample spaces, events, probability, independence, random variables, expectation, correlation, conditional probabilities, and so on.

Many probabilistic and statistical regularities are effectively illustrated even by the simplest random walk generated by Bernoulli trials. In this connection we study both classical results (law of large numbers, local and integral De Moivre and Laplace theorems) and more modern results (for example, the arc sine law).

The first chapter concludes with a discussion of dependent random variables generated by martingales and by Markov chains.

Chapters II–IV form an expanded version of the second part of the course (second semester). Here we present (Chapter II) Kolmogorov's generally accepted axiomatization of probability theory and the mathematical methods that constitute the tools of modern probability theory ( $\sigma$ -algebras, measures and their representations, the Lebesgue integral, random variables and random elements, characteristic functions, conditional expectation with respect to a  $\sigma$ -algebra, Gaussian systems, and so on). Note that two measure-theoretical results—Carathéodory's theorem on the extension of measures and the Radon–Nikodým theorem—are quoted without proof.

The third chapter is devoted to problems about weak convergence of probability distributions and the method of characteristic functions for proving limit theorems. We introduce the concepts of relative compactness and tightness of families of probability distributions, and prove (for the real line) Prohorov's theorem on the equivalence of these concepts.

The same part of the course discusses properties "with probability 1" for sequences and sums of independent random variables (Chapter IV). We give proofs of the "zero or one laws" of Kolmogorov and of Hewitt and Savage, tests for the convergence of series, and conditions for the strong law of large numbers. The law of the iterated logarithm is stated for arbitrary sequences of independent identically distributed random variables with finite second moments, and proved under the assumption that the variables have Gaussian distributions.

Finally, the third part of the book (Chapters V–VIII) is devoted to random processes with discrete parameters (random sequences). Chapters V and VI are devoted to the theory of stationary random sequences, where "stationary" is interpreted either in the strict or the wide sense. The theory of random sequences that are stationary in the strict sense is based on the ideas of ergodic theory: measure preserving transformations, ergodicity, mixing, etc. We reproduce a simple proof (by A. Garsia) of the maximal ergodic theorem; this also lets us give a simple proof of the Birkhoff–Khinchin ergodic theorem.

The discussion of sequences of random variables that are stationary in the wide sense begins with a proof of the spectral representation of the covariance function. Then we introduce orthogonal stochastic measures, and integrals with respect to these, and establish the spectral representation of the sequences themselves. We also discuss a number of statistical problems: estimating the covariance function and the spectral density, extrapolation, interpolation and filtering. The chapter includes material on the Kalman–Bucy filter and its generalizations.

The seventh chapter discusses the basic results of the theory of martingales and related ideas. This material has only rarely been included in traditional courses in probability theory. In the last chapter, which is devoted to Markov chains, the greatest attention is given to problems on the asymptotic behavior of Markov chains with countably many states.

Each section ends with problems of various kinds: some of them ask for proofs of statements made but not proved in the text, some consist of propositions that will be used later, some are intended to give additional information about the circle of ideas that is under discussion, and finally, some are simple exercises.

In designing the course and preparing this text, the author has used a variety of sources on probability theory. The Historical and Bibliographical Notes indicate both the historical sources of the results and supplementary references for the material under consideration.

The numbering system and form of references is the following. Each section has its own enumeration of theorems, lemmas and formulas (with

no indication of chapter or section). For a reference to a result from a different section of the same chapter, we use double numbering, with the first number indicating the number of the section (thus, (2.10) means formula (10) of §2). For references to a different chapter we use triple numbering (thus, formula (II.4.3) means formula (3) of §4 of Chapter II). Works listed in the References at the end of the book have the form  $[L_n]$ , where  $L$  is a letter and  $n$  is a numeral.

The author takes this opportunity to thank his teacher A. N. Kolmogorov, and B. V. Gnedenko and Yu. V. Prokhorov, from whom he learned probability theory and under whose direction he had the opportunity of using it. For discussions and advice, the author also thanks his colleagues in the Departments of Probability Theory and Mathematical Statistics at the Moscow State University, and his colleagues in the Section on probability theory of the Steklov Mathematical Institute of the Academy of Sciences of the U.S.S.R.

*Moscow*  
*Steklov Mathematical Institute*

A. N. SHIRYAEV

*Translator's acknowledgement.* I am grateful both to the author and to my colleague, C. T. Ionescu Tulcea, for advice about terminology.

R. P. B.





# Contents

Preface to the Second Edition	vii
Preface to the First Edition	ix
Introduction	1
CHAPTER I	
Elementary Probability Theory	5
§1. Probabilistic Model of an Experiment with a Finite Number of Outcomes	5
§2. Some Classical Models and Distributions	17
§3. Conditional Probability. Independence	23
§4. Random Variables and Their Properties	32
§5. The Bernoulli Scheme. I. The Law of Large Numbers	45
§6. The Bernoulli Scheme. II. Limit Theorems (Local, De Moivre–Laplace, Poisson)	55
§7. Estimating the Probability of Success in the Bernoulli Scheme	70
§8. Conditional Probabilities and Mathematical Expectations with Respect to Decompositions	76
§9. Random Walk. I. Probabilities of Ruin and Mean Duration in Coin Tossing	83
§10. Random Walk. II. Reflection Principle. Arcsine Law	94
§11. Martingales. Some Applications to the Random Walk	103
§12. Markov Chains. Ergodic Theorem. Strong Markov Property	110
CHAPTER II	
Mathematical Foundations of Probability Theory	131
§1. Probabilistic Model for an Experiment with Infinitely Many Outcomes. Kolmogorov's Axioms	131

§2. Algebras and $\sigma$ -algebras. Measurable Spaces	139
§3. Methods of Introducing Probability Measures on Measurable Spaces	151
§4. Random Variables. I.	170
§5. Random Elements	176
§6. Lebesgue Integral. Expectation	180
§7. Conditional Probabilities and Conditional Expectations with Respect to a $\sigma$ -Algebra	212
§8. Random Variables. II.	234
§9. Construction of a Process with Given Finite-Dimensional Distribution	245
§10. Various Kinds of Convergence of Sequences of Random Variables	252
§11. The Hilbert Space of Random Variables with Finite Second Moment	262
§12. Characteristic Functions	274
§13. Gaussian Systems	297

### CHAPTER III

#### Convergence of Probability Measures. Central Limit Theorem

§1. Weak Convergence of Probability Measures and Distributions	308
§2. Relative Compactness and Tightness of Families of Probability Distributions	317
§3. Proofs of Limit Theorems by the Method of Characteristic Functions	321
§4. Central Limit Theorem for Sums of Independent Random Variables. I. The Lindeberg Condition	328
§5. Central Limit Theorem for Sums of Independent Random Variables. II. Nonclassical Conditions	337
§6. Infinitely Divisible and Stable Distributions	341
§7. Metrizability of Weak Convergence	348
§8. On the Connection of Weak Convergence of Measures with Almost Sure Convergence of Random Elements ("Method of a Single Probability Space")	353
§9. The Distance in Variation between Probability Measures. Kakutani-Hellinger Distance and Hellinger Integrals. Application to Absolute Continuity and Singularity of Measures	359
§10. Contiguity and Entire Asymptotic Separation of Probability Measures	368
§11. Rapidity of Convergence in the Central Limit Theorem	373
§12. Rapidity of Convergence in Poisson's Theorem	376

### CHAPTER IV

Sequences and Sums of Independent Random Variables	379
§1. Zero-or-One Laws	379
§2. Convergence of Series	384
§3. Strong Law of Large Numbers	388
§4. Law of the Iterated Logarithm	395
§5. Rapidity of Convergence in the Strong Law of Large Numbers and in the Probabilities of Large Deviations	400

## CHAPTER V

## Stationary (Strict Sense) Random Sequences and Ergodic Theory

404

- §1. Stationary (Strict Sense) Random Sequences. Measure-Preserving Transformations
- §2. Ergodicity and Mixing
- §3. Ergodic Theorems

404

407

409

## CHAPTER VI

Stationary (Wide Sense) Random Sequences.  $L^2$  Theory

415

- §1. Spectral Representation of the Covariance Function
- §2. Orthogonal Stochastic Measures and Stochastic Integrals
- §3. Spectral Representation of Stationary (Wide Sense) Sequences
- §4. Statistical Estimation of the Covariance Function and the Spectral Density
- §5. Wold's Expansion
- §6. Extrapolation, Interpolation and Filtering
- §7. The Kalman-Bucy Filter and Its Generalizations

415

423

429

440

446

453

464

## CHAPTER VII

## Sequences of Random Variables that Form Martingales

474

- §1. Definitions of Martingales and Related Concepts
- §2. Preservation of the Martingale Property Under Time Change at a Random Time
- §3. Fundamental Inequalities
- §4. General Theorems on the Convergence of Submartingales and Martingales
- §5. Sets of Convergence of Submartingales and Martingales
- §6. Absolute Continuity and Singularity of Probability Distributions
- §7. Asymptotics of the Probability of the Outcome of a Random Walk with Curvilinear Boundary
- §8. Central Limit Theorem for Sums of Dependent Random Variables
- §9. Discrete Version of Itô's Formula
- §10. Applications to Calculations of the Probability of Ruin in Insurance

474

484

492

508

515

524

536

541

554

558

## CHAPTER VIII

## Sequences of Random Variables that Form Markov Chains

564

- §1. Definitions and Basic Properties
- §2. Classification of the States of a Markov Chain in Terms of Arithmetic Properties of the Transition Probabilities  $p_{ij}^{(n)}$
- §3. Classification of the States of a Markov Chain in Terms of Asymptotic Properties of the Probabilities  $p_{ii}^{(n)}$
- §4. On the Existence of Limits and of Stationary Distributions
- §5. Examples

564

569

573

582

587

Historical and Bibliographical Notes	596
References	603
Index of Symbols	609
Index	611

# Introduction

The subject matter of probability theory is the mathematical analysis of random events, i.e., of those empirical phenomena which—under certain circumstance—can be described by saying that:

They do not have *deterministic regularity* (observations of them do not yield the same outcome);

whereas at the same time

They possess some *statistical regularity* (indicated by the statistical stability of their frequency).

We illustrate with the classical example of a “fair” toss of an “unbiased” coin. It is clearly impossible to predict with certainty the outcome of each toss. The results of successive experiments are very irregular (now “head,” now “tail”) and we seem to have no possibility of discovering any regularity in such experiments. However, if we carry out a large number of “independent” experiments with an “unbiased” coin we can observe a very definite statistical regularity, namely that “head” appears with a frequency that is “close” to  $\frac{1}{2}$ .

Statistical stability of a frequency is very likely to suggest a hypothesis about a possible quantitative estimate of the “randomness” of some event  $A$  connected with the results of the experiments. With this starting point, probability theory postulates that corresponding to an event  $A$  there is a definite number  $P(A)$ , called the probability of the event, whose intrinsic property is that as the number of “independent” trials (experiments) increases the frequency of event  $A$  is approximated by  $P(A)$ .

Applied to our example, this means that it is natural to assign the proba-

bility  $\frac{1}{2}$  to the event  $A$  that consists of obtaining "head" in a toss of an "unbiased" coin.

There is no difficulty in multiplying examples in which it is very easy to obtain numerical values intuitively for the probabilities of one or another event. However, these examples are all of a similar nature and involve (so far) undefined concepts such as "fair" toss, "unbiased" coin, "independence," etc.

Having been invented to investigate the quantitative aspects of "randomness," probability theory, like every exact science, became such a science only at the point when the concept of a probabilistic model had been clearly formulated and axiomatized. In this connection it is natural for us to discuss, although only briefly, the fundamental steps in the development of probability theory.

Probability theory, as a science, originated in the middle of the seventeenth century with Pascal (1623–1662), Fermat (1601–1655) and Huygens (1629–1695). Although special calculations of probabilities in games of chance had been made earlier, in the fifteenth and sixteenth centuries, by Italian mathematicians (Cardano, Pacioli, Tartaglia, etc.), the first general methods for solving such problems were apparently given in the famous correspondence between Pascal and Fermat, begun in 1654, and in the first book on probability theory, *De Ratiociniis in Aleae Ludo* (*On Calculations in Games of Chance*), published by Huygens in 1657. It was at this time that the fundamental concept of "mathematical expectation" was developed and theorems on the addition and multiplication of probabilities were established.

The real history of probability theory begins with the work of James Bernoulli (1654–1705), *Ars Conjectandi* (*The Art of Guessing*) published in 1713, in which he proved (quite rigorously) the first limit theorem of probability theory, the law of large numbers; and of De Moivre (1667–1754), *Miscellanea Analytica Supplementum* (a rough translation might be *The Analytic Method* or *Analytic Miscellany*, 1730), in which the central limit theorem was stated and proved for the first time (for symmetric Bernoulli trials).

Bernoulli deserves the credit for introducing the "classical" definition of the concept of the *probability* of an event as the *ratio* of the number of possible outcomes of an experiment, that are favorable to the event, to the number of possible outcomes.

Bernoulli was probably the first to realize the importance of considering infinite sequences of random trials and to make a clear distinction between the probability of an event and the frequency of its realization.

De Moivre deserves the credit for defining such concepts as independence, mathematical expectation, and conditional probability.

In 1812 there appeared Laplace's (1749–1827) great treatise *Théorie Analytique des Probabilités* (*Analytic Theory of Probability*) in which he presented his own results in probability theory as well as those of his predecessors. In particular, he generalized De Moivre's theorem to the general

(unsymmetric) case of Bernoulli trials, and at the same time presented De Moivre's results in a more complete form.

Laplace's most important contribution was the application of probabilistic methods to errors of observation. He formulated the idea of considering errors of observation as the cumulative results of adding a large number of independent elementary errors. From this it followed that under rather general conditions the distribution of errors of observation must be at least approximately normal.

The work of Poisson (1781–1840) and Gauss (1777–1855) belongs to the same epoch in the development of probability theory, when the center of the stage was held by limit theorems. In contemporary probability theory we think of Poisson in connection with the distribution and the process that bear his name. Gauss is credited with originating the theory of errors and, in particular, with creating the fundamental method of least squares.

The next important period in the development of probability theory is connected with the names of P. L. Chebyshev (1821–1894), A. A. Markov (1856–1922), and A. M. Lyapunov (1857–1918), who developed effective methods for proving limit theorems for sums of independent but arbitrarily distributed random variables.

The number of Chebyshev's publications in probability theory is not large—four in all—but it would be hard to overestimate their role in probability theory and in the development of the classical Russian school of that subject.

“On the methodological side, the revolution brought about by Chebyshev was not only his insistence for the first time on complete rigor in the proofs of limit theorems, . . . but also, and principally, that Chebyshev always tried to obtain precise estimates for the deviations from the limiting regularities that are available for large but finite numbers of trials, in the form of inequalities that are valid unconditionally for any number of trials.”

(A. N. KOLMOGOROV [30])

Before Chebyshev the main interest in probability theory had been in the calculation of the probabilities of random events. He, however, was the first to realize clearly and exploit the full strength of the concepts of random variables and their mathematical expectations.

The leading exponent of Chebyshev's ideas was his devoted student Markov, to whom there belongs the indisputable credit of presenting his teacher's results with complete clarity. Among Markov's own significant contributions to probability theory were his pioneering investigations of limit theorems for sums of independent random variables and the creation of a new branch of probability theory, the theory of dependent random variables that form what we now call a Markov chain.

“Markov's classical course in the calculus of probability and his original papers, which are models of precision and clarity, contributed to the greatest extent to the transformation of probability theory into one of the most significant



branches of mathematics and to a wide extension of the ideas and methods of Chebyshev."

(S. N. BERNSTEIN [3])

To prove the central limit theorem of probability theory (the theorem on convergence to the normal distribution), Chebyshev and Markov used what is known as the method of moments. With more general hypotheses and a simpler method, the method of characteristic functions, the theorem was obtained by Lyapunov. The subsequent development of the theory has shown that the method of characteristic functions is a powerful analytic tool for establishing the most diverse limit theorems.

The modern period in the development of probability theory begins with its axiomatization. The first work in this direction was done by S. N. Bernstein (1880–1968), R. von Mises (1883–1953), and E. Borel (1871–1956). A. N. Kolmogorov's book *Foundations of the Theory of Probability* appeared in 1933. Here he presented the axiomatic theory that has become generally accepted and is not only applicable to all the classical branches of probability theory, but also provides a firm foundation for the development of new branches that have arisen from questions in the sciences and involve infinite-dimensional distributions.

The treatment in the present book is based on Kolmogorov's axiomatic approach. However, to prevent formalities and logical subtleties from obscuring the intuitive ideas, our exposition begins with the elementary theory of probability, whose elementariness is merely that in the corresponding probabilistic models we consider only experiments with finitely many outcomes. Thereafter we present the foundations of probability theory in their most general form.

The 1920s and '30s saw a rapid development of one of the new branches of probability theory, the theory of stochastic processes, which studies families of random variables that evolve with time. We have seen the creation of theories of Markov processes, stationary processes, martingales, and limit theorems for stochastic processes. Information theory is a recent addition.

The present book is principally concerned with stochastic processes with discrete parameters: random sequences. However, the material presented in the second chapter provides a solid foundation (particularly of a logical nature) for the study of the general theory of stochastic processes.

It was also in the 1920s and '30s that mathematical statistics became a separate mathematical discipline. In a certain sense mathematical statistics deals with inverses of the problems of probability: If the basic aim of probability theory is to calculate the probabilities of complicated events under a given probabilistic model, mathematical statistics sets itself the inverse problem: to clarify the structure of probabilistic-statistical models by means of observations of various complicated events.

Some of the problems and methods of mathematical statistics are also discussed in this book. However, all that is presented in detail here is probability theory and the theory of stochastic processes with discrete parameters.

## CHAPTER I

# Elementary Probability Theory

### §1. Probabilistic Model of an Experiment with a Finite Number of Outcomes

1. Let us consider an experiment of which all possible results are included in a finite number of outcomes  $\omega_1, \dots, \omega_N$ . We do not need to know the nature of these outcomes, only that there are a finite number  $N$  of them.

We call  $\omega_1, \dots, \omega_N$  *elementary events*, or *sample points*, and the finite set

$$\Omega = \{\omega_1, \dots, \omega_N\},$$

the *space of elementary events* or the *sample space*.

The choice of the space of elementary events is the *first step* in formulating a probabilistic model for an experiment. Let us consider some examples of sample spaces.

EXAMPLE 1. For a single toss of a coin the sample space  $\Omega$  consists of two points:

$$\Omega = \{H, T\},$$

where H = "head" and T = "tail". (We exclude possibilities like "the coin stands on edge," "the coin disappears," etc.)

EXAMPLE 2. For  $n$  tosses of a coin the sample space is

$$\Omega = \{\omega: \omega = (a_1, \dots, a_n), a_i = H \text{ or } T\}$$

and the general number  $N(\Omega)$  of outcomes is  $2^n$ .

EXAMPLE 3. First toss a coin. If it falls "head" then toss a die (with six faces numbered 1, 2, 3, 4, 5, 6); if it falls "tail", toss the coin again. The sample space for this experiment is

$$\Omega = \{H1, H2, H3, H4, H5, H6, TH, TT\}.$$

We now consider some more complicated examples involving the selection of  $n$  balls from an urn containing  $M$  distinguishable balls.

2. EXAMPLE 4 (Sampling with replacement). This is an experiment in which after each step the selected ball is returned again. In this case each sample of  $n$  balls can be presented in the form  $(a_1, \dots, a_n)$ , where  $a_i$  is the label of the ball selected at the  $i$ th step. It is clear that in sampling with replacement each  $a_i$  can have any of the  $M$  values 1, 2,  $\dots$ ,  $M$ . The description of the sample space depends in an essential way on whether we consider samples like, for example, (4, 1, 2, 1) and (1, 4, 2, 1) as different or the same. It is customary to distinguish two cases: *ordered* samples and *unordered* samples. In the first case samples containing the same elements, but arranged differently, are considered to be different. In the second case the order of the elements is disregarded and the two samples are considered to be the same. To emphasize which kind of sample we are considering, we use the notation  $(a_1, \dots, a_n)$  for ordered samples and  $[a_1, \dots, a_n]$  for unordered samples.

Thus for ordered samples the sample space has the form

$$\Omega = \{\omega: \omega = (a_1, \dots, a_n), a_i = 1, \dots, M\}$$

and the number of (different) outcomes is

$$N(\Omega) = M^n. \quad (1)$$

If, however, we consider unordered samples, then

$$\Omega = \{\omega: \omega = [a_1, \dots, a_n], a_i = 1, \dots, M\}.$$

Clearly the number  $N(\Omega)$  of (different) unordered samples is smaller than the number of ordered samples. Let us show that in the present case

$$N(\Omega) = C_{M+n-1}^n, \quad (2)$$

where  $C_k^l \equiv k!/[l!(k-l)!]$  is the number of combinations of  $l$  elements, taken  $k$  at a time.

We prove this by induction. Let  $N(M, n)$  be the number of outcomes of interest. It is clear that when  $k \leq M$  we have

$$N(k, 1) = k = C_k^1.$$

Now suppose that  $N(k, n) = C_{k+n-1}^k$  for  $k \leq M$ ; we show that this formula continues to hold when  $n$  is replaced by  $n + 1$ . For the unordered samples  $[a_1, \dots, a_{n+1}]$  that we are considering, we may suppose that the elements are arranged in nondecreasing order:  $a_1 \leq a_2 \leq \dots \leq a_n$ . It is clear that the number of unordered samples with  $a_1 = 1$  is  $N(M, n)$ , the number with  $a_1 = 2$  is  $N(M - 1, n)$ , etc. Consequently

$$\begin{aligned} N(M, n + 1) &= N(M, n) + N(M - 1, n) + \dots + N(1, n) \\ &= C_{M+n-1}^n + C_{M-1+n-1}^n + \dots + C_n^n \\ &= (C_{M+n}^{n+1} - C_{M+n-1}^{n+1}) + (C_{M-1+n}^{n+1} - C_{M-1+n-1}^{n+1}) \\ &\quad + \dots + (C_{n+1}^{n+1} - C_n^{n+1}) = C_{M+n}^{n+1}; \end{aligned}$$

here we have used the easily verified property

$$C_k^{l-1} + C_k^l = C_{k+1}^l$$

of the binomial coefficients.

**EXAMPLE 5 (Sampling without replacement).** Suppose that  $n \leq M$  and that the selected balls are not returned. In this case we again consider two possibilities, namely ordered and unordered samples.

For ordered samples without replacement the sample space is

$$\Omega = \{\omega: \omega = (a_1, \dots, a_n), a_k \neq a_l, k \neq l, a_i = 1, \dots, M\},$$

and the number of elements of this set (called *permutations*) is  $M(M - 1) \dots (M - n + 1)$ . We denote this by  $(M)_n$  or  $A_M^n$  and call it "the number of permutations of  $M$  things,  $n$  at a time".

For unordered samples (called *combinations*) the sample space

$$\Omega = \{\omega: \omega = [a_1, \dots, a_n], a_k \neq a_l, k \neq l, a_i = 1, \dots, M\}$$

consists of

$$N(\Omega) = C_M^n \quad (3)$$

elements. In fact, from each unordered sample  $[a_1, \dots, a_n]$  consisting of distinct elements we can obtain  $n!$  ordered samples. Consequently

$$N(\Omega) \cdot n! = (M)_n$$

and therefore

$$N(\Omega) = \frac{(M)_n}{n!} = C_M^n.$$

The results on the numbers of samples of  $n$  from an urn with  $M$  balls are presented in Table 1.

Table 1

$M^n$	$C_{M+n-1}^n$	With replacement
$(M)_n$	$C_M^n$	Without replacement
Ordered	Unordered	Sample Type

For the case  $M = 3$  and  $n = 2$ , the corresponding sample spaces are displayed in Table 2.

EXAMPLE 6 (Distribution of objects in cells). We consider the structure of the sample space in the problem of placing  $n$  objects (balls, etc.) in  $M$  cells (boxes, etc.). For example, such problems arise in statistical physics in studying the distribution of  $n$  particles (which might be protons, electrons, ...) among  $M$  states (which might be energy levels).

Let the cells be numbered  $1, 2, \dots, M$ , and suppose first that the objects are distinguishable (numbered  $1, 2, \dots, n$ ). Then a distribution of the  $n$  objects among the  $M$  cells is completely described by an ordered set  $(a_1, \dots, a_n)$ , where  $a_i$  is the index of the cell containing object  $i$ . However, if the objects are indistinguishable their distribution among the  $M$  cells is completely determined by the unordered set  $[a_1, \dots, a_n]$ , where  $a_i$  is the index of the cell into which an object is put at the  $i$ th step.

Comparing this situation with Examples 4 and 5, we have the following correspondences:

(ordered samples)  $\leftrightarrow$  (distinguishable objects),

(unordered samples)  $\leftrightarrow$  (indistinguishable objects),

Table 2

(1, 1) (1, 2) (1, 3) (2, 1) (2, 2) (2, 3) (3, 1) (3, 2) (3, 3)	[1, 1] [2, 2] [3, 3] [1, 2] [1, 3] [2, 3]	With replacement
(1, 2) (1, 3) (2, 1) (2, 3) (3, 1) (3, 2)	[1, 2] [1, 3] [2, 3]	Without replacement
Ordered	Unordered	Sample Type

by which we mean that to an instance of an ordered (unordered) sample of  $n$  balls from an urn containing  $M$  balls there corresponds (one and only one) instance of distributing  $n$  distinguishable (indistinguishable) objects among  $M$  cells.

In a similar sense we have the following correspondences:

(sampling with replacement)  $\leftrightarrow$  (a cell may receive any number of objects)

(sampling without replacement)  $\leftrightarrow$  (a cell may receive at most one object)

These correspondences generate others of the same kind:

(an unordered sample in sampling without replacement)  $\leftrightarrow$  (indistinguishable objects in the problem of distribution among cells when each cell may receive at most one object)

etc.; so that we can use Examples 4 and 5 to describe the sample space for the problem of distributing distinguishable or indistinguishable objects among cells either with exclusion (a cell may receive at most one object) or without exclusion (a cell may receive any number of objects).

Table 3 displays the distributions of two objects among three cells. For distinguishable objects, we denote them by W (white) and B (black). For indistinguishable objects, the presence of an object in a cell is indicated by a +.

Table 3

<div>W B <input type="checkbox"/> <input type="checkbox"/></div> <div>W B <input type="checkbox"/> <input type="checkbox"/></div> <div>W <input type="checkbox"/> B</div> <div>B W <input type="checkbox"/> <input type="checkbox"/></div> <div><input type="checkbox"/> W B <input type="checkbox"/> <input type="checkbox"/></div> <div><input type="checkbox"/> W B</div> <div>B <input type="checkbox"/> W <input type="checkbox"/> <input type="checkbox"/></div> <div>B W <input type="checkbox"/> <input type="checkbox"/></div> <div><input type="checkbox"/> <input type="checkbox"/> W B</div>	<div><input type="checkbox"/> + <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/></div> <div><input type="checkbox"/> <input type="checkbox"/> + <input type="checkbox"/> <input type="checkbox"/></div> <div><input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/> + <input type="checkbox"/></div> <div><input type="checkbox"/> + <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/></div> <div><input type="checkbox"/> + <input type="checkbox"/> <input type="checkbox"/></div> <div><input type="checkbox"/> <input type="checkbox"/> + <input type="checkbox"/></div> <div><input type="checkbox"/> <input type="checkbox"/> + <input type="checkbox"/></div>	Without exclusion
<div>W B <input type="checkbox"/> <input type="checkbox"/></div> <div>W <input type="checkbox"/> <input type="checkbox"/> B</div> <div>B W <input type="checkbox"/> <input type="checkbox"/></div> <div><input type="checkbox"/> <input type="checkbox"/> W B</div> <div>B <input type="checkbox"/> <input type="checkbox"/> W</div> <div><input type="checkbox"/> <input type="checkbox"/> B W</div>	<div><input type="checkbox"/> + <input type="checkbox"/> <input type="checkbox"/> <input type="checkbox"/></div> <div><input type="checkbox"/> + <input type="checkbox"/> <input type="checkbox"/></div> <div><input type="checkbox"/> <input type="checkbox"/> + <input type="checkbox"/></div>	With exclusion
Distinguishable objects	Indistinguishable objects	Distribution Kind of objects

Table 4

$N(\Omega)$ in the problem of placing $n$ objects in $M$ cells			
<div>Kind of objects Distribution</div>	Distinguishable objects	Indistinguishable objects	
Without exclusion	$M^n$ (Maxwell– Boltzmann statistics)	$C_{M+n-1}^n$ (Bose– Einstein statistics)	With replacement
With exclusion	$(M)_n$	$C_M^n$ (Fermi–Dirac statistics)	Without replacement
	Ordered samples	Unordered samples	<div>Sample Type</div>
$N(\Omega)$ in the problem of choosing $n$ balls from an urn containing $M$ balls			

The duality that we have observed between the two problems gives us an obvious way of finding the number of outcomes in the problem of placing objects in cells. The results, which include the results in Table 1, are given in Table 4.

In statistical physics one says that distinguishable (or indistinguishable, respectively) particles that are not subject to the Pauli exclusion principle† obey Maxwell–Boltzmann statistics (or, respectively, Bose–Einstein statistics). If, however, the particles are indistinguishable and are subject to the exclusion principle, they obey Fermi–Dirac statistics (see Table 4). For example, electrons, protons and neutrons obey Fermi–Dirac statistics. Photons and pions obey Bose–Einstein statistics. Distinguishable particles that are subject to the exclusion principle do not occur in physics.

3. In addition to the concept of sample space we now need the fundamental concept of *event*.

Experimenters are ordinarily interested, not in what particular outcome occurs as the result of a trial, but in whether the outcome belongs to some subset of the set of all possible outcomes. We shall describe as *events* all subsets  $A \subset \Omega$  for which, under the conditions of the experiment, it is possible to say either “the outcome  $\omega \in A$ ” or “the outcome  $\omega \notin A$ .”

† At most one particle in each cell. (Translator)

For example, let a coin be tossed three times. The sample space  $\Omega$  consists of the eight points

$$\Omega = \{HHH, HHT, \dots, TTT\}$$

and if we are able to observe (determine, measure, etc.) the results of all three tosses, we say that the set

$$A = \{HHH, HHT, HTH, THH\}$$

is the event consisting of the appearance of at least two heads. If, however, we can determine only the result of the first toss, this set  $A$  cannot be considered to be an event, since there is no way to give either a positive or negative answer to the question of whether a specific outcome  $\omega$  belongs to  $A$ .

Starting from a given collection of sets that are events, we can form new events by means of statements containing the logical connectives "or," "and," and "not," which correspond in the language of set theory to the operations "union," "intersection," and "complement."

If  $A$  and  $B$  are sets, their *union*, denoted by  $A \cup B$ , is the set of points that belong either to  $A$  or to  $B$ :

$$A \cup B = \{\omega \in \Omega: \omega \in A \text{ or } \omega \in B\}.$$

In the language of probability theory,  $A \cup B$  is the event consisting of the realization either of  $A$  or of  $B$ .

The *intersection* of  $A$  and  $B$ , denoted by  $A \cap B$ , or by  $AB$ , is the set of points that belong to both  $A$  and  $B$ :

$$A \cap B = \{\omega \in \Omega: \omega \in A \text{ and } \omega \in B\}.$$

The event  $A \cap B$  consists of the simultaneous realization of both  $A$  and  $B$ .

For example, if  $A = \{HH, HT, TH\}$  and  $B = \{TT, TH, HT\}$  then

$$A \cup B = \{HH, HT, TH, TT\} \quad (= \Omega),$$

$$A \cap B = \{TH, HT\}.$$

If  $A$  is a subset of  $\Omega$ , its *complement*, denoted by  $\bar{A}$ , is the set of points of  $\Omega$  that do not belong to  $A$ .

If  $B \setminus A$  denotes the *difference* of  $B$  and  $A$  (i.e. the set of points that belong to  $B$  but not to  $A$ ) then  $\bar{A} = \Omega \setminus A$ . In the language of probability,  $\bar{A}$  is the event consisting of the nonrealization of  $A$ . For example, if  $A = \{HH, HT, TH\}$  then  $\bar{A} = \{TT\}$ , the event in which two successive tails occur.

The sets  $A$  and  $\bar{A}$  have no points in common and consequently  $A \cap \bar{A}$  is empty. We denote the empty set by  $\emptyset$ . In probability theory,  $\emptyset$  is called an *impossible* event. The set  $\Omega$  is naturally called the *certain* event.

When  $A$  and  $B$  are disjoint ( $AB = \emptyset$ ), the union  $A \cup B$  is called the *sum* of  $A$  and  $B$  and written  $A + B$ .

If we consider a collection  $\mathcal{A}_0$  of sets  $A \subseteq \Omega$  we may use the set-theoretic operators  $\cup$ ,  $\cap$  and  $\setminus$  to form a new collection of sets from the elements of



$\mathcal{A}_0$ ; these sets are again events. If we adjoin the certain and impossible events  $\Omega$  and  $\emptyset$  we obtain a collection  $\mathcal{A}$  of sets which is an *algebra*, i.e. a collection of subsets of  $\Omega$  for which

- (1)  $\Omega \in \mathcal{A}$ ,
- (2) if  $A \in \mathcal{A}$ ,  $B \in \mathcal{A}$ , the sets  $A \cup B$ ,  $A \cap B$ ,  $A \setminus B$  also belong to  $\mathcal{A}$ .

It follows from what we have said that it will be advisable to consider collections of events that form algebras. In the future we shall consider only such collections.

Here are some examples of algebras of events:

- (a)  $\{\Omega, \emptyset\}$ , the collection consisting of  $\Omega$  and the empty set (we call this the *trivial algebra*);
- (b)  $\{A, \bar{A}, \Omega, \emptyset\}$ , the collection generated by  $A$ ;
- (c)  $\mathcal{A} = \{A: A \subseteq \Omega\}$ , the collection consisting of *all* the subsets of  $\Omega$  (including the empty set  $\emptyset$ ).

It is easy to check that all these algebras of events can be obtained from the following principle.

We say that a collection

$$\mathcal{D} = \{D_1, \dots, D_n\}$$

of sets is a *decomposition* of  $\Omega$ , and call the  $D_i$  the *atoms* of the decomposition, if the  $D_i$  are not empty, are pairwise disjoint, and their sum is  $\Omega$ :

$$D_1 + \dots + D_n = \Omega.$$

For example, if  $\Omega$  consists of three points,  $\Omega = \{1, 2, 3\}$ , there are five different decompositions:

$$\begin{array}{ll} \mathcal{D}_1 = \{D_1\} & \text{with } D_1 = \{1, 2, 3\}; \\ \mathcal{D}_2 = \{D_1, D_2\} & \text{with } D_1 = \{1, 2\}, D_2 = \{3\}; \\ \mathcal{D}_3 = \{D_1, D_2\} & \text{with } D_1 = \{1, 3\}, D_2 = \{2\}; \\ \mathcal{D}_4 = \{D_1, D_2\} & \text{with } D_1 = \{2, 3\}, D_2 = \{1\}; \\ \mathcal{D}_5 = \{D_1, D_2, D_3\} & \text{with } D_1 = \{1\}, D_2 = \{2\}, D_3 = \{3\}. \end{array}$$

(For the general number of decompositions of a finite set see Problem 2.)

If we consider all unions of the sets in  $\mathcal{D}$ , the resulting collection of sets, together with the empty set, forms an algebra, called the *algebra induced by*  $\mathcal{D}$ , and denoted by  $\alpha(\mathcal{D})$ . Thus the elements of  $\alpha(\mathcal{D})$  consist of the empty set together with the sums of sets which are atoms of  $\mathcal{D}$ .

Thus if  $\mathcal{D}$  is a decomposition, there is associated with it a specific algebra  $\mathcal{B} = \alpha(\mathcal{D})$ .

The converse is also true. Let  $\mathcal{B}$  be an algebra of subsets of a finite space  $\Omega$ . Then there is a unique decomposition  $\mathcal{D}$  whose atoms are the elements of

$\mathcal{B}$ , with  $\mathcal{B} = \alpha(\mathcal{D})$ . In fact, let  $D \in \mathcal{B}$  and let  $D$  have the property that for every  $B \in \mathcal{B}$  the set  $D \cap B$  either coincides with  $D$  or is empty. Then this collection of sets  $D$  forms a decomposition  $\mathcal{D}$  with the required property  $\alpha(\mathcal{D}) = \mathcal{B}$ . In Example (a),  $\mathcal{D}$  is the trivial decomposition consisting of the single set  $D_1 = \Omega$ ; in (b),  $\mathcal{D} = \{A, \bar{A}\}$ . The most fine-grained decomposition  $\mathcal{D}$ , which consists of the singletons  $\{\omega_i\}$ ,  $\omega_i \in \Omega$ , induces the algebra in Example (c), i.e. the algebra of all subsets of  $\Omega$ .

Let  $\mathcal{D}_1$  and  $\mathcal{D}_2$  be two decompositions. We say that  $\mathcal{D}_2$  is finer than  $\mathcal{D}_1$ , and write  $\mathcal{D}_1 \leq \mathcal{D}_2$ , if  $\alpha(\mathcal{D}_1) \subseteq \alpha(\mathcal{D}_2)$ .

Let us show that if  $\Omega$  consists, as we assumed above, of a finite number of points  $\omega_1, \dots, \omega_N$ , then the number  $N(\mathcal{A})$  of sets in the collection  $\mathcal{A}$  is equal to  $2^N$ . In fact, every nonempty set  $A \in \mathcal{A}$  can be represented as  $A = \{\omega_{i_1}, \dots, \omega_{i_k}\}$ , where  $\omega_{i_j} \in \Omega$ ,  $1 \leq k \leq N$ . With this set we associate the sequence of zeros and ones

$$(0, \dots, 0, 1, 0, \dots, 0, 1, \dots),$$

where there are ones in the positions  $i_1, \dots, i_k$  and zeros elsewhere. Then for a given  $k$  the number of different sets  $A$  of the form  $\{\omega_{i_1}, \dots, \omega_{i_k}\}$  is the same as the number of ways in which  $k$  ones ( $k$  indistinguishable objects) can be placed in  $N$  positions ( $N$  cells). According to Table 4 (see the lower right-hand square) we see that this number is  $C_N^k$ . Hence (counting the empty set) we find that

$$N(\mathcal{A}) = 1 + C_N^1 + \dots + C_N^N = (1 + 1)^N = 2^N.$$

4. We have now taken the first two steps in defining a probabilistic model of an experiment with a finite number of outcomes: we have selected a sample space and a collection  $\mathcal{A}$  of subsets, which form an algebra and are called events. We now take the next step, to assign to each sample point (outcome)  $\omega_i \in \Omega$ ,  $i = 1, \dots, N$ , a *weight*. This is denoted by  $p(\omega_i)$  and called the *probability* of the outcome  $\omega_i$ ; we assume that it has the following properties:

- (a)  $0 \leq p(\omega_i) \leq 1$  (nonnegativity),
- (b)  $p(\omega_1) + \dots + p(\omega_N) = 1$  (normalization).

Starting from the given probabilities  $p(\omega_i)$  of the outcomes  $\omega_i$ , we define the probability  $P(A)$  of any event  $A \in \mathcal{A}$  by

$$P(A) = \sum_{\{i: \omega_i \in A\}} p(\omega_i). \quad (4)$$

Finally, we say that a triple

$$(\Omega, \mathcal{A}, P),$$

where  $\Omega = \{\omega_1, \dots, \omega_N\}$ ,  $\mathcal{A}$  is an algebra of subsets of  $\Omega$  and

$$P = \{P(A); A \in \mathcal{A}\}$$

defines (or assigns) a *probabilistic model*, or a *probability space*, of experiments with a (finite) space  $\Omega$  of outcomes and algebra  $\mathcal{A}$  of events.

The following properties of probability follow from (4):

$$P(\emptyset) = 0, \quad (5)$$

$$P(\Omega) = 1, \quad (6)$$

$$P(A \cup B) = P(A) + P(B) - P(A \cap B). \quad (7)$$

In particular, if  $A \cap B = \emptyset$ , then

$$P(A + B) = P(A) + P(B) \quad (8)$$

and

$$P(\bar{A}) = 1 - P(A). \quad (9)$$

5. In constructing a probabilistic model for a specific situation, the construction of the sample space  $\Omega$  and the algebra  $\mathcal{A}$  of events are ordinarily not difficult. In elementary probability theory one usually takes the algebra  $\mathcal{A}$  to be the algebra of *all* subsets of  $\Omega$ . Any difficulty that may arise is in assigning probabilities to the sample points. In principle, the solution to this problem lies outside the domain of probability theory, and we shall not consider it in detail. We consider that our fundamental problem is not the question of how to assign probabilities, but how to calculate the probabilities of complicated events (elements of  $\mathcal{A}$ ) from the probabilities of the sample points.

It is clear from a mathematical point of view that for finite sample spaces we can obtain all conceivable (finite) probability spaces by assigning non-negative numbers  $p_1, \dots, p_N$ , satisfying the condition  $p_1 + \dots + p_N = 1$ , to the outcomes  $\omega_1, \dots, \omega_N$ .

The validity of the assignments of the numbers  $p_1, \dots, p_N$  can, in specific cases, be checked to a certain extent by using the law of large numbers (which will be discussed later on). It states that in a long series of "independent" experiments, carried out under identical conditions, the frequencies with which the elementary events appear are "close" to their probabilities.

In connection with the difficulty of assigning probabilities to outcomes, we note that there are many actual situations in which for reasons of symmetry it seems reasonable to consider all conceivable outcomes as equally probable. In such cases, if the sample space consists of points  $\omega_1, \dots, \omega_N$ , with  $N < \infty$ , we put

$$p(\omega_1) = \dots = p(\omega_N) = 1/N,$$

and consequently

$$P(A) = N(A)/N \quad (10)$$

for every event  $A \in \mathcal{A}$ , where  $N(A)$  is the number of sample points in  $A$ . This is called the classical method of assigning probabilities. It is clear that in this case the calculation of  $P(A)$  reduces to calculating the number of outcomes belonging to  $A$ . This is usually done by combinatorial methods, so that combinatorics, applied to finite sets, plays a significant role in the calculus of probabilities.

**EXAMPLE 7 (Coincidence problem).** Let an urn contain  $M$  balls numbered  $1, 2, \dots, M$ . We draw an ordered sample of size  $n$  with replacement. It is clear that then

$$\Omega = \{\omega: \omega = (a_1, \dots, a_n), a_i = 1, \dots, M\}$$

and  $N(\Omega) = M^n$ . Using the classical assignment of probabilities, we consider the  $M^n$  outcomes equally probable and ask for the probability of the event

$$A = \{\omega: \omega = (a_1, \dots, a_n), a_i \neq a_j, i \neq j\},$$

i.e., the event in which there is no repetition. Clearly  $N(A) = M(M-1) \dots (M-n+1)$ , and therefore

$$P(A) = \frac{(M)_n}{M^n} = \left(1 - \frac{1}{M}\right) \left(1 - \frac{2}{M}\right) \dots \left(1 - \frac{n-1}{M}\right). \quad (11)$$

This problem has the following striking interpretation. Suppose that there are  $n$  students in a class. Let us suppose that each student's birthday is on one of 365 days and that all days are equally probable. The question is, what is the probability  $P_n$  that there are at least two students in the class whose birthdays coincide? If we interpret selection of birthdays as selection of balls from an urn containing 365 balls, then by (11)

$$P_n = 1 - \frac{(365)_n}{365^n}.$$

The following table lists the values of  $P_n$  for some values of  $n$ :

$n$	4	16	22	23	40	64
$P_n$	0.016	0.284	0.476	0.507	0.891	0.997

It is interesting to note that (unexpectedly!) the size of class in which there is probability  $\frac{1}{2}$  of finding at least two students with the same birthday is not very large: only 23.

**EXAMPLE 8 (Prizes in a lottery).** Consider a lottery that is run in the following way. There are  $M$  tickets numbered  $1, 2, \dots, M$ , of which  $n$ , numbered  $1, \dots, n$ , win prizes ( $M \geq 2n$ ). You buy  $n$  tickets, and ask for the probability ( $P$ , say) of winning at least one prize.

Since the order in which the tickets are drawn plays no role in the presence or absence of winners in your purchase, we may suppose that the sample space has the form

$$\Omega = \{\omega: \omega = [a_1, \dots, a_n], a_k \neq a_l, k \neq l, a_i = 1, \dots, M\}.$$

By Table 1,  $N(\Omega) = C_M^n$ . Now let

$$A_0 = \{\omega: \omega = [a_1, \dots, a_n], a_k \neq a_l, k \neq l, a_i = n+1, \dots, M\}$$

be the event that there is no winner in the set of tickets you bought. Again by Table 1,  $N(A_0) = C_{M-n}^n$ . Therefore

$$\begin{aligned} P(A_0) &= \frac{C_{M-n}^n}{C_M^n} = \frac{(M-n)_n}{(M)_n} \\ &= \left(1 - \frac{n}{M}\right) \left(1 - \frac{n}{M-1}\right) \cdots \left(1 - \frac{n}{M-n+1}\right) \end{aligned}$$

and consequently

$$P = 1 - P(A_0) = 1 - \left(1 - \frac{n}{M}\right) \left(1 - \frac{n}{M-1}\right) \cdots \left(1 - \frac{n}{M-n+1}\right).$$

If  $M = n^2$  and  $n \rightarrow \infty$ , then  $P(A_0) \rightarrow e^{-1}$  and

$$P \rightarrow 1 - e^{-1} \approx 0.632.$$

The convergence is quite fast: for  $n = 10$  the probability is already  $P = 0.670$ .

## 6. PROBLEMS

1. Establish the following properties of the operators  $\cap$  and  $\cup$ :

$$A \cup B = B \cup A, \quad AB = BA \quad (\text{commutativity}),$$

$$A \cup (B \cap C) = (A \cup B) \cap C, \quad A(BC) = (AB)C \quad (\text{associativity}),$$

$$A(B \cup C) = AB \cup AC, \quad A \cap (BC) = (A \cap B)(A \cap C) \quad (\text{distributivity}),$$

$$A \cup A = A, \quad AA = A \quad (\text{idempotency}).$$

Show also that

$$\overline{A \cup B} = \bar{A} \cap \bar{B}, \quad \overline{AB} = \bar{A} \cup \bar{B}.$$

2. Let  $\Omega$  contain  $N$  elements. Show that the number  $d(N)$  of different decompositions of  $\Omega$  is given by the formula

$$d(N) = e^{-1} \sum_{k=0}^{\infty} \frac{k^N}{k!}. \quad (12)$$

(Hint: Show that

$$d(N) = \sum_{k=0}^{N-1} C_{N-1}^k d(k), \quad \text{where } d(0) = 1,$$

and then verify that the series in (12) satisfies the same recurrence relation.)

3. For any finite collection of sets  $A_1, \dots, A_n$ ,

$$P(A_1 \cup \dots \cup A_n) \leq P(A_1) + \dots + P(A_n).$$

4. Let  $A$  and  $B$  be events. Show that  $A\bar{B} \cup B\bar{A}$  is the event in which exactly one of  $A$  and  $B$  occurs. Moreover,

$$P(A\bar{B} \cup B\bar{A}) = P(A) + P(B) - 2P(AB).$$

5. Let  $A_1, \dots, A_n$  be events, and define  $S_0, S_1, \dots, S_n$  as follows:  $S_0 = 1$ ,

$$S_r = \sum_{J_r} P(A_{k_1} \cap \dots \cap A_{k_r}), \quad 1 \leq r \leq n,$$

where the sum is over the unordered subsets  $J_r = [k_1, \dots, k_r]$  of  $\{1, \dots, n\}$ .

Let  $B_m$  be the event in which each of the events  $A_1, \dots, A_n$  occurs exactly  $m$  times. Show that

$$P(B_m) = \sum_{r=m}^n (-1)^{r-m} C_r^m S_r.$$

In particular, for  $m = 0$

$$P(B_0) = 1 - S_1 + S_2 - \dots \pm S_n.$$

Show also that the probability that at least  $m$  of the events  $A_1, \dots, A_n$  occur simultaneously is

$$P(B_1) + \dots + P(B_n) = \sum_{r=m}^n (-1)^{r-m} C_r^{m-1} S_r.$$

In particular, the probability that at least one of the events  $A_1, \dots, A_n$  occurs is

$$P(B_1) + \dots + P(B_n) = S_1 - S_2 + \dots \pm S_n.$$

## §2. Some Classical Models and Distributions

**1. Binomial distribution.** Let a coin be tossed  $n$  times and record the results as an ordered set  $(a_1, \dots, a_n)$ , where  $a_i = 1$  for a head ("success") and  $a_i = 0$  for a tail ("failure"). The sample space is

$$\Omega = \{\omega: \omega = (a_1, \dots, a_n), a_i = 0, 1\}.$$

To each sample point  $\omega = (a_1, \dots, a_n)$  we assign the probability

$$p(\omega) = p^{\sum a_i} q^{n - \sum a_i},$$

where the nonnegative numbers  $p$  and  $q$  satisfy  $p + q = 1$ . In the first place, we verify that this assignment of the weights  $p(\omega)$  is consistent. It is enough to show that  $\sum_{\omega \in \Omega} p(\omega) = 1$ .

We consider all outcomes  $\omega = (a_1, \dots, a_n)$  for which  $\sum_i a_i = k$ , where  $k = 0, 1, \dots, n$ . According to Table 4 (distribution of  $k$  indistinguishable

ones in  $n$  places) the number of these outcomes is  $C_n^k$ . Therefore

$$\sum_{\omega \in \Omega} p(\omega) = \sum_{k=0}^n C_n^k p^k q^{n-k} = (p + q)^n = 1.$$

Thus the space  $\Omega$  together with the collection  $\mathcal{A}$  of all its subsets and the probabilities  $P(A) = \sum_{\omega \in A} p(\omega)$ ,  $A \in \mathcal{A}$ , defines a probabilistic model. It is natural to call this the probabilistic model for  $n$  tosses of a coin.

In the case  $n = 1$ , when the sample space contains just the two points  $\omega = 1$  ("success") and  $\omega = 0$  ("failure"), it is natural to call  $p(1) = p$  the probability of success. We shall see later that this model for  $n$  tosses of a coin can be thought of as the result of  $n$  "independent" experiments with probability  $p$  of success at each trial.

Let us consider the events

$$A_k = \{\omega: \omega = (a_1, \dots, a_n), a_1 + \dots + a_n = k\}, \quad k = 0, 1, \dots, n,$$

consisting of exactly  $k$  successes. It follows from what we said above that

$$P(A_k) = C_n^k p^k q^{n-k}, \quad (1)$$

and  $\sum_{k=0}^n P(A_k) = 1$ .

The set of probabilities  $(P(A_0), \dots, P(A_n))$  is called the *binomial distribution* (the number of successes in a sample of size  $n$ ). This distribution plays an extremely important role in probability theory since it arises in the most diverse probabilistic models. We write  $P_n(k) = P(A_k)$ ,  $k = 0, 1, \dots, n$ . Figure 1 shows the binomial distribution in the case  $p = \frac{1}{2}$  (symmetric coin) for  $n = 5, 10, 20$ .

We now present a different model (in essence, equivalent to the preceding one) which describes the random walk of a "particle."

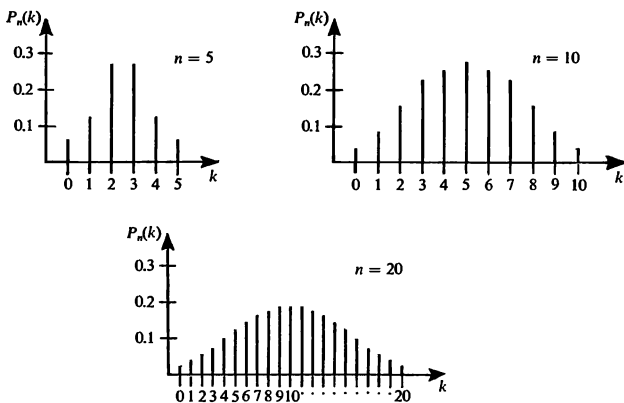
Let the particle start at the origin, and after unit time let it take a unit step upward or downward (Figure 2).

Consequently after  $n$  steps the particle can have moved at most  $n$  units up or  $n$  units down. It is clear that each path  $\omega$  of the particle is completely specified by a set  $(a_1, \dots, a_n)$ , where  $a_i = +1$  if the particle moves up at the  $i$ th step, and  $a_i = -1$  if it moves down. Let us assign to each path  $\omega$  the weight  $p(\omega) = p^{v(\omega)} q^{n-v(\omega)}$ , where  $v(\omega)$  is the number of  $+1$ 's in the sequence  $\omega = (a_1, \dots, a_n)$ , i.e.  $v(\omega) = [(a_1 + \dots + a_n) + n]/2$ , and the nonnegative numbers  $p$  and  $q$  satisfy  $p + q = 1$ .

Since  $\sum_{\omega \in \Omega} p(\omega) = 1$ , the set of probabilities  $p(\omega)$  together with the space  $\Omega$  of paths  $\omega = (a_1, \dots, a_n)$  and its subsets define an acceptable probabilistic model of the motion of the particle for  $n$  steps.

Let us ask the following question: What is the probability of the event  $A_k$  that after  $n$  steps the particle is at a point with ordinate  $k$ ? This condition is satisfied by those paths  $\omega$  for which  $v(\omega) - (n - v(\omega)) = k$ , i.e.

$$v(\omega) = \frac{n + k}{2}.$$

Figure 1. Graph of the binomial probabilities  $P_n(k)$  for  $n = 5, 10, 20$ .

The number of such paths (see Table 4) is  $C_n^{[n+k]/2}$ , and therefore

$$P(A_k) = C_n^{[n+k]/2} p^{[n+k]/2} q^{[n-k]/2}.$$

Consequently the binomial distribution  $(P(A_{-n}), \dots, P(A_0), \dots, P(A_n))$  can be said to describe the probability distribution for the position of the particle after  $n$  steps.

Note that in the symmetric case ( $p = q = \frac{1}{2}$ ) when the probabilities of the individual paths are equal to  $2^{-n}$ ,

$$P(A_k) = C_n^{[n+k]/2} \cdot 2^{-n}.$$

Let us investigate the asymptotic behavior of these probabilities for large  $n$ .

If the number of steps is  $2n$ , it follows from the properties of the binomial coefficients that the largest of the probabilities  $P(A_k)$ ,  $|k| \leq 2n$ , is

$$P(A_0) = C_{2n}^n \cdot 2^{-2n}.$$

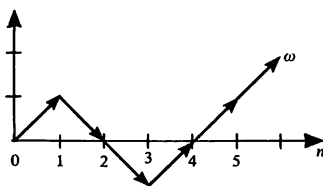


Figure 2



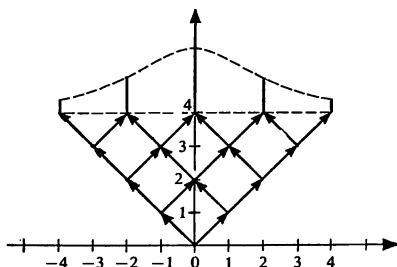


Figure 3. Beginning of the binomial distribution.

From Stirling's formula (see formula (6) in Section 4)

$$n! \sim \sqrt{2\pi n} e^{-n} n^n. \dagger$$

Consequently

$$C_{2n}^n = \frac{(2n)!}{(n!)^2} \sim 2^{2n} \cdot \frac{1}{\sqrt{\pi n}}$$

and therefore for large  $n$

$$P(A_0) \sim \frac{1}{\sqrt{\pi n}}.$$

Figure 3 represents the beginning of the binomial distribution for  $2n$  steps of a random walk (in contrast to Figure 2, the time axis is now directed upward).

**2. Multinomial distribution.** Generalizing the preceding model, we now suppose that the sample space is

$$\Omega = \{\omega: \omega = (a_1, \dots, a_n), a_i = b_1, \dots, b_r\},$$

where  $b_1, \dots, b_r$  are given numbers. Let  $v_i(\omega)$  be the number of elements of  $\omega = (a_1, \dots, a_n)$  that are equal to  $b_i$ ,  $i = 1, \dots, r$ , and define the probability of  $\omega$  by

$$p(\omega) = p_1^{v_1(\omega)} \cdots p_r^{v_r(\omega)},$$

where  $p_i \geq 0$  and  $p_1 + \cdots + p_r = 1$ . Note that

$$\sum_{\omega \in \Omega} p(\omega) = \sum_{\substack{n_1 \geq 0, \dots, n_r \geq 0, \\ n_1 + \dots + n_r = n}} C_n(n_1, \dots, n_r) p_1^{n_1} \cdots p_r^{n_r},$$

where  $C_n(n_1, \dots, n_r)$  is the number of (ordered) sequences  $(a_1, \dots, a_n)$  in which  $b_1$  occurs  $n_1$  times,  $\dots$ ,  $b_r$  occurs  $n_r$  times. Since  $n_1$  elements  $b_1$  can

$\dagger$  The notation  $f(n) \sim g(n)$  means that  $f(n)/g(n) \rightarrow 1$  as  $n \rightarrow \infty$ .

be distributed into  $n$  positions in  $C_n^{n_1}$  ways;  $n_2$  elements  $b_2$  into  $n - n_1$  positions in  $C_{n-n_1}^{n_2}$  ways, etc., we have

$$\begin{aligned} C_n(n_1, \dots, n_r) &= C_n^{n_1} \cdot C_{n-n_1}^{n_2} \cdots C_{n-(n_1+\dots+n_{r-1})}^{n_r} \\ &= \frac{n!}{n_1! (n-n_1)!} \cdot \frac{(n-n_1)!}{n_2! (n-n_1-n_2)!} \cdots 1 \\ &= \frac{n!}{n_1! \cdots n_r!}. \end{aligned}$$

Therefore

$$\sum_{\omega \in \Omega} p(\omega) = \sum_{\substack{n_1 \geq 0, \dots, n_r \geq 0, \\ n_1 + \dots + n_r = n}} \frac{n!}{n_1! \cdots n_r!} p_1^{n_1} \cdots p_r^{n_r} = (p_1 + \cdots + p_r)^n = 1,$$

and consequently we have defined an acceptable method of assigning probabilities.

Let

$$A_{n_1, \dots, n_r} = \{\omega: v_1(\omega) = n_1, \dots, v_r(\omega) = n_r\}.$$

Then

$$P(A_{n_1, \dots, n_r}) = C_n(n_1, \dots, n_r) p_1^{n_1} \cdots p_r^{n_r}. \quad (2)$$

The set of probabilities

$$\{P(A_{n_1, \dots, n_r})\}$$

is called the *multinomial* (or polynomial) distribution.

We emphasize that both this distribution and its special case, the binomial distribution, originate from problems about sampling *with replacement*.

**3. The multidimensional hypergeometric distribution** occurs in problems that involve sampling *without replacement*.

Consider, for example, an urn containing  $M$  balls numbered  $1, 2, \dots, M$ , where  $M_1$  balls have the color  $b_1, \dots, M_r$  balls have the color  $b_r$ , and  $M_1 + \dots + M_r = M$ . Suppose that we draw a sample of size  $n < M$  without replacement. The sample space is

$$\Omega = \{\omega: \omega = (a_1, \dots, a_n), a_k \neq a_l, k \neq l, a_i = 1, \dots, M\}$$

and  $N(\Omega) = (M)_n$ . Let us suppose that the sample points are equiprobable, and find the probability of the event  $B_{n_1, \dots, n_r}$  in which  $n_1$  balls have color  $b_1, \dots, n_r$  balls have color  $b_r$ , where  $n_1 + \dots + n_r = n$ . It is easy to show that

$$N(B_{n_1, \dots, n_r}) = C_n(n_1, \dots, n_r) (M_1)_{n_1} \cdots (M_r)_{n_r},$$

and therefore

$$P(B_{n_1, \dots, n_r}) = \frac{N(B_{n_1, \dots, n_r})}{N(\Omega)} = \frac{C_{M_1}^{n_1} \cdots C_{M_r}^{n_r}}{C_M^n}. \quad (3)$$

The set of probabilities  $\{P(B_{n_1, \dots, n_r})\}$  is called the *multidimensional hypergeometric distribution*. When  $r = 2$  it is simply called the *hypergeometric distribution* because its “generating function” is a hypergeometric function.

The structure of the multidimensional hypergeometric distribution is rather complicated. For example, the probability

$$P(B_{n_1, n_2}) = \frac{C_{M_1}^{n_1} C_{M_2}^{n_2}}{C_M^n}, \quad n_1 + n_2 = n, \quad M_1 + M_2 = M, \quad (4)$$

contains nine factorials. However, it is easily established that if  $M \rightarrow \infty$  and  $M_1 \rightarrow \infty$  in such a way that  $M_1/M \rightarrow p$  (and therefore  $M_2/M \rightarrow 1 - p$ ) then

$$P(B_{n_1, n_2}) \rightarrow C_{n_1+n_2}^{n_1} p^{n_1} (1-p)^{n_2}. \quad (5)$$

In other words, under the present hypotheses the hypergeometric distribution is approximated by the binomial; this is intuitively clear since when  $M$  and  $M_1$  are large (but finite), sampling without replacement ought to give almost the same result as sampling with replacement.

**EXAMPLE.** Let us use (4) to find the probability of picking six “lucky” numbers in a lottery of the following kind (this is an abstract formulation of the “sportloto,” which is well known in Russia):

There are 49 balls numbered from 1 to 49; six of them are lucky (colored red, say, whereas the rest are white). We draw a sample of six balls, without replacement. The question is, What is the probability that all six of these balls are lucky? Taking  $M = 49$ ,  $M_1 = 6$ ,  $n_1 = 6$ ,  $n_2 = 0$ , we see that the event of interest, namely

$$B_{6,0} = \{6 \text{ balls, all lucky}\}$$

has, by (4), probability

$$P(B_{6,0}) = \frac{1}{C_{49}^6} \approx 7.2 \times 10^{-8}.$$

4. The numbers  $n!$  increase extremely rapidly with  $n$ . For example,

$$10! = 3,628,800,$$

$$15! = 1,307,674,368,000,$$

and  $100!$  has 158 digits. Hence from either the theoretical or the computational point of view, it is important to know Stirling's formula,

$$n! = \sqrt{2\pi n} \left(\frac{n}{e}\right)^n \exp\left(\frac{\theta_n}{12n}\right), \quad 0 < \theta_n < 1, \quad (6)$$

whose proof can be found in most textbooks on mathematical analysis (see also [69]).

## 5. PROBLEMS

1. Prove formula (5).
2. Show that for the multinomial distribution  $\{P(A_{n_1}, \dots, A_{n_r})\}$  the maximum probability is attained at a point  $(k_1, \dots, k_r)$  that satisfies the inequalities  $np_i - 1 < k_i \leq (n + r - 1)p_i$ ,  $i = 1, \dots, r$ .
3. *One-dimensional Ising model.* Consider  $n$  particles located at the points  $1, 2, \dots, n$ . Suppose that each particle is of one of two types, and that there are  $n_1$  particles of the first type and  $n_2$  of the second ( $n_1 + n_2 = n$ ). We suppose that all  $n!$  arrangements of the particles are equally probable.

Construct a corresponding probabilistic model and find the probability of the event  $A_n(m_{11}, m_{12}, m_{21}, m_{22}) = \{v_{11} = m_{11}, \dots, v_{22} = m_{22}\}$ , where  $v_{ij}$  is the number of particles of type  $i$  following particles of type  $j$  ( $i, j = 1, 2$ ).

4. Prove the following inequalities by probabilistic reasoning:

$$\sum_{k=0}^n C_n^k = 2^n,$$

$$\sum_{k=0}^n (C_n^k)^2 = C_{2n}^n,$$

$$\sum_{k=0}^n (-1)^{n-k} C_m^k = C_{m-1}^n, \quad m \geq n + 1,$$

$$\sum_{k=0}^n k(k-1)C_m^k = m(m-1)2^{m-2}, \quad m \geq 2.$$

## §3. Conditional Probability. Independence

1. The concept of probabilities of events lets us answer questions of the following kind: If there are  $M$  balls in an urn,  $M_1$  white and  $M_2$  black, what is the probability  $P(A)$  of the event  $A$  that a selected ball is white? With the classical approach,  $P(A) = M_1/M$ .

The concept of *conditional probability*, which will be introduced below, lets us answer questions of the following kind: What is the probability that the second ball is white (event  $B$ ) under the condition that the first ball was also white (event  $A$ )? (We are thinking of sampling without replacement.)

It is natural to reason as follows: if the first ball is white, then at the second step we have an urn containing  $M - 1$  balls, of which  $M_1 - 1$  are white and  $M_2$  black; hence it seems reasonable to suppose that the (conditional) probability in question is  $(M_1 - 1)/(M - 1)$ .

We now give a definition of conditional probability that is consistent with our intuitive ideas.

Let  $(\Omega, \mathcal{A}, P)$  be a (finite) probability space and  $A$  an event (i.e.  $A \in \mathcal{A}$ ).

**Definition 1.** The *conditional probability* of event  $B$  assuming event  $A$  with  $P(A) > 0$  (denoted by  $P(B|A)$ ) is

$$\frac{P(AB)}{P(A)}. \quad (1)$$

In the classical approach we have  $P(A) = N(A)/N(\Omega)$ ,  $P(AB) = N(AB)/N(\Omega)$ , and therefore

$$P(B|A) = \frac{N(AB)}{N(A)}. \quad (2)$$

From Definition 1 we immediately get the following properties of conditional probability:

$$\begin{aligned} P(A|A) &= 1, \\ P(\emptyset|A) &= 0, \\ P(B|A) &= 1, \quad B \supseteq A, \\ P(B_1 + B_2|A) &= P(B_1|A) + P(B_2|A). \end{aligned}$$

It follows from these properties that for a given set  $A$  the conditional probability  $P(\cdot|A)$  has the same properties on the space  $(\Omega \cap A, \mathcal{A} \cap A)$ , where  $\mathcal{A} \cap A = \{B \cap A : B \in \mathcal{A}\}$ , that the original probability  $P(\cdot)$  has on  $(\Omega, \mathcal{A})$ .

Note that

$$P(B|A) + P(\bar{B}|A) = 1;$$

however in general

$$\begin{aligned} P(B|A) + P(B|\bar{A}) &\neq 1, \\ P(\bar{B}|A) + P(\bar{B}|\bar{A}) &\neq 1. \end{aligned}$$

**EXAMPLE 1.** Consider a family with two children. We ask for the probability that both children are boys, assuming

- (a) that the older child is a boy;
- (b) that at least one of the children is a boy.

The sample space is

$$\Omega = \{BB, BG, GB, GG\},$$

where BG means that the older child is a boy and the younger is a girl, etc.

Let us suppose that all sample points are equally probable:

$$P(BB) = P(BG) = P(GB) = P(GG) = \frac{1}{4}.$$

Let  $A$  be the event that the older child is a boy, and  $B$ , that the younger child is a boy. Then  $A \cup B$  is the event that at least one child is a boy, and  $AB$  is the event that both children are boys. In question (a) we want the conditional probability  $P(AB|A)$ , and in (b), the conditional probability  $P(AB|A \cup B)$ .

It is easy to see that

$$P(AB|A) = \frac{P(AB)}{P(A)} = \frac{\frac{1}{4}}{\frac{1}{2}} = \frac{1}{2},$$

$$P(AB|A \cup B) = \frac{P(AB)}{P(A \cup B)} = \frac{\frac{1}{4}}{\frac{1}{2}} = \frac{1}{3}.$$

2. The simple but important formula (3), below, is called the formula for total probability. It provides the basic means for calculating the probabilities of complicated events by using conditional probabilities.

Consider a decomposition  $\mathcal{D} = \{A_1, \dots, A_n\}$  with  $P(A_i) > 0, i = 1, \dots, n$  (such a decomposition is often called a complete set of disjoint events). It is clear that

$$B = BA_1 + \dots + BA_n$$

and therefore

$$P(B) = \sum_{i=1}^n P(BA_i).$$

But

$$P(BA_i) = P(B|A_i)P(A_i).$$

Hence we have the *formula for total probability*:

$$P(B) = \sum_{i=1}^n P(B|A_i)P(A_i). \quad (3)$$

In particular, if  $0 < P(A) < 1$ , then

$$P(B) = P(B|A)P(A) + P(B|\bar{A})P(\bar{A}). \quad (4)$$

EXAMPLE 2. An urn contains  $M$  balls,  $m$  of which are "lucky." We ask for the probability that the second ball drawn is lucky (assuming that the result of the first draw is unknown, that a sample of size 2 is drawn without replacement, and that all outcomes are equally probable). Let  $A$  be the event that the first ball is lucky,  $B$  the event that the second is lucky. Then

$$P(B|A) = \frac{P(BA)}{P(A)} = \frac{\frac{m(m-1)}{M(M-1)}}{\frac{m}{M}} = \frac{m-1}{M-1},$$

$$P(B|\bar{A}) = \frac{P(B\bar{A})}{P(\bar{A})} = \frac{\frac{m(M-m)}{M(M-1)}}{\frac{M-m}{M}} = \frac{m}{M-1}$$

and

$$\begin{aligned} P(B) &= P(B|A)P(A) + P(B|\bar{A})P(\bar{A}) \\ &= \frac{m-1}{M-1} \cdot \frac{m}{M} + \frac{m}{M-1} \cdot \frac{M-m}{M} = \frac{m}{M}. \end{aligned}$$

It is interesting to observe that  $P(A)$  is precisely  $m/M$ . Hence, when the nature of the first ball is unknown, it does not affect the probability that the second ball is lucky.

By the definition of conditional probability (with  $P(A) > 0$ ),

$$P(AB) = P(B|A)P(A). \quad (5)$$

This formula, the *multiplication formula for probabilities*, can be generalized (by induction) as follows: If  $A_1, \dots, A_{n-1}$  are events with  $P(A_1 \cdots A_{n-1}) > 0$ , then

$$P(A_1 \cdots A_n) = P(A_1)P(A_2|A_1) \cdots P(A_n|A_1 \cdots A_{n-1}) \quad (6)$$

(here  $A_1 \cdots A_n = A_1 \cap A_2 \cap \cdots \cap A_n$ ).

3. Suppose that  $A$  and  $B$  are events with  $P(A) > 0$  and  $P(B) > 0$ . Then along with (5) we have the parallel formula

$$P(AB) = P(A|B)P(B). \quad (7)$$

From (5) and (7) we obtain *Bayes's formula*

$$P(A|B) = \frac{P(A)P(B|A)}{P(B)}. \quad (8)$$

If the events  $A_1, \dots, A_n$  form a decomposition of  $\Omega$ , (3) and (8) imply *Bayes's theorem*:

$$P(A_i|B) = \frac{P(A_i)P(B|A_i)}{\sum_{j=1}^n P(A_j)P(B|A_j)}. \quad (9)$$

In statistical applications,  $A_1, \dots, A_n$  ( $A_1 + \dots + A_n = \Omega$ ) are often called hypotheses, and  $P(A_i)$  is called the *a priori*† probability of  $A_i$ . The conditional probability  $P(A_i|B)$  is considered as the *a posteriori* probability of  $A_i$  after the occurrence of event  $B$ .

**EXAMPLE 3.** Let an urn contain two coins:  $A_1$ , a fair coin with probability  $\frac{1}{2}$  of falling H; and  $A_2$ , a biased coin with probability  $\frac{1}{3}$  of falling H. A coin is drawn at random and tossed. Suppose that it falls head. We ask for the probability that the fair coin was selected.

Let us construct the corresponding probabilistic model. Here it is natural to take the sample space to be the set  $\Omega = \{A_1H, A_1T, A_2H, A_2T\}$ , which describes all possible outcomes of a selection and a toss ( $A_1H$  means that coin  $A_1$  was selected and fell heads, etc.) The probabilities  $p(\omega)$  of the various outcomes have to be assigned so that, according to the statement of the problem,

$$P(A_1) = P(A_2) = \frac{1}{2}$$

and

$$P(H|A_1) = \frac{1}{2}, \quad P(H|A_2) = \frac{1}{3}.$$

With these assignments, the probabilities of the sample points are uniquely determined:

$$P(A_1H) = \frac{1}{4}, \quad P(A_1T) = \frac{1}{4}, \quad P(A_2H) = \frac{1}{6}, \quad P(A_2T) = \frac{1}{3}.$$

Then by Bayes's formula the probability in question is

$$P(A_1|H) = \frac{P(A_1)P(H|A_1)}{P(A_1)P(H|A_1) + P(A_2)P(H|A_2)} = \frac{3}{5},$$

and therefore

$$P(A_2|H) = \frac{2}{5}.$$

4. In certain sense, the concept of *independence*, which we are now going to introduce, plays a central role in probability theory: it is precisely this concept that distinguishes probability theory from the general theory of measure spaces.

† *A priori*: before the experiment; *a posteriori*: after the experiment.



If  $A$  and  $B$  are two events, it is natural to say that  $B$  is independent of  $A$  if knowing that  $A$  has occurred has no effect on the probability of  $B$ . In other words, " $B$  is independent of  $A$ " if

$$P(B|A) = P(B) \quad (10)$$

(we are supposing that  $P(A) > 0$ ).

Since

$$P(B|A) = \frac{P(AB)}{P(A)},$$

it follows from (10) that

$$P(AB) = P(A)P(B). \quad (11)$$

In exactly the same way, if  $P(B) > 0$  it is natural to say that " $A$  is independent of  $B$ " if

$$P(A|B) = P(A).$$

Hence we again obtain (11), which is symmetric in  $A$  and  $B$  and still makes sense when the probabilities of these events are zero.

After these preliminaries, we introduce the following definition.

**Definition 2.** Events  $A$  and  $B$  are called *independent* or *statistically independent* (with respect to the probability  $P$ ) if

$$P(AB) = P(A)P(B).$$

In probability theory it is often convenient to consider not only independence of events (or sets) but also independence of collections of events (or sets).

Accordingly, we introduce the following definition.

**Definition 3.** Two algebras  $\mathcal{A}_1$  and  $\mathcal{A}_2$  of events (or sets) are called *independent* or *statistically independent* (with respect to the probability  $P$ ) if all pairs of sets  $A_1$  and  $A_2$ , belonging respectively to  $\mathcal{A}_1$  and  $\mathcal{A}_2$ , are independent.

For example, let us consider the two algebras

$$\mathcal{A}_1 = \{A_1, \bar{A}_1, \emptyset, \Omega\} \quad \text{and} \quad \mathcal{A}_2 = \{A_2, \bar{A}_2, \emptyset, \Omega\},$$

where  $A_1$  and  $A_2$  are subsets of  $\Omega$ . It is easy to verify that  $\mathcal{A}_1$  and  $\mathcal{A}_2$  are independent if and only if  $A_1$  and  $A_2$  are independent. In fact, the independence of  $\mathcal{A}_1$  and  $\mathcal{A}_2$  means the independence of the 16 events  $A_1$  and  $A_2$ ,  $A_1$  and  $\bar{A}_2$ , ...,  $\Omega$  and  $\Omega$ . Consequently  $A_1$  and  $A_2$  are independent. Conversely, if  $A_1$  and  $A_2$  are independent, we have to show that the other 15

pairs of events are independent. Let us verify, for example, the independence of  $A_1$  and  $\bar{A}_2$ . We have

$$\begin{aligned} P(A_1\bar{A}_2) &= P(A_1) - P(A_1A_2) = P(A_1) - P(A_1)P(A_2) \\ &= P(A_1) \cdot (1 - P(A_2)) = P(A_1)P(\bar{A}_2). \end{aligned}$$

The independence of the other pairs is verified similarly.

5. The concept of independence of two sets or two algebras of sets can be extended to any finite number of sets or algebras of sets.

Thus we say that the sets  $A_1, \dots, A_n$  are collectively *independent* or *statistically independent* (with respect to the probability  $P$ ) if for  $k = 1, \dots, n$  and  $1 \leq i_1 < i_2 < \dots < i_k \leq n$

$$P(A_{i_1} \dots A_{i_k}) = P(A_{i_1}) \dots P(A_{i_k}). \quad (12)$$

The algebras  $\mathcal{A}_1, \dots, \mathcal{A}_n$  of sets are called *independent* or *statistically independent* (with respect to the probability  $P$ ) if all sets  $A_1, \dots, A_n$  belonging respectively to  $\mathcal{A}_1, \dots, \mathcal{A}_n$  are independent.

Note that *pairwise independence* of events *does not imply* their independence. In fact if, for example,  $\Omega = \{\omega_1, \omega_2, \omega_3, \omega_4\}$  and all outcomes are equiprobable, it is easily verified that the events

$$A = \{\omega_1, \omega_2\}, \quad B = \{\omega_1, \omega_3\}, \quad C = \{\omega_1, \omega_4\}$$

are pairwise independent, whereas

$$P(ABC) = \frac{1}{4} \neq \left(\frac{1}{2}\right)^3 = P(A)P(B)P(C).$$

Also note that if

$$P(ABC) = P(A)P(B)P(C)$$

for events  $A, B$  and  $C$ , it by no means follows that these events are pairwise independent. In fact, let  $\Omega$  consist of the 36 ordered pairs  $(i, j)$ , where  $i, j = 1, 2, \dots, 6$  and all the pairs are equiprobable. Then if  $A = \{(i, j): j = 1, 2 \text{ or } 5\}$ ,  $B = \{(i, j): j = 4, 5 \text{ or } 6\}$ ,  $C = \{(i, j): i + j = 9\}$  we have

$$P(AB) = \frac{1}{6} \neq \frac{1}{4} = P(A)P(B),$$

$$P(AC) = \frac{1}{36} \neq \frac{1}{18} = P(A)P(C),$$

$$P(BC) = \frac{1}{12} \neq \frac{1}{18} = P(B)P(C),$$

but also

$$P(ABC) = \frac{1}{36} = P(A)P(B)P(C).$$

6. Let us consider in more detail, from the point of view of independence, the classical model  $(\Omega, \mathcal{A}, P)$  that was introduced in §2 and used as a basis for the binomial distribution.

In this model

$$\Omega = \{\omega: \omega = (a_1, \dots, a_n), a_i = 0, 1\}, \quad \mathcal{A} = \{A: A \subseteq \Omega\}$$

and

$$p(\omega) = p^{\sum a_i} q^{n - \sum a_i}. \quad (13)$$

Consider an event  $A \subseteq \Omega$ . We say that this event depends on a trial at time  $k$  if it is determined by the value  $a_k$  alone. Examples of such events are

$$A_k = \{\omega: a_k = 1\}, \quad \bar{A}_k = \{\omega: a_k = 0\}.$$

Let us consider the sequence of algebras  $\mathcal{A}_1, \mathcal{A}_2, \dots, \mathcal{A}_n$ , where  $\mathcal{A}_k = \{A_k, \bar{A}_k, \emptyset, \Omega\}$  and show that under (13) these algebras are independent.

It is clear that

$$\begin{aligned} P(A_k) &= \sum_{\{\omega: a_k = 1\}} p(\omega) = \sum_{\{\omega: a_k = 1\}} p^{\sum a_i} q^{n - \sum a_i} \\ &= p \sum_{(a_1, \dots, a_{k-1}, a_{k+1}, \dots, a_n)} p^{a_1 + \dots + a_{k-1} + a_{k+1} + \dots + a_n} \\ &\quad \times q^{(n-1) - (a_1 + \dots + a_{k-1} + a_{k+1} + \dots + a_n)} = p \sum_{i=0}^{n-1} C_{n-1}^i p^i q^{(n-1)-i} = p, \end{aligned}$$

and a similar calculation shows that  $P(\bar{A}_k) = q$  and that, for  $k \neq l$ ,

$$P(A_k A_l) = p^2, \quad P(A_k \bar{A}_l) = pq, \quad P(\bar{A}_k A_l) = q^2.$$

It is easy to deduce from this that  $\mathcal{A}_k$  and  $\mathcal{A}_l$  are independent for  $k \neq l$ .

It can be shown in the same way that  $\mathcal{A}_1, \mathcal{A}_2, \dots, \mathcal{A}_n$  are independent. This is the basis for saying that our model  $(\Omega, \mathcal{A}, P)$  corresponds to “ $n$  independent trials with two outcomes and probability  $p$  of success.” James Bernoulli was the first to study this model systematically, and established the law of large numbers (§5) for it. Accordingly, this model is also called the Bernoulli scheme with two outcomes (success and failure) and probability  $p$  of success.

A detailed study of the probability space for the Bernoulli scheme shows that it has the structure of a direct product of probability spaces, defined as follows.

Suppose that we are given a collection  $(\Omega_1, \mathcal{B}_1, P_1), \dots, (\Omega_n, \mathcal{B}_n, P_n)$  of finite probability spaces. Form the space  $\Omega = \Omega_1 \times \Omega_2 \times \dots \times \Omega_n$  of points  $\omega = (a_1, \dots, a_n)$ , where  $a_i \in \Omega_i$ . Let  $\mathcal{A} = \mathcal{B}_1 \otimes \dots \otimes \mathcal{B}_n$  be the algebra of the subsets of  $\Omega$  that consists of sums of sets of the form

$$A = B_1 \times B_2 \times \dots \times B_n$$

with  $B_i \in \mathcal{B}_i$ . Finally, for  $\omega = (a_1, \dots, a_n)$  take  $p(\omega) = p_1(a_1) \cdots p_n(a_n)$  and define  $P(A)$  for the set  $A = B_1 \times B_2 \times \dots \times B_n$  by

$$P(A) = \sum_{\{a_1 \in B_1, \dots, a_n \in B_n\}} p_1(a_1) \cdots p_n(a_n).$$

It is easy to verify that  $P(\Omega) = 1$  and therefore the triple  $(\Omega, \mathcal{A}, P)$  defines a probability space. This space is called the *direct product of the probability spaces*  $(\Omega_1, \mathcal{B}_1, P_1), \dots, (\Omega_n, \mathcal{B}_n, P_n)$ .

We note an easily verified property of the direct product of probability spaces: with respect to  $P$ , the events

$$A_1 = \{\omega: a_1 \in B_1\}, \dots, A_n = \{\omega: a_n \in B_n\},$$

where  $B_i \in \mathcal{B}_i$ , are independent. In the same way, the algebras of subsets of  $\Omega$ ,

$$\mathcal{A}_1 = \{A_1: A_1 = \{\omega: a_1 \in B_1\}, B_1 \in \mathcal{B}_1\},$$

$$\dots\dots\dots$$

$$\mathcal{A}_n = \{A_n: A_n = \{\omega: a_n \in B_n\}, B_n \in \mathcal{B}_n\}$$

are independent.

It is clear from our construction that the Bernoulli scheme

$$(\Omega, \mathcal{A}, P) \quad \text{with} \quad \Omega = \{\omega: \omega = (a_1, \dots, a_n), a_i = 0 \text{ or } 1\}$$

$$\mathcal{A} = \{A: A \subseteq \Omega\} \quad \text{and} \quad p(\omega) = p^{\sum a_i} q^{n - \sum a_i}$$

can be thought of as the direct product of the probability spaces  $(\Omega_i, \mathcal{B}_i, P_i)$ ,  $i = 1, 2, \dots, n$ , where

$$\Omega_i = \{0, 1\}, \quad \mathcal{B}_i = \{\{0\}, \{1\}, \emptyset, \Omega_i\},$$

$$P_i(\{1\}) = p, \quad P_i(\{0\}) = q.$$

## 7. PROBLEMS

1. Give examples to show that in general the equations

$$P(B|A) + P(B|\bar{A}) = 1,$$

$$P(B|A) + P(\bar{B}|\bar{A}) = 1$$

are false.

2. An urn contains  $M$  balls, of which  $M_1$  are white. Consider a sample of size  $n$ . Let  $B_j$  be the event that the ball selected at the  $j$ th step is white, and  $A_k$  the event that a sample of size  $n$  contains exactly  $k$  white balls. Show that

$$P(B_j|A_k) = k/n$$

both for sampling with replacement and for sampling without replacement.

3. Let  $A_1, \dots, A_n$  be independent events. Then

$$P\left(\bigcup_{i=1}^n A_i\right) = 1 - \prod_{i=1}^n P(\bar{A}_i).$$

4. Let  $A_1, \dots, A_n$  be independent events with  $P(A_i) = p_i$ . Then the probability  $P_0$  that neither event occurs is

$$P_0 = \prod_{i=1}^n (1 - p_i).$$

5. Let  $A$  and  $B$  be independent events. In terms of  $P(A)$  and  $P(B)$ , find the probabilities of the events that exactly  $k$ , at least  $k$ , and at most  $k$  of  $A$  and  $B$  occur ( $k = 0, 1, 2$ ).
6. Let event  $A$  be independent of itself, i.e. let  $A$  and  $A$  be independent. Show that  $P(A)$  is either 0 or 1.
7. Let event  $A$  have  $P(A) = 0$  or 1. Show that  $A$  and an arbitrary event  $B$  are independent.
8. Consider the electric circuit shown in Figure 4:

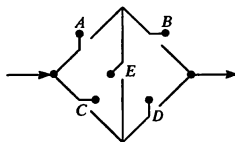


Figure 4

Each of the switches  $A$ ,  $B$ ,  $C$ ,  $D$ , and  $E$  is independently open or closed with probabilities  $p$  and  $q$ , respectively. Find the probability that a signal fed in at "input" will be received at "output". If the signal is received, what is the conditional probability that  $E$  is open?

## §4. Random Variables and Their Properties

1. Let  $(\Omega, \mathcal{A}, P)$  be a probabilistic model of an experiment with a finite number of outcomes,  $N(\Omega) < \infty$ , where  $\mathcal{A}$  is the algebra of all subsets of  $\Omega$ . We observe that in the examples above, where we calculated the probabilities of various events  $A \in \mathcal{A}$ , the specific nature of the sample space  $\Omega$  was of no interest. We were interested only in numerical properties depending on the sample points. For example, we were interested in the probability of some number of successes in a series of  $n$  trials, in the probability distribution for the number of objects in cells, etc.

The concept "random variable," which we now introduce (later it will be given a more general form) serves to define quantities that are subject to "measurement" in random experiments.

**Definition 1.** Any numerical function  $\xi = \xi(\omega)$  defined on a (finite) sample space  $\Omega$  is called a (simple) *random variable*. (The reason for the term "simple" random variable will become clear after the introduction of the general concept of random variable in §4 of Chapter II.)

EXAMPLE 1. In the model of two tosses of a coin with sample space  $\Omega = \{HH, HT, TH, TT\}$ , define a random variable  $\xi = \xi(\omega)$  by the table

$\omega$	HH	HT	TH	TT
$\xi(\omega)$	2	1	1	0

Here, from its very definition,  $\xi(\omega)$  is nothing but the number of heads in the outcome  $\omega$ .

Another extremely simple example of a random variable is the *indicator* (or *characteristic function*) of a set  $A \in \mathcal{A}$ :

$$\xi = I_A(\omega),$$

where†

$$I_A(\omega) = \begin{cases} 1, & \omega \in A, \\ 0, & \omega \notin A. \end{cases}$$

When experimenters are concerned with random variables that describe observations, their main interest is in the probabilities with which the random variables take various values. From this point of view they are interested, not in the distribution of the probability  $P$  over  $(\Omega, \mathcal{A})$ , but in its distribution over the range of a random variable. Since we are considering the case when  $\Omega$  contains only a finite number of points, the range  $X$  of the random variable  $\xi$  is also finite. Let  $X = \{x_1, \dots, x_m\}$ , where the (different) numbers  $x_1, \dots, x_m$  exhaust the values of  $\xi$ .

Let  $\mathcal{X}$  be the collection of all subsets of  $X$ , and let  $B \in \mathcal{X}$ . We can also interpret  $B$  as an event if the sample space is taken to be  $X$ , the set of values of  $\xi$ .

On  $(X, \mathcal{X})$ , consider the probability  $P_\xi(\cdot)$  induced by  $\xi$  according to the formula

$$P_\xi(B) = P\{\omega: \xi(\omega) \in B\}, \quad B \in \mathcal{X}.$$

It is clear that the values of this probability are completely determined by the probabilities

$$P_\xi(x_i) = P\{\omega: \xi(\omega) = x_i\}, \quad x_i \in X.$$

The set of numbers  $\{P_\xi(x_1), \dots, P_\xi(x_m)\}$  is called the *probability distribution of the random variable*  $\xi$ .

† The notation  $I(A)$  is also used. For frequently used properties of indicators see Problem 1.

EXAMPLE 2. A random variable  $\xi$  that takes the two values 1 and 0 with probabilities  $p$  ("success") and  $q$  ("failure"), is called a Bernoulli† random variable. Clearly

$$P_{\xi}(x) = p^x q^{1-x}, \quad x = 0, 1. \quad (1)$$

A *binomial* (or binomially distributed) *random variable*  $\xi$  is a random variable that takes the  $n + 1$  values  $0, 1, \dots, n$  with probabilities

$$P_{\xi}(x) = C_n^x p^x q^{n-x}, \quad x = 0, 1, \dots, n. \quad (2)$$

Note that here and in many subsequent examples we do not specify the sample spaces  $(\Omega, \mathcal{A}, \mathbf{P})$ , but are interested only in the values of the random variables and their probability distributions.

The probabilistic structure of the random variables  $\xi$  is completely specified by the probability distributions  $\{P_{\xi}(x_i), i = 1, \dots, m\}$ . The concept of distribution function, which we now introduce, yields an equivalent description of the probabilistic structure of the random variables.

**Definition 2.** Let  $x \in \mathbf{R}^1$ . The function

$$F_{\xi}(x) = \mathbf{P}\{\omega: \xi(\omega) \leq x\}$$

is called the *distribution function* of the random variable  $\xi$ .

Clearly

$$F_{\xi}(x) = \sum_{\{i: x_i \leq x\}} P_{\xi}(x_i)$$

and

$$P_{\xi}(x_i) = F_{\xi}(x_i) - F_{\xi}(x_{i-1}),$$

where  $F_{\xi}(x-) = \lim_{y \uparrow x} F_{\xi}(y)$ .

If we suppose that  $x_1 < x_2 < \dots < x_m$  and put  $F_{\xi}(x_0) = 0$ , then

$$P_{\xi}(x_i) = F_{\xi}(x_i) - F_{\xi}(x_{i-1}), \quad i = 1, \dots, m.$$

The following diagrams (Figure 5) exhibit  $P_{\xi}(x)$  and  $F_{\xi}(x)$  for a binomial random variable.

It follows immediately from Definition 2 that the distribution  $F_{\xi} = F_{\xi}(x)$  has the following properties:

- (1)  $F_{\xi}(-\infty) = 0, F_{\xi}(+\infty) = 1$ ;
- (2)  $F_{\xi}(x)$  is continuous on the right ( $F_{\xi}(x+) = F_{\xi}(x)$ ) and piecewise constant.

† We use the terms "Bernoulli, binomial, Poisson, Gaussian, . . . , random variables" for what are more usually called random variables with Bernoulli, binomial, Poisson, Gaussian, . . . , distributions.

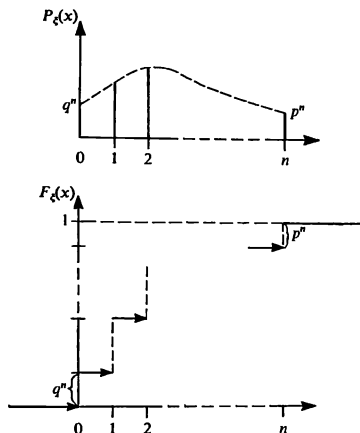


Figure 5

Along with random variables it is often necessary to consider *random vectors*  $\xi = (\xi_1, \dots, \xi_r)$  whose components are random variables. For example, when we considered the multinomial distribution we were dealing with a random vector  $v = (v_1, \dots, v_r)$ , where  $v_i = v_i(\omega)$  is the number of elements equal to  $b_i$ ,  $i = 1, \dots, r$ , in the sequence  $\omega = (a_1, \dots, a_n)$ .

The set of probabilities

$$P_\xi(x_1, \dots, x_r) = \mathbf{P}\{\omega: \xi_1(\omega) = x_1, \dots, \xi_r(\omega) = x_r\},$$

where  $x_i \in X_i$ , the range of  $\xi_i$ , is called the *probability distribution of the random vector*  $\xi$ , and the function

$$F_\xi(x_1, \dots, x_r) = \mathbf{P}\{\omega: \xi_1(\omega) \leq x_1, \dots, \xi_r(\omega) \leq x_r\},$$

where  $x_i \in R^1$ , is called the *distribution function of the random vector*  $\xi = (\xi_1, \dots, \xi_r)$ .

For example, for the random vector  $v = (v_1, \dots, v_r)$  mentioned above,

$$P_v(n_1, \dots, n_r) = C_n(n_1, \dots, n_r) p_1^{n_1} \cdots p_r^{n_r}$$

(see (2.2)).

2. Let  $\xi_1, \dots, \xi_r$  be a set of random variables with values in a (finite) set  $X \subseteq R^1$ . Let  $\mathcal{X}$  be the algebra of subsets of  $X$ .



**Definition 3.** The random variables  $\xi_1, \dots, \xi_r$  are said to be *independent* (collectively independent) if

$$P\{\xi_1 = x_1, \dots, \xi_r = x_r\} = P\{\xi_1 = x_1\} \cdots P\{\xi_r = x_r\}$$

for all  $x_1, \dots, x_r \in X$ ; or, equivalently, if

$$P\{\xi_1 \in B_1, \dots, \xi_r \in B_r\} = P\{\xi_1 \in B_1\} \cdots P\{\xi_r \in B_r\}$$

for all  $B_1, \dots, B_r \in \mathcal{X}$ .

We can get a very simple example of independent random variables from the Bernoulli scheme. Let

$$\Omega = \{\omega: \omega = (a_1, \dots, a_n), a_i = 0, 1\}, \quad p(\omega) = p^{\sum a_i} q^{n - \sum a_i}$$

and  $\xi_i(\omega) = a_i$  for  $\omega = (a_1, \dots, a_n)$ ,  $i = 1, \dots, n$ . Then the random variables  $\xi_1, \xi_2, \dots, \xi_n$  are independent, as follows from the independence of the events

$$A_1 = \{\omega: a_1 = 1\}, \dots, A_n = \{\omega: a_n = 1\},$$

which was established in §3.

3. We shall frequently encounter the problem of finding the probability distributions of random variables that are functions  $f(\xi_1, \dots, \xi_r)$  of random variables  $\xi_1, \dots, \xi_r$ . For the present we consider only the determination of the distribution of a sum  $\zeta = \xi + \eta$  of random variables.

If  $\xi$  and  $\eta$  take values in the respective sets  $X = \{x_1, \dots, x_k\}$  and  $Y = \{y_1, \dots, y_l\}$ , the random variable  $\zeta = \xi + \eta$  takes values in the set  $Z = \{z: z = x_i + y_j, i = 1, \dots, k; j = 1, \dots, l\}$ . Then it is clear that

$$P_\zeta(z) = P\{\zeta = z\} = P\{\xi + \eta = z\} = \sum_{\{(i, j): x_i + y_j = z\}} P\{\xi = x_i, \eta = y_j\}.$$

The case of independent random variables  $\xi$  and  $\eta$  is particularly important. In this case

$$P\{\xi = x_i, \eta = y_j\} = P\{\xi = x_i\}P\{\eta = y_j\},$$

and therefore

$$P_\zeta(z) = \sum_{\{(i, j): x_i + y_j = z\}} P_\xi(x_i)P_\eta(y_j) = \sum_{i=1}^k P_\xi(x_i)P_\eta(z - x_i) \quad (3)$$

for all  $z \in Z$ , where in the last sum  $P_\eta(z - x_i)$  is taken to be zero if  $z - x_i \notin Y$ .

For example, if  $\xi$  and  $\eta$  are independent Bernoulli random variables, taking the values 1 and 0 with respective probabilities  $p$  and  $q$ , then  $Z = \{0, 1, 2\}$  and

$$P_\zeta(0) = P_\xi(0)P_\eta(0) = q^2,$$

$$P_\zeta(1) = P_\xi(0)P_\eta(1) + P_\xi(1)P_\eta(0) = 2pq,$$

$$P_\zeta(2) = P_\xi(1)P_\eta(1) = p^2.$$

It is easy to show by induction that if  $\xi_1, \xi_2, \dots, \xi_n$  are independent Bernoulli random variables with  $P\{\xi_i = 1\} = p$ ,  $P\{\xi_i = 0\} = q$ , then the random variable  $\zeta = \xi_1 + \dots + \xi_n$  has the binomial distribution

$$P_\zeta(k) = C_n^k p^k q^{n-k}, \quad k = 0, 1, \dots, n. \quad (4)$$

4. We now turn to the important concept of the expectation, or mean value, of a random variable.

Let  $(\Omega, \mathcal{A}, P)$  be a (finite) probability space and  $\xi = \xi(\omega)$  a random variable with values in the set  $X = \{x_1, \dots, x_k\}$ . If we put  $A_i = \{\omega: \xi = x_i\}$ ,  $i = 1, \dots, k$ , then  $\xi$  can evidently be represented as

$$\xi(\omega) = \sum_{i=1}^k x_i I(A_i), \quad (5)$$

where the sets  $A_1, \dots, A_k$  form a decomposition of  $\Omega$  (i.e., they are pairwise disjoint and their sum is  $\Omega$ ; see Subsection 3 of §1).

Let  $p_i = P\{\xi = x_i\}$ . It is intuitively plausible that if we observe the values of the random variable  $\xi$  in “ $n$  repetitions of identical experiments”, the value  $x_i$  ought to be encountered about  $p_i n$  times,  $i = 1, \dots, k$ . Hence the mean value calculated from the results of  $n$  experiments is roughly

$$\frac{1}{n} [np_1 x_1 + \dots + np_k x_k] = \sum_{i=1}^k p_i x_i.$$

This discussion provides the motivation for the following definition.

**Definition 4.** The *expectation*<sup>†</sup> or *mean value* of the random variable  $\xi = \sum_{i=1}^k x_i I(A_i)$  is the number

$$E\xi = \sum_{i=1}^k x_i P(A_i). \quad (6)$$

Since  $A_i = \{\omega: \xi(\omega) = x_i\}$  and  $P_\xi(x_i) = P(A_i)$ , we have

$$E\xi = \sum_{i=1}^k x_i P_\xi(x_i). \quad (7)$$

Recalling the definition of  $F_\xi = F_\xi(x)$  and writing

$$\Delta F_\xi(x) = F_\xi(x) - F_\xi(x-),$$

we obtain  $P_\xi(x_i) = \Delta F_\xi(x_i)$  and consequently

$$E\xi = \sum_{i=1}^k x_i \Delta F_\xi(x_i). \quad (8)$$

<sup>†</sup> Also known as mathematical expectation, or expected value, or (especially in physics) expectation value. (Translator)

Before discussing the properties of the expectation, we remark that it is often convenient to use another representation of the random variable  $\xi$ , namely

$$\xi(\omega) = \sum_{j=1}^l x'_j I(B_j),$$

where  $B_1 + \dots + B_l = \Omega$ , but some of the  $x'_j$  may be repeated. In this case  $E\xi$  can be calculated from the formula  $\sum_{j=1}^l x'_j P(B_j)$ , which differs formally from (5) because in (5) the  $x_i$  are all different. In fact,

$$\sum_{\{j: x'_j = x_i\}} x'_j P(B_j) = x_i \sum_{\{j: x'_j = x_i\}} P(B_j) = x_i P(A_i)$$

and therefore

$$\sum_{j=1}^l x'_j P(B_j) = \sum_{i=1}^k x_i P(A_i).$$

5. We list the basic properties of the expectation:

- (1) If  $\xi \geq 0$  then  $E\xi \geq 0$ .
- (2)  $E(a\xi + b\eta) = aE\xi + bE\eta$ , where  $a$  and  $b$  are constants.
- (3) If  $\xi \geq \eta$  then  $E\xi \geq E\eta$ .
- (4)  $|E\xi| \leq E|\xi|$ .
- (5) If  $\xi$  and  $\eta$  are independent, then  $E\xi\eta = E\xi \cdot E\eta$ .
- (6)  $(E|\xi\eta|)^2 \leq E\xi^2 \cdot E\eta^2$  (Cauchy-Bunyakovskii inequality).†
- (7) If  $\xi = I(A)$  then  $E\xi = P(A)$ .

Properties (1) and (7) are evident. To prove (2), let

$$\xi = \sum_i x_i I(A_i), \quad \eta = \sum_j y_j I(B_j).$$

Then

$$\begin{aligned} a\xi + b\eta &= a \sum_{i,j} x_i I(A_i \cap B_j) + b \sum_{i,j} y_j I(A_i \cap B_j) \\ &= \sum_{i,j} (ax_i + by_j) I(A_i \cap B_j) \end{aligned}$$

and

$$\begin{aligned} E(a\xi + b\eta) &= \sum_{i,j} (ax_i + by_j) P(A_i \cap B_j) \\ &= \sum_i ax_i P(A_i) + \sum_j by_j P(B_j) \\ &= a \sum_i x_i P(A_i) + b \sum_j y_j P(B_j) = aE\xi + bE\eta. \end{aligned}$$

† Also known as the Cauchy-Schwarz or Schwarz inequality. (Translator)

Property (3) follows from (1) and (2). Property (4) is evident, since

$$|E\xi| = \left| \sum_i x_i P(A_i) \right| \leq \sum_i |x_i| P(A_i) = E|\xi|.$$

To prove (5) we note that

$$\begin{aligned} E\xi\eta &= E\left(\sum_i x_i I(A_i)\right) \left(\sum_j y_j I(B_j)\right) \\ &= E \sum_{i,j} x_i y_j I(A_i \cap B_j) = \sum_{i,j} x_i y_j P(A_i \cap B_j) \\ &= \sum_{i,j} x_i y_j P(A_i) P(B_j) \\ &= \left(\sum_i x_i P(A_i)\right) \cdot \left(\sum_j y_j P(B_j)\right) = E\xi \cdot E\eta, \end{aligned}$$

where we have used the property that for independent random variables the events

$$A_i = \{\omega: \xi(\omega) = x_i\} \quad \text{and} \quad B_j = \{\omega: \eta(\omega) = y_j\}$$

are independent:  $P(A_i \cap B_j) = P(A_i)P(B_j)$ .

To prove property (6) we observe that

$$\xi^2 = \sum_i x_i^2 I(A_i), \quad \eta^2 = \sum_j y_j^2 I(B_j)$$

and

$$E\xi^2 = \sum_i x_i^2 P(A_i), \quad E\eta^2 = \sum_j y_j^2 P(B_j).$$

Let  $E\xi^2 > 0$ ,  $E\eta^2 > 0$ . Put

$$\tilde{\xi} = \frac{\xi}{\sqrt{E\xi^2}}, \quad \tilde{\eta} = \frac{\eta}{\sqrt{E\eta^2}}.$$

Since  $2|\tilde{\xi}\tilde{\eta}| \leq \tilde{\xi}^2 + \tilde{\eta}^2$ , we have  $2E|\tilde{\xi}\tilde{\eta}| \leq E\tilde{\xi}^2 + E\tilde{\eta}^2 = 2$ . Therefore  $E|\tilde{\xi}\tilde{\eta}| \leq 1$  and  $(E|\xi\eta|)^2 \leq E\xi^2 \cdot E\eta^2$ .

However, if, say,  $E\xi^2 = 0$ , this means that  $\sum_i x_i^2 P(A_i) = 0$  and consequently the mean value of  $\xi$  is 0, and  $P\{\omega: \xi(\omega) = 0\} = 1$ . Therefore if at least one of  $E\xi^2$  or  $E\eta^2$  is zero, it is evident that  $E|\xi\eta| = 0$  and consequently the Cauchy–Bunyakovskii inequality still holds.

**Remark.** Property (5) generalizes in an obvious way to any finite number of random variables: if  $\xi_1, \dots, \xi_r$  are independent, then

$$E\xi_1 \cdots \xi_r = E\xi_1 \cdots E\xi_r.$$

The proof can be given in the same way as for the case  $r = 2$ , or by induction.

EXAMPLE 3. Let  $\xi$  be a Bernoulli random variable, taking the values 1 and 0 with probabilities  $p$  and  $q$ . Then

$$E\xi = 1 \cdot P\{\xi = 1\} + 0 \cdot P\{\xi = 0\} = p.$$

EXAMPLE 4. Let  $\xi_1, \dots, \xi_n$  be  $n$  Bernoulli random variables with  $P\{\xi_i = 1\} = p$ ,  $P\{\xi_i = 0\} = q$ ,  $p + q = 1$ . Then if

$$S_n = \xi_1 + \dots + \xi_n$$

we find that

$$ES_n = np.$$

This result can be obtained in a different way. It is easy to see that  $ES_n$  is not changed if we assume that the Bernoulli random variables  $\xi_1, \dots, \xi_n$  are independent. With this assumption, we have according to (4)

$$P(S_n = k) = C_n^k p^k q^{n-k}, \quad k = 0, 1, \dots, n.$$

Therefore

$$\begin{aligned} ES_n &= \sum_{k=0}^n k P(S_n = k) = \sum_{k=0}^n k C_n^k p^k q^{n-k} \\ &= \sum_{k=0}^n k \cdot \frac{n!}{k! (n-k)!} p^k q^{n-k} \\ &= np \sum_{k=1}^n \frac{(n-1)!}{(k-1)! ((n-1) - (k-1))!} p^{k-1} q^{(n-1) - (k-1)} \\ &= np \sum_{l=0}^n \frac{(n-1)!}{l! ((n-1) - l)!} p^l q^{(n-1) - l} = np. \end{aligned}$$

However, the first method is more direct.

6. Let  $\xi = \sum_i x_i I(A_i)$ , where  $A_i = \{\omega: \xi(\omega) = x_i\}$ , and  $\varphi = \varphi(\xi(\omega))$  is a function of  $\xi(\omega)$ . If  $B_j = \{\omega: \varphi(\xi(\omega)) = y_j\}$ , then

$$\varphi(\xi(\omega)) = \sum_j y_j I(B_j),$$

and consequently

$$E\varphi = \sum_j y_j P(B_j) = \sum_j y_j P_\varphi(y_j). \quad (9)$$

But it is also clear that

$$\varphi(\xi(\omega)) = \sum_i \varphi(x_i) I(A_i).$$

Hence, as in (9), the expectation of the random variable  $\varphi = \varphi(\xi)$  can be calculated as

$$E\varphi(\xi) = \sum_i \varphi(x_i)P_\xi(x_i).$$

7. The important notion of the variance of a random variable  $\xi$  indicates the amount of scatter of the values of  $\xi$  around  $E\xi$ .

**Definition 5.** The *variance* (also called the *dispersion*) of the random variable  $\xi$  (denoted by  $V\xi$ ) is

$$V\xi = E(\xi - E\xi)^2.$$

The number  $\sigma = +\sqrt{V\xi}$  is called the standard deviation.

Since

$$E(\xi - E\xi)^2 = E(\xi^2 - 2\xi \cdot E\xi + (E\xi)^2) = E\xi^2 - (E\xi)^2,$$

we have

$$V\xi = E\xi^2 - (E\xi)^2.$$

Clearly  $V\xi \geq 0$ . It follows from the definition that

$$V(a + b\xi) = b^2V\xi, \quad \text{where } a \text{ and } b \text{ are constants.}$$

In particular,  $Va = 0$ ,  $V(b\xi) = b^2V\xi$ .

Let  $\xi$  and  $\eta$  be random variables. Then

$$\begin{aligned} V(\xi + \eta) &= E((\xi - E\xi) + (\eta - E\eta))^2 \\ &= V\xi + V\eta + 2E(\xi - E\xi)(\eta - E\eta). \end{aligned}$$

Write

$$\text{cov}(\xi, \eta) = E(\xi - E\xi)(\eta - E\eta).$$

This number is called the *covariance* of  $\xi$  and  $\eta$ . If  $V\xi > 0$  and  $V\eta > 0$ , then

$$\rho(\xi, \eta) = \frac{\text{cov}(\xi, \eta)}{\sqrt{V\xi \cdot V\eta}}$$

is called the *correlation coefficient* of  $\xi$  and  $\eta$ . It is easy to show (see Problem 7 below) that if  $\rho(\xi, \eta) = \pm 1$ , then  $\xi$  and  $\eta$  are linearly dependent:

$$\eta = a\xi + b,$$

with  $a > 0$  if  $\rho(\xi, \eta) = 1$  and  $a < 0$  if  $\rho(\xi, \eta) = -1$ .

We observe immediately that if  $\xi$  and  $\eta$  are independent, so are  $\xi - E\xi$  and  $\eta - E\eta$ . Consequently by Property (5) of expectations,

$$\text{cov}(\xi, \eta) = E(\xi - E\xi) \cdot E(\eta - E\eta) = 0.$$

Using the notation that we introduced for covariance, we have

$$V(\xi + \eta) = V\xi + V\eta + 2\text{cov}(\xi, \eta); \quad (10)$$

if  $\xi$  and  $\eta$  are independent, the variance of the sum  $\xi + \eta$  is equal to the sum of the variances,

$$V(\xi + \eta) = V\xi + V\eta. \quad (11)$$

It follows from (10) that (11) is still valid under weaker hypotheses than the independence of  $\xi$  and  $\eta$ . In fact, it is enough to suppose that  $\xi$  and  $\eta$  are uncorrelated, i.e.  $\text{cov}(\xi, \eta) = 0$ .

**Remark.** If  $\xi$  and  $\eta$  are uncorrelated, it does not follow in general that they are independent. Here is a simple example. Let the random variable  $\alpha$  take the values  $0, \pi/2$  and  $\pi$  with probability  $\frac{1}{3}$ . Then  $\xi = \sin \alpha$  and  $\eta = \cos \alpha$  are uncorrelated; however, they are not only stochastically dependent (i.e., not independent with respect to the probability  $P$ ):

$$P\{\xi = 1, \eta = 1\} = 0 \neq \frac{1}{9} = P\{\xi = 1\}P\{\eta = 1\},$$

but even functionally dependent:  $\xi^2 + \eta^2 = 1$ .

Properties (10) and (11) can be extended in the obvious way to any number of random variables:

$$V\left(\sum_{i=1}^n \xi_i\right) = \sum_{i=1}^n V\xi_i + 2 \sum_{i>j} \text{cov}(\xi_i, \xi_j). \quad (12)$$

In particular, if  $\xi_1, \dots, \xi_n$  are pairwise independent (pairwise uncorrelated is sufficient), then

$$V\left(\sum_{i=1}^n \xi_i\right) = \sum_{i=1}^n V\xi_i. \quad (13)$$

**EXAMPLE 5.** If  $\xi$  is a Bernoulli random variable, taking the values 1 and 0 with probabilities  $p$  and  $q$ , then

$$V\xi = E(\xi - E\xi)^2 = (\xi - p)^2 = (1 - p)^2 p + p^2 q = pq.$$

It follows that if  $\xi_1, \dots, \xi_n$  are independent identically distributed Bernoulli random variables, and  $S_n = \xi_1 + \dots + \xi_n$ , then

$$VS_n = npq. \quad (14)$$

**8.** Consider two random variables  $\xi$  and  $\eta$ . Suppose that only  $\xi$  can be observed. If  $\xi$  and  $\eta$  are correlated, we may expect that knowing the value of  $\xi$  allows us to make some inference about the values of the unobserved variable  $\eta$ .

Any function  $f = f(\xi)$  of  $\xi$  is called an *estimator* for  $\eta$ . We say that an estimator  $f^* = f^*(\xi)$  is *best in the mean-square sense* if

$$E(\eta - f^*(\xi))^2 = \inf_f E(\eta - f(\xi))^2.$$

Let us show how to find a best estimator in the class of linear estimators  $\lambda(\xi) = a + b\xi$ . We consider the function  $g(a, b) = E(\eta - (a + b\xi))^2$ . Differentiating  $g(a, b)$  with respect to  $a$  and  $b$ , we obtain

$$\begin{aligned}\frac{\partial g(a, b)}{\partial a} &= -2E[\eta - (a + b\xi)], \\ \frac{\partial g(a, b)}{\partial b} &= -2E[(\eta - (a + b\xi))\xi],\end{aligned}$$

whence, setting the derivatives equal to zero, we find that the best mean-square linear estimator is  $\lambda^*(\xi) = a^* + b^*\xi$ , where

$$a^* = E\eta - b^*E\xi, \quad b^* = \frac{\text{cov}(\xi, \eta)}{V\xi}. \quad (15)$$

In other words,

$$\lambda^*(\xi) = E\eta + \frac{\text{cov}(\xi, \eta)}{V\xi}(\xi - E\xi). \quad (16)$$

The number  $E(\eta - \lambda^*(\xi))^2$  is called the *mean-square error of observation*. An easy calculation shows that it is equal to

$$\Delta^* = E(\eta - \lambda^*(\xi))^2 = V\eta - \frac{\text{cov}^2(\xi, \eta)}{V\xi} = V\eta[1 - \rho^2(\xi, \eta)]. \quad (17)$$

Consequently, the larger (in absolute value) the correlation coefficient  $\rho(\xi, \eta)$  between  $\xi$  and  $\eta$ , the smaller the mean-square error of observation  $\Delta^*$ . In particular, if  $|\rho(\xi, \eta)| = 1$  then  $\Delta^* = 0$  (cf. Problem 7). On the other hand, if  $\xi$  and  $\eta$  are uncorrelated ( $\rho(\xi, \eta) = 0$ ), then  $\lambda^*(\xi) = E\eta$ , i.e. in the absence of correlation between  $\xi$  and  $\eta$  the best estimate of  $\eta$  in terms of  $\xi$  is simply  $E\eta$  (cf. Problem 4).

## 9. PROBLEMS

1. Verify the following properties of indicators  $I_A = I_A(\omega)$ :

$$I_\emptyset = 0, \quad I_\Omega = 1, \quad I_A + I_{\bar{A}} = 1,$$

$$I_{AB} = I_A \cdot I_B,$$

$$I_{A \cup B} = I_A + I_B - I_{AB}.$$

The indicator of  $\bigcup_{i=1}^n A_i$  is  $1 - \prod_{i=1}^n (1 - I_{A_i})$ , the indicator of  $\bigcap_{i=1}^n A_i$  is  $\prod_{i=1}^n I_{A_i}$ , and the indicator of  $\sum_{i=1}^n A_i$  is  $\sum_{i=1}^n I_{A_i}$ .

$$I_{A \Delta B} = (I_A - I_B)^2,$$

where  $A \Delta B$  is the *symmetric difference* of  $A$  and  $B$ , i.e. the set  $(A \setminus B) \cup (B \setminus A)$ .



2. Let  $\xi_1, \dots, \xi_n$  be independent random variables and

$$\xi_{\min} = \min(\xi_1, \dots, \xi_n), \quad \xi_{\max} = \max(\xi_1, \dots, \xi_n).$$

Show that

$$P\{\xi_{\min} \geq x\} = \prod_{i=1}^n P\{\xi_i \geq x\},$$

$$P\{\xi_{\max} < x\} = \prod_{i=1}^n P\{\xi_i < x\}.$$

3. Let  $\xi_1, \dots, \xi_n$  be independent Bernoulli random variables such that

$$P\{\xi_i = 0\} = 1 - \lambda_i \Delta,$$

$$P\{\xi_i = 1\} = \lambda_i \Delta,$$

where  $\Delta$  is a small number,  $\Delta > 0$ ,  $\lambda_i > 0$ .

Show that

$$P\{\xi_1 + \dots + \xi_n = 1\} = \left( \sum_{i=1}^n \lambda_i \right) \Delta + O(\Delta^2),$$

$$P\{\xi_1 + \dots + \xi_n > 1\} = O(\Delta^2).$$

4. Show that  $\inf_{-\infty < a < \infty} E(\xi - a)^2$  is attained for  $a = E\xi$  and consequently

$$\inf_{-\infty < a < \infty} E(\xi - a)^2 = V\xi.$$

5. Let  $\xi$  be a random variable with distribution function  $F_\xi(x)$  and let  $m_e$  be a median of  $F_\xi(x)$ , i.e. a point such that

$$F_\xi(m_e -) \leq \frac{1}{2} \leq F_\xi(m_e).$$

Show that

$$\inf_{-\infty < a < \infty} E|\xi - a| = E|\xi - m_e|.$$

6. Let  $P_\xi(x) = P\{\xi = x\}$  and  $F_\xi(x) = P\{\xi \leq x\}$ . Show that

$$P_{a\xi+b}(x) = P_\xi\left(\frac{x-b}{a}\right),$$

$$F_{a\xi+b}(x) = F_\xi\left(\frac{x-b}{a}\right)$$

for  $a > 0$  and  $-\infty < b < \infty$ . If  $y \geq 0$ , then

$$F_{\xi^+}(y) = F_\xi(+\sqrt{y}) - F_\xi(-\sqrt{y}) + P_\xi(-\sqrt{y}).$$

Let  $\xi^+ = \max(\xi, 0)$ . Then

$$F_{\xi^+}(x) = \begin{cases} 0, & x < 0, \\ F_\xi(0), & x = 0, \\ F_\xi(x), & x > 0. \end{cases}$$

7. Let  $\xi$  and  $\eta$  be random variables with  $V\xi > 0$ ,  $V\eta > 0$ , and let  $\rho = \rho(\xi, \eta)$  be their correlation coefficient. Show that  $|\rho| \leq 1$ . If  $|\rho| = 1$ , find constants  $a$  and  $b$  such that  $\eta = a\xi + b$ . Moreover, if  $\rho = 1$ , then

$$\frac{\eta - E\eta}{\sqrt{V\eta}} = \frac{\xi - E\xi}{\sqrt{V\xi}}$$

(and therefore  $a > 0$ ), whereas if  $\rho = -1$ , then

$$\frac{\eta - E\eta}{\sqrt{V\eta}} = -\frac{\xi - E\xi}{\sqrt{V\xi}}$$

(and therefore  $a < 0$ ).

8. Let  $\xi$  and  $\eta$  be random variables with  $E\xi = E\eta = 0$ ,  $V\xi = V\eta = 1$  and correlation coefficient  $\rho = \rho(\xi, \eta)$ . Show that

$$E \max(\xi^2, \eta^2) \leq 1 + \sqrt{1 - \rho^2}.$$

9. Use the equation

$$\left( \text{Indicator of } \bigcup_{i=1}^n A_i \right) = \prod_{i=1}^n (1 - I_{A_i}),$$

to deduce the formula  $P(B_0) = 1 - S_1 + S_2 - \dots \pm S_n$  from Problem 4 of §1.

10. Let  $\xi_1, \dots, \xi_n$  be independent random variables,  $\varphi_1 = \varphi_1(\xi_1, \dots, \xi_k)$  and  $\varphi_2 = \varphi_2(\xi_{k+1}, \dots, \xi_n)$ , functions respectively of  $\xi_1, \dots, \xi_k$  and  $\xi_{k+1}, \dots, \xi_n$ . Show that the random variables  $\varphi_1$  and  $\varphi_2$  are independent.
11. Show that the random variables  $\xi_1, \dots, \xi_n$  are independent if and only if

$$F_{\xi_1, \dots, \xi_n}(x_1, \dots, x_n) = F_{\xi_1}(x_1) \cdots F_{\xi_n}(x_n)$$

for all  $x_1, \dots, x_n$ , where  $F_{\xi_1, \dots, \xi_n}(x_1, \dots, x_n) = P\{\xi_1 \leq x_1, \dots, \xi_n \leq x_n\}$ .

12. Show that the random variable  $\xi$  is independent of itself (i.e.,  $\xi$  and  $\xi$  are independent) if and only if  $\xi = \text{const}$ .
13. Under what hypotheses on  $\xi$  are the random variables  $\xi$  and  $\sin \xi$  independent?
14. Let  $\xi$  and  $\eta$  be independent random variables and  $\eta \neq 0$ . Express the probabilities of the events  $P\{\xi\eta \leq z\}$  and  $P\{\xi/\eta \leq z\}$  in terms of the probabilities  $P_\xi(x)$  and  $P_\eta(y)$ .

## §5. The Bernoulli Scheme. I. The Law of Large Numbers

1. In accordance with the definitions given above, a triple

$$(\Omega, \mathcal{A}, P) \quad \text{with} \quad \Omega = \{\omega: \omega = (a_1, \dots, a_n), a_i = 0, 1\},$$

$$\mathcal{A} = \{A: A \subseteq \Omega\}, \quad p(\omega) = p^{\sum a_i} q^{n - \sum a_i}$$

is called a probabilistic model of  $n$  independent experiments with two outcomes, or a Bernoulli scheme.

In this and the next section we study some limiting properties (in a sense described below) for Bernoulli schemes. These are best expressed in terms of random variables and of the probabilities of events connected with them.

We introduce random variables  $\xi_1, \dots, \xi_n$  by taking  $\xi_i(\omega) = a_i$ ,  $i = 1, \dots, n$ , where  $\omega = (a_1, \dots, a_n)$ . As we saw above, the Bernoulli variables  $\xi_i(\omega)$  are independent and identically distributed:

$$P\{\xi_i = 1\} = p, \quad P\{\xi_i = 0\} = q, \quad i = 1, \dots, n.$$

It is natural to think of  $\xi_i$  as describing the result of an experiment at the  $i$ th stage (or at time  $i$ ).

Let us put  $S_0(\omega) \equiv 0$  and

$$S_k = \xi_1 + \dots + \xi_k, \quad k = 1, \dots, n.$$

As we found above,  $ES_n = np$  and consequently

$$E \frac{S_n}{n} = p. \quad (1)$$

In other words, the mean value of the frequency of "success", i.e.  $S_n/n$ , coincides with the probability  $p$  of success. Hence we are led to ask how much the frequency  $S_n/n$  of success differs from its probability  $p$ .

We first note that we cannot expect that, for a sufficiently small  $\varepsilon > 0$  and for sufficiently large  $n$ , the deviation of  $S_n/n$  from  $p$  is less than  $\varepsilon$  for all  $\omega$ , i.e. that

$$\left| \frac{S_n(\omega)}{n} - p \right| \leq \varepsilon, \quad \omega \in \Omega. \quad (2)$$

In fact, when  $0 < p < 1$ ,

$$\begin{aligned} P\left\{\frac{S_n}{n} = 1\right\} &= P\{\xi_1 = 1, \dots, \xi_n = 1\} = p^n, \\ P\left\{\frac{S_n}{n} = 0\right\} &= P\{\xi_1 = 0, \dots, \xi_n = 0\} = q^n, \end{aligned}$$

whence it follows that (2) is not satisfied for sufficiently small  $\varepsilon > 0$ .

We observe, however, that when  $n$  is large the probabilities of the events  $\{S_n/n = 1\}$  and  $\{S_n/n = 0\}$  are small. It is therefore natural to expect that the total probability of the events for which  $|(S_n(\omega)/n) - p| > \varepsilon$  will also be small when  $n$  is sufficiently large.

We shall accordingly try to estimate the probability of the event  $\{\omega: |(S_n(\omega)/n) - p| > \varepsilon\}$ . For this purpose we need the following inequality, which was discovered by Chebyshev.

**Chebyshev's inequality.** Let  $(\Omega, \mathcal{A}, P)$  be a probability space and  $\xi = \xi(\omega)$  a nonnegative random variable. Then

$$P\{\xi \geq \varepsilon\} \leq E\xi/\varepsilon \quad (3)$$

for all  $\varepsilon > 0$ .

PROOF. We notice that

$$\xi = \xi I(\xi \geq \varepsilon) + \xi I(\xi < \varepsilon) \geq \xi I(\xi \geq \varepsilon) \geq \varepsilon I(\xi \geq \varepsilon),$$

where  $I(A)$  is the indicator of  $A$ .

Then, by the properties of the expectation,

$$E\xi \geq \varepsilon E I(\xi \geq \varepsilon) = \varepsilon P(\xi \geq \varepsilon),$$

which establishes (3).

**Corollary.** If  $\xi$  is any random variable, we have for  $\varepsilon > 0$ ,

$$\begin{aligned} P\{|\xi| \geq \varepsilon\} &\leq E|\xi|/\varepsilon, \\ P\{|\xi| \geq \varepsilon\} &= P\{\xi^2 \geq \varepsilon^2\} \leq E\xi^2/\varepsilon^2, \\ P\{|\xi - E\xi| \geq \varepsilon\} &\leq V\xi/\varepsilon^2. \end{aligned} \quad (4)$$

In the last of these inequalities, take  $\xi = S_n/n$ . Then using (4.14), we obtain

$$P\left\{\left|\frac{S_n}{n} - p\right| \geq \varepsilon\right\} \leq \frac{V(S_n/n)}{\varepsilon^2} = \frac{VS_n}{n^2\varepsilon^2} = \frac{npq}{n^2\varepsilon^2} = \frac{pq}{n\varepsilon^2}.$$

Therefore

$$P\left\{\left|\frac{S_n}{n} - p\right| \geq \varepsilon\right\} \leq \frac{pq}{n\varepsilon^2} \leq \frac{1}{4n\varepsilon^2}, \quad (5)$$

from which we see that for large  $n$  there is rather small probability that the frequency  $S_n/n$  of success deviates from the probability  $p$  by more than  $\varepsilon$ .

For  $n \geq 1$  and  $0 \leq k \leq n$ , write

$$P_n(k) = C_n^k p^k q^{n-k}.$$

Then

$$P\left\{\left|\frac{S_n}{n} - p\right| \geq \varepsilon\right\} = \sum_{\{k: |(k/n) - p| \geq \varepsilon\}} P_n(k),$$

and we have actually shown that

$$\sum_{\{k: |(k/n) - p| \geq \varepsilon\}} P_n(k) \leq \frac{pq}{n\varepsilon^2} \leq \frac{1}{4n\varepsilon^2}, \quad (6)$$

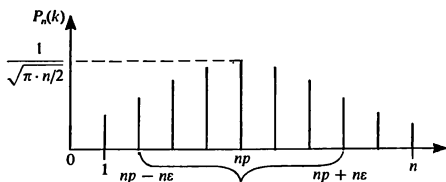


Figure 6

i.e. we have proved an inequality that could also have been obtained analytically, without using the probabilistic interpretation.

It is clear from (6) that

$$\sum_{\{k: |(k/n) - p| \geq \varepsilon\}} P_n(k) \rightarrow 0, \quad n \rightarrow \infty. \quad (7)$$

We can clarify this graphically in the following way. Let us represent the binomial distribution  $\{P_n(k), 0 \leq k \leq n\}$  as in Figure 6.

Then as  $n$  increases the graph spreads out and becomes flatter. At the same time the sum of  $P_n(k)$ , over  $k$  for which  $np - n\varepsilon \leq k < np + n\varepsilon$ , tends to 1.

Let us think of the sequence of random variables  $S_0, S_1, \dots, S_n$  as the path of a wandering particle. Then (7) has the following interpretation.

Let us draw lines from the origin of slopes  $kp$ ,  $k(p + \varepsilon)$ , and  $k(p - \varepsilon)$ . Then on the average the path follows the  $kp$  line, and for every  $\varepsilon > 0$  we can say that when  $n$  is sufficiently large there is a large probability that the point  $S_n$  specifying the position of the particle at time  $n$  lies in the interval  $[n(p - \varepsilon), n(p + \varepsilon)]$ ; see Figure 7.

We would like to write (7) in the following form:

$$P\left\{\left|\frac{S_n}{n} - p\right| \geq \varepsilon\right\} \rightarrow 0, \quad n \rightarrow \infty, \quad (8)$$

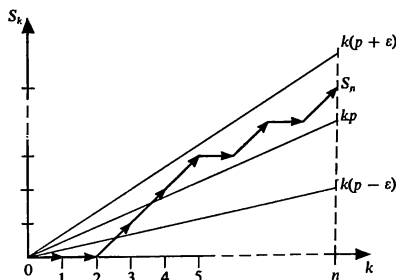


Figure 7

However, we must keep in mind that there is a delicate point involved here. Indeed, the form (8) is really justified only if  $\mathbf{P}$  is a probability on a space  $(\Omega, \mathcal{A})$  on which infinitely many sequences of independent Bernoulli random variables  $\xi_1, \xi_2, \dots$ , are defined. Such spaces can actually be constructed and (8) can be justified in a completely rigorous probabilistic sense (see Corollary 1 below, the end of §4, Chapter II, and Theorem 1, §9, Chapter II). For the time being, if we want to attach a meaning to the analytic statement (7), using the language of probability theory, we have proved only the following.

Let  $(\Omega^{(n)}, \mathcal{A}^{(n)}, \mathbf{P}^{(n)})$ ,  $n \geq 1$ , be a sequence of Bernoulli schemes such that

$$\begin{aligned}\Omega^{(n)} &= \{\omega^{(n)}: \omega^{(n)} = (a_1^{(n)}, \dots, a_n^{(n)}), a_i^{(n)} = 0, 1\}, \\ \mathcal{A}^{(n)} &= \{A: A \subseteq \Omega^{(n)}\}, \\ p^{(n)}(\omega^{(n)}) &= p^{\sum a_i^{(n)}} q^{n - \sum a_i^{(n)}}\end{aligned}$$

and

$$S_k^{(n)}(\omega^{(n)}) = \xi_1^{(n)}(\omega^{(n)}) + \dots + \xi_k^{(n)}(\omega^{(n)}),$$

where, for  $n \leq 1$ ,  $\xi_1^{(n)}, \dots, \xi_n^{(n)}$  are sequences of independent identically distributed Bernoulli random variables.

Then

$$\mathbf{P}^{(n)}\left\{\omega^{(n)}: \left|\frac{S_n^{(n)}(\omega^{(n)})}{n} - p\right| \geq \varepsilon\right\} = \sum_{\{k: |(k/n) - p| \geq \varepsilon\}} P_n(k) \rightarrow 0, \quad n \rightarrow \infty. \quad (9)$$

Statements like (7)–(9) go by the name of **James Bernoulli's law of large numbers**. We may remark that to be precise, Bernoulli's proof consisted in establishing (7), which he did quite rigorously by using estimates for the "tails" of the binomial probabilities  $P_n(k)$  (for the values of  $k$  for which  $|(k/n) - p| \geq \varepsilon$ ). A direct calculation of the sum of the tail probabilities of the binomial distribution  $\sum_{\{k: |(k/n) - p| \geq \varepsilon\}} P_n(k)$  is rather difficult problem for large  $n$ , and the resulting formulas are ill adapted for actual estimates of the probability with which the frequencies  $S_n/n$  differ from  $p$  by less than  $\varepsilon$ . Important progress resulted from the discovery by De Moivre (for  $p = \frac{1}{2}$ ) and then by Laplace (for  $0 < p < 1$ ) of simple asymptotic formulas for  $P_n(k)$ , which led not only to new proofs of the law of large numbers but also to more precise statements of both local and integral limit theorems, the essence of which is that for large  $n$  and at least for  $k \sim np$ ,

$$P_n(k) \sim \frac{1}{\sqrt{2\pi npq}} e^{-(k-np)^2/(2npq)},$$

and

$$\sum_{\{k: |(k/n) - p| \geq \varepsilon\}} P_n(k) \sim \frac{1}{\sqrt{2\pi}} \int_{-\varepsilon\sqrt{n/pq}}^{\varepsilon\sqrt{n/pq}} e^{-x^2/2} dx.$$

2. The next section will be devoted to precise statements and proofs of these results. For the present we consider the question of the real meaning of the law of large numbers, and of its empirical interpretation.

Let us carry out a large number, say  $N$ , of series of experiments, each of which consists of " $n$  independent trials with probability  $p$  of the event  $C$  of interest." Let  $S_n^i/n$  be the frequency of event  $C$  in the  $i$ th series and  $N_\varepsilon$  the number of series in which the frequency deviates from  $p$  by less than  $\varepsilon$ :

$N_\varepsilon$  is the number of  $i$ 's for which  $|(S_n^i/n) - p| \leq \varepsilon$ . Then

$$N_\varepsilon/N \sim P_\varepsilon \quad (10)$$

where  $P_\varepsilon = P\{|(S_n^1/n) - p| \leq \varepsilon\}$ .

It is important to emphasize that an attempt to make (10) precise inevitably involves the introduction of some probability measure, just as an estimate for the deviation of  $S_n/n$  from  $p$  becomes possible only after the introduction of a probability measure  $P$ .

3. Let us consider the estimate obtained above,

$$P\left\{\left|\frac{S_n}{n} - p\right| \geq \varepsilon\right\} = \sum_{\{k: |(k/n) - p| \geq \varepsilon\}} P_n(k) \leq \frac{1}{4n\varepsilon^2}, \quad (11)$$

as an answer to the following question that is typical of mathematical statistics: what is the least number  $n$  of observations that is guaranteed to have (for arbitrary  $0 < p < 1$ )

$$P\left\{\left|\frac{S_n}{n} - p\right| \leq \varepsilon\right\} \geq 1 - \alpha, \quad (12)$$

where  $\alpha$  is a given number (usually small)?

It follows from (11) that this number is the smallest integer  $n$  for which

$$n \geq \frac{1}{4\varepsilon^2\alpha}. \quad (13)$$

For example, if  $\alpha = 0.05$  and  $\varepsilon = 0.02$ , then 12 500 observations guarantee that (12) will hold independently of the value of the unknown parameter  $p$ .

Later (Subsection 5, §6) we shall see that this number is much overstated; this came about because Chebyshev's inequality provides only a very crude upper bound for  $P\{|(S_n/n) - p| \geq \varepsilon\}$ .

4. Let us write

$$C(n, \varepsilon) = \left\{\omega: \left|\frac{S_n(\omega)}{n} - p\right| \leq \varepsilon\right\}.$$

From the law of large numbers that we proved, it follows that for every  $\varepsilon > 0$  and for sufficiently large  $n$ , the probability of the set  $C(n, \varepsilon)$  is close to 1. In this sense it is natural to call paths (realizations)  $\omega$  that are in  $C(n, \varepsilon)$  *typical* (or  $(n, \varepsilon)$ -typical).

We ask the following question: How many typical realizations are there, and what is the weight  $p(\omega)$  of a typical realization?

For this purpose we first notice that the total number  $N(\Omega)$  of points is  $2^n$ , and that if  $p = 0$  or  $1$ , the set of typical paths  $C(n, \varepsilon)$  contains only the single path  $(0, 0, \dots, 0)$  or  $(1, 1, \dots, 1)$ . However, if  $p = \frac{1}{2}$ , it is intuitively clear that "almost all" paths (all except those of the form  $(0, 0, \dots, 0)$  or  $(1, 1, \dots, 1)$ ) are typical and that consequently there should be about  $2^n$  of them.

It turns out that we can give a definitive answer to the question whenever  $0 < p < 1$ ; it will then appear that both the number of typical realizations and the weights  $p(\omega)$  are determined by a function of  $p$  called the entropy.

In order to present the corresponding results in more depth, it will be helpful to consider the somewhat more general scheme of Subsection 2 of §2 instead of the Bernoulli scheme itself.

Let  $(p_1, p_2, \dots, p_r)$  be a finite probability distribution, i.e. a set of nonnegative numbers satisfying  $p_1 + \dots + p_r = 1$ . The *entropy* of this distribution is

$$H = - \sum_{i=1}^r p_i \ln p_i, \quad (14)$$

with  $0 \cdot \ln 0 = 0$ . It is clear that  $H \geq 0$ , and  $H = 0$  if and only if every  $p_i$ , with one exception, is zero. The function  $f(x) = -x \ln x$ ,  $0 \leq x \leq 1$ , is convex upward, so that, as known from the theory of convex functions,

$$\frac{f(x_1) + \dots + f(x_r)}{r} \leq f\left(\frac{x_1 + \dots + x_r}{r}\right).$$

Consequently

$$H = - \sum_{i=1}^r p_i \ln p_i \leq -r \cdot \frac{p_1 + \dots + p_r}{r} \cdot \ln\left(\frac{p_1 + \dots + p_r}{r}\right) = \ln r.$$

In other words, the entropy attains its largest value for  $p_1 = \dots = p_r = 1/r$  (see Figure 8 for  $H = H(p)$  in the case  $r = 2$ ).

If we consider the probability distribution  $(p_1, p_2, \dots, p_r)$  as giving the probabilities for the occurrence of events  $A_1, A_2, \dots, A_r$ , say, then it is quite clear that the "degree of indeterminacy" of an event will be different for

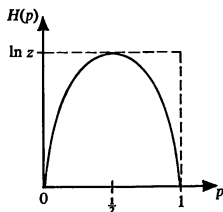


Figure 8. The function  $H(p) = -p \ln p - (1-p) \ln(1-p)$ .



different distributions. If, for example,  $p_1 = 1, p_2 = \dots = p_r = 0$ , it is clear that this distribution does not admit any indeterminacy: we can say with complete certainty that the result of the experiment will be  $A_1$ . On the other hand, if  $p_1 = \dots = p_r = 1/r$ , the distribution has maximal indeterminacy, in the sense that it is impossible to discover any preference for the occurrence of one event rather than another.

Consequently it is important to have a quantitative measure of the indeterminacy of different probability distributions, so that we may compare them in this respect. The entropy successfully provides such a measure of indeterminacy; it plays an important role in statistical mechanics and in many significant problems of coding and communication theory.

Suppose now that the sample space is

$$\Omega = \{\omega: \omega = (a_1, \dots, a_n), a_i = 1, \dots, r\}$$

and that  $p(\omega) = p_1^{v_1(\omega)} \dots p_r^{v_r(\omega)}$ , where  $v_i(\omega)$  is the number of occurrences of  $i$  in the sequence  $\omega$ , and  $(p_1, \dots, p_r)$  is a probability distribution.

For  $\varepsilon > 0$  and  $n = 1, 2, \dots$ , let us put

$$C(n, \varepsilon) = \left\{ \omega: \left| \frac{v_i(\omega)}{n} - p_i \right| < \varepsilon, i = 1, \dots, r \right\}.$$

It is clear that

$$P(C(n, \varepsilon)) \geq 1 - \sum_{i=1}^r P \left\{ \left| \frac{v_i(\omega)}{n} - p_i \right| \geq \varepsilon \right\},$$

and for sufficiently large  $n$  the probabilities  $P\{|(v_i(\omega)/n) - p_i| \geq \varepsilon\}$  are arbitrarily small when  $n$  is sufficiently large, by the law of large numbers applied to the random variables

$$\xi_k(\omega) = \begin{cases} 1, & a_k = i, \\ 0, & a_k \neq i, \end{cases} \quad k = 1, \dots, n.$$

Hence for large  $n$  the probability of the event  $C(n, \varepsilon)$  is close to 1. Thus, as in the case  $n = 2$ , a path in  $C(n, \varepsilon)$  can be said to be typical.

If all  $p_i > 0$ , then for every  $\omega \in \Omega$

$$p(\omega) = \exp \left\{ -n \sum_{k=1}^r \left( -\frac{v_k(\omega)}{n} \ln p_k \right) \right\}.$$

Consequently if  $\omega$  is a typical path, we have

$$\left| \sum_{k=1}^r \left( -\frac{v_k(\omega)}{n} \ln p_k \right) - H \right| \leq - \sum_{k=1}^r \left| \frac{v_k(\omega)}{n} - p_k \right| \ln p_k \leq -\varepsilon \sum_{k=1}^r \ln p_k.$$

It follows that for typical paths the probability  $p(\omega)$  is close to  $e^{-nH}$  and—since, by the law of large numbers, the typical paths “almost” exhaust  $\Omega$  when  $n$  is large—the number of such paths must be of order  $e^{nH}$ . These considerations lead up to the following proposition.

**Theorem (Macmillan).** Let  $p_i > 0$ ,  $i = 1, \dots, r$  and  $0 < \varepsilon < 1$ . Then there is an  $n_0 = n_0(\varepsilon; p_1, \dots, p_r)$  such that for all  $n > n_0$

- (a)  $e^{n(H-\varepsilon)} \leq N(C(n, \varepsilon_1)) \leq e^{n(H+\varepsilon)}$ ;  
 (b)  $e^{-n(H+\varepsilon)} \leq p(\omega) \leq e^{-n(H-\varepsilon)}$ ,  $\omega \in C(n, \varepsilon_1)$ ;  
 (c)  $P(C(n, \varepsilon_1)) = \sum_{\omega \in C(n, \varepsilon_1)} p(\omega) \rightarrow 1$ ,  $n \rightarrow \infty$ ,

where

$$\varepsilon_1 \text{ is the smaller of } \varepsilon \text{ and } \varepsilon / \left\{ -2 \sum_{k=1}^r \ln p_k \right\}.$$

**PROOF.** Conclusion (c) follows from the law of large numbers. To establish the other conclusions, we notice that if  $\omega \in C(n, \varepsilon)$  then

$$np_k - \varepsilon_1 n < v_k(\omega) < np_k + \varepsilon_1 n, \quad k = 1, \dots, r,$$

and therefore

$$\begin{aligned} p(\omega) &= \exp\{-\sum v_k \ln p_k\} < \exp\{-n \sum p_k \ln p_k - \varepsilon_1 n \sum \ln p_k\} \\ &\leq \exp\{-n(H - \tfrac{1}{2}\varepsilon)\}. \end{aligned}$$

Similarly

$$p(\omega) > \exp\{-n(H + \tfrac{1}{2}\varepsilon)\}.$$

Consequently (b) is now established.

Furthermore, since

$$P(C(n, \varepsilon_1)) \geq N(C(n, \varepsilon_1)) \cdot \min_{\omega \in C(n, \varepsilon_1)} p(\omega),$$

we have

$$N(C(n, \varepsilon_1)) \leq \frac{P(C(n, \varepsilon_1))}{\min_{\omega \in C(n, \varepsilon_1)} p(\omega)} < \frac{1}{e^{-n(H + (1/2)\varepsilon)}} = e^{n(H + (1/2)\varepsilon)}$$

and similarly

$$N(C(n, \varepsilon_1)) \geq \frac{P(C(n, \varepsilon_1))}{\max_{\omega \in C(n, \varepsilon_1)} p(\omega)} > P(C(n, \varepsilon_1)) e^{n(H - (1/2)\varepsilon)}.$$

Since  $P(C(n, \varepsilon_1)) \rightarrow 1$ ,  $n \rightarrow \infty$ , there is an  $n_1$  such that  $P(C(n, \varepsilon_1)) > 1 - \varepsilon$  for  $n > n_1$ , and therefore

$$\begin{aligned} N(C(n, \varepsilon_1)) &\geq (1 - \varepsilon) \exp\{n(H - \tfrac{1}{2}\varepsilon)\} \\ &= \exp\{n(H - \varepsilon) + (\tfrac{1}{2}n\varepsilon + \ln(1 - \varepsilon))\}. \end{aligned}$$

Let  $n_2$  be such that

$$\frac{1}{2}n\varepsilon + \ln(1 - \varepsilon) > 0.$$

for  $n > n_2$ . Then when  $n \geq n_0 = \max(n_1, n_2)$  we have

$$N(C(n, \varepsilon_1)) \geq e^{n(H - \varepsilon)}.$$

This completes the proof of the theorem.

5. The law of large numbers for Bernoulli schemes lets us give a simple and elegant proof of Weierstrass's theorem on the approximation of continuous functions by polynomials.

Let  $f = f(p)$  be a continuous function on the interval  $[0, 1]$ . We introduce the polynomials

$$B_n(p) = \sum_{k=0}^n f\left(\frac{k}{n}\right) C_n^k p^k q^{n-k},$$

which are called Bernstein polynomials after the inventor of this proof of Weierstrass's theorem.

If  $\xi_1, \dots, \xi_n$  is a sequence of independent Bernoulli random variables with  $P\{\xi_i = 1\} = p$ ,  $P\{\xi_i = 0\} = q$  and  $S_n = \xi_1 + \dots + \xi_n$ , then

$$Ef\left(\frac{S_n}{n}\right) = B_n(p).$$

Since the function  $f = f(p)$ , being continuous on  $[0, 1]$ , is uniformly continuous, for every  $\varepsilon > 0$  we can find  $\delta > 0$  such that  $|f(x) - f(y)| \leq \varepsilon$  whenever  $|x - y| \leq \delta$ . It is also clear that the function is bounded:  $|f(x)| \leq M < \infty$ .

Using this and (5), we obtain

$$\begin{aligned} |f(p) - B_n(p)| &= \left| \sum_{k=0}^n \left[ f(p) - f\left(\frac{k}{n}\right) \right] C_n^k p^k q^{n-k} \right| \\ &\leq \sum_{\{k: |(k/n) - p| \leq \delta\}} \left| f(p) - f\left(\frac{k}{n}\right) \right| C_n^k p^k q^{n-k} \\ &\quad + \sum_{\{k: |(k/n) - p| > \delta\}} \left| f(p) - f\left(\frac{k}{n}\right) \right| C_n^k p^k q^{n-k} \\ &\leq \varepsilon + 2M \sum_{\{k: |(k/n) - p| > \delta\}} C_n^k p^k q^{n-k} \leq \varepsilon + \frac{2M}{4n\delta^2} = \varepsilon + \frac{M}{2n\delta^2}. \end{aligned}$$

Hence

$$\lim_{n \rightarrow \infty} \max_{0 \leq p \leq 1} |f(p) - B_n(p)| = 0,$$

which is the conclusion of Weierstrass's theorem.

## 6. PROBLEMS

1. Let  $\xi$  and  $\eta$  be random variables with correlation coefficient  $\rho$ . Establish the following two-dimensional analog of Chebyshev's inequality:

$$P\{|\xi - E\xi| \geq \varepsilon\sqrt{V\xi} \text{ or } |\eta - E\eta| \geq \varepsilon\sqrt{V\eta}\} \leq \frac{1}{\varepsilon^2} (1 + \sqrt{1 - \rho^2}).$$

(Hint: Use the result of Problem 8 of §4.)

2. Let  $f = f(x)$  be a nonnegative even function that is nondecreasing for positive  $x$ . Then for a random variable  $\xi$  with  $|\xi(\omega)| \leq C$ ,

$$\frac{Ef(\xi) - f(E\xi)}{f(C)} \leq P\{|\xi - E\xi| \geq \varepsilon\} \leq \frac{Ef(\xi - E\xi)}{f(\varepsilon)}.$$

In particular, if  $f(x) = x^2$ ,

$$\frac{E\xi^2 - \varepsilon^2}{C^2} \leq P\{|\xi - E\xi| \geq \varepsilon\} \leq \frac{V\xi}{\varepsilon^2}.$$

3. Let  $\xi_1, \dots, \xi_n$  be a sequence of independent random variables with  $V\xi_i \leq C$ . Then

$$P\left\{\left|\frac{\xi_1 + \dots + \xi_n}{n} - \frac{E(\xi_1 + \dots + \xi_n)}{n}\right| \geq \varepsilon\right\} \leq \frac{C}{n\varepsilon^2}. \quad (15)$$

(With the same reservations as in (8), inequality (15) implies the validity of the law of large numbers in more general contexts than Bernoulli schemes.)

4. Let  $\xi_1, \dots, \xi_n$  be independent Bernoulli random variables with  $P\{\xi_i = 1\} = p > 0$ ,  $P\{\xi_i = -1\} = 1 - p$ . Derive the following *inequality of Bernstein*: there is a number  $a > 0$  such that

$$P\left\{\left|\frac{S_n}{n} - (2p - 1)\right| \geq \varepsilon\right\} \leq 2e^{-ae^2 n},$$

where  $S_n = \xi_1 + \dots + \xi_n$  and  $\varepsilon > 0$ .

## §6. The Bernoulli Scheme. II. Limit Theorems (Local, De Moivre–Laplace, Poisson)

1. As in the preceding section, let

$$S_n = \xi_1 + \dots + \xi_n.$$

Then

$$E \frac{S_n}{n} = p, \quad (1)$$

and by (4.14)

$$E\left(\frac{S_n}{n} - p\right)^2 = \frac{pq}{n}. \quad (2)$$

It follows from (1) that  $S_n/n \sim p$ , where the equivalence symbol  $\sim$  has been given a precise meaning in the law of large numbers as the assertion  $P\{|(S_n/n) - p| \geq \varepsilon\} \rightarrow 0$ . It is natural to suppose that, in a similar way, the relation

$$\left| \frac{S_n}{n} - p \right| \sim \sqrt{\frac{pq}{n}} \quad (3)$$

which follows from (2), can also be given a precise probabilistic meaning involving, for example, probabilities of the form

$$P\left\{ \left| \frac{S_n}{n} - p \right| \leq x \sqrt{\frac{pq}{n}} \right\}, \quad x \in R^1,$$

or equivalently

$$P\left\{ \left| \frac{S_n - ES_n}{\sqrt{VS_n}} \right| \leq x \right\}$$

(since  $ES_n = np$  and  $VS_n = npq$ ).

If, as before, we write

$$P_n(k) = C_n^k p^k q^{n-k}, \quad 0 \leq k \leq n,$$

for  $n \geq 1$ , then

$$P\left\{ \left| \frac{S_n - ES_n}{\sqrt{VS_n}} \right| \leq x \right\} = \sum_{\{k: |(k - np)|/\sqrt{npq} \leq x\}} P_n(k). \quad (4)$$

We set the problem of finding convenient asymptotic formulas, as  $n \rightarrow \infty$ , for  $P_n(k)$  and for their sum over the values of  $k$  that satisfy the condition on the right-hand side of (4).

The following result provides an answer not only for these values of  $k$  (that is, for those satisfying  $|k - np| = O(\sqrt{npq})$ ) but also for those satisfying  $|k - np| = o(npq)^{2/3}$ .

**Local Limit Theorem.** Let  $0 < p < 1$ ; then

$$P_n(k) \sim \frac{1}{\sqrt{2\pi npq}} e^{-(k - np)^2/(2npq)}, \quad (5)$$

uniformly for  $k$  such that  $|k - np| = o(npq)^{2/3}$ , i.e. as  $n \rightarrow \infty$

$$\sup_{\{k: |k - np| \leq \varphi(n)\}} \left| \frac{P_n(k)}{\frac{1}{\sqrt{2\pi npq}} e^{-(k - np)^2/(2npq)}} - 1 \right| \rightarrow 0,$$

where  $\varphi(n) = o(npq)^{2/3}$ .

The proof depends on Stirling's formula (2.6)

$$n! = \sqrt{2\pi n} e^{-n} n^n (1 + R(n)),$$

where  $R(n) \rightarrow 0$  as  $n \rightarrow \infty$ .

Then if  $n \rightarrow \infty$ ,  $k \rightarrow \infty$ ,  $n - k \rightarrow \infty$ , we have

$$\begin{aligned} C_n^k &= \frac{n!}{k! (n-k)!} \\ &= \frac{\sqrt{2\pi n} e^{-n} n^n}{\sqrt{2\pi k} \cdot \sqrt{2\pi(n-k)} e^{-k} k^k \cdot e^{-(n-k)} (n-k)^{n-k} (1 + R(k))(1 + R(n-k))} \frac{1 + R(n)}{1 + R(k)(1 + R(n-k))} \\ &= \frac{1}{\sqrt{2\pi n} \frac{k}{n} \left(1 - \frac{k}{n}\right)} \cdot \frac{1 + \varepsilon(n, k, n-k)}{\left(\frac{k}{n}\right)^k \left(1 - \frac{k}{n}\right)^{n-k}}, \end{aligned}$$

where  $\varepsilon = \varepsilon(n, k, n-k)$  is defined in an evident way and  $\varepsilon \rightarrow 0$  as  $n \rightarrow \infty$ ,  $k \rightarrow \infty$ ,  $n-k \rightarrow \infty$ .

Therefore

$$P_n(k) = C_n^k p^k q^{n-k} = \frac{1}{\sqrt{2\pi n} \frac{k}{n} \left(1 - \frac{k}{n}\right)} \frac{p^k (1-p)^{n-k}}{\left(\frac{k}{n}\right)^k \left(1 - \frac{k}{n}\right)^{n-k}} (1 + \varepsilon).$$

Write  $\hat{p} = k/n$ . Then

$$\begin{aligned} P_n(k) &= \frac{1}{\sqrt{2\pi n \hat{p}(1-\hat{p})}} \left(\frac{p}{\hat{p}}\right)^k \left(\frac{1-p}{1-\hat{p}}\right)^{n-k} (1 + \varepsilon) \\ &= \frac{1}{\sqrt{2\pi n \hat{p}(1-\hat{p})}} \exp\left\{k \ln \frac{p}{\hat{p}} + (n-k) \ln \frac{1-p}{1-\hat{p}}\right\} (1 + \varepsilon) \\ &= \frac{1}{\sqrt{2\pi n \hat{p}(1-\hat{p})}} \exp\left\{n \left[\frac{k}{n} \ln \frac{p}{\hat{p}} + \left(1 - \frac{k}{n}\right) \ln \frac{1-p}{1-\hat{p}}\right]\right\} (1 + \varepsilon) \\ &= \frac{1}{\sqrt{2\pi n \hat{p}(1-\hat{p})}} \exp\{-nH(\hat{p})\} (1 + \varepsilon), \end{aligned}$$

where

$$H(x) = x \ln \frac{x}{p} + (1-x) \ln \frac{1-x}{1-p}.$$

We are considering values of  $k$  such that  $|k - np| = o(npq)^{2/3}$ , and consequently  $p - \hat{p} \rightarrow 0$ ,  $n \rightarrow \infty$ .

Since, for  $0 < x < 1$ ,

$$H'(x) = \ln \frac{x}{p} - \ln \frac{1-x}{1-p},$$

$$H''(x) = \frac{1}{x} + \frac{1}{1-x},$$

$$H'''(x) = -\frac{1}{x^2} + \frac{1}{(1-x)^2},$$

if we write  $H(\hat{p})$  in the form  $H(p + (\hat{p} - p))$  and use Taylor's formula, we find that for sufficiently large  $n$

$$\begin{aligned} H(\hat{p}) &= H(p) + H'(p)(\hat{p} - p) + \frac{1}{2}H''(p)(\hat{p} - p)^2 + O(|\hat{p} - p|^3) \\ &= \frac{1}{2}\left(\frac{1}{p} + \frac{1}{q}\right)(\hat{p} - p)^2 + O(|\hat{p} - p|^3). \end{aligned}$$

Consequently

$$P_n(k) = \frac{1}{\sqrt{2\pi n\hat{p}(1-\hat{p})}} \exp\left\{-\frac{n}{2pq}(\hat{p} - p)^2 + nO(|\hat{p} - p|^3)\right\} (1 + \varepsilon).$$

Notice that

$$\frac{n}{2pq}(\hat{p} - p)^2 = \frac{n}{2pq}\left(\frac{k}{n} - p\right)^2 = \frac{(k - np)^2}{2npq}.$$

Therefore

$$P_n(k) = \frac{1}{\sqrt{2\pi npq}} e^{-(k - np)^2/(2npq)} (1 + \varepsilon'(n, k, n - k)),$$

where

$$1 + \varepsilon'(n, k, n - k) = (1 + \varepsilon(n, k, n - k)) \exp\{n O(|p - \hat{p}|^3)\} \sqrt{\frac{p(1-p)}{\hat{p}(1-\hat{p})}}$$

and, as is easily seen,

$$\sup |\varepsilon'(n, k, n - k)| \rightarrow 0, \quad n \rightarrow \infty,$$

if the sup is taken over the values of  $k$  for which

$$|k - np| \leq \varphi(n), \quad \varphi(n) = o(npq)^{2/3}.$$

This completes the proof.

**Corollary.** *The conclusion of the local limit theorem can be put in the following equivalent form: For all  $x \in \mathbf{R}^1$  such that  $x = o(npq)^{1/6}$ , and for  $np + x\sqrt{npq}$  an integer from the set  $\{0, 1, \dots, n\}$ ,*

$$P_n(np + x\sqrt{npq}) \sim \frac{1}{\sqrt{2\pi npq}} e^{-x^2/2}, \quad (7)$$

i.e. as  $n \rightarrow \infty$ ,

$$\sup_{\{x: |x| \leq \psi(n)\}} \left| \frac{P_n(np + x\sqrt{npq})}{\frac{1}{\sqrt{2\pi npq}} e^{-x^2/2}} - 1 \right| \rightarrow 0, \quad (8)$$

where  $\psi(n) = o(npq)^{1/6}$ .

With the reservations made in connection with formula (5.8), we can reformulate these results in probabilistic language in the following way:

$$P\{S_n = k\} \sim \frac{1}{\sqrt{2\pi npq}} e^{-(k-np)^2/(2npq)}, \quad |k - np| = o(npq)^{2/3}, \quad (9)$$

$$P\left\{\frac{S_n - np}{\sqrt{npq}} = x\right\} \sim \frac{1}{\sqrt{2\pi npq}} e^{-x^2/2}, \quad x = o(npq)^{1/6}. \quad (10)$$

(In the last formula  $np + x\sqrt{npq}$  is assumed to have one of the values  $0, 1, \dots, n$ .)

If we put  $t_k = (k - np)/\sqrt{npq}$  and  $\Delta t_k = t_{k+1} - t_k = 1/\sqrt{npq}$ , the preceding formula assumes the form

$$P\left\{\frac{S_n - np}{\sqrt{npq}} = t_k\right\} \sim \frac{\Delta t_k}{\sqrt{2\pi}} e^{-t_k^2/2}, \quad t_k = o(npq)^{1/6}. \quad (11)$$

It is clear that  $\Delta t_k = 1/\sqrt{npq} \rightarrow 0$  and the set of points  $\{t_k\}$  as it were "fills" the real line. It is natural to expect that (11) can be used to obtain the integral formula

$$P\left\{a < \frac{S_n - np}{\sqrt{npq}} \leq b\right\} \sim \frac{1}{\sqrt{2\pi}} \int_a^b e^{-x^2/2} dx, \quad -\infty < a \leq b < \infty.$$

Let us now give a precise statement.

2. For  $-\infty < a \leq b < \infty$  let

$$P_n(a, b] = \sum_{a < x \leq b} P_n(np + x\sqrt{npq}),$$

where the summation is over those  $x$  for which  $np + x\sqrt{npq}$  is an integer.



It follows from the local theorem (see also (11)) that for all  $t_k$  defined by  $k = np + t_k\sqrt{npq}$  and satisfying  $|t_k| \leq T < \infty$ ,

$$P_n(np + t_k\sqrt{npq}) = \frac{\Delta t_k}{\sqrt{2\pi}} e^{-t_k^2/2} [1 + \varepsilon(t_k, n)], \quad (12)$$

where

$$\sup_{|t_k| \leq T} |\varepsilon(t_k, n)| \rightarrow 0, \quad n \rightarrow \infty. \quad (13)$$

Consequently, if  $a$  and  $b$  are given so that  $-T \leq a \leq b \leq T$ , then

$$\begin{aligned} \sum_{a < t_k \leq b} P_n(np + t_k\sqrt{npq}) &= \sum_{a < t_k \leq b} \frac{\Delta t_k}{\sqrt{2\pi}} e^{-t_k^2/2} + \sum_{a < t_k \leq b} \varepsilon(t_k, n) \frac{\Delta t_k}{\sqrt{2\pi}} e^{-t_k^2/2} \\ &= \frac{1}{\sqrt{2\pi}} \int_a^b e^{-x^2/2} dx + R_n^{(1)}(a, b) + R_n^{(2)}(a, b), \end{aligned} \quad (14)$$

where

$$\begin{aligned} R_n^{(1)}(a, b) &= \sum_{a < t_k \leq b} \frac{\Delta t_k}{\sqrt{2\pi}} e^{-t_k^2/2} - \frac{1}{\sqrt{2\pi}} \int_a^b e^{-x^2/2} dx, \\ R_n^{(2)}(a, b) &= \sum_{a < t_k \leq b} \varepsilon(t_k, n) \frac{\Delta t_k}{\sqrt{2\pi}} e^{-t_k^2/2}. \end{aligned}$$

From the standard properties of Riemann sums,

$$\sup_{-T \leq a \leq b \leq T} |R_n^{(1)}(a, b)| \rightarrow 0, \quad n \rightarrow \infty. \quad (15)$$

It also clear that

$$\begin{aligned} &\sup_{-T \leq a \leq b \leq T} |R_n^{(2)}(a, b)| \\ &\leq \sup_{|t_k| \leq T} |\varepsilon(t_k, n)| \cdot \sum_{|t_k| \leq T} \frac{\Delta t_k}{\sqrt{2\pi}} e^{-t_k^2/2} \\ &\leq \sup_{|t_k| \leq T} |\varepsilon(t_k, n)| \\ &\quad \times \left[ \frac{1}{\sqrt{2\pi}} \int_{-T}^T e^{-x^2/2} dx + \sup_{-T \leq a \leq b \leq T} |R_n^{(1)}(a, b)| \right] \rightarrow 0, \end{aligned} \quad (16)$$

where the convergence of the right-hand side to zero follows from (15) and from

$$\frac{1}{\sqrt{2\pi}} \int_{-T}^T e^{-x^2/2} dx \leq \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} e^{-x^2/2} dx = 1, \quad (17)$$

the value of the last integral being well known.

We write

$$\Phi(x) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^x e^{-t^2/2} dt.$$

Then it follows from (14)–(16) that

$$\sup_{-T \leq a \leq b \leq T} |P_n(a, b] - (\Phi(b) - \Phi(a))| \rightarrow 0, \quad n \rightarrow \infty. \quad (18)$$

We now show that this result holds for  $T = \infty$  as well as for finite  $T$ . By (17), corresponding to a given  $\varepsilon > 0$  we can find a finite  $T = T(\varepsilon)$  such that

$$\frac{1}{\sqrt{2\pi}} \int_{-T}^T e^{-x^2/2} dx > 1 - \frac{1}{4} \varepsilon. \quad (19)$$

According to (18), we can find an  $N$  such that for all  $n > N$  and  $T = T(\varepsilon)$  we have

$$\sup_{-T \leq a \leq b \leq T} |P_n(a, b] - (\Phi(b) - \Phi(a))| < \frac{1}{4} \varepsilon. \quad (20)$$

It follows from this and (19) that

$$P_n(-T, T] > 1 - \frac{1}{2} \varepsilon,$$

and consequently

$$P_n(-\infty, T] + P_n(T, \infty) \leq \frac{1}{2} \varepsilon,$$

where  $P_n(-\infty, T] = \lim_{S \downarrow -\infty} P_n(S, T]$  and  $P_n(T, \infty) = \lim_{S \uparrow \infty} P_n(T, S]$ .

Therefore for  $-\infty \leq a \leq -T < T \leq b \leq \infty$ ,

$$\begin{aligned} & \left| P_n(a, b] - \frac{1}{\sqrt{2\pi}} \int_a^b e^{-x^2/2} dx \right| \\ & \leq \left| P_n(-T, T] - \frac{1}{\sqrt{2\pi}} \int_{-T}^T e^{-x^2/2} dx \right| \\ & + \left| P_n(a, -T] - \frac{1}{\sqrt{2\pi}} \int_a^{-T} e^{-x^2/2} dx \right| + \left| P_n(T, b] - \frac{1}{\sqrt{2\pi}} \int_T^b e^{-x^2/2} dx \right| \\ & \leq \frac{1}{4} \varepsilon + P_n(-\infty, -T] + \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{-T} e^{-x^2/2} dx + P_n(T, \infty) \\ & + \frac{1}{\sqrt{2\pi}} \int_T^{\infty} e^{-x^2/2} dx \leq \frac{1}{4} \varepsilon + \frac{1}{2} \varepsilon + \frac{1}{8} \varepsilon + \frac{1}{8} \varepsilon = \varepsilon. \end{aligned}$$

By using (18) it is now easy to see that  $P_n(a, b]$  tends uniformly to  $\Phi(b) - \Phi(a)$  for  $-\infty \leq a < b \leq \infty$ .

Thus we have proved the following theorem.

**De Moivre–Laplace Integral Theorem.** Let  $0 < p < 1$ ,

$$P_n(k) = C_n^k p^k q^{n-k}, \quad P_n(a, b] = \sum_{a < x \leq b} P_n(np + x\sqrt{npq}),$$

Then

$$\sup_{-\infty \leq a < b \leq \infty} \left| P_n(a, b] - \frac{1}{\sqrt{2\pi}} \int_a^b e^{-x^2/2} dx \right| \rightarrow 0, \quad n \rightarrow \infty. \quad (21)$$

With the same reservations as in (5.8), (21) can be stated in probabilistic language in the following way:

$$\sup_{-\infty \leq a < b \leq \infty} \left| P \left\{ a < \frac{S_n - ES_n}{\sqrt{VS_n}} \leq b \right\} - \frac{1}{\sqrt{2\pi}} \int_a^b e^{-x^2/2} dx \right| \rightarrow 0, \quad n \rightarrow \infty.$$

It follows at once from this formula that

$$P\{A < S_n \leq B\} - \left[ \Phi\left(\frac{B - np}{\sqrt{npq}}\right) - \Phi\left(\frac{A - np}{\sqrt{npq}}\right) \right] \rightarrow 0, \quad (22)$$

as  $n \rightarrow \infty$ , whenever  $-\infty \leq A < B \leq \infty$ .

**EXAMPLE.** A true die is tossed 12 000 times. We ask for the probability  $P$  that the number of 6's lies in the interval (1800, 2100].

The required probability is

$$P = \sum_{1800 < k \leq 2100} C_{12000}^k \left(\frac{1}{6}\right)^k \left(\frac{5}{6}\right)^{12000-k}.$$

An exact calculation of this sum would obviously be rather difficult. However, if we use the integral theorem we find that the probability  $P$  in question is ( $n = 12\,000$ ,  $p = \frac{1}{6}$ ,  $a = 1800$ ,  $b = 2100$ )

$$\begin{aligned} & \Phi\left(\frac{2100 - 2000}{\sqrt{12\,000 \cdot \frac{1}{6} \cdot \frac{5}{6}}}\right) - \Phi\left(\frac{1800 - 2000}{\sqrt{12\,000 \cdot \frac{1}{6} \cdot \frac{5}{6}}}\right) = \Phi(\sqrt{6}) - \Phi(-2\sqrt{6}) \\ & \approx \Phi(2.449) - \Phi(-4.898) \approx 0.992, \end{aligned}$$

where the values of  $\Phi(2.449)$  and  $\Phi(-4.898)$  were taken from tables of  $\Phi(x)$  (this is the normal distribution function; see Subsection 6 below).

3. We have plotted a graph of  $P_n(np + x\sqrt{npq})$  (with  $x$  assumed such that  $np + x\sqrt{npq}$  is an integer) in Figure 9.

Then the local theorem says that when  $x = o(npq)^{1/6}$ , the curve  $(1/\sqrt{2\pi npq})e^{-x^2/2}$  provides a close fit to  $P_n(np + x\sqrt{npq})$ . On the other hand the integral theorem says that  $P_n(a, b] = P\{a\sqrt{npq} < S_n - np \leq b\sqrt{npq}\} = P\{np + a\sqrt{npq} < S_n \leq np + b\sqrt{npq}\}$  is closely approximated by the integral  $(1/\sqrt{2\pi})\int_a^b e^{-x^2/2} dx$ .

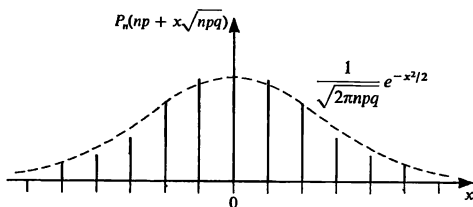


Figure 9

We write

$$F_n(x) = P_n(-\infty, x] \quad \left( = P \left\{ \frac{S_n - np}{\sqrt{npq}} \leq x \right\} \right).$$

Then it follows from (21) that

$$\sup_{-\infty \leq x \leq \infty} |F_n(x) - \Phi(x)| \rightarrow 0, \quad n \rightarrow \infty. \quad (23)$$

It is natural to ask how rapid the approach to zero is in (21) and (23), as  $n \rightarrow \infty$ . We quote a result in this direction (a special case of the Berry-Esseen theorem: see §6 in Chapter III):

$$\sup_{-\infty \leq x \leq \infty} |F_n(x) - \Phi(x)| \leq \frac{p^2 + q^2}{\sqrt{npq}}. \quad (24)$$

It is important to recognize that the order of the estimate  $(1/\sqrt{npq})$  cannot be improved; this means that the approximation of  $F_n(x)$  by  $\Phi(x)$  can be poor for values of  $p$  that are close to 0 or 1, even when  $n$  is large. This suggests the question of whether there is a better method of approximation for the probabilities of interest when  $p$  or  $q$  is small, something better than the normal approximation given by the local and integral theorems. In this connection we note that for  $p = \frac{1}{2}$ , say, the binomial distribution  $\{P_n(k)\}$  is symmetric (Figure 10). However, for small  $p$  the binomial distribution is asymmetric (Figure 10), and hence it is not reasonable to expect that the normal approximation will be satisfactory.

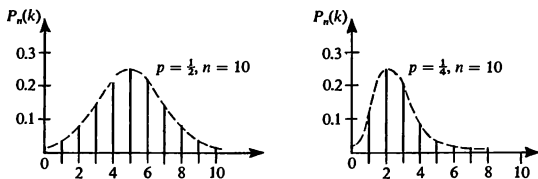


Figure 10

4. It turns out that for small values of  $p$  the distribution known as the Poisson distribution provides a good approximation to  $\{P_n(k)\}$ .

Let

$$P_n(k) = \begin{cases} C_n^k p^k q^{n-k}, & k = 0, 1, \dots, n, \\ 0, & k = n+1, n+2, \dots, \end{cases}$$

and suppose that  $p$  is a function  $p(n)$  of  $n$ .

**Poisson's Theorem.** Let  $p(n) \rightarrow 0$ ,  $n \rightarrow \infty$ , in such a way that  $np(n) \rightarrow \lambda$ , where  $\lambda > 0$ . Then for  $k = 1, 2, \dots$ ,

$$P_n(k) \rightarrow \pi_k, \quad n \rightarrow \infty, \quad (25)$$

where

$$\pi_k = \frac{\lambda^k e^{-\lambda}}{k!}, \quad k = 0, 1, \dots \quad (26)$$

The proof is extremely simple. Since  $p(n) = (\lambda/n) + o(1/n)$  by hypothesis, for a given  $k = 0, 1, \dots$  and sufficiently large  $n$ ,

$$\begin{aligned} P_n(k) &= C_n^k p^k q^{n-k} \\ &= \frac{n(n-1)\cdots(n-k+1)}{k!} \left[ \frac{\lambda}{n} + o\left(\frac{1}{n}\right) \right]^k \cdot \left[ 1 - \frac{\lambda}{n} + o\left(\frac{1}{n}\right) \right]^{n-k}. \end{aligned}$$

But

$$\begin{aligned} &n(n-1)\cdots(n-k+1) \left[ \frac{\lambda}{n} + o\left(\frac{1}{n}\right) \right]^k \\ &= \frac{n(n-1)\cdots(n-k+1)}{n^k} [\lambda + o(1)]^k \rightarrow \lambda^k, \quad n \rightarrow \infty, \end{aligned}$$

and

$$\left[ 1 - \frac{\lambda}{n} + o\left(\frac{1}{n}\right) \right]^{n-k} \rightarrow e^{-\lambda}, \quad n \rightarrow \infty,$$

which establishes (25).

The set of numbers  $\{\pi_k, k = 0, 1, \dots\}$  defines the *Poisson probability distribution* ( $\pi_k \geq 0$ ,  $\sum_{k=0}^{\infty} \pi_k = 1$ ). Notice that all the (discrete) distributions considered previously were concentrated at only a finite number of points. The Poisson distribution is the first example that we have encountered of a (discrete) distribution concentrated at a countable number of points.

The following result of Prokhorov exhibits the rapidity with which  $P_n(k)$  converges to  $\pi_k$  as  $n \rightarrow \infty$ : if  $np(n) = \lambda > 0$ , then

$$\sum_{k=0}^{\infty} |P_n(k) - \pi_k| \leq \frac{2\lambda}{n} \cdot \min(2, \lambda). \quad (27)$$

(A proof of a somewhat weaker result is given in §12, Chapter III.)

5. Let us return to the De Moivre–Laplace limit theorem, and show how it implies the law of large numbers (with the same reservation that was made in connection with (5.8)). Since

$$\mathbf{P}\left\{\left|\frac{S_n}{n} - p\right| \leq \varepsilon\right\} = \mathbf{P}\left\{\left|\frac{S_n - np}{\sqrt{npq}}\right| \leq \varepsilon \sqrt{\frac{n}{pq}}\right\},$$

it is clear from (21) that when  $\varepsilon > 0$

$$\mathbf{P}\left\{\left|\frac{S_n}{n} - p\right| \leq \varepsilon\right\} - \frac{1}{\sqrt{2\pi}} \int_{-\varepsilon\sqrt{n/pq}}^{\varepsilon\sqrt{n/pq}} e^{-x^2/2} dx \rightarrow 0, \quad n \rightarrow \infty, \quad (28)$$

whence

$$\mathbf{P}\left\{\left|\frac{S_n}{n} - p\right| \leq \varepsilon\right\} \rightarrow 1, \quad n \rightarrow \infty,$$

which is the conclusion of the law of large numbers.

From (28)

$$\mathbf{P}\left\{\left|\frac{S_n}{n} - p\right| \leq \varepsilon\right\} \sim \frac{1}{\sqrt{2\pi}} \int_{-\varepsilon\sqrt{n/pq}}^{\varepsilon\sqrt{n/pq}} e^{-x^2/2} dx, \quad n \rightarrow \infty, \quad (29)$$

whereas Chebyshev's inequality yielded only

$$\mathbf{P}\left\{\left|\frac{S_n}{n} - p\right| \leq \varepsilon\right\} \geq 1 - \frac{pq}{n\varepsilon^2}.$$

It was shown at the end of §5 that Chebyshev's inequality yielded the estimate

$$n \geq \frac{1}{4\varepsilon^2\alpha}$$

for the number of observations needed for the validity of the inequality

$$\mathbf{P}\left\{\left|\frac{S_n}{n} - p\right| \leq \varepsilon\right\} \geq 1 - \alpha.$$

Thus with  $\varepsilon = 0.02$  and  $\alpha = 0.05$ , 12 500 observations were needed. We can now solve the same problem by using the approximation (29).

We define the number  $k(\alpha)$  by

$$\frac{1}{\sqrt{2\pi}} \int_{-k(\alpha)}^{k(\alpha)} e^{-x^2/2} dx = 1 - \alpha.$$

Since  $\varepsilon \sqrt{(n/pq)} \geq 2\varepsilon\sqrt{n}$ , if we define  $n$  as the smallest integer satisfying

$$2\varepsilon\sqrt{n} \geq k(\alpha) \quad (30)$$

we find that

$$\mathbf{P}\left\{\left|\frac{S_n}{n} - p\right| \leq \varepsilon\right\} \gtrsim 1 - \alpha. \quad (31)$$

We find from (30) that the smallest integer  $n$  satisfying

$$n \geq \frac{k^2(\alpha)}{4\epsilon^2}$$

guarantees that (31) is satisfied, and the accuracy of the approximation can easily be established by using (24).

Taking  $\epsilon = 0.02$ ,  $\alpha = 0.05$ , we find that in fact 2500 observations suffice, rather than the 12 500 found by using Chebyshev's inequality. The values of  $k(\alpha)$  have been tabulated. We quote a number of values of  $k(\alpha)$  for various values of  $\alpha$ :

$\alpha$	$k(\alpha)$
0.50	0.675
0.3173	1.000
0.10	1.645
0.05	1.960
0.0454	2.000
0.01	2.576
0.0027	3.000

## 6. The function

$$\Phi(x) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^x e^{-t^2/2} dt, \quad (32)$$

which was introduced above and occurs in the De Moivre-Laplace integral theorem, plays an exceptionally important role in probability theory. It is known as the *normal* or *Gaussian distribution* on the real line, with the (normal or Gaussian) density

$$\varphi(x) = \frac{1}{\sqrt{2\pi}} e^{-x^2/2}, \quad x \in \mathbb{R}^1.$$

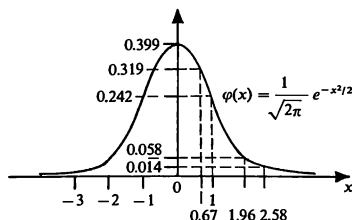
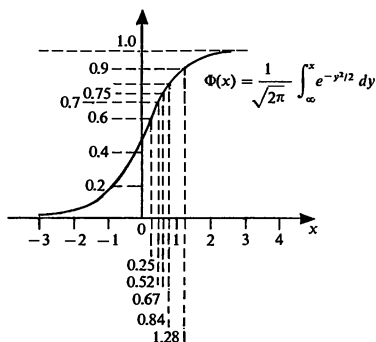


Figure 11. Graph of the normal probability density  $\varphi(x)$ .

Figure 12. Graph of the normal distribution  $\Phi(x)$ .

We have already encountered (discrete) distributions concentrated on a finite or countable set of points. The normal distribution belongs to another important class of distributions that arise in probability theory. We have mentioned its exceptional role; this comes about, first of all, because under rather general hypotheses, sums of a large number of independent random variables (not necessarily Bernoulli variables) are closely approximated by the normal distribution (§4 of Chapter III). For the present we mention only some of the simplest properties of  $\varphi(x)$  and  $\Phi(x)$ , whose graphs are shown in Figures 11 and 12.

The function  $\varphi(x)$  is a symmetric bell-shaped curve, decreasing very rapidly with increasing  $|x|$ : thus  $\varphi(1) = 0.24197$ ,  $\varphi(2) = 0.053991$ ,  $\varphi(3) = 0.004432$ ,  $\varphi(4) = 0.000134$ ,  $\varphi(5) = 0.000016$ . Its maximum is attained at  $x = 0$  and is equal to  $(2\pi)^{-1/2} \approx 0.399$ .

The curve  $\Phi(x) = (1/\sqrt{2\pi}) \int_{-\infty}^x e^{-t^2/2} dt$  approximates 1 very rapidly as  $x$  increases:  $\Phi(1) = 0.841345$ ,  $\Phi(2) = 0.977250$ ,  $\Phi(3) = 0.998650$ ,  $\Phi(4) = 0.999968$ ,  $\Phi(5) = 0.999997$ .

For tables of  $\varphi(x)$  and  $\Phi(x)$ , as well as of other important functions that are used in probability theory and mathematical statistics, see [A1].

7. At the end of subsection 3, §5, we noticed that the upper bound for the probability of the event  $\{\omega: |(S_n/n) - p| \geq \varepsilon\}$ , given by Chebyshev's inequality, was rather crude. That estimate was obtained from Chebyshev's inequality  $P\{X \geq \varepsilon\} \leq EX^2/\varepsilon^2$  for nonnegative random variables  $X \geq 0$ . We may, however, use Chebyshev's inequality in the form

$$P\{X \geq \varepsilon\} = P\{X^{2k} \geq \varepsilon^{2k}\} \leq \frac{EX^{2k}}{\varepsilon^{2k}}. \quad (33)$$

However, we can go further by using the "exponential form" of Chebyshev's



inequality: if  $X \geq 0$  and  $\lambda > 0$ , this states that

$$P\{X \geq \varepsilon\} = P\{e^{\lambda X} \geq e^{\lambda \varepsilon}\} \leq Ee^{\lambda(X-\varepsilon)}. \quad (34)$$

Since the positive number  $\lambda$  is arbitrary, it is clear that

$$P\{X \geq \varepsilon\} \leq \inf_{\lambda > 0} Ee^{\lambda(X-\varepsilon)}. \quad (35)$$

Let us see what the consequences of this approach are in the case when  $X = S_n/n$ ,  $S_n = \xi_1 + \dots + \xi_n$ ,  $P(\xi_i = 1) = p$ ,  $P(\xi_i = 0) = q$ ,  $i \geq 1$ .

Let us set  $\varphi(\lambda) = Ee^{\lambda \xi_1}$ . Then

$$\varphi(\lambda) = 1 - p + pe^{\lambda}$$

and, under the hypothesis of the independence of  $\xi_1, \xi_2, \dots, \xi_n$ ,

$$Ee^{\lambda S_n} = [\varphi(\lambda)]^n.$$

Therefore, ( $0 < a < 1$ )

$$\begin{aligned} P\left\{\frac{S_n}{n} \geq a\right\} &\leq \inf_{\lambda > 0} Ee^{\lambda(S_n/n-a)} = \inf_{\lambda > 0} e^{-n[\lambda a/n - \ln \varphi(\lambda/n)]} \\ &= \inf_{s > 0} e^{-n[as - \ln \varphi(s)]} = e^{-n \sup_{s > 0} [as - \ln \varphi(s)]}. \end{aligned} \quad (36)$$

Similarly,

$$P\left\{\frac{S_n}{n} \leq a\right\} \leq e^{-n \sup_{s < 0} [as - \ln \varphi(s)]}. \quad (37)$$

The function  $f(s) = as - \log[1 - p + pe^s]$  attains its maximum for  $p \leq a \leq 1$  at the point  $s_0$  ( $f'(s_0) = 0$ ) determined by the equation

$$e^{s_0} = \frac{a(1-p)}{p(1-a)}.$$

Consequently,

$$\sup_{s > 0} f(s) = H(a),$$

where

$$H(a) = a \ln \frac{a}{p} + (1-a) \ln \frac{1-a}{1-p}$$

is the function that was previously used in the proof of the local theorem (subsection 1).

Thus, for  $p \leq a \leq 1$

$$P\left\{\frac{S_n}{n} \geq a\right\} \leq e^{-nH(a)}, \quad (38)$$

and therefore, since  $H(p+x) \geq 2x^2$  and  $0 \leq p+x \leq 1$ , we have, for  $\varepsilon > 0$  and  $0 \leq p \leq 1$ ,

$$P \left\{ \frac{S_n}{n} - p \geq \varepsilon \right\} \leq e^{-2n\varepsilon^2}. \quad (39)$$

We can establish similarly that for  $a \leq p \leq 1$

$$P \left\{ \frac{S_n}{n} \leq a \right\} \leq e^{-nH(a)}, \quad (40)$$

and consequently, for every  $\varepsilon > 0$  and  $0 \leq p \leq 1$ ,

$$P \left\{ \frac{S_n}{n} - p \leq -\varepsilon \right\} \leq e^{-2n\varepsilon^2}. \quad (41)$$

Therefore,

$$P \left\{ \left| \frac{S_n}{n} - p \right| \geq \varepsilon \right\} \leq 2e^{-2n\varepsilon^2}. \quad (42)$$

Hence, it follows that the number  $n_3(\alpha)$  of observations of the inequality

$$P \left\{ \left| \frac{S_n}{n} - p \right| \leq \varepsilon \right\} \geq 1 - \alpha, \quad (43)$$

that are guaranteed to be satisfied for every  $p$ ,  $0 < p < 1$ , is determined by the formula

$$n_3(\alpha) = \left\lceil \frac{\ln(2/\alpha)}{2\varepsilon^2} \right\rceil, \quad (44)$$

where  $[x]$  is the integral part of  $x$ . If we neglect "integral parts" and compare  $n_3(\alpha)$  with  $n_1(\alpha) = \lceil (4\alpha\varepsilon^2)^{-1} \rceil$ , we find that

$$\frac{n_1(\alpha)}{n_3(\alpha)} = \frac{1}{2\alpha \ln \frac{2}{\alpha}} \uparrow \infty, \quad \alpha \downarrow 0.$$

It is clear from this that when  $\alpha \downarrow 0$ , an estimate of the smallest number of observations needed that can be obtained from the exponential Chebyshev inequality is more precise than the estimate obtained from the ordinary Chebyshev inequality, especially for small  $\alpha$ .

There is no difficulty in applying the formula

$$\frac{1}{\sqrt{2\pi}} \int_x^\infty e^{-y^2/2} dy \sim \frac{1}{\sqrt{2\pi}x} e^{-x^2/2}, \quad x \rightarrow \infty,$$

to show that  $k^2(\alpha) \sim 2 \log(2/\alpha)$  when  $\alpha \downarrow 0$ . Therefore,

$$\frac{n_2(\alpha)}{n_3(\alpha)} \rightarrow 1, \quad \alpha \downarrow 0.$$

Inequalities like (38)–(42) are known as *inequalities for the probability of large deviations*. This terminology can be explained in the following way.

The De Moivre–Laplace integral theorem makes it possible to estimate in a simple way the probabilities of the events  $\{|S_n - np| \leq x\sqrt{n}\}$  characterizing the “standard” deviation (up to order  $\sqrt{n}$ ) of  $S_n$  from  $np$ . Even the inequalities (39), (41), and (42) provide an estimate of the probabilities of the events  $\{\omega: |S_n - np| \leq x\}$ , describing deviations of order greater than  $\sqrt{n}$ , in fact of order  $n$ .

We shall continue the discussion of probabilities of large deviations, in more general situations, in §5, chap. IV.

## 8. PROBLEMS

1. Let  $n = 100$ ,  $p = \frac{1}{10}, \frac{2}{10}, \frac{3}{10}, \frac{4}{10}, \frac{5}{10}$ . Using tables (for example, those in [A1]) of the binomial and Poisson distributions, compare the values of the probabilities

$$P\{10 < S_{100} \leq 12\}, \quad P\{20 < S_{100} \leq 22\},$$

$$P\{33 < S_{100} \leq 35\}, \quad P\{40 < S_{100} \leq 42\},$$

$$P\{50 < S_{100} \leq 52\}$$

with the corresponding values given by the normal and Poisson approximations.

2. Let  $p = \frac{1}{2}$  and  $Z_n = 2S_n - n$  (the excess of 1's over 0's in  $n$  trials). Show that

$$\sup_j |\sqrt{\pi n} P\{Z_{2n} = j\} - e^{-j^2/4n}| \rightarrow 0, \quad n \rightarrow \infty.$$

3. Show that the rate of convergence in Poisson's theorem is given by

$$\sup_k \left| P_n(k) - \frac{\lambda^k e^{-\lambda}}{k!} \right| \leq \frac{2\lambda^2}{n}.$$

## §7. Estimating the Probability of Success in the Bernoulli Scheme

1. In the Bernoulli scheme  $(\Omega, \mathcal{A}, P)$  with  $\Omega = \{\omega: \omega = (x_1, \dots, x_n), x_i = 0, 1\}$ ,  $\mathcal{A} = A: A \subseteq \Omega$ ,

$$p(\omega) = p^{\sum x_i} q^{n - \sum x_i},$$

we supposed that  $p$  (the probability of success) was known.

Let us now suppose that  $p$  is not known in advance and that we want to determine it by observing the outcomes of experiments; or, what amounts to the same thing, by observations of the random variables  $\xi_1, \dots, \xi_n$ , where  $\xi_i(\omega) = x_i$ . This is a typical problem of mathematical statistics, and can be formulated in various ways. We shall consider two of the possible formulations: the problem of *estimation* and the problem of *constructing confidence intervals*.

In the notation used in mathematical statistics, the unknown parameter is denoted by  $\theta$ , assuming *a priori* that  $\theta$  belongs to the set  $\Theta = [0, 1]$ . We say that the set  $(\Omega, \mathcal{A}, P_\theta; \theta \in \Theta)$  with  $p_\theta(\omega) = \theta^{\sum x_i} (1 - \theta)^{n - \sum x_i}$  is a probabil-

istic-statistical model (corresponding to “ $n$  independent trials” with probability of “success”  $\theta \in \Theta$ ), and any function  $T_n = T_n(\omega)$  with values in  $\Theta$  is called an *estimator*.

If  $S_n = \xi_1 + \dots + \xi_n$  and  $T_n^* = S_n/n$ , it follows from the law of large numbers that  $T_n^*$  is *consistent*, in the sense that ( $\varepsilon > 0$ )

$$P_\theta\{|T_n^* - \theta| \geq \varepsilon\} \rightarrow 0, \quad n \rightarrow \infty. \quad (1)$$

Moreover, this estimator is *unbiased*: for every  $\theta$

$$E_\theta T_n^* = \theta, \quad (2)$$

where  $E_\theta$  is the expectation corresponding to the probability  $P_\theta$ .

The property of being unbiased is quite natural: it expresses the fact that any reasonable estimate ought, at least “on the average,” to lead to the desired result. However, it is easy to see that  $T_n^*$  is not the only unbiased estimator. For example, the same property is possessed by every estimator

$$T_n = \frac{b_1 x_1 + \dots + b_n x_n}{n},$$

where  $b_1 + \dots + b_n = n$ . Moreover, the law of large numbers (1) is also satisfied by such estimators (at least if  $|b_i| \leq K < \infty$ ; see Problem 2, §3, Chapter III) and so these estimators  $T_n$  are just as “good” as  $T_n^*$ .

In this connection there arises the question of how to compare different unbiased estimators, and which of them to describe as best, or optimal.

With the same meaning of “estimator,” it is natural to suppose that an estimator is better, the smaller its deviation from the parameter that is being estimated. On this basis, we call an estimator  $\tilde{T}_n$  *efficient* (in the class of unbiased estimators  $T_n$ ) if,

$$V_\theta \tilde{T}_n = \inf_{T_n} V_\theta T_n, \quad \theta \in \Theta, \quad (3)$$

where  $V_\theta T_n$  is the dispersion of  $T_n$ , i.e.  $E_\theta(T_n - \theta)^2$ .

Let us show that the estimator  $T_n^*$ , considered above, is efficient. We have

$$V_\theta T_n^* = V_\theta \left( \frac{S_n}{n} \right) = \frac{V_\theta S_n}{n^2} = \frac{n\theta(1-\theta)}{n^2} = \frac{\theta(1-\theta)}{n}. \quad (4)$$

Hence to establish that  $T_n^*$  is efficient, we have only to show that

$$\inf_{T_n} V_\theta T_n \geq \frac{\theta(1-\theta)}{n}. \quad (5)$$

This is obvious for  $\theta = 0$  or  $1$ . Let  $\theta \in (0, 1)$  and

$$p_\theta(x_i) = \theta^{x_i}(1-\theta)^{1-x_i}.$$

It is clear that

$$p_\theta(\omega) = \prod_{i=1}^n p_\theta(x_i).$$

Let us write

$$L_{\theta}(\omega) = \ln p_{\theta}(\omega).$$

Then

$$L_{\theta}(\omega) = \ln \theta \cdot \sum x_i + \ln(1 - \theta) \sum (1 - x_i)$$

and

$$\frac{\partial L_{\theta}(\omega)}{\partial \theta} = \frac{\sum (x_i - \theta)}{\theta(1 - \theta)}.$$

Since

$$1 = \mathbf{E}_{\theta} 1 = \sum_{\omega} p_{\theta}(\omega),$$

and since  $T_n$  is unbiased,

$$\theta \equiv \mathbf{E}_{\theta} T_n = \sum_{\omega} T_n(\omega) p_{\theta}(\omega).$$

After differentiating with respect to  $\theta$ , we find that

$$\begin{aligned} 0 &= \sum_{\omega} \frac{\partial p_{\theta}(\omega)}{\partial \theta} = \sum_{\omega} \frac{\left( \frac{\partial p_{\theta}(\omega)}{\partial \theta} \right)}{p_{\theta}(\omega)} p_{\theta}(\omega) = \mathbf{E}_{\theta} \left[ \frac{\partial L_{\theta}(\omega)}{\partial \theta} \right], \\ 1 &= \sum_{\omega} T_n \left( \frac{\partial p_{\theta}(\omega)}{\partial \theta} \right) p_{\theta}(\omega) = \mathbf{E}_{\theta} \left[ T_n \frac{\partial L_{\theta}(\omega)}{\partial \theta} \right]. \end{aligned}$$

Therefore

$$1 = \mathbf{E}_{\theta} \left[ (T_n - \theta) \frac{\partial L_{\theta}(\omega)}{\partial \theta} \right]$$

and by the Cauchy–Bunyakovskii inequality,

$$1 \leq \mathbf{E}_{\theta} [T_n - \theta]^2 \cdot \mathbf{E}_{\theta} \left[ \frac{\partial L_{\theta}(\omega)}{\partial \theta} \right]^2,$$

whence

$$\mathbf{E}_{\theta} [T_n - \theta]^2 \geq \frac{1}{I_n(\theta)}, \quad (6)$$

where

$$I_n(\theta) = \left[ \frac{\partial L_{\theta}(\omega)}{\partial \theta} \right]^2$$

is known as *Fisher's information*.

From (6) we can obtain a special case of the Rao–Cramér inequality for unbiased estimators  $T_n$ :

$$\inf_{T_n} \mathbf{V}_{\theta} T_n \geq \frac{1}{I_n(\theta)}. \quad (7)$$

In the present case

$$I_n(\theta) = E_\theta \left[ \frac{\partial L_\theta(\omega)}{\partial \theta} \right]^2 = E_\theta \left[ \frac{\sum (\xi_i - \theta)}{\theta(1 - \theta)} \right]^2 = \frac{n\theta(1 - \theta)}{[\theta(1 - \theta)]^2} = \frac{n}{\theta(1 - \theta)},$$

which also establishes (5), from which, as we already noticed, there follows the efficiency of the unbiased estimator  $T_n^* = S_n/n$  for the unknown parameter  $\theta$ .

2. It is evident that, in considering  $T_n^*$  as a pointwise estimator for  $\theta$ , we have introduced a certain amount of inaccuracy. It can even happen that the numerical value of  $T_n^*$  calculated from observations of  $x_1, \dots, x_n$  differs rather severely from the true value  $\theta$ . Hence it would be advisable to determine the size of the error.

It would be too much to hope that  $T_n^*(\omega)$  differs little from the true value  $\theta$  for all sample points  $\omega$ . However, we know from the law of large numbers that for every  $\delta > 0$  and for sufficiently large  $n$ , the probability of the event  $\{|\theta - T_n^*(\omega)| > \delta\}$  will be arbitrarily small.

By Chebyshev's inequality

$$P_\theta\{|\theta - T_n^*| > \delta\} \leq \frac{V_\theta T_n^*}{\delta^2} = \frac{\theta(1 - \theta)}{n\delta^2}$$

and therefore, for every  $\lambda > 0$ ,

$$P_\theta\left\{|\theta - T_n^*| \leq \lambda \sqrt{\frac{\theta(1 - \theta)}{n}}\right\} \geq 1 - \frac{1}{\lambda^2}.$$

If we take, for example,  $\lambda = 3$ , then with  $P_\theta$ -probability greater than 0.888 ( $1 - (1/3^2) = \frac{8}{9} \approx 0.8889$ ) the event

$$|\theta - T_n^*| \leq 3 \sqrt{\frac{\theta(1 - \theta)}{n}}$$

will be realized, and a fortiori the event

$$|\theta - T_n^*| \leq \frac{3}{2\sqrt{n}},$$

since  $\theta(1 - \theta) \leq \frac{1}{4}$ .

Therefore

$$P_\theta\left\{|\theta - T_n^*| \leq \frac{3}{2\sqrt{n}}\right\} = P_\theta\left\{T_n^* - \frac{3}{2\sqrt{n}} \leq \theta \leq T_n^* + \frac{3}{2\sqrt{n}}\right\} \geq 0.8888.$$

In other words, we can say with probability greater than 0.8888 that the exact value of  $\theta$  is in the interval  $[T_n^* - (3/2\sqrt{n}), T_n^* + (3/2\sqrt{n})]$ . This statement is sometimes written in the symbolic form

$$\theta \simeq T_n^* \pm \frac{3}{2\sqrt{n}} \quad (\geq 88\%),$$

where " $\geq 88\%$ " means "in more than 88% of all cases."

The interval  $[T_n^* - (3/2\sqrt{n}), T_n^* + (3/2\sqrt{n})]$  is an example of what are called confidence intervals for the unknown parameter.

**Definition.** An interval of the form

$$[\psi_1(\omega), \psi_2(\omega)]$$

where  $\psi_1(\omega)$  and  $\psi_2(\omega)$  are functions of sample points, is called a *confidence interval of reliability*  $1 - \delta$  (or of *significance level*  $\delta$ ) if

$$P_\theta\{\psi_1(\omega) \leq \theta \leq \psi_2(\omega)\} \geq 1 - \delta.$$

for all  $\theta \in \Theta$ .

The preceding discussion shows that the interval

$$\left[ T_n^* - \frac{\lambda}{2\sqrt{n}}, T_n^* + \frac{\lambda}{2\sqrt{n}} \right]$$

has reliability  $1 - (1/\lambda^2)$ . In point of fact, the reliability of this confidence interval is considerably higher, since Chebyshev's inequality gives only crude estimates of the probabilities of events.

To obtain more precise results we notice that

$$\left\{ \omega: |\theta - T_n^*| \leq \lambda \sqrt{\frac{\theta(1-\theta)}{n}} \right\} = \{ \omega: \psi_1(T_n^*, n) \leq \theta \leq \psi_2(T_n^*, n) \},$$

where  $\psi_1 = \psi_1(T_n^*, n)$  and  $\psi_2 = \psi_2(T_n^*, n)$  are the roots of the quadratic equation

$$(\theta - T_n^*)^2 = \frac{\lambda^2}{n} \theta(1 - \theta),$$

which describes an ellipse situated as shown in Figure 13.

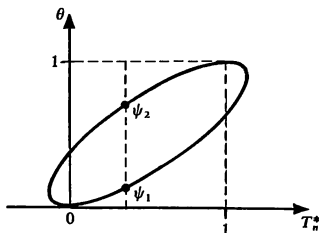


Figure 13

Now let

$$F_{\theta}^n(x) = P_{\theta} \left\{ \frac{S_n - n\theta}{\sqrt{n\theta(1-\theta)}} \leq x \right\}.$$

Then by (6.24)

$$\sup_x |F_{\theta}^n(x) - \Phi(x)| \leq \frac{1}{\sqrt{n\theta(1-\theta)}}.$$

Therefore if we know *a priori* that

$$0 < \Delta \leq \theta \leq 1 - \Delta < 1,$$

where  $\Delta$  is a constant, then

$$\sup_x |F_{\theta}^n(x) - \Phi(x)| \leq \frac{1}{\Delta\sqrt{n}}$$

and consequently

$$\begin{aligned} P_{\theta} \{ \psi_1(T_n^*, n) \leq \theta \leq \psi_2(T_n^*, n) \} &= P_{\theta} \left\{ |\theta - T_n^*| \leq \lambda \sqrt{\frac{\theta(1-\theta)}{n}} \right\} \\ &= P_{\theta} \left\{ \frac{|S_n - n\theta|}{\sqrt{n\theta(1-\theta)}} \leq \lambda \right\} \\ &\geq (2\Phi(\lambda) - 1) - \frac{2}{\Delta\sqrt{n}}. \end{aligned}$$

Let  $\lambda^*$  be the smallest  $\lambda$  for which

$$(2\Phi(\lambda) - 1) - \frac{2}{\Delta\sqrt{n}} \geq 1 - \delta^*,$$

where  $\delta^*$  is a given significance level. Putting  $\delta = \delta^* - (2/\Delta\sqrt{n})$ , we find that  $\lambda^*$  satisfies the equation

$$\Phi(\lambda) = 1 - \frac{1}{2}\delta.$$

For large  $n$  we may neglect the term  $2/\Delta\sqrt{n}$  and assume that  $\lambda^*$  satisfies

$$\Phi(\lambda^*) = 1 - \frac{1}{2}\delta^*.$$

In particular, if  $\lambda^* = 3$  then  $1 - \delta^* = 0.9973 \dots$ . Then with probability approximately 0.9973

$$T_n^* - 3\sqrt{\frac{\theta(1-\theta)}{n}} \leq \theta \leq T_n^* + 3\sqrt{\frac{\theta(1-\theta)}{n}} \quad (8)$$

or, after iterating and then suppressing terms of order  $O(n^{-3/4})$ , we obtain



$$T_n^* - 3\sqrt{\frac{T_n^*(1 - T_n^*)}{n}} \leq \theta \leq T_n^* + 3\sqrt{\frac{T_n^*(1 - T_n^*)}{n}}. \quad (9)$$

Hence it follows that the confidence interval

$$\left[ T_n^* - \frac{3}{2\sqrt{n}}, T_n^* + \frac{3}{2\sqrt{n}} \right] \quad (10)$$

has (for large  $n$ ) reliability 0.9973 (whereas Chebyshev's inequality only provided reliability approximately 0.8889).

Thus we can make the following practical application. Let us carry out a large number  $N$  of series of experiments, in each of which we estimate the parameter  $\theta$  after  $n$  observations. Then in about 99.73% of the  $N$  cases, in each series the estimate will differ from the true value of the parameter by at most  $3/2\sqrt{n}$ . (On this topic see also the end of §5.)

### 3. PROBLEMS

1. Let it be known *a priori* that  $\theta$  has a value in the set  $\Theta_0 \subseteq [0, 1]$ . Construct an unbiased estimator for  $\theta$ , taking values only in  $\Theta_0$ .
2. Under the hypotheses of the preceding problem, find an analog of the Rao-Cramér inequality and discuss the problem of efficient estimators.
3. Under the hypotheses of the first problem, discuss the construction of confidence intervals for  $\theta$ .

## §8. Conditional Probabilities and Mathematical Expectations with Respect to Decompositions

1. Let  $(\Omega, \mathcal{A}, P)$  be a finite probability space and

$$\mathcal{D} = \{D_1, \dots, D_k\}$$

a decomposition of  $\Omega$  ( $D_i \in \mathcal{A}$ ,  $P(D_i) > 0$ ,  $i = 1, \dots, k$ , and  $D_1 + \dots + D_k = \Omega$ ). Also let  $A$  be an event from  $\mathcal{A}$  and  $P(A|D_i)$  the conditional probability of  $A$  with respect to  $D_i$ .

With a set of conditional probabilities  $\{P(A|D_i), i = 1, \dots, k\}$  we may associate the random variable

$$\pi(\omega) = \sum_{i=1}^k P(A|D_i)I_{D_i}(\omega) \quad (1)$$

(cf. (4.5)), that takes the values  $P(A|D_i)$  on the atoms of  $D_i$ . To emphasize that this *random variable* is associated specifically with the decomposition  $\mathcal{D}$ , we denote it by

$$P(A|\mathcal{D}) \quad \text{or} \quad P(A|\mathcal{D})(\omega)$$

and call it the *conditional probability of the event  $A$  with respect to the decomposition  $\mathcal{D}$* .

This concept, as well as the more general concept of conditional probabilities with respect to a  $\sigma$ -algebra, which will be introduced later, plays an important role in probability theory, a role that will be developed progressively as we proceed.

We mention some of the simplest properties of conditional probabilities:

$$P(A + B|\mathcal{D}) = P(A|\mathcal{D}) + P(B|\mathcal{D}); \quad (2)$$

if  $\mathcal{D}$  is the trivial decomposition consisting of the single set  $\Omega$  then

$$P(A|\Omega) = P(A). \quad (3)$$

The definition of  $P(A|\mathcal{D})$  as a random variable lets us speak of its expectation; by using this, we can write the formula (3.3) for total probability in the following compact form:

$$EP(A|\mathcal{D}) = P(A). \quad (4)$$

In fact, since

$$P(A|\mathcal{D}) = \sum_{i=1}^k P(A|D_i)I_{D_i}(\omega),$$

then by the definition of expectation (see (4.5) and (4.6))

$$EP(A|\mathcal{D}) = \sum_{i=1}^k P(A|D_i)P(D_i) = \sum_{i=1}^k P(AD_i) = P(A).$$

Now let  $\eta = \eta(\omega)$  be a random variable that takes the values  $y_1, \dots, y_k$  with positive probabilities:

$$\eta(\omega) = \sum_{j=1}^k y_j I_{D_j}(\omega),$$

where  $D_j = \{\omega: \eta(\omega) = y_j\}$ . The decomposition  $\mathcal{D}_\eta = \{D_1, \dots, D_k\}$  is called the decomposition induced by  $\eta$ . The conditional probability  $P(A|\mathcal{D}_\eta)$  will be denoted by  $P(A|\eta)$  or  $P(A|\eta)(\omega)$ , and called the *conditional probability of  $A$  with respect to the random variable  $\eta$* . We also denote by  $P(A|\eta = y_j)$  the conditional probability  $P(A|D_j)$ , where  $D_j = \{\omega: \eta(\omega) = y_j\}$ .

Similarly, if  $\eta_1, \eta_2, \dots, \eta_m$  are random variables and  $\mathcal{D}_{\eta_1, \eta_2, \dots, \eta_m}$  is the decomposition induced by  $\eta_1, \eta_2, \dots, \eta_m$  with atoms

$$D_{y_1, y_2, \dots, y_m} = \{\omega: \eta_1(\omega) = y_1, \dots, \eta_m(\omega) = y_m\},$$

then  $P(A|D_{\eta_1, \eta_2, \dots, \eta_m})$  will be denoted by  $P(A|\eta_1, \eta_2, \dots, \eta_m)$  and called the conditional probability of  $A$  with respect to  $\eta_1, \eta_2, \dots, \eta_m$ .

**EXAMPLE 1.** Let  $\xi$  and  $\eta$  be independent identically distributed random variables, each taking the values 1 and 0 with probabilities  $p$  and  $q$ . For  $k = 0, 1, 2$ , let us find the conditional probability  $P(\xi + \eta = k|\eta)$  of the event  $A = \{\omega: \xi + \eta = k\}$  with respect to  $\eta$ .

To do this, we first notice the following useful general fact: if  $\xi$  and  $\eta$  are independent random variables with respective values  $x$  and  $y$ , then

$$P(\xi + \eta = z | \eta = y) = P(\xi + y = z). \quad (5)$$

In fact,

$$\begin{aligned} P(\xi + \eta = f | \eta = y) &= \frac{P(\xi + \eta = z, \eta = y)}{P(\eta = y)} \\ &= \frac{P(\xi + y = z, \eta = y)}{P(\eta = y)} = \frac{P(\xi + y = z)P(\eta = y)}{P(\eta = y)} \\ &= P(\xi + y = z). \end{aligned}$$

Using this formula for the case at hand, we find that

$$\begin{aligned} P(\xi + \eta = k | \eta) &= P(\xi + \eta = k | \eta = 0)I_{\{\eta=0\}}(\omega) \\ &\quad + P(\xi + \eta = k | \eta = 1)I_{\{\eta=1\}}(\omega) \\ &= P(\xi = k)I_{\{\eta=0\}}(\omega) + P\{\xi = k - 1\}I_{\{\eta=1\}}(\omega). \end{aligned}$$

Thus

$$P(\xi + \eta = k | \eta) = \begin{cases} qI_{\{\eta=0\}}(\omega), & k = 0, \\ pI_{\{\eta=0\}}(\omega) + qI_{\{\eta=1\}}(\omega), & k = 1, \\ pI_{\{\eta=1\}}(\omega), & k = 2, \end{cases} \quad (6)$$

or equivalently

$$P(\xi + \eta = k | \eta) = \begin{cases} q(1 - \eta), & k = 0, \\ p(1 - \eta) + q\eta, & k = 1, \\ p\eta, & k = 2, \end{cases} \quad (7)$$

2. Let  $\xi = \xi(\omega)$  be a random variable with values in the set  $X = \{x_1, \dots, x_n\}$ :

$$\xi = \sum_{j=1}^l x_j I_{A_j}(\omega), \quad A_j = \{\omega: \xi = x_j\}$$

and let  $\mathcal{D} = \{D_1, \dots, D_k\}$  be a decomposition. Just as we defined the expectation of  $\xi$  with respect to the probabilities  $P(A_j), j = 1, \dots, l$

$$E\xi = \sum_{j=1}^l x_j P(A_j), \quad (8)$$

it is now natural to define the *conditional expectation of  $\xi$  with respect to  $\mathcal{D}$*  by using the conditional probabilities  $P(A_j | \mathcal{D}), j = 1, \dots, l$ . We denote this expectation by  $E(\xi | \mathcal{D})$  or  $E(\xi | \mathcal{D})(\omega)$ , and define it by the formula

$$E(\xi | \mathcal{D}) = \sum_{j=1}^l x_j P(A_j | \mathcal{D}). \quad (9)$$

According to this definition the conditional expectation  $E(\xi | \mathcal{D})(\omega)$  is a random variable which, at all sample points  $\omega$  belonging to the same atom

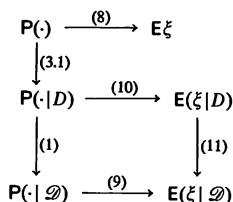


Figure 14

$D_i$ , takes the same value  $\sum_{j=1}^l x_j P(A_j|D_i)$ . This observation shows that the definition of  $E(\xi|\mathcal{D})$  could have been expressed differently. In fact, we could first define  $E(\xi|D_i)$ , the conditional expectation of  $\xi$  with respect to  $D_i$ , by

$$E(\xi|D_i) = \sum_{j=1}^l x_j P(A_j|D_i) \left( = \frac{E[\xi I_{D_i}]}{P(D_i)} \right), \quad (10)$$

and then define

$$E(\xi|\mathcal{D})(\omega) = \sum_{i=1}^k E(\xi|D_i) I_{D_i}(\omega) \quad (11)$$

(see the diagram in Figure 14).

It is also useful to notice that  $E(\xi|D)$  and  $E(\xi|\mathcal{D})$  are independent of the representation of  $\xi$ .

The following properties of conditional expectations follow immediately from the definitions:

$$E(a\xi + b\eta|\mathcal{D}) = aE(\xi|\mathcal{D}) + bE(\eta|\mathcal{D}), \quad a \text{ and } b \text{ constants}; \quad (12)$$

$$E(\xi|\Omega) = E\xi; \quad (13)$$

$$E(C|\mathcal{D}) = C, \quad C \text{ constant}; \quad (14)$$

if  $\xi = I_A(\omega)$  then

$$E(\xi|\mathcal{D}) = P(A|\mathcal{D}). \quad (15)$$

The last equation shows, in particular, that properties of conditional probabilities can be deduced directly from properties of conditional expectations.

The following important property generalizes the *formula for total probability* (5):

$$EE(\xi|\mathcal{D}) = E\xi. \quad (16)$$

For the proof, it is enough to notice that by (5)

$$EE(\xi|\mathcal{D}) = E \sum_{j=1}^l x_j P(A_j|\mathcal{D}) = \sum_{j=1}^l x_j EP(A_j|\mathcal{D}) = \sum_{j=1}^l x_j P(A_j) = E\xi.$$

Let  $\mathcal{D} = \{D_1, \dots, D_k\}$  be a decomposition and  $\eta = \eta(\omega)$  a random variable. We say that  $\eta$  is measurable with respect to this decomposition,

or  $\mathcal{D}$ -measurable, if  $\mathcal{D}_\eta \preceq \mathcal{D}$ , i.e.  $\eta = \eta(\omega)$  can be represented in the form

$$\eta(\omega) = \sum_{i=1}^k y_i I_{D_i}(\omega),$$

where some  $y_i$  might be equal. In other words, a random variable is  $\mathcal{D}$ -measurable if and only if it takes constant values on the atoms of  $\mathcal{D}$ .

EXAMPLE 2. If  $\mathcal{D}$  is the trivial decomposition,  $\mathcal{D} = \{\Omega\}$ , then  $\eta$  is  $\mathcal{D}$ -measurable if and only if  $\eta \equiv C$ , where  $C$  is a constant. Every random variable  $\eta$  is measurable with respect to  $\mathcal{D}_\eta$ .

Suppose that the random variable  $\eta$  is  $\mathcal{D}$ -measurable. Then

$$\mathbf{E}(\xi\eta | \mathcal{D}) = \eta \mathbf{E}(\xi | \mathcal{D}) \quad (17)$$

and in particular

$$\mathbf{E}(\eta | \mathcal{D}) = \eta \quad (\mathbf{E}(\eta | \mathcal{D}_\eta) = \eta). \quad (18)$$

To establish (17) we observe that if  $\xi = \sum_{j=1}^l x_j I_{A_j}$ , then

$$\xi\eta = \sum_{j=1}^l \sum_{i=1}^k x_j y_i I_{A_j D_i}$$

and therefore

$$\begin{aligned} \mathbf{E}(\xi\eta | \mathcal{D}) &= \sum_{j=1}^l \sum_{i=1}^k x_j y_i \mathbf{P}(A_j D_i | \mathcal{D}) \\ &= \sum_{j=1}^l \sum_{i=1}^k x_j y_i \sum_{m=1}^k \mathbf{P}(A_j D_i | D_m) I_{D_m}(\omega) \\ &= \sum_{j=1}^l \sum_{i=1}^k x_j y_i \mathbf{P}(A_j D_i | D_i) I_{D_i}(\omega) \\ &= \sum_{j=1}^l \sum_{i=1}^k x_j y_i \mathbf{P}(A_j | D_i) I_{D_i}(\omega). \end{aligned} \quad (19)$$

On the other hand, since  $I_{D_i}^2 = I_{D_i}$  and  $I_{D_i} \cdot I_{D_m} = 0$ ,  $i \neq m$ , we obtain

$$\begin{aligned} \eta \mathbf{E}(\xi | \mathcal{D}) &= \left[ \sum_{i=1}^k y_i I_{D_i}(\omega) \right] \cdot \left[ \sum_{j=1}^l x_j \mathbf{P}(A_j | \mathcal{D}) \right] \\ &= \left[ \sum_{i=1}^k y_i I_{D_i}(\omega) \right] \cdot \sum_{m=1}^k \left[ \sum_{j=1}^l x_j \mathbf{P}(A_j | D_m) \right] \cdot I_{D_m}(\omega) \\ &= \sum_{i=1}^k \sum_{j=1}^l y_i x_j \mathbf{P}(A_j | D_i) \cdot I_{D_i}(\omega), \end{aligned}$$

which, with (19), establishes (17).

We shall establish another important property of conditional expectations. Let  $\mathcal{D}_1$  and  $\mathcal{D}_2$  be two decompositions, with  $\mathcal{D}_1 \preceq \mathcal{D}_2$  ( $\mathcal{D}_2$  is "finer")

than  $\mathcal{D}_1$ ). Then

$$E[E(\xi | \mathcal{D}_2) | \mathcal{D}_1] = E(\xi | \mathcal{D}_1). \quad (20)$$

For the proof, suppose that

$$\mathcal{D}_1 = \{D_{11}, \dots, D_{1m}\}, \quad \mathcal{D}_2 = \{D_{21}, \dots, D_{2n}\}.$$

Then if  $\xi = \sum_{j=1}^l x_j I_{A_j}$ , we have

$$E(\xi | \mathcal{D}_2) = \sum_{j=1}^l x_j P(A_j | \mathcal{D}_2),$$

and it is sufficient to establish that

$$E[P(A_j | \mathcal{D}_2) | \mathcal{D}_1] = P(A_j | \mathcal{D}_1). \quad (21)$$

Since

$$P(A_j | \mathcal{D}_2) = \sum_{q=1}^n P(A_j | D_{2q}) I_{D_{2q}},$$

we have

$$\begin{aligned} E[P(A_j | \mathcal{D}_2) | \mathcal{D}_1] &= \sum_{q=1}^n P(A_j | D_{2q}) P(D_{2q} | \mathcal{D}_1) \\ &= \sum_{q=1}^n P(A_j | D_{2q}) \left[ \sum_{p=1}^m P(D_{2q} | D_{1p}) I_{D_{1p}} \right] \\ &= \sum_{p=1}^m I_{D_{1p}} \cdot \sum_{q=1}^n P(A_j | D_{2q}) P(D_{2q} | D_{1p}) \\ &= \sum_{p=1}^m I_{D_{1p}} \cdot \sum_{\{q: D_{2q} \subseteq D_{1p}\}} P(A_j | D_{2q}) P(D_{2q} | D_{1p}) \\ &= \sum_{p=1}^m I_{D_{1p}} \cdot \sum_{\{q: D_{2q} \subseteq D_{1p}\}} \frac{P(A_j D_{2q})}{P(D_{2q})} \cdot \frac{P(D_{2q})}{P(D_{1p})} \\ &= \sum_{p=1}^m I_{D_{1p}} \cdot P(A_j | D_{1p}) = P(A_j | \mathcal{D}_1), \end{aligned}$$

which establishes (21).

When  $\mathcal{D}$  is induced by the random variables  $\eta_1, \dots, \eta_k$  ( $\mathcal{D} = \mathcal{D}_{\eta_1, \dots, \eta_k}$ ), the conditional expectation  $E(\xi | \mathcal{D}_{\eta_1, \dots, \eta_k})$  is denoted by  $E(\xi | \eta_1, \dots, \eta_k)$ , or  $E(\xi | \eta_1, \dots, \eta_k)(\omega)$ , and is called the *conditional expectation of  $\xi$  with respect to  $\eta_1, \dots, \eta_k$* .

It follows immediately from the definition of  $E(\xi | \eta)$  that if  $\xi$  and  $\eta$  are *independent*, then

$$E(\xi | \eta) = E\xi. \quad (22)$$

From (18) it also follows that

$$E(\eta | \eta) = \eta. \quad (23)$$

Property (22) admits the following generalization. Let  $\xi$  be independent of  $\mathcal{D}$  (i.e. for each  $D_i \in \mathcal{D}$  the random variables  $\xi$  and  $I_{D_i}$  are independent). Then

$$E(\xi | \mathcal{D}) = E\xi. \quad (24)$$

As a special case of (20) we obtain the following useful formula:

$$E[E(\xi | \eta_1, \eta_2) | \eta_1] = E(\xi | \eta_1). \quad (25)$$

EXAMPLE 3. Let us find  $E(\xi + \eta | \eta)$  for the random variables  $\xi$  and  $\eta$  considered in Example 1. By (22) and (23),

$$E(\xi + \eta | \eta) = E\xi + \eta = p + \eta.$$

This result can also be obtained by starting from (8):

$$E(\xi + \eta | \eta) = \sum_{k=0}^2 kP(\xi + \eta = k | \eta) = p(1 - \eta) + q\eta + 2p\eta = p + \eta.$$

EXAMPLE 4. Let  $\xi$  and  $\eta$  be independent and identically distributed random variables. Then

$$E(\xi | \xi + \eta) = E(\eta | \xi + \eta) = \frac{\xi + \eta}{2}. \quad (26)$$

In fact, if we assume for simplicity that  $\xi$  and  $\eta$  take the values  $1, 2, \dots, m$ , we find ( $1 \leq k \leq m, 2 \leq l \leq 2m$ )

$$\begin{aligned} P(\xi = k | \xi + \eta = l) &= \frac{P(\xi = k, \xi + \eta = l)}{P(\xi + \eta = l)} = \frac{P(\xi = k, \eta = l - k)}{P(\xi + \eta = l)} \\ &= \frac{P(\xi = k)P(\eta = l - k)}{P(\xi + \eta = l)} = \frac{P(\eta = k)P(\xi = l - k)}{P(\xi + \eta = l)} \\ &= P(\eta = k | \xi + \eta = l). \end{aligned}$$

This establishes the first equation in (26). To prove the second, it is enough to notice that

$$2E(\xi | \xi + \eta) = E(\xi | \xi + \eta) + E(\eta | \xi + \eta) = E(\xi + \eta | \xi + \eta) = \xi + \eta.$$

3. We have already noticed in §1 that to each decomposition  $\mathcal{D} = \{D_1, \dots, D_k\}$  of the finite set  $\Omega$  there corresponds an algebra  $\alpha(\mathcal{D})$  of subsets of  $\Omega$ . The converse is also true: every algebra  $\mathcal{B}$  of subsets of the finite space  $\Omega$  generates a decomposition  $\mathcal{D}$  ( $\mathcal{B} = \alpha(\mathcal{D})$ ). Consequently there is a one-to-one correspondence between algebras and decompositions of a finite space  $\Omega$ . This should be kept in mind in connection with the concept, which will be introduced later, of conditional expectation with respect to the special systems of sets called  $\sigma$ -algebras.

For finite spaces, the concepts of algebra and  $\sigma$ -algebra coincide. It will

turn out that if  $\mathcal{B}$  is an algebra, the conditional expectation  $E(\xi|\mathcal{B})$  of a random variable  $\xi$  with respect to  $\mathcal{B}$  (to be introduced in §7 of Chapter II) simply coincides with  $E(\xi|\mathcal{D})$ , the expectation of  $\xi$  with respect to the decomposition  $\mathcal{D}$  such that  $\mathcal{B} = \alpha(\mathcal{D})$ . In this sense we can, in dealing with finite spaces in the future, not distinguish between  $E(\xi|\mathcal{B})$  and  $E(\xi|\mathcal{D})$ , understanding in each case that  $E(\xi|\mathcal{B})$  is simply defined to be  $E(\xi|\mathcal{D})$ .

#### 4. PROBLEMS

1. Give an example of random variables  $\xi$  and  $\eta$  which are not independent but for which

$$E(\xi|\eta) = E\xi.$$

(Cf. (22).)

2. The conditional variance of  $\xi$  with respect to  $\mathcal{D}$  is the random variable

$$V(\xi|\mathcal{D}) = E[(\xi - E(\xi|\mathcal{D}))^2|\mathcal{D}].$$

Show that

$$V\xi = EV(\xi|\mathcal{D}) + VE(\xi|\mathcal{D}).$$

3. Starting from (17), show that for every function  $f = f(\eta)$  the conditional expectation  $E(\xi|\eta)$  has the property

$$E[f(\eta)E(\xi|\eta)] = E[\xi f(\eta)].$$

4. Let  $\xi$  and  $\eta$  be random variables. Show that  $\inf_f E(\eta - f(\xi))^2$  is attained for  $f^*(\xi) = E(\eta|\xi)$ . (Consequently, the best estimator for  $\eta$  in terms of  $\xi$ , in the mean-square sense, is the conditional expectation  $E(\eta|\xi)$ ).
5. Let  $\xi_1, \dots, \xi_n, \tau$  be independent random variables, where  $\xi_1, \dots, \xi_n$  are identically distributed and  $\tau$  takes the values  $1, 2, \dots, n$ . Show that if  $S_\tau = \xi_1 + \dots + \xi_\tau$  is the sum of a random number of the random variables,

$$E(S_\tau|\tau) = \tau E\xi_1, \quad V(S_\tau|\tau) = \tau V\xi_1$$

and

$$ES_\tau = E\tau \cdot E\xi_1, \quad VS_\tau = E\tau \cdot V\xi_1 + V\tau \cdot (E\xi_1)^2.$$

6. Establish equation (24).

## §9. Random Walk. I. Probabilities of Ruin and Mean Duration in Coin Tossing

1. The value of the limit theorems of §6 for Bernoulli schemes is not just that they provide convenient formulas for calculating probabilities  $P(S_n = k)$  and  $P(A < S_n \leq B)$ . They have the additional significance of being of a



universal nature, i.e. they remain useful not only for independent Bernoulli random variables that have only two values, but also for variables of much more general character. In this sense the Bernoulli scheme appears as the simplest model, on the basis of which we can recognize many probabilistic regularities which are inherent also in much more general models.

In this and the next section we shall discuss a number of new probabilistic regularities, some of which are quite surprising. The ones that we discuss are again based on the Bernoulli scheme, although many results on the nature of random oscillations remain valid for random walks of a more general kind.

2. Consider the Bernoulli scheme  $(\Omega, \mathcal{A}, P)$ , where  $\Omega = \{\omega: \omega = (x_1, \dots, x_n), x_i = \pm 1\}$ ,  $\mathcal{A}$  consists of all subsets of  $\Omega$ , and  $p(\omega) = p^{v(\omega)}q^{n-v(\omega)}$ ,  $v(\omega) = (\sum x_i + n)/2$ . Let  $\xi_i(\omega) = x_i$ ,  $i = 1, \dots, n$ . Then, as we know, the sequence  $\xi_1, \dots, \xi_n$  is a sequence of independent Bernoulli random variables,

$$P(\xi_i = 1) = p, \quad P(\xi_i = -1) = q, \quad p + q = 1.$$

Let us put  $S_0 = 0$ ,  $S_k = \xi_1 + \dots + \xi_k$ ,  $1 \leq k \leq n$ . The sequence  $S_0, S_1, \dots, S_n$  can be considered as the path of the random motion of a particle starting at zero. Here  $S_{k+1} = S_k + \xi_{k+1}$ , i.e. if the particle has reached the point  $S_k$  at time  $k$ , then at time  $k+1$  it is displaced either one unit up (with probability  $p$ ) or one unit down (with probability  $q$ ).

Let  $A$  and  $B$  be integers,  $A \leq 0 \leq B$ . An interesting problem about this random walk is to find the probability that after  $n$  steps the moving particle has left the interval  $(A, B)$ . It is also of interest to ask with what probability the particle leaves  $(A, B)$  at  $A$  or at  $B$ .

That these are natural questions to ask becomes particularly clear if we interpret them in terms of a gambling game. Consider two players (first and second) who start with respective bankrolls  $(-A)$  and  $B$ . If  $\xi_i = +1$ , we suppose that the second player pays one unit to the first; if  $\xi_i = -1$ , the first pays the second. Then  $S_k = \xi_1 + \dots + \xi_k$  can be interpreted as the amount won by the first player from the second (if  $S_k < 0$ , this is actually the amount lost by the first player to the second) after  $k$  turns.

At the instant  $k \leq n$  at which for the first time  $S_k = B$  ( $S_k = A$ ) the bankroll of the second (first) player is reduced to zero; in other words, that player is ruined. (If  $k < n$ , we suppose that the game ends at time  $k$ , although the random walk itself is well defined up to time  $n$ , inclusive.)

Before we turn to a precise formulation, let us introduce some notation.

Let  $x$  be an integer in the interval  $[A, B]$  and for  $0 \leq k \leq n$  let  $S_k^x = x + S_k$ ,

$$\tau_k^x = \min\{0 \leq l \leq k: S_l^x = A \text{ or } B\}, \quad (1)$$

where we agree to take  $\tau_k^x = k$  if  $A < S_l^x < B$  for all  $0 \leq l \leq k$ .

For each  $k$  in  $0 \leq k \leq n$  and  $x \in [A, B]$ , the instant  $\tau_k^x$ , called a *stopping time* (see §11), is an integer-valued random variable defined on the sample space  $\Omega$  (the dependence of  $\tau_k^x$  on  $\Omega$  is not explicitly indicated).

It is clear that for all  $l < k$  the set  $\{\omega: \tau_k^x = l\}$  is the event that the random walk  $\{S_i^x, 0 \leq i \leq k\}$ , starting at time zero at the point  $x$ , leaves the interval  $(A, B)$  at time  $l$ . It is also clear that when  $l \leq k$  the sets  $\{\omega: \tau_k^x = l, S_l^x = A\}$  and  $\{\omega: \tau_k^x = l, S_l^x = B\}$  represent the events that the wandering particle leaves the interval  $(A, B)$  at time  $l$  through  $A$  or  $B$  respectively.

For  $0 \leq k \leq n$ , we write

$$\begin{aligned}\mathcal{A}_k^x &= \sum_{0 \leq l \leq k} \{\omega: \tau_k^x = l, S_l^x = A\}, \\ \mathcal{B}_k^x &= \sum_{0 \leq l \leq k} \{\omega: \tau_k^x = l, S_l^x = B\},\end{aligned}\quad (2)$$

and let

$$\alpha_k(x) = P(\mathcal{A}_k^x), \quad \beta_k(x) = P(\mathcal{B}_k^x)$$

be the probabilities that the particle leaves  $(A, B)$ , through  $A$  or  $B$  respectively, during the time interval  $[0, k]$ . For these probabilities we can find recurrent relations from which we can successively determine  $\alpha_1(x), \dots, \alpha_n(x)$  and  $\beta_1(x), \dots, \beta_n(x)$ .

Let, then,  $A < x < B$ . It is clear that  $\alpha_0(x) = \beta_0(x) = 0$ . Now suppose  $1 \leq k \leq n$ . Then by (8.5),

$$\begin{aligned}\beta_k(x) &= P(\mathcal{B}_k^x) = P(\mathcal{B}_k^x | S_1^x = x+1)P(\xi_1 = 1) \\ &\quad + P(\mathcal{B}_k^x | S_1^x = x-1)P(\xi_1 = -1) \\ &= pP(\mathcal{B}_k^x | S_1^x = x+1) + qP(\mathcal{B}_k^x | S_1^x = x-1).\end{aligned}\quad (3)$$

We now show that

$$P(\mathcal{B}_k^x | S_1^x = x+1) = P(\mathcal{B}_{k-1}^{x+1}), \quad P(\mathcal{B}_k^x | S_1^x = x-1) = P(\mathcal{B}_{k-1}^{x-1}).$$

To do this, we notice that  $\mathcal{B}_k^x$  can be represented in the form

$$\mathcal{B}_k^x = \{\omega: (x, x + \xi_1, \dots, x + \xi_1 + \dots + \xi_k) \in B_k^x\},$$

where  $B_k^x$  is the set of paths of the form

$$(x, x + x_1, \dots, x + x_1 + \dots + x_k)$$

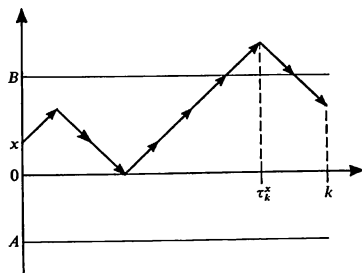
with  $x_1 = \pm 1$ , which during the time  $[0, k]$  first leave  $(A, B)$  at  $B$  (Figure 15).

We represent  $B_k^x$  in the form  $B_k^{x, x+1} + B_k^{x, x-1}$ , where  $B_k^{x, x+1}$  and  $B_k^{x, x-1}$  are the paths in  $B_k^x$  for which  $x_1 = +1$  or  $x_1 = -1$ , respectively.

Notice that the paths  $(x, x+1, x+1+x_2, \dots, x+1+x_2+\dots+x_k)$  in  $B_k^{x, x+1}$  are in one-to-one correspondence with the paths

$$(x+1, x+1+x_2, \dots, x+1+x_2, \dots, x+1+x_2+\dots+x_k)$$

in  $B_{k-1}^{x+1}$ . The same is true for the paths in  $B_k^{x, x-1}$ . Using these facts, together with independence, the identical distribution of  $\xi_1, \dots, \xi_k$ , and (8.6), we obtain

Figure 15. Example of a path from the set  $B_k^x$ .

$$\begin{aligned}
 &P(\mathcal{B}_k^x | S_1^x = x + 1) \\
 &= P(\mathcal{B}_k^x | \xi_1 = 1) \\
 &= P\{(x, x + \xi_1, \dots, x + \xi_1 + \dots + \xi_k) \in B_k^x | \xi_1 = 1\} \\
 &= P\{(x + 1, x + 1 + \xi_2, \dots, x + 1 + \xi_2 + \dots + \xi_k) \in B_{k-1}^{x+1}\} \\
 &= P\{(x + 1, x + 1 + \xi_1, \dots, x + 1 + \xi_1 + \dots + \xi_{k-1}) \in B_{k-1}^{x+1}\} \\
 &= P(\mathcal{B}_{k-1}^{x+1}).
 \end{aligned}$$

In the same way,

$$P(\mathcal{B}_k^x | S_1^x = x - 1) = P(\mathcal{B}_{k-1}^{x-1}).$$

Consequently, by (3) with  $x \in (A, B)$  and  $k \leq n$ ,

$$\beta_k(x) = p\beta_{k-1}(x + 1) + q\beta_{k-1}(x - 1), \quad (4)$$

where

$$\beta_l(B) = 1, \quad \beta_l(A) = 0, \quad 0 \leq l \leq n. \quad (5)$$

Similarly

$$\alpha_k(x) = p\alpha_{k-1}(x + 1) + q\alpha_{k-1}(x - 1) \quad (6)$$

with

$$\alpha_1(A) = 1, \quad \alpha_l(B) = 0, \quad 0 \leq l \leq n.$$

Since  $\alpha_0(x) = \beta_0(x) = 0$ ,  $x \in (A, B)$ , these recurrent relations can (at least in principle) be solved for the probabilities

$$\alpha_1(x), \dots, \alpha_n(x) \quad \text{and} \quad \beta_1(x), \dots, \beta_n(x).$$

Putting aside any explicit calculation of the probabilities, we ask for their values for large  $n$ .

For this purpose we notice that since  $\mathcal{B}_{k-1}^x \subset \mathcal{B}_k^x$ ,  $k \leq n$ , we have  $\beta_{k-1}(x) \leq \beta_k(x) \leq 1$ . It is therefore natural to expect (and this is actually

the case; see Subsection 3) that for sufficiently large  $n$  the probability  $\beta_n(x)$  will be close to the solution  $\beta(x)$  of the equation

$$\beta(x) = p\beta(x+1) + q\beta(x-1) \quad (7)$$

with the boundary conditions

$$\beta(B) = 1, \quad \beta(A) = 0, \quad (8)$$

that result from a formal approach to the limit in (4) and (5).

To solve the problem in (7) and (8), we first suppose that  $p \neq q$ . We see easily that the equation has the two particular solutions  $a$  and  $b(q/p)^x$ , where  $a$  and  $b$  are constants. Hence we look for a solution of the form

$$\beta(x) = a + b(q/p)^x. \quad (9)$$

Taking account of (8), we find that for  $A \leq x \leq B$

$$\beta(x) = \frac{(q/p)^x - (q/p)^A}{(q/p)^B - (q/p)^A}. \quad (10)$$

Let us show that this is the *only* solution of our problem. It is enough to show that all solutions of the problem in (7) and (8) admit the representation (9).

Let  $\tilde{\beta}(x)$  be a solution with  $\tilde{\beta}(A) = 0$ ,  $\tilde{\beta}(B) = 1$ . We can always find constants  $\tilde{a}$  and  $\tilde{b}$  such that

$$\tilde{a} + \tilde{b}(q/p)^A = \tilde{\beta}(A), \quad \tilde{a} + \tilde{b}(q/p)^{A+1} = \tilde{\beta}(A+1).$$

Then it follows from (7) that

$$\tilde{\beta}(A+2) = \tilde{a} + \tilde{b}(q/p)^{A+2}$$

and generally

$$\tilde{\beta}(x) = \tilde{a} + \tilde{b}(q/p)^x.$$

Consequently the solution (10) is the only solution of our problem.

A similar discussion shows that the only solution of

$$\alpha(x) = p\alpha(x+1) + q\alpha(x-1), \quad x \in (A, B) \quad (11)$$

with the boundary conditions

$$\alpha(A) = 1, \quad \alpha(B) = 0 \quad (12)$$

is given by the formula

$$\alpha(x) = \frac{(p/q)^B - (q/p)^x}{(p/q)^B - (p/q)^A}, \quad A \leq x \leq B. \quad (13)$$

If  $p = q = \frac{1}{2}$ , the only solutions  $\beta(x)$  and  $\alpha(x)$  of (7), (8) and (11), (12) are respectively

$$\beta(x) = \frac{x - A}{B - A} \quad (14)$$

and

$$\alpha(x) = \frac{B - x}{B - A}. \quad (15)$$

We note that

$$\alpha(x) + \beta(x) = 1 \quad (16)$$

for  $0 \leq p \leq 1$ .

We call  $\alpha(x)$  and  $\beta(x)$  the *probabilities of ruin for the first and second players*, respectively (when the first player's bankroll is  $x - A$ , and the second player's is  $B - x$ ) under the assumption of infinitely many turns, which of course presupposes an infinite sequence of independent Bernoulli random variables  $\xi_1, \xi_2, \dots$ , where  $\xi_i = +1$  is treated as a gain for the first player, and  $\xi_i = -1$  as a loss. The probability space  $(\Omega, \mathcal{A}, P)$  considered at the beginning of this section turns out to be too small to allow such an infinite sequence of independent variables. We shall see later that such a sequence can actually be constructed and that  $\beta(x)$  and  $\alpha(x)$  are in fact the probabilities of ruin in an unbounded number of steps.

We now take up some corollaries of the preceding formulas.

If we take  $A = 0$ ,  $0 \leq x \leq B$ , then the definition of  $\beta(x)$  implies that this is the probability that a particle starting at  $x$  arrives at  $B$  before it reaches 0. It follows from (10) and (14) (Figure 16) that

$$\beta(x) = \begin{cases} x/B, & p = q = \frac{1}{2}, \\ \frac{(q/p)^x - 1}{(q/p)^B - 1}, & p \neq q. \end{cases} \quad (17)$$

Now let  $q > p$ , which means that the game is unfavorable for the first player, whose limiting probability of being ruined, namely  $\alpha = \alpha(0)$ , is given by

$$\alpha = \frac{(q/p)^B - 1}{(q/p)^B - (q/p)^A}.$$

Next suppose that the rules of the game are changed: the original bankrolls of the players are still  $(-A)$  and  $B$ , but the payoff for each player is now  $\frac{1}{2}$ ,

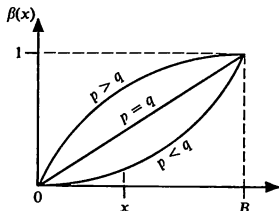


Figure 16. Graph of  $\beta(x)$ , the probability that a particle starting from  $x$  reaches  $B$  before reaching 0.

rather than 1 as before. In other words, now let  $P(\xi_i = \frac{1}{2}) = p$ ,  $P(\xi_i = -\frac{1}{2}) = q$ . In this case let us denote the limiting probability of ruin for the first player by  $\alpha_{1/2}$ . Then

$$\alpha_{1/2} = \frac{(q/p)^{2B} - 1}{(q/p)^{2B} - (q/p)^{2A}},$$

and therefore

$$\alpha_{1/2} = \alpha \cdot \frac{(q/p)^B + 1}{(q/p)^B + (q/p)^A} < \alpha,$$

if  $q > p$ .

Hence we can draw the following conclusion: *if the game is unfavorable to the first player (i.e.,  $q > p$ ) then doubling the stake decreases the probability of ruin.*

3. We now turn to the question of how fast  $\alpha_n(x)$  and  $\beta_n(x)$  approach their limiting values  $\alpha(x)$  and  $\beta(x)$ .

Let us suppose for simplicity that  $x = 0$  and put

$$\alpha_n = \alpha_n(0), \quad \beta_n = \beta_n(0), \quad \gamma_n = 1 - (\alpha_n + \beta_n).$$

It is clear that

$$\gamma_n = P\{A < S_k < B, 0 \leq k \leq n\},$$

where  $\{A < S_k < B, 0 \leq k \leq n\}$  denotes the event

$$\bigcap_{0 \leq k \leq n} \{A < S_k < B\}.$$

Let  $n = rm$ , where  $r$  and  $m$  are integers and

$$\begin{aligned} \zeta_1 &= \xi_1 + \cdots + \xi_m, \\ \zeta_2 &= \xi_{m+1} + \cdots + \xi_{2m}, \\ &\dots\dots\dots \\ \zeta_r &= \xi_{m(r-1)+1} + \cdots + \xi_{rm}. \end{aligned}$$

Then if  $C = |A| + B$ , it is easy to see that

$$\{A < S_k < B, 1 \leq k \leq rm\} \subseteq \{|\zeta_1| < C, \dots, |\zeta_r| < C\},$$

and therefore, since  $\zeta_1, \dots, \zeta_r$  are independent and identically distributed,

$$\gamma_n \leq P\{|\zeta_1| < C, \dots, |\zeta_r| < C\} = \prod_{i=1}^r P\{|\zeta_i| < C\} = (P\{|\zeta_1| < C\})^r. \quad (18)$$

We notice that  $V\zeta_1 = m[1 - (p - q)^2]$ . Hence, for  $0 < p < 1$  and sufficiently large  $m$ ,

$$P\{|\zeta_1| < C\} \leq \varepsilon_1, \quad (19)$$

where  $\varepsilon_1 < 1$ , since  $V\zeta_1 \leq C^2$  if  $P\{|\zeta_1| \leq C\} = 1$ .

If  $p = 0$  or  $p = 1$ , then  $P\{|\zeta_1| < C\} = 0$  for sufficiently large  $m$ , and consequently (19) is satisfied for  $0 \leq p \leq 1$ .

It follows from (18) and (19) that for sufficiently large  $n$

$$\gamma_n \leq \varepsilon^n, \quad (20)$$

where  $\varepsilon = \varepsilon_1^{1/m} < 1$ .

According to (16),  $\alpha + \beta = 1$ . Therefore

$$(\alpha - \alpha_n) + (\beta - \beta_n) = \gamma_n,$$

and since  $\alpha \geq \alpha_n$ ,  $\beta \geq \beta_n$ , we have

$$0 \leq \alpha - \alpha_n \leq \gamma_n \leq \varepsilon^n,$$

$$0 \leq \beta - \beta_n \leq \gamma_n \leq \varepsilon^n, \quad \varepsilon < 1.$$

There are similar inequalities for the differences  $\alpha(x) - \alpha_n(x)$  and  $\beta(x) - \beta_n(x)$ .

4. We now consider the question of the *mean duration* of the random walk.

Let  $m_k(x) = E\tau_k^x$  be the expectation of the stopping time  $\tau_k^x$ ,  $k \leq n$ . Proceeding as in the derivation of the recurrent relations for  $\beta_k(x)$ , we find that, for  $x \in (A, B)$ ,

$$\begin{aligned} m_k(x) &= E\tau_k^x = \sum_{1 \leq l \leq k} l P(\tau_k^x = l) \\ &= \sum_{1 \leq l \leq k} l \cdot [pP(\tau_k^x = l | \xi_1 = 1) + qP(\tau_k^x = l | \xi_1 = -1)] \\ &= \sum_{1 \leq l \leq k} l \cdot [pP(\tau_{k-1}^{x+1} = l-1) + qP(\tau_{k-1}^{x-1} = l-1)] \\ &= \sum_{0 \leq l \leq k-1} (l+1) [pP(\tau_{k-1}^{x+1} = l) + qP(\tau_{k-1}^{x-1} = l)] \\ &= pm_{k-1}(x+1) + qm_{k-1}(x-1) \\ &\quad + \sum_{0 \leq l \leq k-1} [pP(\tau_{k-1}^{x+1} = l) + qP(\tau_{k-1}^{x-1} = l)] \\ &= pm_{k-1}(x+1) + qm_{k-1}(x-1) + 1. \end{aligned}$$

Thus, for  $x \in (A, B)$  and  $0 \leq k \leq n$ , the functions  $m_k(x)$  satisfy the recurrent relations

$$m_k(x) = 1 + pm_{k-1}(x+1) + qm_{k-1}(x-1), \quad (21)$$

with  $m_0(x) = 0$ . From these equations together with the boundary conditions

$$m_k(A) = m_k(B) = 0, \quad (22)$$

we can successively find  $m_1(x), \dots, m_n(x)$ .

Since  $m_k(x) \leq m_{k+1}(x)$ , the limit

$$m(x) = \lim_{n \rightarrow \infty} m_n(x)$$

exists, and by (21) it satisfies the equation

$$m(x) = 1 + pm(x+1) + qm(x-1) \quad (23)$$

with the boundary conditions

$$m(A) = m(B) = 0. \quad (24)$$

To solve this equation, we first suppose that

$$m(x) < \infty, \quad x \in (A, B). \quad (25)$$

Then if  $p \neq q$  there is a particular solution of the form  $x/(q-p)$  and the general solution (see (9)) can be written in the form

$$m(x) = \frac{x}{q-p} + a + b\left(\frac{q}{p}\right)^x.$$

Then by using the boundary conditions  $m(A) = m(B) = 0$  we find that

$$m(x) = \frac{1}{p-q} (B\beta(x) + A\alpha(x) - x), \quad (26)$$

where  $\beta(x)$  and  $\alpha(x)$  are defined by (10) and (13). If  $p = q = \frac{1}{2}$ , the general solution of (23) has the form

$$m(x) = a + bx - x^2,$$

and since  $m(A) = m(B) = 0$  we have

$$m(x) = (B-x)(x-A). \quad (27)$$

It follows, in particular, that if the players start with equal bankrolls ( $B = -A$ ), then

$$m(0) = B^2.$$

If we take  $B = 10$ , and suppose that each turn takes a second, then the (limiting) time to the ruin of one player is rather long: 100 seconds.

We obtained (26) and (27) under the assumption that  $m(x) < \infty$ ,  $x \in (A, B)$ . Let us now show that in fact  $m(x)$  is finite for all  $x \in (A, B)$ . We consider only the case  $x = 0$ ; the general case can be analyzed similarly.

Let  $p = q = \frac{1}{2}$ . We introduce the random variable  $S_{\tau_n}$  defined in terms of the sequence  $S_0, S_1, \dots, S_n$  and the stopping time  $\tau_n = \tau_n^0$  by the equation

$$S_{\tau_n} = \sum_{k=0}^n S_k I_{\{\tau_n = k\}}(\omega). \quad (28)$$

The descriptive meaning of  $S_{\tau_n}$  is clear: it is the position reached by the random walk at the stopping time  $\tau_n$ . Thus, if  $\tau_n < n$ , then  $S_{\tau_n} = A$  or  $B$ ; if  $\tau_n = n$ , then  $A \leq S_{\tau_n} \leq B$ .



Let us show that when  $p = q = \frac{1}{2}$ ,

$$ES_{\tau_n} = 0, \quad (29)$$

$$ES_{\tau_n}^2 = E\tau_n. \quad (30)$$

To establish the first equation we notice that

$$\begin{aligned} ES_{\tau_n} &= \sum_{k=0}^n E[S_k I_{\{\tau_n=k\}}(\omega)] \\ &= \sum_{k=0}^n E[S_n I_{\{\tau_n=k\}}(\omega)] + \sum_{k=0}^n E[(S_k - S_n) I_{\{\tau_n=k\}}(\omega)] \\ &= ES_n + \sum_{k=0}^n E[(S_k - S_n) I_{\{\tau_n=k\}}(\omega)], \end{aligned} \quad (31)$$

where we evidently have  $ES_n = 0$ . Let us show that

$$\sum_{k=0}^n E[(S_k - S_n) I_{\{\tau_n=k\}}(\omega)] = 0.$$

To do this, we notice that  $\{\tau_n > k\} = \{A < S_1 < B, \dots, A < S_k < B\}$  when  $0 \leq k < n$ . The event  $\{A < S_1 < B, \dots, A < S_k < B\}$  can evidently be written in the form

$$\{\omega: (\xi_1, \dots, \xi_k) \in A_k\}, \quad (32)$$

where  $A_k$  is a subset of  $\{-1, +1\}^k$ . In other words, this set is determined by just the values of  $\xi_1, \dots, \xi_k$  and does not depend on  $\xi_{k+1}, \dots, \xi_n$ . Since

$$\{\tau_n = k\} = \{\tau_n > k-1\} \setminus \{\tau_n > k\},$$

this is also a set of the form (32). It then follows from the independence of  $\xi_1, \dots, \xi_n$  and from Problem 10 of §4 that the random variables  $S_n - S_k$  and  $I_{\{\tau_n=k\}}$  are independent, and therefore

$$E[(S_n - S_k) I_{\{\tau_n=k\}}] = E[S_n - S_k] \cdot E I_{\{\tau_n=k\}} = 0.$$

Hence we have established (29).

We can prove (30) by the same method:

$$\begin{aligned} ES_{\tau_n}^2 &= \sum_{k=0}^n ES_k^2 I_{\{\tau_n=k\}} = \sum_{k=0}^n E[(S_n + (S_k - S_n))^2 I_{\{\tau_n=k\}}] \\ &= \sum_{k=0}^n [ES_n^2 I_{\{\tau_n=k\}} + 2ES_n(S_k - S_n) I_{\{\tau_n=k\}} \\ &\quad + E(S_n - S_k)^2 I_{\{\tau_n=k\}}] = ES_n^2 - \sum_{k=0}^n E(S_n - S_k)^2 I_{\{\tau_n=k\}} \\ &= n - \sum_{k=0}^n (n-k)P(\tau_n = k) = \sum_{k=0}^n kP(\tau_n = k) = E\tau_n. \end{aligned}$$

Thus we have (29) and (30) when  $p = q = \frac{1}{2}$ . For general  $p$  and  $q$  ( $p + q = 1$ ) it can be shown similarly that

$$ES_{\tau_n} = (p - q) \cdot E\tau_n, \quad (33)$$

$$E[S_{\tau_n} - \tau_n \cdot E\xi_1]^2 = V\xi_1 \cdot E\tau_n, \quad (34)$$

where  $E\xi_1 = p - q$ ,  $V\xi_1 = 1 - (p - q)^2$ .

With the aid of the results obtained so far we can now show that  $\lim_{n \rightarrow \infty} m_n(0) = m(0) < \infty$ .

If  $p = q = \frac{1}{2}$ , then by (30)

$$E\tau_n \leq \max(A^2, B^2). \quad (35)$$

If  $p \neq q$ , then by (33),

$$E\tau_n \leq \frac{\max(|A|, B)}{|p - q|}, \quad (36)$$

from which it is clear that  $m(0) < \infty$ .

We also notice that when  $p = q = \frac{1}{2}$

$$E\tau_n = ES_{\tau_n}^2 = A^2 \cdot \alpha_n + B^2 \cdot \beta_n + E[S_n^2 I_{\{A < S_n < B\}}]$$

and therefore

$$A^2 \cdot \alpha_n + B^2 \cdot \beta_n \leq E\tau_n \leq A^2 \cdot \alpha_n + B^2 \cdot \beta_n + \max(A^2, B^2) \cdot \gamma_n.$$

It follows from this and (20) that as  $n \rightarrow \infty$ ,  $E\tau_n$  converges with exponential rapidity to

$$m(0) = A^2\alpha + B^2\beta = A^2 \cdot \frac{B}{B - A} - B^2 \cdot \frac{A}{B - A} = |AB|.$$

There is a similar result when  $p \neq q$ :

$$E\tau_n \rightarrow m(0) = \frac{\alpha A + \beta B}{p - q}, \quad \text{exponentially fast.}$$

## 5. PROBLEMS

1. Establish the following generalizations of (33) and (34):

$$ES_{\tau_n}^x = x + (p - q)E\tau_n^x,$$

$$E[S_{\tau_n}^x - \tau_n^x \cdot E\xi_1]^2 = V\xi_1 \cdot E\tau_n^x.$$

2. Investigate the limits of  $\alpha(x)$ ,  $\beta(x)$ , and  $m(x)$  when the level  $A \downarrow -\infty$ .

3. Let  $p = q = \frac{1}{2}$  in the Bernoulli scheme. What is the order of  $E|S_n|$  for large  $n$ ?

4. Two players each toss their own symmetric coins, independently. Show that the probability that each has the same number of heads after  $n$  tosses is  $2^{-2n} \sum_{k=0}^n (C_n^k)^2$ . Hence deduce the equation  $\sum_{k=0}^n (C_n^k)^2 = C_{2n}^n$ .

Let  $\sigma_n$  be the first time when the number of heads for the first player coincides with the number of heads for the second player (if this happens within  $n$  tosses;  $\sigma_n = n + 1$  if there is no such time). Find  $E\sigma_n$ .

## §10. Random Walk. II. Reflection Principle. Arcsine Law

1. As in the preceding section, we suppose that  $\xi_1, \xi_2, \dots, \xi_{2n}$  is a sequence of independent identically distributed Bernoulli random variables with

$$P(\xi_i = 1) = p, \quad P(\xi_i = -1) = q,$$

$$S_k = \xi_1 + \dots + \xi_k, \quad 1 \leq k \leq 2n; \quad S_0 = 0.$$

We define

$$\sigma_{2n} = \min\{1 \leq k \leq 2n: S_k = 0\},$$

putting  $\sigma_{2n} = \infty$  if  $S_k \neq 0$  for  $1 \leq k \leq 2n$ .

The descriptive meaning of  $\sigma_{2n}$  is clear: it is the time of first return to zero. Properties of this time are studied in the present section, where we assume that the random walk is symmetric, i.e.  $p = q = \frac{1}{2}$ .

For  $0 \leq k \leq n$  we write

$$u_{2k} = P(S_{2k} = 0), \quad f_{2k} = P(\sigma_{2n} = 2k). \quad (1)$$

It is clear that  $u_0 = 1$  and

$$u_{2k} = C_{2k}^k \cdot 2^{-2k}.$$

Our immediate aim is to show that for  $1 \leq k \leq n$  the probability  $f_{2k}$  is given by

$$f_{2k} = \frac{1}{2k} u_{2(k-1)}. \quad (2)$$

It is clear that

$$\{\sigma_{2n} = 2k\} = \{S_1 \neq 0, S_2 \neq 0, \dots, S_{2k-1} \neq 0, S_{2k} = 0\}$$

for  $1 \leq k \leq n$ , and by symmetry

$$f_{2k} = P\{S_1 \neq 0, \dots, S_{2k-1} \neq 0, S_{2k} = 0\}$$

$$= 2P\{S_1 > 0, \dots, S_{2k-1} > 0, S_{2k} = 0\}. \quad (3)$$

A sequence  $(S_0, \dots, S_k)$  is called a *path* of length  $k$ ; we denote by  $L_k(A)$  the number of paths of length  $k$  having some specified property  $A$ . Then

$$\begin{aligned} f_{2k} &= 2 \sum_{(a_{2k+1}, \dots, a_n)} L_{2n}(S_1 > 0, \dots, S_{2k-1} > 0, S_{2k} = 0, \\ &\quad \text{and } S_{2k+1} = a_{2k+1}, \dots, S_{2n} = a_{2k+1} + \dots + a_{2n}) \cdot 2^{-2n} \\ &= 2L_{2k}(S_1 > 0, \dots, S_{2k-1} > 0, S_{2k} = 0) \cdot 2^{-2k}, \end{aligned} \quad (4)$$

where the summation is over all sets  $(a_{2k+1}, \dots, a_{2n})$  with  $a_i = \pm 1$ .

Consequently the determination of the probability  $f_{2k}$  reduces to calculating the number of paths  $L_{2k}(S_1 > 0, \dots, S_{2k-1} > 0, S_{2k} = 0)$ .

**Lemma 1.** Let  $a$  and  $b$  be nonnegative integers,  $a - b > 0$  and  $k = a + b$ . Then

$$L_k(S_1 > 0, \dots, S_{k-1} > 0, S_k = a - b) = \frac{a - b}{k} C_k^a. \quad (5)$$

PROOF. In fact,

$$\begin{aligned} L_k(S_1 > 0, \dots, S_{k-1} > 0, S_k = a - b) \\ &= L_k(S_1 = 1, S_2 > 0, \dots, S_{k-1} > 0, S_k = a - b) \\ &= L_k(S_1 = 1, S_k = a - b) - L_k(S_1 = 1, S_k = a - b; \\ &\quad \text{and } \exists i, 2 \leq i \leq k - 1, \text{ such that } S_i \leq 0). \end{aligned} \quad (6)$$

In other words, the number of positive paths  $(S_1, S_2, \dots, S_k)$  that originate at  $(1, 1)$  and terminate at  $(k, a - b)$  is the same as the total number of paths from  $(1, 1)$  to  $(k, a - b)$  after excluding the paths that touch or intersect the time axis.\*

We now notice that

$$\begin{aligned} L_k(S_1 = 1, S_k = a - b; \exists i, 2 \leq i \leq k - 1, \text{ such that } S_i \leq 0) \\ = L_k(S_1 = -1, S_k = a - b), \end{aligned} \quad (7)$$

i.e. the number of paths from  $\alpha = (1, 1)$  to  $\beta = (k, a - b)$ , neither touching nor intersecting the time axis, is equal to the total number of paths that connect  $\alpha^* = (1, -1)$  with  $\beta$ . The proof of this statement, known as the *reflection principle*, follows from the easily established one-to-one correspondence between the paths  $A = (S_1, \dots, S_a, S_{a+1}, \dots, S_k)$  joining  $\alpha$  and  $\beta$ , and paths  $B = (-S_1, \dots, -S_a, S_{a+1}, \dots, S_k)$  joining  $\alpha^*$  and  $\beta$  (Figure 17);  $a$  is the first point where  $A$  and  $B$  reach zero.

\* A path  $(S_1, \dots, S_k)$  is called *positive* (or *nonnegative*) if all  $S_i > 0$  ( $S_i \geq 0$ ); a path is said to *touch* the time axis if  $S_j \geq 0$  or else  $S_j \leq 0$ , for  $1 \leq j \leq k$ , and there is an  $i$ ,  $1 \leq i \leq k$ , such that  $S_i = 0$ ; and a path is said to *intersect* the time axis if there are two times  $i$  and  $j$  such that  $S_i > 0$  and  $S_j < 0$ .

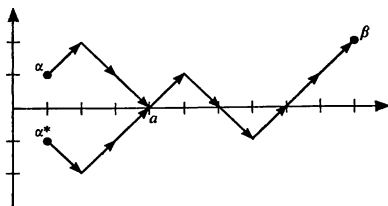


Figure 17. The reflection principle.

From (6) and (7) we find

$$\begin{aligned}
 L_k(S_1 > 0, \dots, S_{k-1} > 0, S_k = a - b) \\
 &= L_k(S_1 = 1, S_k = a - b) - L_k(S_1 = -1, S_k = a - b) \\
 &= C_{k-1}^{a-1} - C_{k-1}^a = \frac{a-b}{k} C_k^a,
 \end{aligned}$$

which establishes (5).

Turning to the calculation of  $f_{2k}$ , we find that by (4) and (5) (with  $a = k$ ,  $b = k - 1$ ),

$$\begin{aligned}
 f_{2k} &= 2L_{2k}(S_1 > 0, \dots, S_{2k-1} > 0, S_{2k} = 0) \cdot 2^{-2k} \\
 &= 2L_{2k-1}(S_1 > 0, \dots, S_{2k-1} = 1) \cdot 2^{-2k} \\
 &= 2 \cdot 2^{-2k} \cdot \frac{1}{2k-1} C_{2k-1}^k = \frac{1}{2k} u_{2(k-1)}.
 \end{aligned}$$

Hence (2) is established.

We present an alternative proof of this formula, based on the following observation. A straightforward verification shows that

$$\frac{1}{2k} u_{2(k-1)} = u_{2(k-1)} - u_{2k}, \quad (8)$$

At the same time, it is clear that

$$\begin{aligned}
 \{\sigma_{2n} = 2k\} &= \{\sigma_{2n} > 2(k-1)\} \setminus \{\sigma_{2n} > 2k\}, \\
 \{\sigma_{2n} > 2l\} &= \{S_1 \neq 0, \dots, S_{2l} \neq 0\}
 \end{aligned}$$

and therefore

$$\{\sigma_{2n} = 2k\} = \{S_1 \neq 0, \dots, S_{2(k-1)} \neq 0\} \setminus \{S_1 \neq 0, \dots, S_{2k} \neq 0\}.$$

Hence

$$f_{2k} = P\{S_1 \neq 0, \dots, S_{2(k-1)} \neq 0\} - P\{S_1 \neq 0, \dots, S_{2k} \neq 0\},$$

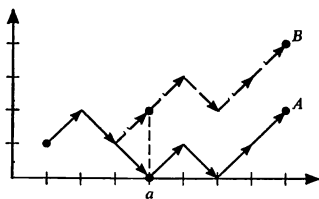


Figure 18

and consequently, because of (8), in order to show that  $f_{2k} = (1/2k)u_{2(k-1)}$  it is enough to show only that

$$L_{2k}(S_1 \neq 0, \dots, S_{2k} \neq 0) = L_{2k}(S_{2k} = 0). \quad (9)$$

For this purpose we notice that evidently

$$L_{2k}(S_1 \neq 0, \dots, S_{2k} \neq 0) = 2L_{2k}(S_1 > 0, \dots, S_{2k} > 0).$$

Hence to verify (9) we need only establish that

$$2L_{2k}(S_1 > 0, \dots, S_{2k} > 0) = L_{2k}(S_1 \geq 0, \dots, S_{2k} \geq 0) \quad (10)$$

and

$$L_{2k}(S_1 \geq 0, \dots, S_{2k} \geq 0) = L_{2k}(S_{2k} = 0). \quad (11)$$

Now (10) will be established if we show that we can establish a one-to-one correspondence between the paths  $A = (S_1, \dots, S_{2k})$  for which at least one  $S_i = 0$ , and the positive paths  $B = (S_1, \dots, S_{2k})$ .

Let  $A = (S_1, \dots, S_{2k})$  be a nonnegative path for which the first zero occurs at the point  $a$  (i.e.,  $S_a = 0$ ). Let us construct the path, starting at  $(a, 2)$ ,  $(S_a + 2, S_{a+1} + 2, \dots, S_{2k} + 2)$  (indicated by the broken lines in Figure 18). Then the path  $B = (S_1, \dots, S_{a-1}, S_a + 2, \dots, S_{2k} + 2)$  is positive.

Conversely, let  $B = (S_1, \dots, S_{2k})$  be a positive path and  $b$  the last instant at which  $S_b = 1$  (Figure 19). Then the path

$$A = (S_1, \dots, S_b, S_{b+1} - 2, \dots, S_{2k} - 2)$$

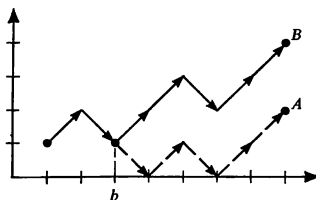


Figure 19

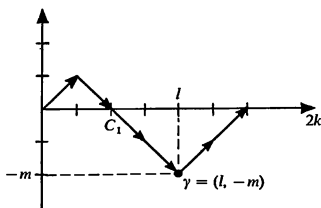


Figure 20

is nonnegative. It follows from these constructions that there is a one-to-one correspondence between the positive paths and the nonnegative paths with at least one  $S_i = 0$ . Therefore formula (10) is established.

We now establish (11). From symmetry and (10) it is enough to show that

$$L_{2k}(S_1 > 0, \dots, S_{2k} > 0) + L_{2k}(S_1 \geq 0, \dots, S_{2k} \geq 0 \text{ and } \exists i, \\ 1 \leq i \leq 2k, \text{ such that } S_i = 0) = L_{2k}(S_{2k} = 0).$$

The set of paths ( $S_{2k} = 0$ ) can be represented as the sum of the two sets  $\mathcal{C}_1$  and  $\mathcal{C}_2$ , where  $\mathcal{C}_1$  contains the paths  $(S_0, \dots, S_{2k})$  that have just one minimum, and  $\mathcal{C}_2$  contains those for which the minimum is attained at at least two points.

Let  $C_1 \in \mathcal{C}_1$  (Figure 20) and let  $\gamma$  be the minimum point. We put the path  $C_1 = (S_0, S_1, \dots, S_{2k})$  in correspondence with the path  $C_1^*$  obtained in the following way (Figure 21). We reflect  $(S_0, S_1, \dots, S_l)$  around the vertical line through the point  $l$ , and displace the resulting path to the right and upward, thus releasing it from the point  $(2k, 0)$ . Then we move the origin to the point  $(l, -m)$ . The resulting path  $C_1^*$  will be positive.

In the same way, if  $C_2 \in \mathcal{C}_2$  we can use the same device to put it into correspondence with a nonnegative path  $C_2^*$ .

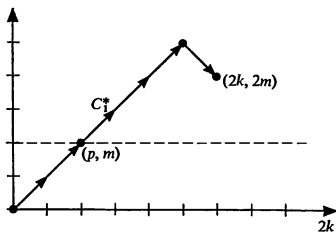


Figure 21

Conversely, let  $C_1^* = (S_1 > 0, \dots, S_{2k} > 0)$  be a positive path with  $S_{2k} = 2m$  (see Figure 21). We make it correspond to the path  $C_1$  that is obtained in the following way. Let  $p$  be the last point at which  $S_p = m$ . Reflect  $(S_p, \dots, S_{2m})$  with respect to the vertical line  $x = p$  and displace the resulting path downward and to the left until its right-hand end coincides with the point  $(0, 0)$ . Then we move the origin to the left-hand end of the resulting path (this is just the path drawn in Figure 20). The resulting path  $C_1 = (S_0, \dots, S_{2k})$  has a minimum at  $S_{2k} = 0$ . A similar construction applied to paths  $(S_1 \geq 0, \dots, S_{2k} \geq 0$  and  $\exists i, 1 \leq i \leq 2k$ , with  $S_i = 0)$  leads to paths for which there are at least two minima and  $S_{2k} = 0$ . Hence we have established a one-to-one correspondence, which establishes (11).

Therefore we have established (9) and consequently also the formula  $f_{2k} = u_{2(k-1)} - u_{2k} = (1/2k)u_{2(k-1)}$ .

By Stirling's formula

$$u_{2k} = C_{2k}^k \cdot 2^{-2k} \sim \frac{1}{\sqrt{\pi k}}, \quad k \rightarrow \infty.$$

Therefore

$$f_{2k} \sim \frac{1}{2\sqrt{\pi k^{3/2}}}, \quad k \rightarrow \infty.$$

Hence it follows that the expectation of the first time when zero is reached, namely

$$\begin{aligned} \text{Emin}(\sigma_{2n}, 2n) &= \sum_{k=1}^n 2kP(\sigma_{2n} = 2k) + 2nu_{2n} \\ &= \sum_{k=1}^n u_{2(k-1)} + 2nu_{2n}, \end{aligned}$$

can be arbitrarily large.

In addition,  $\sum_{k=1}^{\infty} u_{2(k-1)} = \infty$ , and consequently the limiting value of the mean time for the walk to reach zero (in an unbounded number of steps) is  $\infty$ .

This property accounts for many of the unexpected properties of the symmetric random walk that we have been discussing. For example, it would be natural to suppose that after time  $2n$  the number of zero net scores in a game between two equally matched players ( $p = q = \frac{1}{2}$ ), i.e. the number of instants  $i$  at which  $S_i = 0$ , would be proportional to  $2n$ . However, in fact the number of zeros has order  $\sqrt{2n}$  (see [F1] and (15) in §9, Chapter VII). Hence it follows, in particular, that, contrary to intuition, the "typical" walk  $(S_0, S_1, \dots, S_n)$  does not have a sinusoidal character (so that roughly half the time the particle would be on the positive side and half the time on the negative side), but instead must resemble a stretched-out wave. The precise formulation of this statement is given by the arcsine law, which we proceed to investigate.



2. Let  $P_{2k, 2n}$  be the probability that during the interval  $[0, 2n]$  the particle spends  $2k$  units of time on the positive side.\*

**Lemma 2.** Let  $u_0 = 1$  and  $0 \leq k \leq n$ . Then

$$P_{2k, 2n} = u_{2k} \cdot u_{2n-2k} \quad (12)$$

PROOF. It was shown above that  $f_{2k} = u_{2(k-1)} - u_{2k}$ . Let us show that

$$u_{2k} = \sum_{r=1}^k f_{2r} \cdot u_{2(k-r)}. \quad (13)$$

Since  $\{S_{2k} = 0\} \subseteq \{\sigma_{2n} \leq 2k\}$ , we have

$$\{S_{2k} = 0\} = \{S_{2k} = 0\} \cap \{\sigma_{2n} \leq 2k\} = \sum_{1 \leq l \leq k} \{S_{2k} = 0\} \cap \{\sigma_{2n} = 2l\}.$$

Consequently

$$\begin{aligned} u_{2k} &= P(S_{2k} = 0) = \sum_{1 \leq l \leq k} P(S_{2k} = 0, \sigma_{2n} = 2l) \\ &= \sum_{1 \leq l \leq k} P(S_{2k} = 0 | \sigma_{2k} = 2l) P(\sigma_{2n} = 2l). \end{aligned}$$

But

$$\begin{aligned} P(S_{2k} = 0 | \sigma_{2n} = 2l) &= P(S_{2k} = 0 | S_1 \neq 0, \dots, S_{2l-1} \neq 0, S_{2l} = 0) \\ &= P(S_{2l} + (\xi_{2l+1} + \dots + \xi_{2k}) = 0 | S_1 \neq 0, \dots, S_{2l-1} \neq 0, S_{2l} = 0) \\ &= P(S_{2l} + (\xi_{2l+1} + \dots + \xi_{2k}) = 0 | S_{2l} = 0) \\ &= P(\xi_{2l+1} + \dots + \xi_{2k} = 0) = P(S_{2(k-l)} = 0). \end{aligned}$$

Therefore

$$u_{2k} = \sum_{1 \leq l \leq k} P(S_{2(k-l)} = 0) P(\sigma_{2n} = 2l),$$

which establishes (13).

We turn now to the proof of (12). It is obviously true for  $k = 0$  and  $k = n$ . Now let  $1 \leq k \leq n - 1$ . If the particle is on the positive side for exactly  $2k$  instants, it must pass through zero. Let  $2r$  be the time of first passage through zero. There are two possibilities: either  $S_k \geq 0$ ,  $k \leq 2r$ , or  $S_k \leq 0$ ,  $k \leq 2r$ .

The number of paths of the first kind is easily seen to be

$$\left(\frac{1}{2} 2^{2r} f_{2r}\right) \cdot 2^{2(n-r)} P_{2(k-r), 2(n-r)} = \frac{1}{2} \cdot 2^{2n} \cdot f_{2r} \cdot P_{2(k-r), 2(n-r)}.$$

\* We say that the particle is on the positive side in the interval  $[m - 1, m]$  if one, at least, of the values  $S_{m-1}$  and  $S_m$  is positive.

The corresponding number of paths of the second kind is

$$\frac{1}{2} \cdot 2^{2n} \cdot f_{2r} \cdot P_{2k, 2(n-r)}.$$

Consequently, for  $1 \leq k \leq n-1$ ,

$$P_{2k, 2n} = \frac{1}{2} \sum_{r=1}^k f_{2r} \cdot P_{2(k-r), 2(n-r)} + \frac{1}{2} \sum_{r=1}^k f_{2r} \cdot P_{2k, 2(n-r)}. \quad (14)$$

Let us suppose that  $P_{2k, 2m} = u_{2k} \cdot u_{2m-2k}$  holds for  $m = 1, \dots, n-1$ . Then we find from (13) and (14) that

$$\begin{aligned} P_{2k, 2n} &= \frac{1}{2} u_{2n-2k} \cdot \sum_{r=1}^k f_{2r} \cdot u_{2k-2r} + \frac{1}{2} u_{2k} \cdot \sum_{r=1}^k f_{2r} \cdot u_{2n-2r-2k} \\ &= \frac{1}{2} u_{2n-2k} \cdot u_{2k} + \frac{1}{2} u_{2k} \cdot u_{2n-2k} = u_{2k} \cdot u_{2n-2k}. \end{aligned}$$

This completes the proof of the lemma.

Now let  $\gamma(2n)$  be the number of time units that the particle spends on the positive axis in the interval  $[0, 2n]$ . Then, when  $x < 1$ ,

$$\mathbf{P} \left\{ \frac{1}{2} < \frac{\gamma(2n)}{2n} \leq x \right\} = \sum_{\{k: 1/2 < (2k/2n) \leq x\}} P_{2k, 2n}.$$

Since

$$u_{2k} \sim \frac{1}{\sqrt{\pi k}}$$

as  $k \rightarrow \infty$ , we have

$$P_{2k, 2n} = u_{2k} \cdot u_{2(n-k)} \sim \frac{1}{\pi \sqrt{k(n-k)}},$$

as  $k \rightarrow \infty$  and  $n-k \rightarrow \infty$ .

Therefore

$$\sum_{\{k: 1/2 < (2k/2n) \leq x\}} P_{2k, 2n} - \sum_{\{k: 1/2 < (2k/2n) \leq x\}} \frac{1}{\pi n} \cdot \left[ \frac{k}{n} \left( 1 - \frac{k}{n} \right) \right]^{-1/2} \rightarrow 0, \quad n \rightarrow \infty,$$

whence

$$\sum_{\{k: 1/2 < (2k/2n) \leq x\}} P_{2k, 2n} - \frac{1}{\pi} \int_{1/2}^x \frac{dt}{\sqrt{t(1-t)}} \rightarrow 0, \quad n \rightarrow \infty.$$

But, by symmetry,

$$\sum_{\{k: k/n \leq 1/2\}} P_{2k, 2n} \rightarrow \frac{1}{2}$$

and

$$\frac{1}{\pi} \int_{1/2}^x \frac{dt}{\sqrt{t(1-t)}} = \frac{2}{\pi} \arcsin \sqrt{x} - \frac{1}{2}.$$

Consequently we have proved the following theorem.

**Theorem (Arcsine Law).** *The probability that the fraction of the time spent by the particle on the positive side is at most  $x$  tends to  $2\pi^{-1} \arcsin \sqrt{x}$ :*

$$\sum_{\{k: k/n \leq x\}} P_{2k, 2n} \rightarrow 2\pi^{-1} \arcsin \sqrt{x}. \quad (15)$$

We remark that the integrand  $p(t)$  in the integral

$$\frac{1}{\pi} \int_0^x \frac{dt}{\sqrt{t(1-t)}}$$

represents a U-shaped curve that tends to infinity as  $t \rightarrow 0$  or  $1$ .

Hence it follows that, for large  $n$ ,

$$P\left\{0 < \frac{\gamma(2n)}{2n} \leq \Delta\right\} > P\left\{\frac{1}{2} < \frac{\gamma(2n)}{2n} \leq \frac{1}{2} + \Delta\right\},$$

i.e., it is more likely that the fraction of the time spent by the particle on the positive side is close to zero or one, than to the intuitive value  $\frac{1}{2}$ .

Using a table of arcsines and noting that the convergence in (15) is indeed quite rapid, we find that

$$P\left\{\frac{\gamma(2n)}{2n} \leq 0.024\right\} \approx 0.1,$$

$$P\left\{\frac{\gamma(2n)}{2n} \leq 0.1\right\} \approx 0.2,$$

$$P\left\{\frac{\gamma(2n)}{2n} \leq 0.2\right\} \approx 0.3,$$

$$P\left\{\frac{\gamma(2n)}{2n} \leq 0.65\right\} \approx 0.6.$$

Hence if, say,  $n = 1000$ , then in about one case in ten, the particle spends only 24 units of time on the positive axis and therefore spends the greatest amount of time, 976 units, on the negative axis.

### 3. PROBLEMS

1. How fast does  $E\min(\sigma_{2n}, 2n) \rightarrow \infty$  as  $n \rightarrow \infty$ ?
2. Let  $\tau_n = \min\{1 \leq k \leq n: S_k = 1\}$ , where we take  $\tau_n = \infty$  if  $S_k < 1$  for  $1 \leq k \leq n$ . What is the limit of  $E\min(\tau_n, n)$  as  $n \rightarrow \infty$  for symmetric ( $p = q = \frac{1}{2}$ ) and for unsymmetric ( $p \neq q$ ) walks?

## §11. Martingales. Some Applications to the Random Walk

1. The Bernoulli random walk discussed above was generated by a sequence  $\xi_1, \dots, \xi_n$  of *independent* random variables. In this and the next section we introduce two important classes of *dependent* random variables, those that constitute martingales and Markov chains.

The theory of martingales will be developed in detail in Chapter VII. Here we shall present only the essential definitions, prove a theorem on the preservation of the martingale property for stopping times, and apply this to deduce the "ballot theorem." In turn, the latter theorem will be used for another proof of proposition (10.5), which was obtained above by applying the reflection principle.

2. Let  $(\Omega, \mathcal{A}, P)$  be a finite probability space and  $\mathcal{D}_1 \leq \mathcal{D}_2 \leq \dots \leq \mathcal{D}_n$  a sequence of decompositions.

**Definition 1.** A sequence of random variables  $\xi_1, \dots, \xi_n$  is called a *martingale* (with respect to the decomposition  $\mathcal{D}_1 \leq \mathcal{D}_2 \leq \dots \leq \mathcal{D}_n$ ) if

- (1)  $\xi_k$  is  $\mathcal{D}_k$ -measurable,
- (2)  $E(\xi_{k+1} | \mathcal{D}_k) = \xi_k$ ,  $1 \leq k \leq n-1$ .

In order to emphasize the system of decompositions with respect to which the random variables form a martingale, we shall use the notation

$$\xi = (\xi_k, \mathcal{D}_k)_{1 \leq k \leq n}, \quad (1)$$

where for the sake of simplicity we often do not mention explicitly that  $1 \leq k \leq n$ .

When  $\mathcal{D}_k$  is induced by  $\xi_1, \dots, \xi_n$ , i.e.

$$\mathcal{D}_k = \mathcal{D}_{\xi_1, \dots, \xi_k},$$

instead of saying that  $\xi = (\xi_k, \mathcal{D}_k)$  is a martingale, we simply say that the sequence  $\xi = (\xi_k)$  is a martingale.

Here are some examples of martingales.

**EXAMPLE 1.** Let  $\eta_1, \dots, \eta_n$  be independent Bernoulli random variables with

$$P(\eta_k = 1) = P(\eta_k = -1) = \frac{1}{2},$$

$$S_k = \eta_1 + \dots + \eta_k \quad \text{and} \quad \mathcal{D}_k = \mathcal{D}_{\eta_1, \dots, \eta_k}.$$

We observe that the decompositions  $\mathcal{D}_k$  have a simple structure:

$$\mathcal{D}_1 = \{D^+, D^-\},$$

where

$$D^+ = \{\omega: \eta_1 = +1\}, \quad D^- = \{\omega: \eta_1 = -1\}, \\ \mathcal{D}_2 = \{D^{++}, D^{+-}, D^{-+}, D^{--}\},$$

where

$$D^{++} = \{\omega: \eta_1 = +1, \eta_2 = +1\}, \dots, D^{--} = \{\omega: \eta_1 = -1, \eta_2 = -1\},$$

etc.

It is also easy to see that  $\mathcal{D}_{\eta_1, \dots, \eta_k} = \mathcal{D}_{S_1, \dots, S_k}$ .

Let us show that  $(S_k, \mathcal{D}_k)$  forms a martingale. In fact,  $S_k$  is  $\mathcal{D}_k$ -measurable, and by (8.12), (8.18) and (8.24),

$$\begin{aligned} E(S_{k+1} | \mathcal{D}_k) &= E(S_k + \eta_{k+1} | \mathcal{D}_k) \\ &= E(S_k | \mathcal{D}_k) + E(\eta_{k+1} | \mathcal{D}_k) = S_k + E\eta_{k+1} = S_k. \end{aligned}$$

If we put  $S_0 = 0$  and take  $D_0 = \{\Omega\}$ , the trivial decomposition, then the sequence  $(S_k, \mathcal{D}_k)_{0 \leq k \leq n}$  also forms a martingale.

EXAMPLE 2. Let  $\eta_1, \dots, \eta_n$  be independent Bernoulli random variables with  $P(\eta_i = 1) = p$ ,  $P(\eta_i = -1) = q$ . If  $p \neq q$ , each of the sequences  $\xi = (\xi_k)$  with

$$\xi_k = \left(\frac{q}{p}\right)^{S_k}, \quad \xi_k = S_k - k(p - q), \quad \text{where } S_k = \eta_1 + \dots + \eta_k,$$

is a martingale.

EXAMPLE 3. Let  $\eta$  be a random variable,  $\mathcal{D}_1 \preceq \dots \preceq \mathcal{D}_n$ , and

$$\xi_k = E(\eta | \mathcal{D}_k). \quad (2)$$

Then the sequence  $\xi = (\xi_k, \mathcal{D}_k)$  is a martingale. In fact, it is evident that  $E(\eta | \mathcal{D}_k)$  is  $\mathcal{D}_k$ -measurable, and by (8.20)

$$E(\xi_{k+1} | \mathcal{D}_k) = E[E(\eta | \mathcal{D}_{k+1}) | \mathcal{D}_k] = E(\eta | \mathcal{D}_k) = \xi_k.$$

In this connection we notice that if  $\xi = (\xi_k, \mathcal{D}_k)$  is any martingale, then by (8.20)

$$\begin{aligned} \xi_k &= E(\xi_{k+1} | \mathcal{D}_k) = E[E(\xi_{k+2} | \mathcal{D}_{k+1}) | \mathcal{D}_k] \\ &= E(\xi_{k+2} | \mathcal{D}_k) = \dots = E(\xi_n | \mathcal{D}_k). \end{aligned} \quad (3)$$

Consequently the set of martingales  $\xi = (\xi_k, \mathcal{D}_k)$  is exhausted by the martingales of the form (2). (We note that for infinite sequences  $\xi = (\xi_k, \mathcal{D}_k)_{k \geq 1}$  this is, in general, no longer the case; see Problem 7 in §1 of Chapter VII.)

**EXAMPLE 4.** Let  $\eta_1, \dots, \eta_n$  be a sequence of independent identically distributed random variables,  $S_k = \eta_1 + \dots + \eta_k$ , and  $\mathcal{D}_1 = \mathcal{D}_{S_n}$ ,  $\mathcal{D}_2 = \mathcal{D}_{S_n, S_{n-1}}, \dots$ ,  $\mathcal{D}_n = \mathcal{D}_{S_n, S_{n-1}, \dots, S_1}$ . Let us show that the sequence  $\xi = (\xi_k, \mathcal{D}_k)$  with

$$\xi_1 = \frac{S_n}{n}, \xi_2 = \frac{S_{n-1}}{n-1}, \dots, \xi_k = \frac{S_{n+1-k}}{n+1-k}, \dots, \xi_n = S_1$$

is a martingale. In the first place, it is clear that  $\mathcal{D}_k \leq \mathcal{D}_{k+1}$  and  $\xi_k$  is  $\mathcal{D}_k$ -measurable. Moreover, we have by symmetry, for  $j \leq n - k + 1$ ,

$$\mathbf{E}(\eta_j | \mathcal{D}_k) = \mathbf{E}(\eta_1 | \mathcal{D}_k) \quad (4)$$

(compare (8.26)). Therefore

$$(n - k + 1)\mathbf{E}(\eta_1 | \mathcal{D}_k) = \sum_{j=1}^{n-k+1} \mathbf{E}(\eta_j | \mathcal{D}_k) = \mathbf{E}(S_{n-k+1} | \mathcal{D}_k) = S_{n-k+1},$$

and consequently

$$\xi_k = \frac{S_{n-k+1}}{n - k + 1} = \mathbf{E}(\eta_1 | \mathcal{D}_k),$$

and it follows from Example 3 that  $\xi = (\xi_k, \mathcal{D}_k)$  is a martingale.

**Remark.** From this martingale property of the sequence  $\xi = (\xi_k, \mathcal{D}_k)_{1 \leq k \leq n}$ , it is clear why we will sometimes say that the sequence  $(S_k/k)_{1 \leq k \leq n}$  forms a *reversed martingale*. (Compare problem 6 in §1 of Chapter VII.)

**EXAMPLE 5.** Let  $\eta_1, \dots, \eta_n$  be independent Bernoulli random variables with

$$\mathbf{P}(\eta_i = +1) = \mathbf{P}(\eta_i = -1) = \frac{1}{2},$$

$S_k = \eta_1 + \dots + \eta_k$ . Let  $A$  and  $B$  be integers,  $A < 0 < B$ . Then with  $0 < \lambda < \pi/2$ , the sequence  $\xi = (\xi_k, \mathcal{D}_k)$  with  $\mathcal{D}_k = \mathcal{D}_{S_1, \dots, S_k}$  and

$$\xi_k = (\cos \lambda)^{-k} \exp \left\{ i\lambda \left( S_k - \frac{B+A}{2} \right) \right\} \quad (5)$$

is a complex martingale (i.e., the real and imaginary parts of  $\xi_k$  form martingales).

3. It follows from the definition of a martingale that the expectation  $\mathbf{E}\xi_k$  is the same for every  $k$ :

$$\mathbf{E}\xi_k = \mathbf{E}\xi_1.$$

It turns out that this property persists if time  $k$  is replaced by a random time.

In order to formulate this property we introduce the following definition.

**Definition 2.** A random variable  $\tau = \tau(\omega)$  that takes the values  $1, 2, \dots, n$  is called a *stopping time* (with respect to a decomposition  $(\mathcal{D}_k)_{1 \leq k \leq n}$ ,  $\mathcal{D}_1 \leq \mathcal{D}_2 \leq \dots \leq \mathcal{D}_n$ ) if, for  $k = 1, \dots, n$ , the random variable  $I_{\{\tau \leq k\}}(\omega)$  is  $\mathcal{D}_k$ -measurable.

If we consider  $\mathcal{D}_k$  as the decomposition induced by observations for  $k$  steps (for example,  $\mathcal{D}_k = \mathcal{D}_{\eta_1, \dots, \eta_k}$ , the decomposition induced by the variables  $\eta_1, \dots, \eta_k$ ), then the  $\mathcal{D}_k$ -measurability of  $I_{\{\tau=k\}}(\omega)$  means that the realization or nonrealization of the event  $\{\tau=k\}$  is determined only by observations for  $k$  steps (and is independent of the "future").

If  $\mathcal{B}_k = \alpha(\mathcal{D}_k)$ , then the  $\mathcal{D}_k$ -measurability of  $I_{\{\tau=k\}}(\omega)$  is equivalent to the assumption that

$$\{\tau=k\} \in \mathcal{B}_k. \quad (6)$$

We have already introduced specific examples of stopping times: the times  $\tau_k^*$ ,  $\sigma_{2n}$  introduced in §§9 and 10. Those times are special cases of stopping times of the form

$$\begin{aligned} \tau^A &= \min\{0 < k \leq n: \xi_k \in A\}, \\ \sigma^A &= \min\{0 \leq k \leq n: \xi_k \in A\}, \end{aligned} \quad (7)$$

which are the times (respectively the first time after zero and the first time) for a sequence  $\xi_0, \xi_1, \dots, \xi_n$  to attain a point of the set  $A$ .

**4. Theorem 1.** *Let  $\xi = (\xi_k, \mathcal{D}_k)_{1 \leq k \leq n}$  be a martingale and  $\tau$  a stopping time with respect to the decomposition  $(\mathcal{D}_k)_{1 \leq k \leq n}$ . Then*

$$\mathbf{E}(\xi_\tau | \mathcal{D}_1) = \xi_1, \quad (8)$$

where

$$\xi_\tau = \sum_{k=1}^n \xi_k I_{\{\tau \geq k\}}(\omega) \quad (9)$$

and

$$\mathbf{E} \xi_\tau = \mathbf{E} \xi_1. \quad (10)$$

**PROOF** (compare the proof of (9.29)). Let  $D \in \mathcal{D}_1$ . Using (3) and the properties of conditional expectations, we find that

$$\begin{aligned} \mathbf{E}(\xi_\tau | D) &= \frac{\mathbf{E}(\xi_\tau I_D)}{\mathbf{P}(D)} \\ &= \frac{1}{\mathbf{P}(D)} \cdot \sum_{l=1}^n \mathbf{E}(\xi_l \cdot I_{\{\tau \geq l\}} \cdot I_D) \\ &= \frac{1}{\mathbf{P}(D)} \sum_{l=1}^n \mathbf{E}[\mathbf{E}(\xi_l | \mathcal{D}_l) \cdot I_{\{\tau \geq l\}} \cdot I_D] \\ &= \frac{1}{\mathbf{P}(D)} \sum_{l=1}^n \mathbf{E}[\mathbf{E}(\xi_l I_{\{\tau \geq l\}} \cdot I_D | \mathcal{D}_l)] \\ &= \frac{1}{\mathbf{P}(D)} \sum_{l=1}^n \mathbf{E}[\xi_l I_{\{\tau \geq l\}} \cdot I_D] \\ &= \frac{1}{\mathbf{P}(D)} \mathbf{E}(\xi_n I_D) = \mathbf{E}(\xi_n | D), \end{aligned}$$

and consequently

$$E(\xi_\tau | \mathcal{D}_1) = E(\xi_n | \mathcal{D}_1) = \xi_1.$$

The equation  $E\xi_\tau = E\xi_1$  then follows in an obvious way.

This completes the proof of the theorem.

**Corollary.** For the martingale  $(S_k, \mathcal{D}_k)_{1 \leq k \leq n}$  of Example 1, and any stopping time  $\tau$  (with respect to  $(\mathcal{D}_k)$ ) we have the formulas

$$ES_\tau = 0, \quad ES_\tau^2 = E\tau, \quad (11)$$

known as Wald's identities (cf. (9.29) and (9.30); see also Problem 1 and Theorem 3 in §2 of Chapter VII).

5. Let us use Theorem 1 to establish the following proposition.

**Theorem 2 (Ballot Theorem).** Let  $\eta_1, \dots, \eta_n$  be a sequence of independent identically distributed random variables whose values are nonnegative integers,  $S_k = \eta_1 + \dots + \eta_k$ ,  $1 \leq k \leq n$ . Then

$$P\{S_k < k \text{ for all } k, 1 \leq k \leq n | S_n\} = \left(1 - \frac{S_n}{n}\right)^+, \quad (12)$$

where  $a^+ = \max(a, 0)$ .

**PROOF.** On the set  $\{\omega: S_n \geq n\}$  the formula is evident. We therefore prove (12) for the sample points at which  $S_n < n$ .

Let us consider the martingale  $\xi = (\xi_k, \mathcal{D}_k)_{1 \leq k \leq n}$  introduced in Example 4, with  $\xi_k = S_{n+1-k}/(n+1-k)$  and  $\mathcal{D}_k = \mathcal{D}_{S_{n+1-k}, \dots, S_n}$ .

We define

$$\tau = \min\{1 \leq k \leq n: \xi_k \geq 1\},$$

taking  $\tau = n$  on the set  $\{\xi_k < 1 \text{ for all } k \text{ such that } 1 \leq k \leq n\} = \{\max_{1 \leq l \leq n} (S_l/l) < 1\}$ . It is clear that  $\xi_\tau = \xi_n = S_1 = 0$  on this set, and therefore

$$\left\{ \max_{1 \leq l \leq n} \frac{S_l}{l} < 1 \right\} = \left\{ \max_{1 \leq l \leq n} \frac{S_l}{l} < 1, S_n < n \right\} \subseteq \{\xi_\tau = 0\}. \quad (13)$$

Now let us consider those outcomes for which simultaneously  $\max_{1 \leq l \leq n} (S_l/l) \geq 1$  and  $S_n < n$ . Write  $\sigma = n + 1 - \tau$ . It is easy to see that

$$\sigma = \max\{1 \leq k \leq n: S_k \geq k\}$$

and therefore (since  $S_n < n$ ) we have  $\sigma < n$ ,  $S_\sigma \geq \sigma$ , and  $S_{\sigma+1} < \sigma + 1$ . Consequently  $\eta_{\sigma+1} = S_{\sigma+1} - S_\sigma < (\sigma + 1) - \sigma = 1$ , i.e.  $\eta_{\sigma+1} = 0$ . Therefore  $\sigma \leq S_\sigma = S_{\sigma+1} < \sigma + 1$ , and consequently  $S_\sigma = \sigma$  and

$$\xi_\tau = \frac{S_{n+1-\tau}}{n+1-\tau} = \frac{S_\sigma}{\sigma} = 1.$$



Therefore

$$\left\{ \max_{1 \leq l \leq n} \frac{S_l}{l} \geq 1, S_n < n \right\} \subseteq \{\xi_\tau = 1\}. \quad (14)$$

From (13) and (14) we find that

$$\left\{ \max_{1 \leq l \leq n} \frac{S_l}{l} \geq 1, S_n < n \right\} = \{\xi_\tau = 1\} \cap \{S_n < n\}.$$

Therefore, on the set  $\{S_n < n\}$ , we have

$$\mathbf{P}\left\{ \max_{1 \leq l \leq n} \frac{S_l}{l} \geq 1 | S_n \right\} = \mathbf{P}\{\xi_\tau = 1 | S_n\} = \mathbf{E}(\xi_\tau | S_n),$$

where the last equation follows because  $\xi_\tau$  takes only the two values 0 and 1.

Let us notice now that  $\mathbf{E}(\xi_\tau | S_n) = \mathbf{E}(\xi_\tau | \mathcal{D}_1)$ , and (by Theorem 1)  $\mathbf{E}(\xi_\tau | \mathcal{D}_1) = \xi_1 = S_n/n$ . Consequently, on the set  $\{S_n < n\}$  we have  $\mathbf{P}\{S_k < k \text{ for all } k \text{ such that } 1 \leq k \leq n | S_n\} = 1 - (S_n/n)$ .

This completes the proof of the theorem.

We now apply this theorem to obtain a different proof of Lemma 1 of §10, and explain why it is called the ballot theorem.

Let  $\xi_1, \dots, \xi_n$  be independent Bernoulli random variables with

$$\mathbf{P}(\xi_1 = 1) = \mathbf{P}(\xi_i = -1) = \frac{1}{2},$$

$S_k = \xi_1 + \dots + \xi_k$  and  $a, b$  nonnegative integers such that  $a - b > 0$ ,  $a + b = n$ . We are going to show that

$$\mathbf{P}\{S_1 > 0, \dots, S_n > 0 | S_n = a - b\} = \frac{a - b}{a + b}. \quad (15)$$

In fact, by symmetry,

$$\begin{aligned} & \mathbf{P}\{S_1 > 0, \dots, S_n > 0 | S_n = a - b\} \\ &= \mathbf{P}\{S_1 < 0, \dots, S_n < 0 | S_n = -(a - b)\} \\ &= \mathbf{P}\{S_1 + 1 < 1, \dots, S_n + n < n | S_n + n = n - (a - b)\} \\ &= \mathbf{P}\{\eta_1 < 1, \dots, \eta_1 + \dots + \eta_n < n | \eta_1 + \dots + \eta_n = n - (a - b)\} \\ &= \left[ 1 - \frac{n - (a - b)}{n} \right]^+ = \frac{a - b}{n} = \frac{a - b}{a + b}, \end{aligned}$$

where we have put  $\eta_k = \xi_k + 1$  and applied (12).

Now formula (10.5) follows from (15) in an evident way; the formula was also established in Lemma 1 of §10 by using the reflection principle.

Let us interpret  $\xi_i = +1$  as a vote for candidate  $A$  and  $\xi_i = -1$  as a vote for  $B$ . Then  $S_k$  is the difference between the numbers of votes cast for  $A$  and  $B$  at the time when  $k$  votes have been recorded, and

$$P\{S_1 > 0, \dots, S_n > 0 | S_n = a - b\}$$

is the probability that  $A$  was always ahead of  $B$ , with the understanding that  $A$  received  $a$  votes in all,  $B$  received  $b$  votes, and  $a - b > 0$ ,  $a + b = n$ . According to (15) this probability is  $(a - b)/n$ .

## 6. PROBLEMS

1. Let  $\mathcal{D}_0 \leq \mathcal{D}_1 \leq \dots \leq \mathcal{D}_n$  be a sequence of decompositions with  $\mathcal{D}_0 = \{\Omega\}$ , and let  $\eta_k$  be  $\mathcal{D}_k$ -measurable variables,  $1 \leq k \leq n$ . Show that the sequence  $\xi = (\xi_k, \mathcal{D}_k)$  with

$$\xi_k = \sum_{i=1}^k [\eta_i - E(\eta_i | \mathcal{D}_{i-1})]$$

is a martingale.

2. Let the random variables  $\eta_1, \dots, \eta_k$  satisfy  $E(\eta_k | \eta_1, \dots, \eta_{k-1}) = 0$ . Show that the sequence  $\xi = (\xi_k)_{1 \leq k \leq n}$  with  $\xi_1 = \eta_1$  and

$$\xi_{k+1} = \sum_{i=1}^k \eta_{i+1} f_i(\eta_1, \dots, \eta_i),$$

where  $f_i$  are given functions, is a martingale.

3. Show that every martingale  $\xi = (\xi_i, \mathcal{D}_k)$  has uncorrelated increments: if  $a < b < c < d$  then

$$\text{cov}(\xi_d - \xi_c, \xi_b - \xi_a) = 0.$$

4. Let  $\xi = (\xi_1, \dots, \xi_n)$  be a random sequence such that  $\xi_k$  is  $\mathcal{D}_k$ -measurable ( $\mathcal{D} \leq \mathcal{D}_2 \leq \dots \leq \mathcal{D}_n$ ). Show that a necessary and sufficient condition for this sequence to be a martingale (with respect to the system  $(\mathcal{D}_k)$ ) is that  $E\xi_\tau = E\xi_1$  for every stopping time  $\tau$  (with respect to  $(\mathcal{D}_k)$ ). (The phrase "for every stopping time" can be replaced by "for every stopping time that assumes two values.")

5. Show that if  $\xi = (\xi_k, \mathcal{D}_k)_{1 \leq k \leq n}$  is a martingale and  $\tau$  is a stopping time, then

$$E[\xi_n I_{\{\tau=k\}}] = E[\xi_k I_{\{\tau=k\}}]$$

for every  $k$ .

6. Let  $\xi = (\xi_k, \mathcal{D}_k)$  and  $\eta = (\eta_k, \mathcal{D}_k)$  be two martingales,  $\xi_1 = \eta_1 = 0$ . Show that

$$E\xi_n \eta_n = \sum_{k=2}^n E(\xi_k - \xi_{k-1})(\eta_k - \eta_{k-1})$$

and in particular that

$$E\xi_n^2 = \sum_{k=2}^n E(\xi_k - \xi_{k-1})^2.$$

7. Let  $\eta_1, \dots, \eta_n$  be a sequence of independent identically distributed random variables with  $E\eta_i = 0$ . Show that the sequence  $\xi = (\xi_k)$  with

$$\xi_k = \left( \sum_{i=1}^k \eta_i \right)^2 - kE\eta_i^2,$$

$$\xi_k = \frac{\exp \lambda(\eta_1 + \dots + \eta_k)}{(E \exp \lambda \eta_1)^k}$$

is a martingale.

8. Let  $\eta_1, \dots, \eta_n$  be a sequence of independent identically distributed random variables taking values in a finite set  $Y$ . Let  $f_0(y) = P(\eta_1 = y)$ ,  $y \in Y$ , and let  $f_1(y)$  be a non-negative function with  $\sum_{y \in Y} f_1(y) = 1$ . Show that the sequence  $\xi = (\xi_k, \mathcal{D}_k^n)$  with  $\mathcal{D}_k^n = D_{\eta_1, \dots, \eta_k}$ ,

$$\xi_k = \frac{f_1(\eta_1) \cdots f_1(\eta_k)}{f_0(\eta_1) \cdots f_0(\eta_k)},$$

is a martingale. (The variables  $\xi_k$ , known as *likelihood ratios*, are extremely important in mathematical statistics.)

## §12. Markov Chains. Ergodic Theorem. Strong Markov Property

1. We have discussed the Bernoulli scheme with

$$\Omega = \{\omega: \omega = (x_1, \dots, x_n), x_i = 0, 1\},$$

where the probability  $p(\omega)$  of each outcome is given by

$$p(\omega) = p(x_1) \cdots p(x_n), \quad (1)$$

with  $p(x) = p^x q^{1-x}$ . With these hypotheses, the variables  $\xi_1, \dots, \xi_n$  with  $\xi_i(\omega) = x_i$  are *independent and identically distributed* with

$$P(\xi_1 = x) = \cdots = P(\xi_n = x) = p(x), \quad x = 0, 1.$$

If we replace (1) by

$$p(\omega) = p_1(x_1) \cdots p_n(x_n),$$

where  $p_i(x) = p_i^x(1 - p_i)$ ,  $0 \leq p_i \leq 1$ , the random variables  $\xi_1, \dots, \xi_n$  are still *independent*, but in general are *differently distributed*:

$$P(\xi_1 = x) = p_1(x), \dots, P(\xi_n = x) = p_n(x).$$

We now consider a generalization that leads to *dependent* random variables that form what is known as a Markov chain.

Let us suppose that

$$\Omega = \{\omega: \omega = (x_0, x_1, \dots, x_n), x_i \in X\},$$

where  $X$  is a finite set. Let there be given nonnegative functions  $p_0(x)$ ,  $p_1(x, y)$ ,  $\dots$ ,  $p_n(x, y)$  such that

$$\begin{aligned}\sum_{x \in X} p_0(x) &= 1, \\ \sum_{y \in X} p_k(x, y) &= 1, \quad k = 1, \dots, n; \quad y \in X.\end{aligned}\quad (2)$$

For each  $\omega = (x_0, x_1, \dots, x_n)$ , put

$$p(\omega) = p_0(x_0)p_1(x_0, x_1) \cdots p_n(x_{n-1}, x_n). \quad (3)$$

It is easily verified that  $\sum_{\omega \in \Omega} p(\omega) = 1$ , and consequently the set of numbers  $p(\omega)$  together with the space  $\Omega$  and the collection of its subsets defines a probabilistic model, which it is usual to call a *model of experiments that form a Markov chain*.

Let us introduce the random variables  $\xi_0, \xi_1, \dots, \xi_n$  with  $\xi_i(\omega) = x_i$ . A simple calculation shows that

$$\begin{aligned}P(\xi_0 = a) &= p_0(a), \\ P(\xi_0 = a_0, \dots, \xi_k = a_k) &= p_0(a_0)p_1(a_0, a_1) \cdots p_k(a_{k-1}, a_k).\end{aligned}\quad (4)$$

We now establish the validity of the following fundamental property of conditional probabilities:

$$P\{\xi_{k+1} = a_{k+1} | \xi_k = a_k, \dots, \xi_0 = a_0\} = P\{\xi_{k+1} = a_{k+1} | \xi_k = a_k\} \quad (5)$$

(under the assumption that  $P(\xi_k = a_k, \dots, \xi_0 = a_0) > 0$ ).

By (4),

$$\begin{aligned}P\{\xi_{k+1} = a_{k+1} | \xi_k = a_k, \dots, \xi_0 = a_0\} \\ &= \frac{P\{\xi_{k+1} = a_{k+1}, \dots, \xi_0 = a_0\}}{P\{\xi_k = a_k, \dots, \xi_0 = a_0\}} \\ &= \frac{p_0(a_0)p_1(a_0, a_1) \cdots p_{k+1}(a_k, a_{k+1})}{p_0(a_0) \cdots p_k(a_{k-1}, a_k)} = p_{k+1}(a_k, a_{k+1}).\end{aligned}$$

In a similar way we verify

$$P\{\xi_{k+1} = a_{k+1} | \xi_k = a_k\} = p_{k+1}(a_k, a_{k+1}), \quad (6)$$

which establishes (5).

Let  $\mathcal{D}_k^\xi = \mathcal{D}_{\xi_0, \dots, \xi_k}$  be the decomposition induced by  $\xi_0, \dots, \xi_k$ , and  $\mathcal{B}_k^\xi = \sigma(\mathcal{D}_k^\xi)$ .

Then, in the notation introduced in §8, it follows from (5) that

$$P\{\xi_{k+1} = a_{k+1} | \mathcal{B}_k^\xi\} = P\{\xi_{k+1} = a_{k+1} | \xi_k\} \quad (7)$$

or

$$P\{\xi_{k+1} = a_{k+1} | \xi_0, \dots, \xi_k\} = P\{\xi_{k+1} = a_{k+1} | \xi_k\}.$$

If we use the evident equation

$$P(AB|C) = P(A|BC)P(B|C),$$

we find from (7) that

$$P\{\xi_n = a_n, \dots, \xi_{k+1} = a_{k+1} | \mathcal{B}_k^{\xi}\} = P\{\xi_n = a_n, \dots, \xi_{k+1} = a_{k+1} | \xi_k\} \quad (8)$$

or

$$P\{\xi_n = a_n, \dots, \xi_{k+1} = a_{k+1} | \xi_0, \dots, \xi_k\} = P\{\xi_n = a_n, \dots, \xi_{k+1} = a_{k+1} | \xi_k\}. \quad (9)$$

This equation admits the following intuitive interpretation. Let us think of  $\xi_k$  as the position of a particle "at present,"  $(\xi_0, \dots, \xi_{k-1})$  as the "past," and  $(\xi_{k+1}, \dots, \xi_n)$  as the "future." Then (9) says that if the past and the present are given, the future depends only on the present and is independent of how the particle arrived at  $\xi_k$ , i.e. is independent of the past  $(\xi_0, \dots, \xi_{k-1})$ .

Let  $F = (\xi_n = a_n, \dots, \xi_{k+1} = a_{k+1})$ ,  $N = \{\xi_k = a_k\}$ ,

$$B = \{\xi_{k-1} = a_{k-1}, \dots, \xi_0 = a_0\}.$$

Then it follows from (9) that

$$P(F|NB) = P(F|N),$$

from which we easily find that

$$P(FB|N) = P(F|N)P(B|N). \quad (10)$$

In other words, it follows from (7) that for a given present  $N$ , the future  $F$  and the past  $B$  are independent. It is easily shown that the converse also holds: if (10) holds for all  $k = 0, 1, \dots, n-1$ , then (7) holds for every  $k = 0, 1, \dots, n-1$ .

The property of the independence of future and past, or, what is the same thing, the lack of dependence of the future on the past when the present is given, is called the *Markov property*, and the corresponding sequence of random variables  $\xi_0, \dots, \xi_n$  is a *Markov chain*.

Consequently if the probabilities  $p(\omega)$  of the sample points are given by (3), the sequence  $(\xi_0, \dots, \xi_n)$  with  $\xi_i(\omega) = x_i$  forms a Markov chain.

We give the following formal definition.

**Definition.** Let  $(\Omega, \mathcal{A}, P)$  be a (finite) probability space and let  $\xi = (\xi_0, \dots, \xi_n)$  be a sequence of random variables with values in a (finite) set  $X$ . If (7) is satisfied, the sequence  $\xi = (\xi_0, \dots, \xi_n)$  is called a (finite) *Markov chain*.

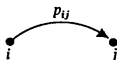
The set  $X$  is called the *phase space* or *state space* of the chain. The set of probabilities  $(p_n(x))$ ,  $x \in X$ , with  $p_0(x) = P(\xi_0 = x)$  is the *initial distribution*, and the matrix  $\|p_k(x, y)\|$ ,  $x, y \in X$ , with  $p(x, y) = P\{\xi_k = y | \xi_{k-1} = x\}$  is the *matrix of transition probabilities* (from state  $x$  to state  $y$ ) at time  $k = 1, \dots, n$ .

When the transition probabilities  $p_k(x, y)$  are independent of  $k$ , that is,  $p_k(x, y) = p(x, y)$ , the sequence  $\xi = (\xi_0, \dots, \xi_n)$  is called a *homogeneous* Markov chain with transition matrix  $\|p(x, y)\|$ .

Notice that the matrix  $\|p(x, y)\|$  is *stochastic*: its elements are nonnegative and the sum of the elements in each row is 1:  $\sum_y p(x, y) = 1, x \in X$ .

We shall suppose that the phase space  $X$  is a finite set of integers ( $X = \{0, 1, \dots, N\}$ ,  $X = \{0, \pm 1, \dots, \pm N\}$ , etc.), and use the traditional notation  $p_i = p_0(i)$  and  $p_{ij} = p(i, j)$ .

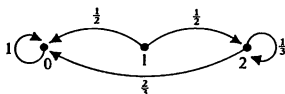
It is clear that the properties of homogeneous Markov chains completely determine the initial distributions  $p_i$  and the transition probabilities  $p_{ij}$ . In specific cases we describe the evolution of the chain, not by writing out the matrix  $\|p_{ij}\|$  explicitly, but by a (directed) graph whose vertices are the states in  $X$ , and an arrow from state  $i$  to state  $j$  with the number  $p_{ij}$  over it indicates that it is possible to pass from point  $i$  to point  $j$  with probability  $p_{ij}$ . When  $p_{ij} = 0$ , the corresponding arrow is omitted.



EXAMPLE 1. Let  $X = \{0, 1, 2\}$  and

$$\|p_{ij}\| = \begin{pmatrix} 1 & 0 & 0 \\ \frac{1}{2} & 0 & \frac{1}{2} \\ \frac{2}{3} & 0 & \frac{1}{3} \end{pmatrix}.$$

The following graph corresponds to this matrix:

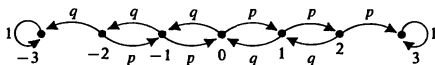


Here state 0 is said to be *absorbing*: if the particle gets into this state it remains there, since  $p_{00} = 1$ . From state 1 the particle goes to the adjacent states 0 or 2 with equal probabilities; state 2 has the property that the particle remains there with probability  $\frac{1}{3}$  and goes to state 0 with probability  $\frac{2}{3}$ .

EXAMPLE 2. Let  $X = \{0, \pm 1, \dots, \pm N\}$ ,  $p_0 = 1$ ,  $p_{NN} = p_{(-N)(-N)} = 1$ , and, for  $|i| < N$ ,

$$p_{ij} = \begin{cases} p, & j = i + 1, \\ q, & j = i - 1, \\ 0 & \text{otherwise.} \end{cases} \quad (11)$$

The transitions corresponding to this chain can be presented graphically in the following way ( $N = 3$ ):



This chain corresponds to the two-player game discussed earlier, when each player has a bankroll  $N$  and at each turn the first player wins  $+1$  from the second with probability  $p$ , and loses (wins  $-1$ ) with probability  $q$ . If we think of state  $i$  as the amount won by the first player from the second, then reaching state  $N$  or  $-N$  means the ruin of the second or first player, respectively.

In fact, if  $\eta_1, \eta_2, \dots, \eta_n$  are independent Bernoulli random variables with  $P(\eta_i = +1) = p$ ,  $P(\eta_i = -1) = q$ ,  $S_0 = 0$  and  $S_k = \eta_1 + \dots + \eta_k$  the amounts won by the first player from the second, then the sequence  $S_0, S_1, \dots, S_n$  is a Markov chain with  $p_0 = 1$  and transition matrix (11), since

$$\begin{aligned}
 P\{S_{k+1} = j | S_k = i_k, S_{k-1} = i_{k-1}, \dots\} \\
 &= P\{S_k + \eta_{k+1} = j | S_k = i_k, S_{k-1} = i_{k-1}, \dots\} \\
 &= P\{S_k + \eta_{k+1} = j | S_k = i_k\} = P\{\eta_{k+1} = j - i_k\}.
 \end{aligned}$$

This Markov chain has a very simple structure:

$$S_{k+1} = S_k + \eta_{k+1}, \quad 0 \leq k \leq n-1,$$

where  $\eta_1, \eta_2, \dots, \eta_n$  is a sequence of independent random variables.

The same considerations show that if  $\xi_0, \eta_1, \dots, \eta_n$  are independent random variables then the sequence  $\xi_0, \xi_1, \dots, \xi_n$  with

$$\xi_{k+1} = f_k(\xi_k, \eta_{k+1}), \quad 0 \leq k \leq n-1, \quad (12)$$

is also a Markov chain.

It is worth noting in this connection that a Markov chain constructed in this way can be considered as a natural probabilistic analog of a (deterministic) sequence  $x = (x_0, \dots, x_n)$  generated by the recurrent equations

$$x_{k+1} = f_k(x_k).$$

We now give another example of a Markov chain of the form (12); this example arises in queueing theory.

**EXAMPLE 3.** At a taxi stand let taxis arrive at unit intervals of time (one at a time). If no one is waiting at the stand, the taxi leaves immediately. Let  $\eta_k$  be the number of passengers who arrive at the stand at time  $k$ , and suppose that  $\eta_1, \dots, \eta_n$  are independent random variables. Let  $\xi_k$  be the length of the

waiting line at time  $k$ ,  $\xi_0 = 0$ . Then if  $\xi_k = i$ , at the next time  $k + 1$  the length  $\xi_{k+1}$  of the waiting line is equal to

$$j = \begin{cases} \eta_{k+1} & \text{if } i = 0, \\ i - 1 + \eta_{k+1} & \text{if } i \geq 1. \end{cases}$$

In other words,

$$\xi_{k+1} = (\xi_k - 1)^+ + \eta_{k+1}, \quad 0 \leq k \leq n - 1,$$

where  $a^+ = \max(a, 0)$ , and therefore the sequence  $\xi = (\xi_0, \dots, \xi_n)$  is a Markov chain.

**EXAMPLE 4.** This example comes from the theory of *branching processes*. A branching process with discrete times is a sequence of random variables  $\xi_0, \xi_1, \dots, \xi_n$ , where  $\xi_k$  is interpreted as the number of particles in existence at time  $k$ , and the process of creation and annihilation of particles is as follows: each particle, independently of the other particles and of the "pre-history" of the process, is transformed into  $j$  particles with probability  $p_j$ ,  $j = 0, 1, \dots, M$ .

We suppose that at the initial time there is just one particle,  $\xi_0 = 1$ . If at time  $k$  there are  $\xi_k$  particles (numbered  $1, 2, \dots, \xi_k$ ), then by assumption  $\xi_{k+1}$  is given as a random sum of random variables,

$$\xi_{k+1} = \eta_1^{(k)} + \dots + \eta_{\xi_k}^{(k)},$$

where  $\eta_i^{(k)}$  is the number of particles produced by particle number  $i$ . It is clear that if  $\xi_k = 0$  then  $\xi_{k+1} = 0$ . If we suppose that all the random variables  $\eta_j^{(k)}$ ,  $k \geq 0$ , are independent of each other, we obtain

$$\begin{aligned} \mathbf{P}\{\xi_{k+1} = i_{k+1} | \xi_k = i_k, \xi_{k-1} = i_{k-1}, \dots\} &= \mathbf{P}\{\xi_{k+1} = i_{k+1} | \xi_k = i_k\} \\ &= \mathbf{P}\{\eta_1^{(k)} + \dots + \eta_{i_k}^{(k)} = i_{k+1}\}. \end{aligned}$$

It is evident from this that the sequence  $\xi_0, \xi_1, \dots, \xi_n$  is a Markov chain.

A particularly interesting case is that in which each particle either vanishes with probability  $q$  or divides in two with probability  $p$ ,  $p + q = 1$ . In this case it is easy to calculate that

$$p_{ij} = \mathbf{P}\{\xi_{k+1} = j | \xi_k = i\}$$

is given by the formula

$$p_{ij} = \begin{cases} C_i^{j/2} p^{j/2} q^{i-j/2}, & j = 0, \dots, 2i, \\ 0 & \text{in all other cases.} \end{cases}$$

2. Let  $\xi = (\xi_k, \Pi, \mathbb{P})$  be a homogeneous Markov chain with starting vectors (rows)  $\Pi = (p_i)$  and transition matrix  $\Pi = \|p_{ij}\|$ . It is clear that

$$p_{ij} = \mathbf{P}\{\xi_1 = j | \xi_0 = i\} = \dots = \mathbf{P}\{\xi_n = j | \xi_{n-1} = i\}.$$



We shall use the notation

$$p_{ij}^{(k)} = \mathbf{P}\{\xi_k = j | \xi_0 = i\} \quad (= \mathbf{P}\{\xi_{k+l} = j | \xi_l = i\})$$

for the probability of a transition from state  $i$  to state  $j$  in  $k$  steps, and

$$p_j^{(k)} = \mathbf{P}\{\xi_k = j\}$$

for the probability of finding the particle at point  $j$  at time  $k$ . Also let

$$\mathbb{P}^{(k)} = \|p_i^{(k)}\|, \quad \mathbb{P}^{(k)} = \|p_{ij}^{(k)}\|.$$

Let us show that the transition probabilities  $p_{ij}^{(k)}$  satisfy the *Kolmogorov–Chapman equation*

$$p_{ij}^{(k+l)} = \sum_{\alpha} p_{i\alpha}^{(k)} p_{\alpha j}^{(l)}, \quad (13)$$

or, in matrix form,

$$\mathbb{P}^{(k+l)} = \mathbb{P}^{(k)} \cdot \mathbb{P}^{(l)} \quad (14)$$

The proof is extremely simple: using the formula for total probability and the Markov property, we obtain

$$\begin{aligned} p_{ij}^{(k+l)} &= \mathbf{P}(\xi_{k+l} = j | \xi_0 = i) = \sum_{\alpha} \mathbf{P}(\xi_{k+l} = j, \xi_k = \alpha | \xi_0 = i) \\ &= \sum_{\alpha} \mathbf{P}(\xi_{k+l} = j | \xi_k = \alpha) \mathbf{P}(\xi_k = \alpha | \xi_0 = i) = \sum_{\alpha} p_{\alpha j}^{(l)} p_{i\alpha}^{(k)}. \end{aligned}$$

The following two cases of (13) are particularly important:

the *backward equation*

$$p_{ij}^{(l+1)} = \sum_{\alpha} p_{i\alpha} p_{\alpha j}^{(l)} \quad (15)$$

and the *forward equation*

$$p_{ij}^{(k+1)} = \sum_{\alpha} p_{i\alpha}^{(k)} p_{\alpha j} \quad (16)$$

(see Figures 22 and 23). The forward and backward equations can be written in the following matrix forms

$$\mathbb{P}^{(k+1)} = \mathbb{P}^{(k)} \cdot \mathbb{P}, \quad (17)$$

$$\mathbb{P}^{(k+1)} = \mathbb{P} \cdot \mathbb{P}^{(k)}. \quad (18)$$

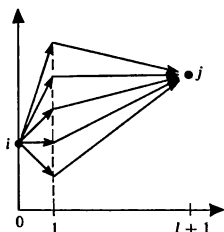


Figure 22. For the backward equation.

Similarly, we find for the (unconditional) probabilities  $p_j^{(k)}$  that

$$p_j^{(k+l)} = \sum_x p_x^{(k)} p_{xj}^{(l)}, \quad (19)$$

or in matrix form

$$\Pi^{(k+l)} = \Pi^{(k)} \cdot \Pi^{(l)}.$$

In particular,

$$\Pi^{(k+1)} = \Pi^{(k)} \cdot \Pi$$

(forward equation) and

$$\Pi^{(k+1)} = \Pi^{(1)} \cdot \Pi^{(k)}$$

(backward equation). Since  $\Pi^{(1)} = \Pi$ ,  $\Pi^{(1)} = \Pi$ , it follows from these equations that

$$\Pi^{(k)} = \Pi^k, \quad \Pi^{(k)} = \Pi^k.$$

Consequently for homogeneous Markov chains the  $k$ -step transition probabilities  $p_{ij}^{(k)}$  are the elements of the  $k$ th powers of the matrix  $\Pi$ , so that many properties of such chains can be investigated by the methods of matrix analysis.

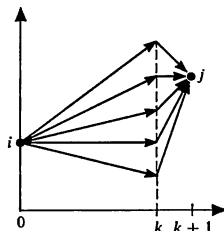


Figure 23. For the forward equation.

EXAMPLE 5. Consider a homogeneous Markov chain with the two states 0 and 1 and the matrix

$$\mathbb{P} = \begin{pmatrix} p_{00} & p_{01} \\ p_{10} & p_{11} \end{pmatrix}.$$

It is easy to calculate that

$$\mathbb{P}^2 = \begin{pmatrix} p_{00}^2 + p_{01}p_{10} & p_{01}(p_{00} + p_{11}) \\ p_{10}(p_{00} + p_{11}) & p_{11}^2 + p_{01}p_{10} \end{pmatrix}$$

and (by induction)

$$\begin{aligned} \mathbb{P}^n = & \frac{1}{2 - p_{00} - p_{11}} \begin{pmatrix} 1 - p_{11} & 1 - p_{00} \\ 1 - p_{11} & 1 - p_{00} \end{pmatrix} \\ & + \frac{(p_{00} + p_{11} - 1)^n}{2 - p_{00} - p_{11}} \begin{pmatrix} 1 - p_{00} & -(1 - p_{00}) \\ -(1 - p_{11}) & 1 - p_{11} \end{pmatrix} \end{aligned}$$

(under the hypothesis that  $|p_{00} + p_{11} - 1| < 1$ ).

Hence it is clear that if the elements of  $\mathbb{P}$  satisfy  $|p_{00} + p_{11} - 1| < 1$  (in particular, if all the transition probabilities  $p_{ij}$  are positive), then as  $n \rightarrow \infty$

$$\mathbb{P}^n \rightarrow \frac{1}{2 - p_{00} - p_{11}} \begin{pmatrix} 1 - p_{11} & 1 - p_{00} \\ 1 - p_{11} & 1 - p_{00} \end{pmatrix}, \quad (20)$$

and therefore

$$\lim_n p_{i0}^{(n)} = \frac{1 - p_{11}}{2 - p_{00} - p_{11}}, \quad \lim_n p_{i1}^{(n)} = \frac{1 - p_{00}}{2 - p_{00} - p_{11}}.$$

Consequently if  $|p_{00} + p_{11} - 1| < 1$ , such a Markov chain exhibits regular behavior of the following kind: the influence of the initial state on the probability of finding the particle in one state or another eventually becomes negligible ( $p_{ij}^{(n)}$  approach limits  $\pi_j$ , independent of  $i$  and forming a probability distribution:  $\pi_0 \geq 0$ ,  $\pi_1 \geq 0$ ,  $\pi_0 + \pi_1 = 1$ ); if also all  $p_{ij} > 0$  then  $\pi_0 > 0$  and  $\pi_1 > 0$ .

3. The following theorem describes a wide class of Markov chains that have the property called *ergodicity*: the limits  $\pi_j = \lim_n p_{ij}$  not only exist, are independent of  $i$ , and form a probability distribution ( $\pi_j \geq 0$ ,  $\sum_j \pi_j = 1$ ), but also  $\pi_j > 0$  for all  $j$  (such a distribution  $\pi_j$  is said to be *ergodic*).

**Theorem 1 (Ergodic Theorem).** Let  $\mathbb{P} = \|p_{ij}\|$  be the transition matrix of a chain with a finite state space  $X = \{1, 2, \dots, N\}$ .

(a) If there is an  $n_0$  such that

$$\min_{i,j} p_{ij}^{(n_0)} > 0, \quad (21)$$

then there are numbers  $\pi_1, \dots, \pi_N$  such that

$$\pi_j > 0, \quad \sum_j \pi_j = 1 \quad (22)$$

and

$$p_{ij}^{(n)} \rightarrow \pi_j, \quad n \rightarrow \infty \quad (23)$$

for every  $i \in X$ .

(b) Conversely, if there are numbers  $\pi_1, \dots, \pi_N$  satisfying (22) and (23), there is an  $n_0$  such that (21) holds.

(c) The numbers  $(\pi_1, \dots, \pi_N)$  satisfy the equations

$$\pi_j = \sum_{\alpha} \pi_{\alpha} p_{\alpha j}, \quad j = 1, \dots, N. \quad (24)$$

PROOF. (a) Let

$$m_j^{(n)} = \min_i p_{ij}^{(n)}, \quad M_j^{(n)} = \max_i p_{ij}^{(n)}.$$

Since

$$p_{ij}^{(n+1)} = \sum_{\alpha} p_{i\alpha} p_{\alpha j}^{(n)}, \quad (25)$$

we have

$$m_j^{(n+1)} = \min_i p_{ij}^{(n+1)} = \min_i \sum_{\alpha} p_{i\alpha} p_{\alpha j}^{(n)} \geq \min_i \sum_{\alpha} p_{i\alpha} \min_{\alpha} p_{\alpha j}^{(n)} = m_j^{(n)},$$

whence  $m_j^{(n)} \leq m_j^{(n+1)}$  and similarly  $M_j^{(n)} \geq M_j^{(n+1)}$ . Consequently, to establish (23) it will be enough to prove that

$$M_j^{(n)} - m_j^{(n)} \rightarrow 0, \quad n \rightarrow \infty, \quad j = 1, \dots, N.$$

Let  $\varepsilon = \min_{i,j} p_{ij}^{(n_0)} > 0$ . Then

$$\begin{aligned} p_{ij}^{(n_0+n)} &= \sum_{\alpha} p_{i\alpha}^{(n_0)} p_{\alpha j}^{(n)} = \sum_{\alpha} [p_{i\alpha}^{(n_0)} - \varepsilon p_{j\alpha}^{(n)}] p_{\alpha j}^{(n)} + \varepsilon \sum_{\alpha} p_{j\alpha}^{(n)} p_{\alpha j}^{(n)} \\ &= \sum_{\alpha} [p_{i\alpha}^{(n_0)} - \varepsilon p_{j\alpha}^{(n)}] p_{\alpha j}^{(n)} + \varepsilon p_{jj}^{(2n)}. \end{aligned}$$

But  $p_{i\alpha}^{(n_0)} - \varepsilon p_{j\alpha}^{(n)} \geq 0$ ; therefore

$$p_{ij}^{(n_0+n)} \geq m_j^{(n)} \cdot \sum_{\alpha} [p_{i\alpha}^{(n_0)} - \varepsilon p_{j\alpha}^{(n)}] + \varepsilon p_{jj}^{(2n)} = m_j^{(n)}(1 - \varepsilon) + \varepsilon p_{jj}^{(2n)},$$

and consequently

$$m_j^{(n_0+n)} \geq m_j^{(n)}(1 - \varepsilon) + \varepsilon p_{jj}^{(2n)}.$$

In a similar way

$$M_j^{(n_0+n)} \leq M_j^{(n)}(1 - \varepsilon) + \varepsilon p_{jj}^{(2n)}.$$

Combining these inequalities, we obtain

$$M_j^{(n_0+n)} - m_j^{(n_0+n)} \leq (M_j^{(n)} - m_j^{(n)}) \cdot (1 - \varepsilon)$$

and consequently

$$M_j^{(kn_0+n)} - m_j^{(kn_0+n)} \leq (M_j^{(n)} - m_j^{(n)})(1 - \varepsilon)^k \downarrow 0, \quad k \rightarrow \infty.$$

Thus  $M_j^{(n_{\beta})} - m_j^{(n_{\beta})} \rightarrow 0$  for some subsequence  $n_{\beta}$ ,  $n_{\beta} \rightarrow \infty$ . But the difference  $M_j^{(n)} - m_j^{(n)}$  is monotonic in  $n$ , and therefore  $M_j^{(n)} - m_j^{(n)} \rightarrow 0$ ,  $n \rightarrow \infty$ .

If we put  $\pi_j = \lim_n m_j^{(n)}$ , it follows from the preceding inequalities that

$$|p_{ij}^{(n)} - \pi_j| \leq M_j^{(n)} - m_j^{(n)} \leq (1 - \varepsilon)^{[n/n_0] - 1}$$

for  $n \geq n_0$ , that is,  $p_{ij}^{(n)}$  converges to its limit  $\pi_j$  geometrically (i.e., as fast as a geometric progression).

It is also clear that  $m_j^{(n)} \geq m_j^{(n_0)} \geq \varepsilon > 0$  for  $n \geq n_0$ , and therefore  $\pi_j > 0$ .

(b) Inequality (21) follows from (23) and (25).

(c) Equation (24) follows from (23) and (25).

This completes the proof of the theorem.

4. Equations (24) play a major role in the theory of Markov chains. A nonnegative solution  $(\pi_1, \dots, \pi_N)$  satisfying  $\sum_{\alpha} \pi_{\alpha} = 1$  is said to be a *stationary* or *invariant* probability distribution for the Markov chain with transition matrix  $\|p_{ij}\|$ . The reason for this terminology is as follows.

Let us select an initial distribution  $(\pi_1, \dots, \pi_N)$  and take  $p_j = \pi_j$ . Then

$$p_j^{(1)} = \sum_{\alpha} \pi_{\alpha} p_{\alpha j} = \pi_j$$

and in general  $p_j^{(n)} = \pi_j$ . In other words, if we take  $(\pi_1, \dots, \pi_N)$  as the initial distribution, this distribution is unchanged as time goes on, i.e. for any  $k$

$$P(\xi_k = j) = P(\xi_0 = j), \quad j = 1, \dots, N.$$

Moreover, with this initial distribution the Markov chain  $\xi = (\xi, \Pi, \mathbb{P})$  is really *stationary*: the joint distribution of the vector  $(\xi_k, \xi_{k+1}, \dots, \xi_{k+l})$  is independent of  $k$  for all  $l$  (assuming that  $k + l \leq n$ ).

Property (21) guarantees both the existence of limits  $\pi_j = \lim p_{ij}^{(n)}$ , which are independent of  $i$ , and the existence of an ergodic distribution, i.e. one with  $\pi_j > 0$ . The distribution  $(\pi_1, \dots, \pi_N)$  is also a *stationary* distribution. Let us now show that the set  $(\pi_1, \dots, \pi_N)$  is the *only* stationary distribution.

In fact, let  $(\tilde{\pi}_1, \dots, \tilde{\pi}_N)$  be another stationary distribution. Then

$$\tilde{\pi}_j = \sum_{\alpha} \tilde{\pi}_{\alpha} p_{\alpha j} = \dots = \sum_{\alpha} \tilde{\pi}_{\alpha} p_{\alpha j}^{(n)},$$

and since  $p_{\alpha j}^{(n)} \rightarrow \pi_j$  we have

$$\tilde{\pi}_j = \sum_{\alpha} (\tilde{\pi}_{\alpha} \cdot \pi_j) = \pi_j.$$

These problems will be investigated in detail in Chapter VIII for Markov chains with countably many states as well as with finitely many states.

We note that a stationary probability distribution (even unique) may exist for a nonergodic chain. In fact, if

$$\mathbb{P} = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix},$$

then

$$\mathbb{P}^{2n} = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}, \quad \mathbb{P}^{2n+1} = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix},$$

and consequently the limits  $\lim p_{ij}^{(n)}$  do not exist. At the same time, the system

$$\pi_j = \sum_{\alpha} \pi_{\alpha} p_{\alpha j}, \quad j = 1, 2,$$

reduces to

$$\pi_1 = \pi_2,$$

$$\pi_2 = \pi_1,$$

of which the unique solution satisfying  $\pi_1 + \pi_2 = 1$  is  $(\frac{1}{2}, \frac{1}{2})$ .

We also notice that for this example the system (24) has the form

$$\pi_0 = \pi_0 p_{00} + \pi_1 p_{10},$$

$$\pi_1 = \pi_0 p_{01} + \pi_1 p_{11},$$

from which, by the condition  $\pi_0 + \pi_1 = 1$ , we find that the unique stationary distribution  $(\pi_0, \pi_1)$  coincides with the one obtained above:

$$\pi_0 = \frac{1 - p_{11}}{2 - p_{00} - p_{11}}, \quad \pi_1 = \frac{1 - p_{00}}{2 - p_{00} - p_{11}}.$$

We now consider some corollaries of the ergodic theorem.

Let  $A$  be a set of states,  $A \subseteq X$  and

$$I_A(x) = \begin{cases} 1, & x \in A, \\ 0, & x \notin A. \end{cases}$$

Consider

$$v_A(n) = \frac{I_A(\xi_0) + \cdots + I_A(\xi_n)}{n+1}$$

which is the fraction of the time spent by the particle in the set  $A$ . Since

$$\mathbb{E}[I_A(\xi_k) | \xi_0 = i] = \mathbb{P}(\xi_k \in A | \xi_0 = i) = \sum_{j \in A} p_{ij}^{(k)} (= p_i^{(k)}(A)),$$

we have

$$\mathbb{E}[v_A(n)|\xi_0 = i] = \frac{1}{n+1} \sum_{k=0}^n p_i^{(k)}(A)$$

and in particular

$$\mathbb{E}[v_{\{j\}}(n)|\xi_0 = i] = \frac{1}{n+1} \sum_{k=0}^n p_{ij}^{(k)}.$$

It is known from analysis (see also Lemma 1 in §3 of Chapter IV) that if  $a_n \rightarrow a$  then  $(a_0 + \dots + a_n)/(n+1) \rightarrow a$ ,  $n \rightarrow \infty$ . Hence if  $p_{ij}^{(k)} \rightarrow \pi_j$ ,  $k \rightarrow \infty$ , then

$$\mathbb{E}v_{\{j\}}(n) \rightarrow \pi_j, \quad \mathbb{E}v_A(n) \rightarrow \pi_A, \quad \text{where } \pi_A = \sum_{j \in A} \pi_j.$$

For ergodic chains one can in fact prove more, namely that the following result holds for  $I_A(\xi_0), \dots, I_A(\xi_n), \dots$ .

**Law of Large Numbers.** *If  $\xi_0, \xi_1, \dots$  form a finite ergodic Markov chain, then*

$$\mathbb{P}\{|v_A(n) - \pi_A| > \varepsilon\} \rightarrow 0, \quad n \rightarrow \infty, \quad (26)$$

*for every  $\varepsilon > 0$  and every initial distribution.*

Before we undertake the proof, let us notice that we cannot apply the results of §5 directly to  $I_A(\xi_0), \dots, I_A(\xi_n), \dots$ , since these variables are, in general, dependent. However, the proof can be carried through along the same lines as for independent variables if we again use Chebyshev's inequality, and apply the fact that for an ergodic chain with finitely many states there is a number  $\rho$ ,  $0 < \rho < 1$ , such that

$$|p_{ij}^{(n)} - \pi_j| \leq C \cdot \rho^n. \quad (27)$$

Let us consider states  $i$  and  $j$  (which might be the same) and show that, for  $\varepsilon > 0$ ,

$$\mathbb{P}\{|v_{\{j\}}(n) - \pi_j| > \varepsilon | \xi_0 = i\} \rightarrow 0, \quad n \rightarrow \infty. \quad (28)$$

By Chebyshev's inequality,

$$\mathbb{P}\{|v_{\{j\}}(n) - \pi_j| > \varepsilon | \xi_0 = i\} < \frac{\mathbb{E}\{|v_{\{j\}}(n) - \pi_j|^2 | \xi_0 = i\}}{\varepsilon^2}.$$

Hence we have only to show that

$$\mathbb{E}\{|v_{\{j\}}(n) - \pi_j|^2 | \xi_0 = i\} \rightarrow 0, \quad n \rightarrow \infty.$$

A simple calculation shows that

$$\begin{aligned} \mathbf{E}\{|v_{ij}(n) - \pi_j|^2 | \xi_0 = i\} &= \frac{1}{(n+1)^2} \cdot \mathbf{E}\left\{\left[\sum_{k=0}^n (I_{ij}(\xi_k) - \pi_j)\right]^2 \middle| \xi_0 = i\right\} \\ &= \frac{1}{(n+1)^2} \sum_{k=0}^n \sum_{l=0}^n m_{ij}^{(k,l)}, \end{aligned}$$

where

$$\begin{aligned} m_{ij}^{(k,l)} &= \mathbf{E}\{[I_{ij}(\xi_k)I_{ij}(\xi_l)] | \xi_0 = i\} \\ &\quad - \pi_j \cdot \mathbf{E}[I_{ij}(\xi_k) | \xi_0 = i] - \pi_j \cdot \mathbf{E}[I_{ij}(\xi_l) | \xi_0 = l] + \pi_j^2 \\ &= p_{ij}^{(s)} \cdot p_{ij}^{(t)} - \pi_j \cdot p_{ij}^{(k)} - \pi_j \cdot p_{ij}^{(l)} + \pi_j^2, \\ &\quad s = \min(k, l) \quad \text{and} \quad t = |k - l|. \end{aligned}$$

By (27),

$$p_{ij}^{(n)} = \pi_j + \varepsilon_{ij}^{(n)}, \quad |\varepsilon_{ij}^{(n)}| \leq C\rho^n.$$

Therefore

$$|m_{ij}^{(k,l)}| \leq C_1[\rho^s + \rho^t + \rho^k + \rho^l],$$

where  $C_1$  is a constant. Consequently

$$\begin{aligned} \frac{1}{(n+1)^2} \sum_{k=0}^n \sum_{l=0}^n m_{ij}^{(k,l)} &\leq \frac{C_1}{(n+1)^2} \sum_{k=0}^n \sum_{l=0}^n [\rho^s + \rho^t + \rho^k + \rho^l] \\ &\leq \frac{4C_1}{(n+1)^2} \cdot \frac{2(n+1)}{1-\rho} = \frac{8C_1}{(n+1)(1-\rho)} \rightarrow 0, \quad n \rightarrow \infty. \end{aligned}$$

Then (28) follows from this, and we obtain (26) in an obvious way.

5. In §9 we gave, for a random walk  $S_0, S_1, \dots$  generated by a Bernoulli scheme, recurrent equations for the probability and the expectation of the exit time at either boundary. We now derive similar equations for Markov chains.

Let  $\xi = (\xi_0, \dots, \xi_n)$  be a Markov chain with transition matrix  $\|p_{ij}\|$  and phase space  $X = \{0, \pm 1, \dots, \pm N\}$ . Let  $A$  and  $B$  be two integers,  $-N \leq A \leq 0 \leq B \leq N$ , and  $x \in X$ . Let  $\mathcal{B}_{k+1}$  be the set of paths  $(x_0, x_1, \dots, x_k)$ ,  $x_i \in X$ , that leave the interval  $(A, B)$  for the first time at the upper end, i.e. leave  $(A, B)$  by going into the set  $(B, B+1, \dots, N)$ .

For  $A \leq x \leq B$ , put

$$\beta_k(x) = \mathbf{P}\{(\xi_0, \dots, \xi_k) \in \mathcal{B}_{k+1} | \xi_0 = x\}.$$

In order to find these probabilities (for the first exit of the Markov chain from  $(A, B)$  through the upper boundary) we use the method that was applied in the deduction of the backward equations.



We have

$$\begin{aligned}\beta_k(x) &= \mathbf{P}\{(\xi_0, \dots, \xi_k) \in \mathcal{B}_{k+1} | \xi_0 = x\} \\ &= \sum_y p_{xy} \cdot \mathbf{P}\{(\xi_0, \dots, \xi_k) \in \mathcal{B}_{k+1} | \xi_0 = x, \xi_1 = y\},\end{aligned}$$

where, as is easily seen by using the Markov property and the homogeneity of the chain,

$$\begin{aligned}\mathbf{P}\{(\xi_0, \dots, \xi_k) \in \mathcal{B}_{k+1} | \xi_0 = x, \xi_1 = y\} \\ &= \mathbf{P}\{(x, y, \xi_2, \dots, \xi_k) \in \mathcal{B}_{k+1} | \xi_0 = x, \xi_1 = y\} \\ &= \mathbf{P}\{(y, \xi_2, \dots, \xi_k) \in \mathcal{B}_k | \xi_1 = y\} \\ &= \mathbf{P}\{(y, \xi_1, \dots, \xi_{k-1}) \in \mathcal{B}_k | \xi_0 = y\} = \beta_{k-1}(y).\end{aligned}$$

Therefore

$$\beta_k(x) = \sum_y p_{xy} \beta_{k-1}(y)$$

for  $A < x < B$  and  $1 \leq k \leq n$ . Moreover, it is clear that

$$\beta_k(x) = 1, \quad x = B, B+1, \dots, N,$$

and

$$\beta_k(x) = 0, \quad x = -N, \dots, A.$$

In a similar way we can find equations for  $\alpha_k(x)$ , the probabilities for first exit from  $(A, B)$  through the lower boundary.

Let  $\tau_k = \min\{0 \leq l \leq k: \xi_l \notin (A, B)\}$ , where  $\tau_k = k$  if the set  $\{\cdot\} = \emptyset$ . Then the same method, applied to  $m_k(x) = \mathbf{E}(\tau_k | \xi_0 = x)$ , leads to the following recurrent equations:

$$m_k(x) = 1 + \sum_y m_{k-1}(y) p_{xy}$$

(here  $1 \leq k \leq n$ ,  $A < x < B$ ). We define

$$m_k(x) = 0, \quad x \notin (A, B).$$

It is clear that if the transition matrix is given by (11) the equations for  $\alpha_k(x)$ ,  $\beta_k(x)$  and  $m_k(x)$  become the corresponding equations from §9, where they were obtained by essentially the same method that was used here.

These equations have the most interesting applications in the limiting case when the walk continues for an unbounded length of time. Just as in §9, the corresponding equations can be obtained by a formal limiting process ( $k \rightarrow \infty$ ).

By way of example, we consider the Markov chain with states  $\{0, 1, \dots, B\}$  and transition probabilities

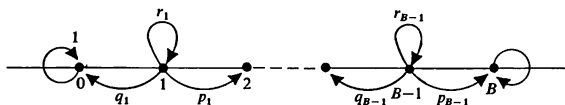
$$p_{00} = 1, \quad p_{BB} = 1,$$

and

$$p_{ij} = \begin{cases} p_i > 0, & j = i + 1, \\ r_i, & j = i, \\ q_i > 0, & j = i - 1, \end{cases}$$

for  $1 \leq i \leq B - 1$ , where  $p_i + q_i + r_i = 1$ .

For this chain, the corresponding graph is



It is clear that states 0 and  $B$  are absorbing, whereas for every other state  $i$  the particle stays there with probability  $r_i$ , moves one step to the right with probability  $p_i$ , and to the left with probability  $q_i$ .

Let us find  $\alpha(x) = \lim_{k \rightarrow \infty} \alpha_k(x)$ , the limit of the probability that a particle starting at the point  $x$  arrives at state zero before reaching state  $B$ . Taking limits as  $k \rightarrow \infty$  in the equations for  $\alpha_k(x)$ , we find that

$$\alpha(j) = q_j \alpha(j-1) + r_j \alpha(j) + p_j \alpha(j+1)$$

when  $0 < j < B$ , with the boundary conditions

$$\alpha(0) = 1, \quad \alpha(B) = 0.$$

Since  $r_j = 1 - q_j - p_j$ , we have

$$p_j(\alpha(j+1) - \alpha(j)) = q_j(\alpha(j) - \alpha(j-1))$$

and consequently

$$\alpha(j+1) - \alpha(j) = \rho_j(\alpha(1) - 1),$$

where

$$\rho_j = \frac{q_1 \cdots q_j}{p_1 \cdots p_j}, \quad \rho_0 = 1.$$

But

$$\alpha(j+1) - 1 = \sum_{i=0}^j (\alpha(i+1) - \alpha(i)).$$

Therefore

$$\alpha(j+1) - 1 = (\alpha(1) - 1) \cdot \sum_{i=0}^j \rho_i.$$

If  $j = B - 1$ , we have  $\alpha(j + 1) = \alpha(B) = 0$ , and therefore

$$\alpha(1) = 1 = -\frac{1}{\sum_{i=1}^{B-1} \rho_i},$$

whence

$$\alpha(1) = \frac{\sum_{i=1}^{B-1} \rho_i}{\sum_{i=0}^{B-1} \rho_i} \quad \text{and} \quad \alpha(j) = \frac{\sum_{i=j}^{B-1} \rho_i}{\sum_{i=1}^{B-1} \rho_i}, \quad j = 1, \dots, B.$$

(This should be compared with the results of §9.)

Now let  $m(x) = \lim_k m_k(x)$ , the limiting value of the average time taken to arrive at one of the states 0 or  $B$ . Then  $m(0) = m(B) = 0$ ,

$$m(x) = 1 + \sum_y m(y) p_{xy}$$

and consequently for the example that we are considering,

$$m(j) = 1 + q_j m(j-1) + r_j m(j) + p_j m(j+1)$$

for  $j = 1, 2, \dots, B-1$ . To find  $m(j)$  we put

$$M(j) = m(j) - m(j-1), \quad j = 0, 1, \dots, B.$$

Then

$$p_j M(j+1) = q_j M(j) - 1, \quad j = 1, \dots, B-1,$$

and consequently we find that

$$M(j+1) = \rho_j M(1) - R_j,$$

where

$$\rho_j = \frac{q_1 \cdots q_j}{p_1 \cdots p_j}, \quad R_j = \frac{1}{p_j} \left[ 1 + \frac{q_j}{p_{j-1}} + \cdots + \frac{q_j \cdots q_2}{p_j \cdots p_1} \right].$$

Therefore

$$\begin{aligned} m(i) = m(j) - m(0) &= \sum_{i=0}^{j-1} M(i+1) \\ &= \sum_{i=0}^{j-1} (\rho_i m(1) - R_i) = m(1) \sum_{i=0}^{j-1} \rho_i - \sum_{i=0}^{j-1} R_i. \end{aligned}$$

It remains only to determine  $m(1)$ . But  $m(B) = 0$ , and therefore

$$m(1) = \frac{\sum_{i=0}^{B-1} R_i}{\sum_{i=0}^{B-1} \rho_i},$$

and for  $1 < j \leq B$ ,

$$m(j) = \sum_{i=0}^{j-1} \rho_i \cdot \frac{\sum_{i=0}^{B-1} R_i}{\sum_{i=0}^{B-1} \rho_i} - \sum_{i=0}^{j-1} R_i.$$

(This should be compared with the results in §9 for the case  $r_i = 0$ ,  $p_i = p$ ,  $q_i = q$ .)

6. In this subsection we consider a stronger version of the Markov property (8), namely that it remains valid if time  $k$  is replaced by a random time (see also Theorem 2). The significance of this, the *strong Markov property*, can be illustrated in particular by the example of the derivation of the recurrent relations (38), which play an important role in the classification of the states of Markov chains (Chapter VIII).

Let  $\xi = (\xi_1, \dots, \xi_n)$  be a homogeneous Markov chain with transition matrix  $\|p_{ij}\|$ ; let  $\mathcal{D}^\xi = (\mathcal{D}_k^\xi)_{0 \leq k \leq n}$  be a system of decompositions,  $\mathcal{D}_k^\xi = \mathcal{D}_{\xi_0, \dots, \xi_k}$ . Let  $\mathcal{B}_k^\xi$  denote the algebra  $\alpha(\mathcal{D}_k^\xi)$  generated by the decomposition  $\mathcal{D}_k^\xi$ .

We first put the Markov property (8) into a somewhat different form. Let  $B \in \mathcal{B}_k^\xi$ . Let us show that then

$$\begin{aligned} \mathbf{P}\{\xi_n = a_n, \dots, \xi_{k+1} = a_{k+1} | B \cap (\xi_k = a_k)\} \\ = \mathbf{P}\{\xi_n = a_n, \dots, \xi_{k+1} = a_{k+1} | \xi_k = a_k\} \end{aligned} \quad (29)$$

(assuming that  $\mathbf{P}\{B \cap (\xi_k = a_k)\} > 0$ ). In fact,  $B$  can be represented in the form

$$B = \sum^* \{\xi_0 = a_0^*, \dots, \xi_k = a_k^*\},$$

where  $\sum^*$  extends over some set  $(a_0^*, \dots, a_k^*)$ . Consequently

$$\begin{aligned} \mathbf{P}\{\xi_n = a_n, \dots, \xi_{k+1} = a_{k+1} | B \cap (\xi_k = a_k)\} \\ = \frac{\mathbf{P}\{(\xi_n = a_n, \dots, \xi_k = a_k) \cap B\}}{\mathbf{P}\{(\xi_k = a_k) \cap B\}} \\ = \frac{\sum^* \mathbf{P}\{(\xi_n = a_n, \dots, \xi_k = a_k) \cap (\xi_0 = a_0^*, \dots, \xi_k = a_k^*)\}}{\mathbf{P}\{(\xi_k = a_k) \cap B\}}. \end{aligned} \quad (30)$$

But, by the Markov property,

$$\begin{aligned} \mathbf{P}\{(\xi_n = a_n, \dots, \xi_k = a_k) \cap (\xi_0 = a_0^*, \dots, \xi_k = a_k^*)\} \\ = \begin{cases} \mathbf{P}\{\xi_n = a_n, \dots, \xi_{k+1} = a_{k+1} | \xi_0 = a_0^*, \dots, \xi_k = a_k^*\} \\ \quad \times \mathbf{P}\{\xi_0 = a_0^*, \dots, \xi_k = a_k^*\} & \text{if } a_k = a_k^*, \\ 0 & \text{if } a_k \neq a_k^*, \end{cases} \\ = \begin{cases} \mathbf{P}\{\xi_n = a_n, \dots, \xi_{k+1} = a_{k+1} | \xi_k = a_k\} \mathbf{P}\{\xi_0 = a_0^*, \dots, \xi_k = a_k^*\} \\ & \text{if } a_k = a_k^*, \\ 0 & \text{if } a_k \neq a_k^*, \end{cases} \\ = \begin{cases} \mathbf{P}\{\xi_n = a_n, \dots, \xi_{k+1} = a_{k+1} | \xi_k = a_k\} \mathbf{P}\{(\xi_k = a_k) \cap B\} \\ & \text{if } a_k = a_k^*, \\ 0 & \text{if } a_k \neq a_k^*. \end{cases} \end{aligned}$$

Therefore the sum  $\sum^*$  in (30) is equal to

$$P\{\xi_n = a_n, \dots, \xi_{k+1} = a_{k+1} | \xi_k = a_k\} P\{(\xi_k = a_k) \cap B\},$$

This establishes (29).

Let  $\tau$  be a stopping time (with respect to the system  $D^\xi = (D_k^\xi)_{0 \leq k \leq n}$ ; see Definition 2 in §11).

**Definition.** We say that a set  $B$  in the algebra  $\mathcal{B}_n^\xi$  belongs to the system of sets  $\mathcal{B}_\tau^\xi$  if, for each  $k$ ,  $0 \leq k \leq n$ ,

$$B \cap \{\tau = k\} \in \mathcal{B}_k^\xi. \quad (31)$$

It is easily verified that the collection of such sets  $B$  forms an algebra (called the algebra of events observed at time  $\tau$ ).

**Theorem 2.** Let  $\xi = (\xi_0, \dots, \xi_n)$  be a homogeneous Markov chain with transition matrix  $\|p_{ij}\|$ ,  $\tau$  a stopping time (with respect to  $\mathcal{D}^\xi$ ),  $B \in \mathcal{B}_\tau^\xi$  and  $A = \{\omega: \tau + l \leq n\}$ . Then if  $P\{A \cap B \cap (\xi_\tau = a_0)\} > 0$ , we have

$$\begin{aligned} P\{\xi_{\tau+l} = a_l, \dots, \xi_{\tau+1} = a_1 | A \cap B \cap (\xi_\tau = a_0)\} \\ = P\{\xi_{\tau+l} = a_l, \dots, \xi_{\tau+1} = a_1 | A \cap (\xi_\tau = a_0)\}, \end{aligned} \quad (32)$$

and if  $P\{A \cap (\xi_\tau = a_0)\} > 0$  then

$$P\{\xi_{\tau+l} = a_l, \dots, \xi_{\tau+1} = a_1 | A \cap (\xi_\tau = a_0)\} = p_{a_0 a_1} \dots p_{a_{l-1} a_l}. \quad (33)$$

For the sake of simplicity, we give the proof only for the case  $l = 1$ . Since  $B \cap (\tau = k) \in \mathcal{B}_k^\xi$ , we have, according to (29),

$$\begin{aligned} P\{\xi_{\tau+1} = a_1, A \cap B \cap (\xi_\tau = a_0)\} \\ = \sum_{k \leq n-1} P\{\xi_{k+1} = a_1, \xi_k = a_0, \tau = k, B\} \\ = \sum_{k \leq n-1} P\{\xi_{k+1} = a_1 | \xi_k = a_0, \tau = k, B\} P\{\xi_k = a_0, \tau = k, B\} \\ = \sum_{k \leq n-1} P\{\xi_{k+1} = a_1 | \xi_k = a_0\} P\{\xi_k = a_0, \tau = k, B\} \\ = p_{a_0 a_1} \cdot \sum_{k \leq n-1} P\{\xi_k = a_0, \tau = k, B\} = p_{a_0 a_1} \cdot P\{A \cap B \cap (\xi_\tau = a_0)\}, \end{aligned}$$

which simultaneously establishes (32) and (33) (for (33) we have to take  $B = \Omega$ ).

**Remark.** When  $l = 1$ , the strong Markov property (32), (33) is evidently equivalent to the property that

$$P\{\xi_{\tau+1} \in C | A \cap B \cap (\xi_\tau = a_0)\} = P_{a_0}(C), \quad (34)$$

for every  $C \subseteq X$ , where

$$P_{a_0}(C) = \sum_{a_1 \in C} p_{a_0 a_1}.$$

In turn, (34) can be restated as follows: on the set  $A = \{\tau \leq n-1\}$ ,

$$P\{\xi_{\tau+1} \in C | \mathcal{B}_\tau^\xi\} = P_{\xi_\tau}(C), \quad (35)$$

which is a form of the strong Markov property that is commonly used in the general theory of homogeneous Markov processes.

7. Let  $\xi = (\xi_0, \dots, \xi_n)$  be a homogeneous Markov chain with transition matrix  $\|p_{ij}\|$ , and let

$$f_{ii}^{(k)} = P\{\xi_k = i, \xi_l \neq i, 1 \leq l \leq k-1 | \xi_0 = i\} \quad (36)$$

and

$$f_{ij}^{(k)} = P\{\xi_k = j, \xi_l \neq j, 1 \leq l \leq k-1 | \xi_0 = i\} \quad (37)$$

for  $i \neq j$  be respectively the probability of first return to state  $i$  at time  $k$  and the probability of first arrival at state  $j$  at time  $k$ .

Let us show that

$$p_{ij}^{(n)} = \sum_{k=1}^n f_{ij}^{(k)} p_{jj}^{(n-k)}, \quad \text{where } p_{jj}^{(0)} = 1. \quad (38)$$

The intuitive meaning of the formula is clear: to go from state  $i$  to state  $j$  in  $n$  steps, it is necessary to reach state  $j$  for the first time in  $k$  steps ( $1 \leq k \leq n$ ) and then to go from state  $j$  to state  $j$  in  $n-k$  steps. We now give a rigorous derivation.

Let  $j$  be given and

$$\tau = \min\{1 \leq k \leq n: \xi_k = j\},$$

assuming that  $\tau = n+1$  if  $\{\cdot\} = \emptyset$ . Then  $f_{ij}^{(k)} = P\{\tau = k | \xi_0 = i\}$  and

$$\begin{aligned} p_{ij}^{(n)} &= P\{\xi_n = j | \xi_0 = i\} \\ &= \sum_{1 \leq k \leq n} P\{\xi_n = j, \tau = k | \xi_0 = i\} \\ &= \sum_{1 \leq k \leq n} P\{\xi_{\tau+n-k} = j, \tau = k | \xi_0 = i\}, \end{aligned} \quad (39)$$

where the last equation follows because  $\xi_{\tau+n-k} = \xi_n$  on the set  $\{\tau = k\}$ . Moreover, the set  $\{\tau = k\} = \{\tau = k, \xi_\tau = j\}$  for every  $k, 1 \leq k \leq n$ . Therefore if  $P\{\xi_0 = i, \tau = k\} > 0$ , it follows from Theorem 2 that

$$\begin{aligned} P\{\xi_{\tau+n-k} = j | \xi_0 = i, \tau = k\} &= P\{\xi_{\tau+n-k} = j | \xi_0 = i, \tau = k, \xi_\tau = j\} \\ &= P\{\xi_{\tau+n-k} = j | \xi_\tau = j\} = p_{jj}^{(n-k)} \end{aligned}$$

and by (37)

$$\begin{aligned} p_{ij}^{(n)} &= \sum_{k=1}^n \mathbf{P}\{\xi_{\tau+n-k} = j | \xi_0 = i, \tau = k\} \mathbf{P}\{\tau = k | \xi_0 = i\} \\ &= \sum_{k=1}^n p_{jj}^{(n-k)} f_{ij}^{(k)}, \end{aligned}$$

which establishes (38).

## 8. PROBLEMS

1. Let  $\xi = (\xi_0, \dots, \xi_n)$  be a Markov chain with values in  $X$  and  $f = f(x)$  ( $x \in X$ ) a function. Will the sequence  $(f(\xi_0), \dots, f(\xi_n))$  form a Markov chain? Will the "reversed" sequence

$$(\xi_n, \xi_{n-1}, \dots, \xi_0)$$

form a Markov chain?

2. Let  $\mathbb{P} = \|p_{ij}\|$ ,  $1 \leq i, j \leq r$ , be a stochastic matrix and  $\lambda$  an eigenvalue of the matrix, i.e. a root of the characteristic equation  $\det\|\mathbb{P} - \lambda E\| = 0$ . Show that  $\lambda_0 = 1$  is an eigenvalue and that all the other eigenvalues have moduli not exceeding 1. If all the eigenvalues  $\lambda_1, \dots, \lambda_r$  are distinct, then  $p_{ij}^{(k)}$  admits the representation

$$p_{ij}^{(k)} = \pi_j + a_{ij}(1)\lambda_1^k + \dots + a_{ij}(r)\lambda_r^k,$$

where  $\pi_j, a_{ij}(1), \dots, a_{ij}(r)$  can be expressed in terms of the elements of  $\mathbb{P}$ . (It follows from this algebraic approach to the study of Markov chains that, in particular, when  $|\lambda_1| < 1, \dots, |\lambda_r| < 1$ , the limit  $\lim p_{ij}^{(k)}$  exists for every  $j$  and is independent of  $i$ .)

3. Let  $\xi = (\xi_0, \dots, \xi_n)$  be a homogeneous Markov chain with state space  $X$  and transition matrix  $\mathbb{P} = \|p_{xy}\|$ . Let

$$T\varphi(x) = \mathbf{E}[\varphi(\xi_1) | \xi_0 = x] \quad \left( = \sum_y \varphi(y) p_{xy} \right).$$

Let the nonnegative function  $\varphi$  satisfy

$$T\varphi(x) = \varphi(x), \quad x \in X.$$

Show that the sequence of random variables

$$\zeta = (\zeta_k, \mathcal{D}_k^{\zeta}) \quad \text{with} \quad \zeta_k = \varphi(\xi_k)$$

is a martingale.

4. Let  $\xi = (\xi_n, \Pi, \mathbb{P})$  and  $\tilde{\xi} = (\xi_n, \tilde{\Pi}, \tilde{\mathbb{P}})$  be two Markov chains with different initial distributions  $\Pi = (p_1, \dots, p_r)$  and  $\tilde{\Pi} = (\tilde{p}_1, \dots, \tilde{p}_r)$ . Show that if  $\min_{i,j} p_{ij} \geq \varepsilon > 0$  then

$$\sum_{i=1}^r |\tilde{p}_i^{(n)} - p_i^{(n)}| \leq 2(1 - \varepsilon)^n.$$

## CHAPTER II

# Mathematical Foundations of Probability Theory

### §1. Probabilistic Model for an Experiment with Infinitely Many Outcomes. Kolmogorov's Axioms

1. The models introduced in the preceding chapter enabled us to give a probabilistic-statistical description of experiments with a finite number of outcomes. For example, the triple  $(\Omega, \mathcal{A}, \mathbf{P})$  with

$$\Omega = \{\omega: \omega = (a_1, \dots, a_n), a_i = 0, 1\}, \mathcal{A} = \{A: A \subseteq \Omega\}$$

and  $p(\omega) = p^{\sum a_i} q^{n - \sum a_i}$  is a model for the experiment in which a coin is tossed  $n$  times "independently" with probability  $p$  of falling head. In this model the number  $N(\Omega)$  of outcomes, i.e. the number of points in  $\Omega$ , is the finite number  $2^n$ .

We now consider the problem of constructing a probabilistic model for the experiment consisting of an infinite number of independent tosses of a coin when at each step the probability of falling head is  $p$ .

It is natural to take the set of outcomes to be the set

$$\Omega = \{\omega: \omega = (a_1, a_2, \dots), a_i = 0, 1\},$$

i.e. the space of sequences  $\omega = (a_1, a_2, \dots)$  whose elements are 0 or 1.

What is the cardinality  $N(\Omega)$  of  $\Omega$ ? It is well known that every number  $a \in [0, 1)$  has a unique binary expansion (containing an infinite number of zeros)

$$a = \frac{a_1}{2} + \frac{a_2}{2^2} + \dots \quad (a_i = 0, 1).$$



Hence it is clear that there is a one-to-one correspondence between the points  $\omega$  of  $\Omega$  and the points  $a$  of the set  $[0, 1)$ , and therefore  $\Omega$  has the cardinality of the continuum.

Consequently if we wish to construct a probabilistic model to describe experiments like tossing a coin infinitely often, we must consider spaces  $\Omega$  of a rather complicated nature.

We shall now try to see what probabilities ought reasonably to be assigned (or assumed) in a model of infinitely many independent tosses of a fair coin ( $p + q = \frac{1}{2}$ ).

Since we may take  $\Omega$  to be the set  $[0, 1)$ , our problem can be considered as the problem of choosing points at random from this set. For reasons of symmetry, it is clear that all outcomes ought to be equiprobable. But the set  $[0, 1)$  is uncountable, and if we suppose that its probability is 1, then it follows that the probability  $p(\omega)$  of each outcome certainly must equal zero. However, this assignment of probabilities ( $p(\omega) = 0, \omega \in [0, 1)$ ) does not lead very far. The fact is that we are ordinarily not interested in the probability of one outcome or another, but in the probability that the result of the experiment is in one or another specified set  $A$  of outcomes (an event). In elementary probability theory we use the probabilities  $p(\omega)$  to find the probability  $P(A)$  of the event  $A$ :  $P(A) = \sum_{\omega \in A} p(\omega)$ . In the present case, with  $p(\omega) = 0, \omega \in [0, 1)$ , we cannot define, for example, the probability that a point chosen at random from  $[0, 1)$  belongs to the set  $[0, \frac{1}{2})$ . At the same time, it is intuitively clear that this probability should be  $\frac{1}{2}$ .

These remarks should suggest that in constructing probabilistic models for uncountable spaces  $\Omega$  we must assign probabilities, not to individual outcomes but to subsets of  $\Omega$ . The same reasoning as in the first chapter shows that the collection of sets to which probabilities are assigned must be closed with respect to unions, intersections and complements. Here the following definition is useful.

**Definition 1.** Let  $\Omega$  be a set of points  $\omega$ . A system  $\mathcal{A}$  of subsets of  $\Omega$  is called an *algebra* if

- (a)  $\Omega \in \mathcal{A}$ ,
- (b)  $A, B \in \mathcal{A} \Rightarrow A \cup B \in \mathcal{A}, \quad A \cap B \in \mathcal{A}$ ,
- (c)  $A \in \mathcal{A} \Rightarrow \bar{A} \in \mathcal{A}$

(Notice that in condition (b) it is sufficient to require only that either  $A \cup B \in \mathcal{A}$  or that  $A \cap B \in \mathcal{A}$ , since  $A \cup B = \overline{\bar{A} \cap \bar{B}}$  and  $A \cap B = \overline{\bar{A} \cup \bar{B}}$ .)

The next definition is needed in formulating the concept of a probabilistic model.

**Definition 2.** Let  $\mathcal{A}$  be an algebra of subsets of  $\Omega$ . A set function  $\mu = \mu(A)$ ,  $A \in \mathcal{A}$ , taking values in  $[0, \infty]$ , is called a *finitely additive measure* defined

on  $\mathcal{A}$  if

$$\mu(A + B) = \mu(A) + \mu(B). \quad (1)$$

for every pair of disjoint sets  $A$  and  $B$  in  $\mathcal{A}$ .

A finitely additive measure  $\mu$  with  $\mu(\Omega) < \infty$  is called *finite*, and when  $\mu(\Omega) = 1$  it is called a *finitely additive probability measure*, or a *finitely additive probability*.

2. We now define a probabilistic model (in the extended sense).

**Definition 3.** An ordered triple  $(\Omega, \mathcal{A}, \mathbf{P})$ , where

- (a)  $\Omega$  is a set of points  $\omega$ ;
- (b)  $\mathcal{A}$  is an algebra of subsets of  $\Omega$ ;
- (c)  $\mathbf{P}$  is a finitely additive probability on  $\mathcal{A}$ ,

is a *probabilistic model in the extended sense*.

It turns out, however, that this model is too broad to lead to a fruitful mathematical theory. Consequently we must restrict both the class of subsets of  $\Omega$  that we consider, and the class of admissible probability measures.

**Definition 4.** A system  $\mathcal{F}$  of subsets of  $\Omega$  is a  $\sigma$ -*algebra* if it is an algebra and satisfies the following additional condition (stronger than (b) of Definition 1):

- (b\*) if  $A_n \in \mathcal{F}$ ,  $n = 1, 2, \dots$ , then

$$\bigcup A_n \in \mathcal{F}, \quad \bigcap A_n \in \mathcal{F}$$

(it is sufficient to require either that  $\bigcup A_n \in \mathcal{F}$  or that  $\bigcap A_n \in \mathcal{F}$ ).

**Definition 5.** The space  $\Omega$  together with a  $\sigma$ -algebra  $\mathcal{F}$  of its subsets is a *measurable space*, and is denoted by  $(\Omega, \mathcal{F})$ .

**Definition 6.** A finitely additive measure  $\mu$  defined on an algebra  $\mathcal{A}$  of subsets of  $\Omega$  is *countably additive* (or  $\sigma$ -*additive*), or simply a *measure*, if, for all pairwise disjoint subsets  $A_1, A_2, \dots$  of  $\mathcal{A}$  with  $\sum A_n \in \mathcal{A}$

$$\mu\left(\sum_{n=1}^{\infty} A_n\right) = \sum_{n=1}^{\infty} \mu(A_n).$$

A finitely additive measure  $\mu$  is said to be  $\sigma$ -*finite* if  $\Omega$  can be represented in the form

$$\Omega = \sum_{n=1}^{\infty} \Omega_n, \quad \Omega_n \in \mathcal{A},$$

with  $\mu(\Omega_n) < \infty$ ,  $n = 1, 2, \dots$ .

If a countably additive measure  $P$  on the algebra  $\mathcal{A}$  satisfies  $P(\Omega) = 1$ , it is called a *probability measure* or a *probability* (defined on the sets that belong to the algebra  $\mathcal{A}$ ).

Probability measures have the following properties.

If  $\emptyset$  is the empty set then

$$P(\emptyset) = 0.$$

If  $A, B \in \mathcal{A}$  then

$$P(A \cup B) = P(A) + P(B) - P(A \cap B).$$

If  $A, B \in \mathcal{A}$  and  $B \subseteq A$  then

$$P(B) \leq P(A).$$

If  $A_n \in \mathcal{A}$ ,  $n = 1, 2, \dots$ , and  $\bigcup A_n \in \mathcal{A}$ , then

$$P(A_1 \cup A_2 \cup \dots) \leq P(A_1) + P(A_2) + \dots$$

The first three properties are evident. To establish the last one it is enough to observe that  $\bigcup_{n=1}^{\infty} A_n = \sum_{n=1}^{\infty} B_n$ , where  $B_1 = A_1$ ,  $B_n = \bar{A}_1 \cap \dots \cap \bar{A}_{n-1} \cap A_n$ ,  $n \geq 2$ ,  $B_i \cap B_j = \emptyset$ ,  $i \neq j$ , and therefore

$$P\left(\bigcup_{n=1}^{\infty} A_n\right) = P\left(\sum_{n=1}^{\infty} B_n\right) = \sum_{n=1}^{\infty} P(B_n) \leq \sum_{n=1}^{\infty} P(A_n).$$

The next theorem, which has many applications, provides conditions under which a finitely additive set function is actually countably additive.

**Theorem.** Let  $P$  be a finitely additive set function defined over the algebra  $\mathcal{A}$ , with  $P(\Omega) = 1$ . The following four conditions are equivalent:

- (1)  $P$  is  $\sigma$ -additive ( $P$  is a probability);
- (2)  $P$  is continuous from below, i.e. for any sets  $A_1, A_2, \dots \in \mathcal{A}$  such that  $A_n \subseteq A_{n+1}$  and  $\bigcup_{n=1}^{\infty} A_n \in \mathcal{A}$ ,

$$\lim_n P(A_n) = P\left(\bigcup_{n=1}^{\infty} A_n\right);$$

- (3)  $P$  is continuous from above, i.e. for any sets  $A_1, A_2, \dots$  such that  $A_n \supseteq A_{n+1}$  and  $\bigcap_{n=1}^{\infty} A_n \in \mathcal{A}$ ,

$$\lim_n P(A_n) = P\left(\bigcap_{n=1}^{\infty} A_n\right);$$

(4)  $P$  is continuous at  $\emptyset$ , i.e. for any sets  $A_1, A_2, \dots \in \mathcal{A}$  such that  $A_{n+1} \subseteq A_n$  and  $\bigcap_{n=1}^{\infty} A_n = \emptyset$ ,

$$\lim_n P(A_n) = 0.$$

PROOF. (1)  $\Rightarrow$  (2). Since

$$\bigcup_{n=1}^{\infty} A_n = A_1 + (A_2 \setminus A_1) + (A_3 \setminus A_2) + \dots,$$

we have

$$\begin{aligned} P\left(\bigcup_{n=1}^{\infty} A_n\right) &= P(A_1) + P(A_2 \setminus A_1) + P(A_3 \setminus A_2) + \dots \\ &= P(A_1) + P(A_2) - P(A_1) + P(A_3) - P(A_2) + \dots \\ &= \lim_n P(A_n). \end{aligned}$$

(2)  $\Rightarrow$  (3). Let  $n \geq 1$ ; then

$$P(A_n) = P(A_1 \setminus (A_1 \setminus A_n)) = P(A_1) - P(A_1 \setminus A_n).$$

The sequence  $\{A_1 \setminus A_n\}_{n \geq 1}$  of sets is nondecreasing (see the table in Subsection 3 below) and

$$\bigcup_{n=1}^{\infty} (A_1 \setminus A_n) = A_1 \setminus \bigcap_{n=1}^{\infty} A_n.$$

Then, by (2)

$$\lim_n P(A_1 \setminus A_n) = P\left(\bigcup_{n=1}^{\infty} (A_1 \setminus A_n)\right)$$

and therefore

$$\begin{aligned} \lim_n P(A_n) &= P(A_1) - \lim_n P(A_1 \setminus A_n) \\ &= P(A_1) - P\left(\bigcup_{n=1}^{\infty} (A_1 \setminus A_n)\right) = P(A_1) - P\left(A_1 \setminus \bigcap_{n=1}^{\infty} A_n\right) \\ &= P(A_1) - P(A_1) + P\left(\bigcap_{n=1}^{\infty} A_n\right) = P\left(\bigcap_{n=1}^{\infty} A_n\right) \end{aligned}$$

(3)  $\Rightarrow$  (4). Obvious.

(4)  $\Rightarrow$  (1). Let  $A_1, A_2, \dots \in \mathcal{A}$  be pairwise disjoint and let  $\sum_{n=1}^{\infty} A_n \in \mathcal{A}$ . Then

$$P\left(\sum_{i=1}^{\infty} A_i\right) = P\left(\sum_{i=1}^n A_i\right) + P\left(\sum_{i=n+1}^{\infty} A_i\right),$$

Table

Notation	Set-theoretic interpretation	Interpretation in probability theory
$\omega$	element or point	outcome, sample point, elementary event
$\Omega$	set of points	sample space; certain event
$\mathcal{F}$	$\sigma$ -algebra of subsets	$\sigma$ -algebra of events
$A \in \mathcal{F}$	set of points	event (if $\omega \in A$ , we say that event $A$ occurs)
$\bar{A} = \Omega \setminus A$	complement of $A$ , i.e. the set of points $\omega$ that are not in $A$	event that $A$ does not occur
$A \cup B$	union of $A$ and $B$ , i.e. the set of points $\omega$ belonging either to $A$ or to $B$	event that either $A$ or $B$ occurs
$A \cap B$ (or $AB$ )	intersection of $A$ and $B$ , i.e. the set of points $\omega$ belonging to both $A$ and $B$	event that both $A$ and $B$ occur
$\emptyset$	empty set	impossible event
$A \cap B = \emptyset$	$A$ and $B$ are disjoint	events $A$ and $B$ are mutually exclusive, i.e. cannot occur simultaneously
$A + B$	sum of sets, i.e. union of disjoint sets	event that one of two mutually exclusive events occurs
$A \setminus B$	difference of $A$ and $B$ , i.e. the set of points that belong to $A$ but not to $B$	event that $A$ occurs and $B$ does not
$A \triangle B$	symmetric difference of sets, i.e. $(A \setminus B) \cup (B \setminus A)$	event that $A$ or $B$ occurs, but not both
$\bigcup_{n=1}^{\infty} A_n$	union of the sets $A_1, A_2, \dots$	event that at least one of $A_1, A_2, \dots$ occurs

$\sum_{n=1}^{\infty} A_n$	sum, i.e. union of pairwise disjoint sets $A_1, A_2, \dots$	event that one of the mutually exclusive events $A_1, A_2, \dots$ occurs
$\bigcap_{n=1}^{\infty} A_n$	intersection of $A_1, A_2, \dots$	event that all the events $A_1, A_2, \dots$ occur
$A_n \uparrow A$ $\left( \text{or } A = \lim_n \uparrow A_n \right)$	the increasing sequence of sets $A_n$ converges to $A$ , i.e. $A_1 \subseteq A_2 \subseteq \dots$ and $A = \bigcup_{n=1}^{\infty} A_n$	the increasing sequence of events converges to event $A$
$A_n \downarrow A$ $\left( \text{or } A = \lim_n \downarrow A_n \right)$	the decreasing sequence of sets $A_n$ converges to $A$ , i.e. $A_1 \supseteq A_2 \supseteq \dots$ and $A = \bigcap_{n=1}^{\infty} A_n$	the decreasing sequence of events converges to event $A$
$\varlimsup A_n$ (or $\limsup A_n$ or* $\{A_n \text{ i.o.}\}$ )	the set $\bigcap_{n=1}^{\infty} \bigcup_{k=n}^{\infty} A_k$	event that infinitely many of events $A_1, A_2, \dots$ occur
$\varliminf A_n$ (or $\liminf A_n$ )	the set $\bigcup_{n=1}^{\infty} \bigcap_{k=n}^{\infty} A_k$	event that all the events $A_1, A_2, \dots$ occur with the possible exception of a finite number of them

\* i.o. = infinitely often.

and since  $\sum_{i=n+1}^{\infty} A_i \downarrow \emptyset, n \rightarrow \infty$ , we have

$$\begin{aligned} \sum_{i=1}^{\infty} P(A_i) &= \lim_n \sum_{i=1}^n P(A_i) = \lim_n P\left(\sum_{i=1}^n A_i\right) \\ &= \lim_n \left[ P\left(\sum_{i=1}^{\infty} A_i\right) - P\left(\sum_{i=n+1}^{\infty} A_i\right) \right] \\ &= P\left(\sum_{i=1}^{\infty} A_i\right) - \lim_n P\left(\sum_{i=n+1}^{\infty} A_i\right) = P\left(\sum_{i=1}^{\infty} A_i\right). \end{aligned}$$

3. We can now formulate Kolmogorov's generally accepted axiom system, which forms the basis for the concept of a probability space.

**Fundamental Definition.** An ordered triple  $(\Omega, \mathcal{F}, P)$  where

- (a)  $\Omega$  is a set of points  $\omega$ ,
- (b)  $\mathcal{F}$  is a  $\sigma$ -algebra of subsets of  $\Omega$ ,
- (c)  $P$  is a probability on  $\mathcal{F}$ ,

is called a probabilistic model or a probability space. Here  $\Omega$  is the sample space or space of elementary events, the sets  $A$  in  $\mathcal{F}$  are events, and  $P(A)$  is the probability of the event  $A$ .

It is clear from the definition that the axiomatic formulation of probability theory is based on set theory and measure theory. Accordingly, it is useful to have a table (pp. 136–137) displaying the ways in which various concepts are interpreted in the two theories. In the next two sections we shall give examples of the measurable spaces that are most important for probability theory and of how probabilities are assigned on them.

#### 4. PROBLEMS

1. Let  $\Omega = \{r: r \in [0, 1]\}$  be the set of rational points of  $[0, 1]$ ,  $\mathcal{A}$  the algebra of sets each of which is a finite sum of disjoint sets  $A$  of one of the forms  $\{r: a < r < b\}$ ,  $\{r: a \leq r < b\}$ ,  $\{r: a < r \leq b\}$ ,  $\{r: a \leq r \leq b\}$ , and  $P(A) = b - a$ . Show that  $P(A)$ ,  $A \in \mathcal{A}$ , is finitely additive set function but not countably additive.
2. Let  $\Omega$  be a countable set and  $\mathcal{F}$  the collection of all its subsets. Put  $\mu(A) = 0$  if  $A$  is finite and  $\mu(A) = \infty$  if  $A$  is infinite. Show that the set function  $\mu$  is finitely additive but not countably additive.
3. Let  $\mu$  be a finite measure on a  $\sigma$ -algebra  $\mathcal{F}$ ,  $A_n \in \mathcal{F}$ ,  $n = 1, 2, \dots$ , and  $A = \lim_n A_n$  (i.e.,  $A = \underline{\lim}_n A_n = \overline{\lim}_n A_n$ ). Show that  $\mu(A) = \lim_n \mu(A_n)$ .
4. Prove that  $P(A \triangle B) = P(A) + P(B) - 2P(A \cap B)$ .

5. Show that the “distances”  $\rho_1(A, B)$  and  $\rho_2(A, B)$  defined by

$$\rho_1(A, B) = P(A \triangle B),$$

$$\rho_2(A, B) = \begin{cases} \frac{P(A \triangle B)}{P(A \cup B)} & \text{if } P(A \cup B) \neq 0, \\ 0 & \text{if } P(A \cup B) = 0 \end{cases}$$

satisfy the triangle inequality.

6. Let  $\mu$  be a finitely additive measure on an algebra  $\mathcal{A}$ , and let the sets  $A_1, A_2, \dots \in \mathcal{A}$  be pairwise disjoint and satisfy  $A = \sum_{i=1}^{\infty} A_i \in \mathcal{A}$ . Then  $\mu(A) \geq \sum_{i=1}^{\infty} \mu(A_i)$ .

7. Prove that

$$\overline{\lim \sup A_n} = \lim \inf A_n, \quad \overline{\lim \inf A_n} = \lim \sup \bar{A}_n,$$

$$\lim \inf A_n \subseteq \lim \sup A_n, \quad \lim \sup(A_n \cup B_n) = \lim \sup A_n \cup \lim \sup B_n,$$

$$\lim \sup A_n \cap \lim \inf B_n \subseteq \lim \sup(A_n \cap B_n) \subseteq \lim \sup A_n \cap \lim \sup B_n.$$

If  $A_n \uparrow A$  or  $A_n \downarrow A$ , then

$$\lim \inf A_n = \lim \sup A_n.$$

8. Let  $\{x_n\}$  be a sequence of numbers and  $A_n = (-\infty, x_n)$ . Show that  $x = \lim \sup x_n$  and  $A = \lim \sup A_n$  are related in the following way:  $(-\infty, x) \subseteq A \subseteq (-\infty, x]$ . In other words,  $A$  is equal to either  $(-\infty, x)$  or to  $(-\infty, x]$ .

9. Give an example to show that if a measure takes the value  $+\infty$ , it does not follow in general that countable additivity implies continuity at  $\emptyset$ .

## §2. Algebras and $\sigma$ -Algebras. Measurable Spaces

**1. Algebras and  $\sigma$ -algebras** are the components out of which probabilistic models are constructed. We shall present some examples and a number of results for these systems.

Let  $\Omega$  be a sample space. Evidently each of the collections of sets

$$\mathcal{F}_* = \{\emptyset, \Omega\}, \quad \mathcal{F}^* = \{A: A \subseteq \Omega\}$$

is both an algebra and a  $\sigma$ -algebra. In fact,  $\mathcal{F}_*$  is trivial, the “poorest”  $\sigma$ -algebra, whereas  $\mathcal{F}^*$  is the “richest”  $\sigma$ -algebra, consisting of all subsets of  $\Omega$ .

When  $\Omega$  is a finite space, the  $\sigma$ -algebra  $\mathcal{F}^*$  is fully surveyable, and commonly serves as the system of events in the elementary theory. However, when the space is uncountable the class  $\mathcal{F}^*$  is much too large, since it is impossible to define “probability” on such a system of sets in any consistent way.

If  $A \subseteq \Omega$ , the system

$$\mathcal{F}_A = \{A, \bar{A}, \emptyset, \Omega\}$$



is another example of an algebra (and a  $\sigma$ -algebra), the algebra (or  $\sigma$ -algebra) generated by  $A$ .

This system of sets is a special case of the systems generated by decompositions. In fact, let

$$\mathcal{D} = \{D_1, D_2, \dots\}$$

be a *countable* decomposition of  $\Omega$  into nonempty sets:

$$\Omega = D_1 + D_2 + \dots; \quad D_i \cap D_j = \emptyset, \quad i \neq j.$$

Then the system  $\mathcal{A} = \alpha(\mathcal{D})$ , formed by the sets that are unions of finite numbers of elements of the decomposition, is an algebra.

The following lemma is particularly useful since it establishes the important principle that there is a smallest algebra, or  $\sigma$ -algebra, containing a given collection of sets.

**Lemma 1.** *Let  $\mathcal{E}$  be a collection of subsets of  $\Omega$ . Then there are a smallest algebra  $\alpha(\mathcal{E})$  and a smallest  $\sigma$ -algebra  $\sigma(\mathcal{E})$  containing all the sets that are in  $\mathcal{E}$ .*

**PROOF.** The class  $\mathcal{F}^*$  of all subsets of  $\Omega$  is a  $\sigma$ -algebra. Therefore there are at least one algebra and one  $\sigma$ -algebra containing  $\mathcal{E}$ . We now define  $\alpha(\mathcal{E})$  (or  $\sigma(\mathcal{E})$ ) to consist of all sets that belong to every algebra (or  $\sigma$ -algebra) containing  $\mathcal{E}$ . It is easy to verify that this system is an algebra (or  $\sigma$ -algebra) and indeed the smallest.

**Remark.** The algebra  $\alpha(E)$  (or  $\sigma(E)$ , respectively) is often referred to as the smallest algebra (or  $\sigma$ -algebra) generated by  $\mathcal{E}$ .

We often need to know what additional conditions will make an algebra, or some other system of sets, into a  $\sigma$ -algebra. We shall present several results of this kind.

**Definition 1.** A collection  $\mathcal{M}$  of subsets of  $\Omega$  is a *monotonic class* if  $A_n \in \mathcal{M}$ ,  $n = 1, 2, \dots$ , together with  $A_n \uparrow A$  or  $A_n \downarrow A$ , implies that  $A \in \mathcal{M}$ .

Let  $\mathcal{E}$  be a system of sets. Let  $\mu(\mathcal{E})$  be the smallest monotonic class containing  $\mathcal{E}$ . (The proof of the existence of this class is like the proof of Lemma 1.)

**Lemma 2.** *A necessary and sufficient condition for an algebra  $\mathcal{A}$  to be a  $\sigma$ -algebra is that it is a monotonic class.*

**PROOF.** A  $\sigma$ -algebra is evidently a monotonic class. Now let  $\mathcal{A}$  be a monotonic class and  $A_n \in \mathcal{A}$ ,  $n = 1, 2, \dots$ . It is clear that  $B_n = \bigcup_{i=1}^n A_i \in \mathcal{A}$  and  $B_n \subseteq B_{n+1}$ . Consequently, by the definition of a monotonic class,  $B_n \uparrow \bigcup_{i=1}^{\infty} A_i \in \mathcal{A}$ . Similarly we could show that  $\bigcap_{i=1}^{\infty} A_i \in \mathcal{A}$ .

By using this lemma, we can prove that, starting with an algebra  $\mathcal{A}$ , we can construct the  $\sigma$ -algebra  $\sigma(\mathcal{A})$  by means of monotonic limiting processes.

**Theorem 1.** *Let  $\mathcal{A}$  be an algebra. Then*

$$\mu(\mathcal{A}) = \sigma(\mathcal{A}). \quad (1)$$

PROOF. By Lemma 2,  $\mu(\mathcal{A}) \subseteq \sigma(\mathcal{A})$ . Hence it is enough to show that  $\mu(\mathcal{A})$  is a  $\sigma$ -algebra. But  $\mathcal{M} = \mu(\mathcal{A})$  is a monotonic class, and therefore, by Lemma 2 again, it is enough to show that  $\mu(\mathcal{A})$  is an algebra.

Let  $A \in \mathcal{M}$ ; we show that  $\bar{A} \in \mathcal{M}$ . For this purpose, we shall apply a principle that will often be used in the future, the *principle of appropriate sets*, which we now illustrate.

Let

$$\tilde{\mathcal{M}} = \{B: B \in \mathcal{M}, \bar{B} \in \mathcal{M}\}$$

be the sets that have the property that concerns us. It is evident that  $\mathcal{A} \subseteq \tilde{\mathcal{M}} \subseteq \mathcal{M}$ . Let us show that  $\tilde{\mathcal{M}}$  is a monotonic class.

Let  $B_n \in \tilde{\mathcal{M}}$ ; then  $B_n \in \mathcal{M}$ ,  $\bar{B}_n \in \mathcal{M}$ , and therefore

$$\lim \uparrow B_n \in \mathcal{M}, \quad \lim \uparrow \bar{B}_n \in \mathcal{M}, \quad \lim \downarrow B_n \in \mathcal{M}, \quad \lim \downarrow \bar{B}_n \in \mathcal{M}.$$

Consequently

$$\begin{aligned} \overline{\lim \uparrow B_n} &= \lim \downarrow \bar{B}_n \in \mathcal{M}, & \overline{\lim \downarrow B_n} &= \lim \uparrow \bar{B}_n \in \mathcal{M}, \\ \overline{\lim \uparrow \bar{B}_n} &= \lim \downarrow B_n \in \mathcal{M}, & \overline{\lim \downarrow \bar{B}_n} &= \lim \uparrow B_n \in \mathcal{M}, \end{aligned}$$

and therefore  $\tilde{\mathcal{M}}$  is a monotonic class. But  $\tilde{\mathcal{M}} \subseteq \mathcal{M}$  and  $\mathcal{M}$  is the smallest monotonic class. Therefore  $\tilde{\mathcal{M}} = \mathcal{M}$ , and if  $A \in \mathcal{M} = \mu(\mathcal{A})$ , then we also have  $\bar{A} \in \mathcal{M}$ , i.e.  $\mathcal{M}$  is closed under the operation of taking complements.

Let us now show that  $\mathcal{M}$  is closed under intersections.

Let  $A \in \mathcal{M}$  and

$$\mathcal{M}_A = \{B: B \in \mathcal{M}, A \cap B \in \mathcal{M}\}.$$

From the equations

$$\begin{aligned} \lim \downarrow (A \cap B_n) &= A \cap \lim \downarrow B_n, \\ \lim \uparrow (A \cap B_n) &= A \cap \lim \uparrow B_n \end{aligned}$$

it follows that  $\mathcal{M}_A$  is a monotonic class.

Moreover, it is easily verified that

$$(A \in \mathcal{M}_B) \Leftrightarrow (B \in \mathcal{M}_A). \quad (2)$$

Now let  $A \in \mathcal{A}$ ; then since  $\mathcal{A}$  is an algebra, for every  $B \in \mathcal{A}$  the set  $A \cap B \in \mathcal{A}$  and therefore

$$\mathcal{A} \subseteq \mathcal{M}_A \subseteq \mathcal{M}.$$

But  $\mathcal{M}_A$  is a monotonic class (since  $\lim \uparrow AB_n = A \lim \uparrow B_n$  and  $\lim \downarrow AB_n = A \lim \downarrow B_n$ ), and  $\mathcal{M}$  is the smallest monotonic class. Therefore  $\mathcal{M}_A = \mathcal{M}$  for all  $A \in \mathcal{A}$ . But then it follows from (2) that

$$(A \in \mathcal{M}_B) \Leftrightarrow (B \in \mathcal{M}_A = \mathcal{M}).$$

whenever  $A \in \mathcal{A}$  and  $B \in \mathcal{M}$ . Consequently if  $A \in \mathcal{A}$  then

$$A \in \mathcal{M}_B$$

for every  $B \in \mathcal{M}$ . Since  $A$  is any set in  $\mathcal{A}$ , it follows that

$$\mathcal{A} \subseteq \mathcal{M}_B \subseteq \mathcal{M}.$$

Therefore for every  $B \in \mathcal{M}$

$$\mathcal{M}_B = \mathcal{M},$$

i.e. if  $B \in \mathcal{M}$  and  $C \in \mathcal{M}$  then  $C \cap B \in \mathcal{M}$ .

Thus  $\mathcal{M}$  is closed under complementation and intersection (and therefore under unions). Consequently  $\mathcal{M}$  is an algebra, and the theorem is established.

**Definition 2.** Let  $\Omega$  be a space. A class  $\mathcal{D}$  of subsets of  $\Omega$  is a *d-system* if

- (a)  $\Omega \in \mathcal{D}$ ;
- (b)  $A, B \in \mathcal{D}, A \subseteq B \Rightarrow B \setminus A \in \mathcal{D}$ ;
- (c)  $A_n \in \mathcal{D}, A_n \subseteq A_{n+1} \Rightarrow \bigcup A_n \in \mathcal{D}$ .

If  $\mathcal{E}$  is a collection of sets then  $d(\mathcal{E})$  denotes the smallest *d-system* containing  $\mathcal{E}$ .

**Theorem 2.** If the collection  $\mathcal{E}$  of sets is closed under intersections, then

$$d(\mathcal{E}) = \sigma(\mathcal{E}) \quad (3)$$

**PROOF.** Every  $\sigma$ -algebra is a *d-system*, and consequently  $d(\mathcal{E}) \subseteq \sigma(\mathcal{E})$ . Hence if we prove that  $d(\mathcal{E})$  is closed under intersections,  $d(\mathcal{E})$  must be a  $\sigma$ -algebra and then, of course, the opposite inclusion  $\sigma(\mathcal{E}) \subseteq d(\mathcal{E})$  is valid.

The proof once again uses the principle of appropriate sets.

Let

$$\mathcal{E}_1 = \{B \in d(\mathcal{E}) : B \cap A \in d(\mathcal{E}) \text{ for all } A \in \mathcal{E}\}.$$

If  $B \in \mathcal{E}$  then  $B \cap A \in \mathcal{E}$  for all  $A \in \mathcal{E}$  and therefore  $\mathcal{E} \subseteq \mathcal{E}_1$ . But  $\mathcal{E}_1$  is a *d-system*. Hence  $d(\mathcal{E}) \subseteq \mathcal{E}_1$ . On the other hand,  $\mathcal{E}_1 \subseteq d(\mathcal{E})$  by definition. Consequently

$$\mathcal{E}_1 = d(\mathcal{E}).$$

Now let

$$\mathcal{E}_2 = \{B \in d(\mathcal{E}) : B \cap A \in d(\mathcal{E}) \text{ for all } A \in d(\mathcal{E})\}.$$

Again it is easily verified that  $\mathcal{E}_2$  is a *d-system*. If  $B \in \mathcal{E}$ , then by the definition of  $\mathcal{E}_1$  we obtain that  $B \cap A \in d(\mathcal{E})$  for all  $A \in \mathcal{E}_1 = d(\mathcal{E})$ . Consequently  $\mathcal{E} \subseteq \mathcal{E}_2$  and  $d(\mathcal{E}) \subseteq \mathcal{E}_2$ . But  $d(\mathcal{E}) \supseteq \mathcal{E}_2$ ; hence  $d(\mathcal{E}) = \mathcal{E}_2$ , and therefore

whenever  $A$  and  $B$  are in  $d(\mathcal{E})$ , the set  $A \cap B$  also belongs to  $d(\mathcal{E})$ , i.e.  $d(\mathcal{E})$  is closed under intersections.

This completes the proof of the theorem.

We next consider some measurable spaces  $(\Omega, \mathcal{F})$  which are extremely important for probability theory.

**2. The measurable space  $(R, \mathcal{B}(R))$ .** Let  $R = (-\infty, \infty)$  be the real line and

$$(a, b] = \{x \in R: a < x \leq b\}$$

for all  $a$  and  $b$ ,  $-\infty \leq a < b < \infty$ . The interval  $(a, \infty]$  is taken to be  $(a, \infty)$ . (This convention is required if the complement of an interval  $(-\infty, b]$  is to be an interval of the same form, i.e. open on the left and closed on the right.)

Let  $\mathcal{A}$  be the system of subsets of  $R$  which are finite sums of disjoint intervals of the form  $(a, b]$ :

$$A \in \mathcal{A} \text{ if } A = \sum_{i=1}^n (a_i, b_i], \quad n < \infty.$$

It is easily verified that this system of sets, in which we also include the empty set  $\emptyset$ , is an algebra. However, it is not a  $\sigma$ -algebra, since if  $A_n = (0, 1 - 1/n] \in \mathcal{A}$ , we have  $\bigcup_n A_n = (0, 1) \notin \mathcal{A}$ .

Let  $\mathcal{B}(R)$  be the smallest  $\sigma$ -algebra  $\sigma(\mathcal{A})$  containing  $\mathcal{A}$ . This  $\sigma$ -algebra, which plays an important role in analysis, is called the *Borel algebra* of subsets of the real line, and its sets are called *Borel sets*.

If  $\mathcal{J}$  is the system of intervals  $\mathcal{J}$  of the form  $(a, b]$ , and  $\sigma(\mathcal{J})$  is the smallest  $\sigma$ -algebra containing  $\mathcal{J}$ , it is easily verified that  $\sigma(\mathcal{J})$  is the Borel algebra. In other words, we can obtain the Borel algebra from  $\mathcal{J}$  without going through the algebra  $\mathcal{A}$ , since  $\sigma(\mathcal{J}) = \sigma(\mathcal{A}(\mathcal{J}))$ .

We observe that

$$\begin{aligned} (a, b) &= \bigcup_{n=1}^{\infty} \left( a, b - \frac{1}{n} \right], & a < b, \\ [a, b] &= \bigcap_{n=1}^{\infty} \left( a - \frac{1}{n}, b \right], & a < b, \\ \{a\} &= \bigcap_{n=1}^{\infty} \left( a - \frac{1}{n}, a \right]. \end{aligned}$$

Thus the Borel algebra contains not only intervals  $(a, b]$  but also the singletons  $\{a\}$  and all sets of the six forms

$$(a, b), \quad [a, b], \quad [a, b), \quad (-\infty, b), \quad (-\infty, b], \quad (a, \infty). \quad (4)$$

Let us also notice that the construction of  $\mathcal{B}(R)$  could have been based on any of the six kinds of intervals instead of on  $(a, b]$ , since all the minimal  $\sigma$ -algebras generated by systems of intervals of any of the forms (4) are the same as  $\mathcal{B}(R)$ .

Sometimes it is useful to deal with the  $\sigma$ -algebra  $\mathcal{B}(\bar{R})$  of subsets of the extended real line  $\bar{R} = [-\infty, \infty]$ . This is the smallest  $\sigma$ -algebra generated by intervals of the form

$$(a, b] = \{x \in \bar{R} : a < x \leq b\}, \quad -\infty \leq a < b \leq \infty,$$

where  $(-\infty, b]$  is to stand for the set  $\{x \in \bar{R} : -\infty \leq x \leq b\}$ .

**Remark 1.** The measurable space  $(R, \mathcal{B}(R))$  is often denoted by  $(R, \mathcal{B})$  or  $(R^1, \mathcal{B}_1)$ .

**Remark 2.** Let us introduce the metric

$$\rho_1(x, y) = \frac{|x - y|}{1 + |x - y|}$$

on the real line  $R$  (this is equivalent to the usual metric  $|x - y|$ ) and let  $\mathcal{B}_0(R)$  be the smallest  $\sigma$ -algebra generated by the open sets  $S_\rho(x^0) = \{x \in R : \rho_1(x, x^0) < \rho\}$ ,  $\rho > 0$ ,  $x^0 \in R$ . Then  $\mathcal{B}_0(R) = \mathcal{B}(R)$  (see Problem 7).

**3. The measurable space  $(R^n, \mathcal{B}(R^n))$ .** Let  $R^n = R \times \cdots \times R$  be the direct, or Cartesian, product of  $n$  copies of the real line, i.e. the set of ordered  $n$ -tuples  $x = (x_1, \dots, x_n)$ , where  $-\infty < x_k < \infty$ ,  $k = 1, \dots, n$ . The set

$$I = I_1 \times \cdots \times I_n,$$

where  $I_k = (a_k, b_k]$ , i.e. the set  $\{x \in R^n : x_k \in I_k, k = 1, \dots, n\}$ , is called a rectangle, and  $I_k$  is a side of the rectangle. Let  $\mathcal{J}$  be the set of all rectangles  $I$ . The smallest  $\sigma$ -algebra  $\sigma(\mathcal{J})$  generated by the system  $\mathcal{J}$  is the *Borel algebra* of subsets of  $R^n$  and is denoted by  $\mathcal{B}(R^n)$ . Let us show that we can arrive at this Borel algebra by starting in a different way.

Instead of the rectangles  $I = I_1 \times \cdots \times I_n$  let us consider the rectangles  $B = B_1 \times \cdots \times B_n$  with Borel sides ( $B_k$  is the Borel subset of the real line that appears in the  $k$ th place in the direct product  $R \times \cdots \times R$ ). The smallest  $\sigma$ -algebra containing all rectangles with Borel sides is denoted by

$$\mathcal{B}(R) \otimes \cdots \otimes \mathcal{B}(R)$$

and called the *direct product* of the  $\sigma$ -algebras  $\mathcal{B}(R)$ . Let us show that in fact

$$\mathcal{B}(R^n) = \mathcal{B}(R) \otimes \cdots \otimes \mathcal{B}(R).$$

In other words, the smallest  $\sigma$ -algebra generated by the rectangles  $I = I_1 \times \cdots \times I_n$  and the (broader) class of rectangles  $B = B_1 \times \cdots \times B_n$  with Borel sides are actually the same.

The proof depends on the following proposition.

**Lemma 3.** *Let  $\mathcal{E}$  be a class of subsets of  $\Omega$ , let  $B \subseteq \Omega$ , and define*

$$\mathcal{E} \cap B = \{A \cap B : A \in \mathcal{E}\}. \quad (5)$$

*Then*

$$\sigma(\mathcal{E} \cap B) = \sigma(\mathcal{E}) \cap B. \quad (6)$$

PROOF. Since  $\mathcal{E} \subseteq \sigma(\mathcal{E})$ , we have

$$\mathcal{E} \cap B \subseteq \sigma(\mathcal{E}) \cap B. \quad (7)$$

But  $\sigma(\mathcal{E}) \cap B$  is a  $\sigma$ -algebra; hence it follows from (7) that

$$\sigma(\mathcal{E} \cap B) \subseteq \sigma(\mathcal{E}) \cap B.$$

To prove the conclusion in the opposite direction, we again use the principle of appropriate sets.

Define

$$\mathcal{C}_B = \{A \in \sigma(\mathcal{E}) : A \cap B \in \sigma(\mathcal{E} \cap B)\}.$$

Since  $\sigma(\mathcal{E})$  and  $\sigma(\mathcal{E} \cap B)$  are  $\sigma$ -algebras,  $\mathcal{C}_B$  is also a  $\sigma$ -algebra, and evidently

$$\mathcal{E} \subseteq \mathcal{C}_B \subseteq \sigma(\mathcal{E}),$$

whence  $\sigma(\mathcal{E}) \subseteq \sigma(\mathcal{C}_B) = \mathcal{C}_B \subseteq \sigma(\mathcal{E})$  and therefore  $\sigma(\mathcal{E}) = \mathcal{C}_B$ . Therefore

$$A \cap B \in \sigma(\mathcal{E} \cap B)$$

for every  $A \subseteq \sigma(\mathcal{E})$ , and consequently  $\sigma(\mathcal{E}) \cap B \subseteq \sigma(\mathcal{E} \cap B)$ .

This completes the proof of the lemma.

Proof that  $\mathcal{B}(R^n)$  and  $\mathcal{B} \otimes \cdots \otimes \mathcal{B}$  are the same. This is obvious for  $n = 1$ . We now show that it is true for  $n = 2$ .

Since  $\mathcal{B}(R^2) \subseteq \mathcal{B} \otimes \mathcal{B}$ , it is enough to show that the Borel rectangle  $B_1 \times B_2$  belongs to  $\mathcal{B}(R^2)$ .

Let  $R^2 = R_1 \times R_2$ , where  $R_1$  and  $R_2$  are the "first" and "second" real lines,  $\tilde{\mathcal{B}}_1 = \mathcal{B}_1 \times R_2$ ,  $\tilde{\mathcal{B}}_2 = R_1 \times \mathcal{B}_2$ , where  $\mathcal{B}_1 \times R_2$  (or  $R_1 \times \mathcal{B}_2$ ) is the collection of sets of the form  $B_1 \times R_2$  (or  $R_1 \times B_2$ ), with  $B_1 \in \mathcal{B}_1$  (or  $B_2 \in \mathcal{B}_2$ ). Also let  $\mathcal{I}_1$  and  $\mathcal{I}_2$  be the sets of intervals in  $R_1$  and  $R_2$ , and  $\tilde{\mathcal{I}}_1 = \mathcal{I}_1 \times R_2$ ,  $\tilde{\mathcal{I}}_2 = R_1 \times \mathcal{I}_2$ . Then, by (6),

$$\begin{aligned} B_1 \times B_2 &= \tilde{B}_1 \cap \tilde{B}_2 \in \tilde{\mathcal{B}}_1 \cap \tilde{\mathcal{B}}_2 = \sigma(\tilde{\mathcal{I}}_1) \cap \tilde{B}_2 \\ &= \sigma(\tilde{\mathcal{I}}_1 \cap \tilde{B}_2) \subseteq \sigma(\tilde{\mathcal{I}}_1 \cap \tilde{\mathcal{I}}_2) \\ &= \sigma(\mathcal{I}_1 \times \mathcal{I}_2), \end{aligned}$$

as was to be proved.

The case of any  $n, n > 2$ , can be discussed in the same way.

**Remark.** Let  $\mathcal{B}_0(R^n)$  be the smallest  $\sigma$ -algebra generated by the open sets

$$S_\rho(x^0) = \{x \in R^n: \rho_n(x, x^0) < \rho\}, \quad x^0 \in R^n, \quad \rho > 0,$$

in the metric

$$\rho_n(x, x^0) = \sum_{k=1}^n 2^{-k} \rho_1(x_k, x_k^0),$$

where  $x = (x_1, \dots, x_n)$ ,  $x^0 = (x_1^0, \dots, x_n^0)$ .

Then  $\mathcal{B}_0(R_n) = \mathcal{B}(R^n)$  (Problem 7).

**4. The measurable space  $(R^\infty, \mathcal{B}(R^\infty))$**  plays a significant role in probability theory, since it is used as the basis for constructing probabilistic models of experiments with infinitely many steps.

The space  $R^\infty$  is the space of *ordered* sequences of numbers,

$$x = (x_1, x_2, \dots), \quad -\infty < x_k < \infty, \quad k = 1, 2, \dots$$

Let  $I_k$  and  $B_k$  denote, respectively, the intervals  $(a_k, b_k]$  and the Borel subsets of the  $k$ th line (with coordinate  $x_k$ ). We consider the *cylinder sets*

$$\mathcal{I}(I_1 \times \dots \times I_n) = \{x: x = (x_1, x_2, \dots), x_1 \in I_1, \dots, x_n \in I_n\}, \quad (8)$$

$$\mathcal{I}(B_1 \times \dots \times B_n) = \{x: x = (x_1, x_2, \dots), x_1 \in B_1, \dots, x_n \in B_n\}, \quad (9)$$

$$\mathcal{I}(B^n) = \{x: (x_1, \dots, x_n) \in B^n\}, \quad (10)$$

where  $B^n$  is a Borel set in  $\mathcal{B}(R^n)$ . Each cylinder  $\mathcal{I}(B_1 \times \dots \times B_n)$ , or  $\mathcal{I}(B^n)$ , can also be thought of as a cylinder with base in  $R^{n+1}, R^{n+2}, \dots$ , since

$$\mathcal{I}(B_1 \times \dots \times B_n) = \mathcal{I}(B_1 \times \dots \times B_n \times R),$$

$$\mathcal{I}(B^n) = \mathcal{I}(B^{n+1}),$$

where  $B^{n+1} = B^n \times R$ .

It follows that both systems of cylinders  $\mathcal{I}(B_1 \times \dots \times B_n)$  and  $\mathcal{I}(B^n)$  are algebras. It is easy to verify that the unions of disjoint cylinders

$$\mathcal{I}(I_1 \times \dots \times I_n)$$

also form an algebra. Let  $\mathcal{B}(R^\infty)$ ,  $\mathcal{B}_1(R^\infty)$  and  $\mathcal{B}_2(R^\infty)$  be the smallest  $\sigma$ -algebras containing all the sets (8), (9) or (10), respectively. (The  $\sigma$ -algebra  $\mathcal{B}_1(R^\infty)$  is often denoted by  $\mathcal{B}(R) \otimes \mathcal{B}(R) \times \dots$ .) It is clear that  $\mathcal{B}(R^\infty) \subseteq \mathcal{B}_1(R^\infty) \subseteq \mathcal{B}_2(R^\infty)$ . As a matter of fact, all three  $\sigma$ -algebras are the same.

To prove this, we put

$$\mathcal{C}_n = \{A \in R^n: \{x: (x_1, \dots, x_n) \in A\} \in \mathcal{B}(R^\infty)\}$$

for  $n = 1, 2, \dots$ . Let  $B^n \in \mathcal{B}(R^n)$ . Then

$$B^n \in \mathcal{C}_n \subseteq \mathcal{B}(R^\infty).$$

But  $\mathcal{C}_n$  is a  $\sigma$ -algebra, and therefore

$$\mathcal{B}(R^n) \subseteq \sigma(\mathcal{C}_n) = \mathcal{C}_n \subseteq \mathcal{B}(R^\infty);$$

consequently

$$\mathcal{B}_2(R^\infty) \subseteq \mathcal{B}(R^\infty).$$

Thus  $\mathcal{B}(R^\infty) = \mathcal{B}_1(R^\infty) = \mathcal{B}_2(R^\infty)$ .

From now on we shall describe sets in  $\mathcal{B}(R^\infty)$  as Borel sets (in  $R^\infty$ ).

**Remark.** Let  $\mathcal{B}_0(R^\infty)$  be the smallest  $\sigma$ -algebra generated by the open sets

$$S_\rho(x^0) = \{x \in R^\infty: \rho_\infty(x, x^0) < \rho\}, \quad x^0 \in R^\infty, \quad \rho > 0,$$

in the metric

$$\rho_\infty(x, x^0) = \sum_{k=1}^{\infty} 2^{-k} \rho_1(x_k, x_k^0),$$

where  $x = (x_1, x_2, \dots)$ ,  $x^0 = (x_1^0, x_2^0, \dots)$ . Then  $\mathcal{B}(R^\infty) = \mathcal{B}_0(R^\infty)$  (Problem 7).

Here are some examples of Borel sets in  $R^\infty$ :

(a)  $\{x \in R^\infty: \sup x_n > a\},$

$$\{x \in R^\infty: \inf x_n < a\};$$

(b)  $\{x \in R^\infty: \overline{\lim} x_n \leq a\},$

$$\{x \in R^\infty: \underline{\lim} x_n > a\},$$

where, as usual,

$$\overline{\lim} x_n = \inf_{n} \sup_{m \geq n} x_m, \quad \underline{\lim} x_n = \sup_{n} \inf_{m \geq n} x_m;$$

(c)  $\{x \in R^\infty: x_n \rightarrow \}$ , the set of  $x \in R^\infty$  for which  $\lim x_n$  exists and is finite;

(d)  $\{x \in R^\infty: \lim x_n > a\};$

(e)  $\{x \in R^\infty: \sum_{n=1}^{\infty} |x_n| > a\};$

(f)  $\{x \in R^\infty: \sum_{k=1}^n x_k = 0 \text{ for at least one } n \geq 1\}.$

To be convinced, for example, that sets in (a) belong to the system  $\mathcal{B}(R^\infty)$ , it is enough to observe that

$$\{x: \sup x_n > a\} = \bigcup_n \{x: x_n > a\} \in \mathcal{B}(R^\infty),$$

$$\{x: \inf x_n < a\} = \bigcup_n \{x: x_n < a\} \in \mathcal{B}(R^\infty).$$

**5. The measurable space  $(R^T, \mathcal{B}(R^T))$ ,** where  $T$  is an arbitrary set. The space  $R^T$  is the collection of real functions  $x = (x_t)$  defined for  $t \in T$ . In general we shall be interested in the case when  $T$  is an uncountable subset of the real

† We shall also use the notations  $x = (x_t)_{t \in R^T}$  and  $x = (x_t), t \in R^T$ , for elements of  $R^T$ .



line. For simplicity and definiteness we shall suppose for the present that  $T = [0, \infty)$ .

We shall consider three types of cylinder sets

$$\mathcal{J}_{t_1, \dots, t_n}(I_1 \times \dots \times I_n) = \{x: x_{t_1} \in I_1, \dots, x_{t_n} \in I_n\}, \quad (11)$$

$$\mathcal{J}_{t_1, \dots, t_n}(B_1 \times \dots \times B_n) = \{x: x_{t_1} \in B_1, \dots, x_{t_n} \in B_n\}, \quad (12)$$

$$\mathcal{J}_{t_1, \dots, t_n}(B^n) = \{x: (x_{t_1}, \dots, x_{t_n}) \in B^n\}, \quad (13)$$

where  $I_k$  is a set of the form  $(a_k, b_k]$ ,  $B_k$  is a Borel set on the line, and  $B^n$  is a Borel set in  $R^n$ .

The set  $\mathcal{J}_{t_1, \dots, t_n}(I_1 \times \dots \times I_n)$  is just the set of functions that, at times  $t_1, \dots, t_n$ , "get through the windows"  $I_1, \dots, I_n$  and at other times have arbitrary values (Figure 24).

Let  $\mathcal{B}(R^T)$ ,  $\mathcal{B}_1(R^T)$  and  $\mathcal{B}_2(R^T)$  be the smallest  $\sigma$ -algebras corresponding respectively to the cylinder sets (11), (12) and (13). It is clear that

$$\mathcal{B}(R^T) \subseteq \mathcal{B}_1(R^T) \subseteq \mathcal{B}_2(R^T). \quad (14)$$

As a matter of fact, all three of these  $\sigma$ -algebras are the same. Moreover, we can give a complete description of the structure of their sets.

**Theorem 3.** *Let  $T$  be any uncountable set. Then  $\mathcal{B}(R^T) = \mathcal{B}_1(R^T) = \mathcal{B}_2(R^T)$ , and every set  $A \in \mathcal{B}(R^T)$  has the following structure: there are a countable set of points  $t_1, t_2, \dots$  of  $T$  and a Borel set  $B$  in  $\mathcal{B}(R^\omega)$  such that*

$$A = \{x: (x_{t_1}, x_{t_2}, \dots) \in B\}. \quad (15)$$

**PROOF.** Let  $\mathcal{E}$  denote the collection of sets of the form (15) (for various aggregates  $(t_1, t_2, \dots)$  and Borel sets  $B$  in  $\mathcal{B}(R^\omega)$ ). If  $A_1, A_2, \dots \in \mathcal{E}$  and the corresponding aggregates are  $T^{(1)} = (t_1^{(1)}, t_2^{(1)}, \dots)$ ,  $T^{(2)} = (t_1^{(2)}, t_2^{(2)}, \dots)$ ,  $\dots$ ,

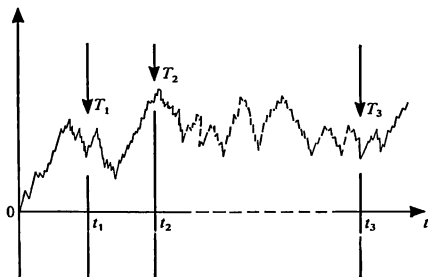


Figure 24

then the set  $T^{(\infty)} = \bigcup_k T^{(k)}$  can be taken as a basis, so that every  $A^{(i)}$  has a representation

$$A_i = \{x: (x_{\tau_1}, x_{\tau_2}, \dots) \in B_i\},$$

where  $B_i$  is a set in one and the same  $\sigma$ -algebra  $\mathcal{B}(R^\infty)$ , and  $\tau_i \in T^{(\infty)}$ .

Hence it follows that the system  $\mathcal{E}$  is a  $\sigma$ -algebra. Clearly this  $\sigma$ -algebra contains all cylinder sets of the form (1) and, since  $\mathcal{B}_2(R^T)$  is the smallest  $\sigma$ -algebra containing these sets, and since we have (14), we obtain

$$\mathcal{B}(R^T) \subseteq \mathcal{B}_1(R^T) \subseteq \mathcal{B}_2(R^T) \subseteq \mathcal{E}. \quad (16)$$

Let us consider a set  $A$  from  $\mathcal{E}$ , represented in the form (15). For a given aggregate  $(t_1, t_2, \dots)$ , the same reasoning as for the space  $(R^\infty, \mathcal{B}(R^\infty))$  shows that  $A$  is an element of the  $\sigma$ -algebra generated by the cylinder sets (11). But this  $\sigma$ -algebra evidently belongs to the  $\sigma$ -algebra  $\mathcal{B}(R^T)$ ; together with (16), this established both conclusions of the theorem.

Thus every Borel set  $A$  in the  $\sigma$ -algebra  $\mathcal{B}(R^T)$  is determined by restrictions imposed on the functions  $x = (x_t)$ ,  $t \in T$ , on an at most countable set of points  $t_1, t_2, \dots$ . Hence it follows, in particular, that the sets

$$A_1 = \{x: \sup_t x_t < C \text{ for all } t \in [0, 1]\},$$

$$A_2 = \{x: x_t = 0 \text{ for at least one } t \in [0, 1]\},$$

$$A_3 = \{x: x_t \text{ is continuous at a given point } t_0 \in [0, 1]\},$$

which depend on the behavior of the function on an uncountable set of points, cannot be Borel sets. And indeed *none of these three sets belongs to  $\mathcal{B}(R^{[0, 1]})$ .*

Let us establish this for  $A_1$ . If  $A_1 \in \mathcal{B}(R^{[0, 1]})$ , then by our theorem there are a point  $(t_1^0, t_2^0, \dots)$  and a set  $B^0 \in \mathcal{B}(R^\infty)$  such that

$$\left\{x: \sup_t x_t < C, t \in [0, 1]\right\} = \{x: (x_{t_1^0}, x_{t_2^0}, \dots) \in B^0\}.$$

It is clear that the function  $y_t \equiv C - 1$  belongs to  $A_1$ , and consequently  $(y_{t_1^0}, \dots) \in B^0$ . Now form the function

$$z_t = \begin{cases} C - 1, & t \in (t_1^0, t_2^0, \dots), \\ C + 1, & t \notin (t_1^0, t_2^0, \dots). \end{cases}$$

It is clear that

$$(y_{t_1^0}, y_{t_2^0}, \dots) = (z_{t_1^0}, z_{t_2^0}, \dots),$$

and consequently the function  $z = (z_t)$  belongs to the set  $\{x: (x_{t_1^0}, \dots) \in B^0\}$ . But at the same time it is clear that it does not belong to the set  $\{x: \sup_t x_t < C\}$ . This contradiction shows that  $A_1 \notin \mathcal{B}(R^{[0, 1]})$ .

Since the sets  $A_1$ ,  $A_2$  and  $A_3$  are nonmeasurable with respect to the  $\sigma$ -algebra  $\mathcal{B}[R^{[0,1]}]$  in the space of all functions  $x = (x_t)$ ,  $t \in [0, 1]$ , it is natural to consider a smaller class of functions for which these sets are measurable. It is intuitively clear that this will be the case if we take the initial space to be, for example, the space of continuous functions.

**6. The measurable space  $(C, \mathcal{B}(C))$ .** Let  $T = [0, 1]$  and let  $C$  be the space of continuous functions  $x = (x_t)$ ,  $0 \leq t \leq 1$ . This is a metric space with the metric  $\rho(x, y) = \sup_{t \in T} |x_t - y_t|$ . We introduce two  $\sigma$ -algebras in  $C$ :  $\mathcal{B}(C)$  is the  $\sigma$ -algebra generated by the cylinder sets, and  $\mathcal{B}_0(C)$  is generated by the open sets (open with respect to the metric  $\rho(x, y)$ ). Let us show that in fact these  $\sigma$ -algebras are the same:  $\mathcal{B}(C) = \mathcal{B}_0(C)$ .

Let  $B = \{x: x_{t_0} < b\}$  be a cylinder set. It is easy to see that this set is open. Hence it follows that  $\{x: x_{t_1} < b_1, \dots, x_{t_n} < b_n\} \in \mathcal{B}_0(C)$ , and therefore  $\mathcal{B}(C) \subseteq \mathcal{B}_0(C)$ .

Conversely, consider a set  $B_\rho = \{y: y \in S_\rho(x^0)\}$  where  $x^0$  is an element of  $C$  and  $S_\rho(x^0) = \{x \in C: \sup_{t \in T} |x_t - x_t^0| < \rho\}$  is an open ball with center at  $x^0$ . Since the functions in  $C$  are continuous,

$$\begin{aligned} B_\rho &= \{y \in C: y \in S_\rho(x^0)\} = \left\{y \in C: \max_t |y_t - x_t^0| < \rho\right\} \\ &= \bigcap_{t_k} \{y \in C: |y_{t_k} - x_{t_k}^0| < \rho\} \in \mathcal{B}(C), \quad (17) \end{aligned}$$

where  $t_k$  are the rational points of  $[0, 1]$ . Therefore  $\mathcal{B}_0(C) \subseteq \mathcal{B}(C)$ .

The following example is fundamental.

**7. The measurable space  $(D, \mathcal{B}(D))$ ,** where  $D$  is the space of functions  $x = (x_t)$ ,  $t \in [0, 1]$ , that are continuous on the right ( $x_t = x_{t+}$  for all  $t < 1$ ) and have limits from the left (at every  $t > 0$ ).

Just as for  $C$ , we can introduce a metric  $d(x, y)$  on  $D$  such that the  $\sigma$ -algebra  $\mathcal{B}_0(D)$  generated by the open sets will coincide with the  $\sigma$ -algebra  $\mathcal{B}(D)$  generated by the cylinder sets. This metric  $d(x, y)$ , which was introduced by Skorohod, is defined as follows:

$$d(x, y) = \inf\{\varepsilon > 0: \exists \lambda \in \Lambda: \sup_t |x_t - y_{\lambda(t)}| + \sup_t |t - \lambda(t)| \leq \varepsilon\}, \quad (18)$$

where  $\Lambda$  is the set of strictly increasing functions  $\lambda = \lambda(t)$  that are continuous on  $[0, 1]$  and have  $\lambda(0) = 0$ ,  $\lambda(1) = 1$ .

**8. The measurable space  $(\prod_{t \in T} \Omega_t, \bigotimes_{t \in T} \mathcal{F}_t)$ .** Along with the space  $(R^T, \mathcal{B}(R^T))$ , which is the direct product of  $T$  copies of the real line together with the system of Borel sets, probability theory also uses the measurable space  $(\prod_{t \in T} \Omega_t, \bigotimes_{t \in T} \mathcal{F}_t)$ , which is defined in the following way.

Let  $T$  be any set of indices and  $(\Omega_t, \mathcal{F}_t)$  a measurable space,  $t \in T$ . Let  $\Omega = \prod_{t \in T} \Omega_t$ , the set of functions  $\omega = (\omega_t)$ ,  $t \in T$ , such that  $\omega_t \in \Omega_t$  for each  $t \in T$ .

The collection of cylinder sets

$$\mathcal{J}_{t_1, \dots, t_n}(B_1 \times \dots \times B_n) = \{\omega: \omega_{t_1} \in B_1, \dots, \omega_{t_n} \in B_n\},$$

where  $B_{t_i} \in \mathcal{F}_{t_i}$ , is easily shown to be an algebra. The smallest  $\sigma$ -algebra containing all these cylinder sets is denoted by  $\bigotimes_{t \in T} \mathcal{F}_t$ , and the measurable space  $(\prod_{t \in T} \Omega_t, \bigotimes_{t \in T} \mathcal{F}_t)$  is called the *direct product* of the measurable spaces  $(\Omega_t, \mathcal{F}_t)$ ,  $t \in T$ .

## 9. PROBLEMS

1. Let  $\mathcal{B}_1$  and  $\mathcal{B}_2$  be  $\sigma$ -algebras of subsets of  $\Omega$ . Are the following systems of sets  $\sigma$ -algebras?

$$\mathcal{B}_1 \cap \mathcal{B}_2 \equiv \{A: A \in \mathcal{B}_1 \text{ and } A \in \mathcal{B}_2\},$$

$$\mathcal{B}_1 \cup \mathcal{B}_2 \equiv \{A: A \in \mathcal{B}_1 \text{ or } A \in \mathcal{B}_2\}.$$

2. Let  $\mathcal{D} = \{D_1, D_2, \dots\}$  be a countable decomposition of  $\Omega$  and  $\mathcal{B} = \sigma(\mathcal{D})$ . Are there also only countably many sets in  $\mathcal{B}$ ?

3. Show that

$$\mathcal{B}(R^n) \otimes \mathcal{B}(R) = \mathcal{B}(R^{n+1}).$$

4. Prove that the sets (b)–(f) (see Subsection 4) belong to  $\mathcal{B}(R^\infty)$ .
5. Prove that the sets  $A_2$  and  $A_3$  (see Subsection 5) do not belong to  $\mathcal{B}(R^{[0,1]})$ .
6. Prove that the function (15) actually defines a metric.
7. Prove that  $\mathcal{B}_0(R^n) = \mathcal{B}(R^n)$ ,  $n \geq 1$ , and  $\mathcal{B}_0(R^\infty) = \mathcal{B}(R^\infty)$ .
8. Let  $C = C[0, \infty)$  be the space of continuous functions  $x = (x_t)$  defined for  $t \geq 0$ . Show that with the metric

$$\rho(x, y) = \sum_{n=1}^{\infty} 2^{-n} \min \left[ \sup_{0 \leq t \leq n} |x_t - y_t|, 1 \right], \quad x, y \in C,$$

this is a complete separable metric space and that the  $\sigma$ -algebra  $\mathcal{B}_0(C)$  generated by the open sets coincides with the  $\sigma$ -algebra  $\mathcal{B}(C)$  generated by the cylinder sets.

## §3. Methods of Introducing Probability Measures on Measurable Spaces

1. **The measurable space  $(R, \mathcal{B}(R))$ .** Let  $P = P(A)$  be a probability measure defined on the Borel subsets  $A$  of the real line. Take  $A = (-\infty, x]$  and put

$$F(x) = P(-\infty, x], \quad x \in R. \quad (1)$$

This function has the following properties:

- (1)  $F(x)$  is nondecreasing;
- (2)  $F(-\infty) = 0, F(+\infty) = 1$ , where

$$F(-\infty) = \lim_{x \downarrow -\infty} F(x), \quad F(+\infty) = \lim_{x \uparrow \infty} F(x);$$

- (3)  $F(x)$  is continuous on the right and has a limit on the left at each  $x \in R$ .

The first property is evident, and the other two follow from the continuity properties of probability measures.

**Definition 1.** Every function  $F = F(x)$  satisfying conditions (1)–(3) is called a *distribution function* (on the real line  $R$ ).

Thus to every probability measure  $P$  on  $(R, \mathcal{B}(R))$  there corresponds (by (1)) a distribution function. It turns out that the converse is also true.

**Theorem 1.** Let  $F = F(x)$  be a distribution function on the real line  $R$ . There exists a unique probability measure  $P$  on  $(R, \mathcal{B}(R))$  such that

$$P(a, b] = F(b) - F(a) \quad (2)$$

for all  $a, b, -\infty \leq a < b < \infty$ .

**PROOF.** Let  $\mathcal{A}$  be the algebra of the subsets  $A$  of  $R$  that are finite sums of disjoint intervals of the form  $(a, b]$ :

$$A = \sum_{k=1}^n (a_k, b_k].$$

On these sets we define a set function  $P_0$  by putting

$$P_0(A) = \sum_{k=1}^n [F(b_k) - F(a_k)], \quad A \in \mathcal{A}. \quad (3)$$

This formula defines, evidently uniquely, a finitely additive set function on  $\mathcal{A}$ . Therefore if we show that this function is also countably additive on this algebra, the existence and uniqueness of the required measure  $P$  on  $\mathcal{B}(R)$  will follow immediately from a general result of measure theory (which we quote without proof).

**Carathéodory's Theorem.** Let  $\Omega$  be a space,  $\mathcal{A}$  an algebra of its subsets, and  $\mathcal{B} = \sigma(\mathcal{A})$  the smallest  $\sigma$ -algebra containing  $\mathcal{A}$ . Let  $\mu_0$  be a  $\sigma$ -additive measure on  $(\Omega, \mathcal{A})$ . Then there is a unique measure  $\mu$  on  $(\Omega, \sigma(\mathcal{A}))$  which is an extension of  $\mu_0$ , i.e. satisfies

$$\mu(A) = \mu_0(A), \quad A \in \mathcal{A}.$$

We are now to show that  $P_0$  is countably additive on  $\mathcal{A}$ . By a theorem from §1 it is enough to show that  $P_0$  is continuous at  $\emptyset$ , i.e. to verify that

$$P_0(A_n) \downarrow 0, \quad A_n \downarrow \emptyset, \quad A_n \in \mathcal{A}.$$

Let  $A_1, A_2, \dots$  be a sequence of sets from  $\mathcal{A}$  with the property  $A_n \downarrow \emptyset$ . Let us suppose first that the sets  $A_n$  belong to a closed interval  $[-N, N]$ ,  $N < \infty$ . Since  $A$  is the sum of finitely many intervals of the form  $(a, b]$  and since

$$P_0(a', b] = F(b) - F(a') \rightarrow F(b) - F(a) = P_0(a, b]$$

as  $a' \downarrow a$ , because  $F(x)$  is continuous on the right, we can find, for every  $A_n$ , a set  $B_n \in \mathcal{A}$  such that its closure  $[B_n] \subseteq A_n$  and

$$P_0(A_n) - P_0(B_n) \leq \varepsilon \cdot 2^{-n},$$

where  $\varepsilon$  is a preassigned positive number.

By hypothesis,  $\bigcap A_n = \emptyset$  and therefore  $\bigcap [B_n] = \emptyset$ . But the sets  $[B_n]$  are closed, and therefore there is a finite  $n_0 = n_0(\varepsilon)$  such that

$$\bigcap_{n=1}^{n_0} [B_n] = \emptyset. \quad (4)$$

(In fact,  $[-N, N]$  is compact, and the collection of sets  $\{[-N, N] \setminus [B_n]\}_{n \geq 1}$  is an open covering of this compact set. By the Heine-Borel theorem there is a finite subcovering:

$$\bigcup_{n=1}^{n_0} ([-N, N] \setminus [B_n]) = [-N, N]$$

and therefore  $\bigcap_{n=1}^{n_0} [B_n] = \emptyset$ ).

Using (4) and the inclusions  $A_{n_0} \subseteq A_{n_0-1} \subseteq \dots \subseteq A_1$ , we obtain

$$\begin{aligned} P_0(A_{n_0}) &= P_0\left(A_{n_0} \setminus \bigcap_{k=1}^{n_0} B_k\right) + P_0\left(\bigcap_{k=1}^{n_0} B_k\right) \\ &= P_0\left(A_{n_0} \setminus \bigcap_{k=1}^{n_0} B_k\right) \leq P_0\left(\bigcup_{k=1}^{n_0} (A_k \setminus B_k)\right) \\ &\leq \sum_{k=1}^{n_0} P_0(A_k \setminus B_k) \leq \sum_{k=1}^{n_0} \varepsilon \cdot 2^{-k} \leq \varepsilon. \end{aligned}$$

Therefore  $P_0(A_n) \downarrow 0$ ,  $n \rightarrow \infty$ .

We now abandon the assumption that  $A_n \subseteq [-N, N]$  for some  $N$ . Take an  $\varepsilon > 0$  and choose  $N$  so that  $P_0[-N, N] > 1 - \varepsilon/2$ . Then, since

$$A_n = A_n \cap [-N, N] + A_n \cap \overline{[-N, N]},$$

we have

$$\begin{aligned} P_0(A_n) &= P_0(A_n \cap [-N, N]) + P_0(A_n \cap \overline{[-N, N]}) \\ &\leq P_0(A_n \cap [-N, N]) + \varepsilon/2 \end{aligned}$$

and, applying the preceding reasoning (replacing  $A_n$  by  $A_n \cap [-N, N]$ ), we find that  $P_0(A_n \cap [-N, N]) \leq \varepsilon/2$  for sufficiently large  $n$ . Hence once again  $P_0(A_n) \downarrow 0, n \rightarrow \infty$ . This completes the proof of the theorem.

Thus there is a one-to-one correspondence between probability measures  $P$  on  $(R, \mathcal{B}(R))$  and distribution functions  $F$  on the real line  $R$ . The measure  $P$  constructed from the function  $F$  is usually called the Lebesgue-Stieltjes probability measure corresponding to the distribution function  $F$ .

The case when

$$F(x) = \begin{cases} 0, & x < 0, \\ x, & 0 \leq x \leq 1, \\ 1, & x > 1. \end{cases}$$

is particularly important. In this case the corresponding probability measure (denoted by  $\lambda$ ) is *Lebesgue measure* on  $[0, 1]$ . Clearly  $\lambda(a, b] = b - a$ . In other words, the Lebesgue measure of  $(a, b]$  (as well as of any of the intervals  $(a, b)$ ,  $[a, b]$  or  $[a, b)$ ) is simply its length  $b - a$ .

Let

$$\mathcal{B}([0, 1]) = \{A \cap [0, 1] : A \in \mathcal{B}(R)\}$$

be the collection of Borel subsets of  $[0, 1]$ . It is often necessary to consider, besides these sets, the Lebesgue measurable subsets of  $[0, 1]$ . We say that a set  $\Lambda \subseteq [0, 1]$  belongs to  $\mathcal{B}([0, 1])$  if there are Borel sets  $A$  and  $B$  such that  $A \subseteq \Lambda \subseteq B$  and  $\lambda(B \setminus A) = 0$ . It is easily verified that  $\mathcal{B}([0, 1])$  is a  $\sigma$ -algebra. It is known as the system of *Lebesgue measurable subsets of  $[0, 1]$* . Clearly  $\mathcal{B}([0, 1]) \subseteq \mathcal{B}([0, 1])$ .

The measure  $\lambda$ , defined so far only for sets in  $\mathcal{B}([0, 1])$ , extends in a natural way to the system  $\mathcal{B}([0, 1])$  of Lebesgue measurable sets. Specifically, if  $\Lambda \in \mathcal{B}([0, 1])$  and  $A \subseteq \Lambda \subseteq B$ , where  $A$  and  $B \in \mathcal{B}([0, 1])$  and  $\lambda(B \setminus A) = 0$ , we define  $\lambda(\Lambda) = \lambda(A)$ . The set function  $\bar{\lambda} = \lambda(\Lambda), \Lambda \in \mathcal{B}([0, 1])$ , is easily seen to be a probability measure on  $([0, 1], \mathcal{B}([0, 1]))$ . It is usually called *Lebesgue measure* (on the system of Lebesgue-measurable sets).

**Remark.** This process of completing (or extending) a measure can be applied, and is useful, in other situations. For example, let  $(\Omega, \mathcal{F}, P)$  be a probability space. Let  $\mathcal{F}^P$  be the collection of all the subsets  $A$  of  $\Omega$  for which there are sets  $B_1$  and  $B_2$  of  $\mathcal{F}$  such that  $B_1 \subseteq A \subseteq B_2$  and  $P(B_2 \setminus B_1) = 0$ . The probability measure can be defined for sets  $A \in \mathcal{F}^P$  in a natural way (by  $P(A) = P(B_1)$ ). The resulting probability space is the completion of  $(\Omega, \mathcal{F}, P)$  with respect to  $P$ .

A probability measure such that  $\mathcal{F}^P = \mathcal{F}$  is called *complete*, and the corresponding space  $(\Omega, \mathcal{F}, P)$  is a *complete probability space*.

The correspondence between probability measures  $P$  and distribution functions  $F$  established by the equation  $P(a, b] = F(b) - F(a)$  makes it

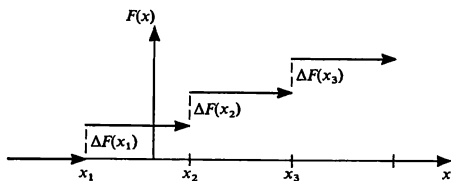


Figure 25

possible to construct various probability measures by obtaining the corresponding distribution functions.

**Discrete measures** are measures  $P$  for which the corresponding distributions  $F = F(x)$  are piecewise constant (Figure 25), changing their values at the points  $x_1, x_2, \dots$  ( $\Delta F(x_i) > 0$ , where  $\Delta F(x) = F(x) - F(x-)$ ). In this case the measure is concentrated at the points  $x_1, x_2, \dots$ :

$$P(\{x_k\}) = \Delta F(x_k) > 0, \quad \sum_k P(\{x_k\}) = 1.$$

The set of numbers  $(p_1, p_2, \dots)$ , where  $p_k = P(\{x_k\})$ , is called a *discrete probability distribution* and the corresponding distribution function  $F = F(x)$  is called *discrete*.

We present a table of the commonest types of discrete probability distribution, with their names.

Table 1

Distribution	Probabilities $p_k$	Parameters
Discrete uniform	$1/N, \quad k = 1, 2, \dots, N$	$N = 1, 2, \dots$
Bernoulli	$p_1 = p, \quad p_0 = q$	$0 \leq p \leq 1, \quad q = 1 - p$
Binomial	$C_n^k p^k q^{n-k}, \quad k = 0, 1, \dots, n$	$0 \leq p \leq 1, \quad q = 1 - p,$ $n = 1, 2, \dots$
Poisson	$e^{-k}/k!, \quad k = 0, 1, \dots$	$\lambda > 0$
Geometric	$q^{k-1}p, \quad k = 0, 1, \dots$	$0 \leq p \leq 1, \quad q = 1 - p$
Negative binomial	$C_{k-1}^{r-1} p^r q^{k-r}, \quad k = r, r+1, \dots$	$0 \leq p \leq 1, \quad q = 1 - p,$ $r = 1, 2, \dots$

**Absolutely continuous measures.** These are measures for which the corresponding distribution functions are such that

$$F(x) = \int_{-\infty}^x f(t) dt, \quad (5)$$



where  $f = f(t)$  are nonnegative functions and the integral is at first taken in the Riemann sense, but later (see §6) in that of Lebesgue.

The function  $f = f(x)$ ,  $x \in R$ , is the *density* of the distribution function  $F = F(x)$  (or the density of the probability distribution, or simply the density) and  $F = F(x)$  is called absolutely continuous.

It is clear that every nonnegative  $f = f(x)$  that is Riemann integrable and such that  $\int_{-\infty}^{\infty} f(x) dx = 1$  defines a distribution function by (5). Table 2 presents some important examples of various kinds of densities  $f = f(x)$  with their names and parameters (a density  $f(x)$  is taken to be zero for values of  $x$  not listed in the table).

Table 2

Distribution	Density	Parameters
Uniform on $[a, b]$	$1/(b - a)$ , $a \leq x \leq b$	$a, b \in R$ ; $a < b$
Normal or Gaussian	$(2\pi\sigma^2)^{-1/2} e^{-(x-m)^2/(2\sigma^2)}$ , $x \in R$	$m \in R$ , $\sigma > 0$
Gamma	$\frac{x^{\alpha-1} e^{-x/\beta}}{\Gamma(\alpha)\beta^\alpha}$ , $x \geq 0$	$\alpha > 0$ , $\beta > 0$
Beta	$\frac{x^{r-1}(1-x)^{s-1}}{B(r, s)}$ , $0 \leq x \leq 1$	$r > 0$ , $s > 0$
Exponential (gamma with $\alpha = 1$ , $\beta = 1/\lambda$ )	$\lambda e^{-\lambda x}$ , $x \geq 0$	$\lambda > 0$
Bilateral exponential	$\frac{1}{2} \lambda e^{-\lambda  x }$ , $x \in R$	$\lambda > 0$
Chi-squared, $\chi^2$ (gamma with $\alpha = n/2$ , $\beta = 2$ )	$2^{-n/2} x^{n/2-1} e^{-x/2} / \Gamma(n/2)$ , $x \geq 0$	$n = 1, 2, \dots$
Student, $t$	$\frac{\Gamma(\frac{1}{2}(n+1))}{(n\pi)^{1/2} \Gamma(n/2)} \left(1 + \frac{x^2}{n}\right)^{-(n+1)/2}$ , $x \in R$	$n = 1, 2, \dots$
$F$	$\frac{(m/n)^{m/2}}{B(m/2, n/2)} \frac{x^{m/2-1}}{(1 + mx/n)^{(m+n)/2}}$	$m, n = 1, 2, \dots$
Cauchy	$\frac{\theta}{\pi(x^2 + \theta^2)}$ , $x \in R$	$\theta > 0$

**Singular measures.** These are measures whose distribution functions are continuous but have all their points of increases on sets of zero *Lebesgue measure*. We do not discuss this case in detail; we merely give an example of such a function.

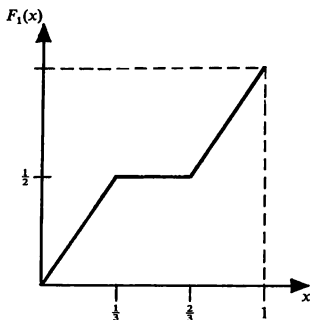


Figure 26

We consider the interval  $[0, 1]$  and construct  $F(x)$  by the following procedure originated by Cantor.

We divide  $[0, 1]$  into thirds and put (Figure 26)

$$F_2(x) = \begin{cases} \frac{1}{2}, & x \in (\frac{1}{3}, \frac{2}{3}), \\ \frac{1}{4}, & x \in (\frac{1}{9}, \frac{2}{9}), \\ \frac{3}{4}, & x \in (\frac{7}{9}, \frac{8}{9}), \\ 0, & x = 0, \\ 1, & x = 1 \end{cases}$$

defining it in the intermediate intervals by linear interpolation.

Then we divide each of the intervals  $[0, \frac{1}{3}]$  and  $[\frac{2}{3}, 1]$  into three parts and define the function (Figure 27) with its values at other points determined by linear interpolation.

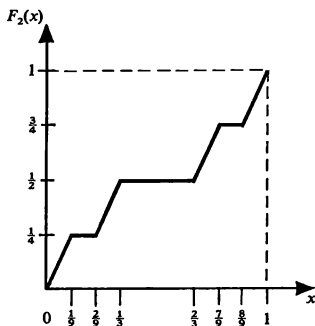


Figure 27

Continuing this process, we construct a sequence of functions  $F_n(x)$ ,  $n = 1, 2, \dots$ , which converges to a nondecreasing continuous function  $F(x)$  (the Cantor function), whose points of increase ( $x$  is a point of increase of  $F(x)$  if  $F(x + \varepsilon) - F(x - \varepsilon) > 0$  for every  $\varepsilon > 0$ ) form a set of Lebesgue measure zero. In fact, it is clear from the construction of  $F(x)$  that the total length of the intervals  $(\frac{1}{3}, \frac{2}{3})$ ,  $(\frac{1}{9}, \frac{2}{9})$ ,  $(\frac{7}{9}, \frac{8}{9})$ ,  $\dots$  on which the function is constant is

$$\frac{1}{3} + \frac{2}{9} + \frac{4}{27} + \dots = \frac{1}{3} \sum_{n=0}^{\infty} \left(\frac{2}{3}\right)^n = 1. \quad (6)$$

Let  $\mathcal{N}$  be the set of points of increase of the Cantor function  $F(x)$ . It follows from (6) that  $\lambda(\mathcal{N}) = 0$ . At the same time, if  $\mu$  is the measure corresponding to the Cantor function  $F(x)$ , we have  $\mu(\mathcal{N}) = 1$ . (We then say that the measure is *singular* with respect to Lebesgue measure  $\lambda$ .)

Without any further discussion of possible types of distribution functions, we merely observe that in fact the *three types* that have been mentioned cover all possibilities. More precisely, every distribution function can be represented in the form  $p_1 F_1 + p_2 F_2 + p_3 F_3$ , where  $F_1$  is discrete,  $F_2$  is absolutely continuous, and  $F_3$  is singular, and  $p_i$  are nonnegative numbers,  $p_1 + p_2 + p_3 = 1$ .

2. Theorem 1 establishes a one-to-one correspondence between probability measures on  $(R, \mathcal{B}(R))$  and distribution functions on  $R$ . An analysis of the proof of the theorem shows that in fact a stronger theorem is true, one that in particular lets us introduce Lebesgue measure on the real line.

Let  $\mu$  be a  $\sigma$ -finite measure on  $(\Omega, \mathcal{A})$ , where  $\mathcal{A}$  is an algebra of subsets of  $\Omega$ . It turns out that the conclusion of Carathéodory's theorem on the extension of a measure and an algebra  $\mathcal{A}$  to a minimal  $\sigma$ -algebra  $\sigma(\mathcal{A})$  remains valid with a  $\sigma$ -finite measure; this makes it possible to generalize Theorem 1.

A *Lebesgue-Stieltjes measure* on  $(R, \mathcal{B}(R))$  is a (countably additive) measure  $\mu$  such that the measure  $\mu(I)$  of every bounded interval  $I$  is finite. A *generalized distribution function* on the real line  $R$  is a nondecreasing function  $G = G(x)$ , with values on  $(-\infty, \infty)$ , that is continuous on the right.

Theorem 1 can be generalized to the statement that the formula

$$\mu(a, b] = G(b) - G(a), \quad a < b,$$

again establishes a one-to-one correspondence between Lebesgue-Stieltjes measures  $\mu$  and generalized distribution functions  $G$ .

In fact, if  $G(+\infty) - G(-\infty) < \infty$ , the proof of Theorem 1 can be taken over without any change, since this case reduces to the case when  $G(+\infty) - G(-\infty) = 1$  and  $G(-\infty) = 0$ .

Now let  $G(+\infty) - G(-\infty) = \infty$ . Put

$$G_n(x) = \begin{cases} G(x), & |x| \leq n, \\ G(n) & x = n, \\ G(-n), & x = -n. \end{cases}$$

On the algebra  $\mathcal{A}$  let us define a finitely additive measure  $\mu_0$  such that  $\mu_0(a, b] = G(b) - G(a)$ , and let  $\mu_n$  be the finitely additive measure previously constructed (by Theorem 1) from  $G_n(x)$ .

Evidently  $\mu_n \uparrow \mu_0$  on  $\mathcal{A}$ . Now let  $A_1, A_2, \dots$  be disjoint sets in  $\mathcal{A}$  and  $A \equiv \sum A_n \in \mathcal{A}$ . Then (Problem 6 of §1)

$$\mu_0(A) \geq \sum_{n=1}^{\infty} \mu_0(A_n).$$

If  $\sum_{n=1}^{\infty} \mu_0(A_n) = \infty$  then  $\mu_0(A) = \sum_{n=1}^{\infty} \mu_0(A_n)$ . Let us suppose that  $\sum \mu_0(A_n) < \infty$ . Then

$$\mu_0(A) = \lim_n \mu_n(A) = \lim_n \sum_{k=1}^{\infty} \mu_n(A_k).$$

By hypothesis,  $\sum \mu_0(A_n) < \infty$ . Therefore

$$0 \leq \mu_0(A) - \sum_{k=1}^{\infty} \mu_0(A_k) = \lim_n \left[ \sum_{k=1}^{\infty} (\mu_n(A_k) - \mu_0(A_k)) \right] \leq 0,$$

since  $\mu_n \leq \mu_0$ .

Thus a  $\sigma$ -finite finitely additive measure  $\mu_0$  is countably additive on  $\mathcal{A}$ , and therefore (by Carathéodory's theorem) it can be extended to a countably additive measure  $\mu$  on  $\sigma(\mathcal{A})$ .

The case  $G(x) = x$  is particularly important. The measure  $\lambda$  corresponding to this generalized distribution function is Lebesgue measure on  $(R, \mathcal{B}(R))$ . As for the interval  $[0, 1]$  of the real line, we can define the system  $\mathcal{B}(R)$  by writing  $A \in \mathcal{B}(R)$  if there are Borel sets  $A$  and  $B$  such that  $A \subseteq \Lambda \subseteq B$ ,  $\lambda(B \setminus A) = 0$ . Then Lebesgue measure  $\lambda$  on  $\mathcal{B}(R)$  is defined by  $\lambda(\Lambda) = \lambda(A)$  if  $A \subseteq \Lambda \subseteq B$ ,  $\lambda \in \mathcal{B}(R)$  and  $\lambda(B \setminus A) = 0$ .

**3. The measurable space  $(R^n, \mathcal{B}(R^n))$ .** Let us suppose, as for the real line, that  $P$  is a probability measure on  $(R^n, \mathcal{B}(R^n))$ .

Let us write

$$F_n(x_1, \dots, x_n) = P((-\infty, x_1] \times \dots \times (-\infty, x_n]),$$

or, in a more compact form,

$$F_n(x) = P(-\infty, x],$$

where  $x = (x_1, \dots, x_n)$ ,  $(-\infty, x] = (-\infty, x_1] \times \dots \times (-\infty, x_n]$ .

Let us introduce the difference operator  $\Delta_{a_i, b_i}: R^n \rightarrow R$ , defined by the formula

$$\begin{aligned} \Delta_{a_i, b_i} F_n(x_1, \dots, x_n) &= F_n(x_1, \dots, x_{i-1}, b_i, x_{i+1}, \dots) \\ &\quad - F_n(x_1, \dots, x_{i-1}, a_i, x_{i+1}, \dots) \end{aligned}$$

where  $a_i \leq b_i$ . A simple calculation shows that

$$\Delta_{a_1 b_1} \cdots \Delta_{a_n b_n} F_n(x_1 \cdots x_n) = P(a, b], \quad (7)$$

where  $(a, b] = (a_1, b_1] \times \cdots \times (a_n, b_n]$ . Hence it is clear, in particular, that (in contrast to the one-dimensional case)  $P(a, b]$  is in general not equal to  $F_n(b) - F_n(a)$ .

Since  $P(a, b] \geq 0$ , it follows from (7) that

$$\Delta_{a_1 b_1} \cdots \Delta_{a_n b_n} F_n(x_1, \dots, x_n) \geq 0 \quad (8)$$

for arbitrary  $a = (a_1, \dots, a_n)$ ,  $b = (b_1, \dots, b_n)$ .

It also follows from the continuity of  $P$  that  $F_n(x_1, \dots, x_n)$  is continuous on the right with respect to the variables collectively, i.e. if  $x^{(k)} \downarrow x$ ,  $x^{(k)} = (x_1^{(k)}, \dots, x_n^{(k)})$ , then

$$F_n(x^{(k)}) \downarrow F_n(x), \quad k \rightarrow \infty. \quad (9)$$

It is also clear that

$$F_n(+\infty, \dots, +\infty) = 1 \quad (10)$$

and

$$\lim_{x \downarrow y} F_n(x_1, \dots, x_n) = 0, \quad (11)$$

if at least one coordinate of  $y$  is  $-\infty$ .

**Definition 2.** An  $n$ -dimensional distribution function (on  $R^n$ ) is a function  $F = F(x_1, \dots, x_n)$  with properties (8)–(11).

The following result can be established by the same reasoning as in Theorem 1.

**Theorem 2.** Let  $F = F_n(x_1, \dots, x_n)$  be a distribution function on  $R^n$ . Then there is a unique probability measure  $P$  on  $(R^n, \mathcal{B}(R^n))$  such that

$$P(a, b] = \Delta_{a_1 b_1} \cdots \Delta_{a_n b_n} F_n(x_1, \dots, x_n). \quad (12)$$

Here are some examples of  $n$ -dimensional distribution functions.

Let  $F^1, \dots, F^n$  be one-dimensional distribution functions (on  $R$ ) and

$$F_n(x_1, \dots, x_n) = F^1(x_1) \cdots F^n(x_n).$$

It is clear that this function is continuous on the right and satisfies (10) and (11). It is also easy to verify that

$$\Delta_{a_1 b_1} \cdots \Delta_{a_n b_n} F_n(x_1, \dots, x_n) = \prod [F^k(b_k) - F^k(a_k)] \geq 0.$$

Consequently  $F_n(x_1, \dots, x_n)$  is a distribution function.

The case when

$$F^k(x_k) = \begin{cases} 0 & x_k < 0, \\ x_k, & 0 \leq x_k \leq 1, \\ 1, & x_k > 1 \end{cases}$$

is particularly important. In this case

$$F_n(x_1, \dots, x_n) = x_1 \cdots x_n.$$

The probability measure corresponding to this  $n$ -dimensional distribution function is  $n$ -dimensional Lebesgue measure on  $[0, 1]^n$ .

Many  $n$ -dimensional distribution functions appear in the form

$$F_n(x_1, \dots, x_n) = \int_{-\infty}^{x_1} \cdots \int_{-\infty}^{x_n} f_n(t_1, \dots, t_n) dt_1 \cdots dt_n,$$

where  $f_n(t_1, \dots, t_n)$  is a nonnegative function such that

$$\int_{-\infty}^{\infty} \cdots \int_{-\infty}^{\infty} f_n(t_1, \dots, t_n) dt_1 \cdots dt_n = 1,$$

and the integrals are Riemann (more generally, Lebesgue) integrals. The function  $f = f_n(t_1, \dots, t_n)$  is called the *density* of the  $n$ -dimensional distribution function, the density of the  $n$ -dimensional probability distribution, or simply an  $n$ -dimensional density.

When  $n = 1$ , the function

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-(x-m)^2/(2\sigma^2)}, \quad x \in R,$$

with  $\sigma > 0$  is the density of the (nondegenerate) *Gaussian* or *normal distribution*. There are natural analogs of this density when  $n > 1$ .

Let  $\mathbb{R} = \|r_{ij}\|$  be a nonnegative definite symmetric  $n \times n$  matrix:

$$\sum_{i,j=1}^n r_{ij} \lambda_i \lambda_j \geq 0, \quad \lambda_i \in R, \quad i = 1, \dots, n, \quad r_{ij} = r_{ji}.$$

When  $\mathbb{R}$  is a positive definite matrix,  $|\mathbb{R}| = \det \mathbb{R} > 0$  and consequently there is an inverse matrix  $A = \|a_{ij}\|$ .

$$f_n(x_1, \dots, x_n) = \frac{|A|^{1/2}}{(2\pi)^{n/2}} \exp\left\{-\frac{1}{2} \sum a_{ij}(x_i - m_i)(x_j - m_j)\right\}, \quad (13)$$

where  $m_i \in R$ ,  $i = 1, \dots, n$ , has the property that its (Riemann) integral over the whole space equals 1 (this will be proved in §13) and therefore, since it is also positive, it is a density.

This function is the *density of the  $n$ -dimensional (nondegenerate) Gaussian or normal distribution* (with vector mean  $m = (m_1, \dots, m_n)$  and covariance matrix  $\mathbb{R} = A^{-1}$ ).

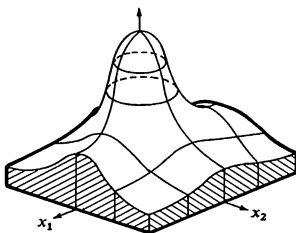


Figure 28. Density of the two-dimensional Gaussian distribution.

When  $n = 2$  the density  $f_2(x_1, x_2)$  can be put in the form

$$f_2(x_1, x_2) = \frac{1}{2\pi\sigma_1\sigma_2\sqrt{1-\rho^2}} \times \exp\left\{-\frac{1}{2(1-\rho^2)}\left[\frac{(x_1-m_1)^2}{\sigma_1^2} - 2\rho\frac{(x_1-m_1)(x_2-m_2)}{\sigma_1\sigma_2} + \frac{(x_2-m_2)^2}{\sigma_2^2}\right]\right\}, \quad (14)$$

where  $\sigma_i > 0$ ,  $|\rho| < 1$ . (The meanings of the parameters  $m_i$ ,  $\sigma_i$  and  $\rho$  will be explained in §8.)

Figure 28 indicates the form of the two-dimensional Gaussian density.

**Remark.** As in the case  $n = 1$ , Theorem 2 can be generalized to (similarly defined) Lebesgue–Stieltjes measures on  $(R^n, \mathcal{B}(R^n))$  and generalized distribution functions on  $R^n$ . When the generalized distribution function  $G_n(x_1, \dots, x_n)$  is  $x_1 \cdots x_n$ , the corresponding measure is Lebesgue measure on the Borel sets of  $R^n$ . It clearly satisfies

$$\lambda(a, b) = \prod_{i=1}^n (b_i - a_i),$$

i.e. the Lebesgue measure of the “rectangle”

$$(a, b] = (a_1, b_1] \times \cdots \times (a_n, b_n]$$

is its “content.”

**4. The measurable space  $(R^\infty, \mathcal{B}(R^\infty))$ .** For the spaces  $R^n$ ,  $n \geq 1$ , the probability measures were constructed in the following way: first for elementary sets (rectangles  $(a, b]$ ), then, in a natural way, for sets  $A = \sum (a_i, b_i]$ , and finally, by using Carathéodory’s theorem, for sets in  $\mathcal{B}(R^n)$ .

A similar construction for probability measures also works for the space  $(R^\infty, \mathcal{B}(R^\infty))$ .

Let

$$\mathcal{J}_n(B) = \{x \in R^\infty : (x_1, \dots, x_n) \in B\}, \quad B \in \mathcal{B}(R^n),$$

denote a cylinder set in  $R^\infty$  with base  $B \in \mathcal{B}(R^n)$ . We see at once that it is natural to take the cylinder sets as *elementary sets* in  $R^\infty$ , with their probabilities defined by the probability measure on the sets of  $\mathcal{B}(R^\infty)$ .

Let  $P$  be a probability measure on  $(R^\infty, \mathcal{B}(R^\infty))$ . For  $n = 1, 2, \dots$ , we take

$$P_n(B) = P(\mathcal{J}_n(B)), \quad B \in \mathcal{B}(R^n). \quad (15)$$

The sequence of probability measures  $P_1, P_2, \dots$  defined respectively on  $(R, \mathcal{B}(R)), (R^2, \mathcal{B}(R^2)), \dots$ , has the following evident consistency property: for  $n = 1, 2, \dots$  and  $B \in \mathcal{B}(R^n)$ ,

$$P_{n+1}(B \times R) = P_n(B). \quad (16)$$

It is noteworthy that the converse also holds.

**Theorem 3** (Kolmogorov's Theorem on the Extension of Measures in  $(R^\infty, \mathcal{B}(R^\infty))$ ). *Let  $P_1, P_2, \dots$  be a sequence of probability measures on  $(R, \mathcal{B}(R)), (R^2, \mathcal{B}(R^2)), \dots$ , possessing the consistency property (16). Then there is a unique probability measure  $P$  on  $(R^\infty, \mathcal{B}(R^\infty))$  such that*

$$P(\mathcal{J}_n(B)) = P_n(B), \quad B \in \mathcal{B}(R^n). \quad (17)$$

for  $n = 1, 2, \dots$ .

**PROOF.** Let  $B^n \in \mathcal{B}(R^n)$  and let  $\mathcal{J}_n(B^n)$  be the cylinder with base  $B^n$ . We assign the measure  $P(\mathcal{J}_n(B^n))$  to this cylinder by taking  $P(\mathcal{J}_n(B^n)) = P_n(B^n)$ .

Let us show that, in virtue of the consistency condition, this definition is consistent, i.e. the value of  $P(\mathcal{J}_n(B^n))$  is independent of the representation of the set  $\mathcal{J}_n(B^n)$ . In fact, let the same cylinder be represented in two way:

$$\mathcal{J}_n(B^n) = \mathcal{J}_{n+k}(B^{n+k}).$$

It follows that, if  $(x_1, \dots, x_{n+k}) \in R^{n+k}$ , we have

$$(x_1, \dots, x_n) \in B^n \Leftrightarrow (x_1, \dots, x_{n+k}) \in B^{n+k}, \quad (18)$$

and therefore, by (16) and (18),

$$\begin{aligned} P_n(B^n) &= P_{n+1}((x_1, \dots, x_{n+1}) : (x_1, \dots, x_n) \in B^n) \\ &= \dots = P_{n+k}((x_1, \dots, x_{n+k}) : (x_1, \dots, x_n) \in B^n) \\ &= P_{n+k}(B^{n+k}). \end{aligned}$$

Let  $\mathcal{A}(R^\infty)$  denote the collection of all cylinder sets  $\hat{B}^n = \mathcal{J}_n(B^n)$ ,  $B^n \in \mathcal{B}(R^n)$ ,  $n = 1, 2, \dots$ .



Now let  $\hat{B}_1, \dots, \hat{B}_k$  be disjoint sets in  $\mathcal{A}(R^\infty)$ . We may suppose without loss of generality that  $\hat{B}_i = \mathcal{J}_n(B_i^n)$ ,  $i = 1, \dots, k$ , for some  $n$ , where  $B_1^n, \dots, B_k^n$  are disjoint sets in  $\mathcal{B}(R^n)$ . Then

$$P\left(\sum_{i=1}^k \hat{B}_i\right) = P\left(\sum_{i=1}^k \mathcal{J}_n(B_i^n)\right) = P_n\left(\sum_{i=1}^k B_i^n\right) = \sum_{i=1}^n P_n(B_i^n) = \sum_{i=1}^n P(\hat{B}_i),$$

i.e. the set function  $P$  is finitely additive on the algebra  $\mathcal{A}(R^\infty)$ .

Let us show that  $P$  is "continuous at zero," i.e. if the sequence of sets  $\hat{B}_n \downarrow \emptyset$ ,  $n \rightarrow \infty$ , then  $P(\hat{B}_n) \rightarrow 0$ ,  $n \rightarrow \infty$ . Suppose the contrary, i.e. let  $\lim P(\hat{B}_n) = \delta > 0$ . We may suppose without loss of generality that  $\{\hat{B}_n\}$  has the form

$$\hat{B}_n = \{x: (x_1, \dots, x_n) \in B_n\}, \quad B_n \in \mathcal{B}(R^n).$$

We use the following property of probability measures  $P_n$  on  $(R^n, \mathcal{B}(R^n))$  (see Problem 9): if  $B_n \in \mathcal{B}(R^n)$ , for a given  $\delta > 0$  we can find a compact set  $A_n \in \mathcal{B}(R^n)$  such that  $A_n \subseteq B_n$  and

$$P_n(B_n \setminus A_n) \leq \delta/2^{n+1}.$$

Therefore if

$$\hat{A}_n = \{x: (x_1, \dots, x_n) \in A_n\},$$

we have

$$P(\hat{B}_n \setminus \hat{A}_n) = P_n(B_n \setminus A_n) \leq \delta/2^{n+1}.$$

Form the set  $\hat{C}_n = \bigcap_{k=1}^n \hat{A}_k$  and let  $C_n$  be such that

$$\hat{C}_n = \{x: (x_1, \dots, x_n) \in C_n\}.$$

Then, since the sets  $\hat{B}_n$  decrease, we obtain

$$P(\hat{B}_n \setminus \hat{C}_n) \leq \sum_{k=1}^n P(\hat{B}_n \setminus \hat{A}_k) \leq \sum_{k=1}^n P(\hat{B}_k \setminus A_k) \leq \delta/2.$$

But by assumption  $\lim_n P(\hat{B}_n) = \delta > 0$ , and therefore  $\lim_n P(\hat{C}_n) \geq \delta/2 > 0$ . Let us show that this contradicts the condition  $\hat{C}_n \downarrow \emptyset$ .

Let us choose a point  $\hat{x}^{(n)} = (x_1^{(n)}, x_2^{(n)}, \dots)$  in  $\hat{C}_n$ . Then  $(x_1^{(n)}, \dots, x_n^{(n)}) \in C_n$  for  $n \geq 1$ .

Let  $(n_1)$  be a subsequence of  $(n)$  such that  $x_1^{(n_1)} \rightarrow x_1^0$ , where  $x_1^0$  is a point in  $C_1$ . (Such a sequence exists since  $x_1^{(n)} \in C_1$  and  $C_1$  is compact.) Then select a subsequence  $(n_2)$  of  $(n_1)$  such that  $(x_1^{(n_2)}, x_2^{(n_2)}) \rightarrow (x_1^0, x_2^0) \in C_2$ . Similarly let  $(x_1^{(m_k)}, \dots, x_k^{(m_k)}) \rightarrow (x_1^0, \dots, x_k^0) \in C_k$ . Finally form the diagonal sequence  $(m_k)$ , where  $m_k$  is the  $k$ th term of  $(n_k)$ . Then  $x_i^{(m_k)} \rightarrow x_i^0$  as  $m_k \rightarrow \infty$  for  $i = 1, 2, \dots$ ; and  $(x_1^0, x_2^0, \dots) \in \hat{C}_n$  for  $n = 1, 2, \dots$ , which evidently contradicts the assumption that  $\hat{C}_n \downarrow \emptyset$ ,  $n \rightarrow \infty$ . This completes the proof of the theorem.

**Remark.** In the present case, the space  $R^\infty$  is a countable product of lines,  $R^\infty = R \times R \times \dots$ . It is natural to ask whether Theorem 3 remains true if  $(R^\infty, \mathcal{B}(R^\infty))$  is replaced by a direct product of measurable spaces  $(\Omega_i, \mathcal{F}_i)$ ,  $i = 1, 2, \dots$ .

We may notice that in the preceding proof the only topological property of the real line that was used was that every set in  $\mathcal{B}(R^n)$  contains a compact subset whose probability measure is arbitrarily close to the probability measure of the whole set. It is known, however, that this is a property not only of spaces  $(R^n, \mathcal{B}(R^n))$ , but also of arbitrary complete separable metric spaces with  $\sigma$ -algebras generated by the open sets.

Consequently Theorem 3 remains valid if we suppose that  $P_1, P_2, \dots$  is a sequence of consistent probability measures on  $(\Omega_1, \mathcal{F}_1)$ ,

$$(\Omega_1 \times \Omega_2, \mathcal{F}_1 \otimes \mathcal{F}_2), \dots,$$

where  $(\Omega_i, \mathcal{F}_i)$  are complete separable metric spaces with  $\sigma$ -algebras  $\mathcal{F}_i$  generated by open sets, and  $(R^\infty, \mathcal{B}(R^\infty))$  is replaced by

$$(\Omega_1 \times \Omega_2 \times \dots, \mathcal{F}_1 \otimes \mathcal{F}_2 \otimes \dots).$$

In §9 (Theorem 2) it will be shown that the result of Theorem 3 remains valid for arbitrary measurable spaces  $(\Omega_i, \mathcal{F}_i)$  if the measures  $P_n$  are concentrated in a particular way. However, Theorem 3 may fail in the general case (without any hypotheses on the topological nature of the measurable spaces or on the structure of the family of measures  $\{P_n\}$ ). This is shown by the following example.

Let us consider the space  $\Omega = (0, 1]$ , which is evidently not complete, and construct a sequence  $\mathcal{F}_1 \subseteq \mathcal{F}_2 \subseteq \dots$  of  $\sigma$ -algebras in the following way. For  $n = 1, 2, \dots$ , let

$$\varphi_n(\omega) = \begin{cases} 1, & 0 < \omega < 1/n, \\ 0, & 1/n \leq \omega \leq 1, \end{cases}$$

$$\mathcal{C}_n = \{A \in \Omega: A = \{\omega: \varphi_n(\omega) \in B\}, B \in \mathcal{B}(R)\}$$

and let  $\mathcal{F}_n = \sigma\{\mathcal{C}_1, \dots, \mathcal{C}_n\}$  be the smallest  $\sigma$ -algebra containing the sets  $\mathcal{C}_1, \dots, \mathcal{C}_n$ . Clearly  $\mathcal{F}_1 \subseteq \mathcal{F}_2 \subseteq \dots$ . Let  $\mathcal{F} = \sigma(\bigcup \mathcal{F}_n)$  be the smallest  $\sigma$ -algebra containing all the  $\mathcal{F}_n$ . Consider the measurable space  $(\Omega, \mathcal{F})$  and define a probability measure  $P_n$  on it as follows:

$$P_n\{\omega: (\varphi_1(\omega), \dots, \varphi_n(\omega)) \in B^n\} = \begin{cases} 1 & \text{if } (1, \dots, 1) \in B^n, \\ 0 & \text{otherwise,} \end{cases}$$

where  $B^n = \mathcal{B}(R^n)$ . It is easy to see that the family  $\{P_n\}$  is consistent: if  $A \in \mathcal{F}_n$  then  $P_{n+1}(A) = P_n(A)$ . However, we claim that there is no probability measure  $P$  on  $(\Omega, \mathcal{F})$  such that its restriction  $P|_{\mathcal{F}_n}$  (i.e., the measure  $P$

considered only on sets in  $\mathcal{F}_n$  coincides with  $P_n$  for  $n = 1, 2, \dots$ . In fact, let us suppose that such a probability measure  $P$  exists. Then

$$P\{\omega: \varphi_1(\omega) = \dots = \varphi_n(\omega) = 1\} = P_n\{\omega: \varphi_1(\omega) = \dots = \varphi_n(\omega) = 1\} = 1 \quad (19)$$

for  $n = 1, 2, \dots$ . But

$$\{\omega: \varphi_1(\omega) = \dots = \varphi_n(\omega) = 1\} = (0, 1/n) \downarrow \emptyset,$$

which contradicts (19) and the hypothesis of countable additivity (and therefore continuity at the "zero"  $\emptyset$ ) of the set function  $P$ .

We now give an example of a probability measure on  $(R^\infty, \mathcal{B}(R^\infty))$ . Let  $F_1(x), F_2(x), \dots$  be a sequence of one-dimensional distribution functions. Define the functions  $G(x) = F_1(x)$ ,  $G_2(x_1, x_2) = F_1(x_1)F_2(x_2), \dots$ , and denote the corresponding probability measures on  $(R, \mathcal{B}(R)), (R^2, \mathcal{B}(R^2)), \dots$  by  $P_1, P_2, \dots$ . Then it follows from Theorem 3 that there is a measure  $P$  on  $(R^\infty, \mathcal{B}(R^\infty))$  such that

$$P\{x \in R^\infty: (x_1, \dots, x_n) \in B\} = P_n(B), \quad B \in \mathcal{B}(R^n)$$

and, in particular,

$$P\{x \in R^\infty: x_1 \leq a_1, \dots, x_n \leq a_n\} = F_1(a_1) \cdots F_n(a_n).$$

Let us take  $F_i(x)$  to be a Bernoulli distribution,

$$F_i(x) = \begin{cases} 0, & x < 0, \\ q, & 0 \leq x < 1, \\ 1, & x \geq 1. \end{cases}$$

Then we can say that there is a probability measure  $P$  on the space  $\Omega$  of sequences of numbers  $x = (x_1, x_2, \dots), x_i = 0$  or  $1$ , together with the  $\sigma$ -algebra of its Borel subsets, such that

$$P\{x: x_1 = a_1, \dots, x_n = a_n\} = p^{\sum a_i} q^{n - \sum a_i}.$$

This is precisely the result that was not available in the first chapter for stating the law of large numbers in the form (I.5.8).

**5. The measurable space  $(R^T, \mathcal{B}(R^T))$ .** Let  $T$  be a set of indices  $t \in T$  and  $R_t$  a real line corresponding to the index  $t$ . We consider a finite unordered set  $\tau = [t_1, \dots, t_n]$  of distinct indices  $t_i, t_i \in T, n \geq 1$ , and let  $P_\tau$  be a probability measure on  $(R^\tau, \mathcal{B}(R^\tau))$ , where  $R^\tau = R_{t_1} \times \dots \times R_{t_n}$ .

We say that the family  $\{P_\tau\}$  of probability measures, where  $\tau$  runs through all finite unordered sets, is *consistent* if, for all sets  $\tau = [t_1, \dots, t_n]$  and  $\sigma = [s_1, \dots, s_k]$  such that  $\sigma \subseteq \tau$  we have

$$P_\sigma\{(x_{s_1}, \dots, x_{s_k}): (x_{s_1}, \dots, x_{s_k}) \in B\} = P_\tau\{(x_{t_1}, \dots, x_{t_n}): (x_{s_1}, \dots, x_{s_k}) \in B\} \quad (20)$$

for every  $B \in \mathcal{B}(R^\sigma)$ .

**Theorem 4** (Kolmogorov's Theorem on the Extension of Measures in  $(R^T, \mathcal{B}(R^T))$ ). Let  $\{P_\tau\}$  be a consistent family of probability measures on  $(R^\tau, \mathcal{B}(R^\tau))$ . Then there is a unique probability measure  $P$  on  $(R^T, \mathcal{B}(R^T))$  such that

$$P(\mathcal{I}_\tau(B)) = P_\tau(B) \quad (21)$$

for all unordered sets  $\tau = [t_1, \dots, t_n]$  of different indices  $t_i \in T$ ,  $B \in \mathcal{B}(R^\tau)$  and  $\mathcal{I}_\tau(B) = \{x \in R^T: (x_{t_1}, \dots, x_{t_n}) \in B\}$ .

PROOF. Let the set  $\hat{B} \in \mathcal{B}(R^T)$ . By the theorem of §2 there is an at most countable set  $S = \{s_1, s_2, \dots\} \subseteq T$  such that  $\hat{B} = \{x: (x_{s_1}, x_{s_2}, \dots) \in B\}$ , where  $B \in \mathcal{B}(R^S)$ ,  $R^S = R_{s_1} \times R_{s_2} \times \dots$ . In other words,  $\hat{B} = \mathcal{I}_S(B)$  is a cylinder set with base  $B \in \mathcal{B}(R^S)$ .

We can define a set function  $P$  on such cylinder sets by putting

$$P(\mathcal{I}_S(B)) = P_S(B), \quad (22)$$

where  $P_S$  is the probability measure whose existence is guaranteed by Theorem 3. We claim that  $P$  is in fact the measure whose existence is asserted in the theorem. To establish this we first verify that the definition (22) is consistent, i.e. that it leads to a unique value of  $P(\hat{B})$  for all possible representations of  $\hat{B}$ ; and second, that this set function is countably additive.

Let  $\hat{B} = \mathcal{I}_{S_1}(B_1)$  and  $\hat{B} = \mathcal{I}_{S_2}(B_2)$ . It is clear that then  $\hat{B} = \mathcal{I}_{S_1 \cup S_2}(B_3)$  with some  $B_3 \in \mathcal{B}(R^{S_1 \cup S_2})$ ; therefore it is enough to show that if  $S \subseteq S'$  and  $B \in \mathcal{B}(R^S)$ , then  $P_{S'}(B') = P_S(B)$ , where

$$B' = \{(x_{s'_1}, x_{s'_2}, \dots): (x_{s_1}, x_{s_2}, \dots) \in B\}$$

with  $S' = \{s'_1, s'_2, \dots\}$ ,  $S = \{s_1, s_2, \dots\}$ . But by the assumed consistency of (20) this equation follows immediately from Theorem 3. This establishes that the value of  $P(\hat{B})$  is independent of the representation of  $\hat{B}$ .

To verify the countable additivity of  $P$ , let us suppose that  $\{\hat{B}_n\}$  is a sequence of pairwise disjoint sets in  $\mathcal{B}(R^T)$ . Then there is an at most countable set  $S \subseteq T$  such that  $\hat{B}_n = \mathcal{I}_S(B_n)$  for all  $n \geq 1$ , where  $B_n \in \mathcal{B}(R^S)$ . Since  $P_S$  is a probability measure, we have

$$\begin{aligned} P(\sum \hat{B}_n) &= P(\sum \mathcal{I}_S(B_n)) = P_S(\sum B_n) = \sum P_S(B_n) \\ &= \sum P(\mathcal{I}_S(B_n)) = \sum P(\hat{B}_n). \end{aligned}$$

Finally, property (21) follows immediately from the way in which  $P$  was constructed.

This completes the proof.

**Remark 1.** We emphasize that  $T$  is any set of indices. Hence, by the remark after Theorem 3, the present theorem remains valid if we replace the real lines  $R_i$  by arbitrary complete separable metric spaces  $\Omega_i$  (with  $\sigma$ -algebras generated by open sets).

**Remark 2.** The original probability measures  $\{P_\tau\}$  were assumed defined on *unordered* sets  $\tau = [t_1, \dots, t_n]$  of different indices. It is also possible to start from a family of probability measures  $\{P_\tau\}$  where  $\tau$  runs through all *ordered* sets  $\tau = (t_1, \dots, t_n)$  of different indices. In this case, in order to have Theorem 4 hold we have to adjoin to (20) a further *consistency condition*:

$$P_{(i_1, \dots, i_n)}(A_{t_{i_1}} \times \dots \times A_{t_{i_n}}) = P_{(t_{i_1}, \dots, t_{i_n})}(A_{t_{i_1}} \times \dots \times A_{t_{i_n}}), \quad (23)$$

where  $(i_1, \dots, i_n)$  is an arbitrary permutation of  $(1, \dots, n)$  and  $A_{t_i} \in \mathcal{B}(R_{t_i})$ . As a necessary condition for the existence of  $P$  this follows from (21) (with  $P_{[t_1, \dots, t_n]}(B)$  replaced by  $P_{(t_1, \dots, t_n)}(B)$ ).

From now on we shall assume that the sets  $\tau$  under consideration are *unordered*. If  $T$  is a subset of the real line (or some completely ordered set), we may assume without loss of generality that the set  $\tau = [t_1, \dots, t_n]$  satisfies  $t_1 < t_2 < \dots < t_n$ . Consequently it is enough to define "finite-dimensional" probabilities only for sets  $\tau = [t_1, \dots, t_n]$  for which  $t_1 < t_2 < \dots < t_n$ .

Now consider the case  $T = [0, \infty)$ . Then  $R^T$  is the space of all real functions  $x = (x_t)_{t \geq 0}$ . A fundamental example of a probability measure on  $(R^{[0, \infty)}, \mathcal{B}(R^{[0, \infty)}))$  is Wiener measure, constructed as follows.

Consider the family  $\{\varphi_t(y|x)\}_{t \geq 0}$  of Gaussian densities (as functions of  $y$  for fixed  $x$ ):

$$\varphi_t(y|x) = \frac{1}{\sqrt{2\pi t}} e^{-(y-x)^2/2t}, \quad y \in R,$$

and for each  $\tau = [t_1, \dots, t_n]$ ,  $t_1 < t_2 < \dots < t_n$ , and each set

$$B = I_1 \times \dots \times I_n, \quad I_k = (a_k, b_k),$$

construct the measure  $P_\tau(B)$  according to the formula

$$P_\tau(I_1 \times \dots \times I_n) = \int_{I_1} \dots \int_{I_n} \varphi_{t_1}(a_1|0) \varphi_{t_2-t_1}(a_2|a_1) \dots \varphi_{t_n-t_{n-1}}(a_n|a_{n-1}) da_1 \dots da_n \quad (24)$$

(integration in the Riemann sense). Now we define the set function  $P$  for each cylinder set  $\mathcal{J}_{t_1 \dots t_n}(I_1 \times \dots \times I_n) = \{x \in R^T: x_{t_1} \in I_1, \dots, x_{t_n} \in I_n\}$  by taking

$$P(\mathcal{J}_{t_1 \dots t_n}(I_1 \times \dots \times I_n)) = P_{[t_1 \dots t_n]}(I_1 \times \dots \times I_n).$$

The intuitive meaning of this method of assigning a measure to the cylinder set  $\mathcal{J}_{t_1 \dots t_n}(I_1 \times \dots \times I_n)$  is as follows.

The set  $\mathcal{J}_{t_1 \dots t_n}(I_1 \times \dots \times I_n)$  is the set of functions that at times  $t_1, \dots, t_n$  pass through the "windows"  $I_1, \dots, I_n$  (see Figure 24 in §2). We shall interpret

$\varphi_{t_k - t_{k-1}}(a_k | a_{k-1})$  as the probability that a particle, starting at  $a_{k-1}$  at time  $t_k - t_{k-1}$ , arrives in a neighborhood of  $a_k$ . Then the product of densities that appears in (24) describes a certain independence of the increments of the displacements of the moving "particle" in the time intervals

$$[0, t_1], [t_1, t_2], \dots, [t_{n-1}, t_n].$$

The family of measures  $\{P_\tau\}$  constructed in this way is easily seen to be consistent, and therefore can be extended to a measure on  $(R^{[0, \infty)}, \mathcal{B}(R^{[0, \infty)}))$ . The measure so obtained plays an important role in probability theory. It was introduced by N. Wiener and is known as *Wiener measure*.

## 6. PROBLEMS

1. Let  $F(x) = P(-\infty, x]$ . Verify the following formulas:

$$P(a, b] = F(b) - F(a), \quad P(a, b) = F(b-) - F(a),$$

$$P[a, b] = F(b) - F(a-), \quad P[a, b) = F(b-) - F(a-),$$

$$P\{x\} = F(x) - F(x-),$$

where  $F(x-) = \lim_{y \uparrow x} F(y)$ .

2. Verify (7).
3. Prove Theorem 2.
4. Show that a distribution function  $F = F(x)$  on  $R$  has at most a countable set of points of discontinuity. Does a corresponding result hold for distribution functions on  $R^n$ ?
5. Show that each of the functions

$$G(x, y) = \begin{cases} 1, & x + y \geq 0, \\ 0, & x + y < 0, \end{cases}$$

$$G(x, y) = [x + y], \text{ the integral part of } x + y,$$

is continuous on the right, and continuous in each argument, but is not a (generalized) distribution function on  $R^2$ .

6. Let  $\mu$  be the Lebesgue-Stieltjes measure generated by a continuous distribution function. Show that if the set  $A$  is at most countable, then  $\mu(A) = 0$ .
7. Let  $c$  be the cardinal number of the continuum. Show that the cardinal number of the collection of Borel sets in  $R^n$  is  $c$ , whereas that of the collection of Lebesgue measurable sets is  $2^c$ .
8. Let  $(\Omega, \mathcal{F}, P)$  be a probability space and  $\mathcal{A}$  an algebra of subsets of  $\Omega$  such that  $\sigma(\mathcal{A}) = \mathcal{F}$ . Using the principle of appropriate sets, prove that for every  $\varepsilon > 0$  and  $B \in \mathcal{F}$  there is a set  $A \in \mathcal{A}$  such that

$$P(A \triangle B) \leq \varepsilon.$$

9. Let  $P$  be a probability measure on  $(R^n, \mathcal{B}(R^n))$ . Using Problem 8, show that, for every  $\varepsilon > 0$  and  $B \in \mathcal{B}(R^n)$ , there is a compact subset  $A$  of  $\mathcal{B}(R^n)$  such that  $A \subseteq B$  and

$$P(B \setminus A) \leq \varepsilon.$$

(This was used in the proof of Theorem 1.)

10. Verify the consistency of the measure defined by (21).

## §4. Random Variables. I

1. Let  $(\Omega, \mathcal{F})$  be a measurable space and let  $(R, \mathcal{B}(R))$  be the real line with the system  $\mathcal{B}(R)$  of Borel sets.

**Definition 1.** A real function  $\xi = \xi(\omega)$  defined on  $(\Omega, \mathcal{F})$  is an  $\mathcal{F}$ -measurable function, or a *random variable*, if

$$\{\omega: \xi(\omega) \in B\} \in \mathcal{F} \quad (1)$$

for every  $B \in \mathcal{B}(R)$ ; or, equivalently, if the inverse image

$$\xi^{-1}(B) \equiv \{\omega: \xi(\omega) \in B\}$$

is a measurable set in  $\Omega$ .

When  $(\Omega, \mathcal{F}) = (R^n, \mathcal{B}(R^n))$ , the  $\mathcal{B}(R^n)$ -measurable functions are called *Borel functions*.

The simplest example of a random variable is the indicator  $I_A(\omega)$  of an arbitrary (measurable) set  $A \in \mathcal{F}$ .

A random variable  $\xi$  that has a representation

$$\xi(\omega) = \sum_{i=1}^{\infty} x_i I_{A_i}(\omega), \quad (2)$$

where  $\sum A_i = \Omega$ ,  $A_i \in \mathcal{F}$ , is called *discrete*. If the sum in (2) is finite, the random variable is called *simple*.

With the same interpretation as in §4 of Chapter I, we may say that a random variable is a numerical property of an experiment, with a value depending on "chance." Here the requirement (1) of measurability is fundamental, for the following reason. If a probability measure  $P$  is defined on  $(\Omega, \mathcal{F})$ , it then makes sense to speak of the probability of the event  $\{\xi(\omega) \in B\}$  that the value of the random variable belongs to a Borel set  $B$ .

We introduce the following definitions.

**Definition 2.** A probability measure  $P_\xi$  on  $(R, \mathcal{B}(R))$  with

$$P_\xi(B) = P\{\omega: \xi(\omega) \in B\}, \quad B \in \mathcal{B}(R),$$

is called the *probability distribution of  $\xi$*  on  $(R, \mathcal{B}(R))$ .

**Definition 3.** The function

$$F_{\xi}(x) = P(\omega: \xi(\omega) \leq x), \quad x \in R,$$

is called the *distribution function* of  $\xi$ .

For a discrete random variable the measure  $P_{\xi}$  is concentrated on an at most countable set and can be represented in the form

$$P_{\xi}(B) = \sum_{\{k: x_k \in B\}} p(x_k), \quad (3)$$

where  $p(x_k) = P\{\xi = x_k\} = \Delta F_{\xi}(x_k)$ .

The converse is evidently true: If  $P_{\xi}$  is represented in the form (3) then  $\xi$  is a *discrete* random variable.

A random variable  $\xi$  is called *continuous* if its distribution function  $F_{\xi}(x)$  is continuous for  $x \in R$ .

A random variable  $\xi$  is called *absolutely continuous* if there is a nonnegative function  $f = f_{\xi}(x)$ , called its *density*, such that

$$F_{\xi}(x) = \int_{-\infty}^x f_{\xi}(y) dy, \quad x \in R, \quad (4)$$

(the integral can be taken in the Riemann sense, or more generally in that of Lebesgue; see §6 below).

**2.** To establish that a function  $\xi = \xi(\omega)$  is a random variable, we have to verify property (1) for all sets  $B \in \mathcal{F}$ . The following lemma shows that the class of such "test" sets can be considerably narrowed.

**Lemma 1.** Let  $\mathcal{E}$  be a system of sets such that  $\sigma(\mathcal{E}) = \mathcal{B}(R)$ . A necessary and sufficient condition that a function  $\xi = \xi(\omega)$  is  $\mathcal{F}$ -measurable is that

$$\{\omega: \xi(\omega) \in E\} \in \mathcal{F} \quad (5)$$

for all  $E \in \mathcal{E}$ .

**PROOF.** The necessity is evident. To prove the sufficiency we again use the principle of appropriate sets.

Let  $\mathcal{D}$  be the system of those Borel sets  $D$  in  $\mathcal{B}(R)$  for which  $\xi^{-1}(D) \in \mathcal{F}$ . The operation "form the inverse image" is easily shown to preserve the set-theoretic operations of union, intersection and complement:

$$\begin{aligned} \xi^{-1}\left(\bigcup_{\alpha} B_{\alpha}\right) &= \bigcup_{\alpha} \xi^{-1}(B_{\alpha}), \\ \xi^{-1}\left(\bigcap_{\alpha} B_{\alpha}\right) &= \bigcap_{\alpha} \xi^{-1}(B_{\alpha}), \\ \overline{\xi^{-1}(B_{\alpha})} &= \xi^{-1}(\overline{B_{\alpha}}). \end{aligned} \quad (6)$$



It follows that  $\mathcal{D}$  is a  $\sigma$ -algebra. Therefore

$$\mathcal{E} \subseteq \mathcal{D} \subseteq \mathcal{B}(R)$$

and

$$\sigma(\mathcal{E}) \subseteq \sigma(\mathcal{D}) = \mathcal{D} \subseteq \mathcal{B}(R).$$

But  $\sigma(E) = \mathcal{B}(R)$  and consequently  $\mathcal{D} = \mathcal{B}(R)$ .

**Corollary.** *A necessary and sufficient condition for  $\xi = \xi(\omega)$  to be a random variable is that*

$$\{\omega: \xi(\omega) < x\} \in \mathcal{F}$$

for every  $x \in R$ , or that

$$\{\omega: \xi(\omega) \leq x\} \in \mathcal{F}$$

for every  $x \in R$ .

The proof is immediate, since each of the systems

$$\mathcal{E}_1 = \{x: x < c, c \in R\},$$

$$\mathcal{E}_2 = \{x: x \leq c, c \in R\}$$

generates the  $\sigma$ -algebra  $\mathcal{B}(R)$ :  $\sigma(E_1) = \sigma(E_2) = \mathcal{B}(R)$  (see §2).

The following lemma makes it possible to construct random variables as functions of other random variables.

**Lemma 2.** *Let  $\varphi = \varphi(x)$  be a Borel function and  $\xi = \xi(\omega)$  a random variable. Then the composition  $\eta = \varphi \circ \xi$ , i.e. the function  $\eta(\omega) = \varphi(\xi(\omega))$ , is also a random variable.*

The proof follows from the equations

$$\{\omega: \eta(\omega) \in B\} = \{\omega: \varphi(\xi(\omega)) \in B\} = \{\omega: \xi(\omega) \in \varphi^{-1}(B)\} \in \mathcal{F} \quad (7)$$

for  $B \in \mathcal{B}(R)$ , since  $\varphi^{-1}(B) \in \mathcal{B}(R)$ .

Therefore if  $\xi$  is a random variable, so are, for examples,  $\xi^n$ ,  $\xi^+ = \max(\xi, 0)$ ,  $\xi^- = -\min(\xi, 0)$ , and  $|\xi|$ , since the functions  $x^n$ ,  $x^+$ ,  $x^-$  and  $|x|$  are Borel functions (Problem 4).

**3.** Starting from a given collection of random variables  $\{\xi_n\}$ , we can construct new functions, for example,  $\sum_{k=1}^{\infty} |\xi_k|$ ,  $\overline{\lim} \xi_n$ ,  $\underline{\lim} \xi_n$ , etc. Notice that in general such functions take values on the extended real line  $\bar{R} = [-\infty, \infty]$ . Hence it is advisable to extend the class of  $\mathcal{F}$ -measurable functions somewhat by allowing them to take the values  $\pm \infty$ .

**Definition 4.** A function  $\xi = \xi(\omega)$  defined on  $(\Omega, \mathcal{F})$  with values in  $\bar{R} = [-\infty, \infty]$  will be called an *extended random variable* if condition (1) is satisfied for every Borel set  $B \in \mathcal{B}(R)$ .

The following theorem, despite its simplicity, is the key to the construction of the Lebesgue integral (§6).

**Theorem 1.**

- (a) For every random variable  $\xi = \xi(\omega)$  (extended ones included) there is a sequence of simple random variables  $\xi_1, \xi_2, \dots$ , such that  $|\xi_n| \leq |\xi|$  and  $\xi_n(\omega) \rightarrow \xi(\omega)$ ,  $n \rightarrow \infty$ , for all  $\omega \in \Omega$ .
- (b) If also  $\xi(\omega) \geq 0$ , there is a sequence of simple random variables  $\xi_1, \xi_2, \dots$ , such that  $\xi_n(\omega) \uparrow \xi(\omega)$ ,  $n \rightarrow \infty$ , for all  $\omega \in \Omega$ .

**PROOF.** We begin by proving the second statement. For  $n = 1, 2, \dots$ , put

$$\xi_n(\omega) = \sum_{k=1}^{n2^n} \frac{k-1}{2^n} I_{k,n}(\omega) + n I_{\{\xi(\omega) \geq n\}}(\omega).$$

where  $I_{k,n}$  is the indicator of the set  $\{(k-1)/2^n \leq \xi(\omega) < k/2^n\}$ . It is easy to verify that the sequence  $\xi_n(\omega)$  so constructed is such that  $\xi_n(\omega) \uparrow \xi(\omega)$  for all  $\omega \in \Omega$ . The first statement follows from this if we merely observe that  $\xi$  can be represented in the form  $\xi = \xi^+ - \xi^-$ . This completes the proof of the theorem.

We next show that the class of extended random variables is closed under pointwise convergence. For this purpose, we note first that if  $\xi_1, \xi_2, \dots$  is a sequence of extended random variables, then  $\sup \xi_n$ ,  $\inf \xi_n$ ,  $\overline{\lim} \xi_n$  and  $\underline{\lim} \xi_n$  are also random variables (possibly extended). This follows immediately from

$$\{\omega: \sup \xi_n > x\} = \bigcup_n \{\omega: \xi_n > x\} \in \mathcal{F},$$

$$\{\omega: \inf \xi_n < x\} = \bigcup_n \{\omega: \xi_n < x\} \in \mathcal{F},$$

and

$$\overline{\lim} \xi_n = \inf_{n} \sup_{m \geq n} \xi_m, \quad \underline{\lim} \xi_n = \sup_{n} \inf_{m \geq n} \xi_m.$$

**Theorem 2.** Let  $\xi_1, \xi_2, \dots$  be a sequence of extended random variables and  $\xi(\omega) = \lim \xi_n(\omega)$ . Then  $\xi(\omega)$  is also an extended random variable.

The proof follows immediately from the remark above and the fact that

$$\begin{aligned} \{\omega: \xi(\omega) < x\} &= \{\omega: \lim \xi_n(\omega) < x\} \\ &= \{\omega: \overline{\lim} \xi_n(\omega) = \underline{\lim} \xi_n(\omega)\} \cap \{\overline{\lim} \xi_n(\omega) < x\} \\ &= \Omega \cap \{\overline{\lim} \xi_n(\omega) < x\} = \{\overline{\lim} \xi_n(\omega) < x\} \in \mathcal{F}. \end{aligned}$$

4. We mention a few more properties of the simplest functions of random variables considered on the measurable space  $(\Omega, \mathcal{F})$  and possibly taking values on the extended real line  $\bar{R} = [-\infty, \infty]$ .†

If  $\xi$  and  $\eta$  are random variables,  $\xi + \eta$ ,  $\xi - \eta$ ,  $\xi\eta$ , and  $\xi/\eta$  are also random variables (assuming that they are defined, i.e. that no indeterminate forms like  $\infty - \infty$ ,  $\infty/\infty$ ,  $a/0$  occur).

In fact, let  $\{\xi_n\}$  and  $\{\eta_n\}$  be sequences of random variables converging to  $\xi$  and  $\eta$  (see Theorem 1). Then

$$\begin{aligned}\xi_n \pm \eta_n &\rightarrow \xi \pm \eta, \\ \xi_n \eta_n &\rightarrow \xi \eta, \\ \frac{\xi_n}{\eta_n + \frac{1}{n} I_{\{\eta_n=0\}}(\omega)} &\rightarrow \frac{\xi}{\eta}.\end{aligned}$$

The functions on the left-hand sides of these relations are simple random variables. Therefore, by Theorem 2, the limit functions  $\xi \pm \eta$ ,  $\xi\eta$  and  $\xi/\eta$  are also random variables.

5. Let  $\xi$  be a random variable. Let us consider sets from  $\mathcal{F}$  of the form  $\{\omega: \xi(\omega) \in B\}$ ,  $B \in \mathcal{B}(R)$ . It is easily verified that they form a  $\sigma$ -algebra, called the  $\sigma$ -algebra generated by  $\xi$ , and denoted by  $\mathcal{F}_\xi$ .

If  $\varphi$  is a Borel function, it follows from Lemma 2 that the function  $\eta = \varphi \circ \xi$  is also a random variable, and in fact  $\mathcal{F}_\xi$ -measurable, i.e. such that

$$\{\omega: \eta(\omega) \in B\} \in \mathcal{F}_\xi, \quad B \in \mathcal{B}(R)$$

(see (7)). It turns out that the converse is also true.

**Theorem 3.** Let  $\eta$  be a  $\mathcal{F}_\xi$ -measurable random variable. Then there is a Borel function  $\varphi$  such that  $\eta = \varphi \circ \xi$ , i.e.  $\eta(\omega) = \varphi(\xi(\omega))$  for every  $\omega \in \Omega$ .

PROOF. Let  $\Phi_\xi$  be the class of  $\mathcal{F}_\xi$ -measurable functions  $\eta = \eta(\omega)$  and  $\tilde{\Phi}_\xi$  the class of  $\mathcal{F}_\xi$ -measurable functions representable in the form  $\varphi \circ \xi$ , where  $\varphi$  is a Borel function. It is clear that  $\tilde{\Phi}_\xi \subseteq \Phi_\xi$ . The conclusion of the theorem is that in fact  $\tilde{\Phi}_\xi = \Phi_\xi$ .

Let  $A \in \mathcal{F}_\xi$  and  $\eta(\omega) = I_A(\omega)$ . Let us show that  $\eta \in \tilde{\Phi}_\xi$ . In fact, if  $A \in \mathcal{F}_\xi$  there is a  $B \in \mathcal{B}(R)$  such that  $A = \{\omega: \xi(\omega) \in B\}$ . Let

$$\chi_B(x) = \begin{cases} 1, & x \in B, \\ 0, & x \notin B. \end{cases}$$

Then  $I_A(\omega) = \chi_B(\xi(\omega)) \in \tilde{\Phi}_\xi$ . Hence it follows that every simple  $\mathcal{F}_\xi$ -measurable function  $\sum_{i=1}^n c_i I_{A_i}(\omega)$ ,  $A_i \in \mathcal{F}_\xi$ , also belongs to  $\tilde{\Phi}_\xi$ .

† We shall assume the usual conventions about arithmetic operations in  $\bar{R}$ : if  $a \in R$  then  $a \pm \infty = \pm \infty$ ,  $a/\pm \infty = 0$ ;  $a \cdot \infty = \infty$  if  $a > 0$ , and  $a \cdot \infty = -\infty$  if  $a < 0$ ;  $0 \cdot (\pm \infty) = 0$ ,  $\infty + \infty = \infty$ ,  $-\infty - \infty = -\infty$ .

Now let  $\eta$  be an arbitrary  $\mathcal{F}_\xi$ -measurable function. By Theorem 1 there is a sequence of simple  $\mathcal{F}_\xi$ -measurable functions  $\{\eta_n\}$  such that  $\eta_n(\omega) \rightarrow \eta(\omega)$ ,  $n \rightarrow \infty$ ,  $\omega \in \Omega$ . As we just showed, there are Borel functions  $\varphi_n = \varphi_n(x)$  such that  $\eta_n(\omega) = \varphi_n(\xi(\omega))$ . Then  $\varphi_n(\xi(\omega)) \rightarrow \eta(\omega)$ ,  $n \rightarrow \infty$ ,  $\omega \in \Omega$ .

Let  $B$  denote the set  $\{x \in R: \lim_n \varphi_n(x) \text{ exists}\}$ . This is a Borel set. Therefore

$$\varphi(x) = \begin{cases} \lim_n \varphi_n(x), & x \in B, \\ 0, & x \notin B \end{cases}$$

is also a Borel function (see Problem 7).

But then it is evident that  $\eta(\omega) = \lim_n \varphi_n(\xi(\omega)) = \varphi(\xi(\omega))$  for all  $\omega \in \Omega$ . Consequently  $\mathcal{F}_\xi = \mathcal{F}_\eta$ .

**6.** Consider a measurable space  $(\Omega, \mathcal{F})$  and a finite or countably infinite decomposition  $\mathcal{D} = \{D_1, D_2, \dots\}$  of the space  $\Omega$ : namely,  $D_i \in \mathcal{F}$  and  $\sum_i D_i = \Omega$ . We form the algebra  $\mathcal{A}$  containing the empty set  $\emptyset$  and the sets of the form  $\sum_a D_a$ , where the sum is taken in the finite or countably infinite sense. It is evident that the system  $\mathcal{A}$  is a monotonic class, and therefore, according to Lemma 2, §2, chap. II, the algebra  $\mathcal{A}$  is at the same time a  $\sigma$ -algebra, denoted  $\sigma(\mathcal{D})$  and called the  $\sigma$ -algebra generated by the decomposition  $\mathcal{D}$ . Clearly  $\sigma(\mathcal{D}) \subseteq \mathcal{F}$ .

**Lemma 3.** Let  $\xi = \xi(\omega)$  be a  $\sigma(\mathcal{D})$ -measurable random variable. Then  $\xi$  is representable in the form

$$\xi(\omega) = \sum_{k=1}^{\infty} x_k I_{D_k}(\omega), \quad (8)$$

where  $x_k \in R$ ,  $k \geq 1$ , i.e.,  $\xi(\omega)$  is constant on the elements  $D_k$  of the decomposition,  $k \geq 1$ .

**PROOF.** Let us choose a set  $D_k$  and show that the  $\sigma(\mathcal{D})$ -measurable function  $\xi$  has a constant value on that set.

For this purpose, we write

$$x_k = \sup[c: D_k \cap \{\omega: \xi(\omega) < c\} = \emptyset].$$

Since  $\{\omega: \xi(\omega) < x_k\} = \bigcup_{r < x_k, r \text{ rational}} \{\omega: \xi(\omega) < r\}$ , we have  $D_k \cap \{\omega: \xi(\omega) < x_k\} = \emptyset$ .

Now let  $c > x_k$ . Then  $D_k \cap \{\omega: \xi(\omega) < c\} \neq \emptyset$ , and since the set  $\{\omega: \xi(\omega) < c\}$  has the form  $\sum_a D_a$ , where the sum is over a finite or countable collection of indices, we have

$$D_k \cap \{\omega: \xi(\omega) < c\} = D_k.$$

Hence, it follows that, for all  $c > x_k$ ,

$$D_k \cap \{\omega: \xi(\omega) \geq c\} = \emptyset,$$

and since  $\{\omega: \xi(\omega) > x_k\} = \bigcup_{r > x_k, r \text{ rational}} \{\omega: \xi(\omega) \geq r\}$ , we have

$$D_k \cap \{\omega: \xi(\omega) > x_k\} = \emptyset.$$

Consequently,  $D_k \cap \{\omega: \xi(\omega) \neq x_k\} = \emptyset$ , and therefore

$$D_k \subseteq \{\omega: \xi(\omega) = x_k\}$$

as required.

## 7. PROBLEMS

1. Show that the random variable  $\xi$  is continuous if and only if  $P(\xi = x) = 0$  for all  $x \in \mathbb{R}$ .
2. If  $|\xi|$  is  $\mathcal{F}$ -measurable, is it true that  $\xi$  is also  $\mathcal{F}$ -measurable?
3. Show that  $\xi = \xi(\omega)$  is an extended random variable if and only if  $\{\omega: \xi(\omega) \in \bar{B}\} \in \mathcal{F}$  for all  $\bar{B} \in \mathcal{B}(\bar{\mathbb{R}})$ .
4. Prove that  $x^+$ ,  $x^+ = \max(x, 0)$ ,  $x^- = -\min(x, 0)$ , and  $|x| = x^+ + x^-$  are Borel functions.
5. If  $\xi$  and  $\eta$  are  $\mathcal{F}$ -measurable, then  $\{\omega: \xi(\omega) = \eta(\omega)\} \in \mathcal{F}$ .
6. Let  $\xi$  and  $\eta$  be random variables on  $(\Omega, \mathcal{F})$ , and  $A \in \mathcal{F}$ . Then the function

$$\zeta(\omega) = \xi(\omega) \cdot I_A + \eta(\omega) I_{\bar{A}}$$

is also a random variable.

7. Let  $\xi_1, \dots, \xi_n$  be random variables and  $\varphi(x_1, \dots, x_n)$  a Borel function. Show that  $\varphi(\xi_1(\omega), \dots, \xi_n(\omega))$  is also a random variable.
8. Let  $\xi$  and  $\eta$  be random variables, both taking the values  $1, 2, \dots, N$ . Suppose that  $\mathcal{F}_\xi = \mathcal{F}$ . Show that there is a permutation  $(i_1, i_2, \dots, i_N)$  of  $(1, 2, \dots, N)$  such that  $\{\omega: \xi = j\} = \{\omega: \eta = i_j\}$  for  $j = 1, 2, \dots, N$ .

## §5. Random Elements

1. In addition to random variables, probability theory and its applications involve random objects of more general kinds, for example random points, vectors, functions, processes, fields, sets, measures, etc. In this connection it is desirable to have the concept of a random object of any kind.

**Definition 1.** Let  $(\Omega, \mathcal{F})$  and  $(E, \mathcal{E})$  be measurable spaces. We say that a function  $X = X(\omega)$ , defined on  $\Omega$  and taking values in  $E$ , is  $\mathcal{F}/\mathcal{E}$ -measurable, or is a *random element* (with values in  $E$ ), if

$$\{\omega: X(\omega) \in B\} \in \mathcal{F} \quad (1)$$

for every  $B \in \mathcal{E}$ . Random elements (with values in  $E$ ) are sometimes called  $E$ -valued random variables.

Let us consider some special cases.

If  $(E, \mathcal{E}) = (R, \mathcal{B}(R))$ , the definition of a random element is the same as the definition of a random variable (§4).

Let  $(E, \mathcal{E}) = (R^n, \mathcal{B}(R^n))$ . Then a random element  $X(\omega)$  is a "random point" in  $R^n$ . If  $\pi_k$  is the projection of  $R^n$  on the  $k$ th coordinate axis,  $X(\omega)$  can be represented in the form

$$X(\omega) = (\xi_1(\omega), \dots, \xi_n(\omega)), \quad (2)$$

where  $\xi_k = \pi_k \circ X$ .

It follows from (1) that  $\xi_k$  is an ordinary random variable. In fact, for  $B \in \mathcal{B}(R)$  we have

$$\begin{aligned} \{\omega: \xi_k(\omega) \in B\} &= \{\omega: \xi_1(\omega) \in R, \dots, \xi_{k-1}(\omega) \in R, \xi_k(\omega) \in B, \xi_{k+1}(\omega) \in R, \dots\} \\ &= \{\omega: X(\omega) \in (R \times \dots \times R \times B \times R \times \dots \times R)\} \in \mathcal{F}, \end{aligned}$$

since  $R \times \dots \times R \times B \times R \times \dots \times R \in \mathcal{B}(R^n)$ .

**Definition 2.** An ordered set  $(\eta_1(\omega), \dots, \eta_n(\omega))$  of random variables is called an  $n$ -dimensional random vector.

According to this definition, every random element  $X(\omega)$  with values in  $R^n$  is an  $n$ -dimensional random vector. The converse is also true: every random vector  $X(\omega) = (\xi_1(\omega), \dots, \xi_n(\omega))$  is a random element in  $R^n$ . In fact, if  $B_k \in \mathcal{B}(R)$ ,  $k = 1, \dots, n$ , then

$$\{\omega: X(\omega) \in (B_1 \times \dots \times B_n)\} = \bigcap_{k=1}^n \{\omega: \xi_k(\omega) \in B_k\} \in \mathcal{F}.$$

But  $\mathcal{B}(R^n)$  is the smallest  $\sigma$ -algebra containing the sets  $B_1 \times \dots \times B_n$ . Consequently we find immediately, by an evident generalization of Lemma 1 of §4, that whenever  $B \in \mathcal{B}(R^n)$ , the set  $\{\omega: X(\omega) \in B\}$  belongs to  $\mathcal{F}$ .

Let  $(E, \mathcal{E}) = (Z, \mathcal{B}(Z))$ , where  $Z$  is the set of complex numbers  $x + iy$ ,  $x, y \in R$ , and  $\mathcal{B}(Z)$  is the smallest  $\sigma$ -algebra containing the sets  $\{z: z = x + iy, a_1 < x \leq b_1, a_2 < y \leq b_2\}$ . It follows from the discussion above that a complex-valued random variable  $Z(\omega)$  can be represented as  $Z(\omega) = X(\omega) + iY(\omega)$ , where  $X(\omega)$  and  $Y(\omega)$  are random variables. Hence we may also call  $Z(\omega)$  a complex random variable.

Let  $(E, \mathcal{E}) = (R^T, \mathcal{B}(R^T))$ , where  $T$  is a subset of the real line. In this case every random element  $X = X(\omega)$  can evidently be represented as  $X = (\xi_t)_{t \in T}$  with  $\xi_t = \pi_t \circ X$ , and is called a random function with time domain  $T$ .

**Definition 3.** Let  $T$  be a subset of the real line. A set of random variables  $X = (\xi_t)_{t \in T}$  is called a random process† with time domain  $T$ .

† Or stochastic process (Translator).

If  $T = \{1, 2, \dots\}$  we call  $X = (\xi_1, \xi_2, \dots)$  a *random process with discrete time*, or a *random sequence*.

If  $T = [0, 1]$ ,  $(-\infty, \infty)$ ,  $[0, \infty)$ ,  $\dots$ , we call  $X = (\xi_t)_{t \in T}$  a *random process with continuous time*.

It is easy to show, by using the structure of the  $\sigma$ -algebra  $\mathcal{B}(R^T)$  (§2) that every random process  $X = (\xi_t)_{t \in T}$  (in the sense of Definition 3) is also a random function on the space  $(R^T, \mathcal{B}(R^T))$ .

**Definition 4.** Let  $X = (\xi_t)_{t \in T}$  be a random process. For each given  $\omega \in \Omega$  the function  $(\xi_t(\omega))_{t \in T}$  is said to be a *realization* or a *trajectory* of the process, corresponding to the outcome  $\omega$ .

The following definition is a natural generalization of Definition 2 of §4.

**Definition 5.** Let  $X = (\xi_t)_{t \in T}$  be a random process. The probability measure  $P_X$  on  $(R^T, \mathcal{B}(R^T))$  defined by

$$P_X(B) = P\{\omega: X(\omega) \in B\}, \quad B \in \mathcal{B}(R^T),$$

is called the *probability distribution of X*. The probabilities

$$P_{t_1, \dots, t_n}(B) \equiv P\{\omega: (\xi_{t_1}, \dots, \xi_{t_n}) \in B\}$$

with  $t_1 < t_2 < \dots < t_n$ ,  $t_i \in T$ , are called *finite-dimensional probabilities* (or *probability distributions*). The functions

$$F_{t_1, \dots, t_n}(x_1, \dots, x_n) \equiv P\{\omega: \xi_{t_1} \leq x_1, \dots, \xi_{t_n} \leq x_n\}$$

with  $t_1 < t_2 < \dots < t_n$ ,  $t_i \in T$ , are called *finite-dimensional distribution functions*.

Let  $(E, \mathcal{E}) = (C, \mathcal{B}_0(C))$ , where  $C$  is the space of continuous functions  $x = (x_t)_{t \in T}$  on  $T = [0, 1]$  and  $\mathcal{B}_0(C)$  is the  $\sigma$ -algebra generated by the open sets (§2). We show that every random element  $X$  on  $(C, \mathcal{B}_0(C))$  is also a random process with continuous trajectories in the sense of Definition 3.

In fact, according to §2 the set  $A = \{x \in C: x_t < a\}$  is open in  $\mathcal{B}_0(C)$ . Therefore

$$\{\omega: \xi_t(\omega) < a\} = \{\omega: X(\omega) \in A\} \in \mathcal{F}.$$

On the other hand, let  $X = (\xi_t(\omega))_{t \in T}$  be a random process (in the sense of Definition 3) whose trajectories are continuous functions for every  $\omega \in \Omega$ . According to (2.14),

$$\{x \in C: x \in S_\rho(x^0)\} = \bigcap_{t_k} \{x \in C: |x_{t_k} - x_{t_k}^0| < \rho\},$$

where  $t_k$  are the rational points of  $[0, 1]$ . Therefore

$$\{\omega: X(\omega) \in S_\rho(X^0\omega)\} = \bigcap_{t_k} \{\omega: |\xi_{t_k}(\omega) - \xi_{t_k}^0(\omega)| < \rho\} \in \mathcal{F},$$

and therefore we also have  $\{\omega: X(\omega) \in B\} \in \mathcal{F}$  for every  $B \in \mathcal{B}_0(C)$ .

Similar reasoning will show that every random element of the space  $(D, \mathcal{B}_0(D))$  can be considered as a random process with trajectories in the space of functions with no discontinuities of the second kind; and conversely.

2. Let  $(\Omega, \mathcal{F}, \mathbf{P})$  be a probability space and  $(E_\alpha, \mathcal{E}_\alpha)$  measurable spaces, where  $\alpha$  belongs to an (arbitrary) set  $\mathfrak{A}$ .

**Definition 6.** We say that the  $\mathcal{F}/\mathcal{E}_\alpha$ -measurable functions  $(X_\alpha(\omega))$ ,  $\alpha \in \mathfrak{A}$ , are independent (or collectively independent) if, for every finite set of indices  $\alpha_1, \dots, \alpha_n$  the random elements  $X_{\alpha_1}, \dots, X_{\alpha_n}$  are independent, i.e.

$$\mathbf{P}(X_{\alpha_1} \in B_{\alpha_1}, \dots, X_{\alpha_n} \in B_{\alpha_n}) = \mathbf{P}(X_{\alpha_1} \in B_{\alpha_1}) \cdots \mathbf{P}(X_{\alpha_n} \in B_{\alpha_n}), \quad (3)$$

where  $B_\alpha \in \mathcal{E}_\alpha$ .

Let  $\mathfrak{A} = \{1, 2, \dots, n\}$ , let  $\xi_\alpha$  be random variables, let  $\alpha \in \mathfrak{A}$  and let

$$F_\xi(x_1, \dots, x_n) = \mathbf{P}(\xi_1 \leq x_1, \dots, \xi_n \leq x_n)$$

be the  $n$ -dimensional distribution function of the random vector  $\xi = (\xi_1, \dots, \xi_n)$ . Let  $F_{\xi_i}(x_i)$  be the distribution functions of the random variables  $\xi_i$ ,  $i = 1, \dots, n$ .

**Theorem.** A necessary and sufficient condition for the random variables  $\xi_1, \dots, \xi_n$  to be independent is that

$$F_\xi(x_1, \dots, x_n) = F_{\xi_1}(x_1) \cdots F_{\xi_n}(x_n) \quad (4)$$

for all  $(x_1, \dots, x_n) \in \mathbb{R}^n$ .

**PROOF.** The necessity is evident. To prove the sufficiency we put  $a = (a_1, \dots, a_n)$ ,  $b = (b_1, \dots, b_n)$ ,

$$\begin{aligned} P_\xi(a, b] &= \mathbf{P}\{\omega: a_1 < \xi_1 \leq b_1, \dots, a_n < \xi_n \leq b_n\}, \\ P_{\xi_i}(a_i, b_i] &= \mathbf{P}\{a_i < \xi_i \leq b_i\}. \end{aligned}$$

Then

$$P_\xi(a, b] = \prod_{i=1}^n [F_{\xi_i}(b_i) - F_{\xi_i}(a_i)] = \prod_{i=1}^n P_{\xi_i}(a_i, b_i]$$

by (4) and (3.7), and therefore

$$\mathbf{P}\{\xi_1 \in I_1, \dots, \xi_n \in I_n\} = \prod_{i=1}^n \mathbf{P}\{\xi_i \in I_i\}, \quad (5)$$

where  $I_i = (a_i, b_i]$ .

We fix  $I_2, \dots, I_n$  and show that

$$\mathbf{P}\{\xi_1 \in B_1, \xi_2 \in I_2, \dots, \xi_n \in I_n\} = \mathbf{P}\{\xi_1 \in B_1\} \prod_{i=2}^n \mathbf{P}\{\xi_i \in I_i\} \quad (6)$$

for all  $B_1 \in \mathcal{B}(R)$ . Let  $\mathcal{M}$  be the collection of sets in  $\mathcal{B}(R)$  for which (6)



holds. Then  $\mathcal{M}$  evidently contains the algebra  $\mathcal{A}$  of sets consisting of sums of disjoint intervals of the form  $I_1 = (a_1, b_1]$ . Hence  $\mathcal{A} \subseteq \mathcal{M} \subseteq \mathcal{B}(R)$ . From the countable additivity (and therefore continuity) of probability measures it also follows that  $\mathcal{M}$  is a monotonic class. Therefore (see Subsection 1 of §2)

$$\mu(\mathcal{A}) \subseteq \mathcal{M} \subseteq \mathcal{B}(R).$$

But  $\mu(\mathcal{A}) = \sigma(\mathcal{A}) = \mathcal{B}(R)$  by Theorem 1 of §2. Therefore  $\mathcal{M} = \mathcal{B}(R)$ .

Thus (6) is established. Now fix  $B_1, I_3, \dots, I_n$ ; by the same method we can establish (6) with  $I_2$  replaced by the Borel set  $B_2$ . Continuing in this way, we can evidently arrive at the required equation,

$$P(\xi_1 \in B_1, \dots, \xi_n \in B_n) = P(\xi_1 \in B_1) \cdots P(\xi_n \in B_n),$$

where  $B_i \in \mathcal{B}(R)$ . This completes the proof of the theorem.

### 3. PROBLEMS

1. Let  $\xi_1, \dots, \xi_n$  be discrete random variables. Show that they are independent if and only if

$$P(\xi_1 = x_1, \dots, \xi_n = x_n) = \prod_{i=1}^n P(\xi_i = x_i)$$

for all real  $x_1, \dots, x_n$ .

2. Carry out the proof that every random function (in the sense of Definition 1) is a random process (in the sense of Definition 3) and conversely.
3. Let  $X_1, \dots, X_n$  be random elements with values in  $(E_1, \mathcal{E}_1), \dots, (E_n, \mathcal{E}_n)$ , respectively. In addition let  $(E'_1, \mathcal{E}'_1), \dots, (E'_n, \mathcal{E}'_n)$  be measurable spaces and let  $g_1, \dots, g_n$  be  $\mathcal{E}_1/\mathcal{E}'_1, \dots, \mathcal{E}_n/\mathcal{E}'_n$ -measurable functions, respectively. Show that if  $X_1, \dots, X_n$  are independent, the random elements  $g_1 \cdot X_1, \dots, g_n \cdot X_n$  are also independent.

## §6. Lebesgue Integral. Expectation

1. When  $(\Omega, \mathcal{F}, P)$  is a finite probability space and  $\xi = \xi(\omega)$  is a simple random variable,

$$\xi(\omega) = \sum_{k=1}^n x_k I_{A_k}(\omega), \quad (1)$$

the expectation  $E\xi$  was defined in §4 of Chapter I. The same definition of the expectation  $E\xi$  of a simple random variable  $\xi$  can be used for any probability space  $(\Omega, \mathcal{F}, P)$ . That is, we define

$$E\xi = \sum_{k=1}^n x_k P(A_k). \quad (2)$$

This definition is consistent (in the sense that  $E\xi$  is independent of the particular representation of  $\xi$  in the form (1)), as can be shown just as for finite probability spaces. The simplest properties of the expectation can be established similarly (see Subsection 5 of §4 of Chapter I).

In the present section we shall define and study the properties of the expectation  $E\xi$  of an arbitrary random variable. In the language of analysis,  $E\xi$  is merely the Lebesgue integral of the  $\mathcal{F}$ -measurable function  $\xi = \xi(\omega)$  with respect to the measure  $P$ . In addition to  $E\xi$  we shall use the notation  $\int_{\Omega} \xi(\omega)P(d\omega)$  or  $\int_{\Omega} \xi dP$ .

Let  $\xi = \xi(\omega)$  be a nonnegative random variable. We construct a sequence of simple nonnegative random variables  $\{\xi_n\}_{n \geq 1}$  such that  $\xi_n(\omega) \uparrow \xi(\omega)$ ,  $n \rightarrow \infty$ , for each  $\omega \in \Omega$  (see Theorem 1 in §4).

Since  $E\xi_n \leq E\xi_{n+1}$  (cf. Property 3) of Subsection 5, §4, Chapter I), the limit  $\lim_n E\xi_n$  exists, possibly with the value  $+\infty$ .

**Definition 1.** The *Lebesgue integral* of the nonnegative random variable  $\xi = \xi(\omega)$ , or its *expectation*, is

$$E\xi \equiv \lim_n E\xi_n. \quad (3)$$

To see that this definition is consistent, we need to show that the limit is independent of the choice of the approximating sequence  $\{\xi_n\}$ . In other words, we need to show that if  $\xi_n \uparrow \xi$  and  $\eta_m \uparrow \xi$ , where  $\{\eta_m\}$  is a sequence of simple functions, then

$$\lim_n E\xi_n = \lim_m E\eta_m. \quad (4)$$

**Lemma 1.** Let  $\eta$  and  $\xi_n$  be simple random variables,  $n \geq 1$ , and

$$\xi_n \uparrow \xi \geq \eta.$$

Then

$$\lim_n E\xi_n \geq E\eta. \quad (5)$$

PROOF. Let  $\varepsilon > 0$  and

$$A_n = \{\omega: \xi_n \geq \eta - \varepsilon\}.$$

It is clear that  $A_n \uparrow \Omega$  and

$$\xi_n = \xi_n I_{A_n} + \xi_n I_{\bar{A}_n} \geq \xi_n I_{A_n} \geq (\eta - \varepsilon) I_{A_n}.$$

Hence by the properties of the expectations of simple random variables we find that

$$\begin{aligned} E\xi_n &\geq E(\eta - \varepsilon) I_{A_n} = E\eta I_{A_n} - \varepsilon P(A_n) \\ &= E\eta - E\eta I_{\bar{A}_n} - \varepsilon P(A_n) \geq E\eta - CP(\bar{A}_n) - \varepsilon, \end{aligned}$$

where  $C = \max_{\omega} \eta(\omega)$ . Since  $\varepsilon$  is arbitrary, the required inequality (5) follows. It follows from this lemma that  $\lim_n E\xi_n \geq \lim_m E\eta_m$  and by symmetry  $\lim_m E\eta_m \geq \lim_n E\xi_n$ , which proves (4).

The following remark is often useful.

**Remark 1.** The expectation  $E\xi$  of the nonnegative random variable  $\xi$  satisfies

$$E\xi = \sup_{\{s \in S: s \leq \xi\}} E s, \quad (6)$$

where  $S = \{s\}$  is a set of simple random variables (Problem 1).

Thus the expectation is well defined for nonnegative random variables. We now consider the general case.

Let  $\xi$  be a random variable and  $\xi^+ = \max(\xi, 0)$ ,  $\xi^- = -\min(\xi, 0)$ .

**Definition 2.** We say that the expectation  $E\xi$  of the random variable  $\xi$  exists, or is defined, if at least one of  $E\xi^+$  and  $E\xi^-$  is finite:

$$\min(E\xi^+, E\xi^-) < \infty.$$

In this case we define

$$E\xi \equiv E\xi^+ - E\xi^-.$$

The expectation  $E\xi$  is also called the *Lebesgue integral* (of the function  $\xi$  with respect to the probability measure  $P$ ).

**Definition 3.** We say that the expectation of  $\xi$  is finite if  $E\xi^+ < \infty$  and  $E\xi^- < \infty$ .

Since  $|\xi| = \xi^+ + \xi^-$ , the finiteness of  $E\xi$ , or  $|E\xi| < \infty$ , is equivalent to  $E|\xi| < \infty$ . (In this sense one says that the Lebesgue integral is absolutely convergent.)

**Remark 2.** In addition to the expectation  $E\xi$ , significant numerical characteristics of a random variable  $\xi$  are the number  $E\xi^r$  (if defined) and  $E|\xi|^r$ ,  $r > 0$ , which are known as the *moment* of order  $r$  (or *rth moment*) and the *absolute moment* of order  $r$  (or *absolute rth moment*) of  $\xi$ .

**Remark 3.** In the definition of the Lebesgue integral  $\int_{\Omega} \xi(\omega) P(d\omega)$  given above, we suppose that  $P$  was a probability measure ( $P(\Omega) = 1$ ) and that the  $\mathcal{F}$ -measurable functions (random variables)  $\xi$  had values in  $R = (-\infty, \infty)$ . Suppose now that  $\mu$  is any measure defined on a measurable space  $(\Omega, \mathcal{F})$  and possibly taking the value  $+\infty$ , and that  $\xi = \xi(\omega)$  is an  $\mathcal{F}$ -measurable function with values in  $\bar{R} = [-\infty, \infty]$  (an extended random variable). In this case the Lebesgue integral  $\int_{\Omega} \xi(\omega) \mu(d\omega)$  is defined in the

same way: first, for nonnegative simple  $\xi$  (by (2) with  $\mathbf{P}$  replaced by  $\mu$ ), then for arbitrary nonnegative  $\xi$ , and in general by the formula

$$\int_{\Omega} \xi(\omega) \mu(d\omega) = \int_{\Omega} \xi^+ \mu(d\omega) - \int_{\Omega} \xi^- \mu(d\omega),$$

provided that no indeterminacy of the form  $\infty - \infty$  arises.

A case that is particularly important for mathematical analysis is that in which  $(\Omega, \mathcal{F}) = (R, \mathcal{B}(R))$  and  $\mu$  is Lebesgue measure. In this case the integral  $\int_R \xi(x) \mu(dx)$  is written  $\int_R \xi(x) dx$ , or  $\int_{-\infty}^{\infty} \xi(x) dx$ , or  $(L) \int_{-\infty}^{\infty} \xi(x) dx$  to emphasize its difference from the Riemann integral  $(R) \int_{-\infty}^{\infty} \xi(x) dx$ . If the measure  $\mu$  (Lebesgue–Stieltjes) corresponds to a generalized distribution function  $G = G(x)$ , the integral  $\int_R \xi(x) \mu(dx)$  is also called a *Lebesgue–Stieltjes integral* and is denoted by  $(L-S) \int_R \xi(x) G(dx)$ , a notation that distinguishes it from the corresponding Riemann–Stieltjes integral

$$(R-S) \int_R \xi(x) G(dx)$$

(see Subsection 10 below).

It will be clear from what follows (Property D) that if  $E\xi$  is defined then so is the expectation  $E(\xi I_A)$  for every  $A \in \mathcal{F}$ . The notations  $E(\xi; A)$  or  $\int_A \xi dP$  are often used for  $E(\xi I_A)$  or its equivalent,  $\int_{\Omega} \xi I_A dP$ . The integral  $\int_A \xi dP$  is called the *Lebesgue integral of  $\xi$  with respect to  $\mathbf{P}$  over the set  $A$* .

Similarly, we write  $\int_A \xi d\mu$  instead of  $\int_{\Omega} \xi \cdot I_A d\mu$  for an arbitrary measure  $\mu$ . In particular, if  $\mu$  is an  $n$ -dimensional Lebesgue–Stieltjes measure, and  $A = (a_1, b_1] \times \cdots \times (a_n, b_n]$ , we use the notation

$$\int_{a_1}^{b_1} \cdots \int_{a_n}^{b_n} \xi(x_1, \dots, x_n) \mu(dx_1 \cdots dx_n) \quad \text{instead of} \quad \int_A \xi d\mu.$$

If  $\mu$  is Lebesgue measure, we write simply  $dx_1 \cdots dx_n$  instead of  $\mu(dx_1, \dots, dx_n)$ .

## 2. Properties of the expectation $E\xi$ of the random variable $\xi$ .

A. Let  $c$  be a constant and let  $E\xi$  exist. Then  $E(c\xi)$  exists and

$$E(c\xi) = cE\xi.$$

B. Let  $\xi \leq \eta$ ; then

$$E\xi \leq E\eta$$

with the understanding that

$$\text{if } -\infty < E\xi \text{ then } -\infty < E\eta \text{ and } E\xi \leq E\eta$$

or

$$\text{if } E\eta < \infty \text{ then } E\xi < \infty \text{ and } E\xi \leq E\eta.$$

C. If  $E\xi$  exists then

$$|E\xi| \leq E|\xi|.$$

D. If  $E\xi$  exists then  $E(\xi I_A)$  exists for each  $A \in \mathcal{F}$ ; if  $E\xi$  is finite,  $E(\xi I_A)$  is finite.

E. If  $\xi$  and  $\eta$  are nonnegative random variables, or such that  $E|\xi| < \infty$  and  $E|\eta| < \infty$ , then

$$E(\xi + \eta) = E\xi + E\eta.$$

(See Problem 2 for a generalization.)

Let us establish A-E.

A. This is obvious for simple random variables. Let  $\xi \geq 0$ ,  $\xi_n \uparrow \xi$ , where  $\xi_n$  are simple random variables and  $c \geq 0$ . Then  $c\xi_n \uparrow c\xi$  and therefore

$$E(c\xi) = \lim E(c\xi_n) = c \lim E\xi_n = cE\xi.$$

In the general case we need to use the representation  $\xi = \xi^+ - \xi^-$  and notice that  $(c\xi)^+ = c\xi^+$ ,  $(c\xi)^- = c\xi^-$  when  $c \geq 0$ , whereas when  $c < 0$ ,  $(c\xi)^+ = -c\xi^-$ ,  $(c\xi)^- = -c\xi^+$ .

B. If  $0 \leq \xi \leq \eta$ , then  $E\xi$  and  $E\eta$  are defined and the inequality  $E\xi \leq E\eta$  follows directly from (6). Now let  $E\xi > -\infty$ ; then  $E\xi^- < \infty$ . If  $\xi \leq \eta$ , we have  $\xi^+ \leq \eta^+$  and  $\xi^- \geq \eta^-$ . Therefore  $E\eta^- \leq E\xi^- < \infty$ ; consequently  $E\eta$  is defined and  $E\xi = E\xi^+ - E\xi^- \leq E\eta^+ - E\eta^- = E\eta$ . The case when  $E\eta < \infty$  can be discussed similarly.

C. Since  $-|\xi| \leq \xi \leq |\xi|$ , Properties A and B imply

$$-E|\xi| \leq E\xi \leq E|\xi|,$$

i.e.  $|E\xi| \leq E|\xi|$ .

D. This follows from B and

$$(\xi I_A)^+ = \xi^+ I_A \leq \xi^+, \quad (\xi I_A)^- = \xi^- I_A \leq \xi^-.$$

E. Let  $\xi \geq 0$ ,  $\eta \geq 0$ , and let  $\{\xi_n\}$  and  $\{\eta_n\}$  be sequences of simple functions such that  $\xi_n \uparrow \xi$  and  $\eta_n \uparrow \eta$ . Then  $E(\xi_n + \eta_n) = E\xi_n + E\eta_n$  and

$$E(\xi_n + \eta_n) \uparrow E(\xi + \eta), \quad E\xi_n \uparrow E\xi, \quad E\eta_n \uparrow E\eta$$

and therefore  $E(\xi + \eta) = E\xi + E\eta$ . The case when  $E|\xi| < \infty$  and  $E|\eta| < \infty$  reduces to this if we use the facts that

$$\xi = \xi^+ - \xi^-, \quad \eta = \eta^+ - \eta^-, \quad \xi^+ \leq |\xi|, \quad \xi^- \leq |\xi|,$$

and

$$\eta^+ \leq |\eta|, \quad \eta^- \leq |\eta|.$$

The following group of statements about expectations involve the notion of "P-almost surely." We say that a property holds "P-almost surely" if there is a set  $\mathcal{N} \in \mathcal{F}$  with  $P(\mathcal{N}) = 0$  such that the property holds for every point  $\omega$  of  $\Omega \setminus \mathcal{N}$ . Instead of "P-almost surely" we often say "P-almost everywhere" or simply "almost surely" (a.s.) or "almost everywhere" (a.e.).

F. If  $\xi = 0$  (a.s.) then  $E\xi = 0$ .

In fact, if  $\xi$  is a simple random variable,  $\xi = \sum x_k I_{A_k}(\omega)$  and  $x_k \neq 0$ , we have  $P(A_k) = 0$  by hypothesis and therefore  $E\xi = 0$ . If  $\xi \geq 0$  and  $0 \leq s \leq \xi$ , where  $s$  is a simple random variable, then  $s = 0$  (a.s.) and consequently  $Es = 0$  and  $E\xi = \sup_{\{s \in S: s \leq \xi\}} Es = 0$ . The general case follows from this by means of the representation  $\xi = \xi^+ - \xi^-$  and the facts that  $\xi^+ \leq |\xi|$ ,  $\xi^- \leq |\xi|$ , and  $|\xi| = 0$  (a.s.).

G. If  $\xi = \eta$  (a.s.) and  $E|\xi| < \infty$ , then  $E|\eta| < \infty$  and  $E\xi = E\eta$  (see also Problem 3).

In fact, let  $\mathcal{N} = \{\omega: \xi \neq \eta\}$ . Then  $P(\mathcal{N}) = 0$  and  $\xi = \xi I_{\mathcal{N}} + \xi I_{\mathcal{N}^c}$ ,  $\eta = \eta I_{\mathcal{N}} + \eta I_{\mathcal{N}^c} = \eta I_{\mathcal{N}} + \xi I_{\mathcal{N}^c}$ . By properties E and F, we have  $E\xi = E\xi I_{\mathcal{N}} + E\xi I_{\mathcal{N}^c} = E\eta I_{\mathcal{N}} + E\xi I_{\mathcal{N}^c}$ . But  $E\eta I_{\mathcal{N}} = 0$ , and therefore  $E\xi = E\eta I_{\mathcal{N}^c} + E\eta I_{\mathcal{N}} = E\eta$ , by Property E.

H. Let  $\xi \geq 0$  and  $E\xi = 0$ . Then  $\xi = 0$  (a.s.).

For the proof, let  $A = \{\omega: \xi(\omega) > 0\}$ ,  $A_n = \{\omega: \xi(\omega) \geq 1/n\}$ . It is clear that  $A_n \uparrow A$  and  $0 \leq \xi \cdot I_{A_n} \leq \xi \cdot I_A$ . Hence, by Property B,

$$0 \leq E\xi I_{A_n} \leq E\xi = 0.$$

Consequently

$$0 = E\xi I_{A_n} \geq \frac{1}{n} P(A_n)$$

and therefore  $P(A_n) = 0$  for all  $n \geq 1$ . But  $P(A) = \lim P(A_n)$  and therefore  $P(A) = 0$ .

I. Let  $\xi$  and  $\eta$  be such that  $E|\xi| < \infty$ ,  $E|\eta| < \infty$  and  $E(\xi I_A) \leq E(\eta I_A)$  for all  $A \in \mathcal{F}$ . Then  $\xi \leq \eta$  (a.s.).

In fact, let  $B = \{\omega: \xi(\omega) > \eta(\omega)\}$ . Then  $E(\eta I_B) \leq E(\xi I_B) \leq E(\eta I_B)$  and therefore  $E(\xi I_B) = E(\eta I_B)$ . By Property E, we have  $E((\xi - \eta) I_B) = 0$  and by Property H we have  $(\xi - \eta) I_B = 0$  (a.s.), whence  $P(B) = 0$ .

J. Let  $\xi$  be an extended random variable and  $E|\xi| < \infty$ . Then  $|\xi| < \infty$  (a.s.). In fact, let  $A = \{\omega: |\xi(\omega)| = \infty\}$  and  $P(A) > 0$ . Then  $E|\xi| \geq E(|\xi| I_A) = \infty \cdot P(A) = \infty$ , which contradicts the hypothesis  $E|\xi| < \infty$ . (See also Problem 4.)

3. Here we consider the fundamental theorems on *taking limits* under the expectation sign (or the Lebesgue integral sign).

**Theorem 1** (On Monotone Convergence). *Let  $\eta, \xi, \xi_1, \xi_2, \dots$  be random variables.*

(a) *If  $\xi_n \geq \eta$  for all  $n \geq 1$ ,  $E\eta > -\infty$ , and  $\xi_n \uparrow \xi$ , then*

$$E\xi_n \uparrow E\xi.$$

(b) *If  $\xi_n \leq \eta$  for all  $n \geq 1$ ,  $E\eta < \infty$ , and  $\xi_n \downarrow \xi$ , then*

$$E\xi_n \downarrow E\xi.$$

**PROOF.** (a) First suppose that  $\eta \geq 0$ . For each  $k \geq 1$  let  $\{\xi_k^{(n)}\}_{n \geq 1}$  be a sequence of simple functions such that  $\xi_k^{(n)} \uparrow \xi_k$ ,  $n \rightarrow \infty$ . Put  $\zeta^{(n)} = \max_{1 \leq k \leq n} \xi_k^{(n)}$ . Then

$$\zeta^{(n-1)} \leq \zeta^{(n)} = \max_{1 \leq k \leq n} \xi_k^{(n)} \leq \max_{1 \leq k \leq n} \xi_k = \xi_n.$$

Let  $\zeta = \lim_n \zeta^{(n)}$ . Since

$$\xi_k^{(n)} \leq \zeta^{(n)} \leq \xi_n$$

for  $1 \leq k \leq n$ , we find by taking limits as  $n \rightarrow \infty$  that

$$\xi_k \leq \zeta \leq \xi$$

for every  $k \geq 1$  and therefore  $\xi = \zeta$ .

The random variables  $\zeta^{(n)}$  are simple and  $\zeta^{(n)} \uparrow \zeta$ . Therefore

$$E\xi = E\zeta = \lim E\zeta^{(n)} \leq \lim E\xi_n.$$

On the other hand, it is obvious, since  $\xi_n \leq \xi_{n+1} \leq \xi$ , that

$$\lim E\xi_n \leq E\xi.$$

Consequently  $\lim E\xi_n = E\xi$ .

Now let  $\eta$  be any random variable with  $E\eta > -\infty$ .

If  $E\eta = \infty$  then  $E\xi_n = E\xi = \infty$  by Property B, and our proposition is proved. Let  $E\eta < \infty$ . Then instead of  $E\eta > -\infty$  we find  $E|\eta| < \infty$ . It is clear that  $0 \leq \xi_n - \eta \uparrow \xi - \eta$  for all  $\omega \in \Omega$ . Therefore by what has been established,  $E(\xi_n - \eta) \uparrow E(\xi - \eta)$  and therefore (by Property E and Problem 2)

$$E\xi_n - E\eta \uparrow E\xi - E\eta.$$

But  $E|\eta| < \infty$ , and therefore  $E\xi_n \uparrow E\xi$ ,  $n \rightarrow \infty$ .

The proof of (b) follows from (a) if we replace the original variables by their negatives.

**Corollary.** *Let  $\{\eta_n\}_{n \geq 1}$  be a sequence of nonnegative random variables. Then*

$$E \sum_{n=1}^{\infty} \eta_n = \sum_{n=1}^{\infty} E\eta_n.$$

The proof follows from Property E (see also Problem 2), the monotone convergence theorem, and the remark that

$$\sum_{n=1}^k \eta_n \uparrow \sum_{n=1}^{\infty} \eta_n, \quad k \rightarrow \infty.$$

**Theorem 2** (Fatou's Lemma). *Let  $\eta, \xi_1, \xi_2, \dots$  be random variables.*

(a) *If  $\xi_n \geq \eta$  for all  $n \geq 1$  and  $E\eta > -\infty$ , then*

$$E \underline{\lim} \xi_n \leq \underline{\lim} E \xi_n.$$

(b) *If  $\xi_n \leq \eta$  for all  $n \geq 1$  and  $E\eta < \infty$ , then*

$$\overline{\lim} E \xi_n \leq E \overline{\lim} \xi_n.$$

(c) *If  $|\xi_n| \leq \eta$  for all  $n \geq 1$  and  $E\eta < \infty$ , then*

$$E \underline{\lim} \xi_n \leq \underline{\lim} E \xi_n \leq \overline{\lim} E \xi_n \leq E \overline{\lim} \xi_n. \quad (7)$$

PROOF. (a) Let  $\zeta_n = \inf_{m \geq n} \xi_m$ ; then

$$\underline{\lim} \xi_n = \lim_{n \rightarrow \infty} \inf_{m \geq n} \xi_m = \lim_{n \rightarrow \infty} \zeta_n.$$

It is clear that  $\zeta_n \uparrow \underline{\lim} \xi_n$  and  $\zeta_n \geq \eta$  for all  $n \geq 1$ . Then by Theorem 1

$$E \underline{\lim} \xi_n = E \lim_{n \rightarrow \infty} \zeta_n = \lim_{n \rightarrow \infty} E \zeta_n = \underline{\lim}_{n \rightarrow \infty} E \zeta_n \leq \underline{\lim}_{n \rightarrow \infty} E \xi_n,$$

which establishes (a). The second conclusion follows from the first. The third is a corollary of the first two.

**Theorem 3** (Lebesgue's Theorem on Dominated Convergence). *Let  $\eta, \xi, \xi_1, \xi_2, \dots$  be random variables such that  $|\xi_n| \leq \eta$ ,  $E\eta < \infty$  and  $\xi_n \rightarrow \xi$  (a.s.). Then  $E|\xi| < \infty$ ,*

$$E \xi_n \rightarrow E \xi \quad (8)$$

and

$$E|\xi_n - \xi| \rightarrow 0 \quad (9)$$

as  $n \rightarrow \infty$ .

PROOF. Formula (7) is valid by Fatou's lemma. By hypothesis,  $\underline{\lim} \xi_n = \overline{\lim} \xi_n = \xi$  (a.s.). Therefore by Property G,

$$E \underline{\lim} \xi_n = \underline{\lim} E \xi_n = \overline{\lim} E \xi_n = E \overline{\lim} \xi_n = E \xi,$$

which establishes (8). It is also clear that  $|\xi| \leq \eta$ . Hence  $E|\xi| < \infty$ .

Conclusion (9) can be proved in the same way if we observe that  $|\xi_n - \xi| \leq 2\eta$ .



**Corollary.** Let  $\eta, \xi, \xi_1, \dots$  be random variables such that  $|\xi_n| \leq \eta, \xi_n \rightarrow \xi$  (a.s.) and  $E\eta^p < \infty$  for some  $p > 0$ . Then  $E|\xi|^p < \infty$  and  $E|\xi - \xi_n|^p \rightarrow 0, n \rightarrow \infty$ .

For the proof, it is sufficient to observe that

$$|\xi| \leq \eta, |\xi - \xi_n|^p \leq (|\xi| + |\xi_n|)^p \leq (2\eta)^p.$$

The condition " $|\xi_n| \leq \eta, E\eta < \infty$ " that appears in Fatou's lemma and the dominated convergence theorem, and ensures the validity of formulas (7)–(9), can be somewhat weakened. In order to be able to state the corresponding result (Theorem 4), we introduce the following definition.

**Definition 4.** A family  $\{\xi_n\}_{n \geq 1}$  of random variables is said to be *uniformly integrable* if

$$\sup_n \int_{\{|\xi_n| > c\}} |\xi_n| P(d\omega) \rightarrow 0, \quad c \rightarrow \infty, \quad (10)$$

or, in a different notation,

$$\sup_n E[|\xi_n| I_{\{|\xi_n| > c\}}] \rightarrow 0, \quad c \rightarrow \infty. \quad (11)$$

It is clear that if  $\xi_n, n \geq 1$ , satisfy  $|\xi_n| \leq \eta, E\eta < \infty$ , then the family  $\{\xi_n\}_{n \geq 1}$  is uniformly integrable.

**Theorem 4.** Let  $\{\xi_n\}_{n \geq 1}$  be a uniformly integrable family of random variables. Then

(a)  $E \liminf \xi_n \leq \liminf E\xi_n \leq \overline{\lim} E\xi_n \leq E \overline{\lim} \xi_n$ .

(b) If in addition  $\xi_n \rightarrow \xi$  (a.s.) then  $\xi$  is integrable and

$$\begin{aligned} E\xi_n &\rightarrow E\xi, & n &\rightarrow \infty, \\ E|\xi_n - \xi| &\rightarrow 0, & n &\rightarrow \infty. \end{aligned}$$

**PROOF.** (a) For every  $c > 0$

$$E\xi_n = E[\xi_n I_{\{\xi_n < -c\}}] + E[\xi_n I_{\{\xi_n \geq -c\}}]. \quad (12)$$

By uniform integrability, for every  $\varepsilon > 0$  we can take  $c$  so large that

$$\sup_n |E[\xi_n I_{\{\xi_n < -c\}}]| < \varepsilon. \quad (13)$$

By Fatou's lemma,

$$\liminf E[\xi_n I_{\{\xi_n \geq -c\}}] \geq E[\liminf \xi_n I_{\{\xi_n \geq -c\}}].$$

But  $\xi_n I_{\{\xi_n \geq -c\}} \geq \xi_n$  and therefore

$$\liminf E[\xi_n I_{\{\xi_n \geq -c\}}] \geq E[\liminf \xi_n]. \quad (14)$$

From (12)–(14) we obtain

$$\lim E\xi_n \geq E[\lim \xi_n] - \varepsilon.$$

Since  $\varepsilon > 0$  is arbitrary, it follows that  $\lim E\xi_n \geq E \lim \xi_n$ . The inequality with upper limits,  $\overline{\lim} E\xi_n \leq E \overline{\lim} \xi_n$ , is proved similarly.

Conclusion (b) can be deduced from (a) as in Theorem 3.

The deeper significance of the concept of uniform integrability is revealed by the following theorem, which gives a necessary and sufficient condition for taking limits under the expectation sign.

**Theorem 5.** *Let  $0 \leq \xi_n \rightarrow \xi$  and  $E\xi_n < \infty$ . Then  $E\xi_n \rightarrow E\xi < \infty$  if and only if the family  $\{\xi_n\}_{n \geq 1}$  is uniformly integrable.*

PROOF. The sufficiency follows from conclusion (b) of Theorem 4. For the proof of the necessity we consider the (at most countable) set

$$A = \{a: P(\xi = a) > 0\}.$$

Then we have  $\xi_n I_{\{\xi_n < a\}} \rightarrow \xi I_{\{\xi < a\}}$  for each  $a \notin A$ , and the family

$$\{\xi_n I_{\{\xi_n < a\}}\}_{n \geq 1}$$

is uniformly integrable. Hence, by the sufficiency part of the theorem, we have  $E\xi_n I_{\{\xi_n < a\}} \rightarrow E\xi I_{\{\xi < a\}}$ ,  $a \notin A$ , and therefore

$$E\xi_n I_{\{\xi_n \geq a\}} \rightarrow E\xi I_{\{\xi \geq a\}}, \quad a \notin A, \quad n \rightarrow \infty. \quad (15)$$

Take an  $\varepsilon > 0$  and choose  $a_0 \notin A$  so large that  $E\xi I_{\{\xi \geq a_0\}} < \varepsilon/2$ ; then choose  $N_0$  so large that

$$E\xi_n I_{\{\xi_n \geq a_0\}} \leq E\xi I_{\{\xi \geq a_0\}} + \varepsilon/2$$

for all  $n \geq N_0$ , and consequently  $E\xi_n I_{\{\xi_n \geq a_0\}} \leq \varepsilon$ . Then choose  $a_1 \geq a_0$  so large that  $E\xi I_{\{\xi \geq a_1\}} \leq \varepsilon$  for all  $n \leq N_0$ . Then we have

$$\sup_n E\xi_n I_{\{\xi_n \geq a_1\}} \leq \varepsilon,$$

which establishes the uniform integrability of the family  $\{\xi_n\}_{n \geq 1}$  of random variables.

4. Let us notice some tests for uniform integrability.

We first observe that if  $\{\xi_n\}$  is a family of uniformly integrable random variables, then

$$\sup_n E|\xi_n| < \infty. \quad (16)$$

In fact, for a given  $\varepsilon > 0$  and sufficiently large  $c > 0$

$$\begin{aligned} \sup_n E|\xi_n| &= \sup_n [E(|\xi_n| I_{\{|\xi_n| \geq c\}}) + E(|\xi_n| I_{\{|\xi_n| < c\}})] \\ &\leq \sup_n E(|\xi_n| I_{\{|\xi_n| \geq c\}}) + \sup_n E(|\xi_n| I_{\{|\xi_n| < c\}}) \leq \varepsilon + c, \end{aligned}$$

which establishes (16).

It turns out that (16) together with a condition of uniform continuity is necessary and sufficient for uniform integrability.

**Lemma 2.** *A necessary and sufficient condition for a family  $\{\xi_n\}_{n \geq 1}$  of random variables to be uniformly integrable is that  $E|\xi_n|$ ,  $n \geq 1$ , are uniformly bounded (i.e., (16) holds) and that  $E\{|\xi_n| I_A\}$ ,  $n \geq 1$ , are uniformly absolutely continuous (i.e.  $\sup E\{|\xi_n| I_A\} \rightarrow 0$  when  $P(A) \rightarrow 0$ ).*

**PROOF.** *Necessity.* Condition (16) was verified above. Moreover,

$$\begin{aligned} E\{|\xi_n| I_A\} &= E\{|\xi_n| I_{A \cap \{|\xi_n| \geq c\}}\} + E\{|\xi_n| I_{A \cap \{|\xi_n| < c\}}\} \\ &\leq E\{|\xi_n| I_{\{|\xi_n| \geq c\}}\} + cP(A). \end{aligned} \quad (17)$$

Take  $c$  so large that  $\sup_n E\{|\xi_n| I_{\{|\xi_n| \geq c\}}\} \leq \varepsilon/2$ . Then if  $P(A) \leq \varepsilon/2c$ , we have

$$\sup_n E\{|\xi_n| I_A\} \leq \varepsilon$$

by (17). This establishes the uniform absolute continuity.

*Sufficiency.* Let  $\varepsilon > 0$  and  $\delta > 0$  be chosen so that  $P(A) < \delta$  implies that  $E(|\xi_n| I_A) \leq \varepsilon$ , uniformly in  $n$ . Since

$$E|\xi_n| \geq E|\xi_n| I_{\{|\xi_n| \geq c\}} \geq cP\{|\xi_n| \geq c\}$$

for every  $c > 0$  (cf. Chebyshev's inequality), we have

$$\sup_n P\{|\xi_n| \geq c\} \leq \frac{1}{c} \sup_n E|\xi_n| \rightarrow 0, \quad c \rightarrow \infty,$$

and therefore, when  $c$  is sufficiently large, any set  $\{|\xi_n| \geq c\}$ ,  $n \geq 1$ , can be taken as  $A$ . Therefore  $\sup E(|\xi_n| I_{\{|\xi_n| \geq c\}}) \leq \varepsilon$ , which establishes the uniform integrability. This completes the proof of the lemma.

The following proposition provides a simple sufficient condition for uniform integrability.

**Lemma 3.** *Let  $\xi_1, \xi_2, \dots$  be a sequence of integrable random variables and  $G = G(t)$  a nonnegative increasing function, defined for  $t \geq 0$ , such that*

$$\lim_{t \rightarrow \infty} \frac{G(t)}{t} = \infty. \quad (18)$$

$$\sup_n E[G(|\xi_n|)] < \infty. \quad (19)$$

*Then the family  $\{\xi_n\}_{n \geq 1}$  is uniformly integrable.*

PROOF. Let  $\varepsilon > 0$ ,  $M = \sup_n \mathbf{E}[G(|\xi_n|)]$ ,  $a = M/\varepsilon$ . Take  $c$  so large that  $G(t)/t \geq a$  for  $t \geq c$ . Then

$$\mathbf{E}[|\xi_n| I_{\{|\xi_n| \geq c\}}] \leq \frac{1}{a} \mathbf{E}[G(|\xi_n|) \cdot I_{\{|\xi_n| \geq c\}}] \leq \frac{M}{a} = \varepsilon$$

uniformly for  $n \geq 1$ .

5. If  $\xi$  and  $\eta$  are independent simple random variables, we can show, as in Subsection 5 of §4 of Chapter I, that  $\mathbf{E}\xi\eta = \mathbf{E}\xi \cdot \mathbf{E}\eta$ . Let us now establish a similar proposition in the general case (see also Problem 5).

**Theorem 6.** Let  $\xi$  and  $\eta$  be independent random variables with  $\mathbf{E}|\xi| < \infty$ ,  $\mathbf{E}|\eta| < \infty$ . Then  $\mathbf{E}|\xi\eta| < \infty$  and

$$\mathbf{E}\xi\eta = \mathbf{E}\xi \cdot \mathbf{E}\eta. \quad (20)$$

PROOF. First let  $\xi \geq 0$ ,  $\eta \geq 0$ . Put

$$\begin{aligned} \xi_n &= \sum_{k=0}^{\infty} \frac{k}{n} I_{\{k/n \leq \xi(\omega) < (k+1)/n\}}, \\ \eta_n &= \sum_{k=0}^{\infty} \frac{k}{n} I_{\{k/n \leq \eta(\omega) < (k+1)/n\}}. \end{aligned}$$

Then  $\xi_n \leq \xi$ ,  $|\xi_n - \xi| \leq 1/n$  and  $\eta_n \leq \eta$ ,  $|\eta_n - \eta| \leq 1/n$ . Since  $\mathbf{E}\xi < \infty$  and  $\mathbf{E}\eta < \infty$ , it follows from Lebesgue's dominated convergence theorem that

$$\lim \mathbf{E}\xi_n = \mathbf{E}\xi, \quad \lim \mathbf{E}\eta_n = \mathbf{E}\eta.$$

Moreover, since  $\xi$  and  $\eta$  are independent,

$$\begin{aligned} \mathbf{E}\xi_n\eta_n &= \sum_{k,l \geq 0} \frac{kl}{n^2} \mathbf{E} I_{\{k/n \leq \xi < (k+1)/n\}} I_{\{l/n \leq \eta < (l+1)/n\}} \\ &= \sum_{k,l \geq 0} \frac{kl}{n^2} \mathbf{E} I_{\{k/n \leq \xi < (k+1)/n\}} \cdot \mathbf{E} I_{\{l/n \leq \eta < (l+1)/n\}} = \mathbf{E}\xi_n \cdot \mathbf{E}\eta_n. \end{aligned}$$

Now notice that

$$\begin{aligned} |\mathbf{E}\xi\eta - \mathbf{E}\xi_n\eta_n| &\leq \mathbf{E}|\xi\eta - \xi_n\eta_n| \leq \mathbf{E}[|\xi| \cdot |\eta - \eta_n|] \\ &\quad + \mathbf{E}[|\eta_n| \cdot |\xi - \xi_n|] \leq \frac{1}{n} \mathbf{E}\xi + \frac{1}{n} \mathbf{E}\left(\eta + \frac{1}{n}\right) \rightarrow 0, \quad n \rightarrow \infty. \end{aligned}$$

Therefore  $\mathbf{E}\xi\eta = \lim_n \mathbf{E}\xi_n\eta_n = \lim \mathbf{E}\xi_n \cdot \lim \mathbf{E}\eta_n = \mathbf{E}\xi \cdot \mathbf{E}\eta$ , and  $\mathbf{E}\xi\eta < \infty$ .

The general case reduces to this one if we use the representations  $\xi = \xi^+ - \xi^-$ ,  $\eta = \eta^+ - \eta^-$ ,  $\xi\eta = \xi^+\eta^+ - \xi^+\eta^- - \xi^-\eta^+ + \xi^-\eta^-$ . This completes the proof.

6. The inequalities for expectations that we develop in this subsection are regularly used both in probability theory and in analysis.

**Chebyshev's Inequality.** Let  $\xi$  be a nonnegative random variable. Then for every  $\varepsilon > 0$

$$P(\xi \geq \varepsilon) \leq \frac{E\xi}{\varepsilon}. \quad (21)$$

The proof follows immediately from

$$E\xi \geq E[\xi \cdot I_{(\xi \geq \varepsilon)}] \geq \varepsilon E I_{(\xi \geq \varepsilon)} = \varepsilon P(\xi \geq \varepsilon).$$

From (21) we can obtain the following variant of Chebyshev's inequality: If  $\xi$  is any random variable then

$$P(\xi \geq \varepsilon) \leq \frac{E\xi^2}{\varepsilon^2} \quad (22)$$

and

$$P(|\xi - E\xi| \geq \varepsilon) \leq \frac{V\xi}{\varepsilon^2}, \quad (23)$$

where  $V\xi = E(\xi - E\xi)^2$  is the variance of  $\xi$ .

**The Cauchy-Bunyakovskii Inequality.** Let  $\xi$  and  $\eta$  satisfy  $E\xi^2 < \infty, E\eta^2 < \infty$ . Then  $E|\xi\eta| < \infty$  and

$$(E|\xi\eta|)^2 \leq E\xi^2 \cdot E\eta^2. \quad (24)$$

PROOF. Suppose that  $E\xi^2 > 0, E\eta^2 > 0$ . Then, with  $\tilde{\xi} = \xi/\sqrt{E\xi^2}, \tilde{\eta} = \eta/\sqrt{E\eta^2}$ , we find, since  $2|\tilde{\xi}\tilde{\eta}| \leq \tilde{\xi}^2 + \tilde{\eta}^2$ , that

$$2E|\tilde{\xi}\tilde{\eta}| \leq E\tilde{\xi}^2 + E\tilde{\eta}^2 = 2,$$

i.e.  $E|\tilde{\xi}\tilde{\eta}| \leq 1$ , which establishes (24).

On the other hand if, say,  $E\xi^2 \equiv 0$ , then  $\xi = 0$  (a.s.) by Property I, and then  $E\xi\eta = 0$  by Property F, i.e. (24) is still satisfied.

**Jensen's Inequality.** Let the Borel function  $g = g(x)$  be convex downward and  $E|\xi| < \infty$ . Then

$$g(E\xi) \leq Eg(\xi). \quad (25)$$

PROOF. If  $g = g(x)$  is convex downward, for each  $x_0 \in R$  there is a number  $\lambda(x_0)$  such that

$$g(x) \geq g(x_0) + (x - x_0) \cdot \lambda(x_0) \quad (26)$$

for all  $x \in R$ . Putting  $x = \xi$  and  $x_0 = E\xi$ , we find from (26) that

$$g(\xi) \geq g(E\xi) + (\xi - E\xi) \cdot \lambda(E\xi),$$

and consequently  $Eg(\xi) \geq g(E\xi)$ .

A whole series of useful inequalities can be derived from Jensen's inequality. We obtain the following one as an example.

**Lyapunov's Inequality.** If  $0 < s < t$ ,

$$(E|\xi|^s)^{1/s} \leq (E|\xi|^t)^{1/t}. \quad (27)$$

To prove this, let  $r = t/s$ . Then, putting  $\eta = |\xi|^s$  and applying Jensen's inequality to  $g(x) = |x|^r$ , we obtain  $|E\eta|^r \leq E|\eta|^r$ , i.e.

$$(E|\xi|^s)^{t/s} \leq E|\xi|^t,$$

which establishes (27).

The following chain of inequalities among absolute moments in a consequence of Lyapunov's inequality:

$$E|\xi| \leq (E|\xi|^2)^{1/2} \leq \dots \leq (E|\xi|^n)^{1/n}. \quad (28)$$

**Hölder's Inequality.** Let  $1 < p < \infty$ ,  $1 < q < \infty$ , and  $(1/p) + (1/q) = 1$ . If  $E|\xi|^p < \infty$  and  $E|\eta|^q < \infty$ , then  $E|\xi\eta| < \infty$  and

$$E|\xi\eta| \leq (E|\xi|^p)^{1/p}(E|\eta|^q)^{1/q}. \quad (29)$$

If  $E|\xi|^p = 0$  or  $E|\eta|^q = 0$ , (29) follows immediately as for the Cauchy-Bunyakovskii inequality (which is the special case  $p = q = 2$  of Hölder's inequality).

Now let  $E|\xi|^p > 0$ ,  $E|\eta|^q > 0$  and

$$\tilde{\xi} = \frac{\xi}{(E|\xi|^p)^{1/p}}, \quad \tilde{\eta} = \frac{\eta}{(E|\eta|^q)^{1/q}}.$$

We apply the inequality

$$x^a y^b \leq ax + by, \quad (30)$$

which holds for positive  $x, y, a, b$  and  $a + b = 1$ , and follows immediately from the concavity of the logarithm:

$$\ln[ax + by] \geq a \ln x + b \ln y = \ln x^a y^b.$$

Then, putting  $x = |\tilde{\xi}|^p$ ,  $y = |\tilde{\eta}|^q$ ,  $a = 1/p$ ,  $b = 1/q$ , we find that

$$|\tilde{\xi}\tilde{\eta}| \leq \frac{1}{p}|\tilde{\xi}|^p + \frac{1}{q}|\tilde{\eta}|^q,$$

whence

$$E|\tilde{\xi}\tilde{\eta}| \leq \frac{1}{p}E|\tilde{\xi}|^p + \frac{1}{q}E|\tilde{\eta}|^q = \frac{1}{p} + \frac{1}{q} = 1.$$

This establishes (29).

**Minkowski's Inequality.** If  $E|\xi|^p < \infty$ ,  $E|\eta|^p < \infty$ ,  $1 \leq p < \infty$ , then we have  $E|\xi + \eta|^p < \infty$  and

$$(E|\xi + \eta|^p)^{1/p} \leq (E|\xi|^p)^{1/p} + (E|\eta|^p)^{1/p}. \quad (31)$$

We begin by establishing the following inequality: if  $a, b > 0$  and  $p \geq 1$ , then

$$(a + b)^p \leq 2^{p-1}(a^p + b^p). \quad (32)$$

In fact, consider the function  $F(x) = (a + x)^p - 2^{p-1}(a^p + x^p)$ . Then

$$F'(x) = p(a + x)^{p-1} - 2^{p-1}px^{p-1},$$

and since  $p \geq 1$ , we have  $F'(a) = 0$ ,  $F'(x) > 0$  for  $x < a$  and  $F'(x) < 0$  for  $x > a$ . Therefore

$$F(b) \leq \max F(x) = F(a) = 0,$$

from which (32) follows.

According to this inequality,

$$|\xi + \eta|^p \leq (|\xi| + |\eta|)^p \leq 2^{p-1}(|\xi|^p + |\eta|^p) \quad (33)$$

and therefore if  $E|\xi|^p < \infty$  and  $E|\eta|^p < \infty$  it follows that  $E|\xi + \eta|^p < \infty$ .

If  $p = 1$ , inequality (31) follows from (33).

Now suppose that  $p > 1$ . Take  $q > 1$  so that  $(1/p) + (1/q) = 1$ . Then

$$|\xi + \eta|^p = |\xi + \eta| \cdot |\xi + \eta|^{p-1} \leq |\xi| \cdot |\xi + \eta|^{p-1} + |\eta| \cdot |\xi + \eta|^{p-1}. \quad (34)$$

Notice that  $(p-1)q = p$ . Consequently

$$E(|\xi + \eta|^{p-1})^q = E|\xi + \eta|^p < \infty,$$

and therefore by Hölder's inequality

$$\begin{aligned} E(|\xi| |\xi + \eta|^{p-1}) &\leq (E|\xi|^p)^{1/p} (E|\xi + \eta|^{(p-1)q})^{1/q} \\ &= (E|\xi|^p)^{1/p} (E|\xi + \eta|^p)^{1/q} < \infty. \end{aligned}$$

In the same way,

$$E(|\eta| |\xi + \eta|^{p-1}) \leq (E|\eta|^p)^{1/p} (E|\xi + \eta|^p)^{1/q}.$$

Consequently, by (34),

$$E|\xi + \eta|^p \leq (E|\xi + \eta|^p)^{1/q} ((E|\xi|^p)^{1/p} + (E|\eta|^p)^{1/p}). \quad (35)$$

If  $E|\xi + \eta|^p = 0$ , the desired inequality (31) is evident. Now let  $E|\xi + \eta|^p > 0$ . Then we obtain

$$(E|\xi + \eta|^p)^{1-(1/q)} \leq (E|\xi|^p)^{1/p} + (E|\eta|^p)^{1/p}$$

from (35), and (31) follows since  $1 - (1/q) = 1/p$ .

7. Let  $\xi$  be a random variable for which  $E\xi$  is defined. Then, by Property D, the set function

$$Q(A) \equiv \int_A \xi \, dP, \quad A \in \mathcal{F}, \quad (36)$$

is well defined. Let us show that this function is countably additive.

First suppose that  $\xi$  is nonnegative. If  $A_1, A_2, \dots$  are pairwise disjoint sets from  $\mathcal{F}$  and  $A = \sum A_n$ , the corollary to Theorem 1 implies that

$$\begin{aligned} Q(A) &= E(\xi \cdot I_A) = E(\xi \cdot I_{\sum A_n}) = E(\sum \xi \cdot I_{A_n}) \\ &= \sum E(\xi \cdot I_{A_n}) = \sum Q(A_n). \end{aligned}$$

If  $\xi$  is an arbitrary random variable for which  $E\xi$  is defined, the countable additivity of  $Q(A)$  follows from the representation

$$Q(A) = Q^+(A) - Q^-(A), \quad (37)$$

where

$$Q^+(A) = \int_A \xi^+ \, dP, \quad Q^-(A) = \int_A \xi^- \, dP,$$

together with the countable additivity for nonnegative random variables and the fact that  $\min(Q^+(\Omega), Q^-(\Omega)) < \infty$ .

Thus if  $E\xi$  is defined, the set function  $Q = Q(A)$  is a signed measure—a countably additive set function representable as  $Q = Q_1 - Q_2$ , where at least one of the measures  $Q_1$  and  $Q_2$  is finite.

We now show that  $Q = Q(A)$  has the following important property of *absolute continuity* with respect to  $P$ :

$$\text{if } P(A) = 0 \quad \text{then } Q(A) = 0 \quad (A \in \mathcal{F})$$

(this property is denoted by the abbreviation  $Q \ll P$ ).

To prove the sufficiency we consider nonnegative random variables. If  $\xi = \sum_{k=1}^n x_k I_{A_k}$  is a simple nonnegative random variable and  $P(A) = 0$ , then

$$Q(A) = E(\xi \cdot I_A) = \sum_{k=1}^n x_k P(A_k \cap A) = 0.$$

If  $\{\xi_n\}_{n \geq 1}$  is a sequence of nonnegative simple functions such that  $\xi_n \uparrow \xi \geq 0$ , then the theorem on monotone convergence shows that

$$Q(A) = E(\xi \cdot I_A) = \lim E(\xi_n \cdot I_A) = 0,$$

since  $E(\xi_n \cdot I_A) = 0$  for all  $n \geq 1$  and  $A$  with  $P(A) = 0$ .

Thus the Lebesgue integral  $Q(A) = \int_A \xi \, dP$ , considered as a function of sets  $A \in \mathcal{F}$ , is a signed measure that is absolutely continuous with respect to  $P$  ( $Q \ll P$ ). It is quite remarkable that the converse is also valid.



**Radon–Nikodým Theorem.** Let  $(\Omega, \mathcal{F})$  be a measurable space,  $\mu$  a  $\sigma$ -finite measure, and  $\lambda$  a signed measure (i.e.,  $\lambda = \lambda_1 - \lambda_2$ , where at least one of the measures  $\lambda_1$  and  $\lambda_2$  is finite) which is absolutely continuous with respect to  $\mu$ . Then there is an  $\mathcal{F}$ -measurable function  $f = f(\omega)$  with values in  $\bar{R} = [-\infty, \infty]$  such that

$$\lambda(A) = \int_A f(\omega) \mu(d\omega), \quad A \in \mathcal{F}. \quad (38)$$

The function  $f(\omega)$  is unique up to sets of  $\mu$ -measure zero: if  $h = h(\omega)$  is another  $\mathcal{F}$ -measurable function such that  $\lambda(A) = \int_A h(\omega) \mu(d\omega)$ ,  $A \in \mathcal{F}$ , then  $\mu\{\omega: f(\omega) \neq h(\omega)\} = 0$ .

If  $\lambda$  is a measure, then  $f = f(\omega)$  has its values in  $\bar{R}^+ = [0, \infty]$ .

**Remark.** The function  $f = f(\omega)$  in the representation (38) is called the *Radon–Nikodým derivative* or the *density* of the measure  $\lambda$  with respect to  $\mu$ , and denoted by  $d\lambda/d\mu$  or  $(d\lambda/d\mu)(\omega)$ .

The Radon–Nikodým theorem, which we quote without proof, will play a key role in the construction of conditional expectations (§7).

8. If  $\xi = \sum_{i=1}^n x_i I_{A_i}$  is a simple random variable,

$$Eg(\xi) = \sum g(x_i) P(A_i) = \sum g(x_i) \Delta F_\xi(x_i). \quad (39)$$

In other words, in order to calculate the expectation of a function of the (simple) random variable  $\xi$  it is unnecessary to know the probability measure  $P$  completely; it is enough to know the probability distribution  $P_\xi$  or, equivalently, the distribution function  $F_\xi$  of  $\xi$ .

The following important theorem generalizes this property.

**Theorem 7** (Change of Variables in a Lebesgue Integral). Let  $(\Omega, \mathcal{F})$  and  $(E, \mathcal{E})$  be measurable spaces and  $X = X(\omega)$  an  $\mathcal{F}/\mathcal{E}$ -measurable function with values in  $E$ . Let  $P$  be a probability measure on  $(\Omega, \mathcal{F})$  and  $P_X$  the probability measure on  $(E, \mathcal{E})$  induced by  $X = X(\omega)$ :

$$P_X(A) = P\{\omega: X(\omega) \in A\}, \quad A \in \mathcal{E}. \quad (40)$$

Then

$$\int_A g(x) P_X(dx) = \int_{X^{-1}(A)} g(X(\omega)) P(d\omega), \quad A \in \mathcal{E}, \quad (41)$$

for every  $\mathcal{E}$ -measurable function  $g = g(x)$ ,  $x \in E$  (in the sense that if one integral exists, the other is well defined, and the two are equal).

**PROOF.** Let  $A \in \mathcal{E}$  and  $g(x) = I_B(x)$ , where  $B \in \mathcal{E}$ . Then (41) becomes

$$P_X(AB) = P(X^{-1}(A) \cap X^{-1}(B)), \quad (42)$$

which follows from (40) and the observation that  $X^{-1}(A) \cap X^{-1}(B) = X^{-1}(A \cap B)$ .

It follows from (42) that (41) is valid for nonnegative simple functions  $g = g(x)$ , and therefore, by the monotone convergence theorem, also for all nonnegative  $\mathcal{E}$ -measurable functions.

In the general case we need only represent  $g$  as  $g^+ - g^-$ . Then, since (41) is valid for  $g^+$  and  $g^-$ , if (for example)  $\int_A g^+(x)P_X(dx) < \infty$ , we have

$$\int_{X^{-1}(A)} g^+(X(\omega))P(d\omega) < \infty$$

also, and therefore the existence of  $\int_A g(x)P_X(dx)$  implies the existence of  $\int_{X^{-1}(A)} g(X(\omega))P(d\omega)$ .

**Corollary.** Let  $(E, \mathcal{E}) = (R, \mathcal{B}(R))$  and let  $\xi = \xi(\omega)$  be a random variable with probability distribution  $P_\xi$ . Then if  $g = g(x)$  is a Borel function and either of the integrals  $\int_A g(x)P_\xi(dx)$  or  $\int_{\xi^{-1}(A)} g(\xi(\omega))P(d\omega)$  exists, we have

$$\int_A g(x)P_\xi(dx) = \int_{\xi^{-1}(A)} g(\xi(\omega))P(d\omega).$$

In particular, for  $A = R$  we obtain

$$Eg(\xi(\omega)) = \int_{\Omega} g(\xi(\omega))P(d\omega) = \int_R g(x)P_\xi(dx). \quad (43)$$

The measure  $P_\xi$  can be uniquely reconstructed from the distribution function  $F_\xi$  (Theorem 1 of §3). Hence the Lebesgue integral  $\int_R g(x)P_\xi(dx)$  is often denoted by  $\int_R g(x)F_\xi(dx)$  and called a *Lebesgue-Stieltjes integral* (with respect to the measure corresponding to the distribution function  $F_\xi(x)$ ).

Let us consider the case when  $F_\xi(x)$  has a density  $f_\xi(x)$ , i.e. let

$$F_\xi(x) = \int_{-\infty}^x f_\xi(y) dy, \quad (44)$$

where  $f_\xi = f_\xi(x)$  is a nonnegative Borel function and the integral is a Lebesgue integral with respect to Lebesgue measure on the set  $(-\infty, x]$  (see Remark 2 in Subsection 1). With the assumption of (44), formula (43) takes the form

$$Eg(\xi(\omega)) = \int_{-\infty}^{\infty} g(x)f_\xi(x) dx, \quad (45)$$

where the integral is the Lebesgue integral of the function  $g(x)f_\xi(x)$  with respect to Lebesgue measure. In fact, if  $g(x) = I_B(x)$ ,  $B \in \mathcal{B}(R)$ , the formula becomes

$$P_\xi(B) = \int_B f_\xi(x) dx, \quad B \in \mathcal{B}(R); \quad (46)$$

its correctness follows from Theorem 1 of §3 and the formula

$$F_{\xi}(b) - F_{\xi}(a) = \int_a^b f_{\xi}(x) dx.$$

In the general case, the proof is the same as for Theorem 7.

9. Let us consider the special case of measurable spaces  $(\Omega, \mathcal{F})$  with a measure  $\mu$ , where  $\Omega = \Omega_1 \times \Omega_2$ ,  $\mathcal{F} = \mathcal{F}_1 \otimes \mathcal{F}_2$ , and  $\mu = \mu_1 \times \mu_2$  is the direct product of measures  $\mu_1$  and  $\mu_2$  (i.e., the measure on  $\mathcal{F}$  such that

$$\mu_1 \times \mu_2(A \times B) = \mu_1(A)\mu_2(B), \quad A \in \mathcal{F}_1, \quad B \in \mathcal{F}_2;$$

the existence of this measure follows from the proof of Theorem 8).

The following theorem plays the same role as the theorem on the reduction of a double Riemann integral to an iterated integral.

**Theorem 8 (Fubini's Theorem).** *Let  $\xi = \xi(\omega_1, \omega_2)$  be an  $\mathcal{F}_1 \otimes \mathcal{F}_2$ -measurable function, integrable with respect to the measure  $\mu_1 \times \mu_2$ :*

$$\int_{\Omega_1 \times \Omega_2} |\xi(\omega_1, \omega_2)| d(\mu_1 \times \mu_2) < \infty. \quad (47)$$

Then the integrals  $\int_{\Omega_1} \xi(\omega_1, \omega_2) \mu_1(d\omega_1)$  and  $\int_{\Omega_2} \xi(\omega_1, \omega_2) \mu_2(d\omega_2)$

(1) are defined for all  $\omega_1$  and  $\omega_2$ ;

(2) are respectively  $\mathcal{F}_2$ - and  $\mathcal{F}_1$ -measurable functions with

$$\begin{aligned} \mu_2 \left\{ \omega_2 : \int_{\Omega_1} |\xi(\omega_1, \omega_2)| \mu_1(d\omega_1) = \infty \right\} &= 0, \\ \mu_1 \left\{ \omega_1 : \int_{\Omega_2} |\xi(\omega_1, \omega_2)| \mu_2(d\omega_2) = \infty \right\} &= 0 \end{aligned} \quad (48)$$

and (3)

$$\begin{aligned} \int_{\Omega_1 \times \Omega_2} \xi(\omega_1, \omega_2) d(\mu_1 \times \mu_2) &= \int_{\Omega_1} \left[ \int_{\Omega_2} \xi(\omega_1, \omega_2) \mu_2(d\omega_2) \right] \mu_1(d\omega_1) \\ &= \int_{\Omega_2} \left[ \int_{\Omega_1} \xi(\omega_1, \omega_2) \mu_1(d\omega_1) \right] \mu_2(d\omega_2). \end{aligned} \quad (49)$$

**PROOF.** We first show that  $\xi_{\omega_1}(\omega_2) = \xi(\omega_1, \omega_2)$  is  $\mathcal{F}_2$ -measurable with respect to  $\omega_2$ , for each  $\omega_1 \in \Omega_1$ .

Let  $F \in \mathcal{F}_1 \otimes \mathcal{F}_2$  and  $\xi(\omega_1, \omega_2) = I_F(\omega_1, \omega_2)$ . Let

$$F_{\omega_1} = \{\omega_2 \in \Omega_2 : (\omega_1, \omega_2) \in F\}$$

be the cross-section of  $F$  at  $\omega_1$ , and let  $\mathcal{C}_{\omega_1} = \{F \in \mathcal{F} : F_{\omega_1} \in F_2\}$ . We must show that  $\mathcal{C}_{\omega_1} = \mathcal{F}$  for every  $\omega_1$ .

If  $F = A \times B$ ,  $A \in \mathcal{F}_1$ ,  $B \in \mathcal{F}_2$ , then

$$(A \times B)_{\omega_1} = \begin{cases} B & \text{if } \omega_1 \in A, \\ \emptyset & \text{if } \omega_1 \notin A. \end{cases}$$

Hence rectangles with measurable sides belong to  $\mathcal{C}_{\omega_1}$ . In addition, if  $F \in \mathcal{F}$ , then  $(\bar{F})_{\omega_1} = \overline{F_{\omega_1}}$ , and if  $\{F^n\}_{n \geq 1}$  are sets in  $\mathcal{F}$ , then  $(\bigcup F^n)_{\omega_1} = \bigcup F^n_{\omega_1}$ . It follows that  $\mathcal{C}_{\omega_1} = \mathcal{F}$ .

Now let  $\xi(\omega_1, \omega_2) \geq 0$ . Then, since the function  $\xi(\omega_1, \omega_2)$  is  $\mathcal{F}_2$ -measurable for each  $\omega_1$ , the integral  $\int_{\Omega_2} \xi(\omega_1, \omega_2) \mu_2(d\omega_2)$  is defined. Let us show that this integral is an  $\mathcal{F}_1$ -measurable function and

$$\int_{\Omega_1} \left[ \int_{\Omega_2} \xi(\omega_1, \omega_2) \mu_2(d\omega_2) \right] \mu_1(d\omega_1) = \int_{\Omega_1 \times \Omega_2} \xi(\omega_1, \omega_2) d(\mu_1 \times \mu_2). \quad (50)$$

Let us suppose that  $\xi(\omega_1, \omega_2) = I_{A \times B}(\omega_1, \omega_2)$ ,  $A \in \mathcal{F}_1$ ,  $B \in \mathcal{F}_2$ . Then since  $I_{A \times B}(\omega_1, \omega_2) = I_A(\omega_1) I_B(\omega_2)$ , we have

$$\int_{\Omega_2} I_{A \times B}(\omega_1, \omega_2) \mu_2(d\omega_2) = I_A(\omega_1) \int_{\Omega_2} I_B(\omega_2) \mu_2(d\omega_2) \quad (51)$$

and consequently the integral on the left of (51) is an  $\mathcal{F}_1$ -measurable function.

Now let  $\xi(\omega_1, \omega_2) = I_F(\omega_1, \omega_2)$ ,  $F \in \mathcal{F} = \mathcal{F}_1 \otimes \mathcal{F}_2$ . Let us show that the integral  $f(\omega_1) = \int_{\Omega_2} I_F(\omega_1, \omega_2) \mu_2(d\omega_2)$  is  $\mathcal{F}$ -measurable. For this purpose we put  $\mathcal{C} = \{F \in \mathcal{F} : f(\omega_1) \text{ is } \mathcal{F}_1\text{-measurable}\}$ . According to what has been proved, the set  $A \times B$  belongs to  $\mathcal{C}$  ( $A \in \mathcal{F}_1$ ,  $B \in \mathcal{F}_2$ ) and therefore the algebra  $\mathcal{A}$  consisting of finite sums of disjoint sets of this form also belongs to  $\mathcal{C}$ . It follows from the monotone convergence theorem that  $\mathcal{C}$  is a monotonic class,  $\mathcal{C} = \mu(\mathcal{C})$ . Therefore, because of the inclusions  $\mathcal{A} \subseteq \mathcal{C} \subseteq \mathcal{F}$  and Theorem 1 of §2, we have  $\mathcal{F} = \sigma(\mathcal{A}) = \mu(\mathcal{A}) \subseteq \mu(\mathcal{C}) = \mathcal{C} \subseteq \mathcal{F}$ , i.e.  $\mathcal{C} = \mathcal{F}$ .

Finally, if  $\xi(\omega_1, \omega_2)$  is an arbitrary nonnegative  $\mathcal{F}$ -measurable function, the  $\mathcal{F}_1$ -measurability of the integral  $\int_{\Omega_2} \xi(\omega_1, \omega_2) \mu_2(d\omega)$  follows from the monotone convergence theorem and Theorem 2 of §4.

Let us now show that the measure  $\mu = \mu_1 \times \mu_2$  defined on  $\mathcal{F} = \mathcal{F}_1 \otimes \mathcal{F}_2$ , with the property  $(\mu_1 \times \mu_2)(A \times B) = \mu_1(A) \cdot \mu_2(B)$ ,  $A \in \mathcal{F}_1$ ,  $B \in \mathcal{F}_2$ , actually exists and is unique.

For  $F \in \mathcal{F}$  we put

$$\mu(F) = \int_{\Omega_1} \left[ \int_{\Omega_2} I_{F_{\omega_1}}(\omega_2) \mu_2(d\omega_2) \right] \mu_1(d\omega_1).$$

As we have shown, the inner integral is an  $\mathcal{F}_1$ -measurable function, and consequently the set function  $\mu(F)$  is actually defined for  $F \in \mathcal{F}$ . It is clear

that if  $F = A \times B$ , then  $\mu(A \times B) = \mu_1(A)\mu_2(B)$ . Now let  $\{F^n\}$  be disjoint sets from  $\mathcal{F}$ . Then

$$\begin{aligned}\mu(\sum F^n) &= \int_{\Omega_1} \left[ \int_{\Omega_2} I_{(\sum F^n)\omega_1}(\omega_2) \mu_2(d\omega_2) \right] \mu_1(d\omega_1) \\ &= \int_{\Omega_1} \sum_n \left[ \int_{\Omega_2} I_{F^n\omega_1}(\omega_2) \mu_2(d\omega_2) \right] \mu_1(d\omega_1) \\ &= \sum_n \int_{\Omega_1} \left[ \int_{\Omega_2} I_{F^n\omega_1}(\omega_2) \mu_2(d\omega_2) \right] \mu_1(d\omega_1) = \sum_n \mu(F^n),\end{aligned}$$

i.e.  $\mu$  is a ( $\sigma$ -finite) measure on  $\mathcal{F}$ .

It follows from Carathéodory's theorem that this measure  $\mu$  is the unique measure with the property that  $\mu(A \times B) = \mu_1(A)\mu_2(B)$ .

We can now establish (50). If  $\xi(\omega_1, \omega_2) = I_{A \times B}(\omega_1, \omega_2)$ ,  $A \in \mathcal{F}_1$ ,  $B \in \mathcal{F}_2$ , then

$$\int_{\Omega_1 \times \Omega_2} I_{A \times B}(\omega_1, \omega_2) d(\mu_1 \times \mu_2) = (\mu_1 \times \mu_2)(A \times B), \quad (52)$$

and since  $I_{A \times B}(\omega_1, \omega_2) = I_A(\omega_1)I_B(\omega_2)$ , we have

$$\begin{aligned}\int_{\Omega_1} \left[ \int_{\Omega_2} I_{A \times B}(\omega_1, \omega_2) \mu_2(d\omega_2) \right] \mu_1(d\omega_1) \\ = \int_{\Omega_1} \left[ I_A(\omega_1) \int_{\Omega_2} I_B(\omega_1, \omega_2) \mu_2(d\omega_2) \right] \mu_1(d\omega_1) = \mu_1(A)\mu_2(B).\end{aligned} \quad (53)$$

But, by the definition of  $\mu_1 \times \mu_2$ ,

$$(\mu_1 \times \mu_2)(A \times B) = \mu_1(A)\mu_2(B).$$

Hence it follows from (52) and (53) that (50) is valid for  $\xi(\omega_1, \omega_2) = I_{A \times B}(\omega_1, \omega_2)$ .

Now let  $\xi(\omega_1, \omega_2) = I_F(\omega_1, \omega_2)$ ,  $F \in \mathcal{F}$ . The set function

$$\lambda(F) = \int_{\Omega_1 \times \Omega_2} I_F(\omega_1, \omega_2) d(\mu_1 \times \mu_2), \quad F \in \mathcal{F},$$

is evidently a  $\sigma$ -finite measure. It is also easily verified that the set function

$$\nu(F) = \int_{\Omega_1} \left[ \int_{\Omega_2} I_F(\omega_1, \omega_2) \mu_2(d\omega_2) \right] \mu_1(d\omega_1)$$

is a  $\sigma$ -finite measure. It will be shown below that  $\lambda$  and  $\nu$  coincide on sets of the form  $F = A \times B$ , and therefore on the algebra  $\mathcal{F}$ . Hence it follows by Carathéodory's theorem that  $\lambda$  and  $\nu$  coincide for all  $F \in \mathcal{F}$ .

We turn now to the proof of the full conclusion of Fubini's theorem. By (47),

$$\int_{\Omega_1 \times \Omega_2} \xi^+(\omega_1, \omega_2) d(\mu_1 \times \mu_2) < \infty, \quad \int_{\Omega_1 \times \Omega_2} \xi^-(\omega_1, \omega_2) d(\mu_1 \times \mu_2) < \infty.$$

By what has already been proved, the integral  $\int_{\Omega_2} \xi^+(\omega_1, \omega_2) \mu_2(d\omega_2)$  is an  $\mathcal{F}_1$ -measurable function of  $\omega_1$  and

$$\int_{\Omega_1} \left[ \int_{\Omega_2} \xi^+(\omega_1, \omega_2) \mu_2(d\omega_2) \right] \mu_1(d\omega_1) = \int_{\Omega_1 \times \Omega_2} \xi^+(\omega_1, \omega_2) d(\mu_1 \times \mu_2) < \infty.$$

Consequently by Problem 4 (see also Property J in Subsection 2)

$$\int_{\Omega_2} \xi^+(\omega_1, \omega_2) \mu_2(d\omega_2) < \infty \quad (\mu_1\text{-a.s.}).$$

In the same way

$$\int_{\Omega_2} \xi^-(\omega_1, \omega_2) \mu_2(d\omega_2) < \infty \quad (\mu_1\text{-a.s.}),$$

and therefore

$$\int_{\Omega_2} |\xi(\omega_1, \omega_2)| \mu_2(d\omega_2) < \infty \quad (\mu_1\text{-a.s.}).$$

It is clear that, except on a set  $\mathcal{N}$  of  $\mu_1$ -measure zero,

$$\int_{\Omega_2} \xi(\omega_1, \omega_2) \mu_2(d\omega_2) = \int_{\Omega_2} \xi^+(\omega_1, \omega_2) \mu_2(d\omega_2) - \int_{\Omega_2} \xi^-(\omega_1, \omega_2) \mu_2(d\omega_2). \quad (54)$$

Taking the integrals to be zero for  $\omega_1 \in \mathcal{N}$ , we may suppose that (54) holds for all  $\omega \in \Omega_1$ . Then, integrating (54) with respect to  $\mu_1$  and using (50), we obtain

$$\begin{aligned} \int_{\Omega_1} \left[ \int_{\Omega_2} \xi(\omega_1, \omega_2) \mu_2(d\omega_2) \right] \mu_1(d\omega_1) &= \int_{\Omega_1} \left[ \int_{\Omega_2} \xi^+(\omega_1, \omega_2) \mu_2(d\omega_2) \right] \mu_1(d\omega_1) \\ &\quad - \int_{\Omega_1} \left[ \int_{\Omega_2} \xi^-(\omega_1, \omega_2) \mu_2(d\omega_2) \right] \mu_1(d\omega_1) \\ &= \int_{\Omega_1 \times \Omega_2} \xi^+(\omega_1, \omega_2) d(\mu_1 \times \mu_2) \\ &\quad - \int_{\Omega_1 \times \Omega_2} \xi^-(\omega_1, \omega_2) d(\mu_1 \times \mu_2) \\ &= \int_{\Omega_1 \times \Omega_2} \xi(\omega_1, \omega_2) d(\mu_1 \times \mu_2). \end{aligned}$$

Similarly we can establish the first equation in (48) and the equation

$$\int_{\Omega_1 \times \Omega_2} \xi(\omega_1, \omega_2) d(\mu_1 \times \mu_2) = \int_{\Omega_2} \left[ \int_{\Omega_1} \xi(\omega_1, \omega_2) \mu_1(d\omega_1) \right] \mu_2(d\omega_2).$$

This completes the proof of the theorem.

**Corollary.** *If  $\int_{\Omega_1} [\int_{\Omega_2} |\xi(\omega_1, \omega_2)| \mu_2(d\omega_2)] \mu_1(d\omega_1) < \infty$ , the conclusion of Fubini's theorem is still valid.*

In fact, under this hypothesis (47) follows from (50), and consequently the conclusions of Fubini's theorem hold.

**EXAMPLE.** Let  $(\xi, \eta)$  be a pair of random variables whose distribution has a two-dimensional density  $f_{\xi\eta}(x, y)$ , i.e.

$$P((\xi, \eta) \in B) = \int_B f_{\xi\eta}(x, y) dx dy, \quad B \in \mathcal{B}(R^2),$$

where  $f_{\xi\eta}(x, y)$  is a nonnegative  $\mathcal{B}(R^2)$ -measurable function, and the integral is a Lebesgue integral with respect to two-dimensional Lebesgue measure.

Let us show that the one-dimensional distributions for  $\xi$  and  $\eta$  have densities  $f_\xi(x)$  and  $f_\eta(y)$ , and furthermore

$$f_\xi(x) = \int_{-\infty}^{\infty} f_{\xi\eta}(x, y) dy$$

and

(55)

$$f_\eta(y) = \int_{-\infty}^{\infty} f_{\xi\eta}(x, y) dx.$$

In fact, if  $A \in \mathcal{B}(R)$ , then by Fubini's theorem

$$P(\xi \in A) = P((\xi, \eta) \in A \times R) = \int_{A \times R} f_{\xi\eta}(x, y) dx dy = \int_A \left[ \int_R f_{\xi\eta}(x, y) dy \right] dx.$$

This establishes both the existence of a density for the probability distribution of  $\xi$  and the first formula in (55). The second formula is established similarly.

According to the theorem in §5, a necessary and sufficient condition that  $\xi$  and  $\eta$  are independent is that

$$F_{\xi\eta}(x, y) = F_\xi(x)F_\eta(y), \quad (x, y) \in R^2.$$

Let us show that when there is a two-dimensional density  $f_{\xi\eta}(x, y)$ , the variables  $\xi$  and  $\eta$  are independent if and only if

$$f_{\xi\eta}(x, y) = f_\xi(x)f_\eta(y) \quad (56)$$

(where the equation is to be understood in the sense of holding almost surely with respect to two-dimensional Lebesgue measure).

In fact, in (56) holds, then by Fubini's theorem

$$\begin{aligned} F_{\xi\eta}(x, y) &= \int_{(-\infty, x] \times (-\infty, y]} f_{\xi\eta}(u, v) du dv = \int_{(-\infty, x] \times (-\infty, y]} f_{\xi}(u) f_{\eta}(v) du dv \\ &= \int_{(-\infty, x]} f_{\xi}(u) du \left( \int_{(-\infty, y]} f_{\eta}(v) dv \right) = F_{\xi}(x) F_{\eta}(y) \end{aligned}$$

and consequently  $\xi$  and  $\eta$  are independent.

Conversely, if they are independent and have a density  $f_{\xi\eta}(x, y)$ , then again by Fubini's theorem

$$\begin{aligned} \int_{(-\infty, x] \times (-\infty, y]} f_{\xi\eta}(u, v) du dv &= \left( \int_{(-\infty, x]} f_{\xi}(u) du \right) \left( \int_{(-\infty, y]} f_{\eta}(v) dv \right) \\ &= \int_{(-\infty, x] \times (-\infty, y]} f_{\xi}(u) f_{\eta}(v) du dv. \end{aligned}$$

It follows that

$$\int_B f_{\xi\eta}(x, y) dx dy = \int_B f_{\xi}(x) f_{\eta}(y) dx dy$$

for every  $B \in \mathcal{B}(R^2)$ , and it is easily deduced from Property I that (56) holds.

**10.** In this subsection we discuss the relation between the Lebesgue and Riemann integrals.

We first observe that the construction of the Lebesgue integral is independent of the measurable space  $(\Omega, \mathcal{F})$  on which the integrands are given. On the other hand, the Riemann integral is not defined on abstract spaces in general, and for  $\Omega = R^n$  it is defined sequentially: first for  $R^1$ , and then extended, with corresponding changes, to the case  $n > 1$ .

We emphasize that the constructions of the Riemann and Lebesgue integrals are based on different ideas. The first step in the construction of the Riemann integral is to group the points  $x \in R^1$  according to their distances along the  $x$  axis. On the other hand, in Lebesgue's construction (for  $\Omega = R^1$ ) the points  $x \in R^1$  are grouped according to a different principle: by the distances between the values of the integrand. It is a consequence of these different approaches that the Riemann approximating sums have limits only for "mildly" discontinuous functions, whereas the Lebesgue sums converge to limits for a much wider class of functions.

Let us recall the definition of the Riemann–Stieltjes integral. Let  $G = G(x)$  be a generalized distribution function on  $R$  (see subsection 2 of §3) and  $\mu$  its corresponding Lebesgue–Stieltjes measure, and let  $g = g(x)$  be a bounded function that vanishes outside  $[a, b]$ .



Consider a decomposition  $\mathcal{P} = \{x_0, \dots, x_n\}$ ,

$$a = x_0 < x_1 < \dots < x_n = b,$$

of  $[a, b]$ , and form the upper and lower sums

$$\bar{\Sigma}_{\mathcal{P}} = \sum_{i=1}^n \bar{g}_i [G(x_i) - G(x_{i-1})], \quad \underline{\Sigma}_{\mathcal{P}} = \sum_{i=1}^n \underline{g}_i [G(x_i) - G(x_{i-1})]$$

where

$$\bar{g}_i = \sup_{x_{i-1} < y \leq x_i} g(y), \quad \underline{g}_i = \inf_{x_{i-1} < y \leq x_i} g(y).$$

Define simple functions  $\bar{g}_{\mathcal{P}}(x)$  and  $\underline{g}_{\mathcal{P}}(x)$  by taking

$$\bar{g}_{\mathcal{P}}(x) = \bar{g}_i, \quad \underline{g}_{\mathcal{P}}(x) = \underline{g}_i,$$

on  $x_{i-1} < x \leq x_i$ , and define  $\bar{g}_{\mathcal{P}}(a) = \underline{g}_{\mathcal{P}}(a) = g(a)$ . It is clear that then

$$\bar{\Sigma}_{\mathcal{P}} = (\text{L-S}) \int_a^b \bar{g}_{\mathcal{P}}(x) G(dx)$$

and

$$\underline{\Sigma}_{\mathcal{P}} = (\text{L-S}) \int_a^b \underline{g}_{\mathcal{P}}(x) G(dx).$$

Now let  $\{\mathcal{P}_k\}$  be a sequence of decompositions such that  $\mathcal{P}_k \subseteq \mathcal{P}_{k+1}$ . Then

$$\bar{g}_{\mathcal{P}_1} \geq \bar{g}_{\mathcal{P}_2} \geq \dots \geq g \geq \dots \geq \underline{g}_{\mathcal{P}_2} \geq \underline{g}_{\mathcal{P}_1},$$

and if  $|g(x)| \leq C$  we have, by the dominated convergence theorem,

$$\begin{aligned} \lim_{k \rightarrow \infty} \bar{\Sigma}_{\mathcal{P}_k} &= (\text{L-S}) \int_a^b \bar{g}(x) G(dx), \\ \lim_{k \rightarrow \infty} \underline{\Sigma}_{\mathcal{P}_k} &= (\text{L-S}) \int_a^b \underline{g}(x) G(dx), \end{aligned} \tag{57}$$

where  $\bar{g}(x) = \lim_k \bar{g}_{\mathcal{P}_k}(x)$ ,  $\underline{g}(x) = \lim_k \underline{g}_{\mathcal{P}_k}(x)$ .

If the limits  $\lim_k \bar{\Sigma}_{\mathcal{P}_k}$  and  $\lim_k \underline{\Sigma}_{\mathcal{P}_k}$  are finite and equal, and their common value is independent of the sequence of decompositions  $\{\mathcal{P}_k\}$ , we say that  $g = g(x)$  is Riemann–Stieltjes integrable, and the common value of the limits is denoted by

$$(\text{R-S}) \int_a^b g(x) G(dx). \tag{58}$$

When  $G(x) = x$ , the integral is called a Riemann integral and denoted by

$$(\text{R}) \int_a^b g(x) dx.$$

Now let  $(L-S) \int_a^b g(x)G(dx)$  be the corresponding Lebesgue–Stieltjes integral (see Remark 2 in Subsection 2).

**Theorem 9.** *If  $g = g(x)$  is continuous on  $[a, b]$ , it is Riemann–Stieltjes integrable and*

$$(R-S) \int_a^b g(x)G(dx) = (L-S) \int_a^b g(x)G(dx). \quad (59)$$

PROOF. Since  $g(x)$  is continuous, we have  $\bar{g}(x) = g(x) = \underline{g}(x)$ . Hence by (57)  $\lim_{k \rightarrow \infty} \sum \mathcal{P}_k = \lim_{k \rightarrow \infty} \underline{\sum} \mathcal{P}_k$ . Consequently  $g = g(x)$  is Riemann–Stieltjes integral (again by (57)).

Let us consider in more detail the question of the correspondence between the Riemann and Lebesgue integrals for the case of Lebesgue measure on the line  $R$ .

**Theorem 10.** *Let  $g(x)$  be a bounded function on  $[a, b]$ .*

- (a) *The function  $g = g(x)$  is Riemann integrable on  $[a, b]$  if and only if it is continuous almost everywhere (with respect to Lebesgue measure  $\bar{\lambda}$  on  $\mathcal{B}([a, b])$ ).*  
 (b) *If  $g = g(x)$  is Riemann integrable, it is Lebesgue integrable and*

$$(R) \int_a^b g(x) dx = (L) \int_a^b g(x) \bar{\lambda}(dx). \quad (60)$$

PROOF. (a) Let  $g = g(x)$  be Riemann integrable. Then, by (57),

$$(L) \int_a^b \bar{g}(x) \bar{\lambda}(dx) = (L) \int_a^b \underline{g}(x) \bar{\lambda}(dx).$$

But  $\underline{g}(x) \leq g(x) \leq \bar{g}(x)$ , and hence by Property H

$$\bar{g}(x) = g(x) = \underline{g}(x) \quad (\bar{\lambda}\text{-a.s.}), \quad (61)$$

from which it is easy to see that  $g(x)$  is continuous almost everywhere (with respect to  $\bar{\lambda}$ ).

Conversely, let  $g = g(x)$  be continuous almost everywhere (with respect to  $\bar{\lambda}$ ). Then (61) is satisfied and consequently  $g(x)$  differs from the (Borel) measurable function  $\bar{g}(x)$  only on a set  $\mathcal{N}$  with  $\bar{\lambda}(\mathcal{N}) = 0$ . But then

$$\begin{aligned} \{x: g(x) \leq c\} &= \{x: g(x) \leq c\} \cap \bar{\mathcal{N}} + \{x: g(x) \leq c\} \cap \mathcal{N} \\ &= \{x: \bar{g}(x) \leq c\} \cap \bar{\mathcal{N}} + \{x: g(x) \leq c\} \cap \mathcal{N} \end{aligned}$$

It is clear that the set  $\{x: \bar{g}(x) \leq c\} \cap \bar{\mathcal{N}} \in \mathcal{B}([a, b])$ , and that

$$\{x: g(x) \leq c\} \cap \mathcal{N}$$

is a subset of  $\mathcal{N}$  having Lebesgue measure  $\bar{\lambda}$  equal to zero and therefore also belonging to  $\overline{\mathcal{B}}([a, b])$ . Therefore  $g(x)$  is  $\overline{\mathcal{B}}([a, b])$ -measurable and, as a bounded function, is Lebesgue integrable. Therefore by Property G,

$$(L) \int_a^b \bar{g}(x) \bar{\lambda}(dx) = (L) \int_a^b \underline{g}(x) \bar{\lambda}(dx) = (L) \int_a^b g(x) \bar{\lambda}(dx),$$

which completes the proof of (a).

(b) If  $g = g(x)$  is Riemann integrable, then according to (a) it is continuous ( $\bar{\lambda}$ -a.s.). It was shown above that then  $g(x)$  is Lebesgue integrable and its Riemann and Lebesgue integrals are equal.

This completes the proof of the theorem.

**Remark.** Let  $\mu$  be a Lebesgue–Stieltjes measure on  $\mathcal{B}([a, b])$ . Let  $\overline{\mathcal{B}}_\mu([a, b])$  be the system consisting of those subsets  $\Lambda \subseteq [a, b]$  for which there are sets  $A$  and  $B$  in  $\mathcal{B}([a, b])$  such that  $A \subseteq \Lambda \subseteq B$  and  $\mu(B \setminus A) = 0$ . Let  $\bar{\mu}$  be an extension of  $\mu$  to  $\overline{\mathcal{B}}_\mu([a, b])$  ( $\bar{\mu}(\Lambda) = \mu(A)$  for  $\Lambda$  such that  $A \subseteq \Lambda \subseteq B$  and  $\mu(B \setminus A) = 0$ ). Then the conclusion of the theorem remains valid if we consider  $\bar{\mu}$  instead of Lebesgue measure  $\bar{\lambda}$ , and the Riemann–Stieltjes and Lebesgue–Stieltjes measures with respect to  $\bar{\mu}$  instead of the Riemann and Lebesgue integrals.

**11.** In this part we present a useful theorem on integration by parts for the Lebesgue–Stieltjes integral.

Let two generalized distribution functions  $F = F(x)$  and  $G = G(x)$  be given on  $(R, \mathcal{B}(R))$ .

**Theorem 11.** *The following formulas are valid for all real  $a$  and  $b, a < b$ :*

$$F(b)G(b) - F(a)G(a) = \int_a^b F(s-)dG(s) + \int_a^b G(s)dF(s), \quad (62)$$

or equivalently

$$\begin{aligned} F(b)G(b) - F(a)G(a) &= \int_a^b F(s-)dG(s) + \int_a^b G(s-)dF(s) \\ &\quad + \sum_{a < s \leq b} \Delta F(s) \cdot \Delta G(s), \end{aligned} \quad (63)$$

where  $F(s-) = \lim_{t \uparrow s} F(t)$ ,  $\Delta F(s) = F(s) - F(s-)$ .

**Remark 1.** Formula (62) can be written symbolically in “differential” form

$$d(FG) = F_- dG + G dF. \quad (64)$$

**Remark 2.** The conclusion of the theorem remains valid for functions  $F$  and  $G$  of bounded variation on  $[a, b]$ . (Every such function that is continuous on the right and has limits on the left can be represented as the difference of two monotone nondecreasing functions.)

**PROOF.** We first recall that in accordance with Subsection 1 an integral  $\int_a^b (\cdot)$  means  $\int_{[a, b]} (\cdot)$ . Then (see formula (2) in §3)

$$(F(b) - F(a))(G(b) - G(a)) = \int_a^b dF(s) \cdot \int_a^b dG(t).$$

Let  $F \times G$  denote the direct product of the measures corresponding to  $F$  and  $G$ . Then by Fubini's theorem

$$\begin{aligned} (F(b) - F(a))(G(b) - G(a)) &= \int_{[a, b] \times [a, b]} d(F \times G)(s, t) \\ &= \int_{[a, b] \times [a, b]} I_{\{s \geq t\}}(s, t) d(F \times G)(s, t) + \int_{[a, b] \times [a, b]} I_{\{s < t\}}(s, t) d(F \times G)(s, t) \\ &= \int_{[a, b]} (G(s) - G(a)) dF(s) + \int_{[a, b]} (F(t-) - F(a)) dG(t) \\ &= \int_a^b G(s) dF(s) + \int_a^b F(s-) dG(s) - G(a)(F(b) - F(a)) - F(a)(G(b) - G(a)), \end{aligned} \quad (65)$$

where  $I_A$  is the indicator of the set  $A$ .

Formula (62) follows immediately from (65). In turn, (63) follows from (62) if we observe that

$$\int_a^b (G(s) - G(s-)) dF(s) = \sum_{a < s \leq b} \Delta G(s) \cdot \Delta F(s). \quad (66)$$

**Corollary 1.** If  $F(x)$  and  $G(x)$  are distribution functions, then

$$F(x)G(x) = \int_{-\infty}^x F(s-) dG(s) + \int_{-\infty}^x G(s) dF(s). \quad (67)$$

If also

$$F(x) = \int_{-\infty}^x f(s) ds,$$

then

$$F(x)G(x) = \int_{-\infty}^x F(s) dG(s) + \int_{-\infty}^x G(s) f(s) ds. \quad (68)$$

**Corollary 2.** Let  $\xi$  be a random variable with distribution function  $F(x)$  and  $E|\xi|^n < \infty$ . Then

$$\int_0^{\infty} x^n dF(x) = n \int_0^{\infty} x^{n-1} [1 - F(x)] dx, \quad (69)$$

$$\int_{-\infty}^0 |x|^n dF(x) = - \int_0^{\infty} x^n dF(-x) = n \int_0^{\infty} x^{n-1} F(-x) dx \quad (70)$$

and

$$E|\xi|^n = \int_{-\infty}^{\infty} |x|^n dF(x) = n \int_0^{\infty} x^{n-1} [1 - F(x) + F(-x)] dx. \quad (71)$$

To prove (69) we observe that

$$\begin{aligned} \int_0^b x^n dF(x) &= - \int_0^b x^n d(1 - F(x)) \\ &= -b^n(1 - F(b)) + n \int_0^b x^{n-1}(1 - F(x)) dx. \end{aligned} \quad (72)$$

Let us show that since  $E|\xi|^n < \infty$ ,

$$b^n(1 - F(b) + F(-b)) \leq b^n P(|\xi| \geq b) \rightarrow 0. \quad (73)$$

In fact,

$$E|\xi|^n = \sum_{k=1}^{\infty} \int_{k-1}^k |x|^n dF(x) < \infty$$

and therefore

$$\sum_{k \geq b+1} \int_{k-1}^k |x|^n dF(x) \rightarrow 0, \quad n \rightarrow \infty.$$

But

$$\sum_{k \geq b+1} \int_{k-1}^k |x|^n dF(x) \geq b^n P(|\xi| \geq b),$$

which establishes (73).

Taking the limit as  $b \rightarrow \infty$  in (72), we obtain (69).

Formula (70) is proved similarly, and (71) follows from (69) and (70).

**12.** Let  $A(t)$ ,  $t \geq 0$ , be a function of locally bounded variation (i.e., of bounded variation on each finite interval  $[a, b]$ ), which is continuous on the right and has limits on the left. Consider the equation

$$Z_t = 1 + \int_0^t Z_{s-} dA(s), \quad (74)$$

which can be written in differential form as

$$dZ = Z_- dA, \quad Z_0 = 1. \quad (75)$$

The formula that we have proved for integration by parts lets us solve (74) explicitly in the class of functions of bounded variation.

We introduce the function

$$\mathcal{E}_t(A) = e^{A(t) - A(0)} \prod_{0 \leq s \leq t} (1 + \Delta A(s)) e^{-\Delta A(s)}, \quad (76)$$

where  $\Delta A(s) = A(s) - A(s-)$  for  $s > 0$ , and  $\Delta A(0) = 0$ .

The function  $A(s)$ ,  $0 \leq s \leq t$ , has bounded variation and therefore has at most a countable number of discontinuities, and so the series  $\sum_{0 \leq s \leq t} |\Delta A(s)|$  converges. It follows that

$$\prod_{0 \leq s \leq t} (1 + \Delta A(s)) e^{-\Delta A(s)}$$

is a function of locally bounded variation.

If  $A^c(t) = A(t) - \sum_{0 \leq s \leq t} \Delta A(s)$  is the continuous component of  $A(t)$ , we can rewrite (76) in the form

$$\mathcal{E}_t(A) = e^{A^c(t) - A^c(0)} \prod_{0 \leq s \leq t} (1 + \Delta A(s)). \quad (77)$$

Let us write

$$F(t) = e^{A^c(t) - A^c(0)}, \quad G(t) = \prod_{0 \leq s \leq t} (1 + \Delta A(s)).$$

Then by (62)

$$\begin{aligned} \mathcal{E}_t(A) &= F(t)G(t) = 1 + \int_0^t F(s) dG(s) + \int_0^t G(s-) dF(s) \\ &= 1 + \sum_{0 \leq s \leq t} F(s)G(s-) \Delta A(s) + \int_0^t G(s-) F(s) dA^c(s) \\ &= 1 + \int_0^t \mathcal{E}_{s-}(A) dA(s). \end{aligned}$$

Therefore  $\mathcal{E}_t(A)$ ,  $t \geq 0$ , is a (locally bounded) solution of (74). Let us show that this is the only locally bounded solution.

Suppose that there are two such solutions and let  $Y = Y(t)$ ,  $t \geq 0$ , be their difference. Then

$$Y(t) = \int_0^t Y(s-) dA(s).$$

Put

$$T = \inf\{t \geq 0: Y(t) \neq 0\},$$

where we take  $T = \infty$  if  $Y(t) = 0$  for  $t \geq 0$ .

Since  $A(t)$  is a function of locally bounded variation, there are two generalized distribution functions  $A_1(t)$  and  $A_2(t)$  such that  $A(t) = A_1(t) - A_2(t)$ . If we suppose that  $T < \infty$ , we can find a finite  $T'$  such that

$$[A_1(T') + A_2(T')] - [A_1(T) + A_2(T)] \leq \frac{1}{2}.$$

Then it follows from the equation

$$Y(t) = \int_T^t Y(s-) dA(s), \quad t \geq T,$$

that

$$\sup_{t \leq T'} |Y(t)| \leq \frac{1}{2} \sup_{t \leq T'} |Y(t)|$$

and since  $\sup |Y(t)| < \infty$ , we have  $Y(t) = 0$  for  $T < t \leq T'$ , contradicting the assumption that  $T < \infty$ .

Thus we have proved the following theorem.

**Theorem 12.** *There is a unique locally bounded solution of (74), and it is given by (76).*

### 13. PROBLEMS

1. Establish the representation (6).
2. Prove the following extension of Property E. Let  $\xi$  and  $\eta$  be random variables for which  $E\xi$  and  $E\eta$  are defined and the sum  $E\xi + E\eta$  is meaningful (does not have the form  $\infty - \infty$  or  $-\infty + \infty$ ). Then

$$E(\xi + \eta) = E\xi + E\eta.$$

3. Generalize Property G by showing that if  $\xi = \eta$  (a.s.) and  $E\xi$  exists, then  $E\eta$  exists and  $E\xi = E\eta$ .
4. Let  $\xi$  be an extended random variable,  $\mu$  a  $\sigma$ -finite measure, and  $\int_{\Omega} |\xi| d\mu < \infty$ . Show that  $|\xi| < \infty$  ( $\mu$ -a.s.) (cf. Property J).
5. Let  $\mu$  be a  $\sigma$ -finite measure,  $\xi$  and  $\eta$  extended random variables for which  $E\xi$  and  $E\eta$  are defined. If  $\int_A \xi d\mu \leq \int_A \eta d\mu$  for all  $A \in \mathcal{F}$ , then  $\xi \leq \eta$  ( $\mu$ -a.s.). (Cf. Property I.)
6. Let  $\xi$  and  $\eta$  be independent nonnegative random variables. Show that  $E\xi\eta = E\xi \cdot E\eta$ .
7. Using Fatou's lemma, show that

$$P(\lim A_n) \leq \lim P(A_n), \quad P(\overline{\lim A_n}) \geq \overline{\lim P(A_n)}.$$

8. Find an example to show that in general it is impossible to weaken the hypothesis " $|\xi_n| \leq \eta$ ,  $E\eta < \infty$ " in the dominated convergence theorem.

9. Find an example to show that in general the hypothesis " $\xi_n \leq \eta$ ,  $E\eta > -\infty$ " in Fatou's lemma cannot be omitted.
10. Prove the following variants of Fatou's lemma. Let the family  $\{\xi_n^+\}_{n \geq 1}$  of random variables be uniformly integrable and let  $E \overline{\lim} \xi_n$  exist. Then

$$\overline{\lim} E\xi_n \leq E \overline{\lim} \xi_n.$$

Let  $\xi_n \leq \eta_n$ ,  $n \geq 1$ , where the family  $\{\xi_n^+\}_{n \geq 1}$  is uniformly integrable and  $\eta_n$  converges a.s. (or only in probability—see §10 below) to a random variable  $\eta$ . Then  $\overline{\lim} E\xi_n \leq E \overline{\lim} \xi_n$ .

11. Dirichlet's function

$$d(x) = \begin{cases} 1, & x \text{ irrational,} \\ 0, & x \text{ rational,} \end{cases}$$

is defined on  $[0, 1]$ , Lebesgue integrable, but not Riemann integrable. Why?

12. Find an example of a sequence of Riemann integrable functions  $\{f_n\}_{n \geq 1}$ , defined on  $[0, 1]$ , such that  $|f_n| \leq 1$ ,  $f_n \rightarrow f$  almost everywhere (with Lebesgue measure), but  $f$  is not Riemann integrable.
13. Let  $(a_{i,j}; i, j \geq 1)$  be a sequence of real numbers such that  $\sum_{i,j} |a_{i,j}| < \infty$ . Deduce from Fubini's theorem that

$$\sum_{(i,j)} a_{ij} = \sum_i \left( \sum_j a_{ij} \right) = \sum_j \left( \sum_i a_{ij} \right). \quad (78)$$

14. Find an example of a sequence  $(a_{i,j}; i, j \geq 1)$  for which  $\sum_{i,j} |a_{i,j}| = \infty$  and the equation in (78) does not hold.
15. Starting from simple functions and using the theorem on taking limits under the Lebesgue integral sign, prove the following result on *integration by substitution*.

Let  $h = h(y)$  be a nondecreasing continuously differentiable function on  $[a, b]$ , and let  $f(x)$  be (Lebesgue) integrable on  $[h(a), h(b)]$ . Then the function  $f(h(y))h'(y)$  is integrable on  $[a, b]$  and

$$\int_{h(a)}^{h(b)} f(x) dx = \int_a^b f(h(y))h'(y) dy.$$

16. Prove formula (70).
17. Let  $\xi, \xi_1, \xi_2, \dots$  be nonnegative integrable random variables such that  $E\xi_n \rightarrow E\xi$  and  $P(\xi - \xi_n > \varepsilon) \rightarrow 0$  for every  $\varepsilon > 0$ . Show that then  $E|\xi_n - \xi| \rightarrow 0$ ,  $n \rightarrow \infty$ .
18. Let  $\xi, \eta, \zeta$  and  $\xi_n, \eta_n, \zeta_n$ ,  $n \geq 1$ , be random variables such that

$$\begin{aligned} \xi_n &\xrightarrow{P} \xi, & \eta_n &\xrightarrow{P} \eta, & \zeta_n &\xrightarrow{P} \zeta, & \eta_n \leq \xi_n \leq \zeta_n, & n \geq 1, \\ E\xi_n &\rightarrow E\xi, & E\eta_n &\rightarrow E\eta, \end{aligned}$$

and the expectations  $E\xi, E\eta, E\zeta$  are finite. Show that then  $E\xi_n \rightarrow E\xi$  (Pratt's lemma).

If also  $\eta_n \leq 0 \leq \zeta_n$  then  $E|\xi_n - \xi| \rightarrow 0$ ,

Deduce that if  $\xi_n \xrightarrow{P} \xi$ ,  $E|\xi_n| \rightarrow E|\xi|$  and  $E|\xi| < \infty$ , then  $E|\xi_n - \xi| \rightarrow 0$ .



## §7. Conditional Probabilities and Conditional Expectations with Respect to a $\sigma$ -Algebra

1. Let  $(\Omega, \mathcal{F}, P)$  be a probability space, and let  $A \in \mathcal{F}$  be an event such that  $P(A) > 0$ . As for finite probability spaces, the *conditional probability of  $B$  with respect to  $A$*  (denoted by  $P(B|A)$ ) means  $P(BA)/P(A)$ , and the *conditional probability of  $B$  with respect to the finite or countable decomposition  $\mathcal{D} = \{D_1, D_2, \dots\}$  with  $P(D_i) > 0, i \geq 1$*  (denoted by  $P(B|\mathcal{D})$ ) is the random variable equal to  $P(B|D_i)$  for  $\omega \in D_i, i \geq 1$ :

$$P(B|\mathcal{D}) = \sum_{i \geq 1} P(B|D_i)I_{D_i}(\omega).$$

In a similar way, if  $\xi$  is a random variable for which  $E\xi$  is defined, the *conditional expectation of  $\xi$  with respect to the event  $A$  with  $P(A) > 0$*  (denoted by  $E(\xi|A)$ ) is  $E(\xi I_A)/P(A)$  (cf. (I.8.10)).

The random variable  $P(B|\mathcal{D})$  is evidently measurable with respect to the  $\sigma$ -algebra  $\mathcal{G} = \sigma(\mathcal{D})$ , and is consequently also denoted by  $P(B|\mathcal{G})$  (see §8 of Chapter I).

However, in probability theory we may have to consider conditional probabilities with respect to events whose probabilities are zero.

Consider, for example, the following experiment. Let  $\xi$  be a random variable that is uniformly distributed on  $[0, 1]$ . If  $\xi = x$ , toss a coin for which the probability of head is  $x$ , and the probability of tail is  $1 - x$ . Let  $v$  be the number of heads in  $n$  independent tosses of this coin. What is the "conditional probability  $P(v = k|\xi = x)$ "? Since  $P(\xi = x) = 0$ , the conditional probability  $P(v = k|\xi = x)$  is undefined, although it is intuitively plausible that "it ought to be  $C_n^k x^k (1 - x)^{n-k}$ ."

Let us now give a general definition of conditional expectation (and, in particular, of conditional probability) with respect to a  $\sigma$ -algebra  $\mathcal{G}, \mathcal{G} \subseteq \mathcal{F}$ , and compare it with the definition given in §8 of Chapter I for finite probability spaces.

2. Let  $(\Omega, \mathcal{F}, P)$  be a probability space,  $\mathcal{G}$  a  $\sigma$ -algebra,  $\mathcal{G} \subseteq \mathcal{F}$  ( $\mathcal{G}$  is a  $\sigma$ -subalgebra of  $\mathcal{F}$ ), and  $\xi = \xi(\omega)$  a random variable. Recall that, according to §6, the expectation  $E\xi$  was defined in two stages: first for a nonnegative random variable  $\xi$ , then in the general case by

$$E\xi = E\xi^+ - E\xi^-,$$

and only under the assumption that

$$\min(E\xi^-, E\xi^+) < \infty.$$

A similar two-stage construction is also used to define conditional expectations  $E(\xi|\mathcal{G})$ .

**Definition 1.**

(1) The *conditional expectation of a nonnegative random variable  $\xi$  with respect to the  $\sigma$ -algebra  $\mathcal{G}$*  is a nonnegative extended random variable, denoted by  $E(\xi|\mathcal{G})$  or  $E(\xi|\mathcal{G})(\omega)$ , such that

- (a)  $E(\xi|\mathcal{G})$  is  $\mathcal{G}$ -measurable;  
 (b) for every  $A \in \mathcal{G}$

$$\int_A \xi dP = \int_A E(\xi|\mathcal{G}) dP. \quad (1)$$

(2) The *conditional expectation  $E(\xi|\mathcal{G})$ , or  $E(\xi|\mathcal{G})(\omega)$ , of any random variable  $\xi$  with respect to the  $\sigma$ -algebra  $\mathcal{G}$* , is considered to be defined if

$$\min(E(\xi^+|\mathcal{G}), E(\xi^-|\mathcal{G})) < \infty,$$

$P$ -a.s., and it is given by the formula

$$E(\xi|\mathcal{G}) \equiv E(\xi^+|\mathcal{G}) - E(\xi^-|\mathcal{G}),$$

where, on the set (of probability zero) of sample points for which  $E(\xi^+|\mathcal{G}) = E(\xi^-|\mathcal{G}) = \infty$ , the difference  $E(\xi^+|\mathcal{G}) - E(\xi^-|\mathcal{G})$  is given an arbitrary value, for example zero.

We begin by showing that, for nonnegative random variables,  $E(\xi|\mathcal{G})$  actually exists. By (6.36) the set function

$$Q(A) = \int_A \xi dP, \quad A \in \mathcal{G}, \quad (2)$$

is a measure on  $(\Omega, \mathcal{G})$ , and is absolutely continuous with respect to  $P$  (considered on  $(\Omega, \mathcal{G})$ ,  $\mathcal{G} \subseteq \mathcal{F}$ ). Therefore (by the Radon–Nikodým theorem) there is a nonnegative  $\mathcal{G}$ -measurable extended random variable  $E(\xi|\mathcal{G})$  such that

$$Q(A) = \int_A E(\xi|\mathcal{G}) dP. \quad (3)$$

Then (1) follows from (2) and (3).

**Remark 1.** In accordance with the Radon–Nikodým theorem, the conditional expectation  $E(\xi|\mathcal{G})$  is defined only up to sets of  $P$ -measure zero. In other words,  $E(\xi|\mathcal{G})$  can be taken to be any  $\mathcal{G}$ -measurable function  $f(\omega)$  for which  $Q(A) = \int_A f(\omega) dP$ ,  $A \in \mathcal{G}$  (a “variant” of the conditional expectation).

Let us observe that, in accordance with the remark on the Radon–Nikodým theorem,

$$E(\xi|\mathcal{G}) \equiv \frac{dQ}{dP}(\omega), \quad (4)$$

i.e. the conditional expectation is just the derivative of the Radon–Nikodým measure  $Q$  with respect to  $P$  (considered on  $(\Omega, \mathcal{G})$ ).

**Remark 2.** In connection with (1), we observe that we cannot in general put  $E(\xi|\mathcal{G}) = \xi$ , since  $\xi$  is not necessarily  $\mathcal{G}$ -measurable.

**Remark 3.** Suppose that  $\xi$  is a random variable for which  $E\xi$  does not exist. Then  $E(\xi|\mathcal{G})$  may be definable as a  $\mathcal{G}$ -measurable function for which (1) holds. This is usually just what happens. Our definition  $E(\xi|\mathcal{G}) \equiv E(\xi^+|\mathcal{G}) - E(\xi^-|\mathcal{G})$  has the advantage that for the trivial  $\sigma$ -algebra  $\mathcal{G} = \{\emptyset, \Omega\}$  it reduces to the definition of  $E\xi$  but does not presuppose the existence of  $E\xi$ . (For example, if  $\xi$  is a random variable with  $E\xi^+ = \infty$ ,  $E\xi^- = \infty$ , and  $\mathcal{G} = \mathcal{F}$ , then  $E\xi$  is not defined, but in terms of Definition 1,  $E(\xi|\mathcal{G})$  exists and is simply  $\xi = \xi^+ - \xi^-$ .

**Remark 4.** Let the random variable  $\xi$  have a conditional expectation  $E(\xi|\mathcal{G})$  with respect to the  $\sigma$ -algebra  $\mathcal{G}$ . The *conditional variance* (denoted by  $V(\xi|\mathcal{G})$  or  $V(\xi|\mathcal{G})(\omega)$ ) of  $\xi$  is the random variable

$$V(\xi|\mathcal{G}) \equiv E[(\xi - E(\xi|\mathcal{G}))^2|\mathcal{G}].$$

(Cf. the definition of the conditional variance  $V(\xi|\mathcal{D})$  of  $\xi$  with respect to a decomposition  $\mathcal{D}$ , as given in Problem 2, §8, Chapter I.)

**Definition 2.** Let  $B \in \mathcal{F}$ . The conditional expectation  $E(I_B|\mathcal{G})$  is denoted by  $P(B|\mathcal{G})$ , or  $P(B|\mathcal{G})(\omega)$ , and is called the *conditional probability of the event  $B$  with respect to the  $\sigma$ -algebra  $\mathcal{G}$* ,  $\mathcal{G} \subseteq \mathcal{F}$ .

It follows from Definitions 1 and 2 that, for a given  $B \in \mathcal{F}$ ,  $P(B|\mathcal{G})$  is a random variable such that

(a)  $P(B|\mathcal{G})$  is  $\mathcal{G}$ -measurable,

$$(b) \quad P(A \cap B) = \int_A P(B|\mathcal{G}) dP \quad (5)$$

for every  $A \in \mathcal{G}$ .

**Definition 3.** Let  $\xi$  be a random variable and  $\mathcal{G}_\eta$  the  $\sigma$ -algebra generated by a random element  $\eta$ . Then  $E(\xi|\mathcal{G}_\eta)$ , if defined, means  $E(\xi|\eta)$  or  $E(\xi|\eta)(\omega)$ , and is called the *conditional expectation of  $\xi$  with respect to  $\eta$* .

The conditional probability  $P(B|\mathcal{G}_\eta)$  is denoted by  $P(B|\eta)$  or  $P(B|\eta)(\omega)$ , and is called the *conditional probability of  $B$  with respect to  $\eta$* .

3. Let us show that the definition of  $E(\xi|\mathcal{G})$  given here agrees with the definition of conditional expectation in §8 of Chapter I.

Let  $\mathcal{D} = \{D_1, D_2, \dots\}$  be a finite or countable decomposition with atoms  $D_i$  with respect to the probability  $P$  (i.e.  $P(D_i) > 0$ , and if  $A \subseteq D_i$ , then either  $P(A) = 0$  or  $P(D_i \setminus A) = 0$ ).

**Theorem 1.** If  $\mathcal{G} = \sigma(\mathcal{D})$  and  $\xi$  is a random variable for which  $E\xi$  is defined, then

$$E(\xi|\mathcal{G}) = E(\xi|D_i) \quad (P\text{-a.s. on } D_i) \quad (6)$$

or equivalently

$$E(\xi|\mathcal{G}) = \frac{E(\xi I_{D_i})}{P(D_i)} \quad (P\text{-a.s. on } D_i).$$

(The notation " $\xi = \eta$  (P-a.s. on  $A$ )" or

" $\xi = \eta(A; P\text{-a.s.})$ " means that  $P(A \cap \{\xi \neq \eta\}) = 0$ .)

**PROOF.** According to Lemma 3 of §4,  $E(\xi|\mathcal{G}) = K_i$  on  $D_i$ , where  $K_i$  are constants. But

$$\int_{D_i} \xi dP = \int_{D_i} E(\xi|\mathcal{G}) dP = K_i P(D_i),$$

whence

$$K_i = \frac{1}{P(D_i)} \int_{D_i} \xi dP = \frac{E(\xi I_{D_i})}{P(D_i)} = E(\xi|D_i).$$

This completes the proof of the theorem.

Consequently the concept of the conditional expectation  $E(\xi|\mathcal{D})$  with respect to a finite decomposition  $\mathcal{D} = \{D_1, \dots, D_n\}$ , as introduced in Chapter I, is a special case of the concept of conditional expectation with respect to the  $\sigma$ -algebra  $\mathcal{G} = \sigma(\mathcal{D})$ .

**4. Properties of conditional expectations.** (We shall suppose that the expectations are defined for all the random variables that we consider and that  $\mathcal{G} \subseteq \mathcal{F}$ .)

**A\*.** If  $C$  is a constant and  $\xi = C$  (a.s.), then  $E(\xi|\mathcal{G}) = C$  (a.s.).

**B\*.** If  $\xi \leq \eta$  (a.s.) then  $E(\xi|\mathcal{G}) \leq E(\eta|\mathcal{G})$  (a.s.).

**C\*.**  $|E(\xi|\mathcal{G})| \leq E(|\xi||\mathcal{G})$  (a.s.).

**D\*.** If  $a, b$  are constants and  $aE\xi + bE\eta$  is defined, then

$$E(a\xi + b\eta|\mathcal{G}) = aE(\xi|\mathcal{G}) + bE(\eta|\mathcal{G}) \quad (a.s.).$$

**E\*.** Let  $\mathcal{F}_* = \{\varphi, \Omega\}$  be the trivial  $\sigma$ -algebra. Then

$$E(\xi|\mathcal{F}_*) = E\xi \quad (a.s.).$$

**F\*.**  $E(\xi|\mathcal{F}) = \xi$  (a.s.).

**G\*.**  $E(E(\xi|\mathcal{G})) = E\xi$ .

**H\*.** If  $\mathcal{G}_1 \subseteq \mathcal{G}_2$  then

$$E[E(\xi|\mathcal{G}_2)|\mathcal{G}_1] = E(\xi|\mathcal{G}_1) \quad (\text{a.s.}).$$

**I\*.** If  $\mathcal{G}_1 \supseteq \mathcal{G}_2$  then

$$E[E(\xi|\mathcal{G}_2)|\mathcal{G}_1] = E(\xi|\mathcal{G}_2) \quad (\text{a.s.}).$$

**J\*.** Let a random variable  $\xi$  for which  $E\xi$  is defined be independent of the  $\sigma$ -algebra  $\mathcal{G}$  (i.e., independent of  $I_B, B \in \mathcal{G}$ ). Then

$$E(\xi|\mathcal{G}) = E\xi \quad (\text{a.s.}).$$

**K\*.** Let  $\eta$  be a  $\mathcal{G}$ -measurable random variable,  $E|\xi| < \infty$  and  $E|\xi\eta| < \infty$ . Then

$$E(\xi\eta|\mathcal{G}) = \eta E(\xi|\mathcal{G}) \quad (\text{a.s.}).$$

Let us establish these properties.

**A\*.** A constant function is measurable with respect to  $\mathcal{G}$ . Therefore we need only verify that

$$\int_A \xi dP = \int_A C dP, \quad A \in \mathcal{G}.$$

But, by the hypothesis  $\xi = C$  (a.s.) and Property G of §6, this equation is obviously satisfied.

**B\*.** If  $\xi \leq \eta$  (a.s.), then by Property B of §6

$$\int_A \xi dP \leq \int_A \eta dP, \quad A \in \mathcal{G},$$

and therefore

$$\int_A E(\xi|\mathcal{G}) dP \leq \int_A E(\eta|\mathcal{G}) dP, \quad A \in \mathcal{G}.$$

The required inequality now follows from Property I (§6).

**C\*.** This follows from the preceding property if we observe that  $-|\xi| \leq \xi \leq |\xi|$ .

**D\*.** If  $A \in \mathcal{G}$  then by Problem 2 of §6,

$$\begin{aligned} \int_A (a\xi + b\eta) dP &= \int_A a\xi dP + \int_A b\eta dP = \int_A aE(\xi|\mathcal{G}) dP \\ &+ \int_A bE(\eta|\mathcal{G}) dP = \int_A [aE(\xi|\mathcal{G}) + bE(\eta|\mathcal{G})] dP, \end{aligned}$$

which establishes **D\***.

**E\***. This property follows from the remark that  $E\xi$  is an  $\mathcal{F}_*$ -measurable function and the evident fact that if  $A = \Omega$  or  $A = \emptyset$  then

$$\int_A \xi \, dP = \int_A E\xi \, dP.$$

**F\***. Since  $\xi$  is  $\mathcal{F}$ -measurable and

$$\int_A \xi \, dP = \int_A \xi \, dP, \quad A \in \mathcal{F},$$

we have  $E(\xi|\mathcal{F}) = \xi$  (a.s.).

**G\***. This follows from **E\*** and **H\*** by taking  $\mathcal{G}_1 = \{\emptyset, \Omega\}$  and  $\mathcal{G}_2 = \mathcal{G}$ .

**H\***. Let  $A \in \mathcal{G}_1$ ; then

$$\int_A E(\xi|\mathcal{G}_1) \, dP = \int_A \xi \, dP.$$

Since  $\mathcal{G}_1 \subseteq \mathcal{G}_2$ , we have  $A \in \mathcal{G}_2$  and therefore

$$\int_A E[E(\xi|\mathcal{G}_2)|\mathcal{G}_1] \, dP = \int_A E(\xi|\mathcal{G}_2) \, dP = \int_A \xi \, dP.$$

Consequently, when  $A \in \mathcal{G}_1$ ,

$$\int_A E(\xi|\mathcal{G}_1) \, dP = \int_A E[E(\xi|\mathcal{G}_2)|\mathcal{G}_1] \, dP$$

and by Property I (§6) and Problem 5 (§6)

$$E(\xi|\mathcal{G}_1) = E[E(\xi|\mathcal{G}_2)|\mathcal{G}_1] \quad (\text{a.s.}).$$

**I\***. If  $A \in \mathcal{G}_1$ , then by the definition of  $E[E(\xi|\mathcal{G}_2)|\mathcal{G}_1]$

$$\int_A E[E(\xi|\mathcal{G}_2)|\mathcal{G}_1] \, dP = \int_A E(\xi|\mathcal{G}_2) \, dP.$$

The function  $E(\xi|\mathcal{G}_2)$  is  $\mathcal{G}_2$ -measurable and, since  $\mathcal{G}_2 \subseteq \mathcal{G}_1$ , also  $\mathcal{G}_1$ -measurable. It follows that  $E(\xi|\mathcal{G}_2)$  is a variant of the expectation  $E[E(\xi|\mathcal{G}_2)|\mathcal{G}_1]$ , which proves Property **I\***.

**J\***. Since  $E\xi$  is a  $\mathcal{G}$ -measurable function, we have only to verify that

$$\int_B dP = \int_B E\xi \, dP,$$

i.e. that  $E[\xi \cdot I_B] = E\xi \cdot E I_B$ . If  $E|\xi| < \infty$ , this follows immediately from Theorem 6 of §6. The general case can be reduced to this by applying Problem 6 of §6.

The proof of Property **K\*** will be given a little later; it depends on conclusion (a) of the following theorem.

**Theorem 2** (On Taking Limits Under the Expectation Sign). *Let  $\{\xi_n\}_{n \geq 1}$  be a sequence of extended random variables.*

(a) *If  $|\xi_n| \leq \eta$ ,  $E\eta < \infty$  and  $\xi_n \rightarrow \xi$  (a.s.), then*

$$E(\xi_n|\mathcal{G}) \rightarrow E(\xi|\mathcal{G}) \quad (\text{a.s.})$$

*and*

$$E(|\xi_n - \xi||\mathcal{G}) \rightarrow 0 \quad (\text{a.s.}).$$

(b) *If  $\xi_n \geq \eta$ ,  $E\eta > -\infty$  and  $\xi_n \uparrow \xi$  (a.s.), then*

$$E(\xi_n|\mathcal{G}) \uparrow E(\xi|\mathcal{G}) \quad (\text{a.s.}).$$

(c) *If  $\xi_n \leq \eta$ ,  $E\eta < \infty$ , and  $\xi_n \downarrow \xi$  (a.s.), then*

$$E(\xi_n|\mathcal{G}) \downarrow E(\xi|\mathcal{G}) \quad (\text{a.s.}).$$

(d) *If  $\xi_n \geq \eta$ ,  $E\eta > -\infty$ , then*

$$E(\lim \xi_n|\mathcal{G}) \leq \lim E(\xi_n|\mathcal{G}) \quad (\text{a.s.}).$$

(e) *If  $\xi_n \leq \eta$ ,  $E\eta < \infty$ , then*

$$\overline{\lim} E(\xi_n|\mathcal{G}) \leq E(\overline{\lim} \xi_n|\mathcal{G}) \quad (\text{a.s.}).$$

(f) *If  $\xi_n \geq 0$  then*

$$E(\sum \xi_n|\mathcal{G}) = \sum E(\xi_n|\mathcal{G}) \quad (\text{a.s.}).$$

**PROOF.** (a) Let  $\zeta_n = \sup_{m \geq n} |\xi_m - \xi|$ . Since  $\xi_n \rightarrow \xi$  (a.s.), we have  $\zeta_n \downarrow 0$  (a.s.). The expectations  $E\zeta_n$  and  $E\xi$  are finite; therefore by Properties **D\*** and **C\*** (a.s.)

$$|E(\xi_n|\mathcal{G}) - E(\xi|\mathcal{G})| = |E(\xi_n - \xi|\mathcal{G})| \leq E(|\xi_n - \xi||\mathcal{G}) \leq E(\zeta_n|\mathcal{G}).$$

Since  $E(\zeta_{n+1}|\mathcal{G}) \leq E(\zeta_n|\mathcal{G})$  (a.s.), the limit  $h = \lim_n E(\zeta_n|\mathcal{G})$  exists (a.s.). Then

$$0 \leq \int_{\Omega} h \, dP \leq \int_{\Omega} E(\zeta_n|\mathcal{G}) \, dP = \int_{\Omega} \zeta_n \, dP \rightarrow 0, \quad n \rightarrow \infty,$$

where the last statement follows from the dominated convergence theorem, since  $0 \leq \zeta_n \leq 2\eta$ ,  $E\eta < \infty$ . Consequently  $\int_{\Omega} h \, dP = 0$  and then  $h = 0$  (a.s.) by Property **H**.

(b) First let  $\eta \equiv 0$ . Since  $E(\xi_n|\mathcal{G}) \leq E(\xi_{n+1}|\mathcal{G})$  (a.s.) the limit  $\zeta(\omega) = \lim_n E(\xi_n|\mathcal{G})$  exists (a.s.). Then by the equation

$$\int_A \xi_n \, dP = \int_A E(\xi_n|\mathcal{G}) \, dP, \quad A \in \mathcal{G},$$

and the theorem on monotone convergence,

$$\int_A \xi \, dP = \int_A \zeta \, dP, \quad A \in \mathcal{G}.$$

Consequently  $\xi = \zeta$  (a.s.) by Property **I** and Problem 5 of §6.

For the proof in the general case, we observe that  $0 \leq \xi_n^+ \uparrow \xi^+$ , and by what has been proved,

$$E(\xi_n^+ | \mathcal{G}) \uparrow E(\xi^+ | \mathcal{G}) \quad (\text{a.s.}) \quad (7)$$

But  $0 \leq \xi_n^- \leq \xi^-$ ,  $E\xi^- < \infty$ , and therefore by (a)

$$E(\xi_n^- | \mathcal{G}) \rightarrow E(\xi^- | \mathcal{G}),$$

which, with (7), proves (b).

Conclusion (c) follows from (b).

(d) Let  $\zeta_n = \inf_{m \geq n} \zeta_m$ ; then  $\zeta_n \uparrow \zeta$ , where  $\zeta = \lim \zeta_n$ . According to (b),  $E(\zeta_n | \mathcal{G}) \uparrow E(\zeta | \mathcal{G})$  (a.s.). Therefore (a.s.)  $E(\lim \zeta_n | \mathcal{G}) = E(\zeta | \mathcal{G}) = \lim_n E(\zeta_n | \mathcal{G}) = \lim_n E(\zeta_n | \mathcal{G}) \leq \lim_n E(\xi_n | \mathcal{G})$ .

Conclusion (e) follows from (d).

(f) If  $\xi_n \geq 0$ , by Property D\* we have

$$E\left(\sum_{k=1}^n \xi_k | \mathcal{G}\right) = \sum_{k=1}^n E(\xi_k | \mathcal{G}) \quad (\text{a.s.})$$

which, with (b), establishes the required result.

This completes the proof of the theorem.

We can now establish Property K\*. Let  $\eta = I_B$ ,  $B \in \mathcal{G}$ . Then, for every  $A \in \mathcal{G}$ ,

$$\int_A \xi \eta \, dP = \int_{A \cap B} \xi \, dP = \int_{A \cap B} E(\xi | \mathcal{G}) \, dP = \int_A I_B E(\xi | \mathcal{G}) \, dP = \int_A \eta E(\xi | \mathcal{G}) \, dP.$$

By the additivity of the Lebesgue integral, the equation

$$\int_A \xi \eta \, dP = \int_A \eta E(\xi | \mathcal{G}) \, dP, \quad A \in \mathcal{G}, \quad (8)$$

remains valid for the simple random variables  $\eta = \sum_{k=1}^n y_k I_{B_k}$ ,  $B_k \in \mathcal{G}$ . Therefore, by Property I (§6), we have

$$E(\xi \eta | \mathcal{G}) = \eta E(\xi | \mathcal{G}) \quad (\text{a.s.}) \quad (9)$$

for these random variables.

Now let  $\eta$  be any  $\mathcal{G}$ -measurable random variable with  $E|\eta| < \infty$ , and let  $\{\eta_n\}_{n \geq 1}$  be a sequence of simple  $\mathcal{G}$ -measurable random variables such that  $|\eta_n| \leq \eta$  and  $\eta_n \rightarrow \eta$ . Then by (9)

$$E(\xi \eta_n | \mathcal{G}) = \eta_n E(\xi | \mathcal{G}) \quad (\text{a.s.}).$$

It is clear that  $|\xi \eta_n| \leq |\xi \eta|$ , where  $E|\xi \eta| < \infty$ . Therefore  $E(\xi \eta_n | \mathcal{G}) \rightarrow E(\xi \eta | \mathcal{G})$  (a.s.) by Property (a). In addition, since  $E|\xi| < \infty$ , we have  $E(\xi | \mathcal{G})$  finite (a.s.) (see Property C\* and Property J of §6). Therefore  $\eta_n E(\xi | \mathcal{G}) \rightarrow \eta E(\xi | \mathcal{G})$  (a.s.). (The hypothesis that  $E(\xi | \mathcal{G})$  is finite, almost surely, is essential, since, according to the footnote on p. 172,  $0 \cdot \infty = 0$ , but if  $\eta_n = 1/n$ ,  $\eta \equiv 0$ , we have  $1/n \cdot \infty \not\rightarrow 0 \cdot \infty = 0$ .)



5. Here we consider the more detailed structure of conditional expectations  $E(\xi|\mathcal{G}_\eta)$ , which we also denote, as usual, by  $E(\xi|\eta)$ .

Since  $E(\xi|\eta)$  is a  $\mathcal{G}_\eta$ -measurable function, then by Theorem 3 of §4 (more precisely, by its obvious modification for extended random variables) there is a Borel function  $m = m(y)$  from  $\bar{R}$  to  $\bar{R}$  such that

$$m(\eta(\omega)) = E(\xi|\eta)(\omega) \quad (10)$$

for all  $\omega \in \Omega$ . We denote this function  $m(y)$  by  $E(\xi|\eta = y)$  and call it the *conditional expectation of  $\xi$  with respect to the event  $\{\eta = y\}$ , or the conditional expectation of  $\xi$  under the condition that  $\eta = y$ .*

Correspondingly we define

$$\int_A \xi dP = \int_A E(\xi|\eta) dP = \int_A m(\eta) dP, \quad A \in \mathcal{G}_\eta. \quad (11)$$

Therefore by Theorem 7 of §6 (on change of variable under the Lebesgue integral sign)

$$\int_{\{\omega: \eta \in B\}} m(\eta) dP = \int_B m(y) P_\eta(dy), \quad B \in \mathcal{B}(\bar{R}), \quad (12)$$

where  $P_\eta$  is the probability distribution of  $\eta$ . Consequently  $m = m(y)$  is a Borel function such that

$$\int_{\{\omega: \eta \in B\}} \xi dP = \int_B m(y) dP_\eta. \quad (13)$$

for every  $B \in \mathcal{B}(R)$ .

This remark shows that we can give a different definition of the conditional expectation  $E(\xi|\eta = y)$ .

**Definition 4.** Let  $\xi$  and  $\eta$  be random variables (possible, extended) and let  $E\xi$  be defined. The conditional expectation of the random variable  $\xi$  under the condition that  $\eta = y$  is any  $\mathcal{B}(\bar{R})$ -measurable function  $m = m(y)$  for which

$$\int_{\{\omega: \eta \in B\}} \xi dP = \int_B m(y) P_\eta(dy), \quad B \in \mathcal{B}(\bar{R}). \quad (14)$$

That such a function exists follows again from the Radon–Nikodým theorem if we observe that the set function

$$Q(B) = \int_{\{\omega: \eta \in B\}} \xi dP$$

is a signed measure absolutely continuous with respect to the measure  $P_\eta$ .

Now suppose that  $m(y)$  is a conditional expectation in the sense of Definition 4. Then if we again apply the theorem on change of variable under the Lebesgue integral sign, we obtain

$$\int_{\{\omega: \eta \in B\}} \xi = \int_B m(y) P_\eta(dy) = \int_{\{\omega: \eta \in B\}} m(\eta), \quad B \in \mathcal{B}(\bar{R}).$$

The function  $m(\eta)$  is  $\mathcal{G}_\eta$ -measurable, and the sets  $\{\omega: \eta \in B\}$ ,  $B \in \mathcal{B}(\bar{R})$ , exhaust the subsets of  $\mathcal{G}_\eta$ .

Hence it follows that  $m(\eta)$  is the expectation  $E(\xi|\eta)$ . Consequently if we know  $E(\xi|\eta = y)$  we can reconstruct  $E(\xi|\eta)$ , and conversely from  $E(\xi|\eta)$  we can find  $E(\xi|\eta = y)$ .

From an intuitive point of view, the conditional expectation  $E(\xi|\eta = y)$  is simpler and more natural than  $E(\xi|\eta)$ . However,  $E(\xi|\eta)$ , considered as a  $\mathcal{G}_\eta$ -measurable random variable, is more convenient to work with.

Observe that Properties A\*-K\* above and the conclusions of Theorem 2 can easily be transferred to  $E(\xi|\eta = y)$  (replacing "almost surely" by " $P_\eta$ -almost surely"). Thus, for example, Property K\* transforms as follows: if  $E|\xi| < \infty$  and  $E|\xi f(\eta)| < \infty$ , where  $f = f(y)$  is a  $\mathcal{B}(\bar{R})$  measurable function, then

$$E(\xi f(\eta)|\eta = y) = f(y)E(\xi|\eta = y) \quad (P_\eta\text{-a.s.}) \quad (15)$$

In addition (cf. Property J\*), if  $\xi$  and  $\eta$  are independent, then

$$E(\xi|\eta = y) = E\xi \quad (P_\eta\text{-a.s.}).$$

We also observe that if  $B \in \mathcal{B}(R^2)$  and  $\xi$  and  $\eta$  are independent, then

$$E[I_B(\xi, \eta)|\eta = y] = E I_B(\xi, y) \quad (P_\eta\text{-a.s.}), \quad (16)$$

and if  $\varphi = \varphi(x, y)$  is a  $\mathcal{B}(R^2)$ -measurable function such that  $E|\varphi(\xi, \eta)| < \infty$ , then

$$E[\varphi(\xi, \eta)|\eta = y] = E[\varphi(\xi, y)] \quad (P_\eta\text{-a.s.}).$$

To prove (16) we make the following observation. If  $B = B_1 \times B_2$ , the validity of (16) will follow from

$$\int_{\{\omega: \eta \in A\}} I_{B_1 \times B_2}(\xi, \eta) P(d\omega) = \int_{\{y \in A\}} E I_{B_1 \times B_2}(\xi, y) P_\eta(dy).$$

But the left-hand side is  $P\{\xi \in B_1, \eta \in A \cap B_2\}$ , and the right-hand side is  $P(\xi \in B_1)P(\eta \in A \cap B_2)$ ; their equality follows from the independence of  $\xi$  and  $\eta$ . In the general case the proof depends on an application of Theorem 1, §2, on monotone classes (cf. the corresponding part of the proof of Fubini's theorem).

**Definition 5.** The conditional probability of the event  $A \in \mathcal{F}$  under the condition that  $\eta = y$  (notation:  $P(A|\eta = y)$ ) is  $E(I_A|\eta = y)$ .

It is clear that  $P(A|\eta = y)$  can be defined as the  $\mathcal{B}(\bar{R})$ -measurable function such that

$$P(A \cap \{\eta \in B\}) = \int_B P(A|\eta = y) P_\eta(dy), \quad B \in \mathcal{B}(\bar{R}). \quad (17)$$

6. Let us calculate some examples of conditional probabilities and conditional expectations.

EXAMPLE 1. Let  $\eta$  be a discrete random variable with  $P(\eta = y_k) > 0$ ,  $\sum_{k=1}^{\infty} P(\eta = y_k) = 1$ . Then

$$P(A|\eta = y_k) = \frac{P(A \cap \{\eta = y_k\})}{P(\eta = y_k)}, \quad k \geq 1.$$

For  $y \notin \{y_1, y_2, \dots\}$  the conditional probability  $P(A|\eta = y)$  can be defined in any way, for example as zero.

If  $\xi$  is a random variable for which  $E\xi$  exists, then

$$E(\xi|\eta = y_k) = \frac{1}{P(\eta = y_k)} \int_{\{\omega: \eta = y_k\}} \xi dP.$$

When  $y \notin \{y_1, y_2, \dots\}$  the conditional expectation  $E(\xi|\eta = y)$  can be defined in any way (for example, as zero).

EXAMPLE 2. Let  $(\xi, \eta)$  be a pair of random variables whose distribution has a density  $f_{\xi\eta}(x, y)$ :

$$P\{(\xi, \eta) \in B\} = \int_B f_{\xi\eta}(x, y) dx dy, \quad B \in \mathcal{B}(R^2).$$

Let  $f_\xi(x)$  and  $f_\eta(y)$  be the densities of the probability distribution of  $\xi$  and  $\eta$  (see (6.46), (6.55) and (6.56).

Let us put

$$f_{\xi|\eta}(x|y) = \frac{f_{\xi\eta}(x, y)}{f_\eta(y)}, \quad (18)$$

taking  $f_{\xi|\eta}(x|y) = 0$  if  $f_\eta(y) = 0$ .

Then

$$P(\xi \in C|\eta = y) = \int_C f_{\xi|\eta}(x|y) dx, \quad C \in \mathcal{B}(R), \quad (19)$$

i.e.  $f_{\xi|\eta}(x|y)$  is the density of a conditional probability distribution.

In fact, in order to prove (19) it is enough to verify (17) for  $B \in \mathcal{B}(R)$ ,  $A = \{\xi \in C\}$ . By (6.43), (6.45) and Fubini's theorem,

$$\begin{aligned} \int_B \left[ \int_C f_{\xi|\eta}(x|y) dx \right] P_\eta(dy) &= \int_B \left[ \int_C f_{\xi|\eta}(x|y) dx \right] f_\eta(y) dy \\ &= \int_{C \times B} f_{\xi|\eta}(x|y) f_\eta(y) dx dy \\ &= \int_{C \times B} f_{\xi\eta}(x, y) dx dy \\ &= P\{(\xi, \eta) \in C \times B\} = P\{(\xi \in C) \cap (\eta \in B)\}, \end{aligned}$$

which proves (17).

In a similar way we can show that if  $E\xi$  exists, then

$$E(\xi|\eta = y) = \int_{-\infty}^{\infty} x f_{\xi|\eta}(x|y) dx. \quad (20)$$

**EXAMPLE 3.** Let the length of time that a piece of apparatus will continue to operate be described by a nonnegative random variable  $\eta = \eta(\omega)$  whose distribution  $F_\eta(y)$  has a density  $f_\eta(y)$  (naturally,  $F_\eta(y) = f_\eta(y) = 0$  for  $y < 0$ ). Find the conditional expectation  $E(\eta - a|\eta \geq a)$ , i.e. the average time for which the apparatus will continue to operate on the hypothesis that it has already been operating for time  $a$ .

Let  $P(\eta \geq a) > 0$ . Then according to the definition (see Subsection 1) and (6.45),

$$\begin{aligned} E(\eta - a|\eta \geq a) &= \frac{E[(\eta - a)I_{(\eta \geq a)}]}{P(\eta \geq a)} = \frac{\int_{\Omega} (\eta - a)I_{(\eta \geq a)} P(d\omega)}{P(\eta \geq a)} \\ &= \frac{\int_a^{\infty} (y - a)f_\eta(y) dy}{\int_a^{\infty} f_\eta(y) dy}. \end{aligned}$$

It is interesting to observe that if  $\eta$  is exponentially distributed, i.e.

$$f_\eta(y) = \begin{cases} \lambda e^{-\lambda y}, & y \geq 0, \\ 0 & y < 0, \end{cases} \quad (21)$$

then  $E\eta = E(\eta|\eta \geq 0) = 1/\lambda$  and  $E(\eta - a|\eta \geq a) = 1/\lambda$  for every  $a > 0$ . In other words, in this case the average time for which the apparatus continues to operate, assuming that it has already operated for time  $a$ , is independent of  $a$  and simply equals the average time  $E\eta$ .

Under the assumption (21) we can find the conditional distribution  $P(\eta - a \leq x|\eta \geq a)$ .

We have

$$\begin{aligned}
 P(\eta - a \leq x | \eta \geq a) &= \frac{P(a \leq \eta \leq a + x)}{P(\eta \geq a)} \\
 &= \frac{F_\eta(a + x) - F_\eta(a) + P(\eta = a)}{1 - F_\eta(a) + P(\eta = a)} \\
 &= \frac{[1 - e^{-\lambda(a+x)}] - [1 - e^{-\lambda a}]}{1 - [1 - e^{-\lambda a}]} \\
 &= \frac{e^{-\lambda a}[1 - e^{-\lambda x}]}{e^{-\lambda a}} = 1 - e^{-\lambda x}.
 \end{aligned}$$

Therefore the conditional distribution  $P(\eta - a \leq x | \eta \geq a)$  is the same as the unconditional distribution  $P(\eta \leq x)$ . This remarkable property is unique to the exponential distribution: there are no other distributions that have densities and possess the property  $P(\eta - a \leq x | \eta \geq a) = P(\eta \leq x)$ ,  $a \geq 0$ ,  $0 \leq x < \infty$ .

**EXAMPLE 4 (Buffon's needle).** Suppose that we toss a needle of unit length "at random" onto a pair of parallel straight lines, a unit distance apart, in a plane. What is the probability that the needle will intersect at least one of the lines?

To solve this problem we must first define what it means to toss the needle "at random." Let  $\xi$  be the distance from the midpoint of the needle to the left-hand line. We shall suppose that  $\xi$  is uniformly distributed on  $[0, 1]$ , and (see Figure 29) that the angle  $\theta$  is uniformly distributed on  $[-\pi/2, \pi/2]$ . In addition, we shall assume that  $\xi$  and  $\theta$  are independent.

Let  $A$  be the event that the needle intersects one of the lines. It is easy to see that if

$$B = \{(a, x): |a| \leq \frac{\pi}{2}, \quad x \in [0, \frac{1}{2}\cos a] \cup [1 - \frac{1}{2}\cos a, 1]\},$$

then  $A = \{\omega: (\theta, \xi) \in B\}$ , and therefore the probability in question is

$$P(A) = EI_A(\omega) = EI_B(\theta(\omega), \xi(\omega)).$$

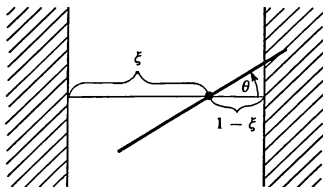


Figure 29

By Property  $G^*$  and formula (16),

$$\begin{aligned} E I_B(\theta(\omega), \xi(\omega)) &= E(E[I_B(\theta(\omega), \xi(\omega)) | \theta(\omega)]) \\ &= \int_{\Omega} E[I_B(\theta(\omega), \xi(\omega)) | \theta(\omega)] P(d\omega) \\ &= \int_{-\pi/2}^{\pi/2} E[I_B(\theta(\omega), \xi(\omega)) | \theta(\omega) = \alpha] P_{\theta}(d\alpha) \\ &= \frac{1}{\pi} \int_{-\pi/2}^{\pi/2} E I_B(a, \xi(\omega)) da = \frac{1}{\pi} \int_{-\pi/2}^{\pi/2} \cos a da = \frac{2}{\pi}, \end{aligned}$$

where we have used the fact that

$$E I_B(a, \xi(\omega)) = P\{\xi \in [0, \frac{1}{2} \cos a] \cup [1 - \frac{1}{2} \cos a]\} = \cos a.$$

Thus the probability that a "random" toss of the needle intersects one of the lines is  $2/\pi$ . This result could be used as the basis for an experimental evaluation of  $\pi$ . In fact, let the needle be tossed  $N$  times independently. Define  $\xi_i$  to be 1 if the needle intersects a line on the  $i$ th toss, and 0 otherwise. Then by the law of large numbers (see, for example, (I.5.6))

$$P\left\{\left|\frac{\xi_1 + \cdots + \xi_N}{N} - P(A)\right| > \varepsilon\right\} \rightarrow 0, \quad N \rightarrow \infty.$$

for every  $\varepsilon > 0$ .

In this sense the frequency satisfies

$$\frac{\xi_1 + \cdots + \xi_N}{N} \approx P(A) = \frac{2}{\pi}$$

and therefore

$$\frac{2N}{\xi_1 + \cdots + \xi_N} \approx \pi.$$

This formula has actually been used for a statistical evaluation of  $\pi$ . In 1850, R. Wolf (an astronomer in Zurich) threw a needle 5000 times and obtained the value 3.1596 for  $\pi$ . Apparently this problem was one of the first applications (now known as Monte Carlo methods) of probabilistic-statistical regularities to numerical analysis.

7. If  $\{\xi_n\}_{n \geq 1}$  is a sequence of nonnegative random variables, then according to conclusion (f) of Theorem 2,

$$E(\sum \xi_n | \mathcal{G}) = \sum E(\xi_n | \mathcal{G}) \quad (\text{a.s.})$$

In particular, if  $B_1, B_2, \dots$  is a sequence of pairwise disjoint sets,

$$P(\sum B_n | \mathcal{G}) = \sum P(B_n | \mathcal{G}) \quad (\text{a.s.}) \quad (22)$$

It must be emphasized that this equation is satisfied only almost surely and that consequently the conditional probability  $P(B|\mathcal{G})(\omega)$  cannot be considered as a measure on  $B$  for given  $\omega$ . One might suppose that, except for a set  $\mathcal{N}$  of measure zero,  $P(\cdot|\mathcal{G})(\omega)$  would still be a measure for  $\omega \in \mathcal{N}$ . However, in general this is not the case, for the following reason. Let  $\mathcal{N}(B_1, B_2, \dots)$  be the set of sample points  $\omega$  such that the countable additivity property (22) fails for these  $B_1, B_2, \dots$ . Then the excluded set  $\mathcal{N}$  is

$$\mathcal{N} = \bigcup \mathcal{N}(B_1, B_2, \dots), \quad (23)$$

where the union is taken over all  $B_1, B_2, \dots$  in  $\mathcal{F}$ . Although the  $P$ -measure of each set  $\mathcal{N}(B_1, B_2, \dots)$  is zero, the  $P$ -measure of  $\mathcal{N}$  can be different from zero (because of an uncountable union in (23)). (Recall that the Lebesgue measure of a single point is zero, but the measure of the set  $\mathcal{N} = [0, 1]$ , which is an uncountable sum of the individual points  $\{x\}$ , is 1).

However, it would be convenient if the conditional probability  $P(\cdot|\mathcal{G})(\omega)$  were a measure for each  $\omega \in \Omega$ , since then, for example, the calculation of conditional probabilities  $E(\xi|\mathcal{G})$  could be carried out (see Theorem 3 below) in a simple way by averaging with respect to the measure  $P(\cdot|\mathcal{G})(\omega)$ :

$$E(\xi|\mathcal{G}) = \int_{\Omega} \xi(\omega) P(d\omega|\mathcal{G}) \quad (\text{a.s.})$$

(cf. (I.8.10)).

We introduce the following definition.

**Definition 6.** A function  $P(\omega; B)$ , defined for all  $\omega \in \Omega$  and  $B \in \mathcal{F}$ , is a *regular conditional probability with respect to  $\mathcal{G}$*  if

- (a)  $P(\omega; \cdot)$  is a probability measure on  $\mathcal{F}$  for every  $\omega \in \Omega$ ;
- (b) For each  $B \in \mathcal{F}$  the function  $P(\omega; B)$ , as a function of  $\omega$ , is a variant of the conditional probability  $P(B|\mathcal{G})(\omega)$ , i.e.  $P(\omega; B) = P(B|\mathcal{G})(\omega)$  (a.s.).

**Theorem 3.** Let  $P(\omega; B)$  be a regular conditional probability with respect to  $\mathcal{G}$  and let  $\xi$  be an integrable random variable. Then

$$E(\xi|\mathcal{G})(\omega) = \int_{\Omega} \xi(\tilde{\omega}) P(\omega; d\tilde{\omega}) \quad (\text{a.s.}) \quad (24)$$

**PROOF.** If  $\xi = I_B$ ,  $B \in \mathcal{F}$ , the required formula (24) becomes

$$P(B|\mathcal{G})(\omega) = P(\omega; B) \quad (\text{a.s.}),$$

which holds by Definition 6(b). Consequently (24) holds for simple functions.

Now let  $\xi \geq 0$  and  $\xi_n \uparrow \xi$ , where  $\xi_n$  are simple functions. Then by (b) of Theorem 2 we have  $E(\xi|\mathcal{G})(\omega) = \lim_n E(\xi_n|\mathcal{G})(\omega)$  (a.s.). But since  $P(\omega; \cdot)$  is a measure for every  $\omega \in \Omega$ , we have

$$\lim_n E(\xi_n|\mathcal{G})(\omega) = \lim_n \int_{\Omega} \xi_n(\tilde{\omega}) P(\omega; d\tilde{\omega}) = \int_{\Omega} \xi(\tilde{\omega}) P(\omega; d\tilde{\omega})$$

by the monotone convergence theorem.

The general case reduces to this one if we use the representation  $\xi = \xi^+ - \xi^-$ .

This completes the proof.

**Corollary.** Let  $\mathcal{G} = \mathcal{G}_\eta$ , where  $\eta$  is a random variable, and let the pair  $(\xi, \eta)$  have a probability distribution with density  $f_{\xi\eta}(x, y)$ . Let  $E|g(\xi)| < \infty$ . Then

$$E(g(\xi)|\eta = y) = \int_{-\infty}^{\infty} g(x) f_{\xi|\eta}(x|y) dx,$$

where  $f_{\xi|\eta}(x|y)$  is the density of the conditional distribution (see (18)).

In order to be able to state the basic result on the existence of regular conditional probabilities, we need the following definitions.

**Definition 7.** Let  $(E, \mathcal{E})$  be a measurable space,  $X = X(\omega)$  a random element with values in  $E$ , and  $\mathcal{G}$  a  $\sigma$ -subalgebra of  $\mathcal{F}$ . A function  $Q(\omega; B)$ , defined for  $\omega \in \Omega$  and  $B \in \mathcal{E}$  is a *regular conditional distribution of  $X$  with respect to  $\mathcal{G}$*  if

- (a) for each  $\omega \in \Omega$  the function  $Q(\omega; B)$  is a probability measure on  $(E, \mathcal{E})$ ;
- (b) for each  $B \in \mathcal{E}$  the function  $Q(\omega; B)$ , as a function of  $\omega$ , is a variant of the conditional probability  $P(X \in B|\mathcal{G})(\omega)$ , i.e.

$$Q(\omega; B) = P(X \in B|\mathcal{G})(\omega) \quad (\text{a.s.}).$$

**Definition 8.** Let  $\xi$  be a random variable. A function  $F = F(\omega; x)$ ,  $\omega \in \Omega$ ,  $x \in R$ , is a *regular distribution function for  $\xi$  with respect to  $\mathcal{G}$*  if :

- (a)  $F(\omega; x)$  is, for each  $\omega \in \Omega$ , a distribution function on  $R$ ;
- (b)  $F(\omega; x) = P(\xi \leq x|\mathcal{G})(\omega)$  (a.s.), for each  $x \in R$ .

**Theorem 4.** A regular distribution function and a regular conditional distribution function always exist for the random variable  $\xi$  with respect to  $\mathcal{G}$ .



PROOF. For each rational number  $r \in R$ , define  $F_r(\omega) = P(\xi \leq r | \mathcal{G})(\omega)$ , where  $P(\xi \leq r | \mathcal{G})(\omega) = E(I_{\{\xi \leq r\}} | \mathcal{G})(\omega)$  is any variant of the conditional probability, with respect to  $\mathcal{G}$ , of the event  $\{\xi \leq r\}$ . Let  $\{r_i\}$  be the set of rational numbers in  $R$ . If  $r_i < r_j$ , Property **B\*** implies that  $P(\xi \leq r_i | \mathcal{G}) \leq P(\xi \leq r_j | \mathcal{G})$  (a.s.), and therefore if  $A_{ij} = \{\omega: F_{r_j}(\omega) < F_{r_i}(\omega)\}$ ,  $A = \bigcup A_{ij}$ , we have  $P(A) = 0$ . In other words, the set of points  $\omega$  at which the distribution function  $F_r(\omega)$ ,  $r \in \{r_i\}$ , fails to be monotonic has measure zero.

Now let

$$B_i = \left\{ \omega: \lim_{n \rightarrow \infty} F_{r_i + (1/n)}(\omega) \neq F_{r_i}(\omega) \right\}, \quad B = \bigcup_{i=1}^{\infty} B_i.$$

It is clear that  $I_{\{\xi \leq r_i + (1/n)\}} \downarrow I_{\{\xi \leq r_i\}}$ ,  $n \rightarrow \infty$ . Therefore, by (a) of Theorem 2,  $F_{r_i + (1/n)}(\omega) \rightarrow F_{r_i}(\omega)$  (a.s.), and therefore the set  $B$  on which continuity on the right fails (with respect to the rational numbers) also has measure zero,  $P(B) = 0$ .

In addition, let

$$C = \left\{ \omega: \lim_{n \rightarrow \infty} F_n(\omega) \neq 1 \right\} \cup \left\{ \omega: \lim_{n \rightarrow -\infty} F_n(\omega) > 0 \right\}.$$

Then, since  $\{\xi \leq n\} \uparrow \Omega$ ,  $n \rightarrow \infty$ , and  $\{\xi \leq n\} \downarrow \emptyset$ ,  $n \rightarrow -\infty$ , we have  $P(C) = 0$ .

Now put

$$F(\omega; x) = \begin{cases} \lim_{r \downarrow x} F_r(\omega), & \omega \notin A \cup B \cup C, \\ G(x), & \omega \in A \cup B \cup C, \end{cases}$$

where  $G(x)$  is any distribution function on  $R$ ; we show that  $F(\omega; x)$  satisfies the conditions of Definition 8.

Let  $\omega \notin A \cup B \cup C$ . Then it is clear that  $F(\omega; x)$  is a nondecreasing function of  $x$ . If  $x < x' \leq r$ , then  $F(\omega; x) \leq F(\omega; x') \leq F(\omega; r) = F_r(\omega) \downarrow F(\omega; x)$  when  $r \downarrow x$ . Consequently  $F(\omega; x)$  is continuous on the right. Similarly  $\lim_{x \rightarrow \infty} F(\omega; x) = 1$ ,  $\lim_{x \rightarrow -\infty} F(\omega; x) = 0$ . Since  $F(\omega; x) = G(x)$  when  $\omega \in A \cup B \cup C$ , it follows that  $F(\omega; x)$  is a distribution function on  $R$  for every  $\omega \in \Omega$ , i.e. condition (a) of Definition 8 is satisfied.

By construction,  $P(\xi \leq r | \mathcal{G})(\omega) = F_r(\omega) = F(\omega; r)$ . If  $r \downarrow x$ , we have  $F(\omega; r) \downarrow F(\omega; x)$  for all  $\omega \in \Omega$  by the continuity on the right that we just established. But by conclusion (a) of Theorem 2, we have  $P(\xi \leq r | \mathcal{G})(\omega) \rightarrow P(\xi \leq x | \mathcal{G})(\omega)$  (a.s.). Therefore  $F(\omega; x) = P(\xi \leq x | \mathcal{G})(\omega)$  (a.s.), which establishes condition (b) of Definition 8.

We now turn to the proof of the existence of a regular conditional distribution of  $\xi$  with respect to  $\mathcal{G}$ .

Let  $F(\omega; x)$  be the function constructed above. Put

$$Q(\omega; B) = \int_B F(\omega; dx),$$

where the integral is a Lebesgue–Stieltjes integral. From the properties of the integral (see §6, Subsection 7), it follows that  $Q(\omega; B)$  is a measure on  $B$  for each given  $\omega \in \Omega$ . To establish that  $Q(\omega; B)$  is a variant of the conditional probability  $P(\xi \in B | \mathcal{G})(\omega)$ , we use the principle of appropriate sets.

Let  $\mathcal{C}$  be the collection of sets  $B$  in  $\mathcal{B}(R)$  for which  $Q(\omega; B) = P(\xi \in B | \mathcal{G})(\omega)$  (a.s.). Since  $F(\omega; x) = P(\xi \leq x | \mathcal{G})(\omega)$  (a.s.), the system  $\mathcal{C}$  contains the sets  $B$  of the form  $B = (-\infty, x]$ ,  $x \in R$ . Therefore  $\mathcal{C}$  also contains the intervals of the form  $(a, b]$ , and the algebra  $\mathcal{A}$  consisting of finite sums of disjoint sets of the form  $(a, b]$ . Then it follows from the continuity properties of  $Q(\omega; B)$  ( $\omega$  fixed) and from conclusion (b) of Theorem 2 that  $\mathcal{C}$  is a monotone class, and since  $\mathcal{A} \subseteq \mathcal{C} \subseteq \mathcal{B}(R)$ , we have, from Theorem 1 of §2,

$$\mathcal{B}(R) = \sigma(\mathcal{A}) \subseteq \sigma(\mathcal{C}) = \mu(\mathcal{C}) = \mathcal{C} \subseteq \mathcal{B}(R),$$

whence  $\mathcal{C} = \mathcal{B}(R)$ .

This completes the proof of the theorem.

By using topological considerations we can extend the conclusion of Theorem 4 on the existence of a regular conditional distribution to random elements with values in what are known as Borel spaces. We need the following definition.

**Definition 9.** A measurable space  $(E, \mathcal{E})$  is a *Borel space* if it is Borel equivalent to a Borel subset of the real line, i.e. there is a one-to-one mapping  $\varphi = \varphi(e): (E, \mathcal{E}) \rightarrow (R, \mathcal{B}(R))$  such that

- (1)  $\varphi(E) \equiv \{\varphi(e): e \in E\}$  is a set in  $\mathcal{B}(R)$ ;
- (2)  $\varphi$  is  $\mathcal{E}$ -measurable ( $\varphi^{-1}(A) \in \mathcal{E}$ ,  $A \in \varphi(E) \cap \mathcal{B}(R)$ ),
- (3)  $\varphi^{-1}$  is  $\mathcal{B}(R)/\mathcal{E}$ -measurable ( $\varphi(B) \in \varphi(E) \cap \mathcal{B}(R)$ ,  $B \in \mathcal{E}$ ).

**Theorem 5.** Let  $X = X(\omega)$  be a random element with values in the Borel space  $(E, \mathcal{E})$ . Then there is a regular conditional distribution of  $X$  with respect to  $\mathcal{G}$ .

**PROOF.** Let  $\varphi = \varphi(e)$  be the function in Definition 9. By (2),  $\varphi(X(\omega))$  is a random variable. Hence, by Theorem 4, we can define the conditional distribution  $Q(\omega; A)$  of  $\varphi(X(\omega))$  with respect to  $\mathcal{G}$ ,  $A \in \varphi(E) \cap \mathcal{B}(R)$ .

We introduce the function  $\tilde{Q}(\omega; B) = Q(\omega; \varphi(B))$ ,  $B \in \mathcal{E}$ . By (3) of Definition 9,  $\varphi(B) \in \varphi(E) \cap \mathcal{B}(R)$  and consequently  $\tilde{Q}(\omega; B)$  is defined. Evidently  $\tilde{Q}(\omega; B)$  is a measure on  $B \in \mathcal{E}$  for every  $\omega$ . Now fix  $B \in \mathcal{E}$ . By the one-to-one character of the mapping  $\varphi = \varphi(e)$ ,

$$\tilde{Q}(\omega; B) = Q(\omega; \varphi(B)) = P\{\varphi(X) \in \varphi(B) | \mathcal{G}\} = P\{X \in B | \mathcal{G}\} \quad (\text{a.s.})$$

Therefore  $\tilde{Q}(\omega; B)$  is a regular conditional distribution of  $X$  with respect to  $\mathcal{G}$ .

This completes the proof of the theorem.

**Corollary.** Let  $X = X(\omega)$  be a random element with values in a complete separable metric space  $(E, \mathcal{E})$ . Then there is a regular conditional distribution of  $X$  with respect to  $\mathcal{G}$ . In particular, such a distribution exists for the spaces  $(R^n, \mathcal{B}(R^n))$  and  $(R^\infty, \mathcal{B}(R^\infty))$ .

The proof follows from Theorem 5 and the well known topological result that such spaces are Borel spaces.

8. The theory of conditional expectations developed above makes it possible to give a generalization of Bayes's theorem; this has applications in statistics.

Recall that if  $\mathcal{D} = \{A_1, \dots, A_n\}$  is a partition of the space  $\Omega$  with  $P(A_i) > 0$ , Bayes's theorem (I.3.9) states that

$$P(A_i|B) = \frac{P(A_i)P(B|A_i)}{\sum_{j=1}^n P(A_j)P(B|A_j)}. \quad (25)$$

for every  $B$  with  $P(B) > 0$ . Therefore if  $\theta = \sum_{i=1}^n a_i I_{A_i}$  is a discrete random variable then, according to (I.8.10),

$$E[g(\theta)|B] = \frac{\sum_{i=1}^n g(a_i)P(A_i)P(B|A_i)}{\sum_{j=1}^n P(A_j)P(B|A_j)}, \quad (26)$$

or

$$E[g(\theta)|B] = \frac{\int_{-\infty}^{\infty} g(a)P(B|\theta = a)P_\theta(da)}{\int_{-\infty}^{\infty} P(B|\theta = a)P_\theta(da)}. \quad (27)$$

On the basis of the definition of  $E[g(\theta)|B]$  given at the beginning of this section, it is easy to establish that (27) holds for all events  $B$  with  $P(B) > 0$ , random variables  $\theta$  and functions  $g = g(a)$  with  $E|g(\theta)| < \infty$ .

We now consider an analog of (27) for conditional expectations  $E[g(\theta)|\mathcal{G}]$  with respect to a  $\sigma$ -algebra  $\mathcal{G}$ ,  $\mathcal{G} \subseteq \mathcal{F}$ .

Let

$$Q(B) = \int_B g(\theta)P(d\omega), \quad B \in \mathcal{G}. \quad (28)$$

Then by (4)

$$E[g(\theta)|\mathcal{G}] = \frac{dQ}{dP}(\omega). \quad (29)$$

We also consider the  $\sigma$ -algebra  $\mathcal{G}_\theta$ . Then, by (5),

$$P(B) = \int_{\Omega} P(B|\mathcal{G}_\theta) dP \quad (30)$$

or, by the formula for change of variable in Lebesgue integrals,

$$P(B) = \int_{-\infty}^{\infty} P(B|\theta = a)P_\theta(da). \quad (31)$$

Since

$$Q(B) = E[g(\theta)I_B] = E[g(\theta) \cdot E(I_B|\mathcal{G}_\theta)],$$

we have

$$Q(B) = \int_{-\infty}^{\infty} g(a)P(B|\theta = a)P_\theta(da). \quad (32)$$

Now suppose that the conditional probability  $P(B|\theta = a)$  is regular and admits the representation

$$P(B|\theta = a) = \int_B \rho(\omega; a)\lambda(d\omega), \quad (33)$$

where  $\rho = \rho(\omega; a)$  is nonnegative and measurable in the two variables jointly, and  $\lambda$  is a  $\sigma$ -finite measure on  $(\Omega, \mathcal{G})$ .

Let  $E|g(\theta)| < \infty$ . Let us show that (P-a.s.)

$$E[g(\theta)|\mathcal{G}] = \frac{\int_{-\infty}^{\infty} g(a)\rho(\omega; a)P_\theta(da)}{\int_{-\infty}^{\infty} \rho(\omega; a)P_\theta(da)} \quad (34)$$

(generalized Bayes theorem).

In proving (34) we shall need the following lemma.

**Lemma.** Let  $(\Omega, \mathcal{F})$  be a measurable space.

(a) Let  $\mu$  and  $\lambda$  be  $\sigma$ -finite measures, and  $f = f(\omega)$  an  $\mathcal{F}$ -measurable function. Then

$$\int_{\Omega} f d\mu = \int_{\Omega} f \frac{d\mu}{d\lambda} \cdot d\lambda \quad (35)$$

(in the sense that if either integral exists, the other exists and they are equal).

(b) If  $\nu$  is a signed measure and  $\mu, \lambda$  are  $\sigma$ -finite measures  $\nu \ll \mu, \mu \ll \lambda$ , then

$$\frac{d\nu}{d\lambda} = \frac{d\nu}{d\mu} \cdot \frac{d\mu}{d\lambda} \quad (\lambda\text{-a.s.}) \quad (36)$$

and

$$\frac{d\nu}{d\mu} = \frac{d\nu}{d\lambda} \bigg/ \frac{d\mu}{d\lambda} \quad (\mu\text{-a.s.}) \quad (37)$$

**PROOF.** (a) Since

$$\mu(A) = \int_A \left( \frac{d\mu}{d\lambda} \right) d\lambda, \quad A \in \mathcal{F},$$

(35) is evidently satisfied for simple functions  $f = \sum f_i I_{A_i}$ . The general case follows from the representation  $f = f^+ - f^-$  and the monotone convergence theorem.

(b) From (a) with  $f = dv/d\mu$  we obtain

$$\nu(A) = \int_A \left( \frac{dv}{d\mu} \right) d\mu = \int_A \left( \frac{dv}{d\mu} \right) \cdot \left( \frac{d\mu}{d\lambda} \right) \cdot d\lambda.$$

Then  $\nu \ll \lambda$  and therefore

$$\nu(A) = \int_A \frac{dv}{d\lambda} d\lambda,$$

whence (36) follows since  $A$  is arbitrary, by Property I (§6).

Property (37) follows from (36) and the remark that

$$\mu \left\{ \omega: \frac{d\mu}{d\lambda} = 0 \right\} = \int_{\{\omega: d\mu/d\lambda = 0\}} \frac{d\mu}{d\lambda} d\lambda = 0$$

(on the set  $\{\omega: d\mu/d\lambda = 0\}$  the right-hand side of (37) can be defined arbitrarily, for example as zero). This completes the proof of the lemma.

To prove (34) we observe that by Fubini's theorem and (33),

$$Q(B) = \int_B \left[ \int_{-\infty}^{\infty} g(a) \rho(\omega; a) P_\theta(da) \right] \lambda(d\omega), \quad (38)$$

$$P(B) = \int_B \left[ \int_{-\infty}^{\infty} \rho(\omega; a) P_\theta(da) \right] \lambda(d\omega). \quad (39)$$

Then by the lemma

$$\frac{dQ}{dP} = \frac{dQ/d\lambda}{dP/d\lambda} \quad (\text{P-a.s.}).$$

Taking account of (38), (39) and (29), we have (34).

**Remark.** Formula (34) remains valid if we replace  $\theta$  by a random element with values in some measurable space  $(E, \mathcal{E})$  (and replace integration over  $R$  by integration over  $E$ ).

Let us consider some special cases of (34).

Let the  $\sigma$ -algebra  $\mathcal{G}$  be generated by the random variable  $\xi$ ,  $\mathcal{G} = \mathcal{G}_\xi$ . Suppose that

$$P(\xi \in A | \theta = a) = \int_A q(x; a) \lambda(dx), \quad A \in \mathcal{B}(R), \quad (40)$$

where  $q = q(x; a)$  is a nonnegative function, measurable with respect to both variables jointly, and  $\lambda$  is a  $\sigma$ -finite measure on  $(R, \mathcal{B}(R))$ . Then we obtain

$$E[g(\theta) | \xi = x] = \frac{\int_{-\infty}^{\infty} g(a) q(x; a) P_\theta(da)}{\int_{-\infty}^{\infty} q(x; a) P_\theta(da)}. \quad (41)$$

In particular, let  $(\theta, \xi)$  be a pair of discrete random variables,  $\theta = \sum a_i I_{A_i}$ ,  $\xi = \sum x_j I_{B_j}$ . Then, taking  $\lambda$  to be the counting measure ( $\lambda(\{x_i\}) = 1, i = 1, 2, \dots$ ) we find from (40) that

$$E[g(\theta)|\xi = x_j] = \frac{\sum_i g(a_i)P(\xi = x_j|\theta = a_i)P(\theta = a_i)}{\sum_i P(\xi = x_j|\theta = a_i)P(\theta = a_i)}. \quad (42)$$

(Compare (26).)

Now let  $(\theta, \xi)$  be a pair of absolutely continuous measures with density  $f_{\theta, \xi}(a, x)$ . Then by (19) the representation (40) applies with  $q(x; a) = f_{\xi|\theta}(x|a)$  and Lebesgue measure  $\lambda$ . Therefore

$$E[g(\theta)|\xi = x] = \frac{\int_{-\infty}^{\infty} g(a)f_{\xi|\theta}(x|a)f_{\theta}(a) da}{\int_{-\infty}^{\infty} f_{\xi|\theta}(x|a)f_{\theta}(a) da}. \quad (43)$$

## 9. PROBLEMS

1. Let  $\xi$  and  $\eta$  be independent identically distributed random variables with  $E\xi$  defined. Show that

$$E(\xi|\xi + \eta) = E(\eta|\xi + \eta) = \frac{\xi + \eta}{2} \quad (\text{a.s.}).$$

2. Let  $\xi_1, \xi_2, \dots$  be independent identically distributed random variables with  $E|\xi_i| < \infty$ . Show that

$$E(\xi_1 | S_n, S_{n+1}, \dots) = \frac{S_n}{n} \quad (\text{a.s.}),$$

where  $S_n = \xi_1 + \dots + \xi_n$ .

3. Suppose that the random elements  $(X, Y)$  are such that there is a regular distribution  $P_x(B) = P(Y \in B | X = x)$ . Show that if  $E|g(X, Y)| < \infty$  then

$$E[g(X, Y)|X = x] = \int g(x, y)P_x(dy) \quad (P_x\text{-a.s.}).$$

4. Let  $\xi$  be a random variable with distribution function  $F_{\xi}(x)$ . Show that

$$E(\xi | \alpha < \xi \leq b) = \frac{\int_{\alpha}^b x dF_{\xi}(x)}{F_{\xi}(b) - F_{\xi}(\alpha)}$$

(assuming that  $F_{\xi}(b) - F_{\xi}(\alpha) > 0$ ).

5. Let  $g = g(x)$  be a convex Borel function with  $E|g(\xi)| < \infty$ . Show that Jensen's inequality

$$g(E(\xi|\mathcal{G})) \leq E(g(\xi)|\mathcal{G})$$

holds for the conditional expectations.

6. Show that a necessary and sufficient condition for the random variable  $\xi$  and the  $\sigma$ -algebra  $\mathcal{G}$  to be independent (i.e., the random variables  $\xi$  and  $I_B(\omega)$  are independent for every  $B \in \mathcal{G}$ ) is that  $E(g(\xi)|\mathcal{G}) = E g(\xi)$  for every Borel function  $g(x)$  with  $E|g(\xi)| < \infty$ .
7. Let  $\xi$  be a nonnegative random variable and  $\mathcal{G}$  a  $\sigma$ -algebra,  $\mathcal{G} \subseteq \mathcal{F}$ . Show that  $E(\xi|\mathcal{G}) < \infty$  (a.s.) if and only if the measure  $Q$ , defined on sets  $A \in \mathcal{G}$  by  $Q(A) = \int_A \xi dP$ , is  $\sigma$ -finite.

## §8. Random Variables. II

1. In the first chapter we introduced characteristics of simple random variables, such as the variance, covariance, and correlation coefficient. These extend similarly to the general case. Let  $(\Omega, \mathcal{F}, P)$  be a probability space and  $\xi = \xi(\omega)$  a random variable for which  $E\xi$  is defined.

The variance of  $\xi$  is

$$V\xi = E(\xi - E\xi)^2.$$

The number  $\sigma = +\sqrt{V\xi}$  is the *standard deviation*.

If  $\xi$  is a random variable with a Gaussian (normal) density

$$f_\xi(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-[(x-m)^2]/2\sigma^2}, \quad \sigma > 0, \quad -\infty < m < \infty, \quad (1)$$

the parameters  $m$  and  $\sigma$  in (1) are very simple:

$$m = E\xi, \quad \sigma^2 = V\xi.$$

Hence the probability distribution of this random variable  $\xi$ , which we call *Gaussian*, or *normally distributed*, is completely determined by its mean value  $m$  and variance  $\sigma^2$ . (It is often convenient to write  $\xi \sim \mathcal{N}(m, \sigma^2)$ .)

Now let  $(\xi, \eta)$  be a pair of random variables. Their covariance is

$$\text{cov}(\xi, \eta) = E(\xi - E\xi)(\eta - E\eta) \quad (2)$$

(assuming that the expectations are defined).

If  $\text{cov}(\xi, \eta) = 0$  we say that  $\xi$  and  $\eta$  are *uncorrelated*.

If  $V\xi > 0$  and  $V\eta > 0$ , the number

$$\rho(\xi, \eta) \equiv \frac{\text{cov}(\xi, \eta)}{\sqrt{V\xi \cdot V\eta}} \quad (3)$$

is the *correlation coefficient* of  $\xi$  and  $\eta$ .

The properties of variance, covariance, and correlation coefficient were investigated in §4 of Chapter I for simple random variables. In the general case these properties can be stated in a completely analogous way.

Let  $\xi = (\xi_1, \dots, \xi_n)$  be a random vector whose components have finite second moments. The *covariance matrix* of  $\xi$  is the  $n \times n$  matrix  $\mathbb{R} = \|R_{ij}\|$ , where  $R_{ij} = \text{cov}(\xi_i, \xi_j)$ . It is clear that  $\mathbb{R}$  is *symmetric*. Moreover, it is *nonnegative definite*, i.e.

$$\sum_{i,j=1}^n R_{ij} \lambda_i \lambda_j \geq 0$$

for all  $\lambda_i \in \mathbb{R}$ ,  $i = 1, \dots, n$ , since

$$\sum_{i,j} R_{ij} \lambda_i \lambda_j = \mathbb{E} \left[ \sum_{i=1}^n (\xi_i - \mathbb{E} \xi_i) \lambda_i \right]^2 \geq 0.$$

The following lemma shows that the converse is also true.

**Lemma.** *A necessary and sufficient condition that an  $n \times n$  matrix  $\mathbb{R}$  is the covariance matrix of a vector  $\xi = (\xi_1, \dots, \xi_n)$  is that the matrix is symmetric and nonnegative definite, or, equivalently, that there is an  $n \times k$  matrix  $A$  ( $1 \leq k \leq n$ ) such that*

$$\mathbb{R} = AA^T,$$

where  $T$  denotes the transpose.

**PROOF.** We showed above that every covariance matrix is symmetric and nonnegative definite.

Conversely, let  $\mathbb{R}$  be a matrix with these properties. We know from matrix theory that corresponding to every symmetric nonnegative definite matrix  $\mathbb{R}$  there is an orthogonal matrix  $\mathcal{O}$  (i.e.,  $\mathcal{O}\mathcal{O}^T = E$ , the unit matrix) such that

$$\mathcal{O}^T \mathbb{R} \mathcal{O} = D,$$

where

$$D = \begin{pmatrix} d_1 & & 0 \\ & \ddots & \\ 0 & & d_n \end{pmatrix}$$

is a diagonal matrix with nonnegative elements  $d_i$ ,  $i = 1, \dots, n$ .

It follows that

$$\mathbb{R} = \mathcal{O} D \mathcal{O}^T = (\mathcal{O} B) (B^T \mathcal{O}^T),$$

where  $B$  is the diagonal matrix with elements  $b_i = +\sqrt{d_i}$ ,  $i = 1, \dots, n$ . Consequently if we put  $A = \mathcal{O} B$  we have the required representation  $\mathbb{R} = AA^T$  for  $\mathbb{R}$ .

It is clear that every matrix  $AA^T$  is symmetric and nonnegative definite. Consequently we have only to show that  $\mathbb{R}$  is the covariance matrix of some random vector.

Let  $\eta_1, \eta_2, \dots, \eta_n$  be a sequence of independent normally distributed random variables,  $\mathcal{N}(0, 1)$ . (The existence of such a sequence follows, for example, from Corollary 1 of Theorem 1, §9, and in principle could easily



be derived from Theorem 2 of §3.) Then the random vector  $\xi = A\eta$  (vectors are thought of as column vectors) has the required properties. In fact,

$$\mathbf{E}\xi\xi^T = \mathbf{E}(A\eta)(A\eta)^T = A \cdot \mathbf{E}\eta\eta^T \cdot A^T = AEA^T = AA^T.$$

(If  $\zeta = \|\zeta_{ij}\|$  is a matrix whose elements are random variables,  $\mathbf{E}\zeta$  means the matrix  $\|\mathbf{E}\zeta_{ij}\|$ ).

This completes the proof of the lemma.

We now turn our attention to the two-dimensional Gaussian (normal) density

$$f_{\xi\eta}(x, y) = \frac{1}{2\pi\sigma_1\sigma_2\sqrt{1-\rho^2}} \exp\left\{-\frac{1}{2(1-\rho^2)}\left[\frac{(x-m_1)^2}{\sigma_1^2} - 2\rho\frac{(x-m_1)(y-m_2)}{\sigma_1\sigma_2} + \frac{(y-m_2)^2}{\sigma_2^2}\right]\right\}, \quad (4)$$

characterized by the five parameters  $m_1, m_2, \sigma_1, \sigma_2$  and  $\rho$  (cf. (3.14)), where  $|m_1| < \infty, |m_2| < \infty, \sigma_1 > 0, \sigma_2 > 0, |\rho| < 1$ . An easy calculation identifies these parameters:

$$m_1 = \mathbf{E}\xi, \quad \sigma_1^2 = \mathbf{V}\xi,$$

$$m_2 = \mathbf{E}\eta, \quad \sigma_2^2 = \mathbf{V}\eta,$$

$$\rho = \rho(\xi, \eta).$$

In §4 of Chapter I we explained that if  $\xi$  and  $\eta$  are uncorrelated ( $\rho(\xi, \eta) = 0$ ), it does not follow that they are independent. However, if the pair  $(\xi, \eta)$  is Gaussian, it does follow that if  $\xi$  and  $\eta$  are uncorrelated then they are independent.

In fact, if  $\rho = 0$  in (4), then

$$f_{\xi\eta}(x, y) = \frac{1}{2\pi\sigma_1\sigma_2} e^{-[(x-m_1)^2]/2\sigma_1^2} \cdot e^{-[(y-m_2)^2]/2\sigma_2^2}.$$

But by (6.55) and (4),

$$f_{\xi}(x) = \int_{-\infty}^{\infty} f_{\xi\eta}(x, y) dy = \frac{1}{\sqrt{2\pi}\sigma_1} e^{-[(x-m_1)^2]/2\sigma_1^2},$$

$$f_{\eta}(y) = \int_{-\infty}^{\infty} f_{\xi\eta}(x, y) dx = \frac{1}{\sqrt{2\pi}\sigma_2} e^{-[(y-m_2)^2]/2\sigma_2^2}.$$

Consequently

$$f_{\xi\eta}(x, y) = f_{\xi}(x) \cdot f_{\eta}(y),$$

from which it follows that  $\xi$  and  $\eta$  are independent (see the end of Subsection 9 of §6).

2. A striking example of the utility of the concept of conditional expectation (introduced in §7) is its application to the solution of the following problem which is connected with *estimation theory* (cf. Subsection 8 of §4 of Chapter I).

Let  $(\xi, \eta)$  be a pair of random variables such that  $\xi$  is observable but  $\eta$  is not. We ask how the unobservable component  $\eta$  can be "estimated" from the knowledge of observations of  $\xi$ .

To state the problem more precisely, we need to define the concept of an *estimator*. Let  $\varphi = \varphi(x)$  be a Borel function. We call the random variable  $\varphi(\xi)$  an estimator of  $\eta$  in terms of  $\xi$ , and  $E[\eta - \varphi(\xi)]^2$  the (mean square) error of this estimator. An estimator  $\varphi^*(\xi)$  is called *optimal* (in the mean-square sense) if

$$\Delta \equiv E[\eta - \varphi^*(\xi)]^2 = \inf_{\varphi} E[\eta - \varphi(\xi)]^2, \quad (5)$$

where  $\inf$  is taken over all Borel functions  $\varphi = \varphi(x)$ .

**Theorem 1.** *Let  $E\eta^2 < \infty$ . Then there is an optimal estimator  $\varphi^* = \varphi^*(\xi)$  and  $\varphi^*(x)$  can be taken to be the function*

$$\varphi^*(x) = E(\eta | \xi = x). \quad (6)$$

**PROOF.** Without loss of generality we may consider only estimators  $\varphi(\xi)$  for which  $E\varphi^2(\xi) < \infty$ . Then if  $\varphi(\xi)$  is such an estimator, and  $\varphi^*(\xi) = E(\eta | \xi)$ , we have

$$\begin{aligned} E[\eta - \varphi(\xi)]^2 &= E[(\eta - \varphi^*(\xi)) + (\varphi^*(\xi) - \varphi(\xi))]^2 \\ &= E[\eta - \varphi^*(\xi)]^2 + E[\varphi^*(\xi) - \varphi(\xi)]^2 \\ &\quad + 2E[(\eta - \varphi^*(\xi))(\varphi^*(\xi) - \varphi(\xi))] \geq E[\eta - \varphi^*(\xi)]^2, \end{aligned}$$

since  $E[\varphi^*(\xi) - \varphi(\xi)]^2 \geq 0$  and, by the properties of conditional expectations,

$$\begin{aligned} E[(\eta - \varphi^*(\xi))(\varphi^*(\xi) - \varphi(\xi))] &= E\{E[(\eta - \varphi^*(\xi))(\varphi^*(\xi) - \varphi(\xi)) | \xi]\} \\ &= E\{(\varphi^*(\xi) - \varphi(\xi))E(\eta - \varphi^*(\xi) | \xi)\} = 0. \end{aligned}$$

This completes the proof of the theorem.

**Remark.** It is clear from the proof that the conclusion of the theorem is still valid when  $\xi$  is not merely a random variable but any random element with values in a measurable space  $(E, \mathcal{E})$ . We would then assume that  $\varphi = \varphi(x)$  is an  $\mathcal{E}/\mathcal{B}(R)$ -measurable function.

Let us consider the form of  $\varphi^*(x)$  on the hypothesis that  $(\xi, \eta)$  is a Gaussian pair with density given by (4).

From (1), (4) and (7.10) we find that the density  $f_{\eta|\xi}(y|x)$  of the conditional probability distribution is given by

$$f_{\eta|\xi}(y|x) = \frac{1}{\sqrt{2\pi(1-\rho^2)\sigma_2}} e^{[(y-m(x))^2]/[2\sigma_2^2(1-\rho^2)]}, \quad (7)$$

where

$$m(x) = m_2 + \frac{\sigma_2}{\sigma_1} \rho \cdot (x - m_1). \quad (8)$$

Then by the Corollary of Theorem 3, §7,

$$E(\eta|\xi = x) = \int_{-\infty}^{\infty} y f_{\eta|\xi}(y|x) dy = m(x) \quad (9)$$

and

$$\begin{aligned} V(\eta|\xi = x) &\equiv E[(\eta - E(\eta|\xi = x))^2 | \xi = x] \\ &= \int_{-\infty}^{\infty} (y - m(x))^2 f_{\eta|\xi}(y|x) dy \\ &= \sigma_2^2(1 - \rho^2). \end{aligned} \quad (10)$$

Notice that the conditional variance  $V(\eta|\xi = x)$  is independent of  $x$  and therefore

$$\Delta = E[\eta - E(\eta|\xi)]^2 = \sigma_2^2(1 - \rho^2). \quad (11)$$

Formulas (9) and (11) were obtained under the assumption that  $V\xi > 0$  and  $V\eta > 0$ . However, if  $V\xi > 0$  and  $V\eta = 0$  they are still evidently valid.

Hence we have the following result (cf. (I.4.16) and (I.4.17)).

**Theorem 2.** *Let  $(\xi, \eta)$  be a Gaussian vector with  $V\xi > 0$ . Then the optimal estimator of  $\eta$  in terms of  $\xi$  is*

$$E(\eta|\xi) = E\eta + \frac{\text{cov}(\xi, \eta)}{V\xi}(\xi - E\xi), \quad (12)$$

and its error is

$$\Delta \equiv E[\eta - E(\eta|\xi)]^2 = V\eta - \frac{\text{cov}^2(\xi, \eta)}{V\xi}. \quad (13)$$

**Remark.** The curve  $y(x) = E(\eta|\xi = x)$  is the *curve of regression of  $\eta$  on  $\xi$*  or of  $\eta$  with respect to  $\xi$ . In the Gaussian case  $E(\eta|\xi = x) = a + bx$  and consequently the regression of  $\eta$  and  $\xi$  is linear. Hence it is not surprising that the right-hand sides of (12) and (13) agree with the corresponding parts of (I.4.6) and (I.4.17) for the optimal linear estimator and its error.

**Corollary.** Let  $\varepsilon_1$  and  $\varepsilon_2$  be independent Gaussian random variables with mean zero and unit variance, and

$$\xi = a_1\varepsilon_1 + a_2\varepsilon_2, \quad \eta = b_1\varepsilon_1 + b_2\varepsilon_2.$$

Then  $E\xi = E\eta = 0$ ,  $V\xi = a_1^2 + a_2^2$ ,  $V\eta = b_1^2 + b_2^2$ ,  $\text{cov}(\xi, \eta) = a_1b_1 + a_2b_2$ , and if  $a_1^2 + a_2^2 > 0$ , then

$$E(\eta|\xi) = \frac{a_1b_1 + a_2b_2}{a_1^2 + a_2^2} \xi, \quad (14)$$

$$\Delta = \frac{(a_1b_2 - a_2b_1)^2}{a_1^2 + a_2^2}. \quad (15)$$

**3.** Let us consider the problem of determining the distribution functions of random variables that are functions of other random variables.

Let  $\xi$  be a random variable with distribution function  $F_\xi(x)$  (and density  $f_\xi(x)$ , if it exists), let  $\varphi = \varphi(x)$  be a Borel function and  $\eta = \varphi(\xi)$ . Letting  $I_y = (-\infty, y)$ , we obtain

$$F_\eta(y) = P(\eta \leq y) = P(\varphi(\xi) \in I_y) = P(\xi \in \varphi^{-1}(I_y)) = \int_{\varphi^{-1}(I_y)} F_\xi(dx), \quad (16)$$

which expresses the distribution function  $F_\eta(y)$  in terms of  $F_\xi(x)$  and  $\varphi$ .

For example, if  $\eta = a\xi + b$ ,  $a > 0$ , we have

$$F_\eta(y) = P\left(\xi \leq \frac{y-b}{a}\right) = F_\xi\left(\frac{y-b}{a}\right). \quad (17)$$

If  $\eta = \xi^2$ , it is evident that  $F_\eta(y) = 0$  for  $y < 0$ , while for  $y \geq 0$

$$\begin{aligned} F_\eta(y) &= P(\xi^2 \leq y) = P(-\sqrt{y} \leq \xi \leq \sqrt{y}) \\ &= F_\xi(\sqrt{y}) - F_\xi(-\sqrt{y}) + P(\xi = -\sqrt{y}). \end{aligned} \quad (18)$$

We now turn to the problem of determining  $f_\eta(y)$ .

Let us suppose that the range of  $\xi$  is a (finite or infinite) open interval  $I = (a, b)$ , and that the function  $\varphi = \varphi(x)$ , with domain  $(a, b)$ , is continuously differentiable and either strictly increasing or strictly decreasing. We also suppose that  $\varphi'(x) \neq 0$ ,  $x \in I$ . Let us write  $h(y) = \varphi^{-1}(y)$  and suppose for definiteness that  $\varphi(x)$  is strictly increasing. Then when  $y \in \varphi(I)$ ,

$$\begin{aligned} F_\eta(y) &= P(\eta \leq y) = P(\varphi(\xi) \leq y) = P(\xi \leq \varphi^{-1}(y)) \\ &= P(\xi \leq h(y)) = \int_{-\infty}^{h(y)} f_\xi(x) dx. \end{aligned} \quad (19)$$

By Problem 15 of §6,

$$\int_{-\infty}^{h(y)} f_{\xi}(x) dx = \int_{-\infty}^y f_{\xi}(h(z))h'(z) dz \quad (20)$$

and therefore

$$f_{\eta}(y) = f_{\xi}(h(y))h'(y). \quad (21)$$

Similarly, if  $\varphi(x)$  is strictly decreasing,

$$f_{\eta}(y) = f_{\xi}(h(y))(-h'(y)).$$

Hence in either case

$$f_{\eta}(y) = f_{\xi}(h(y))|h'(y)|. \quad (22)$$

For example, if  $\eta = a\xi + b$ ,  $a \neq 0$ , we have

$$h(y) = \frac{y-b}{a} \quad \text{and} \quad f_{\eta}(y) = \frac{1}{|a|} f_{\xi}\left(\frac{y-b}{a}\right).$$

If  $\xi \sim \mathcal{N}(m, \sigma^2)$  and  $\eta = e^{\xi}$ , we find from (22) that

$$f_{\eta}(y) = \begin{cases} \frac{1}{\sqrt{2\pi}\sigma y} \exp\left[-\frac{\ln(y/M)^2}{2\sigma^2}\right], & y > 0, \\ 0 & y \leq 0, \end{cases} \quad (23)$$

with  $M = e^m$ .

A probability distribution with the density (23) is said to be *lognormal* (logarithmically normal).

If  $\varphi = \varphi(x)$  is neither strictly increasing nor strictly decreasing, formula (22) is inapplicable. However, the following generalization suffices for many applications.

Let  $\varphi = \varphi(x)$  be defined on the set  $\sum_{k=1}^n [a_k, b_k]$ , continuously differentiable and either strictly increasing or strictly decreasing on each open interval  $I_k = (a_k, b_k)$ , and with  $\varphi'(x) \neq 0$  for  $x \in I_k$ . Let  $h_k = h_k(y)$  be the inverse of  $\varphi(x)$  for  $x \in I_k$ . Then we have the following generalization of (22):

$$f_{\eta}(y) = \sum_{k=1}^n f_{\xi}(h_k(y))|h'_k(y)| \cdot I_{D_k}(y), \quad (24)$$

where  $D_k$  is the domain of  $h_k(y)$ .

For example, if  $\eta = \xi^2$  we can take  $I_1 = (-\infty, 0)$ ,  $I_2 = (0, \infty)$ , and find that  $h_1(y) = -\sqrt{y}$ ,  $h_2(y) = \sqrt{y}$ , and therefore

$$f_{\eta}(y) = \begin{cases} \frac{1}{2\sqrt{y}} [f_{\xi}(\sqrt{y}) + f_{\xi}(-\sqrt{y})], & y > 0, \\ 0, & y \leq 0. \end{cases} \quad (25)$$

We can observe that this result also follows from (18), since  $P(\xi = -\sqrt{y}) = 0$ . In particular, if  $\xi \sim \mathcal{N}(0, 1)$ ,

$$f_{\xi^2}(y) = \begin{cases} \frac{1}{\sqrt{2\pi y}} e^{-y/2}, & y > 0, \\ 0, & y \leq 0. \end{cases} \quad (26)$$

A straightforward calculation shows that

$$f_{|\xi|}(y) = \begin{cases} f_{\xi}(y) + f_{\xi}(-y), & y > 0, \\ 0, & y \leq 0. \end{cases} \quad (27)$$

$$f_{+\sqrt{|\xi|}}(y) = \begin{cases} 2y(f_{\xi}(y^2) + f_{\xi}(-y^2)), & y > 0, \\ 0, & y \leq 0. \end{cases} \quad (28)$$

4. We now consider functions of several random variables.

If  $\xi$  and  $\eta$  are random variables with joint distribution  $F_{\xi\eta}(x, y)$ , and  $\varphi = \varphi(x, y)$  is a Borel function, then if we put  $\zeta = \varphi(\xi, \eta)$  we see at once that

$$F_{\zeta}(z) = \int_{\{x, y: \varphi(x, y) \leq z\}} dF_{\xi\eta}(x, y). \quad (29)$$

For example, if  $\varphi(x, y) = x + y$ , and  $\xi$  and  $\eta$  are independent (and therefore  $F_{\xi\eta}(x, y) = F_{\xi}(x) \cdot F_{\eta}(y)$ ) then Fubini's theorem shows that

$$\begin{aligned} F_{\zeta}(z) &= \int_{\{x, y: x+y \leq z\}} dF_{\xi}(x) \cdot dF_{\eta}(y) \\ &= \int_{\mathbb{R}^2} I_{\{x+y \leq z\}}(x, y) dF_{\xi}(x) dF_{\eta}(y) \\ &= \int_{-\infty}^{\infty} dF_{\xi}(x) \left\{ \int_{-\infty}^{\infty} I_{\{x+y \leq z\}}(x, y) dF_{\eta}(y) \right\} = \int_{-\infty}^{\infty} F_{\eta}(z-x) dF_{\xi}(x) \end{aligned} \quad (30)$$

and similarly

$$F_{\zeta}(z) = \int_{-\infty}^{\infty} F_{\xi}(z-y) dF_{\eta}(y). \quad (31)$$

If  $F$  and  $G$  are distribution functions, the function

$$H(z) = \int_{-\infty}^{\infty} F(z-x) dG(x)$$

is denoted by  $F * G$  and called the *convolution* of  $F$  and  $G$ .

Thus the distribution function  $F_{\zeta}$  of the sum of two independent random variables  $\xi$  and  $\eta$  is the convolution of their distribution functions  $F_{\xi}$  and  $F_{\eta}$ :

$$F_{\zeta} = F_{\xi} * F_{\eta}.$$

It is clear that  $F_{\xi} * F_{\eta} = F_{\eta} * F_{\xi}$ .

Now suppose that the independent random variables  $\xi$  and  $\eta$  have densities  $f_\xi$  and  $f_\eta$ . Then we find from (31), with another application of Fubini's theorem, that

$$\begin{aligned} F_\zeta(z) &= \int_{-\infty}^{\infty} \left[ \int_{-\infty}^{z-y} f_\xi(u) du \right] f_\eta(y) dy \\ &= \int_{-\infty}^{\infty} \left[ \int_{-\infty}^z f_\xi(u-y) du \right] f_\eta(y) dy = \int_{-\infty}^z \left[ \int_{-\infty}^{\infty} f_\xi(u-y) f_\eta(y) dy \right] du, \end{aligned}$$

whence

$$f_\zeta(z) = \int_{-\infty}^{\infty} f_\xi(z-y) f_\eta(y) dy, \quad (32)$$

and similarly

$$f_\zeta(z) = \int_{-\infty}^{\infty} f_\eta(z-x) f_\xi(x) dx. \quad (33)$$

Let us see some examples of the use of these formulas.

Let  $\xi_1, \xi_2, \dots, \xi_n$  be a sequence of independent identically distributed random variables with the uniform density on  $[-1, 1]$ :

$$f(x) = \begin{cases} \frac{1}{2}, & |x| \leq 1, \\ 0, & |x| > 1. \end{cases}$$

Then by (32) we have

$$\begin{aligned} f_{\xi_1 + \xi_2}(x) &= \begin{cases} \frac{2-|x|}{4}, & |x| \leq 2, \\ 0, & |x| > 2, \end{cases} \\ f_{\xi_1 + \xi_2 + \xi_3}(x) &= \begin{cases} \frac{(3-|x|)^2}{16}, & 1 \leq |x| \leq 3, \\ \frac{3-x^2}{8}, & 0 \leq |x| \leq 1, \\ 0, & |x| > 3, \end{cases} \end{aligned}$$

and by induction

$$f_{\xi_1 + \dots + \xi_n}(x) = \begin{cases} \frac{1}{2^n(n-1)!} \sum_{k=0}^{[(n+x)/2]} (-1)^k C_n^k (n+x-2k)^{n-1}, & |x| \leq n, \\ 0, & |x| > n. \end{cases}$$

Now let  $\xi \sim \mathcal{N}(m_1, \sigma_1^2)$  and  $\eta \sim \mathcal{N}(m_2, \sigma_2^2)$ . If we write

$$\varphi(x) = \frac{1}{\sqrt{2\pi}} e^{-x^2/2},$$

then

$$f_{\xi}(x) = \frac{1}{\sigma_1} \varphi\left(\frac{x - m_1}{\sigma_1}\right), \quad f_{\eta}(x) = \frac{1}{\sigma_2} \varphi\left(\frac{x - m_2}{\sigma_2}\right),$$

and the formula

$$f_{\xi+\eta}(x) = \frac{1}{\sqrt{\sigma_1^2 + \sigma_2^2}} \varphi\left(\frac{x - (m_1 + m_2)}{\sqrt{\sigma_1^2 + \sigma_2^2}}\right)$$

follows easily from (32).

Therefore the sum of two independent Gaussian random variables is again a Gaussian random variable with mean  $m_1 + m_2$  and variance  $\sigma_1^2 + \sigma_2^2$ .

Let  $\xi_1, \dots, \xi_n$  be independent random variables each of which is normally distributed with mean 0 and variance 1. Then it follows easily from (26) (by induction) that

$$f_{\xi_1^2 + \dots + \xi_n^2}(x) = \begin{cases} \frac{1}{2^{n/2} \Gamma(n/2)} x^{(n/2)-1} e^{-x/2}, & x > 0, \\ 0, & x \leq 0. \end{cases} \quad (34)$$

The variable  $\xi_1^2 + \dots + \xi_n^2$  is usually denoted by  $\chi_n^2$ , and its distribution (with density (30)) is the  $\chi^2$ -distribution ("chi-square distribution") with  $n$  degrees of freedom (cf. Table 2 in §3).

If we write  $\chi_n = \sqrt{\chi_n^2}$ , it follows from (28) and (34) that

$$f_{\chi_n}(x) = \begin{cases} \frac{2x^{n-1} e^{-x^2/2}}{2^{n/2} \Gamma(n/2)}, & x \geq 0, \\ 0, & x < 0. \end{cases} \quad (35)$$

The probability distribution with this density is the  $\chi$ -distribution (chi-distribution) with  $n$  degrees of freedom.

Again let  $\xi$  and  $\eta$  be independent random variables with densities  $f_{\xi}$  and  $f_{\eta}$ . Then

$$F_{\xi\eta}(z) = \iint_{\{x, y: xy \leq z\}} f_{\xi}(x) f_{\eta}(y) dx dy,$$

$$F_{\xi/\eta}(z) = \iint_{\{x, y: x/y \leq z\}} f_{\xi}(x) f_{\eta}(y) dx dy.$$

Hence we easily obtain

$$f_{\xi\eta}(z) = \int_{-\infty}^{\infty} f_{\xi}\left(\frac{z}{y}\right) f_{\eta}(y) \frac{dy}{|y|} = \int_{-\infty}^{\infty} f_{\eta}\left(\frac{z}{x}\right) f_{\xi}(x) \frac{dx}{|x|} \quad (36)$$

and

$$f_{\xi/\eta}(z) = \int_{-\infty}^{\infty} f_{\xi}(zy) f_{\eta}(y) |y| dy. \quad (37)$$



Putting  $\xi = \xi_0$  and  $\eta = \sqrt{(\xi_1^2 + \dots + \xi_n^2)/n}$ , in (37), where  $\xi_0, \xi_1, \dots, \xi_n$  are independent Gaussian random variables with mean 0 and variance  $\sigma^2 > 0$ , and using (35), we find that

$$f_{\xi_0/\sqrt{(1/n)(\xi_1^2 + \dots + \xi_n^2)}}(x) = \frac{1}{\sqrt{\pi n}} \frac{\Gamma\left(\frac{n+1}{2}\right)}{\Gamma\left(\frac{n}{2}\right)} \frac{1}{\left(1 + \frac{x^2}{n}\right)^{(n+1)/2}}. \quad (38)$$

The variable  $\xi_0/\sqrt{(1/n)(\xi_1^2 + \dots + \xi_n^2)}$  is denoted by  $t$ , and its distribution is the  $t$ -distribution, or *Student's distribution*, with  $n$  degrees of freedom (cf. Table 2 in §3). Observe that this distribution is independent of  $\sigma$ .

## 5. PROBLEMS

1. Verify formulas (9), (10), (24), (27), (28), and (34)–(38).
2. Let  $\xi_1, \dots, \xi_n, n \geq 2$ , be independent identically distributed random variables with distribution function  $F(x)$  (and density  $f(x)$ , if it exists), and let  $\xi = \max(\xi_1, \dots, \xi_n)$ ,  $\underline{\xi} = \min(\xi_1, \dots, \xi_n)$ ,  $\rho = \xi - \underline{\xi}$ . Show that

$$F_{\xi, \underline{\xi}}(y, x) = \begin{cases} (F(y))^n - (F(y) - F(x))^n, & y > x, \\ (F(y))^n, & y \leq x, \end{cases}$$

$$f_{\xi, \underline{\xi}}(y, x) = \begin{cases} n(n-1)[F(y) - F(x)]^{n-2}f(x)f(y), & y > x, \\ 0, & y < x, \end{cases}$$

$$F_\rho(x) = \begin{cases} n \int_{-\infty}^{\infty} [F(y) - F(y-x)]^{n-1}f(y) dy, & x \geq 0, \\ 0, & x < 0, \end{cases}$$

$$f_\rho(x) = \begin{cases} n(n-1) \int_{-\infty}^{\infty} [F(y) - F(y-x)]^{n-2}f(y-x)f(y) dy, & x > 0, \\ 0, & x < 0. \end{cases}$$

3. Let  $\xi_1$  and  $\xi_2$  be independent Poisson random variables with respective parameters  $\lambda_1$  and  $\lambda_2$ . Show that  $\xi_1 + \xi_2$  has a Poisson distribution with parameter  $\lambda_1 + \lambda_2$ .
4. Let  $m_1 = m_2 = 0$  in (4). Show that

$$f_{\xi, \eta}(z) = \frac{\sigma_1 \sigma_2 \sqrt{1 - \rho^2}}{\pi(\sigma_2^2 z - 2\rho\sigma_1\sigma_2 z + \sigma_1^2)}.$$

5. The *maximal correlation coefficient* of  $\xi$  and  $\eta$  is  $\rho^*(\xi, \eta) = \sup_{u, v} \rho(u(\xi), v(\eta))$ , where the supremum is taken over the Borel functions  $u = u(x)$  and  $v = v(x)$  for which the correlation coefficient  $\rho(u(\xi), v(\eta))$  is defined. Show that  $\xi$  and  $\eta$  are independent if and only if  $\rho^*(\xi, \eta) = 0$ .
6. Let  $\tau_1, \tau_2, \dots, \tau_n$  be independent nonnegative identically distributed random variables with the exponential density

$$f(t) = \lambda e^{-\lambda t}, \quad t \geq 0.$$

Show that the distribution of  $\tau_1 + \cdots + \tau_k$  has the density

$$\frac{\lambda^k t^{k-1} e^{-\lambda t}}{(k-1)!}, \quad t \geq 0, \quad 1 \leq k \leq n,$$

and that

$$P(\tau_1 + \cdots + \tau_k > t) = \sum_{i=0}^{k-1} e^{-\lambda t} \frac{(\lambda t)^i}{i!}.$$

7. Let  $\xi \sim \mathcal{N}(0, \sigma^2)$ . Show that, for every  $p \geq 1$ ,

$$E|\xi|^p = C_p \sigma^p,$$

where

$$C_p = \frac{2^{p/2}}{\pi^{1/2}} \Gamma\left(\frac{p+1}{2}\right)$$

and  $\Gamma(s) = \int_0^\infty e^{-x} x^{s-1} dx$  is the gamma function. In particular, for each integer  $n \geq 1$ ,

$$E\xi^{2n} = (2n-1)!! \sigma^{2n}.$$

## §9. Construction of a Process with Given Finite-Dimensional Distribution

1. Let  $\xi = \xi(\omega)$  be a random variable defined on the probability space  $(\Omega, \mathcal{F}, P)$ , and let

$$F_\xi(x) = P\{\omega: \xi(\omega) \leq x\}$$

be its distribution function. It is clear that  $F_\xi(x)$  is a distribution function on the real line in the sense of Definition 1 of §3.

We now ask the following question. Let  $F = F(x)$  be a distribution function on  $R$ . Does there exist a random variable whose distribution function is  $F(x)$ ?

One reason for asking this question is as follows. Many statements in probability theory begin, "Let  $\xi$  be a random variable with the distribution function  $F(x)$ ; then ...". Consequently if a statement of this kind is to be meaningful we need to be certain that the object under consideration actually exists. Since to know a random variable we first have to know its domain  $(\Omega, \mathcal{F})$ , and in order to speak of its distribution we need to have a probability measure  $P$  on  $(\Omega, \mathcal{F})$ , a correct way of phrasing the question of the existence of a random variable with a given distribution function  $F(x)$  is this:

*Do there exist a probability space  $(\Omega, \mathcal{F}, P)$  and a random variable  $\xi = \xi(\omega)$  on it, such that*

$$P\{\omega: \xi(\omega) \leq x\} = F(x)?$$

Let us show that the answer is positive, and essentially contained in Theorem 1 of §1.

In fact, let us put

$$\Omega = R, \quad \mathcal{F} = \mathcal{B}(R).$$

It follows from Theorem 1 of §1 that there is a probability measure  $P$  (and only one) on  $(R, \mathcal{B}(R))$  for which  $P(a, b] = F(b) - F(a)$ ,  $a < b$ .

Put  $\xi(\omega) \equiv \omega$ . Then

$$P\{\omega: \xi(\omega) \leq x\} = P\{\omega: \omega \leq x\} = P(-\infty, x] = F(x).$$

Consequently we have constructed the required probability space and the random variable on it.

2. Let us now ask a similar question for random processes.

Let  $X = (\xi_t)_{t \in T}$  be a random process (in the sense of Definition 3, §5) defined on the probability space  $(\Omega, \mathcal{F}, P)$ , with  $t \in T \subseteq R$ .

From a physical point of view, the most fundamental characteristic of a random process is the set  $\{F_{t_1, \dots, t_n}(x_1, \dots, x_n)\}$  of its *finite-dimensional distribution functions*

$$F_{t_1, \dots, t_n}(x_1, \dots, x_n) = P\{\omega: \xi_{t_1} \leq x_1, \dots, \xi_{t_n} \leq x_n\}, \quad (1)$$

defined for all sets  $t_1, \dots, t_n$  with  $t_1 < t_2 < \dots < t_n$ .

We see from (1) that, for each set  $t_1, \dots, t_n$  with  $t_1 < t_2 < \dots < t_n$  the functions  $F_{t_1, \dots, t_n}(x_1, \dots, x_n)$  are  $n$ -dimensional distribution functions (in the sense of Definition 2, §3) and that the collection  $\{F_{t_1, \dots, t_n}(x_1, \dots, x_n)\}$  has the following *consistency* property:

$$\lim_{x_k \uparrow \infty} F_{t_1, \dots, t_n}(x_1, \dots, x_n) = F_{t_1, \dots, \hat{t}_k, \dots, t_n}(x_1, \dots, \hat{x}_k, \dots, x_n) \quad (2)$$

where  $\hat{\phantom{x}}$  indicates an omitted coordinate.

Now it is natural to ask the following question: under what conditions can a given family  $\{F_{t_1, \dots, t_n}(x_1, \dots, x_n)\}$  of distribution functions  $F_{t_1, \dots, t_n}(x_1, \dots, x_n)$  (in the sense of Definition 2, §3) be the family of finite-dimensional distribution functions of a random process? It is quite remarkable that all such conditions are covered by the consistency condition (2).

**Theorem 1** (Kolmogorov's Theorem on the Existence of a Process). *Let  $\{F_{t_1, \dots, t_n}(x_1, \dots, x_n)\}$ , with  $t_i \in T \subseteq R$ ,  $t_1 < t_2 < \dots < t_n$ ,  $n \geq 1$ , be a given family of finite-dimensional distribution functions, satisfying the consistency condition (2). Then there are a probability space  $(\Omega, \mathcal{F}, P)$  and a random process  $X = (\xi_t)_{t \in T}$  such that*

$$P\{\omega: \xi_{t_1} \leq x_1, \dots, \xi_{t_n} \leq x_n\} = F_{t_1, \dots, t_n}(x_1, \dots, x_n). \quad (3)$$

PROOF. Put

$$\Omega = R^T, \quad \mathcal{F} = \mathcal{B}(R^T),$$

i.e. take  $\Omega$  to be the space of real functions  $\omega = (\omega_t)_{t \in T}$  with the  $\sigma$ -algebra generated by the cylindrical sets.

Let  $\tau = [t_1, \dots, t_n]$ ,  $t_1 < t_2 < \dots < t_n$ . Then by Theorem 2 of §3 we can construct on the space  $(R^n, \mathcal{B}(R^n))$  a unique probability measure  $P_\tau$  such that

$$P_\tau\{(\omega_{t_1}, \dots, \omega_{t_n}): \omega_{t_1} \leq x_1, \dots, \omega_{t_n} \leq x_n\} = F_{t_1 \dots t_n}(x_1, \dots, x_n). \quad (4)$$

It follows from the consistency condition (2) that the family  $\{P_\tau\}$  is also consistent (see (3.20)). According to Theorem 4 of §3 there is a probability measure  $P$  on  $(R^T, \mathcal{B}(R^T))$  such that

$$P\{\omega: (\omega_{t_1}, \dots, \omega_{t_n}) \in B\} = P_\tau(B)$$

for every set  $\tau = [t_1, \dots, t_n]$ ,  $t_1 < \dots < t_n$ .

From this, it also follows that (4) is satisfied. Therefore the required random process  $X = (\xi_t(\omega))_{t \in T}$  can be taken to be the process defined by

$$\xi_t(\omega) = \omega_t, \quad t \in T. \quad (5)$$

This completes the proof of the theorem.

**Remark 1.** The probability space  $(R^T, \mathcal{B}(R^T), P)$  that we have constructed is called *canonical*, and the construction given by (5) is called the *coordinate method* of constructing the process.

**Remark 2.** Let  $(E_\alpha, \mathcal{E}_\alpha)$  be complete separable metric spaces, where  $\alpha$  belongs to some set  $\mathfrak{A}$  of indices. Let  $\{P_\tau\}$  be a set of consistent finite-dimensional distribution functions  $P_\tau$ ,  $\tau = [\alpha_1, \dots, \alpha_n]$  on

$$(E_{\alpha_1} \times \dots \times E_{\alpha_n}, \mathcal{E}_{\alpha_1} \otimes \dots \otimes \mathcal{E}_{\alpha_n}).$$

Then there are a probability space  $(\Omega, \mathcal{F}, P)$  and a family of  $\mathcal{F}/\mathcal{E}_\alpha$ -measurable functions  $(X_\alpha(\omega))_{\alpha \in \mathfrak{A}}$  such that

$$P\{(X_{\alpha_1}, \dots, X_{\alpha_n}) \in B\} = P_\tau(B)$$

for all  $\tau = [\alpha_1, \dots, \alpha_n]$  and  $B \in \mathcal{E}_{\alpha_1} \otimes \dots \otimes \mathcal{E}_{\alpha_n}$ .

This result, which generalizes Theorem 1, follows from Theorem 4 of §3 if we put  $\Omega = \prod_\alpha E_\alpha$ ,  $\mathcal{F} = \bigotimes_\alpha \mathcal{E}_\alpha$  and  $X_\alpha(\omega) = \omega_\alpha$  for each  $\omega = (\omega_\alpha)_{\alpha \in \mathfrak{A}}$ .

**Corollary 1.** Let  $F_1(x), F_2(x), \dots$  be a sequence of one-dimensional distribution functions. Then there exist a probability space  $(\Omega, \mathcal{F}, P)$  and a sequence of independent random variables  $\xi_1, \xi_2, \dots$  such that

$$P\{\omega: \xi_i(\omega) \leq x\} = F_i(x). \quad (6)$$

In particular, there is a probability space  $(\Omega, \mathcal{F}, P)$  on which an infinite sequence of Bernoulli random variables is defined (in this connection see Subsection 2 of §5 of Chapter I). Notice that  $\Omega$  can be taken to be the space

$$\Omega = \{\omega: \omega = (a_1, a_2, \dots), a_i = 0, 1\}$$

(cf. also Theorem 2).

To establish the corollary it is enough to put  $F_{1, \dots, n}(x_1, \dots, x_n) = F_1(x_1) \cdots F_n(x_n)$  and apply Theorem 1.

**Corollary 2.** Let  $T = [0, \infty)$  and let  $\{p(s, x; t, B)\}$  be a family of nonnegative functions defined for  $s, t \in T, t > s, x \in R, B \in \mathcal{B}(R)$ , and satisfying the following conditions:

- (a)  $p(s, x; t, B)$  is a probability measure on  $B$  for given  $s, x$  and  $t$ ;
- (b) for given  $s, t$  and  $B$ , the function  $p(s, x; t, B)$  is a Borel function of  $x$ ;
- (c) for  $0 \leq s < t < \tau$  and  $B \in \mathcal{B}(R)$ , the Kolmogorov–Chapman equation

$$p(s, x; \tau, B) = \int_R p(s, x; t, dy) p(t, y; \tau, B) \quad (7)$$

is satisfied.

Also let  $\pi = \pi(B)$  be a probability measure on  $(R, \mathcal{B}(R))$ . Then there are a probability space  $(\Omega, \mathcal{F}, P)$  and a random process  $X = (\xi_t)_{t \geq 0}$  defined on it, such that

$$\begin{aligned} P\{\xi_{t_0} \leq x_0, \xi_{t_1} \leq x_1, \dots, \xi_{t_n} \leq x_n\} &= \int_{-\infty}^{x_0} \pi(dy_0) \int_{-\infty}^{x_1} p(0, y_0; t_1, dy_1) \\ &\quad \cdots \int_{-\infty}^{x_n} p(t_{n-1}, y_{n-1}; t_n, dy_n) \end{aligned} \quad (8)$$

for  $0 = t_0 < t_1 < \dots < t_n$ .

The process  $X$  so constructed is a *Markov process* with initial distribution  $\pi$  and transition probabilities  $\{p(s, x; t, B)\}$ .

**Corollary 3.** Let  $T = \{0, 1, 2, \dots\}$  and let  $\{P_k(x; B)\}$  be a family of nonnegative functions defined for  $k \geq 1, x \in R, B \in \mathcal{B}(R)$ , and such that  $p_k(x; B)$  is a probability measure on  $B$  (for given  $k$  and  $x$ ) and measurable in  $x$  (for given  $k$  and  $B$ ). In addition, let  $\pi = \pi(B)$  be a probability measure on  $(R, \mathcal{B}(R))$ .

Then there is a probability space  $(\Omega, \mathcal{F}, P)$  with a family of random variables  $X = \{\xi_0, \xi_1, \dots\}$  defined on it, such that

$$\begin{aligned} P\{\xi_{t_0} \leq x_0, \xi_{t_1} \leq x_1, \dots, \xi_{t_n} \leq x_n\} &= \int_{-\infty}^{x_0} \pi(dy_0) \int_{-\infty}^{x_1} p(0, y_0; t_1, dy_1) \\ &\quad \cdots \int_{-\infty}^{x_n} p(t_{n-1}, y_{n-1}; t_n, dy_n) \end{aligned}$$

3. In the situation of Corollary 1, there is a sequence of independent random variables  $\xi_1, \xi_2, \dots$  whose one-dimensional distribution functions are  $F_1, F_2, \dots$ , respectively.

Now let  $(E_1, \mathcal{E}_1), (E_2, \mathcal{E}_2), \dots$  be complete separable metric spaces and let  $P_1, P_2, \dots$  be probability measures on them. Then it follows from Remark 2 that there are a probability space  $(\Omega, \mathcal{F}, P)$  and a sequence of independent elements  $X_1, X_2, \dots$  such that  $X_n$  is  $\mathcal{F}/\mathcal{E}_n$ -measurable and  $P(X_n \in B) = P_n(B), B \in \mathcal{E}_n$ .

It turns out that this result remains valid when the spaces  $(E_n, \mathcal{E}_n)$  are arbitrary measurable spaces.

**Theorem 2** (Ionescu Tulcea's Theorem on Extending a Measure and the Existence of a Random Sequence). *Let  $(\Omega_n, \mathcal{F}_n), n = 1, 2, \dots$ , be arbitrary measurable spaces and  $\Omega = \prod \Omega_n, \mathcal{F} = \bigotimes \mathcal{F}_n$ . Suppose that a probability measure  $P_1$  is given on  $(\Omega_1, \mathcal{F}_1)$  and that, for every set  $(\omega_1, \dots, \omega_n) \in \Omega_1 \times \dots \times \Omega_n, n \geq 1$ , probability measures  $P(\omega_1, \dots, \omega_n; \cdot)$  are given on  $(\Omega_{n+1}, \mathcal{F}_{n+1})$ . Suppose that for every  $B \in \mathcal{F}_{n+1}$  the functions  $P(\omega_1, \dots, \omega_n; B)$  are Borel functions on  $(\omega_1, \dots, \omega_n)$  and let*

$$P_n(A_1 \times \dots \times A_n) = \int_{A_1} P_1(d\omega_1) \int_{A_2} P(\omega_1; d\omega_2) \dots \int_{A_n} P(\omega_1, \dots, \omega_{n-1}; d\omega_n) \quad A_i \in \mathcal{F}_i, \quad n \geq 1. \quad (9)$$

Then there is a unique probability measure  $P$  on  $(\Omega, \mathcal{F})$  such that

$$P\{\omega: \omega_1 \in A_1, \dots, \omega_n \in A_n\} = P_n(A_1 \times \dots \times A_n) \quad (10)$$

for every  $n \geq 1$ , and there is a random sequence  $X = (X_1(\omega), X_2(\omega), \dots)$  such that

$$P\{\omega: X_1(\omega) \in A_1, \dots, X_n(\omega) \in A_n\} = P_n(A_1 \times \dots \times A_n), \quad (11)$$

where  $A_i \in \mathcal{E}_i$ .

**PROOF.** The first step is to establish that for each  $n > 1$  the set function  $P_n$  defined by (9) on the rectangle  $A_1 \times \dots \times A_n$  can be extended to the  $\sigma$ -algebra  $\mathcal{F}_1 \otimes \dots \otimes \mathcal{F}_n$ .

For each  $n \geq 2$  and  $B \in \mathcal{F}_1 \otimes \dots \otimes \mathcal{F}_n$  we put

$$P_n(B) = \int_{\Omega_1} P_1(d\omega_1) \int_{\Omega_2} P(\omega_1; d\omega_2) \int_{\Omega_{n-1}} P(\omega_1, \dots, \omega_{n-2}; d\omega_{n-1}) \times \int_{\Omega_n} I_B(\omega_1, \dots, \omega_n) P(\omega_1, \dots, \omega_{n-1}; d\omega_n). \quad (12)$$

It is easily seen that when  $B = A_1 \times \dots \times A_n$  the right-hand side of (12) is the same as the right-hand side of (9). Moreover, when  $n = 2$  it can be

shown, just as in Theorem 8 of §6, that  $P_2$  is a measure. Consequently it is easily established by induction that  $P_n$  is a measure for all  $n \geq 2$ .

The next step is the same as in Kolmogorov's theorem on the extension of a measure in  $(R^\infty, \mathcal{B}(R^\infty))$  (Theorem 3, §3). Thus for every cylindrical set  $J_n(B) = \{\omega \in \Omega: (\omega_1, \dots, \omega_n) \in B\}$ ,  $B \in \mathcal{F}_1 \otimes \dots \otimes \mathcal{F}_n$ , we define the set function  $P$  by

$$P(J_n(B)) = P_n(B). \quad (13)$$

If we use (12) and the fact that  $P(\omega_1, \dots, \omega_k; \cdot)$  are measures, it is easy to establish that the definition (13) is consistent, in the sense that the value of  $P(J_n(B))$  is independent of the representation of the cylindrical set.

It follows that the set function  $P$  defined in (13) for cylindrical sets, and in an obvious way on the algebra that contains all the cylindrical sets, is a finitely additive measure on this algebra. It remains to verify its countable additivity and apply Carathéodory's theorem.

In Theorem 3 of §3 the corresponding verification was based on the property of  $(R^n, \mathcal{B}(R^n))$  that for every Borel set  $B$  there is a compact set  $A \subseteq B$  whose probability measure is arbitrarily close to the measure of  $B$ . In the present case this part of the proof needs to be modified in the following way.

As in Theorem 3 of §3, let  $\{\hat{B}_n\}_{n \geq 1}$  be a sequence of cylindrical sets

$$\hat{B}_n = \{\omega: (\omega_1, \dots, \omega_n) \in B_n\},$$

that decrease to the empty set  $\emptyset$ , but have

$$\lim_{n \rightarrow \infty} P(\hat{B}_n) > 0. \quad (14)$$

For  $n > 1$ , we have from (12)

$$P(\hat{B}_n) = \int_{\Omega_1} f_n^{(1)}(\omega_1) P_1(d\omega_1),$$

where

$$f_n^{(1)}(\omega_1) = \int_{\Omega_2} P(\omega_1; d\omega_2) \dots \int_{\Omega_n} I_{B_n}(\omega_1, \dots, \omega_n) P(\omega_2, \dots, \omega_{n-1}; d\omega_n).$$

Since  $\hat{B}_{n+1} \subseteq \hat{B}_n$ , we have  $B_{n+1} \subseteq B_n \times \Omega_{n+1}$  and therefore

$$I_{B_{n+1}}(\omega_1, \dots, \omega_{n+1}) \leq I_{B_n}(\omega_1, \dots, \omega_n) I_{\Omega_{n+1}}(\omega_{n+1}).$$

Hence the sequence  $\{f_n^{(1)}(\omega_1)\}_{n \geq 1}$  decreases. Let  $f^{(1)}(\omega_1) = \lim_n f_n^{(1)}(\omega_1)$ . By the dominated convergence theorem

$$\lim_n P(\hat{B}_n) = \lim_n \int_{\Omega_1} f_n^{(1)}(\omega_1) P_1(d\omega_1) = \int_{\Omega_1} f^{(1)}(\omega_1) P_1(d\omega_1).$$

By hypothesis,  $\lim_n P(\hat{B}_n) > 0$ . It follows that there is an  $\omega_1^0 \in B$  such that  $f^{(1)}(\omega_1^0) > 0$ , since if  $\omega_1 \notin B$  then  $f_n^{(1)}(\omega_1) = 0$  for  $n \geq 1$ .

Moreover, for  $n > 2$ ,

$$f_n^{(1)}(\omega_1^0) = \int_{\Omega_2} f_n^{(2)}(\omega_2) P(\omega_1^0; d\omega_2), \quad (15)$$

where

$$\begin{aligned} f_n^{(2)}(\omega_2) &= \int_{\Omega} P(\omega_1^0, \omega_2; d\omega_3) \\ &\quad \cdots \int_{\Omega_n} I_{B_n}(\omega_1^0, \omega_2, \dots, \omega_n) P(\omega_1^0, \omega_2, \dots, \omega_{n-1}, d\omega_n). \end{aligned}$$

We can establish, as for  $\{f_n^{(1)}(\omega_1)\}$ , that  $\{f_n^{(2)}(\omega_2)\}$  is decreasing. Let  $f^{(2)}(\omega_2) = \lim_{n \rightarrow \infty} f_n^{(2)}(\omega_2)$ . Then it follows from (15) that

$$0 < f^{(1)}(\omega_1^0) = \int_{\Omega_2} f^{(2)}(\omega_2) P(\omega_1^0; d\omega_2),$$

and there is a point  $\omega_2^0 \in \Omega_2$  such that  $f^{(2)}(\omega_2^0) > 0$ . Then  $(\omega_1^0, \omega_2^0) \in B_2$ . Continuing this process, we find a point  $(\omega_1^0, \dots, \omega_n^0) \in B_n$  for each  $n$ . Consequently  $(\omega_1^0, \dots, \omega_n^0, \dots) \in \bigcap \hat{B}_n$ , but by hypothesis we have  $\bigcap \hat{B}_n = \emptyset$ . This contradiction shows that  $\lim_n P(\hat{B}_n) = 0$ .

Thus we have proved the part of the theorem about the existence of the probability measure  $P$ . The other part follows from this by putting  $X_n(\omega) = \omega_n, n \geq 1$ .

**Corollary 1.** Let  $(E_n, \mathcal{E}_n)_{n \geq 1}$  be any measurable spaces and  $(P_n)_{n \geq 1}$ , measures on them. Then there are a probability space  $(\Omega, \mathcal{F}, P)$  and a family of independent random elements  $X_1, X_2, \dots$  with values in  $(E_1, \mathcal{E}_1), (E_2, \mathcal{E}_2), \dots$ , respectively, such that

$$P\{\omega: X_n(\omega) \in B\} = P_n(B), \quad B \in \mathcal{E}_n, n \geq 1.$$

**Corollary 2.** Let  $E = \{1, 2, \dots\}$ , and let  $\{p_k(x, y)\}$  be a family of nonnegative functions,  $k \geq 1, x, y \in E$ , such that  $\sum_{y \in E} p_k(x, y) = 1, x \in E, k \geq 1$ . Also let  $\pi = \pi(x)$  be a probability distribution on  $E$  (that is,  $\pi(x) \geq 0, \sum_{x \in E} \pi(x) = 1$ ).

Then there are a probability space  $(\Omega, \mathcal{F}, P)$  and a family  $X = \{\xi_0, \xi_1, \dots\}$  of random variables on it, such that

$$P\{\xi_0 = x_0, \xi_1 = x_1, \dots, \xi_n = x_n\} = \pi(x_0) p_1(x_0, x_1) \cdots p_n(x_{n-1}, x_n) \quad (16)$$

(cf. (I.12.4)) for all  $x_i \in E$  and  $n \geq 1$ . We may take  $\Omega$  to be the space

$$\Omega = \{\omega: \omega = (x_0, x_1, \dots), x_i \in E\}.$$

A sequence  $X = \{\xi_0, \xi_1, \dots\}$  of random variables satisfying (16) is a Markov chain with a countable set  $E$  of states, transition matrix  $\{p_k(x, y)\}$  and initial probability distribution  $\pi$ . (Cf. the definition in §12 of Chapter I.)



## 4. PROBLEMS

1. Let  $\Omega = [0, 1]$ , let  $\mathcal{F}$  be the class of Borel subsets of  $[0, 1]$ , and let  $P$  be Lebesgue measure on  $[0, 1]$ . Show that the space  $(\Omega, \mathcal{F}, P)$  is universal in the following sense. For every distribution function  $F(x)$  on  $(\Omega, \mathcal{F}, P)$  there is a random variable  $\xi = \xi(\omega)$  such that its distribution function  $F_\xi(x) = P(\xi \leq x)$  coincides with  $F(x)$ . (Hint.  $\xi(\omega) = F^{-1}(\omega)$ ,  $0 < \omega < 1$ , where  $F^{-1}(\omega) = \sup\{x: F(x) < \omega\}$ , when  $0 < \omega < 1$ , and  $\xi(0), \xi(1)$  can be chosen arbitrarily.)
2. Verify the consistency of the families of distributions in the corollaries to Theorems 1 and 2.
3. Deduce Corollary 2, Theorem 2, from Theorem 1.

## §10. Various Kinds of Convergence of Sequences of Random Variables

1. Just as in analysis, in probability theory we need to use various kinds of convergence of random variables. Four of these are particularly important: *in probability, with probability one, in mean of order  $p$ , in distribution*.

First some definitions. Let  $\xi, \xi_1, \xi_2, \dots$  be random variables defined on a probability space  $(\Omega, \mathcal{F}, P)$ .

**Definition 1.** The sequence  $\xi_1, \xi_2, \dots$  of random variables converges *in probability* to the random variable  $\xi$  (notation:  $\xi_n \xrightarrow{P} \xi$ ) if for every  $\varepsilon > 0$

$$P\{|\xi_n - \xi| > \varepsilon\} \rightarrow 0, \quad n \rightarrow \infty. \quad (1)$$

We have already encountered this convergence in connection with the law of large numbers for a Bernoulli scheme, which stated that

$$P\left(\left|\frac{S_n}{n} - p\right| > \varepsilon\right) \rightarrow 0, \quad n \rightarrow \infty$$

(see §5 of Chapter I). In analysis this is known as *convergence in measure*.

**Definition 2.** The sequence  $\xi_1, \xi_2, \dots$  of random variables converges *with probability one (almost surely, almost everywhere)* to the random variable  $\xi$  if

$$P\{\omega: \xi_n \not\rightarrow \xi\} = 0, \quad (2)$$

i.e. if the set of sample points  $\omega$  for which  $\xi_n(\omega)$  does not converge to  $\xi$  has probability zero.

This convergence is denoted by  $\xi_n \rightarrow \xi$  (P-a.s.), or  $\xi_n \xrightarrow{a.s.} \xi$  or  $\xi_n \xrightarrow{a.e.} \xi$ .

**Definition 3.** The sequence  $\xi_1, \xi_2, \dots$  of random variables converges in mean of order  $p$ ,  $0 < p < \infty$ , to the random variable  $\xi$  if

$$E|\xi_n - \xi|^p \rightarrow 0, \quad n \rightarrow \infty. \quad (3)$$

In analysis this is known as *convergence in  $L^p$* , and denoted by  $\xi_n \xrightarrow{L^p} \xi$ . In the special case  $p = 2$  it is called *mean square convergence* and denoted by  $\xi = \text{l.i.m. } \xi_n$  (for "limit in the mean").

**Definition 4.** The sequence  $\xi_1, \xi_2, \dots$  of random variables converges in distribution to the random variable  $\xi$  (notation:  $\xi_n \xrightarrow{d} \xi$ ) if

$$Ef(\xi_n) \rightarrow Ef(\xi), \quad n \rightarrow \infty, \quad (4)$$

for every bounded continuous function  $f = f(x)$ . The reason for the terminology is that, according to what will be proved in Chapter III, §1, condition (4) is equivalent to the convergence of the distribution  $F_{\xi_n}(x)$  to  $F_\xi(x)$  at each point  $x$  of continuity of  $F_\xi(x)$ . This convergence is denoted by  $F_{\xi_n} \Rightarrow F_\xi$ .

We emphasize that the convergence of random variables in distribution is defined only in terms of the convergence of their distribution functions. Therefore it makes sense to discuss this mode of convergence even when the random variables are defined on different probability spaces. This convergence will be studied in detail in Chapter III, where, in particular, we shall explain why in the definition of  $F_{\xi_n} \Rightarrow F_\xi$  we require only convergence at points of continuity of  $F_\xi(x)$  and not at all  $x$ .

2. In solving problems of analysis on the convergence (in one sense or another) of a given sequence of functions, it is useful to have the concept of a fundamental sequence (or Cauchy sequence). We can introduce a similar concept for each of the first three kinds of convergence of a sequence of random variables.

Let us say that a sequence  $\{\xi_n\}_{n \geq 1}$  of random variables is *fundamental in probability*, or *with probability 1*, or *in mean of order  $p$* ,  $0 < p < \infty$ , if the corresponding one of the following properties is satisfied:  $P\{|\xi_n - \xi_m| > \varepsilon\} \rightarrow 0$ , as  $m, n \rightarrow \infty$  for every  $\varepsilon > 0$ ; the sequence  $\{\xi_n(\omega)\}_{n \geq 1}$  is fundamental for almost all  $\omega \in \Omega$ ; the sequence  $\{\xi_n(\omega)\}_{n \geq 1}$  is fundamental in  $L^p$ , i.e.  $E|\xi_n - \xi_m|^p \rightarrow 0$  as  $n, m \rightarrow \infty$ .

### 3. Theorem 1.

(a) A necessary and sufficient condition that  $\xi_n \rightarrow \xi$  (P-a.s.) is that

$$P\left\{\sup_{k \geq n} |\xi_k - \xi| \geq \varepsilon\right\} \rightarrow 0, \quad n \rightarrow \infty. \quad (5)$$

for every  $\varepsilon > 0$ .

(b) The sequence  $\{\xi_n\}_{n \geq 1}$  is fundamental with probability 1 if and only if

$$P\left\{\sup_{\substack{k \geq n \\ l \geq n}} |\xi_k - \xi_l| \geq \varepsilon\right\} \rightarrow 0, \quad n \rightarrow \infty, \quad (6)$$

for every  $\varepsilon > 0$ ; or equivalently

$$P\left\{\sup_{k \geq 0} |\xi_{n+k} - \xi_n| \geq \varepsilon\right\} \rightarrow 0, \quad n \rightarrow \infty. \quad (7)$$

PROOF. (a) Let  $A_n^\varepsilon = \{\omega: |\xi_n - \xi| \geq \varepsilon\}$ ,  $A^\varepsilon = \overline{\lim} A_n^\varepsilon \equiv \bigcap_{n=1}^{\infty} \bigcup_{k \geq n} A_k^\varepsilon$ . Then

$$\{\omega: \xi_n \not\rightarrow \xi\} = \bigcup_{\varepsilon > 0} A^\varepsilon = \bigcup_{m=1}^{\infty} A^{1/m}.$$

But

$$P(A^\varepsilon) = \lim_n P\left(\bigcup_{k \geq n} A_k^\varepsilon\right),$$

Hence (a) follows from the following chain of implications:

$$\begin{aligned} 0 = P\{\omega: \xi_n \not\rightarrow \xi\} &= P\left(\bigcup_{\varepsilon > 0} A^\varepsilon\right) \Leftrightarrow P\left(\bigcup_{m=1}^{\infty} A^{1/m}\right) = 0 \\ &\Leftrightarrow P(A^{1/m}) = 0, \quad m \geq 1 \Leftrightarrow P(A^\varepsilon) = 0, \quad \varepsilon > 0, \\ &\Leftrightarrow P\left(\bigcup_{k \geq n} A_k^\varepsilon\right) \rightarrow 0, \quad n \rightarrow \infty \Leftrightarrow P\left(\sup_{k \geq n} |\xi_k - \xi| \geq \varepsilon\right) \rightarrow 0, \\ &\hspace{15em} n \rightarrow \infty. \end{aligned}$$

(b) Let

$$B_{k,l}^\varepsilon = \{\omega: |\xi_k - \xi_l| \geq \varepsilon\}, \quad B^\varepsilon = \bigcap_{n=1}^{\infty} \bigcup_{\substack{k \geq n \\ l \geq n}} B_{k,l}^\varepsilon.$$

Then  $\{\omega: \{\xi_n(\omega)\}_{n \geq 1} \text{ is not fundamental}\} = \bigcup_{\varepsilon > 0} B^\varepsilon$ , and it can be shown as in (a) that  $P\{\omega: \{\xi_n(\omega)\}_{n \geq 1} \text{ is not fundamental}\} = 0 \Leftrightarrow (6)$ . The equivalence of (6) and (7) follows from the obvious inequalities

$$\sup_{k \geq 0} |\xi_{n+k} - \xi_n| \leq \sup_{\substack{k \geq 0 \\ l \geq 0}} |\xi_{n+k} - \xi_{n+l}| \leq 2 \sup_{k \geq 0} |\xi_{n+k} - \xi_n|.$$

This completes the proof of the theorem.

**Corollary.** Since

$$P\left\{\sup_{k \geq n} |\xi_k - \xi| \geq \varepsilon\right\} = P\left\{\bigcup_{k \geq n} (|\xi_k - \xi| \geq \varepsilon)\right\} \leq \sum_{k \geq n} P\{|\xi_k - \xi| \geq \varepsilon\},$$

a sufficient condition for  $\xi_n \xrightarrow{a.s.} \xi$  is that

$$\sum_{k=1}^{\infty} P\{|\xi_k - \xi| \geq \varepsilon\} < \infty \quad (8)$$

is satisfied for every  $\varepsilon > 0$ .

It is appropriate to observe at this point that the reasoning used in obtaining (8) lets us establish the following simple but important result which is essential in studying properties that are satisfied with probability 1.

Let  $A_1, A_2, \dots$  be a sequence of events in  $F$ . Let (see the table in §1)  $\{A_n \text{ i.o.}\}$  denote the event  $\overline{\lim} A_n$  that consists in the realization of infinitely many of  $A_1, A_2, \dots$ .

**Borel–Cantelli Lemma.**

(a) If  $\sum P(A_n) < \infty$  then  $P\{A_n \text{ i.o.}\} = 0$ .

(b) If  $\sum P(A_n) = \infty$  and  $A_1, A_2, \dots$  are independent, then  $P\{A_n \text{ i.o.}\} = 1$ .

PROOF. (a) By definition

$$\{A_n \text{ i.o.}\} = \overline{\lim} A_n = \bigcap_{n=1}^{\infty} \bigcup_{k \geq n} A_k.$$

Consequently

$$P\{A_n \text{ i.o.}\} = P\left\{\bigcap_{n=1}^{\infty} \bigcup_{k \geq n} A_k\right\} = \lim P\left(\bigcup_{k \geq n} A_k\right) \leq \lim \sum_{k \geq n} P(A_k),$$

and (a) follows.

(b) If  $A_1, A_2, \dots$  are independent, so are  $\bar{A}_1, \bar{A}_2, \dots$ . Hence for  $N \geq n$  we have

$$P\left(\bigcap_{k=n}^N A_k\right) = \prod_{k=n}^N P(A_k),$$

and it is then easy to deduce that

$$P\left(\bigcap_{k=n}^{\infty} \bar{A}_k\right) = \prod_{k=n}^{\infty} P(\bar{A}_k). \quad (9)$$

Since  $\log(1 - x) \leq -x$ ,  $0 \leq x < 1$ ,

$$\log \prod_{k=n}^{\infty} [1 - P(A_k)] = \sum_{k=n}^{\infty} \log[1 - P(A_k)] \leq - \sum_{k=n}^{\infty} P(A_k) = -\infty.$$

Consequently

$$P\left(\bigcap_{k=n}^{\infty} \bar{A}_k\right) = 0$$

for all  $n$ , and therefore  $P\{A_n \text{ i.o.}\} = 1$ .

This completes the proof of the lemma.

**Corollary 1.** If  $A_n^\varepsilon = \{\omega: |\xi_n - \xi| \geq \varepsilon\}$  then (8) shows that  $\sum_{n=1}^{\infty} P(A_n) < \infty$ ,  $\varepsilon > 0$ , and then by the Borel-Cantelli lemma we have  $P(A^\varepsilon) = 0$ ,  $\varepsilon > 0$ , where  $A^\varepsilon = \overline{\lim} A_n^\varepsilon$ . Therefore

$$\begin{aligned} \sum P\{|\xi_k - \xi| \geq \varepsilon\} < \infty, \varepsilon > 0 &\Rightarrow P(A^\varepsilon) = 0, \varepsilon > 0 \\ &\Rightarrow P\{\omega: \xi_n \not\rightarrow \xi\} = 0, \end{aligned}$$

as we already observed above.

**Corollary 2.** Let  $(\varepsilon_n)_{n \geq 1}$  be a sequence of positive numbers such that  $\varepsilon_n \downarrow 0$ ,  $n \rightarrow \infty$ . If

$$\sum_{n=1}^{\infty} P\{|\xi_n - \xi| \geq \varepsilon_n\} < \infty, \quad (10)$$

then  $\xi_n \xrightarrow{a.s.} \xi$ .

In fact, let  $A_n = \{|\xi_n - \xi| \geq \varepsilon_n\}$ . Then  $P(A_n \text{ i.o.}) = 0$  by the Borel-Cantelli lemma. This means that, for almost every  $\omega \in \Omega$ , there is an  $N = N(\omega)$  such that  $|\xi_n(\omega) - \xi(\omega)| \leq \varepsilon_n$  for  $n \geq N(\omega)$ . But  $\varepsilon_n \downarrow 0$ , and therefore  $\xi_n(\omega) \rightarrow \xi(\omega)$  for almost every  $\omega \in \Omega$ .

**4. Theorem 2.** We have the following implications:

$$\xi_n \xrightarrow{a.s.} \xi \Rightarrow \xi_n \xrightarrow{P} \xi, \quad (11)$$

$$\xi_n \xrightarrow{L^p} \xi \Rightarrow \xi_n \xrightarrow{P} \xi, \quad p > 0, \quad (12)$$

$$\xi_n \xrightarrow{P} \xi \Rightarrow \xi_n \xrightarrow{d} \xi. \quad (13)$$

**PROOF.** Statement (11) follows from comparing the definition of convergence in probability with (5), and (12) follows from Chebyshev's inequality.

To prove (13), let  $f(x)$  be a continuous function, let  $|f(x)| \leq c$ , let  $\varepsilon > 0$ , and let  $N$  be such that  $P(|\xi| > N) \leq \varepsilon/4c$ . Take  $\delta$  so that  $|f(x) - f(y)| \leq \varepsilon/2c$  for  $|x| < N$  and  $|x - y| \leq \delta$ . Then (cf. the proof of Weierstrass's theorem in Subsection 5, §5, Chapter I)

$$\begin{aligned} E|f(\xi_n) - f(\xi)| &= E(|f(\xi_n) - f(\xi)|; |\xi_n - \xi| \leq \delta, |\xi| \leq N) \\ &\quad + E(|f(\xi_n) - f(\xi)|; |\xi_n - \xi| \leq \delta, |\xi| > N) \\ &\quad + E(|f(\xi_n) - f(\xi)|; |\xi_n - \xi| > \delta) \\ &\leq \varepsilon/2 + \varepsilon/2 + 2cP\{|\xi_n - \xi| > \delta\} \\ &= \varepsilon + 2cP\{|\xi_n - \xi| > \delta\}. \end{aligned}$$

But  $P\{|\xi_n - \xi| > \delta\} \rightarrow 0$ , and hence  $E|f(\xi_n) - f(\xi)| \leq 2\varepsilon$  for sufficiently large  $n$ ; since  $\varepsilon > 0$  is arbitrary, this establishes (13).

This completes the proof of the theorem.

We now present a number of examples which show, in particular, that the converses of (11) and (12) are false in general.

EXAMPLE 1 ( $\xi_n \xrightarrow{P} \xi \not\xrightarrow{a.s.} \xi$ ;  $\xi_n \xrightarrow{L^p} \xi \not\xrightarrow{a.s.} \xi$ ). Let  $\Omega = [0, 1]$ ,  $\mathcal{F} = \mathcal{B}([0, 1])$ ,  $P =$  Lebesgue measure. Put

$$A_n^i = \left[ \frac{i-1}{n}, \frac{i}{n} \right], \quad \xi_n^i = I_{A_n^i}(\omega), \quad i = 1, 2, \dots, n; n \geq 1.$$

Then the sequence

$$\{\xi_1^1; \xi_2^1, \xi_2^2; \xi_3^1, \xi_3^2, \xi_3^3; \dots\}$$

of random variables converges both in probability and in mean of order  $p > 0$ , but does not converge at any point  $\omega \in [0, 1]$ .

EXAMPLE 2 ( $\xi_n \xrightarrow{a.s.} \xi \Rightarrow \xi_n \xrightarrow{P} \xi \not\xrightarrow{L^p} \xi$ ,  $p > 0$ ). Again let  $\Omega = [0, 1]$ ,  $\mathcal{F} = \mathcal{B}([0, 1])$ ,  $P =$  Lebesgue measure, and let

$$\xi_n(\omega) = \begin{cases} e^n, & 0 \leq \omega \leq 1/n, \\ 0, & \omega > 1/n. \end{cases}$$

Then  $\{\xi_n\}$  converges with probability 1 (and therefore in probability) to zero, but

$$E|\xi_n|^p = \frac{e^{np}}{n} \rightarrow \infty, \quad n \rightarrow \infty,$$

for every  $p > 0$ .

EXAMPLE 3 ( $\xi_n \xrightarrow{L^p} \xi \not\xrightarrow{a.s.} \xi$ ). Let  $\{\xi_n\}$  be a sequence of independent random variables with

$$P(\xi_n = 1) = p_n, \quad P(\xi_n = 0) = 1 - p_n.$$

Then it is easy to show that

$$\xi_n \xrightarrow{P} 0 \Leftrightarrow p_n \rightarrow 0, \quad n \rightarrow \infty, \quad (14)$$

$$\xi_n \xrightarrow{L^p} 0 \Leftrightarrow p_n \rightarrow 0, \quad n \rightarrow \infty, \quad (15)$$

$$\xi_n \xrightarrow{a.s.} 0 \Rightarrow \sum_{n=1}^{\infty} p_n < \infty. \quad (16)$$

In particular, if  $p_n = 1/n$  then  $\xi_n \xrightarrow{L^p} 0$  for every  $p > 0$ , but  $\xi_n \not\xrightarrow{a.s.} 0$ .

The following theorem singles out an interesting case when almost sure convergence implies convergence in  $L^1$ .

**Theorem 3.** Let  $(\xi_n)$  be a sequence of nonnegative random variables such that  $\xi_n \xrightarrow{a.s.} \xi$  and  $E\xi_n \rightarrow E\xi < \infty$ . Then

$$E|\xi_n - \xi| \rightarrow 0, \quad n \rightarrow \infty. \quad (17)$$

PROOF. We have  $E\xi_n < \infty$  for sufficiently large  $n$ , and therefore for such  $n$  we have

$$\begin{aligned} E|\xi - \xi_n| &= E(\xi - \xi_n)I_{(\xi \geq \xi_n)} + E(\xi_n - \xi)I_{(\xi_n > \xi)} \\ &= 2E(\xi - \xi_n)I_{(\xi \geq \xi_n)} + E(\xi_n - \xi). \end{aligned}$$

But  $0 \leq (\xi - \xi_n)I_{(\xi \geq \xi_n)} \leq \xi$ . Therefore, by the dominated convergence theorem,  $\lim_n E(\xi - \xi_n)I_{(\xi \geq \xi_n)} = 0$ , which together with  $E\xi_n \rightarrow E\xi$  proves (17).

**Remark.** The dominated convergence theorem also holds when almost sure convergence is replaced by convergence in probability (see Problem 1). Hence in Theorem 3 we may replace " $\xi_n \xrightarrow{a.s.} \xi$ " by " $\xi_n \xrightarrow{P} \xi$ ."

5. It is shown in analysis that every fundamental sequence  $(x_n)$ ,  $x_n \in R$ , is convergent (Cauchy criterion). Let us give a similar result for the convergence of a sequence of random variables.

**Theorem 4** (Cauchy Criterion for Almost Sure Convergence). *A necessary and sufficient condition for the sequence  $(\xi_n)_{n \geq 1}$  of random variables to converge with probability 1 (to a random variable  $\xi$ ) is that it is fundamental with probability 1.*

PROOF. If  $\xi_n \xrightarrow{a.s.} \xi$  then

$$\sup_{\substack{k \geq n \\ l \geq n}} |\xi_k - \xi_l| \leq \sup_{k \geq n} |\xi_k - \xi| + \sup_{l \geq n} |\xi_l - \xi|,$$

whence the necessity follows.

Now let  $(\xi_n)_{n \geq 1}$  be fundamental with probability 1. Let  $\mathcal{N} = \{\omega: (\xi_n(\omega)) \text{ is not fundamental}\}$ . Then whenever  $\omega \in \Omega \setminus \mathcal{N}$  the sequence of numbers  $(\xi_n(\omega))_{n \geq 1}$  is fundamental and, by Cauchy's criterion for sequences of numbers,  $\lim \xi_n(\omega)$  exists. Let

$$\xi(\omega) = \begin{cases} \lim \xi_n(\omega), & \omega \in \Omega \setminus \mathcal{N}, \\ 0, & \omega \in \mathcal{N}. \end{cases} \quad (18)$$

The function so defined is a random variable, and evidently  $\xi_n \xrightarrow{a.s.} \xi$ .

This completes the proof.

Before considering the case of convergence in probability, let us establish the following useful result.

**Theorem 5.** *If the sequence  $(\xi_n)$  is fundamental (or convergent) in probability, it contains a subsequence  $(\xi_{n_k})$  that is fundamental (or convergent) with probability 1.*

PROOF. Let  $(\xi_n)$  be fundamental in probability. By Theorem 4, it is enough to show that it contains a subsequence that converges almost surely.

Take  $n_1 = 1$  and define  $n_k$  inductively as the smallest  $n > n_{k-1}$  for which

$$P\{|\xi_t - \xi_s| > 2^{-k}\} < 2^{-k}.$$

for all  $s \geq n, t \geq n$ . Then

$$\sum_k P\{|\xi_{n_{k+1}} - \xi_{n_k}| > 2^{-k}\} < \sum 2^{-k} < \infty$$

and by the Borel-Cantelli lemma

$$P\{|\xi_{n_{k+1}} - \xi_{n_k}| > 2^{-k} \text{ i.o.}\} = 0.$$

Hence

$$\sum_{k=1}^{\infty} |\xi_{n_{k+1}} - \xi_{n_k}| < \infty$$

with probability 1.

Let  $\mathcal{N} = \{\omega: \sum |\xi_{n_{k+1}} - \xi_{n_k}| = \infty\}$ . Then if we put

$$\xi(\omega) = \begin{cases} \xi_{n_1}(\omega) + \sum_{k=1}^{\infty} (\xi_{n_{k+1}}^{(\omega)} - \xi_{n_k}(\omega)), & \omega \in \Omega \setminus \mathcal{N}, \\ 0, & \omega \in \mathcal{N}, \end{cases}$$

we obtain  $\xi_{n_k} \xrightarrow{\text{a.s.}} \xi$ .

If the original sequence converges in probability, then it is fundamental in probability (see also (19)), and consequently this case reduces to the one already considered.

This completes the proof of the theorem.

**Theorem 6** (Cauchy Criterion for Convergence in Probability). *A necessary and sufficient condition for a sequence  $(\xi_n)_{n \geq 1}$  of random variables to converge in probability is that it is fundamental in probability.*

PROOF. If  $\xi_n \xrightarrow{P} \xi$  then

$$P\{|\xi_n - \xi_m| \geq \varepsilon\} \leq P\{|\xi_n - \xi| \geq \varepsilon/2\} + P\{|\xi_m - \xi| \geq \varepsilon/2\} \quad (19)$$

and consequently  $(\xi_n)$  is fundamental in probability.

Conversely, if  $(\xi_n)$  is fundamental in probability, by Theorem 5 there are a subsequence  $(\xi_{n_k})$  and a random variable  $\xi$  such that  $\xi_{n_k} \xrightarrow{\text{a.s.}} \xi$ . But then

$$P\{|\xi_n - \xi| \geq \varepsilon\} \leq P\{|\xi_n - \xi_{n_k}| \geq \varepsilon/2\} + P\{|\xi_{n_k} - \xi| \geq \varepsilon/2\},$$

from which it is clear that  $\xi_n \xrightarrow{P} \xi$ . This completes the proof.

Before discussing convergence in mean of order  $p$ , we make some observations about  $L^p$  spaces.



We denote by  $L^p = L^p(\Omega, \mathcal{F}, P)$  the space of random variables  $\xi = \xi(\omega)$  with  $E|\xi|^p \equiv \int_{\Omega} |\xi|^p dP < \infty$ . Suppose that  $p \geq 1$  and put

$$\|\xi\|_p = (E|\xi|^p)^{1/p}.$$

It is clear that

$$\|\xi\|_p \geq 0, \quad (20)$$

$$\|c\xi\|_p = |c| \|\xi\|_p, \quad c \text{ constant}, \quad (21)$$

and by Minkowski's inequality (6.31)

$$\|\xi + \eta\|_p \leq \|\xi\|_p + \|\eta\|_p. \quad (22)$$

Hence, in accordance with the usual terminology of functional analysis, the function  $\|\cdot\|_p$ , defined on  $L^p$  and satisfying (20)–(22), is (for  $p \geq 1$ ) a *semi-norm*.

For it to be a *norm*, it must also satisfy

$$\|\xi\|_p = 0 \Rightarrow \xi = 0. \quad (23)$$

This property is, of course, not satisfied, since according to Property H (§6) we can only say that  $\xi = 0$  almost surely.

This fact leads to a somewhat different view of the space  $L^p$ . That is, we connect with every random variable  $\xi \in L^p$  the class  $[\xi]$  of random variables in  $L^p$  that are equivalent to it ( $\xi$  and  $\eta$  are *equivalent* if  $\xi = \eta$  almost surely). It is easily verified that the property of equivalence is *reflexive*, *symmetric*, and *transitive*, and consequently the linear space  $L^p$  can be divided into disjoint equivalence classes of random variables. If we now think of  $[L^p]$  as the collection of the classes  $[\xi]$  of equivalent random variables  $\xi \in L^p$ , and define

$$[\xi] + [\eta] = [\xi + \eta].$$

$$a[\xi] = [a\xi], \quad \text{where } a \text{ is a constant,}$$

$$\|[\xi]\|_p = \|\xi\|_p,$$

then  $[L^p]$  becomes a normed linear space.

In functional analysis, we ordinarily describe elements of a space  $[L^p]$ , not as *equivalence classes of functions*, but simply as *functions*. In the same way we do not actually use the notation  $[L^p]$ . From now on, we no longer think about sets of equivalence classes of functions, but simply about elements, functions, random variables, and so on.

It is a basic result of functional analysis that the spaces  $L^p$ ,  $p \geq 1$ , are *complete*, i.e. that every fundamental sequence has a limit. Let us state and prove this in probabilistic language.

**Theorem 7** (Cauchy Test for Convergence in Mean  $p$ th Power). *A necessary and sufficient condition that a sequence  $(\xi_n)_{n \geq 1}$  of random variables in  $L^p$*

converges in mean of order  $p$  to a random variable in  $L^p$  is that the sequence is fundamental in mean of order  $p$ .

**PROOF.** The necessity follows from Minkowski's inequality. Let  $(\xi_n)$  be fundamental ( $\|\xi_n - \xi_m\|_p \rightarrow 0, n, m \rightarrow \infty$ ). As in the proof of Theorem 5, we select a subsequence  $(\xi_{n_k})$  such that  $\xi_{n_k} \xrightarrow{a.s.} \xi$ , where  $\xi$  is a random variable with  $\|\xi\|_p < \infty$ .

Let  $n_1 = 1$  and define  $n_k$  inductively as the smallest  $n > n_{k-1}$  for which

$$\|\xi_t - \xi_s\|_p < 2^{-2k}$$

for all  $s \geq n, t \geq n$ . Let

$$A_k = \{\omega: |\xi_{n_{k+1}} - \xi_{n_k}| \geq 2^{-k}\}.$$

Then by Chebyshev's inequality

$$P(A_k) \leq \frac{E|\xi_{n_{k+1}} - \xi_{n_k}|^p}{2^{-kr}} \leq \frac{2^{-2kr}}{2^{-kr}} = 2^{-kr} \leq 2^{-k}.$$

As in Theorem 5, we deduce that there is a random variable  $\xi$  such that  $\xi_{n_k} \xrightarrow{a.s.} \xi$ .

We now deduce that  $\|\xi_n - \xi\|_p \rightarrow 0$  as  $n \rightarrow \infty$ . To do this, we fix  $\varepsilon > 0$  and choose  $N = N(\varepsilon)$  so that  $\|\xi_n - \xi_m\|_p^p < \varepsilon$  for all  $n \geq N, m \geq N$ . Then for any fixed  $n \geq N$ , by Fatou's lemma,

$$\begin{aligned} E|\xi_n - \xi|^p &= E\left\{\lim_{n_k \rightarrow \infty} |\xi_n - \xi_{n_k}|^p\right\} = E\left\{\liminf_{n_k \rightarrow \infty} |\xi_n - \xi_{n_k}|^p\right\} \\ &\leq \liminf_{n_k \rightarrow \infty} E|\xi_n - \xi_{n_k}|^p = \liminf_{n_k \rightarrow \infty} \|\xi_n - \xi_{n_k}\|_p^p \leq \varepsilon. \end{aligned}$$

Consequently  $E|\xi_n - \xi|^p \rightarrow 0, n \rightarrow \infty$ . It is also clear that since  $\xi = (\xi - \xi_n) + \xi_n$  we have  $E|\xi|^p < \infty$  by Minkowski's inequality.

This completes the proof of the theorem.

**Remark 1.** In the terminology of functional analysis a complete normed linear space is called a *Banach space*. Thus  $L^p, p \geq 1$ , is a Banach space.

**Remark 2.** If  $0 < p < 1$ , the function  $\|\xi\|_p = (E|\xi|^p)^{1/p}$  does not satisfy the triangle inequality (22) and consequently is not a norm. Nevertheless the space (of equivalence classes)  $L^p, 0 < p < 1$ , is complete in the metric  $d(\xi, \eta) \equiv E|\xi - \eta|^p$ .

**Remark 3.** Let  $L^\infty = L^\infty(\Omega, \mathcal{F}, P)$  be the space (of equivalence classes of) random variables  $\xi = \xi(\omega)$  for which  $\|\xi\|_\infty < \infty$ , where  $\|\xi\|_\infty$ , the *essential supremum* of  $\xi$ , is defined by

$$\|\xi\|_\infty \equiv \text{ess sup} |\xi| \equiv \inf\{0 \leq c < \infty: P(|\xi| > c) = 0\}.$$

The function  $\|\cdot\|_\infty$  is a norm, and  $L^\infty$  is complete in this norm.

## 6. PROBLEMS

1. Use Theorem 5 to show that almost sure convergence can be replaced by convergence in probability in Theorems 3 and 4 of §6.
2. Prove that  $L^\infty$  is complete.
3. Show that if  $\xi_n \xrightarrow{P} \xi$  and also  $\xi_n \xrightarrow{P} \eta$  then  $\xi$  and  $\eta$  are equivalent ( $P(\xi \neq \eta) = 0$ ).
4. Let  $\xi_n \xrightarrow{P} \xi$ ,  $\eta_n \xrightarrow{P} \eta$ , and let  $\xi$  and  $\eta$  be equivalent. Show that

$$P\{|\xi_n - \eta_n| \geq \varepsilon\} \rightarrow 0, \quad n \rightarrow \infty,$$

for every  $\varepsilon > 0$ .

5. Let  $\xi_n \xrightarrow{P} \xi$ ,  $\eta_n \xrightarrow{P} \eta$ . Show that  $a\xi_n + b\eta_n \xrightarrow{P} a\xi + b\eta$  ( $a, b$  constants),  $|\xi_n| \xrightarrow{P} |\xi|$ ,  $\xi_n \eta_n \xrightarrow{P} \xi \eta$ .
6. Let  $(\xi_n - \xi)^2 \rightarrow 0$ . Show that  $\xi_n^2 \rightarrow \xi^2$ .
7. Show that if  $\xi_n \xrightarrow{d} C$ , where  $C$  is a constant, then this sequence converges in probability:

$$\xi_n \xrightarrow{d} C \Rightarrow \xi_n \xrightarrow{P} C.$$

8. Let  $(\xi_n)_{n \geq 1}$  have the property that  $\sum_{n=1}^{\infty} E|\xi_n|^p < \infty$  for some  $p > 0$ . Show that  $\xi_n \rightarrow 0$  (P-a.s.).
9. Let  $(\xi_n)_{n \geq 1}$  be a sequence of independent identically distributed random variables. Show that

$$\begin{aligned} E|\xi_1| < \infty &\Leftrightarrow \sum_{n=1}^{\infty} P\{|\xi_1| > \varepsilon \cdot n\} < \infty \\ &\Leftrightarrow \sum_{n=1}^{\infty} P\left\{\left|\frac{\xi_n}{n}\right| > \varepsilon\right\} < \infty \Rightarrow \frac{\xi_n}{n} \rightarrow 0 \quad (\text{P-a.s.}). \end{aligned}$$

10. Let  $(\xi_n)_{n \geq 1}$  be a sequence of random variables. Suppose that there are a random variable  $\xi$  and a sequence  $\{n_k\}$  such that  $\xi_{n_k} \rightarrow \xi$  (P-a.s.) and  $\max_{n_{k-1} < l \leq n_k} |\xi_l - \xi_{n_{k-1}}| \rightarrow 0$  (P-a.s.) as  $k \rightarrow \infty$ . Show that then  $\xi_n \rightarrow \xi$  (P-a.s.).
11. Let the  $d$ -metric on the set of random variables be defined by

$$d(\xi, \eta) = E \frac{|\xi - \eta|}{1 + |\xi - \eta|}$$

and identify random variables that coincide almost surely. Show that convergence in probability is equivalent to convergence in the  $d$ -metric.

12. Show that there is no metric on the set of random variables such that convergence in that metric is equivalent to almost sure convergence.

## §11. The Hilbert Space of Random Variables with Finite Second Moment

1. An important role among the Banach spaces  $L^p$ ,  $p \geq 1$ , is played by the space  $L^2 = L^2(\Omega, \mathcal{F}, P)$ , the space of (equivalence classes of) random variables with finite second moments.

If  $\xi$  and  $\eta \in L^2$ , we put

$$(\xi, \eta) \equiv E\xi\eta. \quad (1)$$

It is clear that if  $\xi, \eta, \zeta \in L^2$  then

$$\begin{aligned} (a\xi + b\eta, \zeta) &= a(\xi, \zeta) + b(\eta, \zeta), & a, b \in R, \\ (\xi, \xi) &\geq 0 \end{aligned}$$

and

$$(\xi, \xi) = 0 \Leftrightarrow \xi = 0.$$

Consequently  $(\xi, \eta)$  is a *scalar product*. The space  $L^2$  is *complete* with respect to the norm

$$\|\xi\| = (\xi, \xi)^{1/2} \quad (2)$$

induced by this scalar product (as was shown in §10). In accordance with the terminology of functional analysis, a space with the scalar product (1) is a *Hilbert space*.

Hilbert space methods are extensively used in probability theory to study properties that depend only on the first two moments of random variables ("L<sup>2</sup>-theory"). Here we shall introduce the basic concepts and facts that will be needed for an exposition of L<sup>2</sup>-theory (Chapter VI).

**2.** Two random variables  $\xi$  and  $\eta$  in  $L^2$  are said to be *orthogonal* ( $\xi \perp \eta$ ) if  $(\xi, \eta) \equiv E\xi\eta = 0$ . According to §8,  $\xi$  and  $\eta$  are *uncorrelated* if  $\text{cov}(\xi, \eta) = 0$ , i.e. if

$$E\xi\eta = E\xi E\eta.$$

It follows that the properties of being orthogonal and of being uncorrelated coincide for random variables with zero mean values.

A set  $M \subseteq L^2$  is a *system of orthogonal random variables* if  $\xi \perp \eta$  for every  $\xi, \eta \in M$  ( $\xi \neq \eta$ ).

If also  $\|\xi\| = 1$  for every  $\xi \in M$ , then  $M$  is an *orthonormal system*.

**3.** Let  $M = \{\eta_1, \dots, \eta_n\}$  be an orthonormal system and  $\xi$  any random variable in  $L^2$ . Let us find, in the class of linear estimators  $\sum_{i=1}^n a_i \eta_i$ , the best mean-square estimator for  $\xi$  (cf. Subsection 2, §8).

A simple computation shows that

$$\begin{aligned} E \left| \xi - \sum_{i=1}^n a_i \eta_i \right|^2 &\equiv \left\| \xi - \sum_{i=1}^n a_i \eta_i \right\|^2 = \left( \xi - \sum_{i=1}^n a_i \eta_i, \xi - \sum_{i=1}^n a_i \eta_i \right) \\ &= \|\xi\|^2 - 2 \sum_{i=1}^n a_i (\xi, \eta_i) + \left( \sum_{i=1}^n a_i \eta_i, \sum_{i=1}^n a_i \eta_i \right) \\ &= \|\xi\|^2 - 2 \sum_{i=1}^n a_i (\xi, \eta_i) + \sum_{i=1}^n a_i^2 \end{aligned}$$

$$\begin{aligned}
&= \|\xi\|^2 - \sum_{i=1}^n |(\xi, \eta_i)|^2 + \sum_{i=1}^n |a_i - (\xi, \eta_i)|^2 \\
&\geq \|\xi\|^2 - \sum_{i=1}^n |(\xi, \eta_i)|^2,
\end{aligned} \tag{3}$$

where we used the equation

$$a_i^2 - 2a_i(\xi, \eta_i) = |a_i - (\xi, \eta_i)|^2 - |(\xi, \eta_i)|^2.$$

It is now clear that the infimum of  $E|\xi - \sum_{i=1}^n a_i \eta_i|^2$  over all real  $a_1, \dots, a_n$  is attained for  $a_i = (\xi, \eta_i)$ ,  $i = 1, \dots, n$ .

Consequently the best (in the mean-square sense) estimator for  $\xi$  in terms of  $\eta_1, \dots, \eta_n$  is

$$\hat{\xi} = \sum_{i=1}^n (\xi, \eta_i) \eta_i. \tag{4}$$

Here

$$\Delta \equiv \inf E \left| \xi - \sum_{i=1}^n a_i \eta_i \right|^2 = E |\xi - \hat{\xi}|^2 = \|\xi\|^2 - \sum_{i=1}^n |(\xi, \eta_i)|^2 \tag{5}$$

(compare (I.4.17) and (8.13)).

Inequality (3) also implies *Bessel's inequality*: if  $M = \{\eta_1, \eta_2, \dots\}$  is an orthonormal system and  $\xi \in L^2$ , then

$$\sum_{i=1}^{\infty} |(\xi, \eta_i)|^2 \leq \|\xi\|^2; \tag{6}$$

and equality is attained if and only if

$$\xi = \text{l.i.m.}_n \sum_{i=1}^n (\xi, \eta_i) \eta_i. \tag{7}$$

The *best linear estimator* of  $\xi$  is often denoted by  $\hat{E}(\xi|\eta_1, \dots, \eta_n)$  and called the *conditional expectation* (of  $\xi$  with respect to  $\eta_1, \dots, \eta_n$ ) in the *wide sense*.

The reason for the terminology is as follows. If we consider all estimators  $\varphi = \varphi(\eta_1, \dots, \eta_n)$  of  $\xi$  in terms of  $\eta_1, \dots, \eta_n$  (where  $\varphi$  is a Borel function), the best estimator will be  $\varphi^* = E(\xi|\eta_1, \dots, \eta_n)$ , i.e. the conditional expectation of  $\xi$  with respect to  $\eta_1, \dots, \eta_n$  (cf. Theorem 1, §8). Hence the best linear estimator is, by analogy, denoted by  $\hat{E}(\xi|\eta_1, \dots, \eta_n)$  and called the conditional expectation in the wide sense. We note that if  $\eta_1, \dots, \eta_n$  form a Gaussian system (see §13 below), then  $E(\xi|\eta_1, \dots, \eta_n)$  and  $\hat{E}(\xi|\eta_1, \dots, \eta_n)$  are the same.

Let us discuss the *geometric meaning* of  $\hat{\xi} = \hat{E}(\xi|\eta_1, \dots, \eta_n)$ .

Let  $\mathcal{L} = \mathcal{L}\{\eta_1, \dots, \eta_n\}$  denote the *linear manifold* spanned by the orthonormal system of random variables  $\eta_1, \dots, \eta_n$  (i.e., the set of random variables of the form  $\sum_{i=1}^n a_i \eta_i$ ,  $a_i \in R$ ).

Then it follows from the preceding discussion that  $\xi$  admits the “orthogonal decomposition”

$$\xi = \hat{\xi} + (\xi - \hat{\xi}), \quad (8)$$

where  $\hat{\xi} \in \mathcal{L}$  and  $\xi - \hat{\xi} \perp \mathcal{L}$  in the sense that  $\xi - \hat{\xi} \perp \lambda$  for every  $\lambda \in \mathcal{L}$ . It is natural to call  $\hat{\xi}$  the *projection* of  $\xi$  on  $\mathcal{L}$  (the element of  $\mathcal{L}$  “closest” to  $\xi$ ), and to say that  $\xi - \hat{\xi}$  is *perpendicular* to  $\mathcal{L}$ .

4. The concept of orthonormality of the random variables  $\eta_1, \dots, \eta_n$  makes it easy to find the best linear estimator (the projection)  $\hat{\xi}$  of  $\xi$  in terms of  $\eta_1, \dots, \eta_n$ . The situation becomes complicated if we give up the hypothesis of orthonormality. However, the case of arbitrary  $\eta_1, \dots, \eta_n$  can in a certain sense be reduced to the case of orthonormal random variables, as will be shown below. We shall suppose for the sake of simplicity that all our random variables have zero mean values.

We shall say that the random variables  $\eta_1, \dots, \eta_n$  are *linearly independent* if the equation

$$\sum_{i=1}^n a_i \eta_i = 0 \quad (\text{P-a.s.})$$

is satisfied only when all  $a_i$  are zero.

Consider the covariance matrix

$$\mathbb{R} = E\eta\eta^T$$

of the vector  $\eta = (\eta_1, \dots, \eta_n)$ . It is symmetric and nonnegative definite, and as noticed in §8, can be diagonalized by an orthogonal matrix  $\mathcal{O}$ :

$$\mathcal{O}^T \mathbb{R} \mathcal{O} = D,$$

where

$$D = \begin{pmatrix} d_1 & & 0 \\ & \ddots & \\ 0 & & d_n \end{pmatrix}$$

has nonnegative elements  $d_i$ , the eigenvalues of  $\mathbb{R}$ , i.e. the zeros  $\lambda$  of the characteristic equation  $\det(\mathbb{R} - \lambda E) = 0$ .

If  $\eta_1, \dots, \eta_n$  are linearly independent, the Gram determinant ( $\det \mathbb{R}$ ) is not zero and therefore  $d_i > 0$ . Let

$$B = \begin{pmatrix} \sqrt{d_1} & & 0 \\ & \ddots & \\ 0 & & \sqrt{d_n} \end{pmatrix}$$

and

$$\beta = B^{-1} \mathcal{O}^T \eta. \quad (9)$$

Then the covariance matrix of  $\beta$  is

$$E\beta\beta^T = B^{-1} \mathcal{O}^T E\eta\eta^T \mathcal{O} B^{-1} = B^{-1} \mathcal{O}^T \mathbb{R} \mathcal{O} B^{-1} = E,$$

and therefore  $\beta = (\beta_1, \dots, \beta_n)$  consists of uncorrelated random variables.

It is also clear that

$$\eta = (\mathcal{O}B)\beta. \quad (10)$$

Consequently if  $\eta_1, \dots, \eta_n$  are linearly independent there is an orthonormal system such that (9) and (10) hold. Here

$$\mathcal{L}\{\eta_1, \dots, \eta_n\} = \mathcal{L}\{\beta_1, \dots, \beta_n\}.$$

This method of constructing an orthonormal system  $\beta_1, \dots, \beta_n$  is frequently inconvenient. The reason is that if we think of  $\eta_i$  as the value of the random sequence  $(\eta_1, \dots, \eta_n)$  at the instant  $i$ , the value  $\beta_i$  constructed above depends not only on the "past,"  $(\eta_1, \dots, \eta_i)$ , but also on the "future,"  $(\eta_{i+1}, \dots, \eta_n)$ . The *Gram-Schmidt orthogonalization process*, described below, does not have this defect, and moreover has the advantage that it can be applied to an infinite sequence of *linearly independent* random variables (i.e. to a sequence in which every finite set of the variables are linearly independent).

Let  $\eta_1, \eta_2, \dots$  be a sequence of linearly independent random variables in  $L^2$ . We construct a sequence  $\varepsilon_1, \varepsilon_2, \dots$  as follows. Let  $\varepsilon_1 = \eta_1 / \|\eta_1\|$ . If  $\varepsilon_1, \dots, \varepsilon_{n-1}$  have been selected so that they are orthonormal, then

$$\varepsilon_n = \frac{\eta_n - \hat{\eta}_n}{\|\eta_n - \hat{\eta}_n\|}, \quad (11)$$

where  $\hat{\eta}_n$  is the projection of  $\eta_n$  on the linear manifold  $\mathcal{L}(\varepsilon_1, \dots, \varepsilon_{n-1})$  generated by

$$\hat{\eta}_n = \sum_{k=1}^{n-1} (\eta_n, \varepsilon_k) \varepsilon_k. \quad (12)$$

Since  $\eta_1, \dots, \eta_n$  are linearly independent and  $\mathcal{L}\{\eta_1, \dots, \eta_{n-1}\} = \mathcal{L}\{\varepsilon_1, \dots, \varepsilon_{n-1}\}$ , we have  $\|\eta_n - \hat{\eta}_n\| > 0$  and consequently  $\varepsilon_n$  is well defined.

By construction,  $\|\varepsilon_n\| = 1$  for  $n \geq 1$ , and it is clear that  $(\varepsilon_n, \varepsilon_k) = 0$  for  $k < n$ . Hence the sequence  $\varepsilon_1, \varepsilon_2, \dots$  is orthonormal. Moreover, by (11),

$$\eta_n = \hat{\eta}_n + b_n \varepsilon_n,$$

where  $b_n = \|\eta_n - \hat{\eta}_n\|$  and  $\hat{\eta}_n$  is defined by (12).

Now let  $\eta_1, \dots, \eta_n$  be any set of random variables (not necessarily linearly independent). Let  $\det \mathbb{R} = 0$ , where  $\mathbb{R} \equiv \|r_{ij}\|$  is the covariance matrix of  $(\eta_1, \dots, \eta_n)$ , and let

$$\text{rank } \mathbb{R} = r < n.$$

Then, from linear algebra, the quadratic form

$$Q(a) = \sum_{i,j=1}^n r_{ij} a_i a_j, \quad a = (a_1, \dots, a_n),$$

has the property that there are  $n - r$  linearly independent vectors  $a^{(1)}, \dots, a^{(n-r)}$  such that  $Q(a^{(i)}) = 0$ ,  $i = 1, \dots, n - r$ .

But

$$Q(a) = E \left( \sum_{k=1}^n a_k \eta_k \right)^2.$$

Consequently

$$\sum_{k=1}^n a_k^{(i)} \eta_k = 0, \quad i = 1, \dots, n - r,$$

with probability 1.

In other words, there are  $n - r$  linear relations among the variables  $\eta_1, \dots, \eta_n$ . Therefore if, for example,  $\eta_1, \dots, \eta_r$  are linearly independent, the other variables  $\eta_{r+1}, \dots, \eta_n$  can be expressed linearly in terms of them, and consequently  $\mathcal{L}\{\eta_1, \dots, \eta_n\} = \mathcal{L}\{\varepsilon_1, \dots, \varepsilon_r\}$ . Hence it is clear that we can find  $r$  orthonormal random variables  $\varepsilon_1, \dots, \varepsilon_r$  such that  $\eta_1, \dots, \eta_n$  can be expressed linearly in terms of them and  $\mathcal{L}\{\eta_1, \dots, \eta_n\} = \mathcal{L}\{\varepsilon_1, \dots, \varepsilon_r\}$ .

5. Let  $\eta_1, \eta_2, \dots$  be a sequence of random variables in  $L^2$ . Let  $\mathcal{P} = \mathcal{P}\{\eta_1, \eta_2, \dots\}$  be the *linear manifold* spanned by  $\eta_1, \eta_2, \dots$ , i.e. the set of random variables of the form  $\sum_{i=1}^n a_i \eta_i$ ,  $n \geq 1$ ,  $a_i \in R$ . Then  $\overline{\mathcal{P}} = \overline{\mathcal{P}\{\eta_1, \eta_2, \dots\}}$  denotes the *closed linear manifold* spanned by  $\eta_1, \eta_2, \dots$ , i.e. the set of random variables in  $\mathcal{L}$  together with their mean-square limits.

We say that a set  $\eta_1, \eta_2, \dots$  is a *countable orthonormal basis* (or a *complete orthonormal system*) if:

- (a)  $\eta_1, \eta_2, \dots$  is an orthonormal system,
- (b)  $\overline{\mathcal{P}\{\eta_1, \eta_2, \dots\}} = L^2$ .

A Hilbert space with a countable orthonormal basis is said to be *separable*.

By (b), for every  $\xi \in L^2$  and a given  $\varepsilon > 0$  there are numbers  $a_1, \dots, a_n$  such that

$$\left\| \xi - \sum_{i=1}^n a_i \eta_i \right\| \leq \varepsilon.$$

Then by (3)

$$\left\| \xi - \sum_{i=1}^n (\xi, \eta_i) \eta_i \right\| \leq \varepsilon.$$

Consequently every element of a separable Hilbert space  $L^2$  can be represented as

$$\xi = \sum_{i=1}^{\infty} (\xi, \eta_i) \cdot \eta_i, \quad (13)$$

or more precisely as

$$\xi = \text{l.i.m.}_n \sum_{i=1}^n (\xi, \eta_i) \eta_i.$$



We infer from this and (3) that Parseval's equation holds:

$$\|\xi\|^2 = \sum_{i=1}^{\infty} |(\xi, \eta_i)|^2, \quad \xi \in L^2. \quad (14)$$

It is easy to show that the converse is also valid: if  $\eta_1, \eta_2, \dots$  is an orthonormal system and either (13) or (14) is satisfied, then the system is a basis.

We now give some examples of separable Hilbert spaces and their bases.

EXAMPLE 1. Let  $\Omega = R$ ,  $\mathcal{F} = \mathcal{B}(R)$ , and let  $P$  be the Gaussian measure,

$$P(-\infty, a] = \int_{-\infty}^a \varphi(x) dx, \quad \varphi(x) = \frac{1}{\sqrt{2\pi}} e^{-x^2/2}.$$

Let  $D = d/dx$  and

$$H_n(x) = \frac{(-1)^n D^n \varphi(x)}{\varphi(x)}, \quad n \geq 0. \quad (15)$$

We find easily that

$$\begin{aligned} D\varphi(x) &= -x\varphi(x), \\ D^2\varphi(x) &= (x^2 - 1)\varphi(x), \\ D^3\varphi(x) &= (3x - x^3)\varphi(x), \\ &\dots \end{aligned} \quad (16)$$

It follows that  $H_n(x)$  are polynomials (the *Hermite polynomials*). From (15) and (16) we find that

$$\begin{aligned} H_0(x) &= 1, \\ H_1(x) &= x, \\ H_2(x) &= x^2 - 1, \\ H_3(x) &= x^3 - 3x, \\ &\dots \end{aligned}$$

A simple calculation shows that

$$\begin{aligned} (H_m, H_n) &= \int_{-\infty}^{\infty} H_m(x) H_n(x) dP \\ &= \int_{-\infty}^{\infty} H_m(x) H_n(x) \varphi(x) dx = n! \delta_{mn}, \end{aligned}$$

where  $\delta_{mn}$  is the Kronecker delta (0, if  $m \neq n$ , and 1 if  $m = n$ ). Hence if we put

$$h_n(x) = \frac{H_n(x)}{\sqrt{n!}},$$

the system of *normalized Hermite polynomials*  $\{h_n(x)\}_{n \geq 0}$  will be an orthonormal system. We know from functional analysis that if

$$\lim_{\epsilon \downarrow 0} \int_{-\infty}^{\infty} e^{\epsilon |x|} P(dx) < \infty, \quad (17)$$

the system  $\{1, x, x^2, \dots\}$  is complete in  $L^2$ , i.e. every function  $\xi = \xi(x)$  in  $L^2$  can be represented either as  $\sum_{i=1}^n a_i \eta_i(x)$ , where  $\eta_i(x) = x^i$ , or as a limit of these functions (in the mean-square sense). If we apply the Gram-Schmidt orthogonalization process to the sequence  $\eta_1(x), \eta_2(x), \dots$ , with  $\eta_i(x) = x^i$ , the resulting orthonormal system will be precisely the system of normalized Hermite polynomials. In the present case, (17) is satisfied. Hence  $\{h_n(x)\}_{n \geq 0}$  is a basis and therefore every random variable  $\xi = \xi(x)$  on this probability space can be represented in the form

$$\xi(x) = \text{l.i.m.} \sum_{i=0}^n (\xi, h_i) h_i(x). \quad (18)$$

**EXAMPLE 2.** Let  $\Omega = \{0, 1, 2, \dots\}$  and let  $P = \{P_1, P_2, \dots\}$  be the Poisson distribution

$$P_x = \frac{e^{-\lambda} \lambda^x}{x!}, \quad x = 0, 1, \dots; \quad \lambda > 0.$$

Put  $\Delta f(x) = f(x) - f(x-1)$  ( $f(x) = 0, x < 0$ ), and by analogy with (15) define the *Poisson-Charlier polynomials*

$$\Pi_n(x) = \frac{(-1)^n \Delta^n P_x}{P_x}, \quad n \geq 1, \quad \Pi_0 = 1. \quad (19)$$

Since

$$(\Pi_m, \Pi_n) = \sum_{x=0}^{\infty} \Pi_m(x) \Pi_n(x) P_x = c_n \delta_{mn},$$

where  $c_n$  are positive constants, the system of *normalized Poisson-Charlier polynomials*  $\{\pi_n(x)\}_{n \geq 0}$ ,  $\pi_n(x) = \Pi_n(x)/\sqrt{c_n}$ , is an orthonormal system, which is a basis since it satisfies (17).

**EXAMPLE 3.** In this example we describe the Rademacher and Haar systems, which are of interest in function theory as well as in probability theory.

Let  $\Omega = [0, 1]$ ,  $\mathcal{F} = \mathcal{B}([0, 1])$ , and let  $P$  be Lebesgue measure. As we mentioned in §1, every  $x \in [0, 1]$  has a unique binary expansion

$$x = \frac{x_1}{2} + \frac{x_2}{2^2} + \dots,$$

where  $x_i = 0$  or  $1$ . To ensure uniqueness of the expansion, we agree to consider only expansions containing an infinite number of zeros. Thus we

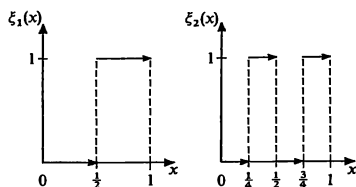


Figure 30

choose the first of the two expansions

$$\frac{1}{2} = \frac{1}{2} + \frac{0}{2^2} + \frac{0}{2^3} + \cdots = \frac{0}{2} + \frac{1}{2^2} + \frac{1}{2^3} + \cdots.$$

We define random variables  $\xi_1(x)$ ,  $\xi_2(x)$ , ... by putting

$$\xi_n(x) = x_n.$$

Then for any numbers  $a_i$ , equal to 0 or 1,

$$\begin{aligned} &P\{x: \xi_1 = a_1, \dots, \xi_n = a_n\} \\ &= P\left\{x: \frac{a_1}{2} + \frac{a_2}{2^2} + \cdots + \frac{a_n}{2^n} \leq x < \frac{a_1}{2} + \frac{a_2}{2^2} + \cdots + \frac{a_n}{2^n} + \frac{1}{2^n}\right\} \\ &= P\left\{x: x \in \left[\frac{a_1}{2} + \cdots + \frac{a_n}{2^n}, \frac{a_1}{2} + \cdots + \frac{a_n}{2^n} + \frac{1}{2^n}\right)\right\} = \frac{1}{2^n}. \end{aligned}$$

It follows immediately that  $\xi_1, \xi_2, \dots$  form a *sequence of independent Bernoulli random variables* (Figure 30 shows the construction of  $\xi_1 = \xi_1(x)$  and  $\xi_2 = \xi_2(x)$ ).

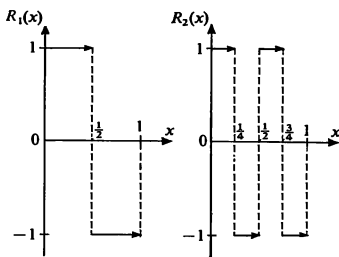


Figure 31. Rademacher functions.

If we now set  $R_n(x) = 1 - 2\xi_n(x)$ ,  $n \geq 1$ , it is easily verified that  $\{R_n\}$  (the *Rademacher functions*, Figure 31) are orthonormal:

$$ER_n R_m = \int_0^1 R_n(x) R_m(x) dx = \delta_{nm}.$$

Notice that  $(1, R_n) \equiv ER_n = 0$ . It follows that this system is not complete.

However, the Rademacher system can be used to construct the *Haar system*, which also has a simple structure and is both *orthonormal* and *complete*.

Again let  $\Omega = [0, 1)$  and  $\mathcal{F} = \mathcal{B}([0, 1))$ . Put

$$H_1(x) = 1,$$

$$H_2(x) = R_1(x),$$

$$H_n(x) = \begin{cases} 2^{j/2} R_j(x) & \text{if } \frac{k-1}{2^j} \leq x < \frac{k}{2^j}, \quad n = 2^j + k, \quad 1 \leq k \leq 2^j, j \geq 1, \\ 0, & \text{otherwise.} \end{cases}$$

It is easy to see that  $H_n(x)$  can also be written in the form

$$H_{2^{m+1}}(x) = \begin{cases} 2^{m/2}, & 0 \leq x < 2^{-(m+1)}, \\ -2^{m/2}, & 2^{-(m+1)} \leq x < 2^{-m}, \quad m = 1, 2, \dots, \\ 0, & \text{otherwise,} \end{cases}$$

$$H_{2^{m+j}}(x) = H_{2^{m+1}}\left(x - \frac{j-1}{2^m}\right), \quad j = 1, \dots, 2^m.$$

Figure 32 shows graphs of the first eight functions, to give an idea of the structure of the Haar functions.

It is easy to see that the Haar system is orthonormal. Moreover, it is complete both in  $L^1$  and in  $L^2$ , i.e. if  $f = f(x) \in L^p$  for  $p = 1$  or  $2$ , then

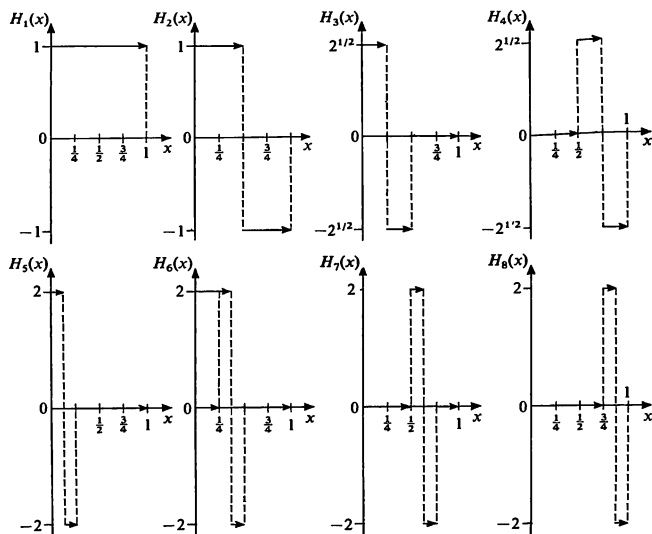
$$\int_0^1 |f(x) - \sum_{k=1}^n (f, H_k) H_k(x)|^p dx \rightarrow 0, \quad n \rightarrow \infty.$$

The system also has the property that

$$\sum_{k=1}^n (f, H_k) H_k(x) \rightarrow f(x), \quad n \rightarrow \infty,$$

with probability 1 (with respect to Lebesgue measure).

In §4, Chapter VII, we shall prove these facts by deriving them from general theorems on the convergence of martingales. This will, in particular, provide a good illustration of the application of martingale methods to the theory of functions.

Figure 32. The Haar functions  $H_1(x), \dots, H_8(x)$ .

6. If  $\eta_1, \dots, \eta_n$  is a finite orthonormal system then, as was shown above, for every random variable  $\xi \in L^2$  there is a random variable  $\hat{\xi}$  in the linear manifold  $\mathcal{L} = \mathcal{L}\{\eta_1, \dots, \eta_n\}$ , namely the projection of  $\xi$  on  $\mathcal{L}$ , such that

$$\|\xi - \hat{\xi}\| = \inf\{\|\xi - \zeta\| : \zeta \in \mathcal{L}\{\eta_1, \dots, \eta_n\}\}.$$

Here  $\hat{\xi} = \sum_{i=1}^n (\xi, \eta_i) \eta_i$ . This result has a natural generalization to the case when  $\eta_1, \eta_2, \dots$  is a countable orthonormal system (not necessarily a basis). In fact, we have the following result.

**Theorem.** Let  $\eta_1, \eta_2, \dots$  be an orthonormal system of random variables, and  $\bar{\mathcal{L}} = \bar{\mathcal{L}}\{\eta_1, \eta_2, \dots\}$  the closed linear manifold spanned by the system. Then there is a unique element  $\hat{\xi} \in \bar{\mathcal{L}}$  such that

$$\|\xi - \hat{\xi}\| = \inf\{\|\xi - \zeta\| : \zeta \in \bar{\mathcal{L}}\}. \quad (20)$$

Moreover,

$$\hat{\xi} = \text{l.i.m.}_n \sum_{i=1}^n (\xi, \eta_i) \eta_i \quad (21)$$

and  $\xi - \hat{\xi} \perp \zeta, \zeta \in \bar{\mathcal{L}}$ .

PROOF. Let  $d = \inf\{\|\xi - \zeta\|: \zeta \in \overline{\mathcal{L}}\}$  and choose a sequence  $\zeta_1, \zeta_2, \dots$  such that  $\|\xi - \zeta_n\| \rightarrow d$ . Let us show that this sequence is fundamental. A simple calculation shows that

$$\|\zeta_n - \zeta_m\|^2 = 2\|\zeta_n - \xi\|^2 + 2\|\zeta_m - \xi\|^2 - 4\left\|\frac{\zeta_n + \zeta_m}{2} - \xi\right\|^2.$$

It is clear that  $(\zeta_n + \zeta_m)/2 \in \overline{\mathcal{L}}$ ; consequently  $\|[(\zeta_n + \zeta_m)/2] - \xi\|^2 \geq d^2$  and therefore  $\|\zeta_n - \zeta_m\|^2 \rightarrow 0$ ,  $n, m \rightarrow \infty$ .

The space  $L^2$  is complete (Theorem 7, §10). Hence there is an element  $\tilde{\xi}$  such that  $\|\zeta_n - \tilde{\xi}\| \rightarrow 0$ . But  $\overline{\mathcal{L}}$  is closed, so  $\tilde{\xi} \in \overline{\mathcal{L}}$ . Moreover,  $\|\zeta_n - \xi\| \rightarrow d$ , and consequently  $\|\xi - \tilde{\xi}\| = d$ , which establishes the existence of the required element.

Let us show that  $\tilde{\xi}$  is the only element of  $\overline{\mathcal{L}}$  with the required property. Let  $\tilde{\xi} \in \overline{\mathcal{L}}$  and let

$$\|\xi - \tilde{\xi}\| = \|\xi - \tilde{\xi}\| = d.$$

Then (by Problem 3)

$$\|\tilde{\xi} + \tilde{\xi} - 2\xi\|^2 + \|\xi - \tilde{\xi}\|^2 = 2\|\tilde{\xi} - \xi\|^2 + 2\|\tilde{\xi} - \xi\|^2 = 4d^2.$$

But

$$\|\tilde{\xi} + \tilde{\xi} - 2\xi\|^2 = 4\|\frac{1}{2}(\tilde{\xi} + \tilde{\xi}) - \xi\|^2 \geq 4d^2.$$

Consequently  $\|\tilde{\xi} - \xi\|^2 = 0$ . This establishes the uniqueness of the element of  $\overline{\mathcal{L}}$  that is closest to  $\xi$ .

Now let us show that  $\xi - \tilde{\xi} \perp \zeta$ ,  $\zeta \in \overline{\mathcal{L}}$ . By (20),

$$\|\xi - \tilde{\xi} - c\zeta\| \geq \|\xi - \tilde{\xi}\|$$

for every  $c \in R$ . But

$$\|\xi - \tilde{\xi} - c\zeta\|^2 = \|\xi - \tilde{\xi}\|^2 + c^2\|\zeta\|^2 - 2(\xi - \tilde{\xi}, c\zeta).$$

Therefore

$$c^2\|\zeta\|^2 \geq 2(\xi - \tilde{\xi}, c\zeta). \quad (22)$$

Take  $c = \lambda(\xi - \tilde{\xi}, \zeta)$ ,  $\lambda \in R$ . Then we find from (22) that

$$(\xi - \tilde{\xi}, \zeta)^2[\lambda^2\|\zeta\|^2 - 2\lambda] \geq 0.$$

We have  $\lambda^2\|\zeta\|^2 - 2\lambda < 0$  if  $\lambda$  is a sufficiently small positive number. Consequently  $(\xi - \tilde{\xi}, \zeta) = 0$ ,  $\zeta \in \overline{\mathcal{L}}$ .

It remains only to prove (21).

The set  $\mathcal{Z} = \mathcal{Z}\{\eta_1, \eta_2, \dots\}$  is a closed subspace of  $L^2$  and therefore a Hilbert space (with the same scalar product). Now the system  $\eta_1, \eta_2, \dots$  is a basis for  $\mathcal{Z}$  (Problem 5), and consequently

$$\tilde{\xi} = \text{l.i.m.} \sum_{k=1}^n (\tilde{\xi}, \eta_k) \eta_k. \quad (23)$$

But  $\xi - \tilde{\xi} \perp \eta_k$ ,  $k \geq 1$ , and therefore  $(\tilde{\xi}, \eta_k) = (\xi, \eta_k)$ ,  $k \geq 0$ . This, with (23) establishes (21).

This completes the proof of the theorem.

**Remark.** As in the finite-dimensional case, we say that  $\xi$  is the projection of  $\xi$  on  $\bar{L} = \bar{L}\{\eta_1, \eta_2, \dots\}$ , that  $\xi - \xi$  is perpendicular to  $\bar{L}$ , and that the representation

$$\xi = \xi + (\xi - \xi)$$

is the orthogonal decomposition of  $\xi$ .

We also denote  $\xi$  by  $\hat{E}(\xi|\eta_1, \eta_2, \dots)$  and call it the *conditional expectation in the wide sense* (of  $\xi$  with respect to  $\eta_1, \eta_2, \dots$ ). From the point of view of estimating  $\xi$  in terms of  $\eta_1, \eta_2, \dots$ , the variable  $\xi$  is the optimal linear estimator, with error

$$\Delta \equiv E|\xi - \xi|^2 \equiv \|\xi - \xi\|^2 = \|\xi\|^2 - \sum_{i=1}^{\infty} |(\xi, \eta_i)|^2,$$

which follows from (5) and (23).

## 7. PROBLEMS

1. Show that if  $\xi = \text{l.i.m. } \xi_n$  then  $\|\xi_n\| \rightarrow \|\xi\|$ .
2. Show that if  $\xi = \text{l.i.m. } \xi_n$  and  $\eta = \text{l.i.m. } \eta_n$  then  $(\xi_n, \eta_n) \rightarrow (\xi, \eta)$ .
3. Show that the norm  $\|\cdot\|$  has the *parallelogram property*

$$\|\xi + \eta\|^2 + \|\xi - \eta\|^2 = 2(\|\xi\|^2 + \|\eta\|^2).$$
4. Let  $(\xi_1, \dots, \xi_n)$  be a family of orthogonal random variables. Show that they have the *Pythagorean property*,

$$\left\| \sum_{i=1}^n \xi_i \right\|^2 = \sum_{i=1}^n \|\xi_i\|^2.$$

5. Let  $\eta_1, \eta_2, \dots$  be an orthonormal system and  $\mathcal{L} = \mathcal{L}\{\eta_1, \eta_2, \dots\}$  the closed linear manifold spanned by  $\eta_1, \eta_2, \dots$ . Show that the system is a basis for the (Hilbert) space  $\mathcal{L}$ .
6. Let  $\xi_1, \xi_2, \dots$  be a sequence of orthogonal random variables and  $S_n = \xi_1 + \dots + \xi_n$ . Show that if  $\sum_{n=1}^{\infty} E\xi_n^2 < \infty$  there is a random variable  $S$  with  $ES^2 < \infty$  such that  $\text{l.i.m. } S_n = S$ , i.e.  $\|S_n - S\|^2 = E|S_n - S|^2 \rightarrow 0, n \rightarrow \infty$ .
7. Show that in the space  $L^2 = L^2([-\pi, \pi], \mathcal{B}([-\pi, \pi]))$  with Lebesgue measure  $\mu$  the system  $\{(1/\sqrt{2\pi})e^{i\lambda n}, n = 0, \pm 1, \dots\}$  is an orthonormal basis.

## §12. Characteristic Functions

1. The method of characteristic functions is one of the main tools of the analytic theory of probability. This will appear very clearly in Chapter III in the proofs of limit theorems and, in particular, in the proof of the central limit theorem, which generalizes the De Moivre-Laplace theorem. In the present section we merely define characteristic functions and present their basic properties.

First we make some general remarks.

Besides random variables which take real values, the theory of characteristic functions requires random variables that take complex values (see Subsection 1 of §5).

Many definitions and properties involving random variables can easily be carried over to the complex case. For example, the expectation  $E\zeta$  of a complex random variable  $\zeta = \xi + i\eta$  will exist if the expectations  $E\xi$  and  $E\eta$  exist. In this case we define  $E\zeta = E\xi + iE\eta$ . It is easy to deduce from the definition of the independence of random elements (Definition 6, §5) that the complex random variables  $\zeta_1 = \xi_1 + i\eta_1$  and  $\zeta_2 = \xi_2 + i\eta_2$  are independent if and only if the pairs  $(\xi_1, \eta_1)$  and  $(\xi_2, \eta_2)$  are independent; or, equivalently, the  $\sigma$ -algebras  $\mathcal{L}_{\xi_1, \eta_1}$  and  $\mathcal{L}_{\xi_2, \eta_2}$  are independent.

Besides the space  $L^2$  of real random variables with finite second moment, we shall consider the Hilbert space of complex random variables  $\zeta = \xi + i\eta$  with  $E|\zeta|^2 < \infty$ , where  $|\zeta|^2 = \xi^2 + \eta^2$  and the scalar product  $(\zeta_1, \zeta_2)$  is defined by  $E\zeta_1\bar{\zeta}_2$ , where  $\bar{\zeta}_2$  is the complex conjugate of  $\zeta$ . The term "random variable" will now be used for both real and complex random variables, with a comment (when necessary) on which is intended.

Let us introduce some notation.

We consider a vector  $a \in R^n$  to be a column vector,

$$a = \begin{pmatrix} a_1 \\ \vdots \\ a_n \end{pmatrix},$$

and  $a^T$  to be a row vector,  $a^T = (a_1, \dots, a_n)$ . If  $a$  and  $b \in R^n$  their scalar product  $(a, b)$  is  $\sum_{i=1}^n a_i b_i$ . Clearly  $(a, b) = a^T b$ .

If  $a \in R^n$  and  $\mathbb{R} = \|r_{ij}\|$  is an  $n$  by  $n$  matrix,

$$(\mathbb{R}a, a) = a^T \mathbb{R}a = \sum_{i,j=1}^n r_{ij} a_i a_j. \quad (1)$$

**2. Definition 1.** Let  $F = F(x_1, \dots, x_n)$  be an  $n$ -dimensional distribution function in  $(R^n, \mathcal{B}(R^n))$ . Its *characteristic function* is

$$\varphi(t) = \int_{R^n} e^{i(t, x)} dF(x), \quad t \in R^n. \quad (2)$$

**Definition 2.** If  $\xi = (\xi_1, \dots, \xi_n)$  is a random vector defined on the probability space  $(\Omega, \mathcal{F}, P)$  with values in  $R^n$ , its *characteristic function* is

$$\varphi_\xi(t) = \int_{R^n} e^{i(t, x)} dF_\xi(x), \quad t \in R^n, \quad (3)$$

where  $F_\xi = F_\xi(x_1, \dots, x_n)$  is the distribution function of the vector  $\xi = (\xi_1, \dots, \xi_n)$ .

If  $F(x)$  has a density  $f = f(x)$  then

$$\varphi(t) = \int_{R^n} e^{i(t, x)} f(x) dx.$$



In other words, in this case the characteristic function is just the Fourier transform of  $f(x)$ .

It follows from (3) and Theorem 6.7 (on change of variable in a Lebesgue integral) that the characteristic function  $\varphi_\xi(t)$  of a random vector can also be defined by

$$\varphi_\xi(t) = \mathbb{E}e^{it(\cdot, \xi)}, \quad t \in R^n. \quad (4)$$

We now present some basic properties of characteristic functions, stated and proved for  $n = 1$ . Further important results for the general case will be given as problems.

Let  $\xi = \xi(\omega)$  be a random variable,  $F_\xi = F_\xi(x)$  its distribution function, and

$$\varphi_\xi(t) = \mathbb{E}e^{it\xi}$$

its characteristic function.

We see at once that if  $\eta = a\xi + b$  then

$$\varphi_\eta(t) = \mathbb{E}e^{it\eta} = \mathbb{E}e^{it(a\xi + b)} = e^{itb}\mathbb{E}e^{iat\xi}.$$

Therefore

$$\varphi_\eta(t) = e^{itb}\varphi_\xi(at). \quad (5)$$

Moreover, if  $\xi_1, \xi_2, \dots, \xi_n$  are independent random variables and  $S_n = \xi_1 + \dots + \xi_n$ , then

$$\varphi_{S_n}(t) = \prod_{j=1}^n \varphi_{\xi_j}(t). \quad (6)$$

In fact,

$$\begin{aligned} \varphi_{S_n} &= \mathbb{E}e^{it(\xi_1 + \dots + \xi_n)} = \mathbb{E}e^{it\xi_1} \dots e^{it\xi_n} \\ &= \mathbb{E}e^{it\xi_1} \dots \mathbb{E}e^{it\xi_n} = \prod_{j=1}^n \varphi_{\xi_j}(t), \end{aligned}$$

where we have used the property that the expectation of a product of independent (bounded) random variables (either real or complex; see Theorem 6 of §6, and Problem 1) is equal to the product of their expectations.

Property (6) is the key to the proofs of limit theorems for sums of independent random variables by the method of characteristic functions (see §3, Chapter III). In this connection we note that the distribution function  $F_{S_n}$  is expressed in terms of the distribution functions of the individual terms in a rather complicated way, namely  $F_{S_n} = F_{\xi_1} * \dots * F_{\xi_n}$  where  $*$  denotes convolution (see §8, Subsection 4).

Here are some examples of characteristic functions.

**EXAMPLE 1.** Let  $\xi$  be a Bernoulli random variable with  $P(\xi = 1) = p$ ,  $P(\xi = 0) = q$ ,  $p + q = 1$ ,  $1 > p > 0$ ; then

$$\varphi_\xi(t) = pe^{it} + q.$$

If  $\xi_1, \dots, \xi_n$  are independent identically distributed random variables like  $\xi$ , then, writing  $T_n = (S_n - np)/\sqrt{npq}$ , we have

$$\begin{aligned}\varphi_{T_n}(t) &= \mathbb{E}e^{iT_n t} = e^{-it\sqrt{npq}}[pe^{it/\sqrt{npq}} + q]^n \\ &= [pe^{it\sqrt{q/(np)}} + qe^{-it\sqrt{p/(nq)}}]^n.\end{aligned}\quad (7)$$

Notice that it follows that as  $n \rightarrow \infty$

$$\varphi_{T_n}(t) \rightarrow e^{-t^2/2}, \quad T_n = \frac{S_n - np}{\sqrt{npq}}. \quad (8)$$

EXAMPLE 2. Let  $\xi \sim \mathcal{N}(m, \sigma^2)$ ,  $|m| < \infty$ ,  $\sigma^2 > 0$ . Let us show that

$$\varphi_\xi(t) = e^{itm - t^2\sigma^2/2} \quad (9)$$

Let  $\eta = (\xi - m)/\sigma$ . Then  $\eta \sim \mathcal{N}(0, 1)$  and, since

$$\varphi_\xi(t) = e^{itm}\varphi_\eta(\sigma t)$$

by (5), it is enough to show that

$$\varphi_\eta(t) = e^{-t^2/2}. \quad (10)$$

We have

$$\begin{aligned}\varphi_\eta(t) &= \mathbb{E}e^{it\eta} = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} e^{itx} e^{-x^2/2} dx \\ &= \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} \sum_{n=0}^{\infty} \frac{(itx)^n}{n!} e^{-x^2/2} dx = \sum_{n=0}^{\infty} \frac{(it)^n}{n!} \frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} x^n e^{-x^2/2} dx \\ &= \sum_{n=0}^{\infty} \frac{(it)^{2n}}{(2n)!} (2n-1)!! = \sum_{n=0}^{\infty} \frac{(it)^{2n} (2n)!}{(2n)! 2^n n!} \\ &= \sum_{n=0}^{\infty} \left(-\frac{t^2}{2}\right)^n \cdot \frac{1}{n!} = e^{-t^2/2},\end{aligned}$$

where we have used the formula (see Problem 7 in §8)

$$\frac{1}{\sqrt{2\pi}} \int_{-\infty}^{\infty} x^{2n} e^{-x^2/2} dx \equiv \mathbb{E}\eta^{2n} = (2n-1)!!.$$

EXAMPLE 3. Let  $\xi$  be a Poisson random variable,

$$P(\xi = k) = \frac{e^{-\lambda} \lambda^k}{k!}, \quad k = 0, 1, \dots$$

Then

$$\mathbb{E}e^{it\xi} = \sum_{k=0}^{\infty} e^{itk} \frac{e^{-\lambda} \lambda^k}{k!} = e^{-\lambda} \sum_{k=0}^{\infty} \frac{(\lambda e^{it})^k}{k!} = \exp\{\lambda(e^{it} - 1)\}. \quad (11)$$

3. As we observed in §9, Subsection 1, with every distribution function in  $(R, \mathcal{B}(R))$  we can associate a random variable of which it is the distribution function. Hence in discussing the properties of characteristic functions (in the sense either of Definition 1 or Definition 2), we may consider only characteristic functions  $\varphi(t) = \varphi_\xi(t)$  of random variables  $\xi = \xi(\omega)$ .

**Theorem 1.** Let  $\xi$  be a random variable with distribution function  $F = F(x)$  and

$$\varphi(t) = \mathbb{E} e^{it\xi}$$

its characteristic function. Then  $\varphi$  has the following properties:

- (1)  $|\varphi(t)| \leq \varphi(0) = 1$ ;
- (2)  $\varphi(t)$  is uniformly continuous for  $t \in R$ ;
- (3)  $\varphi(t) = \overline{\varphi(-t)}$ ;
- (4)  $\varphi(t)$  is real-valued if and only if  $F$  is symmetric ( $\int_B dF(x) = \int_{-B} dF(x)$ ),  $B \in \mathcal{B}(R)$ ,  $-B = \{-x: x \in B\}$ ;
- (5) if  $\mathbb{E}|\xi|^n < \infty$  for some  $n \geq 1$ , then  $\varphi^{(r)}(t)$  exists for every  $r \leq n$ , and

$$\varphi^{(r)}(t) = \int_R (ix)^r e^{itx} dF(x), \quad (12)$$

$$\mathbb{E}\xi^r = \frac{\varphi^{(r)}(0)}{i^r}, \quad (13)$$

$$\varphi(t) = \sum_{r=0}^n \frac{(it)^r}{r!} \mathbb{E}\xi^r + \frac{(it)^n}{n!} \varepsilon_n(t), \quad (14)$$

where  $|\varepsilon_n(t)| \leq 3\mathbb{E}|\xi|^n$  and  $\varepsilon_n(t) \rightarrow 0$ ,  $t \rightarrow 0$ ;

- (6) if  $\varphi^{(2n)}(0)$  exists and is finite then  $\mathbb{E}\xi^{2n} < \infty$ ;
- (7) if  $\mathbb{E}|\xi|^n < \infty$  for all  $n \geq 1$  and

$$\overline{\lim}_n \frac{(\mathbb{E}|\xi|^n)^{1/n}}{n} = \frac{1}{e \cdot R} < \infty,$$

then

$$\varphi(t) = \sum_{n=0}^{\infty} \frac{(it)^n}{n!} \mathbb{E}\xi^n. \quad (15)$$

for all  $|t| < R$ .

**PROOF.** Properties (1) and (3) are evident. Property (2) follows from the inequality

$$|\varphi(t+h) - \varphi(t)| = |\mathbb{E} e^{it\xi}(e^{ih\xi} - 1)| \leq \mathbb{E} |e^{ih\xi} - 1|$$

and the dominated convergence theorem, according to which  $\mathbb{E} |e^{ih\xi} - 1| \rightarrow 0$ ,  $h \rightarrow 0$ .

Property (4). Let  $F$  be symmetric. Then if  $g(x)$  is a bounded odd Borel function, we have  $\int_R g(x) dF(x) = 0$  (observe that for simple odd functions

this follows directly from the definition of the symmetry of  $F$ ). Consequently  $\int_R \sin tx \, dF(x) = 0$  and therefore

$$\varphi(t) = E \cos t\xi.$$

Conversely, let  $\varphi_\xi(t)$  be a real function. Then by (3)

$$\varphi_{-\xi}(t) = \varphi_\xi(-t) = \overline{\varphi_\xi(t)} = \varphi_\xi(t), \quad t \in R.$$

Hence (as will be shown below in Theorem 2) the distribution functions  $F_{-\xi}$  and  $F_\xi$  of the random variables  $-\xi$  and  $\xi$  are the same, and therefore (by Theorem 3.1)

$$P(\xi \in B) = P(-\xi \in B) = P(\xi \in -B)$$

for every  $B \in \mathcal{B}(R)$ .

Property (5). If  $E|\xi|^n < \infty$ , we have  $E|\xi|^r < \infty$  for  $r \leq n$ , by Lyapunov's inequality (6.28).

Consider the difference quotient

$$\frac{\varphi(t+h) - \varphi(t)}{h} = E e^{it\xi} \left( \frac{e^{ih\xi} - 1}{h} \right).$$

Since

$$\left| \frac{e^{ihx} - 1}{h} \right| \leq |x|,$$

and  $E|\xi| < \infty$ , it follows from the dominated convergence theorem that the limit

$$\lim_{h \rightarrow 0} E e^{it\xi} \left( \frac{e^{ih\xi} - 1}{h} \right)$$

exists and equals

$$E e^{it\xi} \lim_{h \rightarrow 0} \left( \frac{e^{ih\xi} - 1}{h} \right) = i E (\xi e^{it\xi}) = i \int_{-\infty}^{\infty} x e^{itx} dF(x). \quad (16)$$

Hence  $\varphi'(t)$  exists and

$$\varphi'(t) = i(E \xi e^{it\xi}) = i \int_{-\infty}^{\infty} x e^{itx} dF(x).$$

The existence of the derivatives  $\varphi^{(r)}(t)$ ,  $1 < r \leq n$ , and the validity of (12), follow by induction.

Formula (13) follows immediately from (12). Let us now establish (14). Since

$$e^{iy} = \cos y + i \sin y = \sum_{k=0}^{n-1} \frac{(iy)^k}{k!} + \frac{(iy)^n}{n!} [\cos \theta_1 y + i \sin \theta_2 y]$$

for real  $y$ , with  $|\theta_1| \leq 1$  and  $|\theta_2| \leq 1$ , we have

$$e^{it\xi} = \sum_{k=0}^{n-1} \frac{(it\xi)^k}{k!} + \frac{(it\xi)^n}{n!} [\cos \theta_1(\omega)t\xi + i \sin \theta_2(\omega)t\xi] \quad (17)$$

and

$$E e^{it\xi} = \sum_{k=0}^{n-1} \frac{(it)^k}{k!} E \xi^k + \frac{(it)^n}{n!} [E \xi^n + \varepsilon_n(t)], \quad (18)$$

where

$$\varepsilon_n(t) = E[\xi^n (\cos \theta_1(\omega)t\xi + i \sin \theta_2(\omega)t\xi - 1)].$$

It is clear that  $|\varepsilon_n(t)| \leq 3E|\xi^n|$ . The theorem on dominated convergence shows that  $\varepsilon_n(t) \rightarrow 0$ ,  $t \rightarrow 0$ .

Property (6). We give a proof by induction. Suppose first that  $\varphi''(0)$  exists and is finite. Let us show that in that case  $E\xi^2 < \infty$ . By L'Hôpital's rule and Fatou's lemma,

$$\begin{aligned} \varphi''(0) &= \lim_{h \rightarrow 0} \frac{1}{2} \left[ \frac{\varphi'(2h) - \varphi'(0)}{2h} + \frac{\varphi'(0) - \varphi'(-2h)}{2h} \right] \\ &= \lim_{h \rightarrow 0} \frac{2\varphi'(2h) - 2\varphi'(-2h)}{8h} = \lim_{h \rightarrow 0} \frac{1}{4h^2} [\varphi(2h) - 2\varphi(0) + \varphi(-2h)] \\ &= \lim_{h \rightarrow 0} \int_{-\infty}^{\infty} \left( \frac{e^{ihx} - e^{-ihx}}{2h} \right)^2 dF(x) \\ &= -\lim_{h \rightarrow 0} \int_{-\infty}^{\infty} \left( \frac{\sin hx}{hx} \right)^2 x^2 dF(x) \leq -\int_{-\infty}^{\infty} \lim_{h \rightarrow 0} \left( \frac{\sin hx}{hx} \right)^2 x^2 dF(x) \\ &= -\int_{-\infty}^{\infty} x^2 dF(x). \end{aligned}$$

Therefore,

$$\int_{-\infty}^{\infty} x^2 dF(x) \leq -\varphi''(0) < \infty.$$

Now let  $\varphi^{(2k+2)}(0)$  exist, finite, and let  $\int_{-\infty}^{\infty} x^{2k} dF(x) < \infty$ . If  $\int_{-\infty}^{\infty} x^{2k} dF(x) = 0$ , then  $\int_{-\infty}^{\infty} x^{2k+2} dF(x) = 0$  also. Hence we may suppose that  $\int_{-\infty}^{\infty} x^{2k} dF(x) > 0$ . Then, by Property (5),

$$\varphi^{(2k)}(t) = \int_{-\infty}^{\infty} (ix)^{2k} e^{itx} dF(x)$$

and therefore,

$$(-1)^k \varphi^{(2k)}(t) = \int_{-\infty}^{\infty} e^{itx} dG(x),$$

where  $G(x) = \int_{-\infty}^x u^{2k} dF(u)$ .

Consequently the function  $(-1)^k \varphi^{(2k)}(t) G(\infty)^{-1}$  is the characteristic function of the probability distribution  $G(x) \cdot G^{-1}(\infty)$  and by what we have proved,

$$G^{-1}(\infty) \int_{-\infty}^{\infty} x^2 dG(x) < \infty.$$

But  $G^{-1}(\infty) > 0$ , and therefore

$$\int_{-\infty}^{\infty} x^{2k+2} dF(x) = \int_{-\infty}^{\infty} x^2 dG(x) < \infty.$$

Property (7). Let  $0 < t_0 < R$ . Then, by Stirling's formula we find that

$$\overline{\lim} \frac{(E|\xi|^n)^{1/n}}{n} < \frac{1}{e \cdot t_0} \Rightarrow \overline{\lim} \frac{(E|\xi|^n t_0^n)^{1/n}}{n} < \frac{1}{e} \Rightarrow \lim \left( \frac{E|\xi|^n t_0^n}{n!} \right)^{1/n} < 1.$$

Consequently the series  $\sum [E|\xi|^n t_0^n / n!]$  converges by Cauchy's test, and therefore the series  $\sum_{r=0}^{\infty} [(it)^r / r!] E\xi^r$  converges for  $|t| \leq t_0$ . But by (14), for  $n \geq 1$ ,

$$\varphi(t) = \sum_{r=0}^n \frac{(it)^r}{r!} E\xi^r + R_n(t),$$

where  $|R_n(t)| \leq 3(|t|^n / n!) E|\xi|^n$ . Therefore

$$\varphi(t) = \sum_{r=0}^{\infty} \frac{(it)^r}{r!} E\xi^r$$

for all  $|t| < R$ . This completes the proof of the theorem.

**Remark 1.** By a method similar to that used for (14), we can establish that if  $E|\xi|^n < \infty$  for some  $n \geq 1$ , then

$$\varphi(t) = \sum_{k=0}^n \frac{i^k (t-s)^k}{k!} \int_{-\infty}^{\infty} x^k e^{isx} dF(x) + \frac{i^n (t-s)^n}{n!} \varepsilon_n(t-s), \quad (19)$$

where  $|\varepsilon_n(t-s)| \leq 3E|\xi|^n$ , and  $\varepsilon_n(t-s) \rightarrow 0$  as  $t-s \rightarrow 0$ .

**Remark 2.** With reference to the condition that appears in Property (7), see also Subsection 9, below, on the "uniqueness of the solution of the moment problem."

4. The following theorem shows that the characteristic function is uniquely determined by the distribution function.

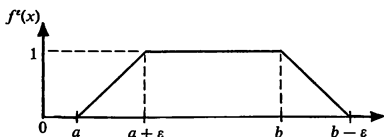


Figure 33

**Theorem 2 (Uniqueness).** Let  $F$  and  $G$  be distribution functions with the same characteristic function, i.e.

$$\int_{-\infty}^{\infty} e^{itx} dF(x) = \int_{-\infty}^{\infty} e^{itx} dG(x) \quad (20)$$

for all  $t \in \mathbb{R}$ . Then  $F(x) \equiv G(x)$ .

**PROOF.** Choose  $a$  and  $b \in \mathbb{R}$ , and  $\epsilon > 0$ , and consider the function  $f^\epsilon = f^\epsilon(x)$  shown in Figure 33. We show that

$$\int_{-\infty}^{\infty} f^\epsilon(x) dF(x) = \int_{-\infty}^{\infty} f^\epsilon(x) dG(x). \quad (21)$$

Let  $n \geq 0$  be large enough so that  $[a - \epsilon, b + \epsilon] \subseteq [-n, n]$ , and let the sequence  $\{\delta_n\}$  be such that  $1 \geq \delta_n \downarrow 0$ ,  $n \rightarrow \infty$ . Like every continuous function on  $[-n, n]$  that has equal values at the endpoints,  $f^\epsilon = f^\epsilon(x)$  can be uniformly approximated by trigonometric polynomials (Weierstrass's theorem), i.e. there is a finite sum

$$f_n^\epsilon(x) = \sum_k a_k \exp\left(i\pi x \frac{k}{n}\right) \quad (22)$$

such that

$$\sup_{-n \leq x \leq n} |f^\epsilon(x) - f_n^\epsilon(x)| \leq \delta_n. \quad (23)$$

Let us extend the periodic function  $f_n^\epsilon(x)$  to all of  $\mathbb{R}$ , and observe that

$$\sup_x |f_n^\epsilon(x)| \leq 2.$$

Then, since by (20)

$$\int_{-\infty}^{\infty} f_n^\epsilon(x) dF(x) = \int_{-\infty}^{\infty} f_n^\epsilon(x) dG(x),$$

we have

$$\begin{aligned}
 \left| \int_{-\infty}^{\infty} f^{\varepsilon}(x) dF(x) - \int_{-\infty}^{\infty} f^{\varepsilon}(x) dG(x) \right| &= \left| \int_{-n}^n f^{\varepsilon} dF - \int_{-n}^n f^{\varepsilon} dG \right| \\
 &\leq \left| \int_{-n}^n f_n^{\varepsilon} dF - \int_{-n}^n f_n^{\varepsilon} dG \right| + 2\delta_n \\
 &\leq \left| \int_{-\infty}^{\infty} f_n^{\varepsilon} dF - \int_{-\infty}^{\infty} f_n^{\varepsilon} dG \right| + 2\delta_n \\
 &\quad + 2F(\overline{[-n, n]}) + 2G(\overline{[-n, n]}),
 \end{aligned} \tag{24}$$

where  $F(A) = \int_A dF(x)$ ,  $G(A) = \int_A dG(x)$ . As  $n \rightarrow \infty$ , the right-hand side of (24) tends to zero, and this establishes (21).

As  $\varepsilon \rightarrow 0$ , we have  $f^{\varepsilon}(x) \rightarrow I_{[a, b]}(x)$ . It follows from (21) by the theorem on distribution functions' being the same.

$$\int_{-\infty}^{\infty} I_{[a, b]}(x) dF(x) = \int_{-\infty}^{\infty} I_{[a, b]}(x) dG(x),$$

i.e.  $F(b) - F(a) = G(b) - G(a)$ . Since  $a$  and  $b$  are arbitrary, it follows that  $F(x) = G(x)$  for all  $x \in R$ .

This completes the proof of the theorem.

5. The preceding theorem says that a distribution function  $F = F(x)$  is uniquely determined by its characteristic function  $\varphi = \varphi(t)$ . The next theorem gives an explicit representation of  $F$  in terms of  $\varphi$ .

**Theorem 3 (Inversion Formula).** *Let  $F = F(x)$  be a distribution function and*

$$\varphi(t) = \int_{-\infty}^{\infty} e^{itx} dF(x)$$

*its characteristic function.*

(a) *For pairs of points  $a$  and  $b$  ( $a < b$ ) at which  $F = F(x)$  is continuous,*

$$F(b) - F(a) = \lim_{\varepsilon \rightarrow \infty} \frac{1}{2\pi} \int_{-\varepsilon}^{\varepsilon} \frac{e^{-ita} - e^{-itb}}{it} \varphi(t) dt; \tag{25}$$

(b) *If  $\int_{-\infty}^{\infty} |\varphi(t)| dt < \infty$ , the distribution function  $F(x)$  has a density  $f(x)$ ,*

$$F(x) = \int_{-\infty}^x f(y) dy \tag{26}$$

*and*

$$f(x) = \frac{1}{2\pi} \int_{-\infty}^{\infty} e^{-itx} \varphi(t) dt. \tag{27}$$



PROOF. We first observe that if  $F(x)$  has density  $f(x)$  then

$$\varphi(t) = \int_{-\infty}^{\infty} e^{itx} f(x) dx, \quad (28)$$

and (27) is just the Fourier transform of the (integrable) function  $\varphi(t)$ . Integrating both sides of (27) and applying Fubini's theorem, we obtain

$$\begin{aligned} F(b) - F(a) &= \int_a^b f(x) dx = \frac{1}{2\pi} \int_a^b \left[ \int_{-\infty}^{\infty} e^{-itx} \varphi(t) dt \right] dx \\ &= \frac{1}{2\pi} \int_{-\infty}^{\infty} \varphi(t) \left[ \int_a^b e^{-itx} dx \right] dt \\ &= \frac{1}{2\pi} \int_{-\infty}^{\infty} \varphi(t) \frac{e^{-ita} - e^{-itb}}{it} dt. \end{aligned}$$

After these remarks, which to some extent clarify (25), we turn to the proof.

(a) We have

$$\begin{aligned} \Phi_c &\equiv \frac{1}{2\pi} \int_{-c}^c \frac{e^{-ita} - e^{-itb}}{it} \varphi(t) dt \\ &= \frac{1}{2\pi} \int_{-c}^c \frac{e^{-ita} - e^{-itb}}{it} \left[ \int_{-\infty}^{\infty} e^{itx} dF(x) \right] dt \\ &= \frac{1}{2\pi} \int_{-\infty}^{\infty} \left[ \int_{-c}^c \frac{e^{-ita} - e^{-itb}}{it} e^{itx} dt \right] dF(x) \\ &= \int_{-\infty}^{\infty} \Psi_c(x) dF(x), \end{aligned} \quad (29)$$

where we have put

$$\Psi_c(x) = \frac{1}{2\pi} \int_{-c}^c \frac{e^{-ita} - e^{-itb}}{it} e^{itx} dt$$

and applied Fubini's theorem, which is applicable in this case because

$$\left| \frac{e^{-ita} - e^{-itb}}{it} \cdot e^{itx} \right| = \left| \frac{e^{-ita} - e^{-itb}}{it} \right| = \left| \int_a^b e^{-itx} dx \right| \leq b - a$$

and

$$\int_{-c}^c \int_{-\infty}^{\infty} (b - a) dF(x) \leq 2c(b - a) < \infty.$$

In addition,

$$\begin{aligned}\Psi_c(x) &= \frac{1}{2\pi} \int_{-c}^c \frac{\sin t(x-a) - \sin t(x-b)}{t} dt \\ &= \frac{1}{2\pi} \int_{-c(x-a)}^{c(x-a)} \frac{\sin v}{v} dv - \frac{1}{2\pi} \int_{-c(x-b)}^{c(x-b)} \frac{\sin u}{u} du.\end{aligned}\quad (30)$$

The function

$$g(s, t) = \int_s^t \frac{\sin v}{v} dv$$

is uniformly continuous in  $s$  and  $t$ , and

$$g(s, t) \rightarrow \pi \quad (31)$$

as  $s \downarrow -\infty$  and  $t \uparrow \infty$ . Hence there is a constant  $C$  such that  $|\Psi_c(x)| < C < \infty$  for all  $c$  and  $x$ . Moreover, it follows from (30) and (31) that

$$\Psi_c(x) \rightarrow \Psi(x), \quad c \rightarrow \infty,$$

where

$$\Psi(x) = \begin{cases} 0, & x < a, x > b, \\ \frac{1}{2}, & x = a, x = b, \\ 1, & a < x < b. \end{cases}$$

Let  $\mu$  be a measure on  $(R, \mathcal{B}(R))$  such that  $\mu(a, b] = F(b) - F(a)$ . Then if we apply the dominated convergence theorem and use the formulas of Problem 1 of §3, we find that, as  $c \rightarrow \infty$ ,

$$\begin{aligned}\Phi_c &= \int_{-\infty}^{\infty} \Psi_c(x) dF(x) \rightarrow \int_{-\infty}^{\infty} \Psi(x) dF(x) \\ &= \mu(a, b) + \frac{1}{2}\mu\{a\} + \frac{1}{2}\mu\{b\} \\ &= F(b-) - F(a) + \frac{1}{2}[F(a) - F(a-) + F(b) - F(b-)] \\ &= \frac{F(b) + F(b-)}{2} - \frac{F(a) + F(a-)}{2} = F(b) - F(a),\end{aligned}$$

where the last equation holds for all points  $a$  and  $b$  of continuity of  $F(x)$ .

Hence (25) is established.

(b) Let  $\int_{-\infty}^{\infty} |\varphi(t)| dt < \infty$ . Write

$$f(x) = \frac{1}{2\pi} \int_{-\infty}^{\infty} e^{-itx} \varphi(t) dt.$$

It follows from the dominated convergence theorem that this is a continuous function of  $x$  and therefore is integrable on  $[a, b]$ . Consequently we find, applying Fubini's theorem again, that

$$\begin{aligned}\int_a^b f(x) dx &= \int_a^b \frac{1}{2\pi} \left( \int_{-\infty}^{\infty} e^{-itx} \varphi(t) dt \right) dx \\ &= \frac{1}{2\pi} \int_{-\infty}^{\infty} \varphi(t) \left[ \int_a^b e^{-itx} dx \right] dt = \lim_{c \rightarrow \infty} \frac{1}{2\pi} \int_{-c}^c \varphi(t) \left[ \int_a^b e^{-itx} dx \right] dt \\ &= \lim_{c \rightarrow \infty} \frac{1}{2\pi} \int_{-c}^c \frac{e^{-ita} - e^{-itb}}{it} \varphi(t) dt = F(b) - F(a)\end{aligned}$$

for all points  $a$  and  $b$  of continuity of  $F(x)$ .

Hence it follows that

$$F(x) = \int_{-\infty}^x f(y) dy, \quad x \in R,$$

and since  $f(x)$  is continuous and  $F(x)$  is nondecreasing,  $f(x)$  is the density of  $F(x)$ .

This completes the proof of the theorem.

**Corollary.** *The inversion formula (25) provides a second proof of Theorem 2.*

**Theorem 4.** *A necessary and sufficient condition for the components of the random vector  $\xi = (\xi_1, \dots, \xi_n)$  to be independent is that its characteristic function is the product of the characteristic functions of the components:*

$$E e^{i(t_1 \xi_1 + \dots + t_n \xi_n)} = \prod_{k=1}^n E e^{it_k \xi_k}, \quad (t_1, \dots, t_n) \in R^n.$$

**PROOF.** The necessity follows from Problem 1. To prove the sufficiency we let  $F(x_1, \dots, x_n)$  be the distribution function of the vector  $\xi = (\xi_1, \dots, \xi_n)$  and  $F_k(x)$ , the distribution functions of the  $\xi_k$ ,  $1 \leq k \leq n$ . Put  $G = G(x_1, \dots, x_n) = F_1(x_1) \cdots F_n(x_n)$ . Then, by Fubini's theorem, for all  $(t_1, \dots, t_n) \in R^n$ ,

$$\begin{aligned}\int_{R^n} e^{i(t_1 x_1 + \dots + t_n x_n)} dG(x_1 \cdots x_n) &= \prod_{k=1}^n \int_R e^{it_k x_k} dF_k(x) \\ &= \prod_{k=1}^n E e^{it_k \xi_k} = E e^{i(t_1 \xi_1 + \dots + t_n \xi_n)} \\ &= \int_{R^n} e^{i(t_1 x_1 + \dots + t_n x_n)} dF(x_1 \cdots x_n).\end{aligned}$$

Therefore by Theorem 2 (or rather, by its multidimensional analog; see Problem 3) we have  $F = G$ , and consequently, by the theorem of §5, the random variables  $\xi_1, \dots, \xi_n$  are independent.

6. Theorem 1 gives us necessary conditions for a function to be a characteristic function. Hence if  $\varphi = \varphi(t)$  fails to satisfy, for example, one of the first three conclusions of the theorem, that function cannot be a characteristic function. We quote without proof some results in the same direction.

**Bochner-Khinchin Theorem.** Let  $\varphi(t)$  be continuous,  $t \in R$ , with  $\varphi(0) = 1$ . A necessary and sufficient condition that  $\varphi(t)$  is a characteristic function is that it is positive semi-definite, i.e. that for all real  $t_1, \dots, t_n$  and all complex  $\lambda_1, \dots, \lambda_n$ ,  $n = 1, 2, \dots$ ,

$$\sum_{i,j=1}^n \varphi(t_i - t_j) \lambda_i \bar{\lambda}_j \geq 0. \quad (32)$$

The necessity of (32) is evident since if  $\varphi(t) = \int_{-\infty}^{\infty} e^{itx} dF(x)$  then

$$\sum_{i,j=1}^n \varphi(t_i - t_j) \lambda_i \bar{\lambda}_j = \int_{-\infty}^{\infty} \left| \sum_{k=1}^n \lambda_k e^{it_k x} \right|^2 dF(x) \geq 0.$$

The proof of the sufficiency of (32) is more difficult.

**Pólya's Theorem.** Let a continuous even function  $\varphi(t)$  satisfy  $\varphi(t) \geq 0$ ,  $\varphi(0) = 1$ ,  $\varphi(t) \rightarrow 0$  as  $t \rightarrow \infty$  and let  $\varphi(t)$  be convex on  $0 \leq t < \infty$ . Then  $\varphi(t)$  is a characteristic function.

This theorem provides a very convenient method of constructing characteristic functions. Examples are

$$\begin{aligned} \varphi_1(t) &= e^{-|t|}, \\ \varphi_2(t) &= \begin{cases} 1 - |t|, & |t| \leq 1, \\ 0, & |t| > 1. \end{cases} \end{aligned}$$

Another is the function  $\varphi_3(t)$  drawn in Figure 34. On  $[-a, a]$ , the function  $\varphi_3(t)$  coincides with  $\varphi_2(t)$ . However, the corresponding distribution functions  $F_2$  and  $F_3$  are evidently different. This example shows that in general two characteristic functions can be the same on a finite interval without their distribution functions' being the same.

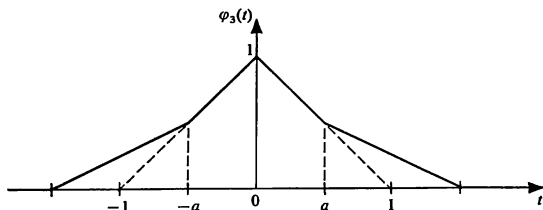


Figure 34

**Marcinkiewicz's Theorem.** *If a characteristic function  $\varphi(t)$  is of the form  $\exp \mathcal{P}(t)$ , where  $\mathcal{P}(t)$  is a polynomial, then this polynomial is of degree at most 2.*

It follows, for example, that  $e^{-t^4}$  is not a characteristic function.

7. The following theorem shows that a property of the characteristic function of a random variable can lead to a nontrivial conclusion about the nature of the random variable.

**Theorem 5.** *Let  $\varphi_\xi(t)$  be the characteristic function of the random variable  $\xi$ .*

(a) *If  $|\varphi_\xi(t_0)| = 1$  for some  $t_0 \neq 0$ , then  $\xi$  is concentrated at the points  $a + nh$ ,  $h = 2\pi/t_0$ , for some  $a$ , that is,*

$$\sum_{n=-\infty}^{\infty} \mathbf{P}\{\xi = a + nh\} = 1, \quad (33)$$

*where  $a$  is a constant.*

(b) *If  $|\varphi_\xi(t)| = |\varphi_\xi(\alpha t)| = 1$  for two different points  $t$  and  $\alpha t$ , where  $\alpha$  is irrational, then  $\xi$  is degenerate:*

$$\mathbf{P}\{\xi = a\} = 1,$$

*where  $a$  is some number.*

(c) *If  $|\varphi_\xi(t)| \equiv 1$ , then  $\xi$  is degenerate.*

**PROOF.** (a) If  $|\varphi_\xi(t_0)| = 1$ ,  $t_0 \neq 0$ , there is a number  $a$  such that  $\varphi(t_0) = e^{it_0 a}$ . Then

$$\begin{aligned} e^{it_0 a} &= \int_{-\infty}^{\infty} e^{it_0 x} dF(x) \Rightarrow 1 = \int_{-\infty}^{\infty} e^{it_0(x-a)} dF(x) \Rightarrow \\ 1 &= \int_{-\infty}^{\infty} \cos t_0(x-a) dF(x) \Rightarrow \int_{-\infty}^{\infty} [1 - \cos t_0(x-a)] dF(x) = 0. \end{aligned}$$

Since  $1 - \cos t_0(x-a) \geq 0$ , it follows from property **H** (Subsection 2 of §6) that

$$1 = \cos t_0(\xi - a) \quad (\text{P-a.s.}),$$

which is equivalent to (33).

(b) It follows from  $|\varphi_\xi(t)| = |\varphi_\xi(\alpha t)| = 1$  and from (33) that

$$\sum_{n=-\infty}^{\infty} \mathbf{P}\left\{\xi = a + \frac{2\pi}{t} n\right\} = \sum_{m=-\infty}^{\infty} \mathbf{P}\left\{\xi = b + \frac{2\pi}{\alpha t} m\right\} = 1.$$

If  $\xi$  is not degenerate, there must be at least two pairs of common points:

$$a + \frac{2\pi}{t} n_1 = b + \frac{2\pi}{\alpha t} m_1, \quad a + \frac{2\pi}{y} n_2 = b + \frac{2\pi}{\alpha t} m_2,$$

in the sets

$$\left\{a + \frac{2\pi}{t}n, n = 0, \pm 1, \dots\right\} \quad \text{and} \quad \left\{b + \frac{2\pi}{\alpha t}m, m = 0, \pm 1, \dots\right\},$$

whence

$$\frac{2\pi}{t}(n_1 - n_2) = \frac{2\pi}{\alpha t}(m_1 - m_2),$$

and this contradicts the assumption that  $\alpha$  is irrational. Conclusion (c) follows from (b).

This completes the proof of the theorem.

8. Let  $\xi = (\xi_1, \dots, \xi_k)$  be a random vector,

$$\varphi_\xi(t) = E e^{i(t, \xi)}, \quad t = (t_1, \dots, t_k),$$

its characteristic function. Let us suppose that  $E|\xi_i|^n < \infty$  for some  $n \geq 1$ ,  $i = 1, \dots, k$ . From the inequalities of Hölder (6.29) and Lyapunov (6.27) it follows that the (mixed) moments  $E(\xi_1^{v_1} \dots \xi_k^{v_k})$  exist for all nonnegative  $v_1, \dots, v_k$  such that  $v_1 + \dots + v_k \leq n$ .

As in Theorem 1, this implies the existence and continuity of the partial derivatives

$$\frac{\partial^{v_1 + \dots + v_k}}{\partial t_1^{v_1} \dots \partial t_k^{v_k}} \varphi_\xi(t_1, \dots, t_k)$$

for  $v_1 + \dots + v_k \leq n$ . Then if we expand  $\varphi_\xi(t_1, \dots, t_k)$  in a Taylor series, we see that

$$\varphi_\xi(t_1, \dots, t_k) = \sum_{v_1 + \dots + v_k \leq n} \frac{i^{v_1 + \dots + v_k}}{v_1! \dots v_k!} m_\xi^{(v_1, \dots, v_k)} t_1^{v_1} \dots t_k^{v_k} + o(|t|^n), \quad (34)$$

where  $|t| = |t_1| + \dots + |t_k|$  and

$$m_\xi^{(v_1, \dots, v_k)} = E \xi_1^{v_1} \dots \xi_k^{v_k}$$

is the *mixed moment of order*  $v = (v_1, \dots, v_k)$ .

Now  $\varphi_\xi(t_1, \dots, t_k)$  is continuous,  $\varphi_\xi(0, \dots, 0) = 1$ , and consequently this function is different from zero in some neighborhood  $|t| < \delta$  of zero. In this neighborhood the partial derivative

$$\frac{\partial^{v_1 + \dots + v_k}}{\partial t_1^{v_1} \dots \partial t_k^{v_k}} \ln \varphi_\xi(t_1, \dots, t_k)$$

exists and is continuous, where  $\ln z$  denotes the principal value of the logarithm (if  $z = re^{i\theta}$ , we take  $\ln z$  to be  $\ln r + i\theta$ ). Hence we can expand  $\ln \varphi_\xi(t_1, \dots, t_k)$  by Taylor's formula,

$$\ln \varphi_\xi(t_1, \dots, t_k) = \sum_{v_1 + \dots + v_k \leq n} \frac{i^{v_1 + \dots + v_k}}{v_1! \dots v_k!} S_\xi^{(v_1, \dots, v_k)} t_1^{v_1} \dots t_k^{v_k} + o(|t|^n), \quad (35)$$

where the coefficients  $s_{\xi}^{(v_1, \dots, v_k)}$  are the (mixed) semi-invariants or cumulants of order  $v = v(v_1, \dots, v_k)$  of  $\xi = \xi_1, \dots, \xi_k$ .

Observe that if  $\xi$  and  $\eta$  are independent, then

$$\ln \varphi_{\xi+\eta}(t) = \ln \varphi_{\xi}(t) + \ln \varphi_{\eta}(t), \quad (36)$$

and therefore

$$s_{\xi+\eta}^{(v_1, \dots, v_k)} = s_{\xi}^{(v_1, \dots, v_k)} + s_{\eta}^{(v_1, \dots, v_k)}. \quad (37)$$

(It is this property that gives rise to the term "semi-invariant" for  $s_{\xi}^{(v_1, \dots, v_k)}$ .)

To simplify the formulas and make (34) and (35) look "one-dimensional," we introduce the following notation.

If  $v = (v_1, \dots, v_k)$  is a vector whose components are nonnegative integers, we put

$$v! = v_1! \cdots v_k!, \quad |v| = v_1 + \cdots + v_k, \quad t^v = t_1^{v_1} \cdots t_k^{v_k}.$$

$$\text{We also put } s_{\xi}^{(v)} = s_{\xi}^{(v_1, \dots, v_k)}, \quad m_{\xi}^{(v)} = m_{\xi}^{(v_1, \dots, v_k)}.$$

Then (34) and (35) can be written

$$\varphi_{\xi}(t) = \sum_{|v| \leq n} \frac{i^{|v|}}{v!} m_{\xi}^{(v)} t^v + o(|t|^n), \quad (38)$$

$$\ln \varphi_{\xi}(t) = \sum_{|v| \leq n} \frac{i^{|v|}}{v!} s_{\xi}^{(v)} t^v + o(|t|^n). \quad (39)$$

The following theorem and its corollaries give formulas that connect moments and semi-invariants.

**Theorem 6.** Let  $\xi = (\xi_1, \dots, \xi_k)$  be a random vector with  $E|\xi_i|^n < \infty$ ,  $i = 1, \dots, k$ ,  $n \geq 1$ . Then for  $v = (v_1, \dots, v_k)$  such that  $|v| \leq n$

$$m_{\xi}^{(v)} = \sum_{\lambda^{(1)} + \dots + \lambda^{(q)} = v} \frac{1}{q!} \frac{v!}{\lambda^{(1)}! \cdots \lambda^{(q)}!} \prod_{p=1}^q s^{(\lambda^{(p)})}, \quad (40)$$

$$s_{\xi}^{(v)} = \sum_{\lambda^{(1)} + \dots + \lambda^{(q)} = v} \frac{(-1)^{q-1}}{q} \frac{v!}{\lambda^{(1)}! \cdots \lambda^{(q)}!} \prod_{p=1}^q m_{\xi}^{(\lambda^{(p)})}, \quad (41)$$

where  $\sum_{\lambda^{(1)} + \dots + \lambda^{(q)} = v}$  indicates summation over all ordered sets of nonnegative integral vectors  $\lambda^{(p)}$ ,  $|\lambda^{(p)}| > 0$ , whose sum is  $v$ .

**PROOF.** Since

$$\varphi_{\xi}(t) = \exp(\ln \varphi_{\xi}(t)),$$

if we expand the function  $\exp$  by Taylor's formula and use (39), we obtain

$$\varphi_{\xi}(t) = 1 + \sum_{q=1}^n \frac{1}{q!} \left( \sum_{1 \leq |\lambda| \leq n} \frac{i^{|\lambda|}}{\lambda!} s_{\xi}^{(\lambda)} t^{\lambda} \right)^q + o(|t|^n). \quad (42)$$

Comparing terms in  $t^\lambda$  on the right-hand sides of (38) and (42), and using  $|\lambda^{(1)}| + \dots + |\lambda^{(q)}| = |\lambda^{(1)} + \dots + \lambda^{(q)}|$ , we obtain (40).

Moreover,

$$\ln \varphi_\xi(t) = \ln \left[ 1 + \sum_{1 \leq |\lambda| \leq n} \frac{i^{|\lambda|}}{\lambda!} m_\xi^{(\lambda)} t^\lambda + o(|t|^n) \right]. \quad (43)$$

For small  $z$  we have the expansion

$$\ln(1 + z) = \sum_{q=1}^n \frac{(-1)^{q-1}}{q} z^q + o(z^q).$$

Using this in (43) and then comparing the coefficients of  $t^\lambda$  with the corresponding coefficients on the right-hand side of (38), we obtain (41).

**Corollary 1.** *The following formulas connect moments and semi-invariants:*

$$m_\xi^{(v)} = \sum_{\{r_1 \lambda^{(1)} + \dots + r_x \lambda^{(x)} = v\}} \frac{1}{r_1! \dots r_x!} \frac{v!}{(\lambda^{(1)})!^{r_1} \dots (\lambda^{(x)})!^{r_x}} \prod_{j=1}^x [s_\xi^{(\lambda^{(j)})}]^{r_j}, \quad (44)$$

$$s_\xi^{(v)} = \sum_{\{r_1 \lambda^{(1)} + \dots + r_x \lambda^{(x)} = v\}} \frac{(-1)^{q-1} (q-1)!}{r_1! \dots r_x!} \frac{v!}{(\lambda^{(1)})!^{r_1} \dots (\lambda^{(x)})!^{r_x}} \prod_{j=1}^x [m_\xi^{(\lambda^{(j)})}]^{r_j}, \quad (45)$$

where  $\sum_{\{r_1 \lambda^{(1)} + \dots + r_x \lambda^{(x)} = v\}}$  denotes summation over all unordered sets of different nonnegative integral vectors  $\lambda^{(j)}$ ,  $|\lambda^{(j)}| > 0$ , and over all ordered sets of positive integral numbers  $r_j$  such that  $r_1 \lambda^{(1)} + \dots + r_x \lambda^{(x)} = v$ .

To establish (44) we suppose that among all the vectors  $\lambda^{(1)}, \dots, \lambda^{(q)}$  that occur in (40), there are  $r_1$  equal to  $\lambda^{(i_1)}, \dots, r_x$  equal to  $\lambda^{(i_x)}$  ( $r_j > 0$ ,  $r_1 + \dots + r_x = q$ ), where all the  $\lambda^{(i_j)}$  are different. There are  $q!/(r_1! \dots r_x!)$  different sets of vectors, corresponding (except for order) with the set  $\{\lambda^{(1)}, \dots, \lambda^{(q)}\}$ . But if two sets, say,  $\{\lambda^{(1)}, \dots, \lambda^{(q)}\}$  and  $\{\bar{\lambda}^{(1)}, \dots, \bar{\lambda}^{(q)}\}$  differ only in order, then  $\prod_{p=1}^q s_\xi^{(\lambda^{(p)})} = \prod_{p=1}^q s_\xi^{(\bar{\lambda}^{(p)})}$ . Hence if we identify sets that differ only in order, we obtain (44) from (40).

Formula (45) can be deduced from (41) in a similar way.

**Corollary 2.** *Let us consider the special case when  $v = (1, \dots, 1)$ . In this case the moments  $m_\xi^{(v)} \equiv E \xi_1 \dots \xi_k$ , and the corresponding semi-invariants, are called simple.*

Formulas connecting simple moments and simple semi-invariants can be read off from the formulas given above. However, it is useful to have them written in a different way.

For this purpose, we introduce the following notation.

Let  $\xi = (\xi_1, \dots, \xi_k)$  be a vector, and  $I_\xi = \{1, 2, \dots, k\}$  its set of indices. If  $I \subseteq I_\xi$ , let  $\xi_I$  denote the vector consisting of the components of  $\xi$  whose



indices belong to  $I$ . Let  $\chi(I)$  be the vector  $\{\chi_1, \dots, \chi_n\}$  for which  $\chi_i = 1$  if  $i \in I$ , and  $\chi_i = 0$  if  $i \notin I$ . These vectors are in one-to-one correspondence with the sets  $I \subseteq I_\xi$ . Hence we can write

$$m_\xi(I) = m_\xi^{\chi(I)}, \quad s_\xi(I) = s_\xi^{\chi(I)}.$$

In other words,  $m_\xi(I)$  and  $s_\xi(I)$  are simple moments and semi-invariants of the subvector  $\xi_I$  of  $\xi$ .

In accordance with the definition given on p. 12, a *decomposition* of a set  $I$  is an unordered collection of disjoint nonempty sets  $I_p$  such that  $\sum_p I_p = I$ .

In terms of these definitions, we have the formulas

$$m_\xi(I) = \sum_{\sum_{p=1}^q I_p = I} \prod_{p=1}^q s_\xi(I_p), \quad (46)$$

$$s_\xi(I) = \sum_{\sum_{p=1}^q I_p = I} (-1)^{q-1} (q-1)! \prod_{p=1}^q m_\xi(I_p). \quad (47)$$

where  $\sum_{\sum_{p=1}^q I_p = I}$  denotes summation over all decompositions of  $I$ ,  $1 \leq q \leq N(I)$ , where  $N(I)$  is the number of elements of the set  $I$ .

We shall derive (46) from (44). If  $v = \chi(I)$  and  $\lambda^{(1)} + \dots + \lambda^{(q)} = v$ , then  $\lambda^{(p)} = \chi(I_p)$ ,  $I_p \subseteq I$ , where the  $\lambda^{(p)}$  are all different,  $\lambda^{(p)}! = v! = 1$ , and every unordered set  $\{\chi(I_1), \dots, \chi(I_q)\}$  is in one-to-one correspondence with the decomposition  $I = \sum_{p=1}^q I_p$ . Consequently (46) follows from (44).

In a similar way, (47) follows from (35).

**EXAMPLE 1.** Let  $\xi$  be a random variable ( $k = 1$ ) and  $m_n = m_\xi^{(n)} = E\xi^n$ ,  $s_n = s_\xi^{(n)}$ . Then (40) and (41) imply the following formulas:

$$\begin{aligned} m_1 &= s_1, \\ m_2 &= s_2 + s_1^2, \\ m_3 &= s_3 + 3s_1s_2 + s_1^3, \\ m_4 &= s_4 + 3s_2^2 + 4s_1s_3 + 6s_1^2s_2 + s_1^4, \\ &\dots \end{aligned} \quad (48)$$

and

$$\begin{aligned} s_1 &= m_1 = E\xi, \\ s_2 &= m_2 - m_1^2 = V\xi, \\ s_3 &= m_3 - 3m_1m_2 + 2m_1^3, \\ s_4 &= m_4 - 3m_2^2 - 4m_1m_3 + 12m_1^2m_2 - 6m_1^4, \\ &\dots \end{aligned} \quad (49)$$

EXAMPLE 2. Let  $\xi \sim \mathcal{N}(m, \sigma^2)$ . Since, by (9),

$$\ln \varphi_\xi(t) = itm - \frac{t^2 \sigma^2}{2},$$

we have  $s_1 = m$ ,  $s_2 = \sigma^2$  by (39), and all the semi-invariants, from the third on, are zero:  $s_n = 0$ ,  $n \geq 3$ .

We may observe that by Marcinkiewicz's theorem a function  $\exp \mathcal{P}(t)$ , where  $\mathcal{P}$  is a polynomial, can be a characteristic function only when the degree of that polynomial is at most 2. It follows, in particular, that the Gaussian distribution is the only distribution with the property that all its semi-invariants  $s_n$  are zero from a certain index onward.

EXAMPLE 3. If  $\xi$  is a Poisson random variable with parameter  $\lambda > 0$ , then by (11)

$$\ln \varphi_\xi(t) = \lambda(e^{it} - 1).$$

It follows that

$$s_n = \lambda \quad (50)$$

for all  $n \geq 1$ .

EXAMPLE 4. Let  $\xi = (\xi_1, \dots, \xi_n)$  be a random vector. Then

$$\begin{aligned} m_\xi(1) &= s_\xi(1), \\ m_\xi(1, 2) &= s_\xi(1, 2) + s_\xi(1)s_\xi(2), \\ m_\xi(1, 2, 3) &= s_\xi(1, 2, 3) + s_\xi(1, 2)s_\xi(3) + \\ &\quad + s_\xi(1, 3)s_\xi(2) + \\ &\quad + s_\xi(2, 3)s_\xi(1) + s_\xi(1)s_\xi(2)s_\xi(3) \\ &\dots\dots\dots \end{aligned} \quad (51)$$

These formulas show that the simple moments can be expressed in terms of the simple semi-invariants in a very *symmetric* way. If we put  $\xi_1 \equiv \xi_2 \equiv \dots \equiv \xi_k$ , we then, of course, obtain (48).

The group-theoretical origin of the coefficients in (48) becomes clear from (51). It also follows from (51) that

$$s_\xi(1, 2) = m_\xi(1, 2) - m_\xi(1)m_\xi(2) = E\xi_1\xi_2 - E\xi_1 E\xi_2, \quad (52)$$

i.e.,  $s_\xi(1, 2)$  is just the *covariance* of  $\xi_1$  and  $\xi_2$ .

9. Let  $\xi$  be a random variable with distribution function  $F = F(x)$  and characteristic function  $\varphi(t)$ . Let us suppose that all the moments  $m_n = E\xi^n$ ,  $n \geq 1$ , exist.

It follows from Theorem 2 that a characteristic function uniquely determines a probability distribution. Let us now ask the following question

(uniqueness for the moment problem): Do the moments  $\{m_n\}_{n \geq 1}$  determine the *probability distribution*?

More precisely, let  $F$  and  $G$  be distribution functions with the same moments, i.e.

$$\int_{-\infty}^{\infty} x^n dF(x) = \int_{-\infty}^{\infty} x^n dG(x) \quad (53)$$

for all integers  $n \geq 0$ . The question is whether  $F$  and  $G$  must be the same.

In general, the answer is "no." To see this, consider the distribution  $F$  with density

$$f(x) = \begin{cases} ke^{-\alpha x^\lambda}, & x > 0, \\ 0, & x \leq 0, \end{cases}$$

where  $\alpha > 0$ ,  $0 < \lambda < \frac{1}{2}$ , and  $k$  is determined by the condition  $\int_0^\infty f(x) dx = 1$ .

Write  $\beta = \alpha \tan \lambda\pi$  and let  $g(x) = 0$  for  $x \leq 0$  and

$$g(x) = ke^{-\alpha x^\lambda} [1 + \varepsilon \sin(\beta x^\lambda)], \quad |\varepsilon| < 1, \quad x > 0.$$

It is evident that  $g(x) \geq 0$ . Let us show that

$$\int_0^\infty x^n e^{-\alpha x^\lambda} \sin \beta x^\lambda dx = 0 \quad (54)$$

for all integers  $n \geq 0$ .

For  $p > 0$  and complex  $q$  with  $\operatorname{Re} q > 0$ , we have

$$\int_0^\infty t^{p-1} e^{-qt} dt = \frac{\Gamma(p)}{q^p}.$$

Take  $p = (n+1)/\lambda$ ,  $q = \alpha + i\beta$ ,  $t = x^\lambda$ . Then

$$\begin{aligned} \int_0^\infty x^{\lambda[(n+1)/\lambda]-1} e^{-(\alpha+i\beta)x^\lambda} \lambda x^{\lambda-1} dx &= \lambda \int_0^\infty x^n e^{-(\alpha+i\beta)x^\lambda} dx \\ &= \lambda \int_0^\infty x^n e^{-\alpha x^\lambda} \cos \beta x^\lambda dx - i\lambda \int_0^\infty x^n e^{-\alpha x^\lambda} \sin \beta x^\lambda dx \\ &= \frac{\Gamma\left(\frac{n+1}{\lambda}\right)}{\alpha^{(n+1)/\lambda} (1 + i \tan \lambda\pi)^{(n+1)/\lambda}}. \end{aligned} \quad (55)$$

But

$$\begin{aligned} (1 + i \tan \lambda\pi)^{(n+1)/\lambda} &= (\cos \lambda\pi + i \sin \lambda\pi)^{(n+1)/\lambda} (\cos \lambda\pi)^{-(n+1)/\lambda} \\ &= e^{i\pi(n+1)} (\cos \lambda\pi)^{-(n+1)/\lambda} \\ &= \cos \pi(n+1) \cdot \cos(\lambda\pi)^{-(n+1)/\lambda}, \end{aligned}$$

since  $\sin \pi(n+1) = 0$ .

Hence right-hand side of (55) is real and therefore (54) is valid for all integral  $n \geq 0$ . Now let  $G(x)$  be the distribution function with density  $g(x)$ . It follows from (54) that the distribution functions  $F$  and  $G$  have equal moments, i.e. (53) holds for all integers  $n \geq 0$ .

We now give some conditions that guarantee the uniqueness of the solution of the moment problem.

**Theorem 7.** Let  $F = F(x)$  be a distribution function and  $\mu_n = \int_{-\infty}^{\infty} |x|^n dF(x)$ . If

$$\overline{\lim}_{n \rightarrow \infty} \frac{\mu_n^{1/n}}{n} < \infty, \quad (56)$$

the moments  $\{m_n\}_{n \geq 1}$ , where  $m_n = \int_{-\infty}^{\infty} x^n dF(x)$ , determine the distribution  $F = F(x)$  uniquely.

**PROOF.** It follows from (56) and conclusion (7) of Theorem 1 that there is a  $t_0 > 0$  such that, for all  $|t| \leq t_0$ , the characteristic function

$$\varphi(t) = \int_{-\infty}^{\infty} e^{itx} dF(x)$$

can be represented in the form

$$\varphi(t) = \sum_{k=0}^{\infty} \frac{(it)^k}{k!} m_k$$

and consequently the moments  $\{m_n\}_{n \geq 1}$  uniquely determine the characteristic function  $\varphi(t)$  for  $|t| \leq t_0$ .

Take a point  $s$  with  $|s| \leq t_0/2$ . Then, as in the proof of (15), we deduce from (56) that

$$\varphi(t) = \sum_{k=0}^{\infty} \frac{i^k (t-s)^k}{k!} \varphi^{(k)}(s)$$

for  $|t-s| \leq t_0$ , where

$$\varphi^{(k)}(s) = i^k \int_{-\infty}^{\infty} x^k e^{isx} dF(x)$$

is uniquely determined by the moments  $\{m_n\}_{n \geq 1}$ . Consequently the moments determine  $\varphi(t)$  uniquely for  $|t| \leq \frac{3}{2}t_0$ . Continuing this process, we see that  $\{m_n\}_{n \geq 1}$  determines  $\varphi(t)$  uniquely for all  $t$ , and therefore also determines  $F(x)$ .

This completes the proof of the theorem.

**Corollary 1.** The moments completely determine the probability distribution if it is concentrated on a finite interval.

**Corollary 2.** *A sufficient condition for the moment problem to have a unique solution is that*

$$\overline{\lim}_{n \rightarrow \infty} \frac{(m_{2n})^{1/2n}}{2n} < \infty. \quad (57)$$

For the proof it is enough to observe that the odd moments can be estimated in terms of the even ones, and then use (56).

**EXAMPLE.** Let  $F(x)$  be the normal distribution function,

$$F(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \int_{-\infty}^x e^{-t^2/2\sigma^2} dt.$$

Then  $m_{2n+1} = 0$ ,  $m_{2n} = [(2n)!/2^n n!] \sigma^{2n}$ , and it follows from (57) that these are the moments only of the normal distribution.

Finally we state, without proof:

Carleman's test for the uniqueness of the moment problem.

(a) Let  $\{m_n\}_{n \geq 1}$  be the moments of a probability distribution, and let

$$\sum_{n=0}^{\infty} \frac{1}{(m_{2n})^{1/2n}} = \infty.$$

Then they determine the probability distribution uniquely.

(b) If  $\{m_n\}_{n \geq 1}$  are the moments of a distribution that is concentrated on  $[0, \infty)$ , then the solution will be unique if we require only that

$$\sum_{n=0}^{\infty} \frac{1}{(m_n)^{1/2n}} = \infty.$$

**10.** Let  $F = F(x)$  and  $G = G(x)$  be distribution functions with characteristic functions  $f = f(t)$  and  $g = g(t)$ , respectively. The following theorem, which we give without proof, makes it possible to estimate how close  $F$  and  $G$  are to each other (in the uniform metric) in terms of the closeness of  $f$  and  $g$ .

**Theorem (Esseen's Inequality).** Let  $G(x)$  have derivative  $G'(x)$  with  $\sup |G'(x)| \leq C$ . Then for every  $T > 0$

$$\sup_x |F(x) - G(x)| \leq \frac{2}{\pi} \int_0^T \left| \frac{f(t) - g(t)}{t} \right| dt + \frac{24}{\pi T} \sup_x |G'(x)|. \quad (58)$$

(This will be used in §6 of Chapter III to prove a theorem on the rapidity of convergence in the central limit theorem.)

## 11. PROBLEMS

1. Let  $\xi$  and  $\eta$  be independent random variables,  $f(x) = f_1(x) + if_2(x)$ ,  $g(x) = g_1(x) + ig_2(x)$ , where  $f_k(x)$  and  $g_k(x)$  are Borel functions,  $k = 1, 2$ . Show that if  $E|f(\xi)| < \infty$  and  $E|g(\eta)| < \infty$ , then

$$E|f(\xi)g(\eta)| < \infty$$

and

$$Ef(\xi)g(\eta) = Ef(\xi) \cdot Eg(\eta).$$

2. Let  $\xi = (\xi_1, \dots, \xi_n)$  and  $E\|\xi\|^n < \infty$ , where  $\|\xi\| = +\sqrt{\sum \xi_i^2}$ . Show that

$$\varphi_\xi(t) = \sum_{k=0}^n \frac{i^k}{k!} E(t, \xi)^k + \varepsilon_n(t) \|t\|^n,$$

where  $t = (t_1, \dots, t_n)$  and  $\varepsilon_n(t) \rightarrow 0$ ,  $t \rightarrow 0$ .

3. Prove Theorem 2 for  $n$ -dimensional distribution functions  $F = F_n(x_1, \dots, x_n)$  and  $G_n(x_1, \dots, x_n)$ .
4. Let  $F = F(x_1, \dots, x_n)$  be an  $n$ -dimensional distribution function and  $\varphi = \varphi(t_1, \dots, t_n)$  its characteristic function. Using the notation of (3.12), establish the inversion formula

$$P(a, b] = \lim_{c \rightarrow \infty} \frac{1}{(2\pi)^n} \int_{-c}^c \prod_{k=1}^n \frac{e^{it_k a_k} - e^{it_k b_k}}{it_k} \varphi(t_1, \dots, t_k) dt_1 \cdots dt_k.$$

(We are to suppose that  $(a, b]$  is an interval of continuity of  $P(a, b]$ , i.e. for  $k = 1, \dots, n$  the points  $a_k, b_k$  are points of continuity of the marginal distribution functions  $F_k(x_k)$  which are obtained from  $F(x_1, \dots, x_n)$  by taking all the variables except  $x_k$  equal to  $+\infty$ .)

5. Let  $\varphi_k(t)$ ,  $k \geq 1$ , be a characteristic function, and let the nonnegative numbers  $\lambda_k$ ,  $k \geq 1$ , satisfy  $\sum \lambda_k = 1$ . Show that  $\sum \lambda_k \varphi_k(t)$  is a characteristic function.
6. If  $\varphi(t)$  is a characteristic function, are  $\operatorname{Re} \varphi(t)$  and  $\operatorname{Im} \varphi(t)$  characteristic functions?
7. Let  $\varphi_1, \varphi_2$  and  $\varphi_3$  be characteristic functions, and  $\varphi_1 \varphi_2 = \varphi_1 \varphi_3$ . Does it follow that  $\varphi_2 = \varphi_3$ ?
8. Construct the characteristic functions of the distributions given in Tables 1 and 2 of §3.
9. Let  $\xi$  be an integral-valued random variable and  $\varphi_\xi(t)$  its characteristic function. Show that

$$P(\xi = k) = \frac{1}{2\pi} \int_{-\pi}^{\pi} e^{-ikt} \varphi_\xi(t) dt, \quad k = 0, \pm 1, \pm 2, \dots$$

## §13. Gaussian Systems

1. Gaussian, or normal, distributions, random variables, processes, and systems play an extremely important role in probability theory and in mathematical statistics. This is explained in the first instance by the central

limit theorem (§4 of Chapter III and §8 of Chapter VII), of which the De Moivre–Laplace limit theorem is a special case (§6, Chapter I). According to this theorem, the normal distribution is universal in the sense that the distribution of the sum of a large number of random variables or random vectors, subject to some not very restrictive conditions, is closely approximated by this distribution.

This is what provides a theoretical explanation of the “law of errors” of applied statistics, which says that errors of measurement that result from large numbers of independent “elementary” errors obey the normal distribution.

A multidimensional Gaussian distribution is specified by a small number of parameters; this is a definite advantage in using it in the construction of simple probabilistic models. Gaussian random variables have finite second moments, and consequently they can be studied by Hilbert space methods. Here it is important that in the Gaussian case “uncorrelated” is equivalent to “independent,” so that the results of  $L^2$ -theory can be significantly strengthened.

2. Let us recall that (see §8) a random variable  $\xi = \xi(\omega)$  is Gaussian, or normally distributed, with parameters  $m$  and  $\sigma^2$  ( $\xi \sim \mathcal{N}(m, \sigma^2)$ ),  $|m| < \infty$ ,  $\sigma^2 > 0$ , if its density  $f_\xi(x)$  has the form

$$f_\xi(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-(x-m)^2/2\sigma^2}, \quad (1)$$

where  $\sigma = +\sqrt{\sigma^2}$ .

As  $\sigma \downarrow 0$ , the density  $f_\xi(x)$  “converges to the  $\delta$ -function supported at  $x = m$ .” It is natural to say that  $\xi$  is normally distributed with mean  $m$  and  $\sigma^2 = 0$  ( $\xi \sim \mathcal{N}(m, 0)$ ) if  $\xi$  has the property that  $P(\xi = m) = 1$ .

We can, however, give a definition that applies both to the *nondegenerate* ( $\sigma^2 > 0$ ) and the *degenerate* ( $\sigma^2 = 0$ ) cases. Let us consider the characteristic function  $\varphi_\xi(t) \equiv E e^{it\xi}$ ,  $t \in R$ .

If  $P(\xi = m) = 1$ , then evidently

$$\varphi_\xi(t) = e^{itm}, \quad (2)$$

whereas if  $\xi \sim \mathcal{N}(m, \sigma^2)$ ,  $\sigma^2 > 0$ ,

$$\varphi_\xi(t) = e^{itm - (1/2)t^2\sigma^2}. \quad (3)$$

It is obvious that when  $\sigma^2 = 0$  the right-hand sides of (2) and (3) are the same. It follows, by Theorem 2 of §12, that the Gaussian random variable with parameters  $m$  and  $\sigma^2$  ( $|m| < \infty$ ,  $\sigma^2 \geq 0$ ) must be the same as the random variable whose characteristic function is given by (3). This is an illustration of the “attraction of characteristic functions,” a very useful technique in the multidimensional case.

Let  $\xi = (\xi_1, \dots, \xi_n)$  be a random vector and

$$\varphi_\xi(t) = E e^{i(t, \xi)}, \quad t = (t_1, \dots, t_n) \in R^n, \quad (4)$$

its characteristic function (see Definition 2, §12).

**Definition 1.** A random vector  $\xi = (\xi_1, \dots, \xi_n)$  is *Gaussian*, or *normally distributed*, if its characteristic function has the form

$$\varphi_\xi(t) = e^{i(t, m) - (1/2)(\mathbb{R}t, t)}, \quad (5)$$

where  $m = (m_1, \dots, m_n)$ ,  $|m_k| < \infty$  and  $\mathbb{R} = \|r_{ki}\|$  is a symmetric nonnegative definite  $n \times n$  matrix; we use the abbreviation  $\xi \sim \mathcal{N}(m, \mathbb{R})$ .

This definition immediately makes us ask whether (5) is in fact a characteristic function. Let us show that it is.

First suppose that  $\mathbb{R}$  is nonsingular. Then we can define the inverse  $A = \mathbb{R}^{-1}$  and the function

$$f(x) = \frac{|A|^{1/2}}{(2\pi)^{n/2}} \exp\{-\frac{1}{2}(A(x - m), (x - m))\}, \quad (6)$$

where  $x = (x_1, \dots, x_n)$  and  $|A| = \det A$ . This function is nonnegative. Let us show that

$$\int_{R^n} e^{i(t, x)} f(x) dx = e^{i(t, m) - (1/2)(\mathbb{R}t, t)},$$

or equivalently that

$$I_n \equiv \int_{R^n} e^{i(t, x - m)} \frac{|A|^{1/2}}{(2\pi)^{n/2}} e^{-(1/2)(A(x - m), (x - m))} dx = e^{-(1/2)(\mathbb{R}t, t)}. \quad (7)$$

Let us make the change of variable

$$x - m = \mathcal{O}u, \quad t = \mathcal{O}v,$$

where  $\mathcal{O}$  is an orthogonal matrix such that

$$\mathcal{O}^T \mathbb{R} \mathcal{O} = D,$$

and

$$D = \begin{pmatrix} d_1 & & 0 \\ & \ddots & \\ 0 & & d_n \end{pmatrix}$$

is a diagonal matrix with  $d_i \geq 0$  (see the proof of the lemma in §8). Since  $|\mathbb{R}| = \det \mathbb{R} \neq 0$ , we have  $d_i > 0$ ,  $i = 1, \dots, n$ . Therefore

$$|A| = |\mathbb{R}^{-1}| = d_1^{-1} \dots d_n^{-1}. \quad (8)$$



Moreover (for notation, see Subsection 1, §12)

$$\begin{aligned} i(t, x - m) - \frac{1}{2}(A(x - m), x - m) &= i(\mathcal{O}v, \mathcal{O}u) - \frac{1}{2}(A\mathcal{O}u, \mathcal{O}u) \\ &= i(\mathcal{O}v)^T \mathcal{O}u - \frac{1}{2}(\mathcal{O}u)^T A(\mathcal{O}u) \\ &= iv^T u - \frac{1}{2}u^T \mathcal{O}^T A \mathcal{O}u \\ &= iv^T u - \frac{1}{2}u^T D^{-1}u. \end{aligned}$$

Together with (8) and (12.9), this yields

$$\begin{aligned} I_n &= (2\pi)^{-n/2} (d_1 \cdots d_n)^{-1/2} \int_{\mathbb{R}^n} \exp(iv^T u - \frac{1}{2}u^T D^{-1}u) du \\ &= \prod_{k=1}^n (2\pi d_k)^{-1/2} \int_{-\infty}^{\infty} \exp\left(iv_k u_k - \frac{u_k^2}{2d_k}\right) du_k = \prod_{k=1}^n \exp(-\frac{1}{2}v_k^2 d_k) \\ &= \exp(-\frac{1}{2}v^T D v) = \exp(-\frac{1}{2}v^T \mathcal{O}^T \mathbb{R} \mathcal{O} v) = \exp(-\frac{1}{2}t^T \mathbb{R} t) = \exp(-\frac{1}{2}(\mathbb{R}t, t)). \end{aligned}$$

It also follows from (6) that

$$\int_{\mathbb{R}^n} f(x) dx = 1. \quad (9)$$

Therefore (5) is the characteristic function of a nondegenerate  $n$ -dimensional Gaussian distribution (see Subsection 3, §3).

Now let  $\mathbb{R}$  be singular. Take  $\varepsilon > 0$  and consider the positive definite symmetric matrix  $\mathbb{R}^\varepsilon \equiv \mathbb{R} + \varepsilon E$ . Then by what has been proved,

$$\varphi^\varepsilon(t) = \exp\{i(t, m) - \frac{1}{2}(\mathbb{R}^\varepsilon t, t)\}$$

is a characteristic function:

$$\varphi^\varepsilon(t) = \int_{\mathbb{R}^n} e^{i(t, x)} dF_\varepsilon(x),$$

where  $F_\varepsilon(x) = F_\varepsilon(x_1, \dots, x_n)$  is an  $n$ -dimensional distribution function.

As  $\varepsilon \rightarrow 0$ ,

$$\varphi^\varepsilon(t) \rightarrow \varphi(t) = \exp\{i(t, m) - \frac{1}{2}(\mathbb{R}^\varepsilon t, t)\}.$$

The limit function  $\varphi(t)$  is continuous at  $(0, \dots, 0)$ . Hence, by Theorem 1 and Problem 1 of §3 of Chapter III, it is a characteristic function.

We have therefore established Theorem 1.

**3.** Let us now discuss the significance of the vector  $m$  and the matrix  $\mathbb{R} = \|r_{kl}\|$  that appear in (5).

Since

$$\ln \varphi_\xi(t) = i(t, m) - \frac{1}{2}(\mathbb{R}t, t) = i \sum_{k=1}^n t_k m_k - \frac{1}{2} \sum_{k,l=1}^n r_{kl} t_k t_l, \quad (10)$$

we find from (12.35) and the formulas that connect the moments and the semi-invariants that

$$m_1 = s_{\xi}^{(1, 0, \dots, 0)} = E \xi_1, \dots, m_k = s_{\xi}^{(0, \dots, 0, 1)} = E \xi_k.$$

Similarly

$$r_{11} = s_{\xi}^{(2, 0, \dots, 0)} = V \xi_1, \quad r_{12} = s_{\xi}^{(1, 1, 0, \dots)} = \text{cov}(\xi_1, \xi_2),$$

and generally

$$r_{kl} = \text{cov}(\xi_k, \xi_l).$$

Consequently  $m$  is the *mean-value vector* of  $\xi$  and  $\mathbb{R}$  is its covariance matrix.

If  $\mathbb{R}$  is nonsingular, we can obtain this result in a different way. In fact, in this case  $\xi$  has a density  $f(x)$  given by (6).

A direct calculation shows that

$$E \xi_k \equiv \int x_k f(x) dx = m_k, \quad (11)$$

$$\text{cov}(\xi_k, \xi_l) = \int (x_k - m_k)(x_l - m_l) f(x) dx = r_{kl}.$$

4. Let us discuss some properties of Gaussian vectors.

### Theorem 1

- (a) *The components of a Gaussian vector are uncorrelated if and only if they are independent.*
- (b) *A vector  $\xi = (\xi_1, \dots, \xi_n)$  is Gaussian if and only if, for every vector  $\lambda = (\lambda_1, \dots, \lambda_n)$ ,  $\lambda_k \in \mathbb{R}$ , the random variable  $(\xi, \lambda) = \lambda_1 \xi_1 + \dots + \lambda_n \xi_n$  has a Gaussian distribution.*

PROOF. (a) If the components of  $\xi = (\xi_1, \dots, \xi_n)$  are uncorrelated, it follows from the form of the characteristic function  $\varphi_{\xi}(t)$  that it is a product of characteristic functions. Therefore, by Theorem 4 of §12, the components are independent.

The converse is evident, since independence always implies lack of correlation.

(b) If  $\xi$  is a Gaussian vector, it follows from (5) that

$$E \exp\{it(\xi_1 \lambda_1 + \dots + \xi_n \lambda_n)\} = \exp\left\{it\left(\sum \lambda_k m_k\right) - \frac{t^2}{2} \left(\sum r_{kl} \lambda_k \lambda_l\right)\right\}, \quad t \in \mathbb{R},$$

and consequently

$$(\xi, \lambda) \sim \mathcal{N}(\sum \lambda_k m_k, \sum r_{kl} \lambda_k \lambda_l).$$

Conversely, to say that the random variable  $(\xi, \lambda) = \xi_1 \lambda_1 + \dots + \xi_n \lambda_n$  is Gaussian means, in particular, that

$$E e^{i(\xi, \lambda)} = \exp \left\{ i E(\xi, \lambda) - \frac{V(\xi, \lambda)}{2} \right\} = \exp \left\{ i \sum \lambda_k E \xi_k - \frac{1}{2} \sum \lambda_k \lambda_l \text{cov}(\xi_k, \xi_l) \right\}.$$

Since  $\lambda_1, \dots, \lambda_n$  are arbitrary it follows from Definition 1 that the vector  $\xi = (\xi_1, \dots, \xi_n)$  is Gaussian.

This completes the proof of the theorem.

**Remark.** Let  $(\theta, \xi)$  be a Gaussian vector with  $\theta = (\theta_1, \dots, \theta_k)$  and  $\xi = (\xi_1, \dots, \xi_l)$ . If  $\theta$  and  $\xi$  are uncorrelated, i.e.  $\text{cov}(\theta_i, \xi_j) = 0$ ,  $i = 1, \dots, k$ ;  $j = 1, \dots, l$ , they are independent.

The proof is the same as for conclusion (a) of the theorem.

Let  $\xi = (\xi_1, \dots, \xi_n)$  be a Gaussian vector; let us suppose, for simplicity, that its mean-value vector is zero. If  $\text{rank } \mathbb{R} = r < n$ , then (as was shown in §11), there are  $n - r$  linear relations connecting  $\xi_1, \dots, \xi_n$ . We may then suppose that, say,  $\xi_1, \dots, \xi_r$  are linearly independent, and the others can be expressed linearly in terms of them. Hence all the basic properties of the vector  $\xi = \xi_1, \dots, \xi_n$  are determined by the first  $r$  components  $(\xi_1, \dots, \xi_r)$  for which the corresponding covariance matrix is already known to be nonsingular.

Thus we may suppose that the original vector  $\xi = (\xi_1, \dots, \xi_n)$  had linearly independent components and therefore that  $|\mathbb{R}| > 0$ .

Let  $\mathcal{O}$  be an orthogonal matrix that diagonalizes  $\mathbb{R}$ ,

$$\mathcal{O}^T \mathbb{R} \mathcal{O} = D.$$

The diagonal elements of  $D$  are positive and therefore determine the inverse matrix. Put  $B^2 = D$  and

$$\beta = B^{-1} \mathcal{O}^T \xi.$$

Then it is easily verified that

$$E e^{i(t, \beta)} = E e^{i \beta^T t} = e^{-(1/2)(Et, t)},$$

i.e. the vector  $\beta = (\beta_1, \dots, \beta_n)$  is a Gaussian vector with components that are uncorrelated and therefore (Theorem 1) independent. Then if we write  $A = \mathcal{O}B$  we find that the original Gaussian vector  $\xi = (\xi_1, \dots, \xi_n)$  can be represented as

$$\xi = A\beta, \quad (12)$$

where  $\beta = (\beta_1, \dots, \beta_n)$  is a Gaussian vector with independent components,  $\beta_k \sim \mathcal{N}(0, 1)$ . Hence we have the following result. Let  $\xi = (\xi_1, \dots, \xi_n)$  be a

vector with linearly independent components such that  $E\xi_k = 0$ ,  $k = 1, \dots, n$ . This vector is Gaussian if and only if there are independent Gaussian variables  $\beta_1, \dots, \beta_n$ ,  $\beta_k \sim \mathcal{N}(0, 1)$ , and a nonsingular matrix  $A$  of order  $n$  such that  $\xi = A\beta$ . Here  $R = AA^T$  is the covariance matrix of  $\xi$ .

If  $|R| \neq 0$ , then by the Gram-Schmidt method (see §11)

$$\xi_k = \hat{\xi}_k + b_k \varepsilon_k, \quad k = 1, \dots, n, \quad (13)$$

where since  $\varepsilon = (\varepsilon_1, \dots, \varepsilon_k) \sim \mathcal{N}(0, E)$  is a Gaussian vector,

$$\hat{\xi}_k = \sum_{l=1}^{k-1} (\xi_k, \varepsilon_l) \varepsilon_l, \quad (14)$$

$$b_k = \|\xi_k - \hat{\xi}_k\| \quad (15)$$

and

$$\mathcal{L}\{\xi_1, \dots, \xi_k\} = \mathcal{L}\{\varepsilon_1, \dots, \varepsilon_k\}. \quad (16)$$

We see immediately from the orthogonal decomposition (13) that

$$\hat{\xi}_k = E(\xi_k | \xi_{k-1}, \dots, \xi_1). \quad (17)$$

From this, with (16) and (14), it follows that in the Gaussian case the conditional expectation  $E(\xi_k | \xi_{k-1}, \dots, \xi_1)$  is a linear function of  $(\xi_1, \dots, \xi_{k-1})$ :

$$E(\xi_k | \xi_{k-1}, \dots, \xi_1) = \sum_{i=1}^{k-1} a_i \xi_i. \quad (18)$$

(This was proved in §8 for the case  $k = 2$ .)

Since, according to a remark made in Theorem 1 of §8,  $E(\xi_k | \xi_{k-1}, \dots, \xi_1)$  is an optimal estimator (in the mean-square sense) for  $\xi_k$  in terms of  $\xi_1, \dots, \xi_{k-1}$ , it follows from (18) that in the Gaussian case the optimal estimator is *linear*.

We shall use these results in looking for optimal estimators of  $\theta = (\theta_1, \dots, \theta_k)$  in terms of  $\xi = (\xi_1, \dots, \xi_l)$  under the hypothesis that  $(\theta, \xi)$  is Gaussian. Let

$$m_\theta = E\theta, \quad m_\xi = E\xi$$

be the column-vector mean values and

$$V_{\theta\theta} \equiv \text{cov}(\theta, \theta) \equiv \|\text{cov}(\theta_i, \theta_j)\|, \quad 1 \leq i, j \leq k,$$

$$V_{\theta\xi} \equiv \text{cov}(\theta, \xi) \equiv \|\text{cov}(\theta_i, \xi_j)\|, \quad 1 \leq i \leq k, 1 \leq j \leq l,$$

$$V_{\xi\xi} \equiv \text{cov}(\xi, \xi) \equiv \|\text{cov}(\xi_i, \xi_j)\|, \quad 1 \leq i, j \leq l$$

the covariance matrices. Let us suppose that  $V_{\xi\xi}$  has an inverse. Then we have the following theorem.

**Theorem 2** (Theorem on Normal Correlation). *For a Gaussian vector  $(\theta, \xi)$ , the optimal estimator  $E(\theta | \xi)$  of  $\theta$  in terms of  $\xi$ , and its error matrix*

$$\Delta = E[\theta - E(\theta | \xi)][\theta - E(\theta | \xi)]^T$$

are given by the formulas

$$E(\theta|\xi) = m_\theta + V_{\theta\xi} V_{\xi\xi}^{-1}(\xi - m_\xi), \quad (19)$$

$$\Delta = V_{\theta\theta} - V_{\theta\xi} V_{\xi\xi}^{-1} (V_{\theta\xi})^T. \quad (20)$$

PROOF. Form the vector

$$\eta = (\theta - m_\theta) - V_{\theta\xi} V_{\xi\xi}^{-1}(\xi - m_\xi). \quad (21)$$

We can verify at once that  $E\eta(\xi - m_\xi)^T = 0$ , i.e.  $\eta$  is not correlated with  $(\xi - m_\xi)$ . But since  $(\theta, \xi)$  is Gaussian, the vector  $(\eta, \xi)$  is also Gaussian. Hence by the remark on Theorem 1,  $\eta$  and  $\xi - m_\xi$  are independent. Therefore  $\eta$  and  $\xi$  are independent, and consequently  $E(\eta|\xi) = E\eta = 0$ . Therefore

$$E[\theta - m_\theta|\xi] - V_{\theta\xi} V_{\xi\xi}^{-1}(\xi - m_\xi) = 0.$$

which establishes (19).

To establish (20) we consider the conditional covariance

$$\text{cov}(\theta, \theta|\xi) \equiv E[(\theta - E(\theta|\xi))(\theta - E(\theta|\xi))^T|\xi]. \quad (22)$$

Since  $\theta - E(\theta|\xi) = \eta$ , and  $\eta$  and  $\xi$  are independent, we find that

$$\begin{aligned} \text{cov}(\theta, \theta|\xi) &= E(\eta\eta^T|\xi) = E\eta\eta^T \\ &= V_{\theta\theta} + V_{\theta\xi}^{-1} V_{\xi\xi} V_{\xi\xi}^{-1} V_{\theta\xi}^T - 2V_{\theta\xi} V_{\xi\xi}^{-1} V_{\xi\xi} V_{\xi\xi}^{-1} V_{\theta\xi}^T \\ &= V_{\theta\theta} - V_{\theta\xi} V_{\xi\xi}^{-1} V_{\theta\xi}^T. \end{aligned}$$

Since  $\text{cov}(\theta, \theta|\xi)$  does not depend on "chance," we have

$$\Delta = E \text{cov}(\theta, \theta|\xi) = \text{cov}(\theta, \theta|\xi),$$

and this establishes (20).

**Corollary.** Let  $(\theta, \xi_1, \dots, \xi_n)$  be an  $(n+1)$ -dimensional Gaussian vector, with  $\xi_1, \dots, \xi_n$  independent. Then

$$\begin{aligned} E(\theta|\xi_1, \dots, \xi_n) &= E\theta + \sum_{i=1}^n \frac{\text{cov}(\theta, \xi_i)}{V_{\xi_i}} (\xi_i - E\xi_i), \\ \Delta &= V_\theta - \sum_{i=1}^n \frac{\text{cov}^2(\theta, \xi_i)}{V_{\xi_i}} \end{aligned}$$

(cf. (8.12) and (8.13)).

5. Let  $\xi_1, \xi_2, \dots$  be a sequence of Gaussian random vectors that converge in probability to  $\xi$ . Let us show that  $\xi$  is also Gaussian.

In accordance with (b) of Theorem 1, it is enough to establish this only for random variables.

Let  $m_n = E\xi_n$ ,  $\sigma_n^2 = V\xi_n$ . Then by Lebesgue's dominated convergence theorem

$$\lim_{n \rightarrow \infty} e^{itm_n - (1/2)\sigma_n^2 t^2} = \lim_{n \rightarrow \infty} E e^{it\xi_n} = E e^{it\xi}.$$

It follows from the existence of the limit on the left-hand side that there are numbers  $m$  and  $\sigma^2$  such that

$$m = \lim_{n \rightarrow \infty} m_n, \quad \sigma^2 = \lim_{n \rightarrow \infty} \sigma_n^2.$$

Consequently

$$E e^{it\xi} = e^{itm - (1/2)\sigma^2 t^2},$$

i.e.  $\xi \sim \mathcal{N}(m, \sigma^2)$ .

It follows, in particular, that the closed linear manifold  $\overline{\mathcal{L}}(\xi_1, \xi_2, \dots)$  generated by the Gaussian variables  $\xi_1, \xi_2, \dots$  (see §11, Subsection 5) consists of Gaussian variables.

6. We now turn to the concept of Gaussian systems in general.

**Definition 2.** A collection of random variables  $\xi = (\xi_\alpha)$ , where  $\alpha$  belongs to some index set  $\mathfrak{A}$ , is a *Gaussian system* if the random vector  $(\xi_{\alpha_1}, \dots, \xi_{\alpha_n})$  is Gaussian for every  $n \geq 1$  and all indices  $\alpha_1, \dots, \alpha_n$  chosen from  $\mathfrak{A}$ .

Let us notice some properties of Gaussian systems.

- (a) If  $\xi = (\xi_\alpha)$ ,  $\alpha \in \mathfrak{A}$ , is a Gaussian system, then every subsystem  $\xi' = (\xi'_{\alpha'})$ ,  $\alpha' \in \mathfrak{A}' \subseteq \mathfrak{A}$ , is also Gaussian.
- (b) If  $\xi_\alpha$ ,  $\alpha \in \mathfrak{A}$ , are independent Gaussian variables, then the system  $\xi = (\xi_\alpha)$ ,  $\alpha \in \mathfrak{A}$ , is Gaussian.
- (c) If  $\xi = (\xi_\alpha)$ ,  $\alpha \in \mathfrak{A}$ , is a Gaussian system, the closed linear manifold  $\overline{\mathcal{L}}(\xi)$ , consisting of all variables of the form  $\sum_{i=1}^n c_{\alpha_i} \xi_{\alpha_i}$ , together with their mean-square limits, forms a Gaussian system.

Let us observe that the converse of (a) is false in general. For example, let  $\xi_1$  and  $\eta_1$  be independent and  $\xi_1 \sim \mathcal{N}(0, 1)$ ,  $\eta_1 \sim \mathcal{N}(0, 1)$ . Define the system

$$(\xi, \eta) = \begin{cases} (\xi_1, |\eta_1|) & \text{if } \xi_1 \geq 0, \\ (\xi_1, -|\eta_1|) & \text{if } \xi_1 < 0. \end{cases} \quad (23)$$

Then it is easily verified that  $\xi$  and  $\eta$  are both Gaussian, but  $(\xi, \eta)$  is not.

Let  $\xi = (\xi_\alpha)_{\alpha \in \mathfrak{A}}$  be a Gaussian system with mean-value vector  $m = (m_\alpha)$ ,  $\alpha \in \mathfrak{A}$ , and covariance matrix  $\mathbb{R} = (r_{\alpha\beta})_{\alpha, \beta \in \mathfrak{A}}$ , where  $m_\alpha = E\xi_\alpha$ . Then  $\mathbb{R}$  is evidently symmetric ( $r_{\alpha\beta} = r_{\beta\alpha}$ ) and nonnegative definite in the sense that for every vector  $c = (c_\alpha)_{\alpha \in \mathfrak{A}}$  with values in  $R^{\mathfrak{A}}$ , and only a finite number of nonzero coordinates  $c_\alpha$ ,

$$(\mathbb{R}c, c) \equiv \sum_{\alpha, \beta} r_{\alpha\beta} c_\alpha c_\beta \geq 0. \quad (24)$$

We now ask the converse question. Suppose that we are given a parameter set  $\mathfrak{U} = \{\alpha\}$ , a vector  $m = (m_\alpha)_{\alpha \in \mathfrak{U}}$  and a symmetric nonnegative definite matrix  $\mathbb{R} = (r_{\alpha\beta})_{\alpha, \beta \in \mathfrak{U}}$ . Do there exist a probability space  $(\Omega, \mathcal{F}, P)$  and a Gaussian system of random variables  $\xi = (\xi_\alpha)_{\alpha \in \mathfrak{U}}$  on it, such that

$$E\xi_\alpha = m_\alpha,$$

$$\text{cov}(\xi_\alpha, \xi_\beta) = r_{\alpha, \beta}, \quad \alpha, \beta \in \mathfrak{U}?$$

If we take a finite set  $\alpha_1, \dots, \alpha_n$ , then for the vector  $\bar{m} = (m_{\alpha_1}, \dots, m_{\alpha_n})$  and the matrix  $\mathbb{R} = (r_{\alpha\beta})$ ,  $\alpha, \beta = \alpha_1, \dots, \alpha_n$ , we can construct in  $R^n$  the Gaussian distribution  $F_{\alpha_1, \dots, \alpha_n}(x_1, \dots, x_n)$  with characteristic function

$$\varphi(t) = \exp\{i(t, \bar{m}) - \frac{1}{2}(\mathbb{R}t, t)\}, \quad t = (t_{\alpha_1}, \dots, t_{\alpha_n}).$$

It is easily verified that the family

$$\{F_{\alpha_1, \dots, \alpha_n}(x_1, \dots, x_n); \alpha_i \in \mathfrak{U}\}$$

is consistent. Consequently by Kolmogorov's theorem (Theorem 1, §9, and Remark 2 on this) the answer to our question is positive.

7. If  $\mathfrak{U} = \{1, 2, \dots\}$ , then in accordance with the terminology of §5 the system of random variables  $\xi = (\xi_\alpha)_{\alpha \in \mathfrak{U}}$  is a *random sequence* and is denoted by  $\xi = (\xi_1, \xi_2, \dots)$ . A Gaussian sequence is completely described by its mean-value vector  $m = (m_1, m_2, \dots)$  and covariance matrix  $\mathbb{R} = \|r_{ij}\|$ ,  $r_{ij} = \text{cov}(\xi_i, \xi_j)$ . In particular, if  $r_{ij} = \sigma_i^2 \delta_{ij}$ , then  $\xi = (\xi_1, \xi_2, \dots)$  is a Gaussian sequence of independent random variables with  $\xi_i \sim \mathcal{N}(m_i, \sigma_i^2)$ ,  $i \geq 1$ .

When  $\mathfrak{U} = [0, 1], [0, \infty), (-\infty, \infty), \dots$ , the system  $\xi = (\xi_t)$ ,  $t \in \mathfrak{U}$ , is a *random process with continuous time*.

Let us mention some examples of Gaussian random processes. If we take their mean values to be zero, their probabilistic properties are completely described by the covariance matrices  $\|r_{st}\|$ . We write  $r(s, t)$  instead of  $r_{st}$  and call it the *covariance function*.

EXAMPLE 1. If  $T = [0, \infty)$  and

$$r(s, t) = \min(s, t), \quad (25)$$

the Gaussian process  $\xi = (\xi_t)_{t \geq 0}$  with this covariance function (see Problem 2) and  $\xi_0 \equiv 0$  is a *Brownian motion* or *Wiener process*.

Observe that this process has *independent increments*; that is, for arbitrary  $t_1 < t_2 < \dots < t_n$  the random variables

$$\xi_{t_2} - \xi_{t_1}, \dots, \xi_{t_n} - \xi_{t_{n-1}}$$

are independent. In fact, because the process is Gaussian it is enough to verify only that the increments are uncorrelated. But if  $s < t < u < v$  then

$$\begin{aligned} E[\xi_t - \xi_s][\xi_v - \xi_u] &= [r(t, v) - r(t, u)] - [r(s, v) - r(s, u)] \\ &= (t - t) - (s - s) = 0. \end{aligned}$$

EXAMPLE 2. The process  $\xi = (\xi_t)$ ,  $0 \leq t \leq 1$ , with  $\xi_0 \equiv 0$  and

$$r(s, t) = \min(s, t) - st \quad (26)$$

is a *conditional Wiener process* (observe that since  $r(1, 1) = 0$  we have  $P(\xi_1 = 0) = 1$ ).

EXAMPLE 3. The process  $\xi = (\xi_t)$ ,  $-\infty < t < \infty$ , with

$$r(s, t) = e^{-|t-s|} \quad (27)$$

is a *Gauss-Markov process*.

## 8. PROBLEMS

1. Let  $\xi_1, \xi_2, \xi_3$  be independent Gaussian random variables,  $\xi_i \sim \mathcal{N}(0, 1)$ . Show that

$$\frac{\xi_1 + \xi_2 \xi_3}{\sqrt{1 + \xi_3^2}} \sim \mathcal{N}(0, 1).$$

(In this case we encounter the interesting problem of describing the nonlinear transformations of independent Gaussian variables  $\xi_1, \dots, \xi_n$  whose distributions are still Gaussian.)

2. Show that (25), (26) and (27) are nonnegative definite (and consequently are actually covariance functions).
3. Let  $A$  be an  $m \times n$  matrix. An  $n \times m$  matrix  $A^\oplus$  is a *pseudoinverse* of  $A$  if there are matrices  $U$  and  $V$  such that

$$AA^\oplus A = A, \quad A^\oplus = UA^\top = A^\top V.$$

Show that  $A^\oplus$  exists and is unique.

4. Show that (19) and (20) in the theorem on normal correlation remains valid when  $V_{\xi\xi}$  is singular provided that  $V_{\xi\xi}^{-1}$  is replaced by  $V_{\xi\xi}^\oplus$ .
5. Let  $(\theta, \xi) = (\theta_1, \dots, \theta_k; \xi_1, \dots, \xi_l)$  be a Gaussian vector with nonsingular matrix  $\Delta \equiv V_{\theta\theta} - V_{\xi\xi}^\oplus V_{\theta\xi}^*$ . Show that the distribution function

$$P(\theta \leq a | \xi) = P(\theta_1 \leq a_1, \dots, \theta_k \leq a_k | \xi)$$

has (P-a.s.) the density  $p(a_1, \dots, a_k | \xi)$  defined by

$$\frac{|\Delta|^{-1/2}}{(2\pi)^{k/2}} \exp\left\{-\frac{1}{2}(a - E(\theta | \xi))^T \Delta^{-1}(a - E(\theta | \xi))\right\}.$$

6. (S. N. Bernstein). Let  $\xi$  and  $\eta$  be independent identically distributed random variables with finite variances. Show that if  $\xi + \eta$  and  $\xi - \eta$  are independent, then  $\xi$  and  $\eta$  are Gaussian.



## CHAPTER III

# Convergence of Probability Measures. Central Limit Theorem

### §1. Weak Convergence of Probability Measures and Distributions

1. Many of the fundamental results in probability theory are formulated as *limit theorems*. Bernoulli's law of large numbers was formulated as a limit theorem; so was the De Moivre–Laplace theorem, which can fairly be called the origin of a genuine theory of probability and, in particular, which led the way to numerous investigations that clarified the conditions for the validity of the central limit theorem. Poisson's theorem on the approximation of the binomial distribution by the "Poisson" distribution in the case of rare events was formulated as a limit theorem. After the example of these propositions, and of results on the rapidity of convergence in the De Moivre–Laplace and Poisson theorems, it became clear that in probability it is necessary to deal with various kinds of convergence of distributions, and to establish the rapidity of convergence connected with the introduction of various "natural" measures of the distance between distributions. In the present chapter we shall discuss some general features of the convergence of probability distributions and of the distance between them. In this section we take up questions in the general theory of weak convergence of probability measures in metric spaces. The De Moivre–Laplace theorem, the progenitor of the central limit theorem, finds a natural place in this theory. From §3, it is clear that the method of characteristic functions is one of the most powerful means for proving limit theorems on the weak convergence of probability distributions in  $R^n$ . In §6, we consider questions of metrizability of weak convergence. Then, in §8, we turn our attention to a different kind of convergence of distributions (stronger than weak convergence), namely convergence in vari-

ation. Proofs of the simplest results on the rapidity of convergence in the central limit theorem and Poisson's theorem will be given in §§10 and 11.

2. We begin by recalling the statement of the law of large numbers (Chapter I, §5) for the Bernoulli scheme.

Let  $\xi_1, \xi_2, \dots$  be a sequence of independent identically distributed random variables with  $P(\xi_i = 1) = p$ ,  $P(\xi_i = 0) = q$ ,  $p + q = 1$ . In terms of the concept of convergence in probability (Chapter II, §10), Bernoulli's law of large numbers can be stated as follows:

$$\frac{S_n}{n} \xrightarrow{P} p, \quad n \rightarrow \infty, \quad (1)$$

where  $S_n = \xi_1 + \dots + \xi_n$ . (It will be shown in Chapter IV that in fact we have convergence with probability 1.)

We put

$$F_n(x) = P\left\{\frac{S_n}{n} \leq x\right\},$$

$$F(x) = \begin{cases} 1, & x \geq p, \\ 0, & x < p, \end{cases} \quad (2)$$

where  $F(x)$  is the distribution function of the degenerate random variable  $\xi \equiv p$ . Also let  $P_n$  and  $P$  be the probability measures on  $(R, \mathcal{B}(R))$  corresponding to the distributions  $F_n$  and  $F$ .

In accordance with Theorem 2 of §10, Chapter II, convergence in probability,  $S_n/n \xrightarrow{P} p$ , implies convergence in distribution,  $S_n/n \xrightarrow{d} p$ , which means that

$$Ef\left(\frac{S_n}{n}\right) \rightarrow Ef(p), \quad n \rightarrow \infty, \quad (3)$$

for every function  $f = f(x)$  belonging to the class  $C(R)$  of bounded continuous functions on  $R$ .

Since

$$Ef\left(\frac{S_n}{n}\right) = \int_R f(x) P_n(dx), \quad Ef(p) = \int_R f(x) P(dx),$$

(3) can be written in the form

$$\int_R f(x) P_n(dx) \rightarrow \int_R f(x) P(dx), \quad f \in C(R), \quad (4)$$

or (in accordance with §6 of Chapter II) in the form

$$\int_R f(x) dF_n(x) \rightarrow \int_R f(x) dF(x), \quad f \in C(R). \quad (5)$$

In analysis, (4) is called *weak convergence* (of  $P_n$  to  $P$ ,  $n \rightarrow \infty$ ) and written  $P_n \xrightarrow{w} P$ . It is also natural to call (5) weak convergence of  $F_n$  to  $F$  and denote it by  $F_n \xrightarrow{w} F$ .

Thus we may say that in a Bernoulli scheme

$$\frac{S_n}{n} \xrightarrow{P} p \Rightarrow F_n \xrightarrow{w} F. \quad (6)$$

It is also easy to see from (1) that, for the distribution functions defined in (2),

$$F_n(x) \rightarrow F(x), \quad n \rightarrow \infty,$$

for all points  $x \in R$  except for the single point  $x = p$ , where  $F(x)$  has a discontinuity.

This shows that weak convergence  $F_n \rightarrow F$  does not imply pointwise convergence of  $F_n(x)$  to  $F(x)$ ,  $n \rightarrow \infty$ , for all points  $x \in R$ . However, it turns out that, both for Bernoulli schemes and for arbitrary distribution functions, weak convergence is equivalent (see Theorem 2 below) to "convergence in general" in the sense of the following definition.

**Definition 1.** A sequence of distribution functions  $\{F_n\}$ , defined on the real line, converges in general to the distribution function  $F$  (notation:  $F_n \Rightarrow F$ ) if as  $n \rightarrow \infty$

$$F_n(x) \rightarrow F(x), \quad x \in P_C(F),$$

where  $P_C(F)$  is the set of points of continuity of  $F = F(x)$ .

For Bernoulli schemes,  $F = F(x)$  is degenerate, and it is easy to see (see Problem 7 of §10, Chapter II) that

$$(F_n \Rightarrow F) \Rightarrow \left( \frac{S_n}{n} \xrightarrow{P} p \right).$$

Therefore, taking account of Theorem 2 below,

$$\left( \frac{S_n}{n} \xrightarrow{P} p \right) \Rightarrow (F_n \xrightarrow{w} F) \Leftrightarrow (F_n \Rightarrow F) \Rightarrow \left( \frac{S_n}{n} \xrightarrow{P} p \right) \quad (7)$$

and consequently the law of large numbers can be considered as a theorem on the weak convergence of the distribution functions defined in (2).

Let us write

$$F_n(x) = P \left\{ \frac{S_n - np}{\sqrt{npq}} \leq x \right\},$$

$$F(x) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^x e^{-u^2/2} du. \quad (8)$$

The De Moivre-Laplace theorem (§6, Chapter I) states that  $F_n(x) \rightarrow F(x)$  for all  $x \in R$ , and consequently  $F_n \Rightarrow F$ . Since, as we have observed, weak convergence  $F_n \xrightarrow{w} F$  and convergence in general,  $F_n \Rightarrow F$ , are equivalent, we may therefore say that the De Moivre-Laplace theorem is also a theorem on the weak convergence of the distribution functions defined by (8).

These examples justify the concept of weak convergence of probability measures that will be introduced below in Definition 2. Although, on the real line, weak convergence is equivalent to convergence in general of the corresponding distribution functions, it is preferable to use weak convergence from the beginning. This is because in the first place it is easier to work with, and in the second place it remains useful in more general spaces than the real line, and in particular for metric spaces, including the especially important spaces  $R^n$ ,  $R^\infty$ ,  $C$ , and  $D$  (see §3 of Chapter II).

3. Let  $(E, \mathcal{E}, \rho)$  be a metric space with metric  $\rho = \rho(x, y)$  and  $\sigma$ -algebra  $\mathcal{E}$  of Borel subsets generated by the open sets, and let  $P, P_1, P_2, \dots$  be probability measures on  $(E, \mathcal{E}, \rho)$ .

**Definition 2.** A sequence of probability measures  $\{P_n\}$  converges weakly to the probability measure  $P$  (notation:  $P_n \xrightarrow{w} P$ ) if

$$\int_E f(x) P_n(dx) \rightarrow \int_E f(x) P(dx) \quad (9)$$

for every function  $f = f(x)$  in the class  $\mathcal{C}(E)$  of continuous bounded functions on  $E$ .

**Definition 3.** A sequence of probability measures  $\{P_n\}$  converges in general to the probability measure  $P$  (notation:  $P_n \Rightarrow P$ ) if

$$P_n(A) \rightarrow P(A) \quad (10)$$

for every set  $A$  of  $\mathcal{E}$  for which

$$P(\partial A) = 0. \quad (11)$$

(Here  $\partial A$  denotes the boundary of  $A$ :

$$\partial A = [A] \cap [\bar{A}],$$

where  $[A]$  is the closure of  $A$ .)

The following fundamental theorem shows the equivalence of the concepts of weak convergence and convergence in general for probability measures, and contains still another equivalent statement.

**Theorem 1.** The following statements are equivalent.

- (I)  $P_n \xrightarrow{w} P$ .
- (II)  $\limsup P_n(A) \leq P(A)$ ,  $A$  closed.
- (III)  $\liminf P_n(A) \geq P(A)$ ,  $A$  open.
- (IV)  $P_n \Rightarrow P$ .

PROOF. (I)  $\Rightarrow$  (II). Let  $A$  be closed,  $f(x) = I_A(x)$  and

$$f_\varepsilon(x) = g\left(\frac{1}{\varepsilon} \rho(x, A)\right), \quad \varepsilon > 0,$$

where

$$\rho(x, A) = \inf\{\rho(x, y): y \in A\},$$

$$g(t) = \begin{cases} 1, & t \leq 0, \\ 1 - t, & 0 \leq t \leq 1, \\ 0, & t \geq 1. \end{cases}$$

Let us also put

$$A_\varepsilon = \{x: \rho(x, A) < \varepsilon\}$$

and observe that  $A_\varepsilon \downarrow A$  as  $\varepsilon \downarrow 0$ .

Since  $f_\varepsilon(x)$  is bounded, continuous, and satisfies

$$P_n(A) = \int_E I_A(x) P_n(dx) \leq \int_E f_\varepsilon(x) P_n(dx),$$

we have

$$\overline{\lim}_n P_n(A) \leq \overline{\lim}_n \int_E f_\varepsilon(x) P_n(dx) = \int_E f_\varepsilon(x) P(dx) \leq P(A_\varepsilon) \downarrow P(A), \quad \varepsilon \downarrow 0,$$

which establishes the required implication.

The implications (II)  $\Rightarrow$  (III) and (III)  $\Rightarrow$  (II) become obvious if we take the complements of the sets concerned.

(III)  $\Rightarrow$  (IV). Let  $A^0 = A \setminus \partial A$  be the interior, and  $[A]$  the closure, of  $A$ . Then from (II), (III), and the hypothesis  $P(\partial A) = 0$ , we have

$$\overline{\lim}_n P_n(A) \leq \overline{\lim}_n P_n([A]) \leq P([A]) = P(A),$$

$$\underline{\lim}_n P_n(A) \geq \underline{\lim}_n P_n(A^0) \geq P(A^0) = P(A),$$

and therefore  $P_n(A) \rightarrow P(A)$  for every  $A$  such that  $P(\partial A) = 0$ .

(IV)  $\rightarrow$  (I). Let  $f = f(x)$  be a bounded continuous function with  $|f(x)| \leq M$ . We put

$$D = \{t \in R: P\{x: f(x) = t\} \neq 0\}$$

and consider a decomposition  $T_k = (t_0, t_1, \dots, t_k)$  of  $[-M, M]$ :

$$-M = t_0 < t_1 < \dots < t_k = M, \quad k \geq 1,$$

with  $t_i \notin D$ ,  $i = 0, 1, \dots, k$ . (Observe that  $D$  is at most countable since the sets  $f^{-1}\{t\}$  are disjoint and  $P$  is finite.)

Let  $B_i = \{x: t_i \leq f(x) < t_{i+1}\}$ . Since  $f(x)$  is continuous and therefore the set  $f^{-1}(t_i, t_{i+1})$  is open, we have  $\partial B_i \subseteq f^{-1}\{t_i\} \cup f^{-1}\{t_{i+1}\}$ . The points  $t_i, t_{i+1} \notin D$ ; therefore  $P(\partial B_i) = 0$  and, by (IV),

$$\sum_{i=0}^{k-1} t_i P_n(B_i) \rightarrow \sum_{i=0}^{k-1} t_i P(B_i). \quad (12)$$

But

$$\begin{aligned} \left| \int_E f(x) P_n(dx) - \int_E f(x) P(dx) \right| &\leq \left| \int_E f(x) P_n(dx) - \sum_{i=0}^{k-1} t_i P_n(B_i) \right| \\ &\quad + \left| \sum_{i=0}^{k-1} t_i P_n(B_i) - \sum_{i=0}^{k-1} t_i P(B_i) \right| \\ &\quad + \left| \sum_{i=0}^{k-1} t_i P(B_i) - \int_E f(x) P(dx) \right| \\ &\leq 2 \max_{0 \leq i \leq k-1} (t_{i+1} - t_i) \\ &\quad + \left| \sum_{i=0}^{k-1} t_i P_n(B_i) - \sum_{i=0}^{k-1} t_i P(B_i) \right|, \end{aligned}$$

whence, by (12), since the  $T_k$  ( $k \geq 1$ ) are arbitrary,

$$\lim_n \int_E f(x) P_n(dx) = \int_E f(x) P(dx).$$

This completes the proof of the theorem.

**Remark 1.** The functions  $f(x) = I_A(x)$  and  $f_\varepsilon(x)$  that appear in the proof that (I)  $\Rightarrow$  (II) are respectively *upper semicontinuous* and *uniformly continuous*. Hence it is easy to show that each of the conditions of the theorem is equivalent to one of the following:

- (V)  $\int_E f(x) P_n(x) dx \rightarrow \int_E f(x) P(dx)$  for all bounded uniformly continuous  $f(x)$ ;
- (VI)  $\lim \int_E f(x) P_n(dx) \leq \int_E f(x) P(dx)$  for all bounded  $f(x)$  that are upper semicontinuous ( $\lim f(x_n) \leq f(x)$ ,  $x_n \rightarrow x$ );
- (VII)  $\lim \int_E f(x) P_n(dx) \geq \int_E f(x) P(dx)$  for all bounded  $f(x)$  that are lower semicontinuous ( $\lim f(x_n) \geq f(x)$ ,  $x_n \rightarrow x$ ).

**Remark 2.** Theorem 1 admits a natural generalization to the case when the probability measures  $P$  and  $P_n$  defined on  $(E, \mathcal{E}, \rho)$  are replaced by arbitrary (not necessarily probability) *finite measures*  $\mu$  and  $\mu_n$ . For such measures we can introduce weak convergence  $\mu_n \xrightarrow{w} \mu$  and convergence in general  $\mu_n \Rightarrow \mu$  and, just as in Theorem 1, we can establish the equivalence of the following conditions:

- (I\*)  $\mu_n \xrightarrow{w} \mu$ ;
- (II\*)  $\lim \mu_n(A) \leq \mu(A)$ , where  $A$  is closed and  $\mu_n(E) \rightarrow \mu(E)$ ;
- (III\*)  $\lim \mu_n(A) \geq \mu(A)$ , where  $A$  is open and  $\mu_n(E) \rightarrow \mu(E)$ ;
- (IV\*)  $\mu_n \Rightarrow \mu$ .

Each of these is equivalent to any of (V\*), (VI\*), and (VII\*), which are (V), (VI), and (VII) with  $P_n$  and  $P$  replaced by  $\mu_n$  and  $\mu$ .

4. Let  $(R, \mathcal{B}(R))$  be the real line with the system  $\mathcal{B}(R)$  of sets generated by the Euclidean metric  $\rho(x, y) = |x - y|$  (compare Remark 2 of subsection 2 of §2 of Chapter II). Let  $P$  and  $P_n, n \geq 1$ , be probability measures on  $(R, \mathcal{B}(R))$  and let  $F$  and  $F_n, n \geq 1$ , be the corresponding distribution functions.

**Theorem 2.** *The following conditions are equivalent:*

- (1)  $P_n \xrightarrow{w} P$ ,
- (2)  $P_n \Rightarrow P$ ,
- (3)  $F_n \xrightarrow{w} F$ ,
- (4)  $F_n \Rightarrow F$ .

PROOF. Since (2)  $\Leftrightarrow$  (1)  $\Leftrightarrow$  (3), it is enough to show that (2)  $\Leftrightarrow$  (4).

If  $P_n \Rightarrow P$ , then in particular

$$P_n(-\infty, x] \rightarrow P(-\infty, x]$$

for all  $x \in R$  such that  $P\{x\} = 0$ . But this means that  $F_n \Rightarrow F$ .

Now let  $F_n \Rightarrow F$ . To prove that  $P_n \Rightarrow P$  it is enough (by Theorem 1) to show that  $\lim_n P_n(A) \geq P(A)$  for every open set  $A$ .

If  $A$  is open, there is a countable collection of disjoint open intervals  $I_1, I_2, \dots$  (of the form  $(a, b)$ ) such that  $A = \sum_{k=1}^{\infty} I_k$ . Choose  $\varepsilon > 0$  and in each interval  $I_k = (a_k, b_k)$  select a subinterval  $I'_k = (a'_k, b'_k]$  such that  $a'_k, b'_k \in P_c(F)$  and  $P(I_k) \leq P(I'_k) + \varepsilon \cdot 2^{-k}$ . (Since  $F(x)$  has at most countably many discontinuities, such intervals  $I'_k, k \geq 1$ , certainly exist.) By Fatou's lemma,

$$\begin{aligned} \lim_n P_n(A) &= \lim_n \sum_{k=1}^{\infty} P_n(I_k) \geq \sum_{k=1}^{\infty} \lim_n P_n(I_k) \\ &\geq \sum_{k=1}^{\infty} \lim_n P_n(I'_k). \end{aligned}$$

But

$$P_n(I'_k) = F_n(b'_k) - F_n(a'_k) \rightarrow F(b'_k) - F(a'_k) = P(I'_k).$$

Therefore

$$\lim_n P_n(A) \geq \sum_{k=1}^{\infty} P(I'_k) \geq \sum_{k=1}^{\infty} (P(I_k) - \varepsilon \cdot 2^{-k}) = P(A) - \varepsilon.$$

Since  $\varepsilon > 0$  is arbitrary, this shows that  $\lim_n P_n(A) \geq P(A)$  if  $A$  is open.

This completes the proof of the theorem.

5. Let  $(E, \mathcal{E})$  be a measurable space. A collection  $\mathcal{K}_0(E) \subseteq \mathcal{E}$  of subsets is a *determining class* whenever two probability measures  $P$  and  $Q$  on  $(E, \mathcal{E})$

satisfy

$$P(A) = Q(A) \quad \text{for all } A \in \mathcal{K}_0(E)$$

it follows that the measures are identical, i.e.,

$$P(A) = Q(A) \quad \text{for all } A \in \mathcal{E}.$$

If  $(E, \mathcal{E}, \rho)$  is a metric space, a collection  $\mathcal{K}_1(E) \subseteq \mathcal{E}$  is a *convergence-determining class* whenever probability measures  $P, P_1, P_2, \dots$  satisfy

$$P_n(A) \rightarrow P(A) \quad \text{for all } A \in \mathcal{K}_1(E) \quad \text{with } P(\partial A) = 0$$

it follows that

$$P_n(A) \rightarrow P(A) \quad \text{for all } A \in E \quad \text{with } P(\partial A) = 0.$$

When  $(E, \mathcal{E}) = (R, \mathcal{B}(R))$ , we can take a determining class  $\mathcal{K}_0(R)$  to be the class of “elementary” sets  $\mathcal{K} = \{(-\infty, x], x \in R\}$  (Theorem 1, §3, Chapter II). It follows from the equivalence of (2) and (4) of Theorem 2 that this class  $\mathcal{K}$  is also a convergence-determining class.

It is natural to ask about such determining classes in more general spaces.

For  $R^n$ ,  $n \geq 2$ , the class  $\mathcal{K}$  of “elementary” sets of the form  $(-\infty, x] = (-\infty, x_1] \times \cdots \times (-\infty, x_n]$ , where  $x = (x_1, \dots, x_n) \in R^n$ , is both a determining class (Theorem 2, §3, Chapter II) and a convergence-determining class (Problem 2).

For  $R^\infty$  the cylindrical sets  $\mathcal{K}_0(R^\infty)$  are the “elementary” sets whose probabilities uniquely determine the probabilities of the Borel sets (Theorem 3, §3, Chapter II). It turns out that in this case the class of cylindrical sets is also the class of convergence-determining sets (Problem 3). Therefore  $\mathcal{K}_1(R^\infty) = \mathcal{K}_0(R^\infty)$ .

We might expect that the cylindrical sets would still constitute determining classes in more general spaces. However, this is, in general, not the case.

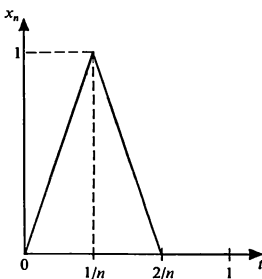


Figure 35



For example, consider the space  $(C, \mathcal{B}_0(C), \rho)$  with the uniform metric  $\rho$  (see subsection 6, §2, Chapter II). Let  $P$  be the probability measure concentrated on the element  $x = x(t) \equiv 0$ ,  $0 \leq t \leq 1$ , and let  $P_n$ ,  $n \geq 1$ , be the probability measures each of which is concentrated on the element  $x = x_n(t)$  shown in Figure 35. It is easy to see that  $P_n(A) \rightarrow P(A)$  for all cylindrical sets  $A$  with  $P(\partial A) = 0$ . But if we consider, for example, the set

$$A = \{\alpha \in C: |\alpha(t)| \leq \frac{1}{2}, 0 \leq t \leq 1\} \in \mathcal{B}_0(C),$$

then  $P(\partial A) = 0$ ,  $P_n(A) = 0$ ,  $P(A) = 1$  and consequently  $P_n \not\Rightarrow P$ .

Therefore  $\mathcal{K}_0(C) = \mathcal{B}_0(C)$  but  $\mathcal{K}_0(C) \subset \mathcal{K}_1(C)$  (with strict inclusion).

## 6. PROBLEMS

1. Let us say that a function  $F = F(x)$ , defined on  $R^n$ , is *continuous at*  $x \in R^n$  provided that, for every  $\varepsilon > 0$ , there is a  $\delta > 0$  such that  $|F(x) - F(y)| < \varepsilon$  for all  $y \in R^n$  that satisfy

$$x - \delta e < y < x + \delta e,$$

where  $e = (1, \dots, 1) \in R^n$ . Let us say that a sequence of distribution functions  $\{F_n\}$  *converges in general* to the distribution function  $F$  ( $F_n \Rightarrow F$ ) if  $F_n(x) \rightarrow F(x)$ , for all points  $x \in R^n$  where  $F = F(x)$  is continuous.

Show that the conclusion of Theorem 2 remains valid for  $R^n$ ,  $n > 1$ . (See the remark on Theorem 2.)

2. Show that the class  $\mathcal{K}$  of "elementary" sets in  $R^n$  is a convergence-determining class.
3. Let  $E$  be one of the spaces  $R^\infty$ ,  $C$ , or  $D$ . Let us say that a sequence  $\{P_n\}$  of probability measures (defined on the  $\sigma$ -algebra  $\mathcal{E}$  of Borel sets generated by the open sets) *converges in general in the sense of finite-dimensional distributions* to the probability measure  $P$  (notation:  $P_n \xrightarrow{f} P$ ) if  $P_n(A) \rightarrow P(A)$ ,  $n \rightarrow \infty$ , for all cylindrical sets  $A$  with  $P(\partial A) = 0$ .

For  $R^\infty$ , show that

$$(P_n \xrightarrow{f} P) \Leftrightarrow (P_n \Rightarrow P).$$

4. Let  $F$  and  $G$  be distribution functions on the real line and let

$$L(F, G) = \inf\{h > 0: F(x-h) - h \leq G(x) \leq F(x+h) + h\}$$

be the *Lévy distance* (between  $F$  and  $G$ ). Show that convergence in general is equivalent to convergence in the Lévy metric:

$$(F_n \Rightarrow F) \Leftrightarrow L(F_n, F) \rightarrow 0.$$

5. Let  $F_n \Rightarrow F$  and let  $F$  be continuous. Show that in this case  $F_n(x)$  converges uniformly to  $F(x)$ :

$$\sup_x |F_n(x) - F(x)| \rightarrow 0, \quad n \rightarrow \infty.$$

6. Prove the statement in Remark 1 on Theorem 1.

7. Establish the equivalence of (I\*)–(IV\*) as stated in Remark 2 on Theorem 1.
8. Show that  $P_n \xrightarrow{w} P$  if and only if every subsequence  $\{P_{n'}\}$  of  $\{P_n\}$  contains a subsequence  $\{P_{n''}\}$  such that  $P_{n''} \xrightarrow{w} P$ .

## §2. Relative Compactness and Tightness of Families of Probability Distributions

1. If we are given a sequence of probability measures, then before we can consider the question of its (weak) convergence to some probability measure, we have of course to establish whether the sequence converges in general to some measure, or has at least one convergent subsequence.

For example, the sequence  $\{P_n\}$ , where  $P_{2n} = P$ ,  $P_{2n+1} = Q$ , and  $P$  and  $Q$  are different probability measures, is evidently not convergent, but has the two convergent subsequences  $\{P_{2n}\}$  and  $\{P_{2n+1}\}$ .

It is easy to construct a sequence  $\{P_n\}$  of probability measures  $P_n$ ,  $n \geq 1$ , that not only fails to converge, but contains no convergent subsequences at all. All that we have to do is to take  $P_n$ ,  $n \geq 1$ , to be concentrated at  $\{n\}$  (that is,  $P_n\{n\} = 1$ ). In fact, since  $\lim_n P_n(a, b] = 0$  whenever  $a < b$ , a limit measure would have to be identically zero, contradicting the fact that  $1 = P_n(R) \not\rightarrow 0$ ,  $n \rightarrow \infty$ . It is interesting to observe that in this example the corresponding sequence  $\{F_n\}$  of distribution functions,

$$F_n(x) = \begin{cases} 1, & x \geq n, \\ 0, & x < n, \end{cases}$$

is evidently convergent: for every  $x \in R$ ,

$$F_n(x) \rightarrow G(x) \equiv 0.$$

However, the limit function  $G = G(x)$  is not a distribution function (in the sense of Definition 1 of §3, Chapter II).

This instructive example shows that the space of distribution functions is not compact. It also shows that if a sequence of distribution functions is to converge to a limit that is also a distribution function, we must have some hypothesis that will prevent mass from "escaping to infinity."

After these introductory remarks, which illustrate the kinds of difficulty that can arise, we turn to the basic definitions.

2. Let us suppose that all measures are defined on the metric space  $(E, \mathcal{E}, \rho)$ .

**Definition 1.** A family of probability measures  $\mathcal{P} = \{P_a; a \in \mathfrak{A}\}$  is *relatively compact* if every sequence of measures from  $\mathcal{P}$  contains a subsequence which converges weakly to a probability measure.

We emphasize that in this definition the limit measure is to be a *probability measure*, although it need not belong to the original class  $\mathcal{P}$ . (This is why the word "relatively" appears in the definition.)

It is often far from simple to verify that a given family of probability measures is relatively compact. Consequently it is desirable to have simple and useable tests for this property. We need the following definitions.

**Definition 2.** A family of probability measures  $\mathcal{P} = \{P_\alpha; \alpha \in \mathfrak{A}\}$  is *tight* if, for every  $\varepsilon > 0$ , there is a compact set  $K \subseteq E$  such that

$$\sup_{\alpha \in \mathfrak{A}} P_\alpha(E \setminus K) \leq \varepsilon. \quad (1)$$

**Definition 3.** A family of distribution functions  $F = \{F_\alpha; \alpha \in \mathfrak{A}\}$  defined on  $R^n$ ,  $n \geq 1$ , is *relatively compact* (or *tight*) if the same property is possessed by the family  $\mathcal{P} = \{P_\alpha; \alpha \in \mathfrak{A}\}$  of probability measures, where  $P_\alpha$  is the measure constructed from  $F_\alpha$ .

3. The following result is fundamental for the study of weak convergence of probability measures.

**Theorem 1** (Prokhorov's Theorem). Let  $\mathcal{P} = \{P_\alpha; \alpha \in \mathfrak{A}\}$  be a family of probability measures defined on a complete separable metric space  $(E, \mathcal{E}, \rho)$ . Then  $\mathcal{P}$  is relatively compact if and only if it is tight.

We shall give the proof only when the space is the real line. (The proof can be carried over, almost unchanged, to arbitrary Euclidean spaces  $R^n$ ,  $n \geq 2$ . Then the theorem can be extended successively to  $R^\infty$ , to  $\sigma$ -compact spaces; and finally to general complete separable metric spaces, by reducing each case to the preceding one.)

*Necessity.* Let the family  $\mathcal{P} = \{P_\alpha; \alpha \in \mathfrak{A}\}$  of probability measures defined on  $(R, \mathcal{B}(R))$  be relatively compact but not tight. Then there is an  $\varepsilon > 0$  such that for every compact  $K \subseteq R$

$$\sup_{\alpha} P_\alpha(R \setminus K) > \varepsilon,$$

and therefore, for each interval  $I = (a, b)$ ,

$$\sup_{\alpha} P_\alpha(R \setminus I) > \varepsilon.$$

It follows that for every interval  $I_n = (-n, n)$ ,  $n \geq 1$ , there is a measure  $P_{\alpha_n}$  such that

$$P_{\alpha_n}(R \setminus I_n) > \varepsilon.$$

Since the original family  $\mathcal{P}$  is relatively compact, we can select from  $\{P_{\alpha_n}\}_{n \geq 1}$  a subsequence  $\{P_{\alpha_{n_k}}\}$  such that  $P_{\alpha_{n_k}} \xrightarrow{w} Q$ , where  $Q$  is a probability measure. Then, by the equivalence of conditions (I) and (II) in Theorem 1 of §1, we have

$$\lim_{k \rightarrow \infty} P_{\alpha_{n_k}}(R \setminus I_n) \leq Q(R \setminus I_n) \quad (2)$$

for every  $n \geq 1$ . But  $Q(R \setminus I_n) \downarrow 0$ ,  $n \rightarrow \infty$ , and the left side of (2) exceeds  $\varepsilon > 0$ . This contradiction shows that relatively compact sets are tight.

To prove the sufficiency we need a general result (Helly's theorem) on the *sequential compactness* of families of generalized distribution functions (Subsection 2 of §3 of Chapter II).

Let  $\mathcal{J} = \{G\}$  be the collection of generalized distribution functions  $G = G(x)$  that satisfy:

- (1)  $G(x)$  is nondecreasing;
- (2)  $0 \leq G(-\infty)$ ,  $G(+\infty) \leq 1$ ;
- (3)  $G(x)$  is continuous on the right.

Then  $\mathcal{J}$  clearly contains the class of distribution functions  $\mathcal{F} = \{F\}$  for which  $F(-\infty) = 0$  and  $F(+\infty) = 1$ .

**Theorem 2 (Helly's Theorem).** *The class  $\mathcal{J} = \{G\}$  of generalized distribution functions is sequentially compact, i.e., for every sequence  $\{G_n\}$  of functions from  $\mathcal{J}$  we can find a function  $G \in \mathcal{J}$  and a sequence  $\{n_k\} \subseteq \{n\}$  such that*

$$G_{n_k}(x) \rightarrow G(x), \quad k \rightarrow \infty,$$

for every point  $x$  of the set  $P_c(G)$  of points of continuity of  $G = G(x)$ .

**PROOF.** Let  $T = \{x_1, x_2, \dots\}$  be a countable dense subset of  $R$ . Since the sequence of numbers  $\{G_n(x_1)\}$  is bounded, there is a subsequence  $N_1 = \{n_1^{(1)}, n_2^{(1)}, \dots\}$  such that  $G_{n_i^{(1)}}(x_1)$  approaches a limit  $g_1$  as  $i \rightarrow \infty$ . Then we extract from  $N_1$  a subsequence  $N_2 = \{n_1^{(2)}, n_2^{(2)}, \dots\}$  such that  $G_{n_i^{(2)}}(x_2)$  approaches a limit  $g_2$  as  $i \rightarrow \infty$ ; and so on.

On the set  $T \subseteq R$  we can define a function  $G_T(x)$  by

$$G_T(x_i) = g_i, \quad x_i \in T,$$

and consider the "Cantor" diagonal sequence  $N = \{n_1^{(1)}, n_2^{(2)}, \dots\}$ . Then, for each  $x_i \in T$ , as  $m \rightarrow \infty$ , we have

$$G_{n_m^{(m)}}(x_i) \rightarrow G_T(x_i).$$

Finally, let us define  $G = G(x)$  for all  $x \in R$  by putting

$$G(x) = \inf\{G_T(y): y \in T, y > x\}. \quad (3)$$

We claim that  $G = G(x)$  is the required function and  $G_{n_m^{(m)}}(x) \rightarrow G(x)$  at all points  $x$  of continuity of  $G$ .

Since all the functions  $G_n$  under consideration are nondecreasing, we have  $G_{n_m^{(m)}}(x) \leq G_{n_m^{(m)}}(y)$  for all  $x$  and  $y$  that belong to  $T$  and satisfy the inequality  $x \leq y$ . Hence for such  $x$  and  $y$ ,

$$G_T(x) \leq G_T(y).$$

It follows from this and (3) that  $G = G(x)$  is nondecreasing.

Now let us show that it is continuous on the right. Let  $x_k \downarrow x$  and  $d = \lim_k G(x_k)$ . Clearly  $G(x) \leq d$ , and we have to show that actually  $G(x) = d$ . Suppose the contrary, that is, let  $G(x) < d$ . It follows from (3) that there is a  $y \in T$ ,  $x < y$ , such that  $G_T(y) < d$ . But  $x < x_k < y$  for sufficiently large  $k$ , and therefore  $G(x_k) \leq G_T(y) < d$  and  $\lim G(x_k) < d$ , which contradicts  $d = \lim_k G(x_k)$ . Thus we have constructed a function  $G$  that belongs to  $\mathcal{J}$ .

We now establish that  $G_{n_{H^m}}(x^0) \rightarrow G(x^0)$  for every  $x^0 \in P_C(G)$ .

If  $x^0 < y \in T$ , then

$$\overline{\lim}_m G_{n_{H^m}}(x^0) \leq \overline{\lim}_m G_{n_{H^m}}(y) = G_T(y),$$

whence

$$\overline{\lim}_m G_{n_{H^m}}(x^0) \leq \inf\{G_T(y) : y > x^0, y \in T\} = G(x^0). \quad (4)$$

On the other hand, let  $x^1 < y < x^0$ ,  $y \in T$ . Then

$$G(x^1) \leq G_T(y) = \lim_m G_{n_{H^m}}(y) = \underline{\lim}_m G_{n_{H^m}}(y) \leq \underline{\lim}_m G_{n_{H^m}}(x^0).$$

Hence if we let  $x^1 \uparrow x^0$  we find that

$$G(x^0 -) \leq \underline{\lim}_m G_{n_{H^m}}(x^0). \quad (5)$$

But if  $G(x^0 -) = G(x^0)$  we can infer from (4) and (5) that  $G_{n_{H^m}}(x^0) \rightarrow G(x^0)$ ,  $m \rightarrow \infty$ .

This completes the proof of the theorem.

We can now complete the proof of Theorem 1.

*Sufficiency.* Let the family  $\mathcal{P}$  be tight and let  $\{P_n\}$  be a sequence of probability measures from  $\mathcal{P}$ . Let  $\{F_n\}$  be the corresponding sequence of distribution functions.

By Helly's theorem, there are a subsequence  $\{F_{n_k}\} \subseteq \{F_n\}$  and a generalized distribution function  $G \in \mathcal{J}$  such that  $F_{n_k}(x) \rightarrow G(x)$  for  $x \in P_C(G)$ . Let us show that because  $\mathcal{P}$  was assumed tight, the function  $G = G(x)$  is in fact a genuine distribution function ( $G(-\infty) = 0$ ,  $G(+\infty) = 1$ ).

Take  $\varepsilon > 0$ , and let  $I = (a, b]$  be the interval for which

$$\sup_n P_n(R \setminus I) < \varepsilon,$$

or, equivalently,

$$1 - \varepsilon \leq P_n(a, b], \quad n \geq 1.$$

Choose points  $a', b' \in P_C(G)$  such that  $a' < a$ ,  $b' > b$ . Then  $1 - \varepsilon \leq P_{n_k}(a, b] \leq P_{n_k}(a', b'] = F_{n_k}(b') - F_{n_k}(a') \rightarrow G(b') - G(a')$ . It follows that  $G(+\infty) - G(-\infty) = 1$ , and since  $0 \leq G(-\infty) \leq G(+\infty) \leq 1$ , we have  $G(-\infty) = 0$  and  $G(+\infty) = 1$ .

Therefore the limit function  $G = G(x)$  is a distribution function and  $F_{n_k} \Rightarrow G$ . Together with Theorem 2 of §1 this shows that  $P_{n_k} \xrightarrow{w} Q$ , where  $Q$  is the probability measure corresponding to the distribution function  $G$ .

This completes the proof of Theorem 1.

#### 4. PROBLEMS

1. Carry out the proofs of Theorems 1 and 2 for  $R^n$ ,  $n \geq 2$ .
2. Let  $P_\alpha$  be a Gaussian measure on the real line, with parameters  $m_\alpha$  and  $\sigma_\alpha^2$ ,  $\alpha \in \mathfrak{A}$ . Show that the family  $\mathscr{P} = \{P_\alpha; \alpha \in \mathfrak{A}\}$  is tight if and only if

$$|m_\alpha| \leq a, \quad \sigma_\alpha^2 \leq b, \quad \alpha \in \mathfrak{A}.$$

3. Construct examples of tight and nontight families  $\mathscr{P} = \{P_\alpha; \alpha \in \mathfrak{A}\}$  of probability measures defined on  $(R^\infty, \mathscr{B}(R^\infty))$ .

### §3. Proofs of Limit Theorems by the Method of Characteristic Functions

1. The proofs of the first limit theorems of probability theory—the law of large numbers, and the De Moivre–Laplace and Poisson theorems for Bernoulli schemes—were based on direct analysis of the limit functions of the distributions  $F_n$ , which are expressed rather simply in terms of binomial probabilities. (In the Bernoulli scheme, we are adding random variables that take only two values, so that in principle we can find  $F_n$  explicitly.) However, it is practically impossible to apply a similar direct method to the study of more complicated random variables.

The first step in proving limit theorems for sums of arbitrarily distributed random variables was taken by Chebyshev. The inequality that he discovered, and which is now known as Chebyshev's inequality, not only makes it possible to give an elementary proof of James Bernoulli's law of large numbers, but also lets us establish very general conditions for this law to hold, when stated in the form

$$P\left\{\left|\frac{S_n}{n} - \frac{ES_n}{n}\right| \geq \varepsilon\right\} \rightarrow 0, \quad n \rightarrow \infty, \quad \text{every } \varepsilon > 0, \quad (1)$$

for sums  $S_n = \xi_1 + \cdots + \xi_n$ ,  $n \geq 1$ , of independent random variables. (See Problem 2.)

Furthermore, Chebyshev created (and Markov perfected) the “method of moments” which made it possible to show that the conclusion of the De Moivre–Laplace theorem, written in the form

$$P\left\{\frac{S_n - ES_n}{\sqrt{VS_n}} \leq x\right\} \rightarrow \frac{1}{\sqrt{2\pi}} \int_{-\infty}^x e^{-u^2/2} du, \quad (2)$$

is "universal," in the sense that it is valid under very general hypotheses concerning the nature of the random variables. For this reason it is known as the *central limit theorem* of probability theory.

Somewhat later Lyapunov proposed a different method for proving the central limit theorem, based on the idea (which goes back to Laplace) of the characteristic function of a probability distribution. Subsequent developments have shown that Lyapunov's method of characteristic functions is extremely effective for proving the most diverse limit theorems. Consequently it has been extensively developed and widely applied.

In essence, the method is as follows.

2. We already know (Chapter II, §12) that there is a one-to-one correspondence between distribution functions and characteristic functions. Hence we can study the properties of distribution functions by using the corresponding characteristic functions. It is a fortunate circumstance that weak convergence  $F_n \xrightarrow{w} F$  of distributions is equivalent to pointwise convergence  $\varphi_n \rightarrow \varphi$  of the corresponding characteristic functions. Moreover, we have the following result, which provides the basic method of proving theorems on weak convergence for distributions on the real line.

**Theorem 1** (Continuity Theorem). *Let  $\{F_n\}$  be a sequence of distribution functions  $F_n = F_n(x)$ ,  $x \in R$ , and let  $\{\varphi_n\}$  be the corresponding sequence of characteristic functions,*

$$\varphi_n(t) = \int_{-\infty}^{\infty} e^{itx} dF_n(x), \quad t \in R.$$

- (1) *If  $F_n \xrightarrow{w} F$ , where  $F = F(x)$  is a distribution function, then  $\varphi_n(t) \rightarrow \varphi(t)$ ,  $t \in R$ , where  $\varphi(t)$  is the characteristic function of  $F = F(x)$ .*
- (2) *If  $\lim_n \varphi_n(t)$  exists for each  $t \in R$  and  $\varphi(t) = \lim_n \varphi_n(t)$  is continuous at  $t = 0$ , then  $\varphi(t)$  is the characteristic function of a probability distribution  $F = F(x)$ , and*

$$F_n \xrightarrow{w} F.$$

The proof of conclusion (1) is an immediate consequence of the definition of weak convergence, applied to the functions  $\operatorname{Re} e^{itx}$  and  $\operatorname{Im} e^{itx}$ .

The proof of (2) requires some preliminary propositions.

**Lemma 1.** *Let  $\{P_n\}$  be a tight family of probability measures. Suppose that every weakly convergent subsequence  $\{P_{n_k}\}$  of  $\{P_n\}$  converges to the same probability measure  $P$ . Then the whole sequence  $\{P_n\}$  converges to  $P$ .*

**PROOF.** Suppose that  $P_n \not\xrightarrow{w} P$ . Then there is a bounded continuous function  $f = f(x)$  such that

$$\int_R f(x) P_n(dx) \not\xrightarrow{w} \int_R f(x) P(dx).$$

It follows that there exist  $\varepsilon > 0$  and an infinite sequence  $\{n'\} \subseteq \{n\}$  such that

$$\left| \int_R f(x) P_{n'}(dx) - \int_R f(x) P(dx) \right| \geq \varepsilon > 0. \quad (3)$$

By Prokhorov's theorem (§2) we can select a subsequence  $\{P_{n''}\}$  of  $\{P_{n'}\}$  such that  $P_{n''} \xrightarrow{w} Q$ , where  $Q$  is a probability measure.

By the hypotheses of the lemma,  $Q = P$ , and therefore

$$\int_R f(x) P_{n''}(dx) \rightarrow \int_R f(x) P(dx),$$

which leads to a contradiction with (3). This completes the proof of the lemma.

**Lemma 2.** Let  $\{P_n\}$  be a tight family of probability measures on  $(R, \mathcal{B}(R))$ . A necessary and sufficient condition for the sequence  $\{P_n\}$  to converge weakly to a probability measure is that for each  $t \in R$  the limit  $\lim_n \varphi_n(t)$  exists, where  $\varphi_n(t)$  is the characteristic function of  $P_n$ :

$$\varphi_n(t) = \int_R e^{itx} P_n(dx).$$

**PROOF.** If  $\{P_n\}$  is tight, by Prokhorov's theorem there is a subsequence  $\{P_{n'}\}$  and a probability measure  $P$  such that  $P_{n'} \xrightarrow{w} P$ . Suppose that the whole sequence  $\{P_n\}$  does not converge to  $P$  ( $P_n \not\xrightarrow{w} P$ ). Then, by Lemma 1, there is a subsequence  $\{P_{n''}\}$  and a probability measure  $Q$  such that  $P_{n''} \xrightarrow{w} Q$ , and  $P \neq Q$ .

Now we use the existence of  $\lim_n \varphi_n(t)$  for each  $t \in R$ . Then

$$\lim_{n'} \int_R e^{itx} P_{n'}(dx) = \lim_{n''} \int_R e^{itx} P_{n''}(dx)$$

and therefore

$$\int_R e^{itx} P(dx) = \int_R e^{itx} Q(dx), \quad t \in R.$$

But the characteristic function determines the distribution uniquely (Theorem 2, §12, Chapter II). Hence  $P = Q$ , which contradicts the assumption that  $P_n \not\xrightarrow{w} P$ .

The converse part of the lemma follows immediately from the definition of weak convergence.

The following lemma estimates the "tails" of a distribution function in terms of the behavior of its characteristic function in a neighborhood of zero.

**Lemma 3.** Let  $F = F(x)$  be a distribution function on the real line and let



$\varphi = \varphi(t)$  be its characteristic function. Then there is a constant  $K > 0$  such that for every  $a > 0$

$$\int_{|x| \geq 1/a} dF(x) \leq \frac{K}{a} \int_0^a [1 - \operatorname{Re} \varphi(t)] dt. \quad (4)$$

PROOF. Since  $\operatorname{Re} \varphi(t) = \int_{-\infty}^{\infty} \cos tx dF(x)$ , we find by Fubini's theorem that

$$\begin{aligned} \frac{1}{a} \int_0^a [1 - \operatorname{Re} \varphi(t)] dt &= \frac{1}{a} \int_0^a \left[ \int_{-\infty}^{\infty} (1 - \cos tx) dF(x) \right] dt \\ &= \int_{-\infty}^{\infty} \left[ \frac{1}{a} \int_0^a (1 - \cos tx) dt \right] dF(x) \\ &= \int_{-\infty}^{\infty} \left( 1 - \frac{\sin ax}{ax} \right) dF(x) \\ &\geq \inf_{|y| \geq 1} \left( 1 - \frac{\sin y}{y} \right) \cdot \int_{|ax| \geq 1} dF(x) \\ &= \frac{1}{K} \int_{|x| \geq 1/a} dF(x), \end{aligned}$$

where

$$\frac{1}{K} = \inf_{|y| \geq 1} \left( 1 - \frac{\sin y}{y} \right) = 1 - \sin 1 \geq \frac{1}{7},$$

so that (4) holds with  $K = 7$ . This establishes the lemma.

Proof of conclusion (2) of Theorem 1. Let  $\varphi_n(t) \rightarrow \varphi(t)$ ,  $n \rightarrow \infty$ , where  $\varphi(t)$  is continuous at 0. Let us show that it follows that the family of probability measures  $\{P_n\}$  is tight, where  $P_n$  is the measure corresponding to  $F_n$ .

By (4) and the dominated convergence theorem,

$$\begin{aligned} P_n \left\{ R \setminus \left( -\frac{1}{a}, \frac{1}{a} \right) \right\} &= \int_{|x| \geq 1/a} dF_n(x) \leq \frac{K}{a} \int_0^a [1 - \operatorname{Re} \varphi_n(t)] dt \\ &\rightarrow \frac{K}{a} \int_0^a [1 - \operatorname{Re} \varphi(t)] dt \end{aligned}$$

as  $n \rightarrow \infty$ .

Since, by hypothesis,  $\varphi(t)$  is continuous at 0 and  $\varphi(0) = 1$ , for every  $\varepsilon > 0$  there is an  $a > 0$  such that

$$P_n \left\{ R \setminus \left( -\frac{1}{a}, \frac{1}{a} \right) \right\} \leq \varepsilon$$

for all  $n \geq 1$ . Consequently  $\{P_n\}$  is tight, and by Lemma 2 there is a prob-

ability measure  $P$  such that

$$P_n \xrightarrow{w} P.$$

Hence

$$\varphi_n(t) = \int_{-\infty}^{\infty} e^{itx} P_n(dx) \rightarrow \int_{-\infty}^{\infty} e^{itx} P(dx),$$

but also  $\varphi_n(t) \rightarrow \varphi(t)$ . Therefore  $\varphi(t)$  is the characteristic function of  $P$ .

This completes the proof of the theorem.

**Corollary.** Let  $\{F_n\}$  be a sequence of distribution functions and  $\{\varphi_n\}$  the corresponding sequence of characteristic functions. Also let  $F$  be a distribution function and  $\varphi$  its characteristic function. Then  $F_n \xrightarrow{w} F$  if and only if  $\varphi_n(t) \rightarrow \varphi(t)$  for all  $t \in R$ .

**Remark.** Let  $\eta, \eta_1, \eta_2, \dots$  be random variables and  $F_{\eta_n} \xrightarrow{w} F_{\eta}$ . In accordance with the definition of §10 of Chapter II, we then say that the random variables  $\eta_1, \eta_2, \dots$  converge to  $\eta$  in distribution, and write  $\eta_n \xrightarrow{d} \eta$ .

Since this notation is self-explanatory, we shall frequently use it instead of  $F_{\eta_n} \xrightarrow{w} F_{\eta}$  when stating limit theorems.

3. In the next section, Theorem 1 will be applied to prove the central limit theorem for independent but not identically distributed random variables. In the present section we shall merely apply the method of characteristic functions to prove some simple limit theorems.

**Theorem 2 (Law of Large Numbers).** Let  $\xi_1, \xi_2, \dots$  be a sequence of independent identically distributed random variables with  $E|\xi_1| < \infty$ ,  $S_n = \xi_1 + \dots + \xi_n$  and  $E\xi_1 = m$ . Then  $S_n/n \xrightarrow{P} m$ , that is, for every  $\varepsilon > 0$

$$P\left\{\left|\frac{S_n}{n} - m\right| \geq \varepsilon\right\} \rightarrow 0, \quad n \rightarrow \infty.$$

**PROOF.** Let  $\varphi(t) = Ee^{it\xi_1}$  and  $\varphi_{S_n/n}(t) = Ee^{itS_n/n}$ . Since the random variables are independent, we have

$$\varphi_{S_n/n}(t) = \left[\varphi\left(\frac{t}{n}\right)\right]^n$$

by (II.12.6). But according to (II.12.14)

$$\varphi(t) = 1 + itm + o(t), \quad t \rightarrow 0.$$

Therefore for each given  $t \in R$

$$\varphi\left(\frac{t}{n}\right) = 1 + i \frac{t}{n} m + o\left(\frac{1}{n}\right), \quad n \rightarrow \infty,$$

and therefore

$$\varphi_{S_n/n}(t) = \left[1 + i \frac{t}{n} m + o\left(\frac{1}{n}\right)\right]^n \rightarrow e^{itm}.$$

The function  $\varphi(t) = e^{itm}$  is continuous at 0 and is the characteristic function of the degenerate probability distribution that is concentrated at  $m$ . Therefore

$$\frac{S_n}{n} \xrightarrow{d} m,$$

and consequently (see Problem 7, §10, Chapter II)

$$\frac{S_n}{n} \xrightarrow{P} m.$$

This completes the proof of the theorem.

**Theorem 3** (Central Limit Theorem for *Independent Identically Distributed* Random Variables). Let  $\xi_1, \xi_2, \dots$  be a sequence of independent identically distributed (nondegenerate) random variables with  $E\xi_1^2 < \infty$  and  $S_n = \xi_1 + \dots + \xi_n$ . Then as  $n \rightarrow \infty$

$$P\left\{\frac{S_n - ES_n}{\sqrt{VS_n}} \leq x\right\} \rightarrow \Phi(x), \quad x \in R, \quad (5)$$

where

$$\Phi(x) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^x e^{-u^2/2} du.$$

PROOF. Let  $E\xi_1 = m$ ,  $V\xi_1 = \sigma^2$  and

$$\varphi(t) = Ee^{it(\xi_1 - m)}.$$

Then if we put

$$\varphi_n(t) = E \exp\left\{it \frac{S_n - ES_n}{\sqrt{VS_n}}\right\},$$

we find that

$$\varphi_n(t) = \left[\varphi\left(\frac{t}{\sigma\sqrt{n}}\right)\right]^n.$$

But by (II.12.14)

$$\varphi(t) = 1 - \frac{\sigma^2 t^2}{2} + o(t^2), \quad t \rightarrow 0.$$

Therefore

$$\varphi_n(t) = \left[ 1 - \frac{\sigma^2 t^2}{2\sigma^2 n} + o\left(\frac{1}{n}\right) \right]^n \rightarrow e^{-t^2/2},$$

as  $n \rightarrow \infty$  for fixed  $t$ .

The function  $e^{-t^2/2}$  is the characteristic function of a random variable (denoted by  $\mathcal{N}(0, 1)$ ) with mean zero and unit variance. This, by Theorem 1, also establishes (5). In accordance with the remark in Theorem 1, this can also be written in the form

$$\frac{S_n - \mathbf{E}S_n}{\sqrt{\mathbf{V}S_n}} \xrightarrow{d} \mathcal{N}(0, 1). \quad (6)$$

This completes the proof of the theorem.

The preceding two theorems have dealt with the behavior of the probabilities of (normalized and symmetrized) sums of independent and identically distributed random variables. However, in order to state Poisson's theorem (§6, Chapter I) we have to use a more general model.

Let us suppose that for each  $n \geq 1$  we are given a sequence of independent random variables  $\xi_{n1}, \dots, \xi_{nn}$ . In other words, let there be given a triangular array

$$\begin{pmatrix} \xi_{11} \\ \xi_{21}, \xi_{22} \\ \xi_{31}, \xi_{32}, \xi_{33} \end{pmatrix}$$

of random variables, those in each row being independent. Put  $S_n = \xi_{n1} + \dots + \xi_{nn}$ .

**Theorem 4 (Poisson's Theorem).** *For each  $n \geq 1$  let the independent random variables  $\xi_{n1}, \dots, \xi_{nn}$  be such that*

$$P(\xi_{nk} = 1) = p_{nk}, \quad P(\xi_{nk} = 0) = q_{nk},$$

$p_{nk} + q_{nk} = 1$ . Suppose that

$$\max_{1 \leq k \leq n} p_{nk} \rightarrow 0, \quad n \rightarrow \infty,$$

and  $\sum_{k=1}^n p_{nk} \rightarrow \lambda > 0, n \rightarrow \infty$ . Then, for each  $m = 0, 1, \dots$ ,

$$P(S_n = m) \rightarrow \frac{e^{-\lambda} \lambda^m}{m!}, \quad n \rightarrow \infty. \quad (7)$$

**PROOF.** Since

$$\mathbf{E}e^{it\xi_{nk}} = p_{nk}e^{it} + q_{nk}$$

for  $1 \leq k \leq n$ , by our assumptions we have

$$\begin{aligned}\varphi_{S_n}(t) &= \mathbf{E} e^{itS_n} = \prod_{k=1}^n (p_{nk} e^{it} + q_{nk}) \\ &= \prod_{k=1}^n (1 + p_{nk}(e^{it} - 1)) \rightarrow \exp\{\lambda(e^{it} - 1)\}, \quad n \rightarrow \infty.\end{aligned}$$

The function  $\varphi(t) = \exp\{\lambda(e^{it} - 1)\}$  is the characteristic function of the Poisson distribution (II.12.11), so that (7) is established.

If  $\pi(\lambda)$  denotes a Poisson random variable with parameter  $\lambda$ , then (7) can be written like (6), in the form

$$S_n \xrightarrow{d} \pi(\lambda).$$

This completes the proof of the theorem.

#### 4. PROBLEMS

1. Prove Theorem 1 for  $R^n$ ,  $n \geq 2$ .
2. Let  $\xi_1, \xi_2, \dots$  be a sequence of independent random variables with finite means  $\mathbf{E}|\xi_n|$  and variances  $\mathbf{V}\xi_n$  such that  $\mathbf{V}\xi_n \leq K < \infty$ , where  $K$  is a constant. Use Chebyshev's inequality to prove the law of large numbers (1).
3. Show, as a corollary to Theorem 1, that the family  $\{\varphi_n\}$  is *uniformly continuous* and that  $\varphi_n \rightarrow \varphi$  uniformly on every finite interval.
4. Let  $\xi_n, n \geq 1$ , be random variables with characteristic functions  $\varphi_{\xi_n}(t)$ ,  $n \geq 1$ . Show that  $\xi_n \xrightarrow{d} 0$  if and only if  $\varphi_{\xi_n}(t) \rightarrow 1, n \rightarrow \infty$ , in some neighborhood of  $t = 0$ .
5. Let  $X_1, X_2, \dots$  be a sequence of independent random vectors (with values in  $R^k$ ) with mean zero and (finite) covariance matrix  $\Gamma$ . Show that

$$\frac{X_1 + \dots + X_n}{\sqrt{n}} \xrightarrow{d} \mathcal{N}(0, \Gamma).$$

(Compare Theorem 3.)

## §4. Central Limit Theorem for Sums of Independent Random Variables.

### I. The Lindeberg Condition

1. In this section, the central limit theorem for (normalized and centralized) sums of independent random variables  $\xi_1, \xi_2, \dots$  will be proved under the traditional hypothesis that the *classical Lindeberg condition* is satisfied. In the next section, we shall consider a more general situation. First, the central limit theorem will be stated in the "series form" and, second, we shall prove it under the so-called *nonclassical hypotheses*.

**Theorem 1.** Let  $\xi_1, \xi_2, \dots$  be a sequence of independent random variables with finite second moments. Let  $m_k = E\xi_k$ ,  $\sigma_k^2 = V\xi_k > 0$ ,  $S_n = \xi_1 + \dots + \xi_n$ ,  $D_n^2 = \sum_{k=1}^n \sigma_k^2$ , and let  $F_k = F_k(x)$  be the distribution function of the random variable  $\xi_k$ .

Let us suppose that the Lindeberg condition is satisfied: for every  $\varepsilon > 0$

$$(L) \quad \frac{1}{D_n^2} \sum_{k=1}^n \int_{\{x: |x - m_k| \geq \varepsilon D_n\}} (x - m_k)^2 dF_k(x) \rightarrow 0, \quad n \rightarrow \infty. \quad (1)$$

Then

$$\frac{S_n - ES_n}{\sqrt{VS_n}} \xrightarrow{d} \mathcal{N}(0, 1). \quad (2)$$

**PROOF.** Without loss of generality we assume that  $m_k = 0$  for  $k \geq 1$ . We set  $\varphi_k(t) = Ee^{it\xi_k}$ ,  $T_n = S_n/\sqrt{VS_n} = S_n/D_n$ ,  $\varphi_{S_n}(t) = Ee^{itS_n}$ ,  $\varphi_{T_n}(t) = Ee^{itT_n}$ .

Then

$$\varphi_{T_n}(t) = Ee^{itT_n} = Ee^{it(t/D_n)S_n} = \varphi_{S_n}\left(\frac{t}{D_n}\right) = \prod_{k=1}^n \varphi_k\left(\frac{t}{D_n}\right), \quad (3)$$

and for the proof of (2) it is sufficient (by Theorem 1 of §3) to establish that, for every  $t \in R$ ,

$$\varphi_{T_n}(t) \rightarrow e^{-t^2/2}, \quad n \rightarrow \infty. \quad (4)$$

We choose a  $t \in R$  and suppose that it is fixed throughout the proof. By the representations

$$e^{iy} = 1 + iy + \frac{\theta_1 y^2}{2},$$

$$e^{iy} = 1 + iy - \frac{y^2}{2} + \frac{\theta_2 |y|^3}{3!},$$

which are valid for all real  $y$ , with  $\theta_1 = \theta_1(y)$  and  $\theta_2 = \theta_2(y)$ , such that  $|\theta_1| \leq 1$ ,  $|\theta_2| \leq 1$ , we obtain

$$\begin{aligned} \varphi_k(t) &= Ee^{it\xi_k} = \int_{-\infty}^{\infty} e^{itx} dF_k(x) = \int_{|x| \geq \varepsilon D_n} \left(1 + itx + \frac{\theta_1 (tx)^2}{2}\right) dF_k(x) \\ &\quad + \int_{|x| < \varepsilon D_n} \left(1 + itx - \frac{t^2 x^2}{2} + \frac{\theta_2 |tx|^3}{6}\right) dF_k(x) \\ &= 1 + \frac{t^2}{2} \int_{|x| \geq \varepsilon D_n} \theta_1 x^2 dF_k(x) - \frac{t^2}{2} \int_{|x| < \varepsilon D_n} x^2 dF_k(x) \\ &\quad + \frac{|t|^3}{6} \int_{|x| < \varepsilon D_n} \theta_2 |x|^3 dF_k(x) \end{aligned}$$

(here we have also used the fact that, by hypothesis,  $m_k = \int_{-\infty}^{\infty} x dF_k(x) = 0$ ).

Consequently,

$$\begin{aligned} \varphi_k\left(\frac{1}{D_n}\right) &= 1 - \frac{t^2}{2D_n^2} \int_{|x| < \varepsilon D_n} x^2 dF_k(x) + \frac{t^2}{2D_n^2} \int_{|x| < \varepsilon D_n} \theta_1 x^2 dF_k(x) \\ &\quad + \frac{|t|^3}{6D_n^3} \int_{|x| < \varepsilon D_n} \theta_2 |x|^3 dF_k(x). \end{aligned} \quad (5)$$

Since

$$\left| \frac{1}{2} \int_{|x| \geq \varepsilon D_n} \theta_1 x^2 dF_k(x) \right| \leq \frac{1}{2} \int_{|x| \geq \varepsilon D_n} x^2 dF_k(x),$$

we have

$$\frac{1}{2} \int_{|x| \geq \varepsilon D_n} \theta_1 x^2 dF_k(x) = \tilde{\theta}_1 \int_{|x| \geq \varepsilon D_n} x^2 dF_k(x), \quad (6)$$

where  $\tilde{\theta}_1 = \tilde{\theta}_1(t, k, n)$  and  $|\tilde{\theta}_1| \leq 1/2$ .

In the same way,

$$\left| \frac{1}{6} \int_{|x| < \varepsilon D_n} \theta_2 |x|^3 dF_k(x) \right| \leq \frac{1}{6} \int_{|x| < \varepsilon D_n} \frac{\varepsilon D_n}{|x|} \cdot |x|^3 dF_k(x) \leq \frac{1}{6} \int_{|x| < \varepsilon D_n} \varepsilon D_n x^2 dF_k(x),$$

and therefore,

$$\frac{1}{6} \int_{|x| < \varepsilon D_n} \theta_2 |x|^3 dF_k(x) = \tilde{\theta}_2 \int_{|x| < \varepsilon D_n} \varepsilon D_n x^2 dF_k(x), \quad (7)$$

where  $\tilde{\theta}_2 = \tilde{\theta}_2(t, k, n)$  and  $|\tilde{\theta}_2| \leq 1/6$ .

We now set

$$\begin{aligned} A_{kn} &= \frac{1}{D_n^2} \int_{|x| < \varepsilon D_n} x^2 dF_k(x), \\ B_{kn} &= \frac{1}{D_n^2} \int_{|x| \geq \varepsilon D_n} x^2 dF_k(x). \end{aligned}$$

Then, by (5)–(7),

$$\varphi_k\left(\frac{t}{D_n}\right) = 1 - \frac{t^2 A_{kn}}{2} + t^2 \tilde{\theta}_1 B_{kn} + |t|^3 \varepsilon \tilde{\theta}_2 A_{kn} = 1 + C_{kn}. \quad (8)$$

We note that

$$\sum_{k=1}^n (A_{kn} + B_{kn}) = 1 \quad (9)$$

and by (1)

$$\sum_{k=1}^n B_{kn} \rightarrow 0, \quad n \rightarrow \infty. \quad (10)$$

Consequently, for sufficiently large  $n$ ,

$$\max_{1 \leq k \leq n} |C_{kn}| \leq t^2 \varepsilon^2 + \varepsilon |t|^3 \quad (11)$$

and

$$\sum_{k=1}^n |C_{kn}| \leq t^2 + \varepsilon |t|^3. \quad (12)$$

We now appeal to the fact that, for any complex numbers  $z$  with  $|z| \leq 1/2$ ,

$$\ln(1+z) = z + \theta |z|^2,$$

where  $\theta = \theta(z)$  with  $|\theta| \leq 1$  and  $\ln$  denotes the *principal value* of the logarithm.\* Then, for sufficiently large  $n$ , it follows from (8) and (11) that, for sufficiently small  $\varepsilon > 0$ ,

$$\ln \varphi_k \left( \frac{t}{D_n} \right) = \ln(1 + C_{kn}) = C_{kn} + \theta_{kn} |C_{kn}|^2,$$

where  $|\theta_{kn}| \leq 1$ . Consequently, by (3),

$$\frac{t^2}{2} + \ln \varphi_{T_n}(t) = \frac{t^2}{2} + \sum_{k=1}^n \ln \varphi_k \left( \frac{t}{D_n} \right) = \frac{t^2}{2} + \sum_{k=1}^n C_{kn} + \sum_{k=1}^n \theta_{kn} |C_{kn}|^2.$$

But

$$\frac{t^2}{2} + \sum_{k=1}^n C_{kn} = \frac{t^2}{2} \left( 1 - \sum_{k=1}^n A_{kn} \right) + t^2 \sum_{k=1}^n \tilde{\theta}_1(t, k, n) B_{kn} + \varepsilon |t|^3 \sum_{k=1}^n \tilde{\theta}_2(t, k, n) A_{kn},$$

and by (9) and (10), for any  $\delta > 0$  we can find numbers  $n_0$  and  $\varepsilon > 0$ , with  $n_0$  so large that for all  $n \geq n_0$

$$\left| \frac{t^2}{2} + \sum_{k=1}^n C_{kn} \right| \leq \frac{\delta}{2}.$$

In addition, by (11) and (12), we can find a positive number  $\varepsilon$  such that

$$\left| \sum_{k=1}^n \theta_{kn} |C_{kn}|^2 \right| \leq \max_{1 \leq k \leq n} |C_{kn}| \sum_{k=1}^n |C_{kn}| \leq (t^2 \varepsilon^2 + \varepsilon |t|^3)(t^2 + \varepsilon |t|^3).$$

Therefore, for sufficiently large  $n$ , we can choose  $\varepsilon > 0$  so that

$$\left| \sum_{k=1}^n \theta_{kn} |C_{kn}|^2 \right| \leq \frac{\delta}{2}$$

and consequently,

$$\left| \frac{t^2}{2} + \ln \varphi_{T_n}(t) \right| \leq \delta.$$

\* The principal value  $\ln z$  of the complex number  $z$  is defined by  $\ln z = \ln |z| + i \arg z$ ,  $-\pi < \arg z \leq \pi$ .



Therefore, for any real  $t$ ,

$$\varphi_{T_n}(t)e^{t^2/2} \rightarrow 1, \quad n \rightarrow \infty$$

and hence,

$$\varphi_{T_n}(t) \rightarrow e^{-t^2/2}, \quad n \rightarrow \infty.$$

This completes the proof of the theorem.

2. We turn our attention to some special cases in which the Lindeberg condition (1) is satisfied and consequently, the central limit theorem is valid.

a) Let the "Lyapunov condition" be satisfied: for some  $\delta > 0$

$$\frac{1}{D_n^{2+\delta}} \sum_{k=1}^n \mathbb{E} |\xi_k - m_k|^{2+\delta} \rightarrow 0, \quad n \rightarrow \infty. \quad (13)$$

Let  $\varepsilon > 0$ ; then

$$\begin{aligned} \mathbb{E} |\xi_k - m_k|^{2+\delta} &= \int_{-\infty}^{\infty} |x - m_k|^{2+\delta} dF_k(x) \\ &\geq \int_{\{x: |x - m_k| \geq \varepsilon D_n\}} |x - m_k|^{2+\delta} dF_k(x) \\ &\geq \varepsilon^\delta D_n^\delta \int_{\{x: |x - m_k| \geq \varepsilon D_n\}} (x - m_k)^2 dF_k(x) \end{aligned}$$

and therefore,

$$\frac{1}{D_n^2} \sum_{k=1}^n \int_{\{x: |x - m_k| \geq \varepsilon D_n\}} (x - m_k)^2 dF_k(x) \leq \frac{1}{\varepsilon^\delta} \cdot \frac{1}{D_n^{2+\delta}} \sum_{k=1}^n \mathbb{E} |\xi_k - m_k|^{2+\delta}.$$

Consequently, the Lyapunov condition implies the Lindeberg condition.

b) Let  $\xi_1, \xi_2, \dots$  be independent identically distributed random variables with  $m = \mathbb{E} \xi_1$  and variance  $0 < \sigma^2 \equiv \mathbb{V} \xi_1 < \infty$ . Then

$$\begin{aligned} &\frac{1}{D_n^2} \sum_{k=1}^n \int_{\{x: |x - m| \geq \varepsilon D_n\}} |x - m|^2 dF_k(x) \\ &= \frac{n}{n\sigma^2} \int_{\{x: |x - m| \geq \varepsilon \sigma^2 \sqrt{n}\}} |x - m|^2 dF_1(x) \rightarrow 0, \end{aligned}$$

since  $\{x: |x - m| \geq \varepsilon \sigma^2 \sqrt{n}\} \downarrow \emptyset, n \rightarrow \infty$ , and  $\sigma^2 = \mathbb{E} |\xi_1 - m|^2 < \infty$ .

Therefore, the Lindeberg condition is satisfied and consequently, Theorem 3 of §3 follows from the proof of Theorem 1.

c) Let  $\xi_1, \xi_2, \dots$  be independent random variables such that for all  $n \geq 1$

$$|\xi_k| \leq K < \infty,$$

where  $K$  is a constant and  $D_n \rightarrow \infty, n \rightarrow \infty$ .

Then by Chebyshev's inequality

$$\begin{aligned} \int_{\{x: |x-m_k| \geq \varepsilon D_n\}} |x - m_k|^2 dF_k(x) &= E[(\xi_k - m_k)^2 I(|\xi_k - m_k| \geq \varepsilon D_n)] \\ &\leq (2K)^2 P\{|\xi_k - m_k| \geq \varepsilon D_n\} \leq (2K)^2 \frac{\sigma_k^2}{\varepsilon^2 D_n^2} \end{aligned}$$

and therefore,

$$\frac{1}{D_n^2} \sum_{k=1}^n \int_{\{x: |x-m_k| \geq \varepsilon D_n\}} |x - m_k|^2 dF_k(x) \leq \frac{(2K)^2}{\varepsilon^2 D_n^2} \rightarrow 0, \quad n \rightarrow \infty.$$

Consequently, the Lindeberg condition is satisfied again and therefore, the central limit theorem is verified.

**3. Remark 1.** Let  $T_n = (S_n - ES_n)/D_n$  and  $F_{T_n}(x) = P(T_n \leq x)$ . Then proposition (2) shows that for all  $x \in R$

$$F_{T_n}(x) \rightarrow \Phi(x), \quad n \rightarrow \infty.$$

Since  $\Phi(x)$  is continuous, the convergence here is actually uniform (problem 5, §1):

$$\sup_{x \in R} |F_{T_n}(x) - \Phi(x)| \rightarrow 0, \quad n \rightarrow \infty. \quad (14)$$

In particular, it follows that

$$P\{S_n \leq x\} - \Phi\left(\frac{x - ES_n}{D_n}\right) \rightarrow 0, \quad n \rightarrow \infty.$$

This proposition is often expressed by the statement that for sufficiently large  $n$  the value  $S_n$  is *approximately normally distributed with mean  $ES_n$  and variance  $D_n^2 \equiv VS_n$* .

**Remark 2.** Since, according to the preceding remarks,  $F_{T_n}(x) \rightarrow \Phi(x)$  as  $n \rightarrow \infty$ , uniformly in  $x$ , it is natural to raise the question of the *rate of convergence* in (14). In the case when the numbers  $\xi_1, \xi_2, \dots$  are independent and uniformly distributed with  $E|\xi_1|^3 < \infty$ , this question is answered by the *Berry-Esseen inequality*:

$$\sup_x |F_{T_n}(x) - \Phi(x)| \leq C \frac{E|\xi_1 - E\xi_1|^3}{\sigma^3 \sqrt{n}}, \quad (15)$$

where the absolute constant  $C$  satisfies the inequality

$$1/\sqrt{(2\pi)} \leq C < 0.8.$$

The proof of (15) will be given in §11.

**Remark 3.** We can state the Lindeberg condition in a somewhat different

(and more compact) version which is especially convenient in the "series form."

Let  $\xi_1, \xi_2, \dots$  be a sequence of independent random variables, with  $m_k = E\xi_k$ ,  $\sigma_k^2 = V\xi_k > 0$ ,  $D_n^2 = \sum_{k=1}^n \sigma_k^2$ , and  $\xi_{nk} = (\xi_k - m_k)/D_n$ . In this notation, condition (1) assumes the following form:

$$(L) \quad \sum_{k=1}^n E[\xi_{nk}^2 I(|\xi_{nk}| \geq \varepsilon)] \rightarrow 0, \quad n \rightarrow \infty. \quad (16)$$

If  $S_n = \xi_{n1} + \dots + \xi_{nn}$ , we have  $VS_n = 1$  and Theorem 1 can be given the following form: if (16) is satisfied, we have

$$S_n \xrightarrow{d} \mathcal{N}(0, 1).$$

In this form the central limit theorem is valid without the assumption that  $\xi_{nk}$  has the special form  $(\xi_k - m_k)/D_n$ . In fact, we have the following result whose proof is word for word the same as that of Theorem 1.

**Theorem 2.** For each  $n \geq 1$  let

$$\xi_{n1}, \xi_{n2}, \dots, \xi_{nn}$$

be a sequence of independent random variables for which  $E\xi_{nk} = 0$  and  $VS_n = 1$ , where  $S_n = \xi_{n1} + \dots + \xi_{nn}$ .

Then the Lindeberg condition (16) is a sufficient condition for the convergence  $S_n \xrightarrow{d} \mathcal{N}(0, 1)$ .

4. Since

$$\max_{1 \leq k \leq n} E\xi_{nk}^2 \leq \varepsilon^2 + \sum_{k=1}^n E[\xi_{nk}^2 I(|\xi_{nk}| \geq \varepsilon)],$$

it is clear that the Lindeberg condition (16) implies that

$$\max_{1 \leq k \leq n} E\xi_{nk}^2 \rightarrow 0, \quad n \rightarrow \infty. \quad (17)$$

It is noteworthy that when this condition is satisfied, it follows automatically from the validity of the central limit theorem that the Lindeberg condition is satisfied (Lindeberg-Feller theorem).

**Theorem 3.** For each  $n \geq 1$  let

$$\xi_{n1}, \xi_{n2}, \dots, \xi_{nn}$$

be a sequence of independent random variables for which  $E\xi_{nk} = 0$  and  $VS_n = 1$ , where  $S_n = \xi_{n1} + \dots + \xi_{nn}$ . Let (17) be satisfied. Then the Lindeberg condition is necessary and sufficient for the validity of the central limit theorem,  $S_n \rightarrow \mathcal{N}(0, 1)$ .

The sufficiency follows from Theorem 2. To establish the necessity we need the following lemma (compare Lemma 3, §3, Chapter III).

**Lemma.** Let  $\xi$  be a random variable with distribution function  $F = F(x)$ ,  $E\xi = 0$ ,  $V\xi = \gamma > 0$ . Then for every  $a > 0$

$$\int_{|x| \geq 1/a} x^2 dF(x) \leq \frac{1}{a^2} [\operatorname{Re} f(\sqrt{6}a) - 1 + 3\gamma a^2], \quad (18)$$

where  $f(t) = Ee^{it\xi}$  is the characteristic function of  $\xi$ .

**PROOF.** We have

$$\begin{aligned} \operatorname{Re} f(t) - 1 + \frac{1}{2}\gamma t^2 &= \frac{1}{2}\gamma t^2 - \int_{-\infty}^{\infty} [1 - \cos tx] dF(x) \\ &= \frac{1}{2}\gamma t^2 - \int_{|x| < 1/a} [1 - \cos tx] dF(x) - \int_{|x| \geq 1/a} [1 - \cos tx] dF(x) \\ &\geq \frac{1}{2}\gamma t^2 - \frac{1}{2}t^2 \int_{|x| < 1/a} x^2 dF(x) - 2a^2 \int_{|x| \geq 1/a} x^2 dF(x) \\ &= (\frac{1}{2}t^2 - 2a^2) \cdot \int_{|x| \geq 1/a} x^2 dF(x). \end{aligned}$$

If we set  $t = \sqrt{6}a$ , we obtain (18), as required.

We now turn to the proof of the necessity in Theorem 3.

Let

$$\begin{aligned} F_{nk}(x) &= P(\xi_{nk} \leq x), \quad f_{nk}(t) = Ee^{it\xi_{nk}}, \\ E\xi_{nk} &= 0, \quad V\xi_{nk} = \gamma_{nk} > 0, \end{aligned} \quad (19)$$

$$\sum_{k=1}^n \gamma_{nk} = 1, \quad \max_{1 \leq k \leq n} \gamma_{nk} \rightarrow 0, \quad n \rightarrow \infty.$$

Let  $\ln z$  denote the *principal value* of the logarithm of the complex number  $z$ .

Then

$$\ln \prod_{k=1}^n f_{nk}(t) = \sum_{k=1}^n \ln f_{nk}(t) + 2\pi im,$$

where  $m = m(n, t)$  is an integer. Consequently,

$$\operatorname{Re} \ln \prod_{k=1}^n f_{nk}(t) = \operatorname{Re} \sum_{k=1}^n \ln f_{nk}(t). \quad (20)$$

Since

$$\prod_{k=1}^n f_{nk}(t) \rightarrow e^{-(1/2)t^2},$$

we have

$$\left| \prod_{k=1}^n f_{nk}(t) \right| \rightarrow e^{-(1/2)t^2}.$$

Therefore,

$$\operatorname{Re} \ln \prod_{k=1}^n f_{nk}(t) = \operatorname{Re} \ln \left| \prod_{k=1}^n f_{nk}(t) \right| \rightarrow -\frac{1}{2}t^2. \quad (21)$$

For  $|z| < 1$

$$\ln(1+z) = z - \frac{z^2}{2} + \frac{z^3}{3} - \cdots \quad (22)$$

and for  $|z| \leq 1/2$

$$|\ln(1+z) - z| \leq |z|^2. \quad (23)$$

By (19), for each fixed  $t$ , all sufficiently large  $n$  and all  $k = 1, 2, \dots, n$ , we have

$$|f_{nk}(t) - 1| \leq \frac{1}{2}\gamma_{nk}t^2 \leq \frac{1}{2}. \quad (24)$$

Hence, we obtain from (23) and (24)

$$\begin{aligned} \left| \sum_{k=1}^n \{ \ln[1 + (f_{nk}(t) - 1)] - (f_{nk}(t) - 1) \} \right| &\leq \sum_{k=1}^n |f_{nk}(t) - 1|^2 \\ &\leq \frac{t^4}{4} \max_{1 \leq k \leq n} \gamma_{nk} \sum_{k=1}^n \gamma_{nk} = \frac{t^4}{4} \max_{1 \leq k \leq n} \gamma_{nk} \rightarrow 0, \quad n \rightarrow \infty, \end{aligned}$$

and consequently,

$$\left| \operatorname{Re} \sum_{k=1}^n \ln f_{nk}(t) - \operatorname{Re} \sum_{k=1}^n (f_{nk}(t) - 1) \right| \rightarrow 0, \quad n \rightarrow \infty. \quad (25)$$

It follows from (20), (21), and (25) that

$$\operatorname{Re} \sum_{k=1}^n (f_{nk}(t) - 1) + \frac{1}{2}t^2 = \sum_{k=1}^n [\operatorname{Re} f_{nk}(t) - 1 + \frac{1}{2}t^2\gamma_{nk}] \rightarrow 0, \quad n \rightarrow \infty.$$

Setting  $t = \sqrt{6a}$ , we find that for each  $a > 0$

$$\sum_{k=1}^n [\operatorname{Re} f_{nk}(\sqrt{6a}) - 1 + 3a^2\gamma_{nk}] \rightarrow 0, \quad n \rightarrow \infty. \quad (26)$$

Finally, from (18) with  $a = 1/\varepsilon$  and (26), we obtain

$$\begin{aligned} \sum_{k=1}^n E[\xi_{nk} I(|\xi_{nk}| \geq \varepsilon)] &= \sum_{k=1}^n \int_{|x| \geq \varepsilon} x^2 dF_{nk}(x) \\ &\leq \varepsilon^2 \sum_{k=1}^n [\operatorname{Re} f_{nk}(\sqrt{6a}) - 1 + 3a^2\gamma_{nk}] \rightarrow 0, \quad n \rightarrow \infty, \end{aligned}$$

which shows that the Lindeberg condition is satisfied.

## 5. PROBLEMS

1. Let  $\xi_1, \xi_2, \dots$  be a sequence of independent identically distributed random variables with  $E\xi_1 = 0$  and  $E\xi_1^2 = 1$ . Show that

$$\max \left( \frac{|\xi_1|}{\sqrt{n}}, \dots, \frac{|\xi_n|}{\sqrt{n}} \right) \xrightarrow{d} 0, \quad n \rightarrow \infty.$$

2. Show that in the Bernoulli scheme the quantity  $\sup_x |F_n(x) - \Phi(x)|$  is of order  $1/\sqrt{n}$ ,  $n \rightarrow \infty$ .

## §5. Central Limit Theorem for Sums of Independent Random Variables

### II. Nonclassical Conditions

1. It was shown in §4 that the Lindeberg condition (4.16) implies that the condition

$$\max_{1 \leq k \leq n} E\xi_{nk}^2 \rightarrow 0,$$

is satisfied. In turn, this implies the so-called condition of being *negligible in the limit* (*asymptotically infinitesimal*), that is, for every  $\varepsilon > 0$ ,

$$\max_{1 \leq k \leq n} P\{|\xi_{nk}| \geq \varepsilon\} \rightarrow 0, \quad n \rightarrow \infty.$$

Consequently, we may say that Theorems 1 and 2 of §4 provide a condition of realization of the central limit theorem for sums of independent random variables under the hypothesis of negligibility in the limit. Limit theorems in which conditions of negligibility in the limit are imposed on individual terms are usually called theorems with a classical formulation. It is easy, however, to give examples of nondegenerate random variables for which neither the Lindeberg condition nor negligibility in the limit is satisfied, but nevertheless the central limit theorem is satisfied. Here is the simplest example.

Let  $\xi_1, \xi_2, \dots$  be a sequence of independent normally distributed random variables with  $E\xi_n = 0$ ,  $V\xi_1 = 1$ ,  $V\xi_k = 2^{k-2}$ ,  $k \geq 2$ . Let  $S_n = \xi_{n1} + \dots + \xi_{nn}$  with

$$\xi_{nk} = \xi_k / \sqrt{\sum_{i=1}^n V\xi_i}.$$

It is easily verified that here neither the Lindeberg condition nor the condition of negligibility in the limit is satisfied, although the validity of the central limit theorem is evident, since  $S_n$  is normally distributed with  $ES_n = 0$  and  $VS_n = 1$ .

Theorem 1 (below) provides a sufficient (and necessary) condition for the central limit theorem without assuming the "classical" condition of negligibility in the limit. In this sense, condition (A), presented below, is an example of "nonclassical" conditions which reflect the title of this section.

2. We shall suppose that for each  $n \geq 1$  there is a given sequence ("series form") of independent random variables

$$\xi_{n1}, \xi_{n2}, \dots, \xi_{nn}$$

with  $E\xi_{nk} = 0$ ,  $V\xi_{nk} = \sigma_{nk}^2 > 0$ ,  $\sum_{k=1}^n \sigma_{nk}^2 = 1$ . Let  $S_n = \xi_{n1} + \dots + \xi_{nn}$ ,

$$F_{nk}(x) = P\{\xi_{nk} \leq x\}, \quad \Phi(x) = (2\pi)^{-1/2} \int_{-\infty}^x e^{-y^2/2} dy, \quad \Phi_{nk}(x) = \Phi\left(\frac{x}{\sigma_{nk}}\right).$$

**Theorem 1.** *To have*

$$S_n \xrightarrow{d} \mathcal{N}(0, 1), \quad (1)$$

*it is sufficient (and necessary) that for every  $\varepsilon > 0$  the condition*

$$(\Lambda) \quad \sum_{k=1}^n \int_{|x| > \varepsilon} |x| |F_{nk}(x) - \Phi_{nk}(x)| dx \rightarrow 0, \quad n \rightarrow \infty \quad (2)$$

*is satisfied.*

The following theorem clarifies the connection between condition  $(\Lambda)$  and the classical Lindeberg condition

$$(L) \quad \sum_{k=1}^n \int_{|x| > s} x^2 dF_{nk}(x) \rightarrow 0, \quad n \rightarrow \infty. \quad (3)$$

**Theorem 2. 1.** *The Lindeberg condition implies that condition  $(\Lambda)$  is satisfied:*

$$(L) \Rightarrow (\Lambda).$$

2. *If  $\max_{1 \leq k \leq n} E\xi_{nk}^2 \rightarrow 0$  as  $n \rightarrow \infty$ , the condition  $(\Lambda)$  implies the Lindeberg condition  $(L)$ :*

$$(\Lambda) \Rightarrow (L).$$

**PROOF OF THEOREM 1.** The proof of the necessity of condition  $(\Lambda)$  is rather complicated. Here we only prove the sufficiency.

Let

$$\begin{aligned} f_{nk}(t) &= Ee^{it\xi_{nk}}, & f_n(t) &= Ee^{itS_n}, \\ \varphi_{nk}(t) &= \int_{-\infty}^{\infty} e^{itx} d\Phi_{nk}(x), & \varphi(t) &= \int_{-\infty}^{\infty} e^{itx} d\Phi(x). \end{aligned}$$

It follows from §12, Chapter II, that

$$\varphi_{nk}(t) = e^{-(t^2\sigma_{nk}^2)/2}, \quad \varphi(t) = e^{-t^2/2}.$$

By the corollary of Theorem 1 of §3, we have  $S_n \xrightarrow{d} \mathcal{N}(0, 1)$  if and only if  $f_n(t) \rightarrow \varphi(t)$  as  $n \rightarrow \infty$ , for every real  $t$ .

We have

$$f_n(t) - \varphi(t) = \prod_{k=1}^n f_{nk}(t) - \prod_{k=1}^n \varphi_{nk}(t).$$

Since  $|f_{nk}(t)| \leq 1$  and  $|\varphi_{nk}(t)| \leq 1$ , we have

$$\begin{aligned} |f_n(t) - \varphi(t)| &= \left| \prod_{k=1}^n f_{nk}(t) - \prod_{k=1}^n \varphi_{nk}(t) \right| \\ &\leq \sum_{k=1}^n |f_{nk}(t) - \varphi_{nk}(t)| = \sum_{k=1}^n \left| \int_{-\infty}^{\infty} e^{itx} d(F_{nk} - \Phi_{nk}) \right| \\ &= \sum_{k=1}^n \left| \int_{-\infty}^{\infty} \left( e^{itx} - itx + \frac{t^2 x^2}{2} \right) d(F_{nk} - \Phi_{nk}) \right|, \end{aligned} \quad (4)$$

where we have used the fact that

$$\int_{-\infty}^{\infty} x^k dF_{nk} = \int_{-\infty}^{\infty} x^k d\Phi_{nk} \quad \text{for } k = 1, 2.$$

If we apply the formula for integration by parts (Theorem 11, §6, Chapter II) to the integral

$$\int_a^b \left( e^{itx} - itx + \frac{t^2 x^2}{2} \right) d(F_{nk} - \Phi_{nk}),$$

we obtain (taking account of the limits  $x^2[1 - F_{nk}(x) + F_{nk}(-x)] \rightarrow 0$ , and  $x^2[1 - \Phi_{nk}(x) + \Phi_{nk}(-x)] \rightarrow 0$ ,  $x \rightarrow \infty$ )

$$\begin{aligned} &\int_{-\infty}^{\infty} \left( e^{itx} - itx + \frac{t^2 x^2}{2} \right) d(F_{nk} - \Phi_{nk}) \\ &= it \int_{-\infty}^{\infty} (e^{itx} - 1 - itx)(F_{nk}(x) - \Phi_{nk}(x)) dx. \end{aligned} \quad (5)$$

From (4) and (5), we obtain

$$\begin{aligned} |f_n(t) - \varphi(t)| &\leq \sum_{k=1}^n \left| t \int_{-\infty}^{\infty} (e^{itx} - 1 - itx)(F_{nk}(x) - \Phi_{nk}(x)) dx \right| \\ &\leq \frac{|t|^3}{2} \varepsilon \sum_{k=1}^n \int_{|x| \geq \varepsilon} |x| |F_{nk}(x) - \Phi_{nk}(x)| dx \\ &\quad + 2t^2 \sum_{k=1}^n \int_{|x| > \varepsilon} |x| |F_{nk}(x) - \Phi_{nk}(x)| dx \\ &\leq \varepsilon |t|^3 \sum_{k=1}^n \sigma_{nk}^2 + 2t^2 \sum_{k=1}^n \int_{|x| > \varepsilon} |x| |F_{nk}(x) - \Phi_{nk}(x)| dx, \end{aligned} \quad (6)$$

where we have used the inequality



$$\int_{|x| \leq \varepsilon} |x| |F_{nk}(x) - \Phi_{nk}(x)| dx \leq 2\sigma_{nk}^2, \quad (7)$$

which is easily established by using (71), §6, Chapter II.

It follows from (6) that  $f_n(t) \rightarrow \varphi(t)$  as  $n \rightarrow \infty$ , because  $\varepsilon$  is an arbitrary positive number and condition (A) is satisfied.

This completes the proof of the theorem.

**PROOF OF THEOREM 2. 1.** By §4 of the Lindeberg condition (L), it follows that  $\max_{1 \leq k \leq n} \sigma_{nk}^2 \rightarrow 0$ . Hence, if we use the fact that  $\sum_{k=1}^n \sigma_{nk}^2 = 1$ , we obtain

$$\sum_{k=1}^n \int_{|x| > \varepsilon} x^2 d\Phi_{nk}(x) \leq \int_{|x| > \varepsilon / \sqrt{\max_{1 \leq k \leq n} \sigma_{nk}^2}} x^2 d\Phi(x) \rightarrow 0, \quad n \rightarrow \infty. \quad (8)$$

Together with Condition (L), this shows that, for every  $\varepsilon > 0$ ,

$$\sum_{k=1}^n \int_{|x| > \varepsilon} x^2 d[F_{nk}(x) + \Phi_{nk}(x)] \rightarrow 0, \quad n \rightarrow \infty. \quad (9)$$

Let us fix  $\varepsilon > 0$ . Then there is a continuous differentiable even function  $h = h(x)$  for which  $|h(x)| \leq x^2$ ,  $|h'(x)| \leq 4x$ , and

$$h(x) = \begin{cases} x^2 & |x| > 2\varepsilon, \\ 0, & |x| \leq \varepsilon. \end{cases}$$

For  $h(x)$ , we have by (9)

$$\sum_{k=1}^n \int_{|x| > \varepsilon} h(x) d[F_{nk}(x) + \Phi_{nk}(x)] \rightarrow 0, \quad n \rightarrow \infty. \quad (10)$$

By integrating by parts in (10), we obtain

$$\begin{aligned} & \sum_{k=1}^n \int_{x \geq \varepsilon} h'(x) [(1 - F_{nk}(x)) + (1 - \Phi_{nk}(x))] dx \\ &= \sum_{k=1}^n \int_{x \geq \varepsilon} h(x) d[F_{nk} + \Phi_{nk}] \rightarrow 0, \\ & \sum_{k=1}^n \int_{x \leq -\varepsilon} h'(x) [F_{nk}(x) + \Phi_{nk}(x)] dx = \sum_{k=1}^n \int_{x \leq -\varepsilon} h(x) d[F_{nk} + \Phi_{nk}] \rightarrow 0. \end{aligned}$$

Since  $h'(x) = 2x$  for  $|x| \geq 2\varepsilon$ , we obtain

$$\sum_{k=1}^n \int_{|x| \geq 2\varepsilon} |x| |F_{nk}(x) - \Phi_{nk}(x)| dx \rightarrow 0, \quad n \rightarrow \infty.$$

Therefore, since  $\varepsilon$  is an arbitrary positive number, we find that (L)  $\Rightarrow$  (A).

2. For the function  $h = h(x)$  introduced above, we find by (8) and the condition  $\max_{1 \leq k \leq n} \sigma_{nk}^2 \rightarrow 0$  that

$$\sum_{k=1}^n \int_{|x| > \varepsilon} h(x) d\Phi_{nk}(x) \leq \sum_{k=1}^n \int_{|x| > \varepsilon} x^2 d\Phi_{nk}(x) \rightarrow 0, \quad n \rightarrow \infty. \quad (11)$$

If we integrate by parts, we obtain

$$\begin{aligned}
 \left| \sum_{k=1}^n \int_{|x| \geq \varepsilon} h(x) d[F_{nk} - \Phi_{nk}] \right| &\leq \left| \sum_{k=1}^n \int_{x \geq \varepsilon} h(x) d[(1 - F_{nk}) - (1 - \Phi_{nk})] \right| \\
 &\quad + \left| \sum_{k=1}^n \int_{x \leq -\varepsilon} h(x) d[F_{nk} - \Phi_{nk}] \right| \\
 &\leq \sum_{k=1}^n \int_{x \geq \varepsilon} |h'(x)| [(1 - F_{nk}) - (1 - \Phi_{nk})] dx \\
 &\quad + \sum_{k=1}^n \int_{x \leq -\varepsilon} |h'(x)| [F_{nk} - \Phi_{nk}] dx \\
 &\leq 4 \sum_{k=1}^n \int_{|x| \geq \varepsilon} |x| |F_{nk}(x) - \Phi_{nk}(x)| dx. \quad (12)
 \end{aligned}$$

It follows from (11) and (12) that

$$\sum_{k=1}^n \int_{|x| \geq 2\varepsilon} x^2 dF_{nk}(x) \leq \sum_{k=1}^n \int_{|x| \geq \varepsilon} h(x) dF_{nk}(x) \rightarrow 0, \quad n \rightarrow \infty,$$

i.e., the Lindeberg condition (L) is satisfied.

This completes the proof of the theorem.

### 3. PROBLEMS

1. Establish formula (5).
2. Verify relations (10) and (12).

## §6. Infinitely Divisible and Stable Distributions

1. In stating Poisson's theorem in §3 we found it necessary to use a triangular array, supposing that for each  $n \geq 1$  there was a sequence of independent random variables  $\{\xi_{n,k}\}$ ,  $1 \leq k \leq n$ .

Put

$$T_n = \xi_{n,1} + \cdots + \xi_{n,n}, \quad n \geq 1. \quad (1)$$

The idea of an infinitely divisible distribution arises in the following problem: how can we determine all the distributions that can be expressed as limits of sequences of distributions of random variables  $T_n$ ,  $n \geq 1$ ?

Generally speaking, the problem of limit distributions is indeterminate in such great generality. Indeed, if  $\xi$  is a random variable and  $\xi_{n,1} = \xi$ ,  $\xi_{n,k} = 0$ ,  $1 < k \leq n$ , then  $T_n \equiv \xi$  and consequently the limit distribution is the distribution of  $\xi$ , which can be arbitrary.

In order to have a more meaningful problem, we shall suppose in the

present section that the variables  $\xi_{n,1}, \dots, \xi_{n,n}$  are, for each  $n \geq 1$ , not only independent, but also identically distributed.

Recall that this was the situation in Poisson's theorem (Theorem 4 of §3). The same framework also includes the central limit theorem (Theorem 3 of §3) for sums  $S_k = \xi_1 + \dots + \xi_n$ ,  $n \geq 1$ , of independent identically distributed random variables  $\xi_1, \xi_2, \dots$ . In fact, if we put

$$\xi_{n,k} = \frac{\xi_k - E\xi_k}{V_n}, \quad V_n^2 = VS_n,$$

then

$$T_n = \sum_{k=1}^n \xi_{n,k} = \frac{S_n - ES_n}{V_n}.$$

Consequently both the normal and the Poisson distributions can be presented as limits in a triangular array. If  $T_n \rightarrow T$ , it is intuitively clear that since  $T_n$  is a sum of independent identically distributed random variables, the limit variable  $T$  must also be a sum of independent identically distributed random variables. With this in mind, we introduce the following definition.

**Definition 1.** A random variable  $T$ , its distribution  $F_T$ , and its characteristic function  $\varphi_T$  are said to be infinitely divisible if, for each  $n \geq 1$ , there are independent identically distributed random variables  $\eta_1, \dots, \eta_n$  such that†  $T \stackrel{d}{=} \eta_1 + \dots + \eta_n$  (or, equivalently,  $F_T = F_{\eta_1} * \dots * F_{\eta_n}$ , or  $\varphi_T = (\varphi_{\eta_1})^n$ ).

**Theorem 1.** A random variable  $T$  can be a limit of sums  $T_n = \sum_{k=1}^n \xi_{n,k}$  if and only if  $T$  is infinitely divisible.

**PROOF.** If  $T$  is infinitely divisible, for each  $n \geq 1$  there are independent identically distributed random variables  $\xi_{n,1}, \dots, \xi_{n,n}$  such that  $T \stackrel{d}{=} \xi_{n,1} + \dots + \xi_{n,n}$ , and this means that  $T \stackrel{d}{=} T_n$ ,  $n \geq 1$ .

Conversely, let  $T_n \xrightarrow{d} T$ . Let us show that  $T$  is infinitely divisible, i.e., for each  $k$  there are independent identically distributed random variables  $\eta_1, \dots, \eta_k$  such that  $T \stackrel{d}{=} \eta_1 + \dots + \eta_k$ .

Choose a  $k \geq 1$  and represent  $T_{nk}$  in the form  $\zeta_n^{(1)} + \dots + \zeta_n^{(k)}$ , where

$$\zeta_n^{(1)} = \xi_{nk,1} + \dots + \xi_{nk,n}, \dots, \zeta_n^{(k)} = \xi_{nk,n(k-1)+1} + \dots + \xi_{nk,nk}.$$

Since  $T_{nk} \xrightarrow{d} T$ ,  $n \rightarrow \infty$ , the sequence of distribution functions corresponding to the random variables  $T_{nk}$ ,  $n \geq 1$ , is relatively compact and therefore, by Prohorov's theorem, is tight. Moreover,

$$[P(\zeta_n^{(1)} > z)]^k = P(\zeta_n^{(1)} > z, \dots, \zeta_n^{(k)} > z) \leq P(T_{nk} > kz)$$

and

$$[P(\zeta_n^{(1)} < -z)]^k = P(\zeta_n^{(1)} < -z, \dots, \zeta_n^{(k)} < -z) \leq P(T_{nk} < -kz).$$

† The notation  $\xi \stackrel{d}{=} \eta$  means that the random variables  $\xi$  and  $\eta$  agree in distribution, i.e.,  $F_\xi(x) = F_\eta(x)$ ,  $x \in R$ .

The family of distributions for  $\zeta_n^{(1)}$ ,  $n \geq 1$ , is tight because of the preceding two inequalities and because the family of distributions for  $T_{nk}$ ,  $n \geq 1$ , is tight. Therefore there is a subsequence  $\{n_i\} \subseteq \{n\}$  and a random variable  $\eta_1$  such that  $\zeta_{n_i}^{(1)} \xrightarrow{d} \eta_1$  as  $n_i \rightarrow \infty$ . Since the variables  $\zeta_n^{(1)}, \dots, \zeta_n^{(k)}$  are identically distributed, we have  $\zeta_{n_i}^{(2)} \xrightarrow{d} \eta_2, \dots, \zeta_{n_i}^{(k)} \xrightarrow{d} \eta_k$ , where  $\eta_1 \stackrel{d}{=} \eta_2 \stackrel{d}{=} \dots = \eta_k$ . Since  $\zeta_n^{(1)}, \dots, \zeta_n^{(k)}$  are independent, it follows from the corollary to Theorem 1 of §3 that  $\eta_1, \dots, \eta_k$  are independent and

$$T_{n_i k} = \zeta_{n_i}^{(1)} + \dots + \zeta_{n_i}^{(k)} \xrightarrow{d} \eta_1 + \dots + \eta_k.$$

But  $T_{n_i k} \xrightarrow{d} T$ , therefore (Problem 1)

$$T \stackrel{d}{=} \eta_1 + \dots + \eta_k.$$

This completes the proof of the theorem.

**Remark.** The conclusion of the theorem remains valid if we replace the hypothesis that  $\xi_{n,1}, \dots, \xi_{n,n}$  are identically distributed for each  $n \geq 1$  by the hypothesis that they are uniformly asymptotically infinitesimal (4.2).

2. To test whether a given random variable  $T$  is infinitely divisible, it is simplest to begin with its characteristic function  $\varphi(t)$ . If we can find characteristic functions  $\varphi_n(t)$  such that  $\varphi(t) = [\varphi_n(t)]^n$  for every  $n \geq 1$ , then  $T$  is infinitely divisible.

In the Gaussian case,

$$\varphi(t) = e^{itm} e^{-(1/2)t^2\sigma^2},$$

and if we put

$$\varphi_n(t) = e^{itm/n} e^{-(1/2)t^2\sigma^2/n},$$

we see at once that  $\varphi(t) = [\varphi_n(t)]^n$ .

In the Poisson case,

$$\varphi(t) = \exp\{\lambda(e^{it} - 1)\},$$

and if we put  $\varphi_n(t) = \exp\{(\lambda/n)(e^{it} - 1)\}$  then  $\varphi(t) = [\varphi_n(t)]^n$ .

If a random variable  $T$  has a  $\Gamma$ -distribution with density

$$f(x) = \begin{cases} \frac{x^{\alpha-1} e^{-x/\beta}}{\Gamma(\alpha)\beta^\alpha}, & x \geq 0, \\ 0, & x < 0, \end{cases}$$

it is easy to show that its characteristic function is

$$\varphi(t) = \frac{1}{(1 - i\beta t)^\alpha}.$$

Consequently  $\varphi(t) = [\varphi_n(t)]^n$  where

$$\varphi_n(t) = \frac{1}{(1 - i\beta t)^{\alpha/n}},$$

and therefore  $T$  is infinitely divisible.

We quote without proof the following result on the general form of the characteristic functions of infinitely divisible distributions.

**Theorem 2 (Lévy-Khinchin Theorem).** *A random variable  $T$  is infinitely divisible if and only if  $\varphi(t) = \exp \psi(t)$  and*

$$\psi(t) = it\beta - \frac{t^2\sigma^2}{2} + \int_{-\infty}^{\infty} \left( e^{itx} - 1 - \frac{itx}{1+x^2} \right) \frac{1+x^2}{x^2} d\lambda(x), \quad (2)$$

where  $\beta \in \mathbb{R}$ ,  $\sigma^2 \geq 0$  and  $\lambda$  is a finite measure on  $(\mathbb{R}, \mathcal{B}(\mathbb{R}))$  with  $\lambda\{0\} = 0$ .

3. Let  $\xi_1, \xi_2, \dots$  be a sequence of independent identically distributed random variables and  $S_n = \xi_1 + \dots + \xi_n$ . Suppose that there are constants  $b_n$  and  $a_n > 0$ , and a random variable  $T$ , such that

$$\frac{S_n - b_n}{a_n} \xrightarrow{d} T. \quad (3)$$

We ask for a description of the distributions (random variables  $T$ ) that can be obtained as limit distributions in (3).

If the independent identically distributed random variables  $\xi_1, \xi_2, \dots$  satisfy  $0 < \sigma^2 \equiv V\xi_1 < \infty$ , then if we put  $b_n = nE\xi_1$  and  $a_n = \sigma\sqrt{n}$ , we find by §4 that  $T$  is the normal distribution  $\mathcal{N}(0, 1)$ .

If  $f(x) = \theta/\pi(x^2 + \theta^2)$  is the Cauchy density (with parameter  $\theta > 0$ ) and  $\xi_1, \xi_2, \dots$  are independent random variables with density  $f(x)$ , the characteristic functions  $\varphi_{\xi_1}(t)$  are equal to  $e^{-\theta|t|}$  and therefore  $\varphi_{S_n/n}(t) = (e^{-\theta|t|/n})^n = e^{-\theta|t|}$ , i.e.,  $S_n/n$  also has a Cauchy distribution (with the same parameter  $\theta$ ).

Consequently there are other limit distributions besides the normal, for example the Cauchy distribution.

If we put  $\xi_{nk} = (\xi_k/a_n) - (b_n/na_n)$ ,  $1 \leq k \leq n$ , we find that

$$\frac{S_n - b_n}{a_n} = \sum_{k=1}^n \xi_{n,k} \quad (= T_n).$$

Therefore all conceivable distributions for  $T$  that can conceivably appear as limits in (3) are necessarily (in agreement with Theorem 1) infinitely divisible. However, the specific characteristics of the variable  $T_n = (S_n - b_n)/a_n$  may make it possible to obtain further information on the structure of the limit distributions that arise.

For this reason we introduce the following definition.

**Definition 2.** A random variable  $T$ , its distribution function  $F(x)$ , and its characteristic function  $\varphi(t)$  are *stable* if, for every  $n \geq 1$ , there are constants  $a_n > 0$ ,  $b_n$ , and independent random variables  $\xi_1, \dots, \xi_n$ , distributed like  $T$ , such that

$$a_n T + b_n \stackrel{d}{=} \xi_1 + \dots + \xi_n \quad (4)$$

or, equivalently,  $F[(x - b_n)/a_n] = \underbrace{F * \dots * F(x)}_{n \text{ times}}$ , or

$$[\varphi(t)]^n = [\varphi(a_n t)]e^{ib_n t}. \quad (5)$$

**Theorem 3.** *A necessary and sufficient condition for the random variable  $T$  to be a limit in distribution of random variables  $(S_n - b_n)/a_n$ ,  $a_n > 0$ , is that  $T$  is stable.*

PROOF. If  $T$  is stable, then by (4)

$$T \stackrel{d}{=} \frac{S_n - b_n}{a_n},$$

where  $S_n = \xi_1 + \dots + \xi_n$ , and consequently  $(S_n - b_n)/a_n \xrightarrow{d} T$ .

Conversely, let  $\xi_1, \xi_2, \dots$  be a sequence of independent identically distributed random variables,  $S_n = \xi_1 + \dots + \xi_n$  and  $(S_n - b_n)/a_n \rightarrow T$ ,  $a_n > 0$ . Let us show that  $T$  is a stable random variable.

If  $T$  is degenerate, it is evidently stable. Let us suppose that  $T$  is non-degenerate.

Choose  $k \geq 1$  and write

$$\begin{aligned} S_n^{(1)} &= \xi_1 + \dots + \xi_n, \dots, & S_n^{(k)} &= \xi_{(k-1)n+1} + \dots + \xi_{kn}, \\ T_n^{(1)} &= \frac{S_n^{(1)} - b_n}{a_n}, \dots, & T_n^{(k)} &= \frac{S_n^{(k)} - b_n}{a_n}. \end{aligned}$$

It is clear that all the variables  $T_n^{(1)}, \dots, T_n^{(k)}$  have the same distribution and

$$T_n^{(i)} \xrightarrow{d} T, \quad n \rightarrow \infty, \quad i = 1, \dots, k.$$

Write

$$U_n^{(k)} = T_n^{(1)} + \dots + T_n^{(k)}.$$

Then

$$U_n^{(k)} \xrightarrow{d} T^{(1)} + \dots + T^{(k)},$$

where  $T^{(1)} \stackrel{d}{=} \dots \stackrel{d}{=} T^{(k)} \stackrel{d}{=} T$ .

On the other hand,

$$\begin{aligned} U_n^{(k)} &= \frac{\xi_1 + \dots + \xi_{kn} - kb_n}{a_n} \\ &= \frac{a_{kn}}{a_n} \left( \frac{\xi_1 + \dots + \xi_{kn} - b_{kn}}{a_{kn}} \right) + \frac{b_{kn} - kb_n}{a_n} \\ &= \alpha_n^{(k)} V_{kn} + \beta_n^{(k)}, \end{aligned} \quad (6)$$

where

$$\alpha_n^{(k)} = \frac{a_{kn}}{a_n}, \quad \beta_n^{(k)} = \frac{b_{kn} - kb_n}{a_n}$$

and

$$V_{kn} = \frac{\xi_1 + \cdots + \xi_{kn} - b_{kn}}{a_{kn}}.$$

It is clear from (6) that

$$V_{kn} = \frac{U_n^{(k)} - \beta_n^{(k)}}{\alpha_n^{(k)}},$$

where  $V_{kn} \xrightarrow{d} T$ ,  $U_n^{(k)} \xrightarrow{d} T^{(1)} + \cdots + T^{(k)}$ ,  $n \rightarrow \infty$ .

It follows from the lemma established below that there are constants  $\alpha^{(k)} > 0$  and  $\beta^{(k)}$  such that  $\alpha_n^{(k)} \rightarrow \alpha^{(k)}$  and  $\beta_n^{(k)} \rightarrow \beta^{(k)}$  as  $n \rightarrow \infty$ . Therefore

$$T \xrightarrow{d} \frac{T^{(1)} + \cdots + T^{(k)} - \beta^{(k)}}{\alpha^{(k)}},$$

which shows that  $T$  is a stable random variable.

This completes the proof of the theorem.

We now state and prove the lemma that we used above.

**Lemma.** Let  $\xi_n \xrightarrow{d} \xi$  and let there be constants  $a_n > 0$  and  $b_n$  such that

$$a_n \xi_n + b_n \xrightarrow{d} \tilde{\xi},$$

where the random variables  $\xi$  and  $\tilde{\xi}$  are not degenerate. Then there are constants  $a > 0$  and  $b$  such that  $\lim a_n = a$ ,  $\lim b_n = b$ , and

$$\tilde{\xi} = a\xi + b.$$

**PROOF.** Let  $\varphi_n$ ,  $\varphi$  and  $\tilde{\varphi}$  be the characteristic functions of  $\xi_n$ ,  $\xi$  and  $\tilde{\xi}$ , respectively. Then  $\varphi_{a_n \xi_n + b_n}(t)$ , the characteristic function of  $a_n \xi_n + b_n$ , is equal to  $e^{itb_n} \varphi_n(a_n t)$  and, by Theorem 1 and Problem 3 of §3,

$$e^{itb_n} \varphi_n(a_n t) \rightarrow \tilde{\varphi}(t), \quad (7)$$

$$\varphi_n(t) \rightarrow \varphi(t) \quad (8)$$

uniformly on every finite interval of length  $t$ .

Let  $\{n_i\}$  be a subsequence of  $\{n\}$  such that  $a_{n_i} \rightarrow a$ . Let us first show that  $a < \infty$ . Suppose that  $a = \infty$ . By (7),

$$\sup_{|t| \leq c} |\varphi_n(a_n t) - \tilde{\varphi}(t)| \rightarrow 0, \quad n \rightarrow \infty$$

for every  $c > 0$ . We replace  $t$  by  $t_0/a_{n_i}$ . Then, since  $a_{n_i} \rightarrow \infty$ , we have

$$\left| \varphi_{n_i} \left( a_{n_i} \frac{t_0}{a_{n_i}} \right) - \tilde{\varphi} \left( \frac{t_0}{a_{n_i}} \right) \right| \rightarrow 0$$

and therefore

$$|\varphi_{n_i}(t_0)| \rightarrow |\tilde{\varphi}(0)| = 1.$$

But  $|\varphi_{n_i}(t_0)| \rightarrow |\varphi(t_0)|$ . Therefore  $|\varphi(t_0)| = 1$  for every  $t_0 \in R$ , and consequently, by Theorem 5, §12, Chapter II, the random variable  $\xi$  must be degenerate, which contradicts the hypotheses of the lemma.

Thus  $a < \infty$ . Now suppose that there are two subsequences  $\{n_i\}$  and  $\{n'_i\}$  such that  $a_{n_i} \rightarrow a$ ,  $a_{n'_i} \rightarrow a'$ , where  $a \neq a'$ ; suppose for definiteness that  $0 \leq a' < a$ . Then by (7) and (8),

$$|\varphi_{n_i}(a_{n_i}t)| \rightarrow |\varphi(at)|, \quad |\varphi_{n_i}(a_{n_i}t)| \rightarrow |\tilde{\varphi}(t)|$$

and

$$|\varphi_{n_i}(a_{n_i}t)| \rightarrow |\varphi(a't)|, \quad |\varphi_{n_i}(a_{n_i}t)| \rightarrow |\tilde{\varphi}(t)|.$$

Consequently

$$|\varphi(at)| = |\varphi(a't)|,$$

and therefore, for all  $t \in R$ ,

$$|\varphi(t)| = \left| \varphi\left(\frac{a'}{a}t\right) \right| = \cdots = \left| \varphi\left(\left(\frac{a'}{a}\right)^n t\right) \right| \rightarrow 1, \quad n \rightarrow \infty.$$

Therefore  $|\varphi(t)| \equiv 1$  and, by Theorem 5 of §12, Chapter II, it follows that  $\xi$  is a degenerate random variable. This contradiction shows that  $a = a'$  and therefore that there is a finite limit  $\lim a_n = a$ , with  $a \geq 0$ .

Let us now show that there is a limit  $\lim b_n = b$ , and that  $a > 0$ . Since (8) is satisfied uniformly on each finite interval, we have

$$\varphi_n(a_n t) \rightarrow \varphi(at),$$

and therefore, by (7), the limit  $\lim_{n \rightarrow \infty} e^{itb_n}$  exists for all  $t$  such that  $\varphi(at) \neq 0$ . Let  $\delta > 0$  be such that  $\varphi(at) \neq 0$  for all  $|t| < \delta$ . For such  $t$ ,  $\lim e^{itb_n}$  exists. Hence we can deduce (Problem 9) that  $\overline{\lim} |b_n| < \infty$ .

Let there be two sequences  $\{n_i\}$  and  $\{n'_i\}$  such that  $\lim b_{n_i} = b$  and  $\lim b_{n'_i} = b'$ . Then

$$e^{itb} = e^{itb'},$$

for  $|t| < \delta$ , and consequently  $b = b'$ . Thus there is a finite limit  $b = \lim b_n$  and, by (7),

$$\tilde{\varphi}(t) = e^{itb}\varphi(at),$$

which means that  $\tilde{\xi} \stackrel{d}{=} a\xi + b$ . Since  $\tilde{\xi}$  is not degenerate, we have  $a > 0$ .

This completes the proof of the lemma.

**4.** We quote without proof a theorem on the general form of the characteristic functions of stable distributions.

**Theorem 4 (Lévy-Khinchin Representation).** *A random variable  $T$  is stable if and only if its characteristic function  $\varphi(t)$  has the form  $\varphi(t) = \exp \psi(t)$ ,*



$$\psi(t) = it\beta - d|t|^\alpha \left( 1 + i\theta \frac{t}{|t|} G(t, \alpha) \right), \quad (9)$$

where  $0 < \alpha < 2$ ,  $\beta \in \mathbb{R}$ ,  $d \geq 0$ ,  $|\theta| \leq 1$ ,  $t/|t| = 0$  for  $t = 0$ , and

$$G(t, \alpha) = \begin{cases} \tan \frac{1}{2}\pi\alpha & \text{if } \alpha \neq 1, \\ (2/\pi) \log |t| & \text{if } \alpha = 1. \end{cases} \quad (10)$$

Observe that it is easy to exhibit characteristic functions of symmetric stable distributions:

$$\varphi(t) = e^{-d|t|^\alpha}, \quad (11)$$

where  $0 < \alpha \leq 2$ ,  $d \geq 0$ .

## 5. PROBLEMS

1. Show that  $\xi \stackrel{d}{=} \eta$  if  $\xi_n \xrightarrow{d} \xi$  and  $\xi_n \xrightarrow{d} \eta$ .
2. Show that if  $\varphi_1$  and  $\varphi_2$  are infinitely divisible characteristic functions, so is  $\varphi_1 \cdot \varphi_2$ .
3. Let  $\varphi_n$  be infinitely divisible characteristic functions and let  $\varphi_n(t) \rightarrow \varphi(t)$  for every  $t \in \mathbb{R}$ , where  $\varphi(t)$  is a characteristic function. Show that  $\varphi(t)$  is infinitely divisible.
4. Show that the characteristic function of an infinitely divisible distribution cannot take the value 0.
5. Give an example of a random variable that is infinitely divisible but not stable.
6. Show that a stable random variable  $\xi$  always satisfies the inequality  $E|\xi|^r < \infty$  for all  $r \in (0, \alpha)$ .
7. Show that if  $\xi$  is a stable random variable with parameter  $0 < \alpha \leq 1$ , then  $\varphi(t)$  is not differentiable at  $t = 0$ .
8. Prove that  $e^{-d|t|^\alpha}$  is a characteristic function provided that  $d \geq 0$ ,  $0 < \alpha \leq 2$ .
9. Let  $(b_n)_{n \geq 1}$  be a sequence of numbers such that  $\lim_n e^{itb_n}$  exists for all  $|t| < \delta$ ,  $\delta > 0$ . Show that  $\lim |b_n| < \infty$ .

## §7. Metrizable of Weak Convergence

1. Let  $(E, \mathcal{E}, \rho)$  be a metric space and  $\mathcal{P}(E) = \{P\}$ , a family of probability measures on  $(E, \mathcal{E})$ . It is natural to raise the question of whether it is possible to "metrize" the weak convergence  $P_n \xrightarrow{w} P$  that was introduced in §1, that is, whether it is possible to introduce a distance  $\mu(P, \tilde{P})$  between any two measures  $P$  and  $\tilde{P}$  in  $\mathcal{P}(E)$  in such a way that the limit  $\mu(P_n, P) \rightarrow 0$  is equivalent to the limit  $P_n \xrightarrow{w} P$ .

In connection with this formulation of the problem, it is useful to recall that convergence of random variables in probability,  $\xi_n \xrightarrow{P} \xi$ , can be metrized

by using, for example, the distance  $d_P(\xi, \eta) = \inf\{\varepsilon > 0: P(|\xi - \eta| \geq \varepsilon) \leq \varepsilon\}$  or the distances  $d(\xi, \eta) = E(|\xi - \eta|/(1 + |\xi - \eta|))$ ,  $d(\xi, \eta) = E \min(1, |\xi - \eta|)$ . (More generally, we can set  $d(\xi, \eta) = E g(|\xi - \eta|)$ , where the function  $g = g(x)$ ,  $x \geq 0$ , can be chosen as any nonnegative increasing Borel function that is continuous at zero and has the properties  $g(x + y) \leq g(x) + g(y)$  for  $x \geq 0$ ,  $y \geq 0$ ,  $g(0) = 0$ , and  $g(x) > 0$  for  $x > 0$ .) However, at the same time there is, in the space of random variables over  $(\Omega, \mathcal{F}, P)$ , no distance  $d(\xi, \eta)$  such that  $d(\xi_n, \xi) \rightarrow 0$  if and only if  $\xi_n$  converges to  $\xi$  with probability one. (In this connection, it is easy to find a sequence of random variables  $\xi_n$ ,  $n \geq 1$ , that converges to  $\xi$  in probability but does not converge with probability one.) In other words, *convergence with probability one is not metrizable*. (See the statements of problems 11 and 12 in §10, Chapter II.)

The aim of this section is to obtain concrete instances of two metrics,  $L(P, \tilde{P})$  and  $\|P - \tilde{P}\|_{BL}^*$  in the space  $\mathcal{P}(E)$  of measures, that metrize weak convergence:

$$P_n \xrightarrow{w} P \Leftrightarrow L(P_n, P) \rightarrow 0 \Leftrightarrow \|P_n - P\|_{BL}^* \rightarrow 0.$$

## 2. The Lévy-Prokhorov metric $L(P, \tilde{P})$ . Let

$$\rho(x, A) = \inf\{\rho(x, y): y \in A\},$$

$$A^\varepsilon = \{x \in E: \rho(x, A) < \varepsilon\}, \quad A \in \mathcal{E}.$$

For any two measures  $P$  and  $\tilde{P} \in \mathcal{P}(E)$ , we set

$$\sigma(P, \tilde{P}) = \inf\{\varepsilon > 0: P(F) \leq \tilde{P}(F^\varepsilon) + \varepsilon \text{ for all closed sets } F \in \mathcal{E}\} \quad (2)$$

and

$$L(P, \tilde{P}) = \max[\sigma(P, \tilde{P}), \sigma(\tilde{P}, P)]. \quad (3)$$

The following lemma shows that the function  $L(P, \tilde{P}) \in \mathcal{P}(E)$ , which is defined in this way, and is called the *Lévy-Prokhorov metric*, actually defines a metric.

**Lemma 1.** *The function  $L(P, \tilde{P})$  has the following properties:*

- (a)  $L(P, \tilde{P}) = L(\tilde{P}, P)$  ( $= \sigma(P, \tilde{P}) = \sigma(\tilde{P}, P)$ ),
- (b)  $L(P, \tilde{P}) \leq L(P, \hat{P}) + L(\hat{P}, \tilde{P})$ ,
- (c)  $L(P, \tilde{P}) = 0$  if and only if  $\tilde{P} = P$ .

**PROOF.** a) It is sufficient to show that (with  $\alpha > 0$  and  $\beta > 0$ )

$$P(F) \leq \tilde{P}(F^\alpha) + \beta \quad \text{for all closed sets } F \in \mathcal{E} \quad (4)$$

if and only if

$$\tilde{P}(F) \leq P(F^\alpha) + \beta \quad \text{for all closed sets } F \in \mathcal{E}. \quad (5)$$

Let  $T$  be a closed subset of  $\mathcal{E}$ . Then the set  $T^\alpha$  is open and it is easy to verify that  $T \subseteq E \setminus (E \setminus T^\alpha)^\alpha$ . If (4) is satisfied, then, in particular,

$$P(E \setminus T^\alpha) \leq \tilde{P}((E \setminus T^\alpha)^c) + \beta$$

and therefore,

$$\tilde{P}(T) \leq \tilde{P}(E \setminus (E \setminus T^\alpha)^c) \leq P(T^\alpha) + \beta,$$

which establishes the equivalence of (4) and (5). Hence, it follows that

$$\sigma(P, \tilde{P}) = \sigma(\tilde{P}, P) \quad (6)$$

and therefore,

$$L(P, \tilde{P}) = \sigma(P, \tilde{P}) = \sigma(\tilde{P}, P) = L(\tilde{P}, P). \quad (7)$$

b) Let  $L(P, \hat{P}) < \delta_1$  and  $L(\hat{P}, \tilde{P}) < \delta_2$ . Then for each closed set  $F \in \mathcal{E}$

$$\tilde{P}(F) \leq \hat{P}(F^{\delta_2}) + \delta_2 \leq P((F^{\delta_2})^{\delta_1}) + \delta_1 + \delta_2 \leq P(F^{\delta_1 + \delta_2}) + \delta_1 + \delta_2$$

and therefore,  $L(P, \tilde{P}) \leq \delta_1 + \delta_2$ . Hence, it follows that

$$L(P, \tilde{P}) \leq L(P, \hat{P}) + L(\hat{P}, \tilde{P}).$$

c) If  $L(P, \tilde{P}) = 0$ , then for every closed set  $F \in \mathcal{E}$  and every  $\alpha > 0$

$$P(F) \leq \tilde{P}(F^\alpha) + \alpha. \quad (8)$$

Since  $F^\alpha \downarrow F$ ,  $\alpha \downarrow 0$ , we find, by taking the limit in (8) as  $\alpha \downarrow 0$ , that  $P(F) \leq \tilde{P}(F)$  and by symmetry  $\tilde{P}(F) \leq P(F)$ . Hence,  $P(F) = \tilde{P}(F)$  for all closed sets  $F \in \mathcal{E}$ . For each Borel set  $A \in \mathcal{E}$  and every  $\varepsilon > 0$ , there is an open set  $G_\varepsilon \supseteq A$  and a closed set  $F_\varepsilon \subseteq A$  such that  $P(G_\varepsilon \setminus F_\varepsilon) \leq \varepsilon$ . Hence, it follows that every probability measure  $P$  on a metric space  $(E, \mathcal{E}, \rho)$  is completely determined by its values on closed sets. Consequently, it follows from the condition  $\tilde{P}(F) = P(F)$  for all closed sets  $F \in \mathcal{E}$  that  $\tilde{P}(A) = P(A)$  for all Borel sets  $A \in \mathcal{E}$ .

**Theorem 1.** *The Lévy-Prokhorov metric  $L(P, \tilde{P})$  metrizes weak convergence:*

$$L(P_n, P) \rightarrow 0 \Leftrightarrow P_n \xrightarrow{w} P. \quad (9)$$

PROOF. ( $\Rightarrow$ ) Let  $L(P_n, P) \rightarrow 0$ ,  $n \rightarrow \infty$ . Then for every specified closed set  $F \in \mathcal{E}$  and every  $\varepsilon > 0$ , we have, by (2) and equation a) of Lemma 1,

$$\overline{\lim}_n P_n(F) \leq P(F^\varepsilon) + \varepsilon. \quad (10)$$

If we then let  $\varepsilon \downarrow 0$ , we find that

$$\overline{\lim}_n P_n(F) \leq P(F).$$

According to Theorem 1 of §1, it follows that

$$P_n \xrightarrow{w} P. \quad (11)$$

The proof of the implication ( $\Leftarrow$ ) will be based on a series of deep and powerful facts that illuminate the content of the concept of weak convergence and the method of establishing it, as well as methods of studying rates of convergence.

Thus, let  $P_n \xrightarrow{w} P$ . This means that for every bounded continuous function  $f = f(x)$ .

$$\int_E f(x) P_n(dx) \rightarrow \int_E f(x) P(dx). \quad (12)$$

Now suppose that  $\mathcal{G}$  is a class of equicontinuous functions  $g = g(x)$  (for every  $\varepsilon > 0$  there is a  $\delta > 0$  such that  $|g(y) - g(x)| < \varepsilon$  if  $\rho(x, y) < \delta$  for all  $g \in \mathcal{G}$ ) and  $|g(x)| \leq C$  for the same constant  $C > 0$  (for all  $x \in E$  and  $g \in \mathcal{G}$ ). By Theorem 3, §8, the following condition, stronger than (12), is valid for  $\mathcal{G}$ :

$$P_n \xrightarrow{w} P \Rightarrow \sup_{g \in \mathcal{G}} \left| \int_E g(x) P_n(dx) - \int_E g(x) P(dx) \right| \rightarrow 0. \quad (13)$$

For each  $A \in \mathcal{E}$  and  $\varepsilon > 0$ , we set (as in Theorem 1, §1)

$$f_A^\varepsilon(x) = \left[ 1 - \frac{\rho(x, A)}{\varepsilon} \right]^+. \quad (14)$$

It is clear that

$$I_A(x) \leq f_A^\varepsilon(x) \leq I_{A^c}(x) \quad (15)$$

and

$$|f_A^\varepsilon(x) - f_A^\varepsilon(y)| \leq \varepsilon^{-1} |\rho(x, A) - \rho(y, A)| \leq \varepsilon^{-1} \rho(x, y).$$

Therefore, we have (13) for the class  $\mathcal{G} = \{f_A^\varepsilon(x), A \in \mathcal{E}\}$ , i.e.,

$$\Delta_n \equiv \sup_{A \in \mathcal{E}} \left| \int_E f_A^\varepsilon(x) P_n(dx) - \int_E f_A^\varepsilon(x) P(dx) \right| \rightarrow 0, \quad n \rightarrow \infty. \quad (16)$$

From this and (15) we conclude that, for every closed set  $A \in \mathcal{E}$  and  $\varepsilon > 0$ ,

$$P(A^c) \geq \int_E f_A^\varepsilon(x) dP \geq \int_E f_A^\varepsilon(x) dP_n - \Delta_n \geq P_n(A) - \Delta_n. \quad (17)$$

We choose  $n(\varepsilon)$  so that  $\Delta_n \leq \varepsilon$  for  $n \geq n(\varepsilon)$ . Then, by (17), for  $n \geq n(\varepsilon)$

$$P(A^c) \geq P_n(A) - \varepsilon. \quad (18)$$

Hence, it follows from definitions (2) and (3) that  $L(P_n, P) \leq \varepsilon$  as soon as  $n \geq n(\varepsilon)$ . Consequently,

$$P_n \xrightarrow{w} P \Rightarrow \Delta_n \rightarrow 0 \Rightarrow L(P_n, P) \rightarrow 0.$$

The theorem is now proved (up to (13)).

**3. The metric  $\|P - \tilde{P}\|_{BL}^*$ .** We denote by  $BL$  the set of bounded continuous functions  $f = f(x)$ ,  $x \in E$  (with  $\|f\|_\infty = \sup_x |f(x)| < \infty$ ) that also satisfy the Lipschitz condition

$$\|f\|_L = \sup_{x \neq y} \frac{|f(x) - f(y)|}{\rho(x, y)} < \infty.$$

We set  $\|f\|_{BL} = \|f\|_{\infty} + \|f\|_L$ . The space  $BL$  with the norm  $\|\cdot\|_{BL}$  is a Banach space.

We define the metric  $\|P - \tilde{P}\|_{BL}^*$  by setting

$$\|P - \tilde{P}\|_{BL}^* = \sup_{f \in BL} \left\{ \left| \int f d(P - \tilde{P}) \right| : \|f\|_{BL} \leq 1 \right\}. \quad (19)$$

(We can verify that  $\|P - \tilde{P}\|_{BL}^*$  actually satisfies the conditions for a metric; Problem 2.)

**Theorem 2.** *The metric  $\|P - \tilde{P}\|_{BL}^*$  metrizes weak convergence:*

$$\|P_n - P\|_{BL}^* \rightarrow 0 \Leftrightarrow P_n \xrightarrow{w} P.$$

**PROOF.** The implication  $(\Leftarrow)$  follows directly from (13). To prove  $(\Rightarrow)$ , it is enough to show that in the definition of weak convergence  $P_n \xrightarrow{w} P$  as given by (12) for every continuous bounded function  $f = f(x)$ , it is enough to restrict consideration to the class of bounded functions that satisfy a Lipschitz condition. In other words, the implication  $(\Rightarrow)$  will be proved if we establish the following result.

**Lemma 2.** *Weak convergence  $P_n \xrightarrow{w} P$  occurs if and only if property (12) is satisfied for every function  $f = f(x)$  of class  $BL$ .*

**PROOF.** The proof is obvious in one direction. Let us now consider the functions  $f_A^{\varepsilon} = f_A^{\varepsilon}(x)$  defined in (14). As was established above in the proof of Theorem 1, for each  $\varepsilon > 0$  the class  $\mathcal{G}^{\varepsilon} = \{f_A^{\varepsilon}(x), A \in \mathcal{G}\} \subseteq BL$ . If we now analyze the proof of the implication (I)  $\Rightarrow$  (II) in Theorem 1 of §1, we can observe that it actually establishes property (12) *not for all* bounded continuous functions but only for functions of class  $\mathcal{G}^{\varepsilon}$ ,  $\varepsilon > 0$ . Since  $\mathcal{G}^{\varepsilon} \subseteq BL$ ,  $\varepsilon > 0$ , it is evidently true that the satisfaction of (12) for functions of class  $BL$  implies proposition II of Theorem 1, §1, which is equivalent (by the same Theorem 1, §1) to the weak convergence  $P_n \xrightarrow{w} P$ .

**Remark.** The conclusion of Theorem 2 can be derived from Theorem 1 (the same as before) if we use the following inequalities between the metrics  $L(P, \tilde{P})$  and  $\|P - \tilde{P}\|_{BL}^*$ , which are valid for the separable metric spaces  $(E, \mathcal{E}, \rho)$ :

$$\|P - \tilde{P}\|_{BL}^* \leq 2L(P, \tilde{P}), \quad (20)$$

$$\varphi(L(P, \tilde{P})) \leq \|P - \tilde{P}\|_{BL}^*, \quad (21)$$

where  $\varphi(x) = 2x^2/(2+x)$ .

We notice that, for  $x \geq 0$ , we have  $0 \leq \varphi \leq 2/3$  if and only if  $x \leq 1$ ; and  $(2/3)x^2 \leq \varphi(x)$  for  $0 \leq x \leq 1$ ; we deduce from (20) and (21) that if  $L(P, \tilde{P}) \leq 1$  or  $\|P - \tilde{P}\|_{BL}^* \leq 2/3$ , we have

$$\frac{2}{3}L^2(P, \tilde{P}) \leq \|P - \tilde{P}\|_{BL}^* \leq 2L(P, \tilde{P}). \quad (22)$$

## 4. PROBLEMS

1. Show that in case  $E = R$  the Lévy-Prokhorov metric between the probability distributions  $P$  and  $\tilde{P}$  becomes the Lévy distance  $L(F, \tilde{F})$  between the distributions  $F$  and  $\tilde{F}$  that correspond to  $P$  and  $\tilde{P}$  (see Problem 4 in §1).
2. Show that formula (19) defines a metric on the space  $BL$ .
3. Establish the inequalities (20), (21), and (22).

## §8. On the Connection of Weak Convergence of Measures with Almost Sure Convergence of Random Elements ("Method of a Single Probability Space")

1. Let us suppose that on the probability space  $(\Omega, \mathcal{F}, P)$  there are given random elements  $X = X(\omega)$ ,  $X_n = X_n(\omega)$ ,  $n \geq 1$ , taking values in the metric space  $(E, \mathcal{E}, \rho)$ ; see §5, Chapter II. We denote by  $P$  and  $P_n$  the probability distributions of  $X$  and  $X_n$ , i.e., let

$$P(A) = P\{\omega: X(\omega) \in A\}, \quad P_n(A) = P\{\omega: X_n(\omega) \in A\}, \quad A \in \mathcal{E}.$$

Generalizing the concept of convergence in distribution of random variables (see §10, chapter II), we introduce the following definition.

**Definition 1.** A sequence of random elements  $X_n$ ,  $n \geq 1$ , is said to converge in distribution, or in law (notation:  $X_n \xrightarrow{\mathcal{D}} X$ , or  $X_n \xrightarrow{w} X$ ), if  $P_n \xrightarrow{w} P$ .

By analogy with the definitions of convergence of random variables in probability or with probability one (§10, Chapter II), it is natural to introduce the following definitions.

**Definition 2.** A sequence of random elements  $X_n$ ,  $n \geq 1$ , is said to converge in probability to  $X$  if

$$P\{\omega: \rho(X_n(\omega), X(\omega)) \geq \varepsilon\} \rightarrow 0, \quad n \rightarrow \infty. \quad (1)$$

**Definition 3.** A sequence of random elements  $X_n$ ,  $n \geq 1$ , is said to converge to  $X$  with probability one (almost surely, almost everywhere) if  $\rho(X_n(\omega), X(\omega)) \xrightarrow{a.s.} 0$ ,  $n \rightarrow \infty$ .

**Remark 1.** Both of the preceding definitions make sense, of course, provided that  $\rho(X_n(\omega), X(\omega))$  are, as functions of  $\omega \in \Omega$ , random variables, i.e.,  $\mathcal{F}$ -measurable functions. This will certainly be the case if the space  $(E, \mathcal{E}, \rho)$  is separable (Problem 1).

**Remark 2.** In connection with Definition 2, we note that our convergence in probability is metrized by the following metric that connects random elements  $X$  and  $Y$  (defined on  $(\Omega, \mathcal{F}, P)$  with values in  $E$ ):

$$d_P(X, Y) = \inf\{\varepsilon > 0 : P\{\rho(X(\omega), Y(\omega)) \geq \varepsilon\} \leq \varepsilon\}. \quad (2)$$

**Remark 3.** If the definitions of convergence in probability and with probability one are defined for random elements on the same probability space, the definition  $X_n \xrightarrow{\mathcal{D}} X$  of convergence in distribution is connected only with the convergence of distributions, and consequently, we may suppose that  $X(\omega)$ ,  $X_1(\omega)$ ,  $X_2(\omega)$ , ... have values in the same space  $E$ , but may be defined on "their own" probability spaces  $(\Omega, \mathcal{F}, P)$ ,  $(\Omega_1, \mathcal{F}_1, P_1)$ ,  $(\Omega_2, \mathcal{F}_2, P_2)$ , ... However, without loss of generality we may always suppose that they are defined on the same probability space, taken as the direct product of the preceding spaces and with the definitions  $X(\omega, \omega_1, \omega_2, \dots) = X(\omega)$ ,  $X_1(\omega, \omega_1, \omega_2, \dots) = X_1(\omega_1)$ , ...

2. By Definition 1 and the theorem on change of variables under the Lebesgue integral sign (Theorem 7, §6, Chapter II)

$$X_n \xrightarrow{\mathcal{D}} X \Leftrightarrow E f(X_n) \rightarrow E f(X) \quad (3)$$

for every bounded continuous function  $f = f(x)$ ,  $x \in E$ .

From (3) it is clear that, by Lebesgue's theorem on dominated convergence (Theorem 3, §6, Chapter II), the limit  $X_n \xrightarrow{a.s.} X$  immediately implies the limit  $X_n \xrightarrow{\mathcal{D}} X$ , which is hardly surprising if we think of the situation when  $X$  and  $X_n$  are random variables (Theorem 2, §10, Chapter II). More unexpectedly, in a certain sense there is a converse result, the precise formulation and application we now turn to.

Preliminarily, we introduce a definition.

**Definition 4.** Random elements  $X = X(\omega')$  and  $Y = Y(\omega'')$ , defined on probability spaces  $(\Omega', \mathcal{F}', P')$  and  $(\Omega'', \mathcal{F}'', P'')$  and with values in the same space  $E$ , are said to be *equivalent in distribution* (notation:  $X \stackrel{\mathcal{D}}{=} Y$ ), if they have congruent probability distributions.

**Theorem 1.** Let  $(E, \mathcal{E}, \rho)$  be a separable metric space.

1. Let random elements  $X$ ,  $X_n$ ,  $n \geq 1$ , defined on a probability space  $(\Omega, \mathcal{F}, P)$ , and with values in  $E$ , have the property that  $X_n \xrightarrow{\mathcal{D}} X$ . Then we can find a probability space  $(\Omega^*, \mathcal{F}^*, P^*)$  and random elements  $X^*$ ,  $X_n^*$ ,  $n \geq 1$ , defined on it, with values in  $E$ , such that

$$X_n^* \xrightarrow{a.s.} X^*$$

and

$$X^* \stackrel{\mathcal{D}}{=} X, \quad X_n^* \stackrel{\mathcal{D}}{=} X_n, \quad n \geq 1.$$

2. Let  $P, P_n, n \geq 1$ , be probability measures on  $(E, \mathcal{E}, \rho)$ . Then there is a probability space  $(\Omega^*, \mathcal{F}^*, P^*)$  and random elements  $X^*, X_n^*, n \geq 1$ , defined on it, with values in  $E$ , such that

$$X_n^* \xrightarrow{a.s.} X^*$$

and

$$P^* = P, \quad P_n^* = P_n, \quad n \geq 1,$$

where  $P^*$  and  $P_n^*$  are the probability distributions of  $X^*$  and  $X_n^*$ .

Before turning to the proof, we first notice that it is enough to prove only the second conclusion, since the first follows from it if we take  $P$  and  $P_n$  to be the distributions of  $X$  and  $X_n$ . Similarly, the second conclusion follows from the first. Second, we notice that a proof of the theorem in full generality is technically rather complicated. For this reason, here we give a proof only of the case  $E = R$ . This proof is rather transparent and moreover, provides a simple, clear construction of the required objectives. (Unfortunately, this construction does not work in the general case, even for  $E = R^2$ .)

PROOF OF THE THEOREM IN THE CASE  $E = R$ . Let  $F = F(x)$  and  $F_n = F_n(x)$  be distribution functions corresponding to the measures  $P$  and  $P_n$  on  $(R, \mathcal{B}(R))$ . We associate with a function  $F = F(x)$  its corresponding *quantile function*  $Q = Q(u)$ , uniquely defined by the formula

$$Q(u) = \inf\{x: F(x) \geq u\}, \quad 0 < u < 1. \quad (4)$$

It is easily verified that

$$F(x) \geq u \Leftrightarrow Q(u) \leq x. \quad (5)$$

We now take  $\Omega^* = (0, 1)$ ,  $\mathcal{F}^* = \mathcal{B}(0, 1)$ ,  $P^*$  to be Lebesgue measure, and  $P^*(dx) = dx$ . We also take  $X^*(\omega^*) = Q(\omega^*)$  and  $\omega^* \in \Omega^*$ . Then

$$P^*\{\omega^*: X^*(\omega^*) \leq x\} = P^*\{\omega^*: Q(\omega^*) \leq x\} = P^*\{\omega^*: \omega^* \leq F(x)\} = F(x),$$

i.e., the distribution of the random variable  $X^*(\omega^*) = Q(\omega^*)$  coincides exactly with  $P$ . Similarly, the distribution of  $X_n^*(\omega^*) = Q_n(\omega^*)$  coincides with  $P_n$ .

In addition, it is not difficult to show that the convergence of  $F_n(x)$  to  $F(x)$  at each point of continuity of the limit function  $F = F(x)$  (equivalent, if  $E = R$ , to the limit  $P_n \xrightarrow{w} P$ ; see Theorem 1 in §1) implies that the sequence of quantiles  $Q_n(u)$ ,  $n \geq 1$ , also converges to  $Q(u)$  at every point of continuity of the limit  $Q = Q(u)$ . Since the set of points of discontinuity of  $Q = Q(u)$ ,  $u \in (0, 1)$ , is at most countable, its Lebesgue measure  $P^*$  is zero and therefore,

$$X_n^*(\omega^*) = Q_n(\omega^*) \xrightarrow{a.s.} X^*(\omega^*) = Q(\omega^*).$$

The theorem is established in the case of  $E = R$ .



This construction in Theorem 1 of a passage from given random elements  $X$  and  $X_n$  to new elements  $X^*$  and  $X_n^*$ , defined on the same probability space, explains the announcement in the heading of this section of the *method of a single probability space*.

We now turn to a number of propositions that are established very simply by using this method.

3. Let us assume that the random elements  $X$  and  $X_n$ ,  $n \geq 1$ , are defined, for example, on a probability space  $(\Omega, \mathcal{F}, P)$  with values in a separable metric space  $(E, \mathcal{E}, \rho)$ , so that  $X_n \xrightarrow{\mathcal{D}} X$ . Also let  $h = h(x)$ ,  $x \in E$ , be a measurable mapping of  $(E, \mathcal{E}, \rho)$  into another separable metric space  $(E', \mathcal{E}', \rho')$ . In probability and mathematical statistics it is often necessary to deal with the search for conditions under which we can say of  $h = h(x)$  that the limit  $X_n \xrightarrow{\mathcal{D}} X$  implies the limit  $h(X_n) \xrightarrow{\mathcal{D}} h(X)$ .

For example, let  $\xi_1, \xi_2, \dots$  be independent identically distributed random variables with  $E\xi_1 = m$ ,  $V\xi_1 = \sigma^2 > 0$ . Let  $\bar{X}_n = (\xi_1 + \dots + \xi_n)/n$ . The central limit theorem shows that

$$\frac{\sqrt{n}(\bar{X}_n - m)}{\sigma} \xrightarrow{d} \mathcal{N}(0, 1).$$

Let us ask, for what functions  $h = h(x)$  can we guarantee that

$$h\left(\frac{\sqrt{n}(\bar{X}_n - m)}{\sigma}\right) \xrightarrow{d} h(\mathcal{N}(0, 1))?$$

(The *Mann-Wald theorem*, which is applicable to the present case, since it is satisfied for continuous functions  $h = h(x)$ , says that  $n(\bar{X} - m)^2/\sigma^2 \xrightarrow{d} \chi_1^2$ , where  $\chi_1^2$  is a random variable with a chi-squared distribution with one degree of freedom; see Table 2 in §3, Chapter I.)

A second example. If  $X = X(t, \omega)$ ,  $X_n = X_n(t, \omega)$ ,  $t \in T$ , are random processes (see §5, Chapter II) and  $h(X) = \sup_{t \in T} |X(t, \omega)|$ ,  $h(X_n) = \sup_{t \in T} |X_n(t, \omega)|$ , our problem amounts to asking under what conditions on the convergence in distribution of the processes  $X_n \xrightarrow{\mathcal{D}} X$  will there follow the convergence in distribution of their suprema,  $h(X_n) \xrightarrow{\mathcal{D}} h(X)$ .

A simple condition that guarantees the validity of the implication

$$X_n \xrightarrow{\mathcal{D}} X \Rightarrow h(X_n) \xrightarrow{\mathcal{D}} h(X),$$

is that the mapping  $h = h(x)$  is continuous. In fact, if  $f = f(x')$  is a *bounded continuous* function on  $E'$ , the function  $f(h(x))$  will also be a bounded continuous function on  $E$ . Consequently,

$$X_n \xrightarrow{\mathcal{D}} X \Rightarrow Ef(h(X_n)) \rightarrow Ef(h(X)).$$

The theorem given below shows that in fact the requirement of continuity of the function  $h = h(x)$  can be somewhat weakened by using the limit properties of the random element  $X$ .

We denote by  $\Delta_h$  the set  $\{x \in E: h(x) \text{ is not } \rho\text{-continuous at } x\}$ ; i.e., let  $\Delta_h$  be the set of points of discontinuity of the function  $h = h(x)$ . We note that  $\Delta_h \in \mathcal{E}$  (problem 4).

**Theorem 2.1.** Let  $(E, \mathcal{E}, \rho)$  and  $(E', \mathcal{E}', \rho')$  be separable metric spaces, and  $X_n \xrightarrow{\mathcal{D}} X$ . Let the mapping  $h = h(x)$ ,  $x \in E$ , have the property that

$$P\{\omega: X(\omega) \in \Delta_h\} = 0. \quad (6)$$

Then  $h(X_n) \xrightarrow{\mathcal{D}} h(X)$ .

2. Let  $P, P_n, n \geq 1$ , be probability distributions on the separable metric space  $(E, \mathcal{E}, \rho)$  and  $h = h(x)$  a measurable mapping of  $(E, \mathcal{E}, \rho)$  on a separable metric space  $(E', \mathcal{E}', \rho')$ . Let

$$P\{x: x \in \Delta_h\} = 0.$$

Then  $P_n^h \xrightarrow{w} P^h$ , where  $P_n^h(A) = P_n\{h(x) \in A\}$ ,  $P^h(A) = P\{h(x) \in A\}$ ,  $A \in \mathcal{E}'$ .

PROOF. As in Theorem 1, it is enough to prove the validity of, for example, the first proposition.

Let  $X^*$  and  $X_n^*$ ,  $n \geq 1$ , be random elements constructed by the "method of a single probability space," so that  $X^* \stackrel{\mathcal{D}}{=} X$ ,  $X_n^* \stackrel{\mathcal{D}}{=} X_n$ ,  $n \geq 1$ , and  $X_n^* \xrightarrow{a.s.} X^*$ . Let  $A^* = \{\omega^*: \rho(X_n^*, X^*) \not\rightarrow 0\}$ ,  $B^* = \{\omega^*: X^*(\omega^*) \in \Delta_n\}$ . Then  $P^*(A^* \cup B^*) = 0$ , and for  $w^* \notin A^* \cup B^*$

$$h(X_n^*(\omega^*)) \rightarrow h(X^*(\omega^*)),$$

which implies that  $h(X_n^*) \xrightarrow{a.s.} h(X^*)$ . As we noticed in subsection 1, it follows that  $h(X_n^*) \xrightarrow{\mathcal{D}} h(X^*)$ . But  $h(X_n^*) \stackrel{\mathcal{D}}{=} h(X_n)$  and  $h(X^*) \stackrel{\mathcal{D}}{=} h(X)$ . Therefore,  $h(X_n) \xrightarrow{\mathcal{D}} h(X)$ .

This completes the proof of the theorem.

4. In §6, in the proof of the implication  $(\Leftarrow)$  in Theorem 1, we used (13). We now give a proof that again relies on the "method of a single probability space."

Let  $(E, \mathcal{E}, \rho)$  be a separable metric space, and  $\mathcal{G}$  a class of equicontinuous functions  $g = g(x)$  for which also  $|g(x)| \leq C$  for all  $x \in E$  and  $g \in \mathcal{G}$ .

**Theorem 3.** Let  $P$  and  $P_n, n \geq 1$ , be probability measures on  $(E, \mathcal{E}, \rho)$  for which  $P_n \xrightarrow{w} P$ . Then

$$\sup_{g \in \mathcal{G}} \left| \int_E g(x) P_n(dx) - \int_E g(x) P(dx) \right| \rightarrow 0, \quad n \rightarrow \infty. \quad (7)$$

PROOF. Let (7) not occur. Then there are an  $a > 0$  and functions  $g_1, g_2, \dots$  from  $\mathcal{G}$  such that

$$\left| \int_E g_n(x) P_n(dx) - \int_E g_n(x) P(dx) \right| \geq a > 0 \quad (8)$$

for infinitely many values of  $n$ . Turning by the "method of a single probability space" to random elements  $X^*$  and  $X_n^*$  (see Theorem 1), we transform (8) to the form

$$|E^*g_n(X_n^*) - E^*g_n(X^*)| \geq a > 0 \quad (9)$$

for infinitely many values of  $n$ . But, by the properties of  $\mathcal{G}$ , for every  $\varepsilon > 0$  there is a  $\delta > 0$  for which  $|g(y) - g(x)| < \varepsilon$  for all  $g \in \mathcal{G}$ , if  $\rho(x, y) < \delta$ . In addition,  $|g(x)| \leq C$  for all  $x \in E$  and  $g \in \mathcal{G}$ . Therefore,

$$\begin{aligned} |E^*g_n(X_n^*) - E^*g_n(X^*)| &\leq E^*\{|g_n(X_n^*) - g_n(X^*)|; \rho(X_n^*, X^*) > \delta\} \\ &\quad + E^*\{|g_n(X_n^*) - g_n(X^*)|; \rho(X_n^*, X^*) \leq \delta\} \\ &\leq 2CP\{\rho(X_n^*, X^*) > \delta\} + \varepsilon. \end{aligned}$$

Since  $X_n^* \xrightarrow{a.s.} X^*$ , we have  $P^*\{\rho(X_n^*, X^*) > \delta\} \rightarrow 0$  as  $n \rightarrow \infty$ . Consequently, since  $\varepsilon > 0$  is arbitrary,

$$\lim_n |E^*g_n(X_n^*) - E^*g_n(X^*)| = 0,$$

which contradicts (9).

This completes the proof of the theorem.

5. In this section the idea of the "method of a single probability space" used in Theorem 1 will be applied to estimating upper bounds of the Lévy-Prokhorov metric  $L(P, \tilde{P})$  between two probability distributions on a separable space  $(E, \mathcal{E}, \rho)$ .

**Theorem 4.** *For each pair  $P, \tilde{P}$  of measures we can find a probability space  $(\Omega^*, \mathcal{F}^*, P^*)$  and random elements  $X$  and  $\tilde{X}$  on it with values in  $E$  such that their distributions coincide respectively with  $P$  and  $\tilde{P}$  and*

$$L(P, \tilde{P}) \leq d_{P^*}(X, \tilde{X}) = \inf\{\varepsilon > 0: P^*(\rho(X, \tilde{X}) \geq \varepsilon) \leq \varepsilon\}. \quad (10)$$

**PROOF.** By Theorem 1, we can actually find a probability space  $(\Omega^*, \mathcal{F}^*, P^*)$  and random elements  $X$  and  $\tilde{X}$  such that  $P^*(X \in A) = P(A)$  and  $P^*(\tilde{X} \in A) = \tilde{P}(A)$ ,  $A \in \mathcal{E}$ .

Let  $\varepsilon > 0$  have the property that

$$P^*(\rho(X, \tilde{X}) \geq \varepsilon) \leq \varepsilon. \quad (11)$$

Then for every  $A \in \mathcal{E}$

$$\begin{aligned} \tilde{P}(A) &= P^*(\tilde{X} \in A) = P^*(\tilde{X} \in A, X \in A^c) + P^*(\tilde{X} \in A, X \in A) \\ &\leq P^*(X \in A^c) + P^*(\rho(X, \tilde{X}) \geq \varepsilon) \leq P(A^c) + \varepsilon. \end{aligned}$$

Hence, by the definition of the Lévy-Prokhorov metric (subsection 2, §6)

$$L(P, \tilde{P}) \leq \varepsilon. \quad (12)$$

From (11) and (12), if we take the infimum for  $\varepsilon > 0$  we obtain the required assertion (10).

**Corollary.** Let  $X$  and  $\tilde{X}$  be random elements defined on a probability space  $(\Omega, \mathcal{F}, P)$  with values in  $E$ . Let  $P_X$  and  $P_{\tilde{X}}$  be their probability distributions. Then

$$L(P_X, P_{\tilde{X}}) \leq d_P(X, \tilde{X}).$$

**Remark 1.** The preceding proof shows that in fact (10) is valid whenever we can exhibit on any probability space  $(\Omega^*, \mathcal{F}^*, P^*)$  random elements  $X$  and  $\tilde{X}$  with values in  $E$  whose distributions coincide with  $P$  and  $\tilde{P}$  and for which the set  $\{\omega^*: \rho(X(\omega^*), \tilde{X}(\omega^*)) \geq \varepsilon\} \in \mathcal{F}^*, \varepsilon > 0$ . Hence, the property of (10) depends in an essential way on how well, with respect to the measures  $P$  and  $\tilde{P}$ , the objects  $(\Omega^*, \mathcal{F}^*, P^*)$  and  $X, \tilde{X}$  are constructed. (The procedure for constructing  $\Omega^*, \mathcal{F}^*, P^*$  and  $X, \tilde{X}$  as well as the measure  $P^*$ , is called *coupling* (joining, linking).) We could, for example, choose  $P^*$  equal to the direct product of the measures  $P$  and  $\tilde{P}$ , but this choice would, as a rule, not lead to a good estimate (10).

**Remark 2.** It is natural to raise the question of when there is equality in (10). In this connection we state the following result without proof: Let  $P$  and  $\tilde{P}$  be two probability measures on a separable metric space  $(E, \mathcal{E}, \rho)$ ; then there are  $(\Omega^*, \mathcal{F}^*, P^*)$  and  $X, \tilde{X}$ , such that

$$L(P, \tilde{P}) = d_{P^*}(X, \tilde{X}) = \inf\{\varepsilon > 0: P^*(\rho(X, \tilde{X}) \geq \varepsilon) \leq \varepsilon\}.$$

## 5. PROBLEMS

1. Prove that in the case of separable metric spaces the real function  $\rho(X(\omega), Y(\omega))$  is a random variable for all random elements  $X(\omega)$  and  $Y(\omega)$  defined on a probability space  $(\Omega, \mathcal{F}, P)$ .
2. Prove that the function  $d_P(X, Y)$  defined in (2) is a metric in the space of random elements with values in  $E$ .
3. Establish (5).
4. Prove that the set  $\Delta_h = \{x \in E: h(x) \text{ is not } \rho\text{-continuous at } x\} \in \mathcal{E}$ .

## §9. The Distance in Variation between Probability Measures. Kakutani–Hellinger Distance and Hellinger Integrals. Application to Absolute Continuity and Singularity of Measures

1. Let  $(\Omega, \mathcal{F})$  be a measurable space and  $\mathcal{P} = \{P\}$  a family of probability measures on it.

**Definition 1.** The *distance in variation* between measures  $P$  and  $\tilde{P}$  in  $\mathcal{P}$  (notation:  $\|P - \tilde{P}\|$ ) is the total (signed) variation of  $P - \tilde{P}$ , i.e.,

$$\text{Var}(P - \tilde{P}) \equiv \sup \left| \int_{\Omega} \varphi(\omega) d(P - \tilde{P}) \right|, \quad (1)$$

where the sup is over the class of all  $\mathcal{F}$ -measurable functions that satisfy the condition that  $|\varphi(\omega)| \leq 1$ .

**Lemma 1.** *The distance in variation is given by*

$$\|P - \tilde{P}\| = 2 \sup_{A \in \mathcal{F}} |P(A) - \tilde{P}(A)|. \quad (2)$$

PROOF. Since, for all  $A \in \mathcal{F}$ ,

$$P(A) - \tilde{P}(A) = \tilde{P}(\bar{A}) - P(\bar{A}),$$

we have

$$2|P(A) - \tilde{P}(A)| = |P(A) - \tilde{P}(A)| + |P(\bar{A}) - \tilde{P}(\bar{A})| \leq \|P - \tilde{P}\|,$$

where the last inequality follows from (1).

For the proof of the converse inequality we turn to the Hahn decomposition (see, for example, [K9] or [H1], p. 121) of a *signed measure*  $\mu \equiv P - \tilde{P}$ . In this decomposition the measure  $\mu$  is represented in the form  $\mu = \mu_+ - \mu_-$ , where the nonnegative measures  $\mu_+$  and  $\mu_-$  (the upper and lower variations of  $\mu$ ) are of the form

$$\mu_+(A) = \int_{A \cap M} d\mu, \quad \mu_-(A) = - \int_{A \cap \bar{M}} d\mu, \quad A \in \mathcal{F},$$

where  $M$  is a subset of  $\mathcal{F}$ . Here

$$\text{Var } \mu = \text{Var } \mu_+ + \text{Var } \mu_- = \mu_+(\Omega) + \mu_-(\Omega).$$

Since

$$\mu_+(\Omega) = P(M) - \tilde{P}(M), \quad \mu_-(\Omega) = \tilde{P}(\bar{M}) - P(\bar{M}),$$

we have

$$\|P - \tilde{P}\| = (P(M) - \tilde{P}(M)) + (\tilde{P}(\bar{M}) - P(\bar{M})) \leq 2 \sup_{A \in \mathcal{F}} |P(A) - \tilde{P}(A)|.$$

This completes the proof of the lemma.

**Definition 2.** A sequence of probability measures  $P_n$ ,  $n \geq 1$ , is said to be *convergent in variation* to the measure  $P$  if

$$\|P_n - P\| \rightarrow 0, \quad n \rightarrow \infty. \quad (3)$$

From this definition and Theorem 1, §1, Chapter III, it is easily seen that convergence in variation of probability measures defined on a metric space  $(\Omega, \mathcal{F}, \rho)$  implies their weak convergence.

*The proximity in variation of distributions* is, perhaps, the strongest form of closeness of probability distributions, since if two distributions are close in

variation, then in practice, in specific situations, they can be considered indistinguishable. In this connection, the impression may be created that the study of distance in variation is not of much probabilistic interest. However, for example, in Poisson's theorem (§6, Chapter I) the convergence of the binomial to the Poisson distribution takes place in the sense of convergence in variation to the zero distribution. (Later, in §11, we shall obtain an upper bound for this distance.)

We also provide an example from the field of mathematical statistics, where the necessity of determining the distance in variation between measures  $P$  and  $\tilde{P}$  arises in a natural way in connection with the problem of discrimination between the results of observations of two statistical hypotheses  $H$  (the true distribution is  $P$ ) and  $\tilde{H}$  (the true distribution is  $\tilde{P}$ ) in connection with the question of whether the measure  $P$  or  $\tilde{P}$ , defined on  $(\Omega, \mathcal{F})$ , is more plausible. If  $\omega \in \Omega$  is treated as the result of an observation, by a test (for different hypotheses  $H$  and  $\tilde{H}$ ) we understand any  $\mathcal{F}$ -measurable function  $\varphi = \varphi(\omega)$  with values in  $[0, 1]$ , the statistical meaning of which is that  $\varphi(\omega)$  is "the probability with which hypothesis  $\tilde{H}$  is accepted if the result of the observation is  $\omega$ ."

We shall characterize the quality of different hypotheses  $H$  and  $\tilde{H}$  by the probabilities of errors of the first and second kind:

$$\alpha(\varphi) = E\varphi(\omega) \quad (= \text{Prob (accepting } \tilde{H} | H \text{ is true)}),$$

$$\beta(\varphi) = \tilde{E}(1 - \varphi(\omega)) \quad (= \text{Prob (accepting } H | \tilde{H} \text{ is true)}).$$

In the case when hypotheses  $H$  and  $\tilde{H}$  are equivalent, the optimum is naturally to consider a test  $\varphi^* = \varphi^*(\omega)$  (if there is such a test) that minimizes the sum  $\alpha(\varphi) + \beta(\varphi)$  of the errors.

We set

$$\mathcal{E}r(P, \tilde{P}) = \inf_{\varphi} [\alpha(\varphi) + \beta(\varphi)]. \quad (4)$$

Let  $Q = (P + \tilde{P})/2$  and  $z = dP/dQ$ ,  $\tilde{z} = d\tilde{P}/dQ$ . Then

$$\begin{aligned} \mathcal{E}r(P, \tilde{P}) &= \inf_{\varphi} [E\varphi + \tilde{E}(1 - \varphi)] \\ &= \inf_{\varphi} E_Q[z\varphi + \tilde{z}(1 - \varphi)] = 1 + \inf_{\varphi} E_Q[\varphi(z - \tilde{z})], \end{aligned}$$

where  $E_Q$  is the expectation of the measure  $Q$ .

It is easy to see that the inf is attained by the function

$$\varphi^*(\omega) = I\{\tilde{z} < z\}$$

and, since  $E_Q(z - \tilde{z}) = 0$ , that

$$\mathcal{E}r(P, \tilde{P}) = 1 - \frac{1}{2}E_Q|z - \tilde{z}| = 1 - \frac{1}{2}\|P - \tilde{P}\|, \quad (5)$$

where the last equation will follow from Lemma 2, below. Therefore, it is clear from (5) that the quality of various hypotheses that characterize the

function  $\mathcal{E}r(P, \tilde{P})$  really depends on the degree of proximity of the measures  $P$  and  $\tilde{P}$  in the sense of *distance in variation*.

**Lemma 2.** Let  $Q$  be a  $\sigma$ -finite measure such that  $P \ll Q$ ,  $\tilde{P} \ll Q$  and  $z = dP/dQ$ ,  $\tilde{z} = d\tilde{P}/dQ$  are Radon–Nikodým measures of  $P$  and  $\tilde{P}$  with respect to  $Q$ . Then

$$\|P - \tilde{P}\| = E_Q|z - \tilde{z}| \quad (6)$$

and if  $Q = (P + \tilde{P})/2$ , we have

$$\|P - \tilde{P}\| = E_Q|z - \tilde{z}| = 2E_Q|1 - z| = 2E_Q|1 - \tilde{z}|. \quad (7)$$

PROOF. For all  $\mathcal{F}$ -measurable functions  $\psi = \psi(\omega)$  with  $|\psi(\omega)| \leq 1$ , we see from the definitions of  $z$  and  $\tilde{z}$  that

$$|E\psi - \tilde{E}\psi| = |E_Q\psi(z - \tilde{z})| \leq E_Q|\psi||z - \tilde{z}| \leq E_Q|z - \tilde{z}|. \quad (8)$$

Therefore,

$$\|P - \tilde{P}\| \leq E_Q|z - \tilde{z}|. \quad (9)$$

However, for the function

$$\psi = \operatorname{sgn}(\tilde{z} - z) = \begin{cases} 1, & \tilde{z} \geq z, \\ -1, & \tilde{z} < z, \end{cases}$$

we have

$$|E\psi - \tilde{E}\psi| = E_Q|z - \tilde{z}|. \quad (10)$$

We obtain the required equation (6) from (9) and (10). Then (7) follows from (6) because  $z + \tilde{z} = 2$  ( $Q$ -a.s.).

**Corollary 1.** Let  $P$  and  $\tilde{P}$  be two probability distributions on  $(R, \mathcal{B}(R))$  with probability densities (with respect to Lebesgue measure  $dx$ )  $p(x)$  and  $\tilde{p}(x)$ ,  $x \in R$ . Then

$$\|P - \tilde{P}\| = \int_{-\infty}^{\infty} |p(x) - \tilde{p}(x)| dx. \quad (11)$$

(As the measure  $Q$ , we are to take Lebesgue measure on  $(R, \mathcal{B}(R))$ .)

**Corollary 2.** Let  $P$  and  $\tilde{P}$  be two discrete measures,  $P = (p_1, p_2, \dots)$ ,  $\tilde{P} = (\tilde{p}_1, \tilde{p}_2, \dots)$ , concentrated on a countable set of points  $x_1, x_2, \dots$ . Then

$$\|P - \tilde{P}\| = \sum_{i=1}^{\infty} |p_i - \tilde{p}_i|. \quad (12)$$

(As the measure  $Q$ , we are to take the counting measure, i.e., that with  $Q(\{x_i\}) = 1$ ,  $i = 1, 2, \dots$ )

2. We now turn to still another measure of the proximity of two probability measures, from among many (as will follow later) related proximities of measures in variation.

Let  $P$  and  $\tilde{P}$  be probability measures on  $(\Omega, \mathcal{F})$  and  $Q$ , the third probability measure, dominating  $P$  and  $\tilde{P}$ , i.e., with the probabilities  $P \ll Q$  and  $\tilde{P} \ll Q$ . We again use the notation

$$z = \frac{dP}{dQ}, \quad \tilde{z} = \frac{d\tilde{P}}{dQ}.$$

**Definition 3.** The *Kakutani–Hellinger distance* between the measures  $P$  and  $\tilde{P}$  is the nonnegative number  $\rho(P, \tilde{P})$  such that

$$\rho^2(P, \tilde{P}) = \frac{1}{2} E_Q[\sqrt{z} - \sqrt{\tilde{z}}]^2. \quad (13)$$

Since

$$E_Q[\sqrt{z} - \sqrt{\tilde{z}}]^2 = \int_{\Omega} \left[ \sqrt{\frac{dP}{dQ}} - \sqrt{\frac{d\tilde{P}}{dQ}} \right]^2 dQ, \quad (14)$$

it is natural to write  $\rho^1(P, \tilde{P})$  symbolically in the form

$$\rho^2(P, \tilde{P}) = \frac{1}{2} \int_{\Omega} [\sqrt{dP} - \sqrt{d\tilde{P}}]^2. \quad (15)$$

If we set

$$H(P, \tilde{P}) = E_Q \sqrt{z\tilde{z}}, \quad (16)$$

then, by analogy with (15), we may write symbolically

$$H(P, \tilde{P}) = \int_{\Omega} \sqrt{dP d\tilde{P}}. \quad (17)$$

From (13) and (16), as well as from (15) and (17), it is clear that

$$\rho^2(P, \tilde{P}) = 1 - H(P, \tilde{P}). \quad (18)$$

The number  $H(P, \tilde{P})$  is called the *Hellinger integral* of the measures  $P$  and  $\tilde{P}$ . It turns out to be convenient, for many purposes, to consider the *Hellinger integrals*  $H(\alpha; P, \tilde{P})$  of order  $\alpha \in (0, 1)$ , defined by the formula

$$H(\alpha; P, \tilde{P}) = E_Q z^{\alpha} \tilde{z}^{1-\alpha}, \quad (19)$$

or, symbolically,

$$H(\alpha; P, \tilde{P}) = \int_{\Omega} (dP)^{\alpha} (d\tilde{P})^{1-\alpha}. \quad (20)$$

It is clear that  $H(1/2; P, \tilde{P}) = H(P, \tilde{P})$ .

For Definition 3 to be reasonable, we need to show that the number  $\rho^2(P, \tilde{P})$  is independent of the choice of the dominating measure and that in fact  $\rho(P, \tilde{P})$  satisfies the requirements of the concept of “distance.”

**Lemma 3.1.** *The Hellinger integral of order  $\alpha \in (0, 1)$  (and consequently also  $\rho(P, \tilde{P})$ ) is independent of the choice of the dominating measure  $Q$ .*



2. The function  $\rho$  defined in (13) is a metric on the set of probability measures.

PROOF. 1. If the measure  $Q'$  dominates  $P$  and  $\tilde{P}$ ,  $Q'$  also dominates  $Q = (P + \tilde{P})/2$ . Hence, it is enough to show that if  $Q \ll Q'$ , we have

$$E_Q(z^\alpha \tilde{z}^{1-\alpha}) = E_{Q'}(z')^\alpha (\tilde{z}')^{1-\alpha},$$

where  $z' = dP/dQ'$  and  $\tilde{z}' = d\tilde{P}/dQ'$ .

Let us set  $v = dQ/dQ'$ . Then  $z' = zv$ ,  $\tilde{z}' = \tilde{z}v$ , and

$$E_Q(z^\alpha \tilde{z}^{1-\alpha}) = E_{Q'}(vz^\alpha \tilde{z}^{1-\alpha}) = E_{Q'}(z')^\alpha (\tilde{z}')^{1-\alpha},$$

which establishes the first assertion.

2. If  $\rho(P, \tilde{P}) = 0$  we have  $z = \tilde{z}$  ( $Q$ -a.e.) and hence,  $P = \tilde{P}$ . By symmetry, we evidently have  $\rho(P, \tilde{P}) = \rho(\tilde{P}, P)$ . Finally, let  $P$ ,  $P'$ , and  $P''$  be three measures,  $P \ll Q$ ,  $P' \ll Q$ , and  $P'' \ll Q$ , with  $z = dP/dQ$ ,  $z' = dP'/dQ$ , and  $z'' = dP''/dQ$ . By using the validity of the triangle inequality for the norm in  $L^2(\Omega, \mathcal{F}, Q)$ , we obtain

$$[E_Q(\sqrt{z} - \sqrt{z''})^2]^{1/2} \leq [E_Q(\sqrt{z} - \sqrt{z'})^2]^{1/2} + [E_Q(\sqrt{z'} - \sqrt{z''})^2]^{1/2},$$

i.e.,

$$\rho(P, P'') \leq \rho(P, P') + \rho(P', P'').$$

This completes the proof of the lemma.

By Definition (19) and Fubini's theorem (§6, Chapter II), it follows immediately that in the case when the measures  $P$  and  $\tilde{P}$  are *direct products* of measures,  $P = P_1 \times \cdots \times P_n$ ,  $\tilde{P} = \tilde{P}_1 \times \cdots \times \tilde{P}_n$  (see subsection 9, §6, Chapter II), the Hellinger integral between the measures  $P$  and  $\tilde{P}$  is equal to the product of the corresponding Hellinger integrals:

$$H(\alpha; P, \tilde{P}) = \prod_{i=1}^n H(\alpha; P_i, \tilde{P}_i).$$

The following theorem shows the connection between distance in variation and Kakutani-Hellinger distance (or, equivalently, the Hellinger integral). In particular, it shows that these distances define the same topology in the space of probability measures on  $(\Omega, \mathcal{F})$ .

**Theorem 1.** *We have the following inequalities:*

$$2[1 - H(P, \tilde{P})] \leq \|P - \tilde{P}\| \leq \sqrt{8[1 - H(P, \tilde{P})]}, \quad (21)$$

$$\|P - \tilde{P}\| \leq 2\sqrt{1 - H^2(P, \tilde{P})}. \quad (22)$$

*In particular,*

$$2\rho^2(P, \tilde{P}) \leq \|P - \tilde{P}\| \leq \sqrt{8\rho(P, \tilde{P})}. \quad (23)$$

PROOF. Since  $H(P, \tilde{P}) \leq 1$  and  $1 - x^2 \leq 2(1 - x)$  for  $0 \leq x \leq 1$ , the right-

hand inequality in (21) follows from (22), the proof of which is provided by the following chain of inequalities (where  $Q = (1/2)(P + \tilde{P})$ ):

$$\begin{aligned}\frac{1}{2}\|P - \tilde{P}\| &= E_Q|1 - z| \leq \sqrt{E_Q|1 - z|^2} = \sqrt{1 - E_Q z(2 - z)} \\ &= \sqrt{1 - E_Q z z} = \sqrt{1 - E_Q(\sqrt{z\tilde{z}})^2} \leq \sqrt{1 - (E_Q\sqrt{z\tilde{z}})^2} \\ &= \sqrt{1 - H^2(P, \tilde{P})}.\end{aligned}$$

Finally, the first inequality in (21) follows from the fact that by the inequality

$$\frac{1}{2}[\sqrt{z} - \sqrt{2 - z}]^2 \leq |z - 1|, \quad z \in [0, 2],$$

and we have (again,  $Q = (1/2)(P + \tilde{P})$ )

$$1 - H(P, \tilde{P}) = \rho^2(P, \tilde{P}) = \frac{1}{2}E_Q[\sqrt{z} - \sqrt{2 - z}]^2 \leq \frac{1}{2}E_Q|z - 1| = \frac{1}{2}\|P - \tilde{P}\|.$$

**Remark.** It can be shown in a similar way that, for every  $\alpha \in (0, 1)$ ,

$$2[1 - H(\alpha; P, \tilde{P})] \leq \|P - \tilde{P}\| \leq \sqrt{c_\alpha(1 - H(\alpha; P, \tilde{P}))}, \quad (24)$$

where  $c_\alpha$  is a constant.

**Corollary 1.** Let  $P$  and  $P^n$ ,  $n \geq 1$ , be probability measures on  $(\Omega, \mathcal{F})$ . Then (as  $n \rightarrow \infty$ )

$$\|P^n - P\| \rightarrow 0 \Leftrightarrow H(P^n, P) \rightarrow 1 \Leftrightarrow \rho(P^n, P) \rightarrow 0,$$

$$\|P^n - P\| \rightarrow 2 \Leftrightarrow H(P^n, P) \rightarrow 0 \Leftrightarrow \rho(P^n, P) \rightarrow 1.$$

**Corollary 2.** Since by (5)

$$\mathcal{E}r(P, \tilde{P}) = 1 - \frac{1}{2}\|P - \tilde{P}\|,$$

we have, by (21) and (22),

$$\frac{1}{2}H^2(P, \tilde{P}) \leq 1 - \sqrt{1 - H^2(P, \tilde{P})} \leq \mathcal{E}r(P, \tilde{P}) \leq H(P, \tilde{P}). \quad (25)$$

In particular, let

$$P^n = \underbrace{P \times \cdots \times P}_n, \quad \tilde{P}^n = \underbrace{\tilde{P} \times \cdots \times \tilde{P}}_n$$

be direct products of measures. Then, since  $H(P^n, \tilde{P}^n) = [H(P, \tilde{P})]^n = e^{-\lambda n}$  with  $\lambda = -\ln H(P, \tilde{P}) \geq \rho^2(P, \tilde{P})$ , we obtain from (25) the inequalities

$$\frac{1}{2}e^{-2\lambda n} \leq \mathcal{E}r(P^n, \tilde{P}^n) \leq e^{-\lambda n} \leq e^{-n\rho^2(P, \tilde{P})}. \quad (26)$$

In connection with the problem, considered above, of *distinguishing two statistical hypotheses* from these inequalities, we have the following result.

Let  $\xi_1, \xi_2, \dots$  be independent identically distributed random elements, that have either the probability distribution  $P$  (Hypothesis  $H$ ) or  $\tilde{P}$  (Hypothe-

sis  $\tilde{H}$ ), with  $\tilde{P} \neq P$ , and therefore,  $\rho^2(P, \tilde{P}) > 0$ . Therefore, when  $n \rightarrow \infty$ , the function  $\mathcal{E}r(P^n, \tilde{P}^n)$ , which describes the quality of optimality of the hypotheses  $H$  and  $\tilde{H}$  as observations of  $\xi_1, \xi_2, \dots$ , decreases exponentially to zero.

4. In using Hellinger integrals of order  $\alpha$  (described above), it will be convenient to introduce the notions of absolute continuity and singularity of probability measures.

Let  $P$  and  $\tilde{P}$  be two probability measures defined on a measurable space  $(\Omega, \mathcal{F})$ . We say that  $\tilde{P}$  is *absolutely continuous* with respect to  $P$  (notation:  $\tilde{P} \ll P$ ) if  $\tilde{P}(A) = 0$  whenever  $P(A) = 0$  for  $A \in \mathcal{F}$ . If  $\tilde{P} \ll P$  and  $P \ll \tilde{P}$ , we say that  $P$  and  $\tilde{P}$  are *equivalent* ( $\tilde{P} \sim P$ ). The measures  $P$  and  $\tilde{P}$  are called *singular* or *orthogonal* ( $\tilde{P} \perp P$ ), if there is an  $A \in \mathcal{F}$  for which  $P(A) = 1$  and  $\tilde{P}(A) = 0$  (i.e.,  $P$  and  $\tilde{P}$  "sit" on different sets).

Let  $Q$  be a probability measure, with  $P \ll Q$ ,  $\tilde{P} \ll Q$ ,  $z = dP/dQ$ ,  $\tilde{z} = d\tilde{P}/dQ$ .

**Theorem 2.** *The following conditions are equivalent:*

- (a)  $\tilde{P} \ll P$ ,
- (b)  $\tilde{P}(z > 0) = 1$ ,
- (c)  $H(\alpha; P, \tilde{P}) \rightarrow 1, \alpha \downarrow 0$ .

**Theorem 3.** *The following conditions are equivalent:*

- (a)  $\tilde{P} \perp P$ ,
- (b)  $\tilde{P}(z > 0) = 0$ ,
- (c)  $H(\alpha; P, \tilde{P}) \rightarrow 0, \alpha \downarrow 0$ ,
- (d)  $H(\alpha; P, \tilde{P}) = 0$  for all  $\alpha \in (0, 1)$ ,
- (e)  $H(\alpha; P, \tilde{P}) = 0$  for some  $\alpha \in (0, 1)$ .

The proofs of these theorems will be given simultaneously. By the definitions of  $z$  and  $\tilde{z}$ ,

$$P(z = 0) = E_Q[zI(z = 0)] = 0, \quad (27)$$

$$\begin{aligned} \tilde{P}(A \cap \{z > 0\}) &= E_Q[\tilde{z}I(A \cap \{z > 0\})] \\ &= E_Q\left[\tilde{z} \frac{z}{z} I(A \cap \{z > 0\})\right] = E\left[\frac{\tilde{z}}{z} I(A \cap \{z > 0\})\right] \\ &= E\left[\frac{\tilde{z}}{z} I(A)\right]. \end{aligned} \quad (28)$$

Consequently, we have the *Lebesgue decomposition*

$$\tilde{P}(A) = E\left[\frac{\tilde{z}}{z} I(A)\right] + \tilde{P}(A \cap \{z = 0\}), \quad A \in \mathcal{F}, \quad (29)$$

in which  $Z = \tilde{z}/z$  is called the *Lebesgue derivative* of  $\tilde{P}$  with respect to  $P$  and

denoted by  $d\tilde{P}/dP$  (compare the remark on the Radon–Nikodým theorem, §6, Chapter III).

Hence, we immediately obtain the equivalence of (a) and (b) in both theorems.

Moreover, since

$$z^\alpha \tilde{z}^{1-\alpha} \rightarrow \tilde{z} I(z > 0), \quad \alpha \downarrow 0,$$

and for  $\alpha \in (0, 1)$

$$0 \leq z^\alpha \tilde{z}^{1-\alpha} \leq \alpha z + (1 - \alpha) \tilde{z} \leq z + \tilde{z}$$

with  $E_Q(z + \tilde{z}) = 2$ , we have, by Lebesgue's dominated convergence theorem,

$$\lim_{\alpha \downarrow 0} H(\alpha; P, \tilde{P}) = E_Q \tilde{z} I(z > 0) = \tilde{P}(z > 0)$$

and therefore, (b)  $\Leftrightarrow$  (c) in both theorems.

Finally, let us show that in the second theorem (c)  $\Leftrightarrow$  (d)  $\Leftrightarrow$  (e). For this, we need only note that  $H(\alpha; P, \tilde{P}) = \tilde{E}(z/\tilde{z})^\alpha I(\tilde{z} > 0)$  and  $\tilde{P}(\tilde{z} > 0) = 1$ . Hence, for each  $\alpha \in (0, 1)$  we have  $\tilde{P}(z > 0) = 0 \Leftrightarrow H(\alpha; P, \tilde{P}) = 0$ , from which there follows the implication (c)  $\Leftrightarrow$  (d)  $\Leftrightarrow$  (e).

**EXAMPLE 1.** Let  $P = P_1 \times P_2 \times \dots$ ,  $\tilde{P} = \tilde{P}_1 \times \tilde{P}_2 \times \dots$ , where  $P_k$  and  $\tilde{P}_k$  are Gaussian measures on  $(R, \mathcal{B}(R))$  with densities

$$p_k(x) = \frac{1}{\sqrt{2\pi}} e^{(x-a_k)^2/2}, \quad \tilde{p}_k(x) = \frac{1}{\sqrt{2\pi}} e^{-(x-\tilde{a}_k)^2/2}.$$

Since

$$H(\alpha; P, \tilde{P}) = \prod_{k=1}^{\infty} H(\alpha; P_k, \tilde{P}_k),$$

where a simple calculation shows that

$$H(\alpha; P_k, \tilde{P}_k) = \int_{-\infty}^{\infty} p_k^\alpha(x) \tilde{p}_k^{1-\alpha}(x) dx = e^{-\alpha(1-\alpha)/2 (a_k - \tilde{a}_k)^2},$$

we have

$$H(\alpha; P, \tilde{P}) = e^{-\alpha(1-\alpha)/2 \sum_{k=1}^{\infty} (a_k - \tilde{a}_k)^2}.$$

From Theorems 2 and 3, we find that

$$\tilde{P} \ll P \Leftrightarrow P \ll \tilde{P} \Leftrightarrow \tilde{P} \sim P \Leftrightarrow \sum_{k=1}^{\infty} (a_k - \tilde{a}_k)^2 < \infty,$$

$$\tilde{P} \perp P \Leftrightarrow \sum_{k=1}^{\infty} (a_k - \tilde{a}_k)^2 = \infty.$$

**EXAMPLE 2.** Again let  $P = P_1 \times P_2 \times \dots$ ,  $\tilde{P} = \tilde{P}_1 \times \tilde{P}_2 \times \dots$ , where  $P_k$  and  $\tilde{P}_k$  are Poisson distributions with respective parameters  $\lambda_k > 0$  and  $\tilde{\lambda}_k > 0$ . Then it is easily shown that

$$\begin{aligned}\tilde{P} \ll P &\Leftrightarrow P \ll \tilde{P} \Leftrightarrow \tilde{P} \sim P \Leftrightarrow \sum_{k=1}^{\infty} (\sqrt{\lambda_k} - \sqrt{\tilde{\lambda}_k})^2 < \infty, \\ \tilde{P} \perp P &\Leftrightarrow \sum_{k=1}^{\infty} (\sqrt{\lambda_k} - \sqrt{\tilde{\lambda}_k})^2 = \infty.\end{aligned}\tag{30}$$

## 5. PROBLEMS

1. In the notation of Lemma 2, set

$$P \wedge \tilde{P} = E_Q(z \wedge \tilde{z}),$$

where  $z \wedge \tilde{z} = \min(z, \tilde{z})$ . Show that

$$\|P - \tilde{P}\| = 2(1 - P \wedge \tilde{P})$$

(and consequently,  $\mathcal{E}r(P, \tilde{P}) = P \wedge \tilde{P}$ ).

2. Let  $P, P_n, n \geq 1$ , be probability measures on  $(R, \mathcal{B}(R))$  with densities (with respect to Lebesgue measure)  $p(x), p_n(x), n \geq 1$ . Let  $p_n(x) \rightarrow p(x)$  for almost all  $x$  (with respect to Lebesgue measure). Show that then

$$\|P - P_n\| = \int_{-\infty}^{\infty} |p(x) - p_n(x)| dx \rightarrow 0, \quad n \rightarrow \infty$$

(compare Problem 17 in §6, Chapter II).

3. Let  $P$  and  $\tilde{P}$  be two probability measures. We define Kullback information  $K(P, \tilde{P})$  as information by using  $P$  against  $\tilde{P}$ , by the equation

$$K(P, \tilde{P}) = \begin{cases} E \ln(dP/d\tilde{P}) & \text{if } P \ll \tilde{P}, \\ \infty & \text{otherwise.} \end{cases}$$

Show that

$$K(P, \tilde{P}) \geq -2 \ln(1 - \rho^2(P, \tilde{P})) \geq 2\rho^2(P, \tilde{P}).$$

4. Establish formulas (11) and (12).  
 5. Prove inequalities (24).  
 6. Let  $P, \tilde{P}$ , and  $Q$  be probability measures on  $(R, \mathcal{B}(R))$ ;  $P * Q$  and  $\tilde{P} * Q$ , their convolutions (see subsection 4, §8, Chapter II). Then
- $$\|P * Q - \tilde{P} * Q\| \leq \|P - \tilde{P}\|.$$
7. Prove (30).

## §10. Contiguity and Entire Asymptotic Separation of Probability Measures

1. These concepts play a fundamental role in the asymptotic theory of mathematical statistics, being natural extensions of the concepts of absolute conti-

nuity and singularity of two measures in the case of *sequences* of pairs of measures.

Let us begin with definitions.

Let  $(\Omega^n, \mathcal{F}^n)_{n \geq 1}$  be a sequence of measurable spaces; let  $(P^n)_{n \geq 1}$  and  $(\tilde{P}^n)_{n \geq 1}$  be sequences of probability measures with  $P^n$  and  $\tilde{P}^n$  defined on  $(\Omega^n, \mathcal{F}^n)$ ,  $n \geq 1$ .

**Definition 1.** We say that a sequence  $(\tilde{P}^n)$  of measures is *contiguous* to the sequence  $(P^n)$  (notation:  $(\tilde{P}^n) \triangleleft (P^n)$ ) if, for all  $A^n \in \mathcal{F}^n$  such that  $P^n(A^n) \rightarrow 0$  as  $n \rightarrow \infty$ , we have  $\tilde{P}^n(A^n) \rightarrow 0$ ,  $n \rightarrow \infty$ .

**Definition 2.** We say that sequences  $(\tilde{P}^n)$  and  $(P^n)$  of measures are *entirely (asymptotically) separated* (or for short:  $(\tilde{P}^n) \triangle (P^n)$ ), if there is a subsequence  $n_k \uparrow \infty$ ,  $k \rightarrow \infty$ , and sets  $A^{n_k} \in \mathcal{F}^{n_k}$  such that

$$P^{n_k}(A^{n_k}) \rightarrow 1 \quad \text{and} \quad \tilde{P}^{n_k}(A^{n_k}) \rightarrow 0, \quad k \rightarrow \infty.$$

We notice immediately that entire separation is a *symmetric* concept:  $(\tilde{P}^n) \triangle (P^n) \Leftrightarrow (P^n) \triangle (\tilde{P}^n)$ . Contiguity does not have this property. If  $(\tilde{P}^n) \triangleleft (P^n)$  and  $(P^n) \triangleleft (\tilde{P}^n)$ , we write  $(\tilde{P}^n) \triangleleft \triangleright (P^n)$  and say that the sequences  $(P^n)$  and  $(\tilde{P}^n)$  of measures are *mutually contiguous*.

We notice that in the case when  $(\Omega^n, \mathcal{F}^n) = (\Omega, \mathcal{F})$ ,  $P' = P$ ,  $\tilde{P}' = \tilde{P}$  for all  $n \geq 1$ , we have

$$(\tilde{P}^n) \triangleleft (P^n) \Leftrightarrow \tilde{P} \ll P, \quad (1)$$

$$(\tilde{P}^n) \triangleleft \triangleright (P^n) \Leftrightarrow \tilde{P} \sim P, \quad (2)$$

$$(\tilde{P}^n) \triangle (P^n) \Leftrightarrow \tilde{P} \perp P. \quad (3)$$

These properties and the definitions given above explain why contiguity and entire asymptotic separation are often thought of as “asymptotic absolute continuity” and “asymptotic singularity” for sequences  $(\tilde{P}^n)$  and  $(P^n)$ .

2. Theorems 1 and 2 presented below are natural extensions of Theorems 2 and 3 of §8 to sequences of measures.

Let  $(\Omega^n, \mathcal{F}^n)_{n \geq 1}$  be a sequence of measurable spaces;  $Q^n$ , a probability measure on  $(\Omega^n, \mathcal{F}^n)$ ; and  $\xi^n$ , a random variable (generally speaking, extended; see §4, Chapter II) on  $(\Omega^n, \mathcal{F}^n)$ ,  $n \geq 1$ .

**Definition 3.** A sequence  $(\xi^n)$  of random variables is *tight* with respect to a sequence of measures  $(Q^n)$  (notation:  $(\xi^n | Q^n)$  is *tight*) if

$$\lim_{N \uparrow \infty} \overline{\lim}_n Q^n(|\xi^n| > N) = 0. \quad (4)$$

(Compare the corresponding definition of tightness of a family of probability measures in §2.)

We shall always set

$$Q^n = \frac{P^n + \tilde{P}^n}{2}, \quad z^n = \frac{dP^n}{dQ^n}, \quad \tilde{z}^n = \frac{d\tilde{P}^n}{dQ^n}.$$

We shall also use the notation

$$Z^n = \tilde{z}^n / z^n \quad (5)$$

for the Lebesgue derivative of  $\tilde{P}^n$  with respect to  $P^n$  (see (29) in §9), taking  $2/0 = \infty$ . We note that if  $\tilde{P}^n \ll P^n$ ,  $Z^n$  is precisely one of the versions of the density  $d\tilde{P}^n/dP^n$  of the measure  $\tilde{P}^n$  with respect to  $P^n$  (see §6, Chapter II).

For later use it is convenient to note that since

$$P^n\left(z^n \leq \frac{1}{N}\right) = E_{Q^n}\left(z^n I\left(z^n \leq \frac{1}{N}\right)\right) \leq \frac{1}{N} \quad (6)$$

and  $Z^n \leq 2/z^n$ , we have

$$((1/z^n)|P^n) \quad \text{tight}, \quad (Z^n|P^n) \quad \text{tight}. \quad (7)$$

**Theorem 1.** *The following statements are equivalent:*

- (a)  $(\tilde{P}^n) \triangleleft (P^n)$ ,
- (b)  $(z^{-n}|\tilde{P}^n)$  is tight,
- (b')  $(Z^n|\tilde{P}^n)$  is tight,
- (c)  $\lim_{\alpha \downarrow 0} \underline{\lim}_n H(\alpha; P^n, \tilde{P}^n) = 1$ .

**Theorem 2.** *The following statements are equivalent:*

- (a)  $(\tilde{P}^n) \triangle (P^n)$ ,
- (b)  $\underline{\lim}_n \tilde{P}^n(z^n \geq \varepsilon) = 0$  for every  $\varepsilon > 0$ ,
- (b')  $\underline{\lim}_n \tilde{P}^n(Z^n \leq N) = 0$  for every  $N > 0$ ,
- (c)  $\lim_{\alpha \downarrow 0} \underline{\lim}_n H(\alpha; P^n, \tilde{P}^n) = 0$ ,
- (d)  $\underline{\lim}_n H(\alpha; P^n, \tilde{P}^n) = 0$  for all  $\alpha \in (0, 1)$ ,
- (e)  $\underline{\lim}_n H(\alpha; P^n, \tilde{P}^n) = 0$  for some  $\alpha \in (0, 1)$ .

**PROOF OF THEOREM 1.**

(a)  $\Rightarrow$  (b). If (b) is not satisfied, there are an  $\varepsilon > 0$  and a sequence  $n \uparrow \infty$  such that  $\tilde{P}^{n_k}(z^{n_k} < 1/n_k) \geq \varepsilon$ . But by (6),  $P^{n_k}(z^{n_k} < 1/n_k) \leq 1/n_k$ ,  $k \rightarrow \infty$ , which contradicts the assumption that  $(\tilde{P}^n) \triangleleft (P^n)$ .

(b)  $\Leftrightarrow$  (b'). We have only to note that  $Z^n = 2z^{-k} - 1$ .

(b)  $\Rightarrow$  (a). Let  $A^n \in \mathcal{F}^n$  and  $P^n(A^n) \rightarrow 0$ ,  $n \rightarrow \infty$ . We have

$$\begin{aligned} \tilde{P}^n(A^n) &\leq \tilde{P}^n(z^n \leq \varepsilon) + E_{Q^n}(z^n I(A^n \cap \{z^n > \varepsilon\})) \\ &\leq \tilde{P}^n(z^n \leq \varepsilon) + \frac{2}{\varepsilon} E_{Q^n}(z^n I(A^n)) = \tilde{P}^n(z^n \leq \varepsilon) + \frac{2}{\varepsilon} P^n(A^n). \end{aligned}$$

Therefore,

$$\lim_n \tilde{P}^n(A^n) \leq \lim_n \tilde{P}^n(z^n \leq \varepsilon), \quad \varepsilon > 0.$$

Proposition (b) is equivalent to saying that  $\lim_{\varepsilon \downarrow 0} \lim_n \tilde{P}^n(z^n \leq \varepsilon) = 0$ . Therefore,  $\tilde{P}^n(A^n) \rightarrow 0$ , i.e., (b)  $\Rightarrow$  (a).

(b)  $\Rightarrow$  (c). Let  $\varepsilon > 0$ . Then

$$\begin{aligned} H(\alpha; P^n, \tilde{P}^n) &= E_{Q^n}[(z^n)^\alpha (\tilde{z}^n)^{1-\alpha}] \\ &\geq E_{Q^n} \left[ \left( \frac{z^n}{\tilde{z}^n} \right)^\alpha I(z^n \geq \varepsilon) I(\tilde{z}^n > 0) \tilde{z}^n \right] \\ &= E_{\tilde{P}^n} \left[ \left( \frac{z^n}{\tilde{z}^n} \right)^\alpha I(z^n \geq \varepsilon) \right] \geq \left( \frac{\varepsilon}{2} \right)^\alpha \tilde{P}^n(z^n \geq \varepsilon), \end{aligned} \quad (8)$$

since  $z^n + \tilde{z}^n = 2$ . Therefore, for  $\varepsilon > 0$ ,

$$\lim_{\alpha \downarrow 0} \lim_n H(\alpha; P^n, \tilde{P}^n) \geq \lim_{\alpha \downarrow 0} \left( \frac{\varepsilon}{2} \right)^\alpha \lim_n \tilde{P}^n(z^n \geq \varepsilon) = \lim_n \tilde{P}^n(z^n \geq \varepsilon). \quad (9)$$

By (b),  $\lim_{\varepsilon \downarrow 0} \lim_n \tilde{P}^n(z^n \geq \varepsilon) = 1$ . Hence, (c) follows from (9) and the fact that  $H(\alpha; P^n, \tilde{P}^n) \leq 1$ .

(c)  $\Rightarrow$  (b). Let  $\delta \in (0, 1)$ . Then

$$\begin{aligned} H(\alpha; P^n, \tilde{P}^n) &= E_{Q^n}[(z^n)^\alpha (\tilde{z}^n)^{1-\alpha} I(z^n < \varepsilon)] \\ &\quad + E_{Q^n}[(z^n)^\alpha (\tilde{z}^n)^{1-\alpha} I(z^n \geq \varepsilon, \tilde{z}^n \leq \delta)] \\ &\quad + E_{Q^n}[(z^n)^\alpha (\tilde{z}^n)^{1-\alpha} I(z^n \geq \varepsilon, \tilde{z}^n > \delta)] \\ &\leq 2\varepsilon^\alpha + 2\delta^{1-\alpha} + E_{Q^n} \left[ \tilde{z}^n \left( \frac{z^n}{\tilde{z}^n} \right)^\alpha I(z^n \geq \varepsilon, \tilde{z}^n > \delta) \right] \\ &\leq 2\varepsilon^\alpha + 2\delta^{1-\alpha} + \left( \frac{2}{\delta} \right)^\alpha \tilde{P}^n(z^n \geq \varepsilon). \end{aligned} \quad (10)$$

Consequently,

$$\lim_{\varepsilon \downarrow 0} \lim_n \tilde{P}^n(z^n \geq \varepsilon) \geq \left( \frac{\delta}{2} \right)^\alpha \lim_n H(\alpha; P^n, \tilde{P}^n) - \frac{2}{2^\alpha} \delta$$

for all  $\alpha \in (0, 1)$ ,  $\delta \in (0, 1)$ . If we first let  $\alpha \downarrow 0$ , use (c), and then let  $\delta \downarrow 0$ , we obtain

$$\lim_{\varepsilon \downarrow 0} \lim_n \tilde{P}^n(z^n \geq \varepsilon) \geq 1,$$

from which (b) follows.

PROOF OF THEOREM 2.

(a)  $\Rightarrow$  (b). Let  $(\tilde{P}^n) \triangle (P^n)$ ,  $n_k \uparrow \infty$ , and let  $A^{n_k} \in \mathcal{F}^{n_k}$  have the property that  $P^{n_k}(A^{n_k}) \rightarrow 1$  and  $\tilde{P}^{n_k}(A^{n_k}) \rightarrow 0$ . Then, since  $z^n + \tilde{z}^n = 2$ , we have



$$\begin{aligned}
\tilde{P}^{n_k}(Z^{n_k} \geq \varepsilon) &\leq \tilde{P}^{n_k}(A^{n_k}) + E_{Q^{n_k}} \left\{ Z^{n_k} \cdot \frac{\tilde{Z}^{n_k}}{Z^{n_k}} I(\bar{A}^{n_k}) I(Z^{n_k} \geq \varepsilon) \right\} \\
&= \tilde{P}^{n_k}(A^{n_k}) + E_{P^{n_k}} \left\{ \frac{\tilde{Z}^{n_k}}{Z^{n_k}} I(\bar{A}^{n_k}) I(Z^{n_k} \geq \varepsilon) \right\} \\
&\leq \tilde{P}^{n_k}(A^{n_k}) + \frac{2}{\varepsilon} P^{n_k}(\bar{A}^{n_k}).
\end{aligned}$$

Consequently,  $\tilde{P}^{n_k}(Z^{n_k} \geq \varepsilon) \rightarrow 0$  and therefore, (b) is satisfied.

(b)  $\Rightarrow$  (a). If (b) is satisfied, there is a sequence  $n_k \uparrow \infty$  such that

$$\tilde{P}^{n_k} \left( Z^{n_k} \geq \frac{1}{k} \right) \leq \frac{1}{k} \rightarrow 0, \quad k \rightarrow \infty.$$

Hence, having observed (see (6)) that  $P^{n_k}(Z^{n_k} \geq 1/k) \geq 1 - (1/k)$ , we obtain (a).

(b)  $\Rightarrow$  (b'). We have only to observe that  $Z^n = (2/Z^n) - 1$ .

(b)  $\Rightarrow$  (d). By (10) and (b),

$$\lim_n H(\alpha; P^n, \tilde{P}^n) \leq 2\varepsilon^\alpha + 2\delta^{1-\alpha}$$

for arbitrary  $\varepsilon$  and  $\delta$  on the interval  $(0, 1)$ . Therefore, (d) is satisfied.

(d)  $\Rightarrow$  (c) and (d)  $\Rightarrow$  (e) are evident.

Finally, from (8) we have

$$\lim_n \tilde{P}^n(Z^n \geq \varepsilon) \leq \left( \frac{2}{\varepsilon} \right)^\alpha \lim_n H(\alpha; P^n, \tilde{P}^n).$$

Therefore, (c)  $\Rightarrow$  (b) and (e)  $\Rightarrow$  (b), since  $(2/\varepsilon)^\alpha \rightarrow 1$ ,  $\alpha \downarrow 0$ .

3. We now consider a special case corresponding to the method of independent observations, where the calculation of the integrals  $H(\alpha; P^n, \tilde{P}^n)$  and application of Theorems 1 and 2 do not present much difficulty.

Let us suppose that the measures  $P^n$  and  $\tilde{P}^n$  are direct products of measures:

$$P^n = P_1 \times \cdots \times P_n, \quad \tilde{P}^n = \tilde{P}_1 \times \cdots \times \tilde{P}_n, \quad n \geq 1,$$

where  $P_k$  and  $\tilde{P}_k$  are given on  $(\Omega_k, \mathcal{F}_k)$ ,  $k \geq 1$ .

Since in this case

$$H(\alpha; P^n, \tilde{P}^n) = \prod_{k=1}^n H(\alpha; P_k, \tilde{P}_k) = e^{\sum_{k=1}^n \ln[1 - (1 - H(\alpha; P_k, \tilde{P}_k))]},$$

we obtain the following result from Theorems 1 and 2:

$$(\tilde{P}^n) \triangleleft (P^n) \Leftrightarrow \lim_{\alpha \downarrow 0} \overline{\lim}_n \sum_{k=1}^n [1 - H(\alpha; P_k, \tilde{P}_k)] = 0, \quad (11)$$

$$(\tilde{P}^n) \triangle (P^n) \Leftrightarrow \overline{\lim}_n \sum_{k=1}^n [1 - H(\alpha; P_k, \tilde{P}_k)] = \infty. \quad (12)$$

EXAMPLE. Let  $(\Omega_k, \mathcal{F}_k) = (R, \mathcal{B}(R))$ ,  $a_k \in [0, 1)$ ,

$$P_k(dx) = I_{[0,1]}(x) dx, \quad \tilde{P}_k(dx) = \frac{1}{1-a_k} I_{[a_k,1]}(x) dx.$$

Since here  $H(\alpha; P_k, \tilde{P}_k) = (1-a_k)^\alpha$ ,  $\alpha \in (0, 1)$ , from (11) and the fact that  $H(\alpha; P_k, \tilde{P}_k) = H(1-\alpha; \tilde{P}_k, P_k)$ , we obtain

$$(\tilde{P}^n) \triangleleft (P^n) \Leftrightarrow \overline{\lim}_n na_n < \infty, \quad \text{i.e.,} \quad a_n = O\left(\frac{1}{n}\right),$$

$$(P^n) \triangleleft (\tilde{P}^n) \Leftrightarrow \overline{\lim}_n na_n = 0, \quad \text{i.e.,} \quad a_n = o\left(\frac{1}{n}\right),$$

$$(\tilde{P}^n) \triangle (\tilde{P}^n) \Leftrightarrow \overline{\lim}_n na_n = \infty.$$

#### 4. PROBLEMS

1. Let  $P^n = P_1^n \times \cdots \times P_n^n$ ,  $\tilde{P}^n = \tilde{P}_1^n \times \cdots \times \tilde{P}_n^n$ ,  $n \geq 1$ , where  $P_k^n$  and  $\tilde{P}_k^n$  are Gaussian measures with parameters  $(a_k^n, 1)$  and  $(\tilde{a}_k^n, 1)$ . Find conditions on  $(a_k^n)$  and  $(\tilde{a}_k^n)$  under which  $(\tilde{P}^n) \triangleleft (P^n)$  and  $(\tilde{P}^n) \triangle (P^n)$ .
2. Let  $P^n = P_1^n \times \cdots \times P_n^n$  and  $\tilde{P}^n = \tilde{P}_1^n \times \cdots \times \tilde{P}_n^n$ , where  $P_k^n$  and  $\tilde{P}_k^n$  are probability measures on  $(R, \mathcal{B}(R))$  for which  $P_k^n(dx) = I_{[0,1]}(x) dx$  and  $\tilde{P}_k^n(dx) = I_{[a_n, 1+a_n]}(dx)$ ,  $0 \leq a_n \leq 1$ . Show that  $H(\alpha; P_k^n, \tilde{P}_k^n) = 1 - a_n$  and

$$(\tilde{P}^n) \triangleleft (P^n) \Leftrightarrow (P^n) \triangleleft (\tilde{P}^n) \Leftrightarrow \overline{\lim}_n na_n = 0, \quad (\tilde{P}^n) \triangle (P^n) \Leftrightarrow \overline{\lim}_n na_n = \infty.$$

## §11. Rapidity of Convergence in the Central Limit Theorem

1. Let  $\xi_{n1}, \dots, \xi_{nn}$  be a sequence of independent random variables,  $S_n = \xi_{n1} + \cdots + \xi_{nn}$ ,  $F_n(x) = P(S_n \leq x)$ . If  $S_n \rightarrow \mathcal{N}(0, 1)$ , then  $F_n(x) \rightarrow \Phi(x)$  for every  $x \in R$ . Since  $\Phi(x)$  is continuous, the convergence here is actually uniform (Problem 5 in §1):

$$\sup_x |F_n(x) - \Phi(x)| \rightarrow 0, \quad n \rightarrow \infty. \quad (1)$$

It is natural to ask how rapid the convergence in (1) is. We shall establish a result for the case when

$$S_n = \frac{\xi_1 + \cdots + \xi_n}{\sigma\sqrt{n}}, \quad n \geq 1,$$

where  $\xi_1, \xi_2, \dots$  is a sequence of independent identically distributed random variables with  $E\xi_k = 0$ ,  $V\xi_k = \sigma^2$  and  $E|\xi_1|^3 < \infty$ .

**Theorem (Berry and Esseen).** *We have the bound*

$$\sup_x |F_n(x) - \Phi(x)| \leq \frac{CE|\xi_1|^3}{\sigma^3\sqrt{n}}, \quad (2)$$

where  $C$  is an absolute constant  $((2\pi)^{-1/2} \leq C < 0.8)$ .

**PROOF.** For simplicity, let  $\sigma^2 = 1$  and  $\beta_3 = E|\xi_1|^3$ . By Esseen's inequality (Subsection 10, §12, Chapter II)

$$\sup_x |F_n(x) - \Phi(x)| \leq \frac{2}{\pi} \int_0^T \left| \frac{f_n(t) - \varphi(t)}{t} \right| dt + \frac{24}{\pi T} \frac{1}{\sqrt{2\pi}} \quad (3)$$

where  $\varphi(t) = e^{-t^2/2}$  and

$$f_n(t) = [f(t/\sqrt{n})]^n,$$

with  $f(t) = Ee^{it\xi_1}$ .

In (3) we may take  $T$  arbitrarily. Let us choose

$$T = \sqrt{n}/(5\beta_3).$$

We are going to show that for this  $T$ ,

$$|f_n(t) - \varphi(t)| \leq \frac{7}{6} \frac{\beta_3}{\sqrt{n}} |t|^3 e^{-t^2/4}, \quad |t| \leq T. \quad (4)$$

The required estimate (2), with  $C$  an absolute constant, will follow immediately from (3) by means of (4). (A more detailed analysis shows that  $C < 0.8$ .)

We now turn to the proof of (4).

By formula (18) from §2, Chapter II ( $n = 3$ ,  $E\xi_1 = 0$ ,  $E\xi_1^2 = 1$ ,  $E|\xi_1|^3 < \infty$ ) we obtain

$$f(t) = Ee^{it\xi_1} = 1 - \frac{t^2}{2} + \frac{(it)^3}{6} [E\xi_1^3(\cos \theta_1 t\xi_1 + i \sin \theta_1 t\xi_1)], \quad (5)$$

where  $|\theta_1| \leq 1$ ,  $|\theta_2| \leq 1$ . Consequently,

$$f\left(\frac{t}{\sqrt{n}}\right) = 1 - \frac{t^2}{2n} + \frac{(it)^3}{6n^{3/2}} \left[ E\xi_1^3 \left( \cos \theta_1 \frac{t}{\sqrt{n}} \xi_1 + i \sin \theta_1 \frac{t}{\sqrt{n}} \xi_1 \right) \right].$$

If  $|t| \leq T = \sqrt{n}/5\beta_3$ , we find, by using the inequality  $\beta_3 \geq \sigma^3 = 1$  (see (28), §6, Chapter II), that

$$1 - \left| f\left(\frac{t}{\sqrt{n}}\right) \right| \leq \left| 1 - f\left(\frac{t}{\sqrt{n}}\right) \right| \leq \frac{t^2}{2n} + \frac{|t|^3 \beta_3}{3n^{3/2}} \leq \frac{1}{25}.$$

Consequently, for  $|t| \leq T$  it is possible to have the representation

$$\left[ f\left(\frac{t}{\sqrt{n}}\right) \right]^n = e^{n \ln f(t/\sqrt{n})}, \quad (6)$$

where  $\ln z$  means the principal value of the logarithm of the complex number  $z$  ( $\ln z = \ln|z| + i \arg z$ ,  $-\pi < \arg z \leq \pi$ ).

Since  $\beta_3 < \infty$ , we obtain from Taylor's theorem with the Lagrange remainder (compare (35) in §12, Chapter II)

$$\begin{aligned} \ln f\left(\frac{t}{\sqrt{n}}\right) &= \frac{it}{\sqrt{n}} s_{\xi_1}^{(1)} + \frac{(it)^2}{2n} s_{\xi_1}^{(2)} + \frac{(it)^3}{6n^{3/2}} (\ln f)''' \left(\theta \frac{t}{\sqrt{n}}\right) \\ &= -\frac{t^2}{2n} + \frac{(it)^3}{6n^{3/2}} (\ln f)''' \left(\theta \frac{t}{\sqrt{n}}\right), \quad |\theta| \leq 1, \end{aligned} \quad (7)$$

since the semi-invariants are  $s_{\xi_1}^{(1)} = E\xi_1 = 0$ ,  $s_{\xi_1}^{(2)} = \sigma^2 = 1$ .

In addition,

$$\begin{aligned} (\ln f(s))''' &= \frac{f'''(s) \cdot f^2(s) - 3f''(s)f'(s)f(s) + 2(f'(s))^3}{f^3(s)} \\ &= \frac{E[(i\xi_1)^3 e^{i\xi_1 s}] f^2(s) - 3E[(i\xi_1)^2 e^{i\xi_1 s}] E[(i\xi_1) e^{i\xi_1 s}] f(s) + 2E[(i\xi_1) e^{i\xi_1 s}]^3}{f^3(s)}. \end{aligned}$$

From this, taking into account that  $|f(t/\sqrt{n})| \geq 24/25$  for  $|t| \leq T$  and  $|f(s)| \leq 1$ , we obtain

$$\left| (\ln f)''' \left(\theta \frac{t}{\sqrt{n}}\right) \right| \leq \frac{\beta_3 + 3\beta_1 \cdot \beta_2 + 2\beta_1^3}{\left(\frac{24}{25}\right)^3} \leq 7\beta_3 \quad (8)$$

( $\beta_k = E|\xi_1|^k$ ,  $k = 1, 2, 3$ ;  $\beta_1 \leq \beta_2^{1/2} \leq \beta_3^{1/3}$ ; see (28), §6, Chapter II).

From (6)–(8), using the inequality  $|e^z - 1| \leq |z|e^{|z|}$ , we find for  $|t| \leq T = \sqrt{n}/5\beta_3$  that

$$\begin{aligned} \left| \left[ f\left(\frac{t}{\sqrt{n}}\right) \right]^n - e^{t^2/2} \right| &= |e^{n \ln f(t/\sqrt{n})} - e^{-t^2/2}| \\ &\leq \left(\frac{7}{6}\right) \frac{\beta_3 |t|^3}{\sqrt{n}} \exp \left\{ -\frac{t^2}{2} + \left(\frac{7}{6}\right) |t|^3 \frac{\beta_3}{\sqrt{n}} \right\} \leq \frac{7}{6} \frac{\beta_3 |t|^3}{\sqrt{n}} e^{-t^2/4}. \end{aligned}$$

This completes the proof of the theorem.

**Remark.** We observe that unless we make some supplementary hypothesis about the behavior of the random variables that are added, (2) cannot be improved. In fact, let  $\xi_1, \xi_2, \dots$  be independent identically distributed Bernoulli random variables with

$$P(\xi_k = 1) = P(\xi_k = -1) = \frac{1}{2}.$$

It is evident by symmetry that

$$2P\left(\sum_{k=1}^{2n} \xi_k < 0\right) + P\left(\sum_{k=1}^{2n} \xi_k = 0\right) = 1,$$

and hence, by Stirling's formula ((6), §2, chap. I)

$$\begin{aligned} \left| P\left(\sum_{k=1}^{2n} \xi_k < 0\right) - \frac{1}{2} \right| &= \frac{1}{2} P\left(\sum_{k=1}^{2n} \xi_k = 0\right) \\ &= \frac{1}{2} C_{2n}^n \cdot 2^{-2n} \sim \frac{1}{2\sqrt{\pi n}} = \frac{1}{\sqrt{(2\pi) \cdot (2n)}}. \end{aligned}$$

It follows, in particular, that the constant  $C$  in (2) cannot be less than  $(2\pi)^{-1/2}$ .

## 2. PROBLEMS

1. Prove (8).

2. Let  $\xi_1, \xi_2, \dots$  be independent identically distributed random variables with  $E\xi_k = 0$ ,  $V\xi_k = \sigma^2$  and  $E|\xi_1|^3 < \infty$ .

It is known that the following *nonuniform inequality* holds: for all  $x \in R$ ,

$$|F_n(x) - \Phi(x)| \leq \frac{CE|\xi_1|^3}{\sigma^3\sqrt{n}} \cdot \frac{1}{(1+|x|)^3}.$$

Prove this, at least for Bernoulli random variables.

## §12. Rapidity of Convergence in Poisson's Theorem

1. Let  $\xi_1, \xi_2, \dots, \xi_n$  be independent Bernoulli random variables that take the values 1 and 0 with probabilities

$$P(\xi_k = 1) = p_k, \quad P(\xi_k = 0) = q_k (= 1 - p_k), \quad 1 \leq k \leq n.$$

We set  $S = \xi_1 + \dots + \xi_n$ ; let  $B = (B_0, B_1, \dots, B_n)$  be the binomial distribution of probabilities of the sum  $S$ , where  $B_k = P(S = k)$ . Also let  $\Pi = (\Pi_0, \Pi_1, \dots)$  be the Poisson distribution with parameter  $\lambda$ , where

$$\Pi_k = \frac{e^{-\lambda} \lambda^k}{k!}, \quad k \geq 0.$$

We noticed in §6, Chapter I, that if

$$p_1 = \dots = p_n, \quad \lambda = np, \quad (1)$$

there is the following estimate (Prokhorov) for the *distance in variation* between the measures  $B$  and  $\Pi$  ( $B_{n+1} = B_{n+2} = \dots = 0$ ):

$$\|B - \Pi\| = \sum_{k=0}^{\infty} |B_k - \Pi_k| \leq C_1(\lambda)p = C_1(\lambda) \cdot \frac{\lambda}{n}, \quad (2)$$

where

$$C_1(\lambda) = 2 \min(2, \lambda). \quad (3)$$

For the case when  $p_k$  are not necessarily equal, but satisfy  $\sum_{k=1}^n p_k = \lambda$ , LeCam showed that

$$\|B - \Pi\| = \sum_{k=0}^{\infty} |B_k - \Pi_k| \leq C_2(\lambda) \max_{1 \leq k \leq n} p_k, \quad (4)$$

where

$$C_2(\lambda) = 2 \min(9, \lambda). \quad (5)$$

A theorem to be presented below will imply the estimate

$$\|B - \Pi\| \leq C_3(\lambda) \max_{1 \leq k \leq n} p_k, \quad (6)$$

in which

$$C_3(\lambda) = 2\lambda. \quad (7)$$

Although  $C_2(\lambda) < C_3(\lambda)$  for  $\lambda > 9$ , i.e., (6) is worse than (4), we nevertheless have preferred to give a proof of (6), since this proof is essentially elementary, whereas an emphasis on obtaining a “good” constant  $C_2(\lambda)$  in (4) greatly complicates the proof.

**2. Theorem.** *Let  $\lambda = \sum_{k=1}^n p_k$ . Then*

$$\|B - \Pi\| = \sum_{k=0}^8 |B_k - \Pi_k| \leq 2 \sum_{k=1}^n p_k^2. \quad (8)$$

**PROOF.** We use the fact that each of the distributions  $B$  and  $\Pi$  is a *convolution* of distributions:

$$\begin{aligned} B &= B(p_1) * B(p_2) * \cdots * B(p_n), \\ \Pi &= \Pi(p_1) * \Pi(p_2) * \cdots * \Pi(p_n), \end{aligned} \quad (9)$$

understood as a convolution of the corresponding distribution functions (see subsection 4, §8, Chapter II), where  $B(p_k) = (1 - p_k, p_k)$  is a Bernoulli distribution on the points 0 and 1, and  $\Pi(p_k)$  is a Poisson distribution supported on the points 0, 1, ... with the parameters  $p_k$ .

It is easy to show that the difference  $B - \Pi$  can be represented in the form

$$B - \Pi = R_1 + \cdots + R_n, \quad (10)$$

where

$$R_k = (B(p_k) - \Pi(p_k)) * F_k \quad (11)$$

with

$$\begin{aligned} F_1 &= \Pi(p_2) * \cdots * \Pi(p_n), \\ &\dots\dots\dots \\ F_k &= B(p_1) * \cdots * B(p_{k-1}) * \Pi(p_{k+1}) * \cdots * \Pi(p_n), \quad 2 \leq k \leq n-1, \\ F_n &= B(p_1) * \cdots * B(p_{n-1}) \Pi. \end{aligned}$$

By problem 6 in §9, we have  $\|R_k\| \leq \|B(p_k) - \Pi(p_k)\|$ . Consequently, we see immediately from (10) that

$$\|B - \Pi\| \leq \sum_{k=1}^n \|B(p_k) - \Pi(p_k)\|. \quad (12)$$

By formula (12) in §9, we see that there is no difficulty in calculating the variation  $\|B(p_k) - \Pi(p_k)\|$ :

$$\begin{aligned} \|B(p_k) - \Pi(p_k)\| &= |(1 - p_k) - e^{-p_k}| + |p_k - p_k e^{-p_k}| + \sum_{j \geq 2} \frac{p_k^j e^{-p_k}}{j!} \\ &= |(1 - p_k) - e^{-p_k}| + |p_k - p_k e^{-p_k}| + 1 - e^{-p_k} - p_k e^{-p_k} \\ &= 2p_k(1 - e^{-p_k}) \leq 2p_k^2. \end{aligned}$$

From this, together with (12), we obtain the required inequality (8).

This completes the proof of the theorem.

**Corollary.** Since  $\sum_{k=1}^n p_k^2 \leq \lambda \max_{1 \leq k \leq n} p_k$ , we obtain (6).

### 3. PROBLEMS

1. Show that, if  $\lambda_k = -\ln(1 - p_k)$ ,

$$\|B(p_k) - \Pi(\lambda_k)\| = 2(1 - e^{-\lambda_k} - \lambda_k e^{-\lambda_k}) \leq \lambda_k^3$$

and consequently,  $\|B - \Pi\| \leq \sum_{k=1}^n \lambda_k^2$ .

2. Establish the representations (9) and (10).

## CHAPTER IV

# Sequences and Sums of Independent Random Variables

### §1. Zero-or-One Laws

1. The series  $\sum_{n=1}^{\infty} (1/n)$  diverges and the series  $\sum_{n=1}^{\infty} (-1)^n(1/n)$  converges. We ask the following question. What can we say about the convergence or divergence of a series  $\sum_{n=1}^{\infty} (\xi_n/n)$ , where  $\xi_1, \xi_2, \dots$  is a sequence of independent identically distributed Bernoulli random variables with  $P(\xi_1 = +1) = P(\xi_1 = -1) = \frac{1}{2}$ ? In other words, what can be said about the convergence of a series whose general term is  $\pm 1/n$ , where the signs are chosen in a random manner, according to the sequence  $\xi_1, \xi_2, \dots$ ?

Let

$$A_1 = \left\{ \omega: \sum_{n=1}^{\infty} \frac{\xi_n}{n} \text{ converges} \right\}$$

be the set of sample points for which  $\sum_{n=1}^{\infty} (\xi_n/n)$  converges (to a finite number) and consider the probability  $P(A_1)$  of this set. It is far from clear, to begin with, what values this probability might have. However, it is a remarkable fact that we are able to say that the probability can have only two values, 0 or 1. This is a corollary of Kolmogorov's "zero-one law," whose statement and proof form the main content of the present section.

2. Let  $(\Omega, \mathcal{F}, P)$  be a probability space, and let  $\xi_1, \xi_2, \dots$  be a sequence of random variables. Let  $\mathcal{F}_n^{\infty} = \sigma(\xi_n, \xi_{n+1}, \dots)$  be the  $\sigma$ -algebra generated by  $\xi_n, \xi_{n+1}, \dots$ , and write

$$\mathcal{X} = \bigcap_{n=1}^{\infty} \mathcal{F}_n^{\infty}.$$



Since an intersection of  $\sigma$ -algebras is again a  $\sigma$ -algebra,  $\mathcal{X}$  is a  $\sigma$ -algebra. It is called a *tail algebra* (or terminal or asymptotic algebra), because every event  $A \in \mathcal{X}$  is independent of the values of  $\xi_1, \dots, \xi_n$  for every finite number  $n$ , and is determined, so to speak, only by the behavior of the infinitely remote values of  $\xi_1, \xi_2, \dots$ .

Since, for every  $k \geq 1$ ,

$$A_1 \equiv \left\{ \sum_{n=1}^{\infty} \frac{\xi_n}{n} \text{ converges} \right\} = \left\{ \sum_{n=k}^{\infty} \frac{\xi_n}{n} \text{ converges} \right\} \in \mathcal{F}_k^{\infty},$$

we have  $A_1 \in \bigcap_k \mathcal{F}_k^{\infty} \equiv \mathcal{X}$ . In the same way, if  $\xi_1, \xi_2, \dots$  is any sequence,

$$A_2 = \left\{ \sum_n \xi_n \text{ converges} \right\} \in \mathcal{X}.$$

The following events are also tail events:

$$A_3 = \{\xi_n \in I_n \text{ for infinitely many } n\},$$

where  $I_n \in \mathcal{B}(R)$ ,  $n \geq 1$ ;

$$A_4 = \left\{ \overline{\lim}_n \xi_n < \infty \right\};$$

$$A_5 = \left\{ \overline{\lim}_n \frac{\xi_1 + \dots + \xi_n}{n} < \infty \right\};$$

$$A_6 = \left\{ \overline{\lim}_n \frac{\xi_1 + \dots + \xi_n}{n} < c \right\};$$

$$A_7 = \left\{ \frac{S_n}{n} \text{ converges} \right\};$$

$$A_8 = \left\{ \overline{\lim}_n \frac{S_n}{\sqrt{2n \log n}} = 1 \right\}.$$

On the other hand,

$$B_1 = \{\xi_n = 0 \text{ for all } n \geq 1\},$$

$$B_2 = \left\{ \lim_n (\xi_1 + \dots + \xi_n) \text{ exists and is less than } c \right\}$$

are examples of events that do not belong to  $\mathcal{X}$ .

Let us now suppose that our random variables are *independent*. Then by the Borel–Cantelli lemma it follows that

$$P(A_3) = 0 \Leftrightarrow \sum P(\xi_n \in I_n) < \infty,$$

$$P(A_3) = 1 \Leftrightarrow \sum P(\xi_n \in I_n) = \infty.$$

Therefore the probability of  $A_3$  can take only the values 0 or 1 according to the convergence or divergence of  $\sum P(\xi_n \in I_n)$ . This is Borel's zero-one law.

**Theorem 1** (Kolmogorov's Zero-One Law). *Let  $\xi_1, \xi_2, \dots$  be a sequence of independent random variables and let  $A \in \mathcal{X}$ . The  $P(A)$  can only have one of the values zero or one.*

**PROOF.** The idea of the proof is to show that every tail event  $A$  is independent of itself and therefore  $P(A \cap A) = P(A) \cdot P(A)$ , i.e.,  $P(A) = P^2(A)$ , so that  $P(A) = 0$  or 1.

If  $A \in \mathcal{X}$  then  $A \in \mathcal{F}_1^\infty = \sigma\{\xi_1, \xi_2, \dots\} = \sigma(\bigcup_n \mathcal{F}_1^n)$ , where  $\mathcal{F}_1^n = \sigma\{\xi_1, \dots, \xi_n\}$ , and we can find (Problem 8, §3, Chapter II) sets  $A_n \in \mathcal{F}_1^n$ ,  $m \geq 1$ , such that  $P(A \triangle A_n) \rightarrow 0$ ,  $n \rightarrow \infty$ . Hence

$$P(A_n) \rightarrow P(A), \quad P(A_n \cap A) \rightarrow P(A). \quad (1)$$

But if  $A \in \mathcal{X}$ , the events  $A_n$  and  $A$  are independent for every  $n \geq 1$ . Hence it follows from (1) that  $P(A) = P^2(A)$  and therefore  $P(A) = 0$  or 1.

This completes the proof of the theorem.

**Corollary.** *Let  $\eta$  be a random variable that is measurable with respect to the tail  $\sigma$ -algebra  $\mathcal{X}$ , i.e.,  $\{\eta \in B\} \in \mathcal{X}$ ,  $B \in \mathcal{B}(R)$ . Then  $\eta$  is degenerate, i.e., there is a constant  $c$  such that  $P(\eta = c) = 1$ .*

3. Theorem 2 below provides an example of a nontrivial application of Kolmogorov's zero-one law.

Let  $\xi_1, \xi_2, \dots$  be a sequence of independent Bernoulli random variables with  $P(\xi_n = 1) = p$ ,  $P(\xi_n = -1) = q$ ,  $p + q = 1$ ,  $n \geq 1$ , and let  $S_n = \xi_1 + \dots + \xi_n$ . It seems intuitively clear that in the symmetric case ( $p = \frac{1}{2}$ ) a "typical" path of the random walk  $S_n$ ,  $n \geq 1$ , will cross zero infinitely often, whereas when  $p \neq \frac{1}{2}$  it will go off to infinity. Let us give a precise formulation.

**Theorem 2.** (a) *If  $p = \frac{1}{2}$  then  $P(S_n = 0 \text{ i.o.}) = 1$ .*

(b) *If  $p \neq \frac{1}{2}$ , then  $P(S_n = 0 \text{ i.o.}) = 0$ .*

**PROOF.** We first observe that the event  $B = (S_n = 0 \text{ i.o.})$  is not a tail event, i.e.,  $B \notin \mathcal{X} = \bigcap \mathcal{F}_n^\infty$ ,  $\mathcal{F}_n^\infty = \sigma\{\xi_n, \xi_{n+1}, \dots\}$ . Consequently it is, in principle, not clear that  $B$  should have only the values 0 or 1.

Statement (b) is easily proved by applying (the first part of) the Borel-Cantelli lemma. In fact, if  $B_{2n} = \{S_{2n} = 0\}$ , then by Stirling's formula

$$P(B_{2n}) = C_{2n} p^n q^n \sim \frac{(4pq)^n}{\sqrt{\pi n}}$$

and therefore  $\sum P(B_{2n}) < \infty$ . Consequently  $P(S_n = 0 \text{ i.o.}) = 0$ .

To prove (a), it is enough to prove that the event

$$A = \left\{ \overline{\lim} \frac{S_n}{\sqrt{n}} = \infty, \underline{\lim} \frac{S_n}{\sqrt{n}} = -\infty \right\}$$

has probability 1 (since  $A \subseteq B$ ).

Let

$$A_c = \left\{ \overline{\lim} \frac{S_n}{\sqrt{n}} > c \right\} \cap \left\{ \lim \frac{S_n}{\sqrt{n}} < -c \right\} (= A'_c \cap A''_c).$$

Then  $A_c \downarrow A$ ,  $c \rightarrow \infty$ , and all the events  $A$ ,  $A_c$ ,  $A'_c$ ,  $A''_c$  are tail events. Let us show that  $P(A'_c) = P(A''_c) = 1$  for each  $c > 0$ . Since  $A'_c \in \mathcal{X}$  and  $A''_c \in \mathcal{X}$ , it is sufficient to show only that  $P(A'_c) > 0$ ,  $P(A''_c) > 0$ . But by Problem 5

$$P\left(\lim \frac{S_n}{\sqrt{n}} < -c\right) = P\left(\overline{\lim} \frac{S_n}{\sqrt{n}} > c\right) \geq \overline{\lim} P\left(\frac{S_n}{\sqrt{n}} > c\right) > 0,$$

where the last inequality follows from the De Moivre-Laplace theorem.

Thus  $P(A_c) = 1$  for all  $c > 0$  and therefore  $P(A) = \lim_{c \rightarrow \infty} P(A_c) = 1$ .

This completes the proof of the theorem.

4. Let us observe again that  $B = \{S_n = 0 \text{ i.o.}\}$  is not a tail event. Nevertheless, it follows from Theorem 2 that, for a Bernoulli scheme, the probability of this event, just as for tail events, takes only the values 0 and 1. This phenomenon is not accidental: it is a corollary of the Hewitt-Savage zero-one law, which for independent identically distributed random variables extends the result of Theorem 1 to the class of "symmetric" events (which includes the class of tail events).

Let us give the essential definitions. A one-to-one mapping  $\pi = (\pi_1, \pi_2, \dots)$  of the set  $(1, 2, \dots)$  on itself is said to be a finite permutation if  $\pi_n = n$  for every  $n$  with a finite number of exceptions.

If  $\xi = \xi_1, \xi_2, \dots$  is a sequence of random variables,  $\pi(\xi)$  denotes the sequence  $(\xi_{\pi_1}, \xi_{\pi_2}, \dots)$ . If  $A$  is the event  $\{\xi \in B\}$ ,  $B \in \mathcal{B}(R^\infty)$ , then  $\pi(A)$  denotes the event  $\{\pi(\xi) \in B\}$ ,  $B \in \mathcal{B}(R^\infty)$ .

We call an event  $A = \{\xi \in B\}$ ,  $B \in \mathcal{B}(R^\infty)$ , *symmetric* if  $\pi(A)$  coincides with  $A$  for every finite permutation  $\pi$ .

An example of a symmetric event is  $A = \{S_n = 0 \text{ i.o.}\}$ , where  $S_n = \xi_1 + \dots + \xi_n$ . Moreover, we may suppose (Problem 4) that every event in the tail  $\sigma$ -algebra  $\mathcal{X}(S) = \bigcap \mathcal{F}_n^\infty(S) = \sigma\{\omega: S_n, S_{n+1}, \dots\}$  generated by  $S_1 = \xi_1, S_2 = \xi_1 + \xi_2, \dots$  is symmetric.

**Theorem 3 (Hewitt-Savage Zero-One Law).** *Let  $\xi_1, \xi_2, \dots$  be a sequence of independent identically distributed random variables, and*

$$A = \{\omega: (\xi_1, \xi_2, \dots) \in B\}$$

*a symmetric event. Then  $P(A) = 0$  or  $1$ .*

**PROOF.** Let  $A = \{\xi \in B\}$  be a symmetric event. Choose sets  $B_n \in \mathcal{B}(R^n)$  such that, for  $A_n = \{\omega: (\xi_1, \dots, \xi_n) \in B_n\}$ ,

$$P(A \triangle A_n) \rightarrow 0, \quad n \rightarrow \infty. \quad (2)$$

Since the random variables  $\xi_1, \xi_2, \dots$  are independent and identically distributed, the probability distributions  $P_\xi(B) = P(\xi \in B)$  and  $P_{\pi_n(\xi)}(B) = P(\pi_n(\xi) \in B)$  coincide. Therefore

$$P(A \triangle A_n) = P_\xi(B \triangle B_n) = P_{\pi_n(\xi)}(B \triangle B_n). \quad (3)$$

Since  $A$  is symmetric, we have

$$A \equiv \{\xi \in B\} = \pi_n(A) \equiv \{\pi_n(\xi) \in B\}.$$

Therefore

$$\begin{aligned} P_{\pi_n(\xi)}(B \triangle B_n) &= P\{\pi_n(\xi) \in B\} \triangle (\pi_n(\xi) \in B_n)\} \\ &= P\{(\xi \in B) \triangle (\pi_n(\xi) \in B_n)\} = P\{A \triangle \pi_n(A_n)\}. \end{aligned} \quad (4)$$

Hence, by (3) and (4),

$$P(A \triangle A_n) = P(A \triangle \pi_n(A_n)). \quad (5)$$

It then follows from (2) that

$$P(A \triangle (A_n \cap \pi_n(A_n))) \rightarrow 0, \quad n \rightarrow \infty. \quad (6)$$

Hence, by (2), (5), and (6), we obtain

$$\begin{aligned} P(A_n) &\rightarrow P(A), & P(\pi_n(A)) &\rightarrow P(A), \\ P(A_n \cap \pi_n(A_n)) &\rightarrow P(A). \end{aligned} \quad (7)$$

Moreover, since  $\xi_1$  and  $\xi_2$  are independent,

$$\begin{aligned} P(A_n \cap \pi_n(A_n)) &= P\{(\xi_1, \dots, \xi_n) \in B_n, (\xi_{n+1}, \dots, \xi_{2n}) \in B_n\} \\ &= P\{(\xi_1, \dots, \xi_n) \in B_n\} \cdot P\{(\xi_{n+1}, \dots, \xi_{2n}) \in B_n\} \\ &= P(A_n)P(\pi_n(A_n)), \end{aligned}$$

whence by (7)

$$P(A) = P^2(A)$$

and therefore  $P(A) = 0$  or  $1$ .

This completes the proof of the theorem.

## 5. PROBLEMS

1. Prove the corollary to Theorem 1.
2. Show that if  $(\xi_n)$  is a sequence of independent random variables, the random variables  $\overline{\lim} \xi_n$  and  $\underline{\lim} \xi_n$  are degenerate.
3. Let  $(\xi_n)$  be a sequence of independent random variables,  $S_n = \xi_1 + \dots + \xi_n$ , and let the constants  $b_n$  satisfy  $0 < b_n \uparrow \infty$ . Show that the random variables  $\overline{\lim}(S_n/b_n)$  and  $\underline{\lim}(S_n/b_n)$  are degenerate.
4. Let  $S_n = \xi_1 + \dots + \xi_n$ ,  $n \geq 1$ , and  $\mathcal{X}(S) = \bigcap \mathcal{F}_n^\infty(S)$ ,  $\mathcal{F}_n^\infty(S) = \sigma\{\omega: S_n, S_{n+1}, \dots\}$ . Show that every event in  $\mathcal{X}(S)$  is symmetric.
5. Let  $(\xi_n)$  be a sequence of random variables. Show that  $\{\overline{\lim} \xi_n > c\} \supseteq \overline{\lim}\{\xi_n > c\}$  for each  $c > 0$ .

## §2. Convergence of Series

1. Let us suppose that  $\xi_1, \xi_2, \dots$  is a sequence of independent random variables,  $S_n = \xi_1 + \dots + \xi_n$ , and let  $A$  be the set of sample points  $\omega$  for which  $\sum \xi_n(\omega)$  converges to a finite limit. It follows from Kolmogorov's zero-one law that  $P(A) = 0$  or  $1$ , i.e., the series  $\sum \xi_n$  converges or diverges with probability  $1$ . The object of the present section is to give criteria that will determine whether a sum of independent random variables converges or diverges.

**Theorem 1** (Kolmogorov and Khinchin).

(a) Let  $E\xi_n = 0$ ,  $n \geq 1$ . Then if

$$\sum E\xi_n^2 < \infty, \quad (1)$$

the series  $\sum \xi_n$  converges with probability  $1$ .

(b) If the random variables  $\xi_n$ ,  $n \geq 1$ , are uniformly bounded (i.e.,  $P(|\xi_n| \leq c) = 1$ ,  $c < \infty$ ), the converse is true: the convergence of  $\sum \xi_n$  with probability  $1$  implies (1).

The proof depends on

### Kolmogorov's Inequality

(a) Let  $\xi_1, \xi_2, \dots, \xi_n$  be independent random variables with  $E\xi_i = 0$ ,  $E\xi_i^2 < \infty$ ,  $i \leq n$ . Then for every  $\varepsilon > 0$

$$P\left\{\max_{1 \leq k \leq n} |S_k| \geq \varepsilon\right\} \leq \frac{ES_n^2}{\varepsilon^2}. \quad (2)$$

(b) If also  $P(|\xi_i| \leq c) = 1$ ,  $i \leq n$ , then

$$P\left\{\max_{1 \leq k \leq n} |S_k| \geq \varepsilon\right\} \geq 1 - \frac{(c + \varepsilon)^2}{ES_n^2}. \quad (3)$$

**PROOF.** (a) Put

$$A = \{\max |S_k| \geq \varepsilon\},$$

$$A_k = \{|S_i| < \varepsilon, i = 1, \dots, k-1, |S_k| \geq \varepsilon\}, \quad 1 \leq k \leq n.$$

Then  $A = \sum A_k$  and

$$ES_n^2 \geq ES_n^2 I_A = \sum ES_n^2 I_{A_k}.$$

But

$$\begin{aligned} ES_n^2 I_{A_k} &= E(S_k + (\xi_{k+1} + \cdots + \xi_n))^2 I_{A_k} \\ &= ES_k^2 I_{A_k} + 2ES_k(\xi_{k+1} + \cdots + \xi_n) I_{A_k} + E(\xi_{k+1} + \cdots + \xi_n)^2 I_{A_k} \\ &\geq ES_k^2 I_{A_k}, \end{aligned}$$

since

$$ES_k(\xi_{k+1} + \cdots + \xi_n) I_{A_k} = ES_k I_{A_k} \cdot E(\xi_{k+1} + \cdots + \xi_n) = 0$$

because of independence and the conditions  $E\xi_i = 0$ ,  $i \leq n$ . Hence

$$ES_n^2 \geq \sum ES_k^2 I_{A_k} \geq \varepsilon^2 \sum P(A_k) = \varepsilon^2 P(A),$$

which proves the first inequality.

(b) To prove (3), we observe that

$$ES_n^2 I_A = ES_n^2 - ES_n^2 I_{\bar{A}} \geq ES_n^2 - \varepsilon^2 P(\bar{A}) = ES_n^2 - \varepsilon^2 + \varepsilon^2 P(A). \quad (4)$$

On the other hand, on the set  $A_k$

$$|S_{k-1}| \leq \varepsilon, \quad |S_k| \leq |S_{k-1}| + |\xi_k| \leq \varepsilon + c$$

and therefore

$$\begin{aligned} ES_n^2 I_A &= \sum_k ES_k^2 I_{A_k} + \sum_k E(I_{A_k} (S_n - S_k)^2) \\ &\leq (\varepsilon + c)^2 \sum_k P(A_k) + \sum_{k=1}^n P(A_k) \sum_{j=k+1}^n E\xi_j^2 \\ &\leq P(A) \left[ (\varepsilon + c)^2 + \sum_{j=1}^n E\xi_j^2 \right] = P(A) [(\varepsilon + c)^2 + ES_n^2]. \quad (5) \end{aligned}$$

From (4) and (5) we obtain

$$P(A) \geq \frac{ES_n^2 - \varepsilon^2}{(\varepsilon + c)^2 + ES_n^2 - \varepsilon^2} = 1 - \frac{(\varepsilon + c)^2}{(\varepsilon + c)^2 + ES_n^2 - \varepsilon^2} \geq 1 - \frac{(\varepsilon + c)^2}{ES_n^2}.$$

This completes the proof of (3).

**PROOF OF THEOREM 1.** (a) By Theorem 4 of §10, Chapter II, the sequence  $(S_n)$ ,  $n \geq 1$ , converges with probability 1, if and only if it is fundamental with probability 1. By Theorem 1 of §10, Chapter II, the sequence  $(S_n)$ ,  $n \geq 1$ , is fundamental (P-a.s.) if and only if

$$P\left\{\sup_{k \geq 1} |S_{n+k} - S_n| \geq \varepsilon\right\} \rightarrow 0, \quad n \rightarrow \infty. \quad (6)$$

By (2),

$$\begin{aligned} P\left\{\sup_{k \geq 1} |S_{n+k} - S_n| \geq \varepsilon\right\} &= \lim_{N \rightarrow \infty} P\left\{\max_{1 \leq k \leq N} |S_{n+k} - S_n| \geq \varepsilon\right\} \\ &\leq \lim_{N \rightarrow \infty} \frac{\sum_{k=n}^{n+N} E\xi_k^2}{\varepsilon^2} = \frac{\sum_{k=n}^{\infty} E\xi_k^2}{\varepsilon^2}. \end{aligned}$$

Therefore (6) is satisfied if  $\sum_{k=1}^{\infty} E\xi_k^2 < \infty$ , and consequently  $\sum \xi_k$  converges with probability 1.

(b) Let  $\sum \xi_k$  converge. Then, by (6), for sufficiently large  $n$ ,

$$P\left\{\sup_{k \geq 1} |S_{n+k} - S_n| \geq \varepsilon\right\} < \frac{1}{2}. \quad (7)$$

By (3),

$$P\left\{\sup_{k \geq 1} |S_{n+k} - S_n| \geq \varepsilon\right\} \geq 1 - \frac{(c + \varepsilon)^2}{\sum_{k=n}^{\infty} E\xi_k^2}.$$

Therefore if we suppose that  $\sum_{k=1}^{\infty} E\xi_k^2 = \infty$ , we obtain

$$P\left\{\sup_{k \geq 1} |S_{n+k} - S_n| \geq \varepsilon\right\} = 1,$$

which contradicts (7).

This completes the proof of the theorem.

**EXAMPLE.** If  $\xi_1, \xi_2, \dots$  is a sequence of independent Bernoulli random variables with  $P(\xi_n = +1) = P(\xi_n = -1) = \frac{1}{2}$ , then the series  $\sum \xi_n a_n$ , with  $|a_n| \leq c$ , converges with probability 1, if and only if  $\sum a_n^2 < \infty$ .

**2. Theorem 2 (Two-Series Theorem).** *A sufficient condition for the convergence of the series  $\sum \xi_n$  of independent random variables, with probability 1, is that both series  $\sum E\xi_n$  and  $\sum V\xi_n$  converge. If  $P(|\xi_n| \leq c) = 1$ , the condition is also necessary.*

**PROOF.** If  $\sum V\xi_n < \infty$ , then by Theorem 1 the series  $\sum (\xi_n - E\xi_n)$  converges (P-a.s.). But by hypothesis the series  $\sum E\xi_n$  converges; hence  $\sum \xi_n$  converges (P-a.s.).

To prove the necessity we use the following symmetrization method. In addition to the sequence  $\xi_1, \xi_2, \dots$  we consider a different sequence  $\tilde{\xi}_1, \tilde{\xi}_2, \dots$  of independent random variables such that  $\tilde{\xi}_n$  has the same distribution as  $\xi_n$ ,  $n \geq 1$ . (When the original sample space is sufficiently rich, the existence of such a sequence follows from Theorem 1 of §9, Chapter II. We can also show that this assumption involves no loss of generality.)

Then if  $\sum \xi_n$  converges (P-a.s.), the series  $\sum \tilde{\xi}_n$  also converges, and hence so does  $\sum (\xi_n - \tilde{\xi}_n)$ . But  $E(\xi_n - \tilde{\xi}_n) = 0$  and  $P(|\xi_n - \tilde{\xi}_n| \leq 2c) = 1$ . Therefore  $\sum V(\xi_n - \tilde{\xi}_n) < \infty$  by Theorem 1. In addition,

$$\sum V\xi_n = \frac{1}{2} \sum V(\xi_n - \tilde{\xi}_n) < \infty.$$

Consequently, by Theorem 1,  $\sum (\xi_n - E\xi_n)$  converges with probability 1, and therefore  $\sum E\xi_n$  converges.

Thus if  $\sum \xi_n$  converges (P-a.s.) (and  $P(|\xi_n| \leq c) = 1, n \geq 1$ ) it follows that both  $\sum E\xi_n$  and  $\sum V\xi_n$  converge.

This completes the proof of the theorem.

**3.** The following theorem provides a necessary and sufficient condition for the convergence of  $\sum \xi_n$  without any boundedness condition on the random variables.

Let  $c$  be a constant and

$$\xi^c = \begin{cases} \xi, & |\xi| \leq c, \\ 0, & |\xi| > c. \end{cases}$$

**Theorem 3** (Kolmogorov's Three-Series Theorem). *Let  $\xi_1, \xi_2, \dots$  be a sequence of independent random variables. A necessary condition for the convergence of  $\sum \xi_n$  with probability 1 is that the series*

$$\sum E \xi_n^c, \quad \sum V \xi_n^c, \quad \sum P(|\xi_n| \geq c)$$

*converge for every  $c > 0$ ; a sufficient condition is that these series converge for some  $c > 0$ .*

**PROOF.** *Sufficiency.* By the two-series theorem,  $\sum \xi_n^c$  converges with probability 1. But if  $\sum P(|\xi_n| \geq c) < \infty$ , then by the Borel-Cantelli lemma we have  $\sum I(|\xi_n| \geq c) < \infty$  with probability 1. Therefore  $\xi_n = \xi_n^c$  for all  $n$  with at most finitely many exceptions. Therefore  $\sum \xi_n$  also converges (P-a.s.).

*Necessity.* If  $\sum \xi_n$  converges (P-a.s.) then  $\xi_n \rightarrow 0$  (P-a.s.), and therefore, for every  $c > 0$ , at most a finite number of the events  $\{|\xi_n| \geq c\}$  can occur (P-a.s.). Therefore  $\sum I(|\xi_n| \geq c) < \infty$  (P-a.s.), and, by the second part of the Borel-Cantelli lemma,  $\sum P(|\xi_n| > c) < \infty$ . Moreover, the convergence of  $\sum \xi_n$  implies the convergence of  $\sum \xi_n^c$ . Therefore, by the two-series theorem, both of the series  $\sum E \xi_n^c$  and  $\sum V \xi_n^c$  converge.

This completes the proof of the theorem.

**Corollary.** *Let  $\xi_1, \xi_2, \dots$  be independent variables with  $E \xi_n = 0$ . Then if*

$$\sum E \frac{\xi_n^2}{1 + |\xi_n|} < \infty,$$

*the series  $\sum \xi_n$  converges with probability 1.*

For the proof we observe that

$$\sum E \frac{\xi_n^2}{1 + |\xi_n|} < \infty \Leftrightarrow \sum E[\xi_n^2 I(|\xi_n| \leq 1) + |\xi_n| I(|\xi_n| > 1)] < \infty.$$

Therefore if  $\xi_n^1 = \xi_n I(|\xi_n| \leq 1)$ , we have

$$\sum E(\xi_n^1)^2 < \infty.$$

Since  $E \xi_n = 0$ , we have

$$\begin{aligned} \sum |E \xi_n^1| &= \sum |E \xi_n I(|\xi_n| \leq 1)| = \sum |E \xi_n I(|\xi_n| > 1)| \\ &\leq \sum E |\xi_n| I(|\xi_n| > 1) < \infty. \end{aligned}$$

Therefore both  $\sum E \xi_n^1$  and  $\sum V \xi_n^1$  converge. Moreover, by Chebyshev's inequality,

$$P\{|\xi_n| > 1\} = P\{|\xi_n| I(|\xi_n| > 1) > 1\} \leq E(|\xi_n| I(|\xi_n| > 1)).$$

Therefore  $\sum P(|\xi_n| > 1) < \infty$ . Hence the convergence of  $\sum \xi_n$  follows from the three-series theorem.



## 4. PROBLEMS

1. Let  $\xi_1, \xi_2, \dots$  be a sequence of independent random variables,  $S_n = \xi_1 + \dots + \xi_n$ . Show, using the three-series theorem, that
- (a) if  $\sum \xi_n^2 < \infty$  (P-a.s.) then  $\sum \xi_n$  converges with probability 1, if and only if  $\sum E \xi_i I(|\xi_i| \leq 1)$  converges;
- (b) if  $\sum \xi_n$  converges (P-a.s.) then  $\sum \xi_n^2 < \infty$  (P-a.s.) if and only if
- $$\sum (E |\xi_n| I(|\xi_n| \leq 1))^2 < \infty.$$
2. Let  $\xi_1, \xi_2, \dots$  be a sequence of independent random variables. Show that  $\sum \xi_n^2 < \infty$  (P-a.s.) if and only if

$$\sum E \frac{\xi_n^2}{1 + \xi_n^2} < \infty.$$

3. Let  $\xi_1, \xi_2, \dots$  be a sequence of independent random variables. Show that  $\sum \xi_n$  converges (P-a.s.) if and only if it converges in probability.

## §3. Strong Law of Large Numbers

1. Let  $\xi_1, \xi_2, \dots$  be a sequence of independent random variables with finite second moments;  $S_n = \xi_1 + \dots + \xi_n$ . By Problem 2, §3, Chapter III, if the numbers  $V\xi_i$  are uniformly bounded, we have the law of large numbers:

$$\frac{S_n - ES_n}{n} \xrightarrow{P} 0, \quad n \rightarrow \infty. \quad (1)$$

A *strong law of large numbers* is a proposition in which convergence in probability is replaced by *convergence with probability 1*.

One of the earliest results in this direction is the following theorem.

**Theorem 1 (Cantelli).** Let  $\xi_1, \xi_2, \dots$  be independent random variables with finite fourth moments and let

$$E |\xi_n - E\xi_n|^4 \leq C, \quad n \geq 1,$$

for some constant  $C$ . Then as  $n \rightarrow \infty$

$$\frac{S_n - ES_n}{n} \rightarrow 0 \quad (\text{P-a.s.}). \quad (2)$$

**PROOF.** Without loss of generality, we may assume that  $E\xi_n = 0$  for  $n \geq 1$ . By the corollary to Theorem 1, §10, Chapter II, we will have  $S_n/n \rightarrow 0$  (P-a.s.) provided that

$$\sum P \left\{ \left| \frac{S_n}{n} \right| \geq \varepsilon \right\} < \infty$$

for every  $\varepsilon > 0$ . In turn, by Chebyshev's inequality, this will follow from

$$\sum E \left| \frac{S_n}{n} \right|^4 < \infty.$$

Let us show that this condition is actually satisfied under our hypotheses.

We have

$$\begin{aligned} S_n^4 &= (\xi_1 + \cdots + \xi_n)^4 = \sum_{i=1}^n \xi_i^4 - \sum_{\substack{i,j \\ i < j}} \frac{4!}{2!2!} \xi_i^2 \xi_j^2 \\ &\quad + \sum_{\substack{i \neq j \\ i \neq k \\ j < k}} \frac{4!}{2!1!1!} \xi_i^2 \xi_j \xi_k + \sum_{i < j < k < l} 4! \xi_i \xi_j \xi_k \xi_l \\ &\quad + \sum_{i \neq j} \frac{4!}{3!1!} \xi_i^3 \xi_j. \end{aligned}$$

Remembering that  $E\xi_k = 0$ ,  $k \leq n$ , we then obtain

$$\begin{aligned} ES_n^4 &= \sum_{i=1}^n E\xi_i^4 + 6 \sum_{i,j=1}^n E\xi_i^2 E\xi_j^2 \leq nC + 6 \sum_{\substack{i,j=1 \\ i < j}} \sqrt{E\xi_i^4 \cdot E\xi_j^4} \\ &\leq nC + \frac{6n(n-1)}{2} C = (3n^2 - 2n)C < 3n^2 C. \end{aligned}$$

Consequently

$$\sum E\left(\frac{S_n}{n}\right)^4 \leq 3C \sum \frac{1}{n^2} < \infty.$$

This completes the proof of the theorem.

**2.** The hypotheses of Theorem 1 can be considerably weakened by the use of more precise methods. In this way we obtain a stronger law of large numbers.

**Theorem 2** (Kolmogorov). *Let  $\xi_1, \xi_2, \dots$  be a sequence of independent random variables with finite second moments, and let there be positive numbers  $b_n$  such that  $b_n \uparrow \infty$  and*

$$\sum \frac{V\xi_n}{b_n^2} < \infty. \quad (3)$$

Then

$$\frac{S_n - ES_n}{b_n} \rightarrow 0 \quad (\text{P-a.s.}). \quad (4)$$

In particular, if

$$\sum \frac{V\xi_n}{n^2} < \infty \quad (5)$$

then

$$\frac{S_n - ES_n}{n} \rightarrow 0 \quad (\text{P-a.s.}). \quad (6)$$

For the proof of this, and of Theorem 2 below, we need two lemmas.

**Lemma 1** (Toeplitz). Let  $\{a_n\}$  be a sequence of nonnegative numbers,  $b_n = \sum_{i=1}^n a_i$ ,  $b_n > 0$  for  $n \geq 1$ , and  $b_n \uparrow \infty$ ,  $n \rightarrow \infty$ . Let  $\{x_n\}$  be a sequence of numbers converging to  $x$ . Then

$$\frac{1}{b_n} \sum_{j=1}^n a_j x_j \rightarrow x. \quad (7)$$

In particular, if  $a_n = 1$  then

$$\frac{x_1 + \cdots + x_n}{n} \rightarrow x. \quad (8)$$

**PROOF.** Let  $\varepsilon > 0$  and let  $n_0 = n_0(\varepsilon)$  be such that  $|x_n - x| \leq \varepsilon/2$  for all  $n \geq n_0$ . Choose  $n_1 > n_0$  so that

$$\frac{1}{b_{n_1}} \sum_{j=1}^{n_0} |x_j - x| < \varepsilon/2.$$

Then, for  $n > n_1$ ,

$$\begin{aligned} \left| \frac{1}{b_n} \sum_{j=1}^n a_j x_j - x \right| &\leq \frac{1}{b_n} \sum_{j=1}^n a_j |x_j - x| \\ &= \frac{1}{b_n} \sum_{j=1}^{n_0} a_j |x_j - x| + \frac{1}{b_n} \sum_{j=n_0+1}^n a_j |x_j - x| \\ &\leq \frac{1}{b_{n_1}} \sum_{j=1}^{n_0} a_j |x_j - x| + \frac{1}{b_n} \sum_{j=n_0+1}^n a_j |x_j - x| \\ &\leq \frac{\varepsilon}{2} + \frac{b_n - b_{n_0}}{b_n} \frac{\varepsilon}{2} \leq \varepsilon. \end{aligned}$$

This completes the proof of the lemma.

**Lemma 2** (Kronecker). Let  $\{b_n\}$  be a sequence of positive increasing numbers,  $b_n \uparrow \infty$ ,  $n \rightarrow \infty$ , and let  $\{x_n\}$  be a sequence of numbers such that  $\sum x_n$  converges. Then

$$\frac{1}{b_n} \sum_{j=1}^n b_j x_j \rightarrow 0, \quad n \rightarrow \infty. \quad (9)$$

In particular, if  $b_n = n$ ,  $x_n = y_n/n$  and  $\sum (y_n/n)$  converges, then

$$\frac{y_1 + \cdots + y_n}{n} \rightarrow 0, \quad n \rightarrow \infty. \quad (10)$$

**PROOF.** Let  $b_0 = 0$ ,  $S_0 = 0$ ,  $S_n = \sum_{j=1}^n x_j$ . Then (by summation by parts)

$$\sum_{j=1}^n b_j x_j = \sum_{j=1}^n b_j (S_j - S_{j-1}) = b_n S_n - b_0 S_0 - \sum_{j=1}^n S_{j-1} (b_j - b_{j-1})$$

and therefore

$$\frac{1}{b_n} \sum_{j=1}^n b_j x_j = S_n - \frac{1}{b_n} \sum_{j=1}^n S_{j-1} a_j \rightarrow 0,$$

since, if  $S_n \rightarrow x$ , then by Toeplitz's lemma,

$$\frac{1}{b_n} \sum_{j=1}^n S_{j-1} a_j \rightarrow x.$$

This establishes the lemma.

**PROOF OF THEOREM 1.** Since

$$\frac{S_n - ES_n}{b_n} = \frac{1}{b_n} \sum_{k=1}^n b_k \left( \frac{\xi_k - E\xi_k}{b_k} \right),$$

a sufficient condition for (4) is, by Kronecker's lemma, that the series  $\sum [(\xi_k - E\xi_k)/b_k]$  converges (P-a.s.). But this series does converge by (3) of Theorem 1, §2.

This completes the proof of the theorem.

**EXAMPLE 1.** Let  $\xi_1, \xi_2, \dots$  be a sequence of independent Bernoulli random variables with  $P(\xi_n=1)=P(\xi_n=-1)=\frac{1}{2}$ . Then, since  $\sum [1/(n \log^2 n)] < \infty$ , we have

$$\frac{S_n}{\sqrt{n \log n}} \rightarrow 0 \quad (\text{P-a.s.}). \quad (11)$$

**3.** In the case when the variables  $\xi_1, \xi_2, \dots$  are not only independent but also identically distributed, we can obtain a strong law of large numbers without requiring (as in Theorem 2) the existence of the second moment, provided that the first absolute moment exists.

**Theorem 3 (Kolmogorov).** Let  $\xi_1, \xi_2, \dots$  be a sequence of independent identically distributed random variables with  $E|\xi_1| < \infty$ . Then

$$\frac{S_n}{n} \rightarrow m \quad (\text{P-a.s.}) \quad (12)$$

where  $m = E\xi_1$ .

We need the following lemma.

**Lemma 3.** Let  $\xi$  be a nonnegative random variable. Then

$$\sum_{n=1}^{\infty} P(\xi \geq n) \leq E\xi \leq 1 + \sum_{n=1}^{\infty} P(\xi \geq n). \quad (13)$$

The proof consists of the following chain of inequalities:

$$\begin{aligned}
 \sum_{n=1}^{\infty} P(\xi \geq n) &= \sum_{n=1}^{\infty} \sum_{k \geq n} P(k \leq \xi < k+1) \\
 &= \sum_{k=1}^{\infty} k P(k \leq \xi < k+1) = \sum_{k=0}^{\infty} E[k I(k \leq \xi < k+1)] \\
 &\leq \sum_{k=0}^{\infty} E[\xi I(k \leq \xi < k+1)] \\
 &= E\xi \leq \sum_{k=0}^{\infty} E[(k+1) I(k \leq \xi < k+1)] \\
 &= \sum_{k=0}^{\infty} (k+1) P(k \leq \xi < k+1) \\
 &= \sum_{n=1}^{\infty} P(\xi \geq n) + \sum_{k=0}^{\infty} P(k \leq \xi < k+1) = \sum_{n=1}^{\infty} P(\xi \geq n) + 1.
 \end{aligned}$$

PROOF OF THEOREM 3. By Lemma 3 and the Borel-Cantelli lemma,

$$\begin{aligned}
 E|\xi_1| < \infty &\Leftrightarrow \sum P\{|\xi_1| \geq n\} < \infty \\
 &\Leftrightarrow \sum P\{|\xi_n| \geq n\} < \infty \Leftrightarrow P\{|\xi_n| \geq n \text{ i.o.}\} = 0.
 \end{aligned}$$

Hence  $|\xi_n| < n$ , except for a finite number of  $n$ , with probability 1.

Let us put

$$\tilde{\xi}_n = \begin{cases} \xi_n, & |\xi_n| < n, \\ 0, & |\xi_n| \geq n, \end{cases}$$

and suppose that  $E\xi_n = 0$ ,  $n \geq 1$ . Then  $(\xi_1 + \dots + \xi_n)/n \rightarrow 0$  (P-a.s.), if and only if  $(\tilde{\xi}_1 + \dots + \tilde{\xi}_n)/n \rightarrow 0$  (P-a.s.). Note that in general  $E\tilde{\xi}_n \neq 0$  but

$$E\tilde{\xi}_n = E\xi_n I(|\xi_n| < n) = E\xi_1 I(|\xi_1| < n) \rightarrow E\xi_1 = 0.$$

Hence by Toeplitz's lemma

$$\frac{1}{n} \sum_{k=1}^n E\tilde{\xi}_k \rightarrow 0, \quad n \rightarrow \infty, \quad \text{ } \quad \text{ }$$

and consequently  $(\xi_1 + \dots + \xi_n)/n \rightarrow 0$  (P-a.s.), if and only if

$$\frac{(\tilde{\xi}_1 - E\tilde{\xi}_1) + \dots + (\tilde{\xi}_n - E\tilde{\xi}_n)}{n} \rightarrow 0, \quad n \rightarrow \infty \quad (\text{P-a.s.}), \quad n \rightarrow \infty. \quad (14)$$

Write  $\bar{\xi}_n = \tilde{\xi}_n - E\tilde{\xi}_n$ . By Kronecker's lemma, (14) will be established if  $\sum (\bar{\xi}_n/n)$  converges (P-a.s.). In turn, by Theorem 1 of §2, this will follow if we show that, when  $E|\xi_1| < \infty$ , the series  $\sum (V\bar{\xi}_n/n^2)$  converges.

We have

$$\begin{aligned}
 \sum \frac{V_{\xi_n}^E}{n^2} &\leq \sum_{n=1}^{\infty} \frac{E \xi_n^2}{n^2} = \sum_{n=1}^{\infty} \frac{1}{n^2} E[\xi_n I(|\xi_n| < n)]^2 \\
 &= \sum_{n=1}^{\infty} \frac{1}{n^2} E[\xi_1^2 I(|\xi_1| < n)] = \sum_{n=1}^{\infty} \frac{1}{n^2} \sum_{k=1}^n E[\xi_1^2 I(k-1 \leq |\xi_1| < k)] \\
 &= \sum_{k=1}^{\infty} E[\xi_1^2 I(k-1 \leq |\xi_1| < k)] \cdot \sum_{n=k}^{\infty} \frac{1}{n^2} \\
 &\leq 2 \sum_{k=1}^{\infty} \frac{1}{k} E[\xi_1^2 I(k-1 \leq |\xi_1| < k)] \\
 &\leq 2 \sum_{k=1}^{\infty} E[|\xi_1| I(k-1 \leq |\xi_1| < k)] = 2E|\xi_1| < \infty.
 \end{aligned}$$

This completes the proof of the theorem.

**Remark 1.** The theorem admits a converse in the following sense. Let  $\xi_1, \xi_2, \dots$  be a sequence of independent identically distributed random variables such that

$$\frac{\xi_1 + \dots + \xi_n}{n} \rightarrow C,$$

with probability 1, where  $C$  is a (finite) constant. Then  $E|\xi_1| < \infty$  and  $C = E\xi_1$ .

In fact, if  $S_n/n \rightarrow C$  (P-a.s.) then

$$\frac{\xi_n}{n} = \frac{S_n}{n} - \left(\frac{n-1}{n}\right) \frac{S_{n-1}}{n-1} \rightarrow 0 \quad (\text{P-a.s.})$$

and therefore  $P(|\xi_n| > n \text{ i.o.}) = 0$ . By the Borel–Cantelli lemma,

$$\sum P(|\xi_1| > n) < \infty,$$

and by Lemma 3 we have  $E|\xi_1| < \infty$ . Then it follows from the theorem that  $C = E\xi_1$ .

Consequently for independent identically distributed random variables the condition  $E|\xi_1| < \infty$  is necessary and sufficient for the convergence (with probability 1) of the ratio  $S_n/n$  to a finite limit.

**Remark 2.** If the expectation  $m = E\xi_1$  exists but is not necessarily finite, the conclusion (10) of the theorem remains valid.

In fact, let, for example,  $E\xi_1^- < \infty$  and  $E\xi_1^+ = \infty$ . With  $C > 0$ , put

$$S_n^C = \sum_{i=1}^n \xi_i I(\xi_i \leq C).$$

Then (P-a.s.).

$$\lim_n \frac{S_n}{n} \geq \lim_n \frac{S_n^C}{n} = E\xi_1 I(\xi_1 \leq C).$$

But as  $C \rightarrow \infty$ ,

$$E\xi_1 I(\xi_1 \leq C) \rightarrow E\xi_1 = \infty;$$

therefore  $S_n/n \rightarrow +\infty$  (P-a.s.).

#### 4. Let us give some applications of the strong law of large numbers.

**EXAMPLE 1** (Application to number theory). Let  $\Omega = [0, 1)$ , let  $\mathcal{B}$  be the algebra of Borel subsets of  $\Omega$  and let  $P$  be Lebesgue measure on  $[0, 1)$ . Consider the binary expansions  $\omega = 0.\omega_1\omega_2\dots$  of numbers  $\omega \in \Omega$  (with infinitely many 0's) and define random variables  $\xi_1(\omega), \xi_2(\omega), \dots$  by putting  $\xi_n(\omega) = \omega_n$ . Since, for all  $n \geq 1$  and all  $x_1, \dots, x_n$  taking the values 0 or 1,

$$\begin{aligned} \{\omega: \xi_1(\omega) = x_1, \dots, \xi_n(\omega) = x_n\} \\ = \left\{ \omega: \frac{x_1}{2} + \frac{x_2}{2^2} + \dots + \frac{x_n}{2^n} \leq \omega < \frac{x_1}{2} + \dots + \frac{x_n}{2^n} + \frac{1}{2^n} \right\}, \end{aligned}$$

the  $P$ -measure of this set is  $1/2^n$ . It follows that  $\xi_1, \xi_2, \dots$  is a sequence of independent identically distributed random variables with

$$P(\xi_1 = 0) = P(\xi_1 = 1) = \frac{1}{2}.$$

Hence, by the strong law of large numbers, we have the following result of Borel: *almost every number in  $[0, 1)$  is normal, in the sense that with probability 1 the proportion of zeros and ones in its binary expansion tends to  $\frac{1}{2}$ , i.e.,*

$$\frac{1}{n} \sum_{k=1}^n I(\xi_k = 1) \rightarrow \frac{1}{2} \quad (\text{P-a.s.}).$$

**EXAMPLE 2** (The Monte Carlo method). Let  $f(x)$  be a continuous function defined on  $[0, 1]$ , with values on  $[0, 1]$ . The following idea is the foundation of the statistical method of calculating  $\int_0^1 f(x) dx$  (the "Monte Carlo method").

Let  $\xi_1, \eta_1, \xi_2, \eta_2, \dots$  be a sequence of independent random variables, uniformly distributed on  $[0, 1]$ . Put

$$\rho_i = \begin{cases} 1 & \text{if } f(\xi_i) > \eta_i, \\ 0 & \text{if } f(\xi_i) < \eta_i. \end{cases}$$

It is clear that

$$E\rho_1 = P\{f(\xi_1) > \eta_1\} = \int_0^1 f(x) dx.$$

By the strong law of large numbers (Theorem 3)

$$\frac{1}{n} \sum_{i=1}^n \rho_i \rightarrow \int_0^1 f(x) dx \quad (\text{P-a.s.}).$$

Consequently we can approximate an integral  $\int_0^1 f(x) dx$  by taking a simulation consisting of a pair of random variables  $(\xi_i, \eta_i)$ ,  $i \geq 1$ , and then calculating  $\rho_i$  and  $(1/n) \sum_{i=1}^n \rho_i$ .

## 5. PROBLEMS

1. Show that  $E\xi^2 < \infty$  if and only if  $\sum_{n=1}^{\infty} nP(|\xi| > n) < \infty$ .
2. Supposing that  $\xi_1, \xi_2, \dots$  are independent and identically distributed, show that if  $E|\xi_1|^\alpha < \infty$  for some  $\alpha$ ,  $0 < \alpha < 1$ , then  $S_n/n^{1/\alpha} \rightarrow 0$  (P-a.s.), and if  $E|\xi_1|^\beta < \infty$  for some  $\beta$ ,  $1 \leq \beta < 2$ , then  $(S_n - nE\xi_1)/n^{1/\beta} \rightarrow 0$  (P-a.s.).
3. Let  $\xi_1, \xi_2, \dots$  be a sequence of independent identically distributed random variables and let  $E|\xi_1| = \infty$ . Show that

$$\overline{\lim} \left| \frac{S_n}{n} - a_n \right| = \infty \quad (\text{P-a.s.})$$

for every sequence of constants  $\{a_n\}$ .

4. Show that a rational number on  $[0, 1)$  is never normal (in the sense of Example 1, Subsection 4).

## §4. Law of the Iterated Logarithm

1. Let  $\xi_1, \xi_2, \dots$  be a sequence of independent Bernoulli random variables with  $P(\xi_n = 1) = P(\xi_n = -1) = \frac{1}{2}$ ; let  $S_n = \xi_1 + \dots + \xi_n$ . It follows from the proof of Theorem 2, §1, that

$$\overline{\lim} \frac{S_n}{\sqrt{n}} = +\infty, \quad \underline{\lim} \frac{S_n}{\sqrt{n}} = -\infty, \quad (1)$$

with probability 1. On the other hand, by (3.11),

$$\frac{S_n}{\sqrt{n \log n}} \rightarrow 0 \quad (\text{P-a.s.}). \quad (2)$$

Let us compare these results.

It follows from (1) that with probability 1 the paths of  $(S_n)_{n \geq 1}$  intersect the “curves”  $\pm \varepsilon \sqrt{n}$  infinitely often for any given  $\varepsilon$ ; but at the same time (2)



shows that they only finitely often leave the region bounded by the curves  $\pm \varepsilon \sqrt{n} \log n$ . These two results yield useful information on the amplitude of the oscillations of the symmetric random walk  $(S_n)_{n \geq 1}$ . The law of the iterated logarithm, which we present below, improves this picture of the amplitude of the oscillations of  $(S_n)_{n \geq 1}$ .

Let us introduce the following definition. We call a function  $\varphi^* = \varphi^*(n)$ ,  $n \geq 1$ , *upper* (for  $(S_n)_{n \geq 1}$ ) if, with probability 1,  $S_n \leq \varphi^*(n)$  for all  $n$  from  $n = n_0(\omega)$  on.

We call a function  $\varphi_* = \varphi_*(n)$ ,  $n \geq 1$ , *lower* (for  $(S_n)_{n \geq 1}$ ) if, with probability 1,  $S_n > \varphi_*(n)$  for infinitely many  $n$ .

Using these definitions, and appealing to (1) and (2), we can say that every function  $\varphi^* = \varepsilon \sqrt{n} \log n$ ,  $\varepsilon > 0$ , is upper, whereas  $\varphi_* = \varepsilon \sqrt{n}$  is lower,  $\varepsilon > 0$ .

Let  $\varphi = \varphi(n)$  be a function and  $\varphi_\varepsilon^* = (1 + \varepsilon)\varphi$ ,  $\varphi_{*\varepsilon} = (1 - \varepsilon)\varphi$ , where  $\varepsilon > 0$ . Then it is easily seen that

$$\begin{aligned} \left\{ \overline{\lim} \frac{S_n}{\varphi(n)} \leq 1 \right\} &= \left\{ \lim_n \left[ \sup_{m \geq n} \frac{S_m}{\varphi(m)} \right] \leq 1 \right\} \\ &\Leftrightarrow \left\{ \sup_{m \geq n_1(\varepsilon)} \frac{S_m}{\varphi(m)} \leq 1 + \varepsilon \text{ for every } \varepsilon > 0, \text{ from some } n_1(\varepsilon) \text{ on} \right\} \\ &\Leftrightarrow \{S_m \leq (1 + \varepsilon)\varphi(m) \text{ for every } \varepsilon > 0, \text{ from some } n_1(\varepsilon) \text{ on}\}. \end{aligned} \quad (3)$$

In the same way,

$$\begin{aligned} \left\{ \overline{\lim} \frac{S_n}{\varphi(n)} \geq 1 \right\} &= \left\{ \lim_n \left[ \sup_{m \geq n} \frac{S_m}{\varphi(m)} \right] \geq 1 \right\} \\ &\Leftrightarrow \left\{ \sup_{m \geq n_2(\varepsilon)} \frac{S_m}{\varphi(m)} \leq 1 + \varepsilon \text{ for every } \varepsilon > 0, \text{ from some } n_1(\varepsilon) \text{ on} \right\} \\ &\Leftrightarrow \{S_m \geq (1 - \varepsilon)\varphi(m) \text{ for every } \varepsilon > 0 \text{ and for infinitely many } m \text{ larger than some } n_3(\varepsilon) \geq n_2(\varepsilon)\}. \end{aligned} \quad (4)$$

It follows from (3) and (4) that in order to verify that each function  $\varphi_\varepsilon^* = (1 + \varepsilon)\varphi$ ,  $\varepsilon > 0$ , is upper, we have to show that

$$\mathbf{P} \left\{ \overline{\lim} \frac{S_n}{\varphi(n)} \leq 1 \right\} = 1. \quad (5)$$

But to show that  $\varphi_{*\varepsilon} = (1 - \varepsilon)\varphi$ ,  $\varepsilon > 0$ , is lower, we have to show that

$$\mathbf{P} \left\{ \overline{\lim} \frac{S_n}{\varphi(n)} \geq 1 \right\} = 1. \quad (6)$$

**2. Theorem 1** (Law of the Iterated Logarithm). *Let  $\xi_1, \xi_2, \dots$  be a sequence of independent identically distributed random variables with  $E\xi_i = 0$  and  $E\xi_i^2 = \sigma^2 > 0$ . Then*

$$P\left\{\overline{\lim} \frac{S_n}{\psi(n)} = 1\right\} = 1, \quad (7)$$

where

$$\psi(n) = \sqrt{2\sigma^2 n \log \log n}. \quad (8)$$

For uniformly bounded random variables, the law of the iterated logarithm was established by Khinchin (1924). In 1929 Kolmogorov generalized this result to a wide class of independent variables. Under the conditions of Theorem 1, the law of the iterated logarithm was established by Hartman and Wintner (1941).

Since the proof of Theorem 1 is rather complicated, we shall confine ourselves to the special case when the random variables  $\xi_n$  are normal,  $\xi_n \sim \mathcal{N}(0, 1)$ ,  $n \geq 1$ .

We begin by proving two auxiliary results.

**Lemma 1.** *Let  $\xi_1, \dots, \xi_n$  be independent random variables that are symmetrically distributed ( $P(\xi_k \in B) = P(-\xi_k \in B)$  for every  $B \in \mathcal{B}(R)$ ,  $k \leq n$ ). Then for every real number  $a$*

$$P\left(\max_{1 \leq k \leq n} S_k > a\right) \leq 2P(S_n > a). \quad (9)$$

**PROOF.** Let  $A = \{\max_{1 \leq k \leq n} S_k > a\}$ ,  $A_k = \{S_i \leq a, i \leq k-1; S_k > a\}$  and  $B = \{S_n > a\}$ . Since  $S_n > a$  on  $A_k$  (because  $S_k \leq S_n$ ), we have

$$\begin{aligned} P(B \cap A_k) &\geq P(A_k \cap \{S_n \geq S_k\}) = P(A_k)P(S_n \geq S_k) \\ &= P(A_k)P(\xi_{k+1} + \dots + \xi_n \geq 0). \end{aligned}$$

By the symmetry of the distributions of the random variables  $\xi_1, \dots, \xi_n$ , we have

$$P(\xi_{k+1} + \dots + \xi_n > 0) = P(\xi_{k+1} + \dots + \xi_n < 0).$$

Hence  $P(\xi_{k+1} + \dots + \xi_n > 0) \geq \frac{1}{2}$ , and therefore

$$P(B) \geq \sum_{k=1}^n P(A_k \cap B) \geq \frac{1}{2} \sum_{k=1}^n P(A_k) = \frac{1}{2} P(A),$$

which establishes (9).

**Lemma 2.** *Let  $S_n \sim \mathcal{N}(0, \sigma^2(n))$ ,  $\sigma^2(n) \uparrow \infty$ , and let  $a(n)$ ,  $n \geq 1$ , satisfy  $a(n)/\sigma(n) \rightarrow \infty$ ,  $n \rightarrow \infty$ . Then*

$$P(S_n > a(n)) \sim \frac{\sigma(n)}{\sqrt{2\pi}a(n)} \exp\{-\frac{1}{2}a^2(n)/\sigma^2(n)\}. \quad (10)$$

The proof follows from the asymptotic formula

$$\frac{1}{\sqrt{2\pi}} \int_x^\infty e^{-y^2/2} dy \sim \frac{1}{\sqrt{2\pi x}} e^{-x^2/2}, \quad x \rightarrow \infty,$$

since  $S_n/\sigma(n) \sim \mathcal{N}(0, 1)$ .

PROOF OF THEOREM 1 (for  $\xi_i \sim \mathcal{N}(0, 1)$ ).

Let us first establish (5). Let  $\varepsilon > 0$ ,  $\lambda = 1 + \varepsilon$ ,  $n_k = \lambda^k$ , where  $k \geq k_0$ , and  $k_0$  is chosen so that  $\ln \ln k_0$  is defined. We also define

$$A_k = \{S_n > \lambda\psi(n) \text{ for some } n \in (n_k, n_{k+1}]\}, \quad (11)$$

and put

$$A = \{A_k \text{ i.o.}\} = \{S_n > \lambda\psi(n) \text{ for infinitely many } n\}.$$

In accordance with (3), we can establish (5) by showing that  $P(A) = 0$ .

Let us show that  $\sum P(A_k) < \infty$ . Then  $P(A) = 0$  by the Borel-Cantelli lemma.

From (11), (9), and (10) we find that

$$\begin{aligned} P(A_k) &\leq P\{S_n > \lambda\psi(n_k) \text{ for some } n \in (n_k, n_{k+1}]\} \\ &\leq P\{S_n > \lambda\psi(n_k) \text{ for some } n \leq n_{k+1}\} \\ &\leq 2P\{S_{n_{k+1}} > \lambda\psi(n_k)\} \sim \frac{2\sqrt{n_k}}{\sqrt{2\pi}\lambda\psi(n_k)} \exp\{-\tfrac{1}{2}\lambda^2[\psi(n_k)/\sqrt{n_k}]^2\} \\ &\leq C_1 \exp(-\lambda \ln \ln \lambda^k) \leq C e^{-\lambda \ln k} = C_2 k^{-\lambda}, \end{aligned}$$

where  $C_1$  and  $C_2$  are constants. But  $\sum_{k=1}^\infty k^{-\lambda} < \infty$ , and therefore

$$\sum P(A_k) < \infty.$$

Consequently (5) is established.

We turn now to the proof of (6). In accordance with (4) we must show that, with  $\lambda = 1 - \varepsilon$ ,  $\varepsilon > 0$ , we have with probability 1 that  $S_n \geq \lambda\psi(n)$  for infinitely many  $n$ .

Let us apply (5), which we just proved, to the sequence  $(-S_n)_{n \geq 1}$ . Then we find that for all  $n$ , with finitely many exceptions,  $-S_n \leq 2\psi(n)$  (P-a.s.). Consequently if  $n_k = N^k$ ,  $N > 1$ , then for sufficiently large  $k$ , either

$$S_{n_{k-1}} \geq -2\psi(n_{k-1})$$

or

$$S_{n_k} \geq Y_k - 2\psi(n_{k-1}), \quad (12)$$

where  $Y_k = S_{n_k} - S_{n_{k-1}}$ .

Hence if we show that for infinitely many  $k$

$$Y_k > \lambda\psi(n_k) + 2\psi(n_{k-1}), \quad (13)$$

this and (12) show that (P-a.s.)  $S_{n_k} > \lambda\psi(n_k)$  for infinitely many  $k$ . Take some  $\lambda' \in (\lambda, 1)$ . Then there is an  $N > 1$  such that for all  $k$

$$\begin{aligned} \lambda'[2(N^k - N^{k-1}) \ln \ln N^k]^{1/2} &> \lambda(2N^k \ln \ln N^k)^{1/2} \\ &+ 2(2N^{k-1} \ln \ln N^{k-1})^{1/2} \equiv \lambda\psi(N^k) + 2\psi(N^{k-1}). \end{aligned}$$

It is now enough to show that

$$Y_k > \lambda'[2(N^k - N^{k-1}) \ln \ln N^k]^{1/2} \quad (14)$$

for infinitely many  $k$ . Evidently  $Y_k \sim \mathcal{N}(0, N^k - N^{k-1})$ . Therefore, by Lemma 2,

$$\begin{aligned} P\{Y_k > \lambda'[2(N^k - N^{k-1}) \ln \ln N^k]^{1/2}\} &\sim \frac{1}{\sqrt{2\pi\lambda'}(2 \ln \ln N^k)^{1/2}} e^{-(\lambda')^2 \ln \ln N^k} \\ &\geq \frac{C_1}{(\ln k)^{1/2}} k^{-(\lambda')^2} \geq \frac{C_2}{k \ln k}. \end{aligned}$$

Since  $\sum (1/k \ln k) = \infty$ , it follows from the second part of the Borel–Cantelli lemma that, with probability 1, inequality (14) is satisfied for infinitely many  $k$ , so that (6) is established.

This completes the proof of the theorem.

**Remark 1.** Applying (7) to the random variables  $(-S_n)_{n \geq 1}$ , we find that

$$\varliminf \frac{S_n}{\varphi(n)} = -1. \quad (15)$$

It follows from (7) and (15) that the law of the iterated logarithm can be put in the form

$$P\left\{\varliminf \frac{|S_n|}{\varphi(n)} = 1\right\} = 1. \quad (16)$$

**Remark 2.** The law of the iterated logarithm says that for every  $\varepsilon > 0$  each function  $\psi_\varepsilon^* = (1 + \varepsilon)\psi$  is upper, and  $\psi_{*\varepsilon} = (1 - \varepsilon)\psi$  is lower.

The conclusion (7) is also equivalent to the statement that, for each  $\varepsilon > 0$ ,

$$P\{|S_n| \geq (1 - \varepsilon)\psi(n) \text{ i.o.}\} = 1,$$

$$P\{|S_n| \geq (1 + \varepsilon)\psi(n) \text{ i.o.}\} = 0.$$

### 3. PROBLEMS

1. Let  $\xi_1, \xi_2, \dots$  be a sequence of independent random variables with  $\xi_n \sim \mathcal{N}(0, 1)$ . Show that

$$P\left\{\varliminf \frac{\xi_n}{\sqrt{2 \ln n}} = 1\right\} = 1.$$

2. Let  $\xi_1, \xi_2, \dots$  be a sequence of independent random variables, distributed according to Poisson's law with parameter  $\lambda > 0$ . Show that (independently of  $\lambda$ )

$$P\left\{\lim_{n \rightarrow \infty} \frac{\xi_n \ln \ln n}{\ln n} = 1\right\} = 1.$$

3. Let  $\xi_1, \xi_2, \dots$  be a sequence of independent identically distributed random variables with

$$E e^{t\xi_1} = e^{-|t|^\alpha}, \quad 0 < \alpha < 2.$$

Show that

$$P\left\{\lim_{n \rightarrow \infty} \left| \frac{S_n}{n^{1/\alpha}} \right|^{1/(1 \ln \ln n)} = e^{1/\alpha}\right\} = 1.$$

4. Establish the following generalization of (9). Let  $\xi_1, \dots, \xi_n$  be independent random variables. *Lévy's inequality*

$$P\left\{\max_{0 \leq k \leq n} [S_k + \mu(S_n - S_k)] > a\right\} \leq 2P(S_n > a), \quad S_0 = 0,$$

holds for every real  $a$ , where  $\mu(\xi)$  is the median of  $\xi$ , i.e. a constant such that

$$P(\xi \geq \mu(\xi)) \geq \frac{1}{2}, \quad P(\xi \leq \mu(\xi)) \geq \frac{1}{2}.$$

## §5. Rapidity of Convergence in the Strong Law of Large Numbers and in the Probabilities of Large Deviations

1. By the results of §6, Chapter I, we have the following estimate for the Bernoulli scheme:

$$P\left\{\left|\frac{S_n}{n} - p\right| \geq \varepsilon\right\} \leq 2e^{-2n\varepsilon^2} \quad (1)$$

(see (42), subsection 7, §6, Chapter I). From this, of course, there follows the inequalities

$$P\left\{\sup_{m \geq n} \left|\frac{S_m}{m} - p\right| \geq \varepsilon\right\} \leq \sum_{m \geq n} P\left\{\left|\frac{S_m}{m} - p\right| \geq \varepsilon\right\} \leq \frac{2}{1 - e^{-2\varepsilon^2}} e^{-2n\varepsilon^2}, \quad (2)$$

which provide an approximation of the rate of convergence to  $p$  by the quantity  $S_n/n$  with probability 1.

We now consider the question of the validity of formulas of the types (1) and (2) in some more general situations, when  $S_n = \xi_1 + \dots + \xi_n$  is a sum of independent identically distributed random variables.

2. Let  $\xi, \xi_1, \xi_2, \dots$  be a sequence of independent random variables. We say that a random variable satisfies *Cramér's condition* if there is a  $\lambda > 0$  for which

$$\varphi(\lambda) = Ee^{\lambda|\xi|} < \infty \quad (3)$$

(it can be shown that this condition is equivalent to an exponential decrease of  $P(|\xi| > x)$ ,  $x \rightarrow \infty$ ).

Let

$$\Lambda = \{\lambda \in R: \varphi(\lambda) < \infty\}, \quad (4)$$

where we suppose that  $\Lambda$  contains a neighborhood of the point  $\lambda = 0$ , i.e., Cramér's condition (3) is satisfied for some  $\lambda > 0$ .

On the set  $\Lambda$  the function

$$\psi(\lambda) = \ln \varphi(\lambda) \quad (5)$$

is convex (from below), and strictly convex if the random variable  $\xi$  is not degenerate. We also notice that

$$\psi(0) = 0, \quad \psi'(0) = m (= E\xi), \quad \psi''(\lambda) \geq 0.$$

We extend  $\psi(\lambda)$  to all  $\lambda \in R$  by setting  $\psi(\lambda) = \infty$  for  $\lambda \notin \Lambda$ .

We define the function

$$H(a) = \sup_{\lambda} [a\lambda - \psi(\lambda)], \quad a \in R, \quad (6)$$

called the *Cramér transform* (of the distribution function  $F = F(x)$  of the random variable  $\xi$ ). The function  $H(a)$  is also convex (from below) and its minimum is zero, attained at  $\lambda = m$ .

If  $a > m$ , we have

$$H(a) = \sup_{\lambda > 0} [a\lambda - \psi(\lambda)].$$

Then

$$P\{\xi \geq a\} \leq \inf_{\lambda > 0} Ee^{\lambda(\xi - a)} = \inf_{\lambda > 0} e^{-[a\lambda - \psi(\lambda)]} = e^{-H(a)}. \quad (7)$$

Similarly, for  $a < m$  we have  $H(a) = \sup_{\lambda < 0} [a\lambda - \psi(\lambda)]$  and

$$P\{\xi \leq a\} \leq e^{-H(a)}. \quad (8)$$

Consequently (compare (42), §6, Chapter 1),

$$P\{|\xi - m| \geq \varepsilon\} \leq e^{-\min\{H(m-\varepsilon), H(m+\varepsilon)\}}. \quad (9)$$

If  $\xi, \xi_1, \dots, \xi_n$  are independent identically distributed random variables that satisfy Cramér's condition (3),  $S_n = \xi_1 + \dots + \xi_n$ ,  $\psi_n(\lambda) = \ln E \exp(\lambda S_n/n)$ ,  $\psi(\lambda) = \ln E e^{\lambda \xi}$ , and

$$H_n(a) = \sup_{\lambda} [a\lambda - \psi_n(\lambda)], \quad (10)$$

then

$$H_n(a) = nH(a) \left( = n \sup_{\lambda} [a\lambda - \psi(\lambda)] \right)$$

and the inequalities (7), (8), and (9) assume the following forms:

$$P\left\{\frac{S_n}{n} \geq a\right\} \leq e^{-nH(a)}, \quad a > m, \quad (11)$$

$$P\left\{\frac{S_n}{n} \leq a\right\} \leq e^{-nH(a)}, \quad a < m, \quad (12)$$

$$P\left\{\left|\frac{S_n}{n} - m\right| \geq \varepsilon\right\} \leq 2e^{-\min\{H(m-\varepsilon), H(m+\varepsilon)\} \cdot n}. \quad (13)$$

**Remark.** Results of the type

$$P\left\{\left|\frac{S_n}{n} - m\right| \geq \varepsilon\right\} \leq ae^{-bn}, \quad (14)$$

where  $a > 0$  and  $b > 0$ , indicate exponential convergence "adjusted" by the constants  $a$  and  $b$ . In the theory of *large deviations*, results are often presented in a somewhat different, "cruder," form:

$$\lim_n \frac{1}{n} \ln P\left\{\left|\frac{S_n}{n} - m\right| \geq \varepsilon\right\} < 0, \quad (15)$$

that clearly arises from (14) and refers to the "exponential" rapidity of convergence, but without specifying the values of the constants  $a$  and  $b$ .

Now we turn to the question of upper bounds for the probabilities

$$P\left\{\sup_{k \geq n} \frac{S_k}{k} > a\right\}, \quad P\left\{\inf_{k \geq n} \frac{S_k}{k} < a\right\}, \quad P\left\{\sup_{k \geq n} \left|\frac{S_k}{k} - m\right| > \varepsilon\right\},$$

which can provide definite bounds on the rapidity of convergence in the strong law of large numbers.

Let us suppose that the independent identically distributed nondegenerate random variables  $\xi, \xi_1, \xi_2, \dots$  satisfy Cramér's condition, i.e.,  $\varphi(\lambda) = Ee^{\lambda|\xi|} < \infty$  for some  $\lambda > 0$ .

We fix  $n \geq 1$  and set

$$\kappa = \inf\left\{k \geq n: \frac{S_k}{k} > a\right\},$$

taking  $\kappa = \infty$  if  $S_k/k < a$  for  $k \geq n$ .

In addition, let  $a$  and  $\lambda > 0$  satisfy

$$\lambda a - \ln \varphi(\lambda) \geq 0. \quad (16)$$

Then

$$\begin{aligned} P\left\{\sup_{k \geq n} \frac{S_k}{k} > a\right\} &= P\left\{\bigcup_{k \geq n} \left\{\frac{S_k}{k} > a\right\}\right\} \\ &= P\left\{\frac{S_\kappa}{\kappa} > a, \kappa < \infty\right\} = P\{e^{\lambda S_\kappa} > e^{\lambda a \kappa}, \kappa < \infty\} \end{aligned}$$

$$\begin{aligned}
&= P\{e^{\lambda S_k - \kappa \ln \varphi(\lambda)} > e^{\kappa(\lambda a - \ln \varphi(\lambda))}, \kappa < \infty\} \\
&\leq P\{e^{\lambda S_k - \kappa \ln \varphi(\lambda)} > e^{n(\lambda a - \ln \varphi(\lambda))}, \kappa < \infty\} \\
&\leq P\left\{\sup_{k \geq n} e^{\lambda S_k - k \ln \varphi(\lambda)} \geq e^{n(\lambda a - \ln \varphi(\lambda))}\right\}. \quad (17)
\end{aligned}$$

To take the final step, we notice that the sequence of random variables

$$e^{\lambda S_k - k \ln \varphi(\lambda)}, \quad k \geq 1,$$

with respect to the flow of  $\sigma$ -algebras  $\mathcal{F}_k = \sigma\{\xi_1, \dots, \xi_k\}$ ,  $k \geq 1$ , forms a *martingale*. (For more details, see Chapter VII and, in particular, Example 1 in §1). Then it follows from inequality (8) in §3, Chapter VII, that

$$P\left\{\sup_{k \geq n} e^{\lambda S_k - k \ln \varphi(\lambda)} \geq e^{n(\lambda a - \ln \varphi(\lambda))}\right\} \leq e^{-n(\lambda a - \ln \varphi(\lambda))},$$

and consequently, (by (16)) we obtain the inequality

$$P\left\{\sup_{k \geq n} \frac{S_k}{k} > a\right\} \leq e^{-n(\lambda a - \ln \varphi(\lambda))}. \quad (18)$$

Let  $a > m$ . Since the function  $f(\lambda) = a\lambda - \ln \varphi(\lambda)$  has the properties  $f(0) = 0$ ,  $f'(0) > 0$ , there is a  $\lambda > 0$  for which (16) is satisfied, and consequently, we obtain from (18) that if  $a > m$  we have

$$P\left\{\sup_{k \geq n} \frac{S_k}{k} > a\right\} \leq e^{-n \sup_{\lambda > 0} [\lambda a - \ln \varphi(\lambda)]} = e^{-nH(a)}. \quad (19)$$

Similarly, if  $a < m$ , we have

$$P\left\{\sup_{k \geq n} \frac{S_k}{k} < a\right\} \leq e^{-n \sup_{\lambda < 0} [\lambda a - \ln \varphi(\lambda)]} = e^{-nH(a)}. \quad (20)$$

From (19) and (20), we obtain

$$P\left\{\sup_{k \geq n} \left|\frac{S_k}{k} - m\right| > \varepsilon\right\} \leq 2e^{-\min[H(m-\varepsilon), H(m+\varepsilon)] \cdot n}. \quad (21)$$

**Remark.** Combining the right-hand sides of the inequalities (11) and (19) leads us to suspect that this situation is not random. In fact, this expectation is concealed in the fact that the sequences  $(S_k/k)_{n \leq k \leq N}$  form, for every  $n \leq N$ , *reversed martingales* (see Problem 6 in §1, Chapter VII, and Example 4 in §11, Chapter I).

## 2. PROBLEMS

1. Carry out the proof of inequalities (8) and (20).
2. Investigate the properties of  $H(a)$ .



## CHAPTER V

# Stationary (Strict Sense) Random Sequences and Ergodic Theory

### §1. Stationary (Strict Sense) Random Sequences. Measure-Preserving Transformations

1. Let  $(\Omega, \mathcal{F}, P)$  be a probability space and  $\xi = (\xi_1, \xi_2, \dots)$  a sequence of random variables or, as we say, a *random sequence*. Let  $\theta_k \xi$  denote the sequence  $(\xi_{k+1}, \xi_{k+2}, \dots)$ .

**Definition 1.** A random sequence  $\xi$  is *stationary (in the strict sense)* if the probability distributions of  $\theta_k \xi$  and  $\xi$  are the same for every  $k \geq 1$ :

$$P((\xi_1, \xi_2, \dots) \in B) = P((\xi_{k+1}, \xi_{k+2}, \dots) \in B), \quad B \in \mathcal{B}(R^\infty).$$

The simplest example is a sequence  $\xi = (\xi_1, \xi_2, \dots)$  of independent identically distributed random variables. Starting from such a sequence, we can construct a broad class of stationary sequences  $\eta = (\eta_1, \eta_2, \dots)$  by choosing any Borel function  $g(x_1, \dots, x_n)$  and setting  $\eta_k = g(\xi_k, \xi_{k+1}, \dots, \xi_{k+1})$ .

If  $\xi = (\xi_1, \xi_2, \dots)$  is a sequence of independent identically distributed random variables with  $E|\xi_1| < \infty$  and  $E\xi_1 = m$ , the law of large numbers tells us that, with probability 1,

$$\frac{\xi_1 + \dots + \xi_n}{n} \rightarrow m, \quad n \rightarrow \infty.$$

In 1931 Birkhoff obtained a remarkable generalization of this fact for the case of stationary sequences. The present chapter consists mainly of a proof of Birkhoff's theorem.

The following presentation is based on the idea of measure-preserving transformations, something that brings us in contact with an interesting

branch of analysis (ergodic theory), and at the same time shows the connection between this theory and stationary random processes.

Let  $(\Omega, \mathcal{F}, P)$  be a probability space.

**Definition 2.** A transformation  $T$  of  $\Omega$  into  $\Omega$  is *measurable* if, for every  $A \in \mathcal{F}$ ,

$$T^{-1}A = \{\omega: T\omega \in A\} \in \mathcal{F}.$$

**Definition 3.** A measurable transformation  $T$  is a *measure-preserving transformation* (or *morphism*) if, for every  $A \in \mathcal{F}$ ,

$$P(T^{-1}A) = P(A).$$

Let  $T$  be a measure-preserving transformation,  $T^n$  its  $n$ th iterate, and  $\xi_1 = \xi_1(\omega)$  a random variable. Put  $\xi_k(\omega) = \xi_1(T^{k-1}\omega)$ ,  $k \geq 2$ , and consider the sequence  $\xi = (\xi_1, \xi_2, \dots)$ . We claim that this sequence is stationary.

In fact, let  $A = \{\omega: \xi \in B\}$  and  $A_1 = \{\omega: \theta_1 \xi \in B\}$ , where  $B \in \mathcal{B}(R^\infty)$ . Since  $A = \{\omega: (\xi_1(\omega), \xi_1(T\omega), \dots) \in B\}$ , and  $A_1 = \{\omega: (\xi_1(T\omega), \xi_1(T^2\omega), \dots) \in B\}$ , we have  $\omega \in A_1$  if and only if either  $T\omega \in A$  or  $A_1 = T^{-1}A$ . But  $P(T^{-1}A) = P(A)$ , and therefore  $P(A_1) = P(A)$ . Similarly  $P(A_k) = P(A)$  for every  $A_k = \{\omega: \theta_k \xi \in B\}$ ,  $k \geq 2$ .

Thus we can use measure-preserving transformations to construct stationary (strict sense) random variables.

In a certain sense, there is a converse result: for every stationary sequence  $\xi$  considered on  $(\Omega, \mathcal{F}, P)$  we can construct a new probability space  $(\tilde{\Omega}, \tilde{\mathcal{F}}, \tilde{P})$ , a random variable  $\tilde{\xi}_1(\tilde{\omega})$  and a measure-preserving transformation  $\tilde{T}$ , such that the distribution of  $\tilde{\xi} = \{\tilde{\xi}_1(\tilde{\omega}), \tilde{\xi}_1(\tilde{T}\tilde{\omega}), \dots\}$  coincides with the distribution of  $\xi$ .

In fact, take  $\tilde{\Omega}$  to be the coordinate space  $R^\infty$  and put  $\tilde{\mathcal{F}} = \mathcal{B}(R^\infty)$ ,  $\tilde{P} = P_\xi$ , where  $P_\xi(B) = P\{\omega: \xi \in B\}$ ,  $B \in \mathcal{B}(R^\infty)$ . The action of  $\tilde{T}$  on  $\tilde{\Omega}$  is given by

$$\tilde{T}(x_1, x_2, \dots) = (x_2, x_3, \dots).$$

If  $\tilde{\omega} = (x_1, x_2, \dots)$ , put

$$\tilde{\xi}_1(\tilde{\omega}) = x_1, \quad \tilde{\xi}_n(\tilde{\omega}) = \tilde{\xi}_1(\tilde{T}^{n-1}\tilde{\omega}), \quad n \geq 2.$$

Now let  $A = \{\tilde{\omega}: (x_1, \dots, x_k) \in B\}$ ,  $B \in \mathcal{B}(R^k)$ , and

$$\tilde{T}^{-1}A = \{\tilde{\omega}: (x_2, \dots, x_{k+1}) \in B\}.$$

Then the property of being stationary means that

$$\tilde{P}(A) = P\{\omega: (\xi_1, \dots, \xi_k) \in B\} = P\{\omega: (\xi_2, \dots, \xi_{k+1}) \in B\} = \tilde{P}(\tilde{T}^{-1}A),$$

i.e.  $\tilde{T}$  is a measure-preserving transformation. Since  $\tilde{P}\{\tilde{\omega}: (\tilde{\xi}_1, \dots, \tilde{\xi}_k) \in B\} = \tilde{P}\{\omega: (\xi_1, \dots, \xi_k) \in B\}$  for every  $k$ , it follows that  $\xi$  and  $\tilde{\xi}$  have the same distribution.

Here are some examples of measure-preserving transformations.

**EXAMPLE 1.** Let  $\Omega = \{\omega_1, \dots, \omega_n\}$  consist of  $n$  points (a finite number),  $n \geq 2$ , let  $\mathcal{F}$  be the collection of its subsets, and let  $T\omega_i = \omega_{i+1}$ ,  $1 \leq i \leq n-1$ , and  $T\omega_n = \omega_1$ . If  $P(\omega_i) = 1/n$ , the transformation  $T$  is measure-preserving.

**EXAMPLE 2.** If  $\Omega = [0, 1)$ ,  $\mathcal{F} = \mathcal{B}([0, 1))$ ,  $P$  is Lebesgue measure,  $\lambda \in [0, 1)$ , then  $Tx = (x + \lambda) \bmod 1$  and  $T = 2x \bmod 1$  are both measure-preserving transformations.

2. Let us pause to consider the physical hypotheses that led to the consideration of measure-preserving transformations.

Let us suppose that  $\Omega$  is the phase space of a system that evolves (in discrete time) according to a given law of motion. If  $\omega$  is the state at instant  $n = 1$ , then  $T^n\omega$ , where  $T$  is the translation operator induced by the given law of motion, is the state attained by the system after  $n$  steps. Moreover, if  $A$  is some set of states  $\omega$  then  $T^{-1}A = \{\omega: T\omega \in A\}$  is, by definition, the set of states  $\omega$  that lead to  $A$  in one step. Therefore if we interpret  $\Omega$  as an incompressible fluid, the condition  $P(T^{-1}A) = P(A)$  can be thought of as the rather natural condition of conservation of volume. (For the classical conservative Hamiltonian systems, Liouville's theorem asserts that the corresponding transformation  $T$  preserves Lebesgue measure.)

3. One of the earliest results on measure-preserving transformations was Poincaré's recurrence theorem (1912).

**Theorem 1.** Let  $(\Omega, \mathcal{F}, P)$  be a probability space, let  $T$  be a measure-preserving transformation, and let  $A \in \mathcal{F}$ . Then, for almost every point  $\omega \in A$ , we have  $T^n\omega \in A$  for infinitely many  $n \geq 1$ .

**PROOF.** Let  $C = \{\omega \in A: T^n\omega \notin A, \text{ for all } n \geq 1\}$ . Since  $C \cap T^{-n}C = \emptyset$  for all  $n \geq 1$ , we have  $T^{-m}C \cap T^{-(m+n)}C = T^{-m}(C \cap T^{-n}C) = \emptyset$ . Therefore the sequence  $\{T^{-n}C\}$  consists of disjoint sets of equal measure. Therefore  $\sum_{n=0}^{\infty} P(C) = \sum_{n=0}^{\infty} P(T^{-n}C) \leq P(\Omega) = 1$  and consequently  $P(C) = 0$ . Therefore, for almost every point  $\omega \in A$ , for at least one  $n \geq 1$ , we have  $T^n\omega \in A$ . It follows that  $T^n\omega \in A$  for infinitely many  $n$ .

Let us apply the preceding result to  $T^k$ ,  $k \geq 1$ . Then for every  $\omega \in A \setminus N$ , where  $N$  is a set of probability zero, the union of the corresponding sets corresponding to the various values of  $k$ , there is an  $n_k$  such that  $(T^k)^{n_k}\omega \in A$ . It is then clear that  $T^n\omega \in A$  for infinitely many  $n$ . This completes the proof of the theorem.

**Corollary.** Let  $\xi(\omega) \geq 0$ . Then

$$\sum_{k=0}^{\infty} \xi(T^k\omega) = \infty \quad (P\text{-a.s.})$$

on the set  $\{\omega: \xi(\omega) > 0\}$ .

In fact, let  $A_n = \{\omega: \xi(\omega) \geq 1/n\}$ . Then, according to the theorem,  $\sum_{k=0}^{\infty} \xi(T^k \omega) = \infty$  (P-a.s.) on  $A_n$ , and the required result follows by letting  $n \rightarrow \infty$ .

**Remark.** The theorem remains valid if we replace the probability measure  $P$  by any finite measure  $\mu$  with  $\mu(\Omega) < \infty$ .

#### 4. PROBLEMS

1. Let  $T$  be a measure-preserving transformation and  $\xi = \xi(\omega)$  a random variable whose expectation  $E\xi(\omega)$  exists. Show that  $E\xi(\omega) = E\xi(T\omega)$ .
2. Show that the transformations in Examples 1 and 2 are measure-preserving.
3. Let  $\Omega = [0, 1]$ ,  $F = \mathcal{B}([0, 1])$  and let  $P$  be a measure whose distribution function is continuous. Show that the transformations  $Tx = \lambda x$ ,  $0 < \lambda < 1$ , and  $Tx = x^2$  are not measure-preserving.

## §2. Ergodicity and Mixing

1. In the present section  $T$  denotes a measure-preserving transformation on the probability space  $(\Omega, \mathcal{F}, P)$ .

**Definition 1.** A set  $A \in \mathcal{F}$  is *invariant* if  $T^{-1}A = A$ . A set  $A \in \mathcal{F}$  is *almost invariant* if  $A$  and  $T^{-1}A$  differ only by a set of measure zero, i.e.  $P(A \Delta T^{-1}A) = 0$ .

It is easily verified that the classes  $\mathcal{I}$  and  $\mathcal{I}^*$  of invariant or almost invariant sets, respectively, are  $\sigma$ -algebras.

**Definition 2.** A measure-preserving transformation  $T$  is *ergodic* (or *metrically transitive*) if every invariant set  $A$  has measure either zero or one.

**Definition 3.** A random variable  $\xi = \xi(\omega)$  is *invariant* (or *almost invariant*) if  $\xi(\omega) = \xi(T\omega)$  for all  $\omega \in \Omega$  (or for almost all  $\omega \in \Omega$ ).

The following lemma establishes a connection between invariant and almost invariant sets.

**Lemma 1.** If  $A$  is almost invariant, there is an invariant set  $B$  such that  $P(A \Delta B) = 0$ .

PROOF. Let  $B = \overline{\lim} T^{-n}A$ . Then  $T^{-1}B = \overline{\lim} T^{-(n+1)}A = B$ , i.e.  $B \in \mathcal{I}$ . It is easily seen that  $A \Delta B \subseteq \bigcup_{k=0}^{\infty} (T^{-k}A \Delta T^{-(k+1)}A)$ . But

$$P(T^{-k}A \Delta T^{-(k+1)}A) = P(A \Delta T^{-1}A) = 0.$$

Hence  $P(A \Delta B) = 0$ .

**Lemma 2.** *A transformation  $T$  is ergodic if and only if every almost invariant set has measure zero or one.*

**PROOF.** Let  $A \in \mathcal{I}^*$ ; then according to Lemma 1 there is an invariant set  $B$  such that  $P(A \triangle B) = 0$ . But  $T$  is ergodic and therefore  $P(B) = 0$  or 1. Therefore  $P(A) = 0$  or 1. The converse is evident, since  $\mathcal{I} \subseteq \mathcal{I}^*$ . This completes the proof of the lemma.

**Theorem 1.** *Let  $T$  be a measure-preserving transformation. Then the following conditions are equivalent:*

- (1)  $T$  is ergodic;
- (2) every almost invariant random variable is (P-a.s.) constant;
- (3) every invariant random variable is (P-a.s.) constant.

**PROOF.** (1)  $\Leftrightarrow$  (2). Let  $T$  be ergodic and  $\xi$  almost invariant, i.e. (P-a.s.)  $\xi(\omega) = \xi(T\omega)$ . Then for every  $c \in R$  we have  $A_c = \{\omega: \xi(\omega) \leq c\} \in \mathcal{I}^*$ , and then  $P(A_c) = 0$  or 1 by Lemma 2. Let  $C = \sup\{c: P(A_c) = 0\}$ . Since  $A_c \uparrow \Omega$  as  $c \uparrow \infty$  and  $A_c \downarrow \emptyset$  as  $c \downarrow -\infty$ , we have  $|C| < \infty$ . Then

$$P\{\omega: \xi(\omega) < C\} = P\left\{\bigcup_{n=1}^{\infty} \left\{\xi(\omega) \leq C - \frac{1}{n}\right\}\right\} = 0$$

and similarly  $P\{\omega: \xi(\omega) > C\} = 0$ . Consequently  $P\{\omega: \xi(\omega) = C\} = 1$ .

(2)  $\Rightarrow$  (3). Evident.

(3)  $\Rightarrow$  (1). Let  $A \in \mathcal{I}$ ; then  $I_A$  is an invariant random variable and therefore, (P-a.s.),  $I_A = 0$  or  $I_A = 1$ , whence  $P(A) = 0$  or 1,

**Remark.** The conclusion of the theorem remains valid in the case when "random variable" is replaced by "bounded-random variable".

We illustrate the theorem with the following example.

**EXAMPLE.** Let  $\Omega = [0, 1)$ ,  $\mathcal{F} = \mathcal{B}([0, 1))$ , let  $P$  be Lebesgue measure and let  $T\omega = (\omega + \lambda) \bmod 1$ . Let us show that  $T$  is ergodic if and only if  $\lambda$  is irrational.

Let  $\xi = \xi(\omega)$  be a random variable with  $E\xi^2(\omega) < \infty$ . Then we know that the Fourier series  $\sum_{n=-\infty}^{\infty} c_n e^{2\pi i n \omega}$  of  $\xi(\omega)$  converges in the mean square sense,  $\sum |c_n|^2 < \infty$ , and, because  $T$  is a measure-preserving transformation (Example 2, §1), we have (Problem 1, §1) that for the random variable  $\xi$

$$\begin{aligned} c_n E\xi(\omega) e^{2\pi i n \xi(\omega)} &= E\xi(T\omega) e^{2\pi i n T\omega} = e^{2\pi i n \lambda} E\xi(T\omega) e^{2\pi i n \omega} \\ &= e^{2\pi i n \lambda} E\xi(\omega) e^{2\pi i n \omega} = c_n e^{2\pi i n \lambda}. \end{aligned}$$

So  $c_n(1 - e^{2\pi i n \lambda}) = 0$ . By hypothesis,  $\lambda$  is irrational and therefore  $e^{2\pi i n \lambda} \neq 1$  for all  $n \neq 0$ . Therefore  $c_n = 0$ ,  $n \neq 0$ ,  $\xi(\omega) = c_0$  (P-a.s.), and  $T$  is ergodic by Theorem 1.

On the other hand, let  $\lambda$  be rational, i.e.  $\lambda = k/m$ , where  $k$  and  $m$  are integers. Consider the set

$$A = \bigcup_{k=0}^{2m-2} \left\{ \omega: \frac{k}{2m} \leq \omega < \frac{k+1}{2m} \right\}.$$

It is clear that this set is invariant; but  $P(A) = \frac{1}{2}$ . Consequently  $T$  is not ergodic.

**2. Definition 4.** A measure-preserving transformation is *mixing* (or has the mixing property) if, for all  $A$  and  $B \in \mathcal{F}$ ,

$$\lim_{n \rightarrow \infty} P(A \cap T^{-n}B) = P(A)P(B). \quad (1)$$

The following theorem establishes a connection between ergodicity and mixing.

**Theorem 2.** Every mixing transformation  $T$  is ergodic.

**PROOF.** Let  $A \in \mathcal{F}$ ,  $B \in \mathcal{J}$ . Then  $B = T^{-n}A$ ,  $n \geq 1$ , and therefore

$$P(A \cap T^{-n}B) = P(A \cap B)$$

for all  $n \geq 1$ . Because of (1),  $P(A \cap B) = P(A)P(B)$ . Hence we find, when  $A = B$ , that  $P(B) = P^2(B)$ , and consequently  $P(B) = 0$  or  $1$ . This completes the proof.

### 3. PROBLEMS

1. Show that a random variable  $\xi$  is invariant if and only if it is  $\mathcal{J}$ -measurable.
2. Show that a set  $A$  is almost invariant if and only if either

$$P(T^{-1}A \setminus A) = 0 \quad \text{or} \quad P(A \setminus T^{-1}A) = 0.$$

3. Show that the transformation considered in the example of Subsection 1 of the present section is not mixing.
4. Show that a transformation is mixing if and only if, for all random variables  $\xi$  and  $\eta$  with  $E\xi^2 < \infty$  and  $E\eta^2 < \infty$ ,

$$E\xi(T^n\omega)\eta(\omega) \rightarrow E\xi(\omega)E\eta(\omega), \quad n \rightarrow \infty.$$

## §3. Ergodic Theorems

**1. Theorem 1 (Birkhoff and Khinchin).** Let  $T$  be a measure-preserving transformation and  $\xi = \xi(\omega)$  a random variable with  $E|\xi| < \infty$ . Then (P-a.s.)

$$\lim_n \frac{1}{n} \sum_{k=0}^{n-1} \xi(T^k\omega) = E(\xi|\mathcal{J}). \quad (1)$$

If also  $T$  is ergodic then (P-a.s.)

$$\lim_n \frac{1}{n} \sum_{k=0}^{n-1} \xi(T^k \omega) = E\xi. \quad (2)$$

The proof given below is based on the following proposition, whose simple proof was given by A. Garsia (1965).

**Lemma (Maximal Ergodic Theorem).** Let  $T$  be a measure-preserving transformation, let  $\xi$  be a random variable with  $E|\xi| < \infty$ , and let

$$S_k(\omega) = \xi(\omega) + \xi(T\omega) + \dots + \xi(T^{k-1}\omega),$$

$$M_k(\omega) = \max\{0, S_1(\omega), \dots, S_k(\omega)\}.$$

Then

$$E[\xi(\omega)I_{\{M_n > 0\}}(\omega)] \geq 0$$

for every  $n \geq 1$ .

PROOF. If  $n \geq k$ , we have  $M_n(T\omega) \geq S_k(T\omega)$  and therefore  $\xi(\omega) + M_n(T\omega) \geq \xi(\omega) + S_k(T\omega) = S_{k+1}(\omega)$ . Since it is evident that  $\xi(\omega) \geq S_1(\omega) - M_n(T\omega)$ , we have

$$\xi(\omega) \geq \max\{S_1(\omega), \dots, S_n(\omega)\} - M_n(T\omega).$$

Therefore

$$E[\xi(\omega)I_{\{M_n > 0\}}(\omega)] \geq E(\max\{S_1(\omega), \dots, S_n(\omega)\} - M_n(T\omega)),$$

But  $\max\{S_1, \dots, S_n\} = M_n$  on the set  $\{M_n > 0\}$ . Consequently,

$$\begin{aligned} E[\xi(\omega)I_{\{M_n > 0\}}(\omega)] &\geq E\{(M_n(\omega) - M_n(T\omega))I_{\{M_n(\omega) > 0\}}\} \\ &\geq E\{M_n(\omega) - M_n(T\omega)\} = 0, \end{aligned}$$

since if  $T$  is a measure-preserving transformation we have  $EM_n(\omega) = EM_n(T\omega)$  (Problem 1, §1).

This completes the proof of the lemma.

PROOF OF THE THEOREM. Let us suppose that  $E(\xi|\mathcal{J}) = 0$  (otherwise replace  $\xi$  by  $\xi - E(\xi|\mathcal{J})$ ).

Let  $\bar{\eta} = \overline{\lim}(S_n/n)$  and  $\underline{\eta} = \underline{\lim}(S_n/n)$ . It will be enough to establish that (P-a.s.)

$$0 \leq \underline{\eta} \leq \bar{\eta} \leq 0.$$

Consider the random variable  $\bar{\eta} = \bar{\eta}(\omega)$ . Since  $\bar{\eta}(\omega) = \bar{\eta}(T\omega)$ , the variable  $\bar{\eta}$  is invariant and consequently, for every  $\varepsilon > 0$ , the set  $A_\varepsilon = \{\bar{\eta}(\omega) > \varepsilon\}$  is also invariant. Let us introduce the new random variable

$$\xi^*(\omega) = (\xi(\omega) - \varepsilon)I_{A_\varepsilon}(\omega),$$

and put

$$S_k^*(\omega) = \xi^*(\omega) + \dots + \xi^*(T^{k-1}\omega), \quad M_k^*(\omega) = \max(0, S_1^*, \dots, S_k^*).$$

Then, by the lemma,

$$E[\xi^* I_{\{M_n^* > 0\}}] \geq 0$$

for every  $n \geq 1$ . But as  $n \rightarrow \infty$ ,

$$\begin{aligned} \{M_n^* > 0\} &= \left\{ \max_{1 \leq k \leq n} S_k^* > 0 \right\} \uparrow \left\{ \sup_{k \geq 1} S_k^* > 0 \right\} = \left\{ \sup_{k \geq 1} \frac{S_k^*}{k} > 0 \right\} \\ &= \left\{ \sup_{k \geq 1} \frac{S_k}{k} > \varepsilon \right\} \cap A_\varepsilon = A_\varepsilon, \end{aligned}$$

where the last equation follows because  $\sup_{k \geq 1} (S_k^*/k) \geq \bar{\eta}$ , and  $A_\varepsilon = \{\omega: \bar{\eta} > \varepsilon\}$ .

Moreover,  $E|\xi^*| \leq E|\xi| + \varepsilon$ . Hence, by the dominated convergence theorem,

$$0 \leq E[\xi^* I_{\{M_n^* > 0\}}] \rightarrow E[\xi^* I_A].$$

Thus

$$\begin{aligned} 0 &\leq E[\xi^* I_{A_\varepsilon}] = E[(\xi - \varepsilon) I_{A_\varepsilon}] = E[\xi I_{A_\varepsilon}] - \varepsilon P(A_\varepsilon) \\ &= E[E(\xi | \mathcal{J}) I_{A_\varepsilon}] - \varepsilon P(A_\varepsilon) = -\varepsilon P(A_\varepsilon), \end{aligned}$$

so that  $P(A_\varepsilon) = 0$  and therefore  $P(\bar{\eta} \leq 0) = 1$ .

Similarly, if we consider  $-\xi(\omega)$  instead of  $\xi(\omega)$ , we find that

$$\varliminf \left( -\frac{S_n}{n} \right) = -\varlimsup \frac{S_n}{n} = -\eta$$

and  $P(-\eta \leq 0) = 1$ , i.e.  $P(\eta \geq 0) = 1$ . Therefore  $0 \leq \underline{\eta} \leq \bar{\eta} \leq 0$  (P-a.s.) and the first part of the theorem is established.

To prove the second part, we observe that since  $E(\xi | \mathcal{J})$  is an invariant random variable, we have  $E(\xi | \mathcal{J}) = E\xi$  (P-a.s.) in the ergodic case.

This completes the proof of the theorem.

**Corollary.** A measure-preserving transformation  $T$  is ergodic if and only if, for all  $A$  and  $B \in \mathcal{F}$ ,

$$\lim_n \frac{1}{n} \sum_{k=0}^{n-1} P(A \cap T^{-k}B) = P(A)P(B). \quad (3)$$

To prove the ergodicity of  $T$  we use  $A = B \in \mathcal{J}$  in (3). Then  $A \cap T^{-k}B = B$  and therefore  $P(B) = P^2(B)$ , i.e.  $P(B) = 0$  or  $1$ . Conversely, let  $T$  be ergodic. Then if we apply (2) to the random variable  $\xi = I_B(\omega)$ , where  $B \in \mathcal{F}$ , we find that (P-a.s.)

$$\lim_n \frac{1}{n} \sum_{k=0}^{n-1} I_{T^{-k}B}(\omega) = P(B).$$



If we now integrate both sides over  $A \in \mathcal{F}$  and use the dominated convergence theorem, we obtain (3) as required.

2. We now show that, under the hypotheses of Theorem 1, there is not only almost sure convergence in (1) and (2), but also convergence in mean. (This result will be used below in the proof of Theorem 3.)

**Theorem 2.** *Let  $T$  be a measure-preserving transformation and let  $\xi = \xi(\omega)$  be a random variable with  $E|\xi| < \infty$ . Then*

$$E \left| \frac{1}{n} \sum_{k=0}^{n-1} \xi(T^k \omega) - E(\xi | \mathcal{J}) \right| \rightarrow 0, \quad n \rightarrow \infty. \quad (4)$$

If also  $T$  is ergodic, then

$$E \left| \frac{1}{n} \sum_{k=0}^{n-1} \xi(T^k \omega) - E\xi \right| \rightarrow 0, \quad n \rightarrow \infty. \quad (5)$$

PROOF. For every  $\varepsilon > 0$  there is a bounded random variable  $\eta(|\eta(\omega)| \leq M)$  such that  $E|\xi - \eta| \leq \varepsilon$ . Then

$$\begin{aligned} E \left| \frac{1}{n} \sum_{k=0}^{n-1} \xi(T^k \omega) - E(\xi | \mathcal{J}) \right| &\leq E \left| \frac{1}{n} \sum_{k=0}^{n-1} (\xi(T^k \omega) - \eta(T^k \omega)) \right| \\ &\quad + E \left| \frac{1}{n} \sum_{k=0}^{n-1} \eta(T^k \omega) - E(\eta | \mathcal{J}) \right| + E|E(\xi | \mathcal{J}) - E(\eta | \mathcal{J})|. \end{aligned} \quad (6)$$

Since  $|\eta| \leq M$ , then by the dominated convergence theorem and by using (1) we find that the second term on the right of (6) tends to zero as  $n \rightarrow \infty$ . The first and third terms are each at most  $\varepsilon$ . Hence for sufficiently large  $n$  the left-hand side of (6) is less than  $2\varepsilon$ , so that (4) is proved. Finally, if  $T$  is ergodic, then (5) follows from (4) and the remark that  $E(\xi | I) = E\xi$  (P-a.s.).

This completes the proof of the theorem.

3. We now turn to the question of the validity of the ergodic theorem for stationary (strict sense) random sequences  $\xi = (\xi_1, \xi_2, \dots)$  defined on a probability space  $(\Omega, \mathcal{F}, P)$ . In general,  $(\Omega, \mathcal{F}, P)$  need not carry any measure-preserving transformations, so that it is not possible to apply Theorem 1 directly. However, as we observed in §1, we can construct a coordinate probability space  $(\tilde{\Omega}, \tilde{\mathcal{F}}, \tilde{P})$ , random variables  $\tilde{\xi} = (\tilde{\xi}_1, \tilde{\xi}_2, \dots)$ , and a measure-preserving transformation  $\tilde{T}$  such that  $\tilde{\xi}_n(\tilde{\omega}) = \tilde{\xi}_1(\tilde{T}^{n-1}\tilde{\omega})$  and the distributions of  $\xi$  and  $\tilde{\xi}$  are the same. Since such properties as almost sure convergence and convergence in the mean are defined only for probability distributions, from the convergence of  $(1/n) \sum_{k=1}^n \tilde{\xi}_1(\tilde{T}^{k-1}\tilde{\omega})$  (P-a.s. and in mean) to a random variable  $\tilde{\eta}$  it follows that  $(1/n) \sum_{k=1}^n \xi_k(\omega)$  also converges (P-a.s. and in mean) to a random variable  $\eta$  such that  $\eta \stackrel{d}{=} \tilde{\eta}$ . It

follows from Theorem 1 that if  $\tilde{E}|\xi_1| < \infty$  then  $\tilde{E}(\xi_1|\tilde{\mathcal{F}}) = \tilde{E}(\xi_1|\tilde{\mathcal{F}})$ , where  $\tilde{\mathcal{F}}$  is a collection of invariant sets ( $\tilde{E}$  is the average with respect to the measure  $\tilde{P}$ ). We now describe the structure of  $\eta$ .

**Definition 1.** A set  $A \in \mathcal{F}$  is *invariant* with respect to the sequence  $\xi$  if there is a set  $B \in \mathcal{B}(R^\infty)$  such that for  $n \geq 1$

$$A = \{\omega: (\xi_n, \xi_{n+1}, \dots) \in B\}.$$

The collection of all such invariant sets is a  $\sigma$ -algebra, denoted by  $\mathcal{I}_\xi$ .

**Definition 2.** A stationary sequence  $\xi$  is *ergodic* if the measure of every invariant set is either 0 or 1.

Let us now show that the random variable  $\eta$  can be taken equal to  $E(\xi_1|\mathcal{I}_\xi)$ . In fact, let  $A \in \mathcal{I}_\xi$ . Then since

$$E\left|\frac{1}{n} \sum_{k=1}^{n-1} \xi_k - \eta\right| \rightarrow 0,$$

we have

$$\frac{1}{n} \sum_{k=1}^n \int_A \xi_k dP \rightarrow \int_A \eta dP. \quad (7)$$

Let  $B \in \mathcal{B}(R^\infty)$  be such that  $A = \{\omega: (\xi_k, \xi_{k+1}, \dots) \in B\}$  for all  $k \geq 1$ . Then since  $\xi$  is stationary,

$$\int_A \xi_k dP = \int_{\{\omega: (\xi_k, \xi_{k+1}, \dots) \in B\}} \xi_k dP = \int_{\{\omega: (\xi_1, \xi_2, \dots) \in B\}} \xi_1 dP = \int_A \xi_1 dP.$$

Hence it follows from (7) that for all  $A \in \mathcal{I}_\xi$ , which implies (see §7, Chapter II) that  $\eta = E(\xi_1|\mathcal{I}_\xi)$ . Here  $E(\xi_1|\mathcal{I}_\xi) = E\xi_1$  if  $\xi$  is ergodic.

Therefore we have proved the following theorem.

**Theorem 3 (Ergodic Theorem).** Let  $\xi = (\xi_1, \xi_2, \dots)$  be a stationary (strict sense) random sequence with  $E|\xi_1| < \infty$ . Then (P-a.s., and in the mean)

$$\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{k=1}^n \xi_k(\omega) = E(\xi_1|\mathcal{I}_\xi).$$

If  $\xi$  is also an ergodic sequence, then (P-a.s., and in the mean)

$$\lim_{n \rightarrow \infty} \frac{1}{n} \sum_{k=1}^n \xi_k(\omega) = E\xi_1.$$

#### 4. PROBLEMS

1. Let  $\xi = (\xi_1, \xi_2, \dots)$  be a Gaussian stationary sequence with  $E\xi_n = 0$  and covariance function  $R(n) = E\xi_{k+n}\xi_k$ . Show that  $R(n) \rightarrow 0$  is a sufficient condition for  $\xi$  to be ergodic.

2. Show that every sequence  $\xi = (\xi_1, \xi_2, \dots)$  of independent identically distributed random variables is ergodic.
3. Show that a stationary sequence  $\xi$  is ergodic if and only if

$$\frac{1}{n} \sum_{i=1}^n I_B(\xi_i, \dots, \xi_{i+k}) \rightarrow \mathbf{P}((\xi_1, \dots, \xi_{1+k}) \in B) \quad (\text{P-a.s.})$$

for every  $B \in \mathcal{B}(R)$ ,  $k = 1, 2, \dots$

## CHAPTER VI

# Stationary (Wide Sense) Random Sequences. $L^2$ -Theory

### §1. Spectral Representation of the Covariance Function

1. According to the definition given in the preceding chapter, a random sequence  $\xi = (\xi_1, \xi_2, \dots)$  is stationary in the strict sense if, for every set  $B \in \mathcal{B}(R^\infty)$  and every  $n \geq 1$ ,

$$P\{(\xi_1, \xi_2, \dots) \in B\} = P\{(\xi_{n+1}, \xi_{n+2}, \dots) \in B\}. \quad (1)$$

It follows, in particular, that if  $E\xi_1^2 < \infty$  then  $E\xi_n$  is independent of  $n$ :

$$E\xi_n = E\xi_1, \quad (2)$$

and the covariance  $\text{cov}(\xi_{n+m}, \xi_n) = E(\xi_{n+m} - E\xi_{n+m})(\xi_n - E\xi_n)$  depends only on  $m$ :

$$\text{cov}(\xi_{n+m}, \xi_n) = \text{cov}(\xi_{1+m}, \xi_1). \quad (3)$$

In the present chapter we study sequences that are stationary in the wide sense (and have finite second moments), namely those for which (1) is replaced by the (weaker) conditions (2) and (3).

The random variables  $\xi_n$  are understood to be defined for  $n \in \mathbb{Z} = \{0, \pm 1, \dots\}$  and to be complex-valued. The latter assumption not only does not complicate the theory, but makes it more elegant. It is also clear that results for real random variables can easily be obtained as special cases of the corresponding results for complex random variables.

Let  $H^2 = H^2(\Omega, \mathcal{F}, P)$  be the space of (complex) random variables  $\xi = \alpha + i\beta$ ,  $\alpha, \beta \in R$ , with  $E|\xi|^2 < \infty$ , where  $|\xi|^2 = \alpha^2 + \beta^2$ . If  $\xi$  and  $\eta \in H^2$ , we put

$$(\xi, \eta) = E\xi\bar{\eta}, \quad (4)$$

where  $\bar{\eta} = \alpha - i\beta$  is the complex conjugate of  $\eta = \alpha + i\beta$  and

$$\|\xi\| = (\xi, \xi)^{1/2}. \quad (5)$$

As for real random variables, the space  $H^2$  (more precisely, the space of equivalence classes of random variables; compare §§10 and 11 of Chapter II) is complete under the scalar product  $(\xi, \eta)$  and norm  $\|\xi\|$ . In accordance with the terminology of functional analysis,  $H^2$  is called the complex (or unitary) Hilbert space (of random variables considered on the probability space  $(\Omega, \mathcal{F}, P)$ ).

If  $\xi, \eta \in H^2$  their covariance is

$$\text{cov}(\xi, \eta) = E(\xi - E\xi)(\overline{\eta - E\eta}). \quad (6)$$

It follows from (4) and (6) that if  $E\xi = E\eta = 0$  then

$$\text{cov}(\xi, \eta) = (\xi, \eta). \quad (7)$$

**Definition.** A sequence of complex random variables  $\xi = (\xi_n)_{n \in \mathbb{Z}}$  with  $E|\xi_n|^2 < \infty, n \in \mathbb{Z}$ , is *stationary (in the wide sense)* if, for all  $n \in \mathbb{Z}$ ,

$$\begin{aligned} E\xi_n &= E\xi_0, \\ \text{cov}(\xi_{k+n}, \xi_k) &= \text{cov}(\xi_n, \xi_0), \quad k \in \mathbb{Z}. \end{aligned} \quad (8)$$

As a matter of convenience, we shall always suppose that  $E\xi_0 = 0$ . This involves no loss of generality, but does make it possible (by (7)) to identify the covariance with the scalar product and hence to apply the methods and results of the theory of Hilbert spaces.

Let us write

$$R(n) = \text{cov}(\xi_n, \xi_0), \quad n \in \mathbb{Z}, \quad (9)$$

and (assuming  $R(0) = E|\xi_0|^2 \neq 0$ )

$$\rho(n) = \frac{R(n)}{R(0)}, \quad n \in \mathbb{Z}. \quad (10)$$

We call  $R(n)$  the *covariance function*, and  $\rho(n)$ , the *correlation function*, of the sequence  $\xi$  (assumed stationary in the wide sense).

It follows immediately from (9) that  $R(n)$  is nonnegative-definite, i.e. for all complex numbers  $a_1, \dots, a_m$  and  $t_1, \dots, t_m \in \mathbb{Z}, m \geq 1$ , we have

$$\sum_{i,j=1}^m a_i \bar{a}_j R(t_i - t_j) \geq 0. \quad (11)$$

It is then easy to deduce (either from (11) or directly from (9)) the following properties of the covariance function (see Problem 1):

$$\begin{aligned} R(0) &\geq 0, \quad R(-n) = \overline{R(n)}, \quad |R(n)| \leq R(0), \\ |R(n) - R(m)|^2 &\leq 2R(0)[R(0) - \text{Re } R(n - m)]. \end{aligned} \quad (12)$$

2. Let us give some examples of stationary sequences  $\xi = (\xi_n)_{n \in \mathbb{Z}}$ . (From now on, the words “in the wide sense” and the statement  $n \in \mathbb{Z}$  will both be omitted.)

EXAMPLE 1. Let  $\xi_n = \xi_0 \cdot g(n)$ , where  $E\xi_0 = 0$ ,  $E\xi_0^2 = 1$  and  $g = g(n)$  is a function. The sequence  $\xi = (\xi_n)$  will be stationary if and only if  $g(k+n)\overline{g(k)}$  depends only on  $n$ . Hence it is easy to see that there is a  $\lambda$  such that

$$g(n) = g(0)e^{i\lambda n}.$$

Consequently the sequence of random variables

$$\xi_n = \xi_0 \cdot g(0)e^{i\lambda n}$$

is stationary with

$$R(n) = |g(0)|^2 e^{i\lambda n}.$$

In particular, the random “constant”  $\xi \equiv \xi_0$  is a stationary sequence.

EXAMPLE 2. An almost periodic sequence. Let

$$\xi_n = \sum_{k=1}^N z_k e^{i\lambda_k n}, \quad (13)$$

where  $z_1, \dots, z_N$  are orthogonal ( $E z_i \bar{z}_j = 0$ ,  $i \neq j$ ) random variables with zero means and  $E|z_k|^2 = \sigma_k^2 > 0$ ;  $-\pi \leq \lambda_k < \pi$ ,  $k = 1, \dots, N$ ;  $\lambda_i \neq \lambda_j$ ,  $i \neq j$ . The sequence  $\xi = (\xi_n)$  is stationary with

$$R(n) = \sum_{k=1}^N \sigma_k^2 e^{i\lambda_k n}. \quad (14)$$

As a generalization of (13) we now suppose that

$$\xi_n = \sum_{k=-\infty}^{\infty} z_k e^{i\lambda_k n}, \quad (15)$$

where  $z_k$ ,  $k \in \mathbb{Z}$ , have the same properties as in (13). If we suppose that  $\sum_{k=-\infty}^{\infty} \sigma_k^2 < \infty$ , the series on the right of (15) converges in mean-square and

$$R(n) = \sum_{k=-\infty}^{\infty} \sigma_k^2 e^{i\lambda_k n}. \quad (16)$$

Let us introduce the function

$$F(\lambda) = \sum_{\{k: \lambda_k \leq \lambda\}} \sigma_k^2. \quad (17)$$

Then the covariance function (16) can be written as a Lebesgue–Stieltjes integral,

$$R(n) = \int_{-\pi}^{\pi} e^{i\lambda n} dF(\lambda). \quad (18)$$

The stationary sequence (15) is represented as a sum of "harmonics"  $e^{i\lambda_k n}$  with "frequencies"  $\lambda_k$  and random "amplitudes"  $z_k$  of "intensities"  $\sigma_k^2 = E|z_k|^2$ . Consequently the values of  $F(\lambda)$  provide complete information on the "spectrum" of the sequence  $\xi$ , i.e. on the intensity with which each frequency appears in (15). By (18), the values of  $F(\lambda)$  also completely determine the structure of the covariance function  $R(n)$ .

Up to a constant multiple, a (nondegenerate)  $F(\lambda)$  is evidently a distribution function, which in the examples considered so far has been piecewise constant. It is quite remarkable that the covariance function of every stationary (wide sense) random sequence can be represented (see the theorem in Subsection 3) in the form (18), where  $F(\lambda)$  is a distribution function (up to normalization), whose support is concentrated on  $[-\pi, \pi)$ , i.e.  $F(\lambda) = 0$  for  $\lambda < -\pi$  and  $F(\lambda) = F(\pi)$  for  $\lambda > \pi$ .

The result on the integral representation of the covariance function, if compared with (15) and (16), suggests that every stationary sequence also admits an "integral" representation. This is in fact the case, as will be shown in §3 by using what we shall learn to call stochastic integrals with respect to orthogonal stochastic measures (§2).

**EXAMPLE 3 (White noise).** Let  $\varepsilon = (\varepsilon_n)$  be an orthonormal sequence of random variables,  $E\varepsilon_n = 0$ ,  $E\varepsilon_i \varepsilon_j = \delta_{ij}$ , where  $\delta_{ij}$  is the Kronecker delta. Such a sequence is evidently stationary, and

$$R(n) = \begin{cases} 1 & n = 0, \\ 0, & n \neq 0. \end{cases}$$

Observe that  $R(n)$  can be represented in the form

$$R(n) = \int_{-\pi}^{\pi} e^{i\lambda n} dF(\lambda), \quad (19)$$

where

$$F(\lambda) = \int_{-\pi}^{\lambda} f(v) dv; \quad f(\lambda) = \frac{1}{2\pi}, \quad -\pi \leq \lambda < \pi. \quad (20)$$

Comparison of the spectral functions (17) and (20) shows that whereas the spectrum in Example 2 is discrete, in the present example it is absolutely continuous with constant "spectral density"  $f(\lambda) \equiv \frac{1}{2\pi}$ . In this sense we can say that the sequence  $\varepsilon = (\varepsilon_n)$  "consists of harmonics of equal intensities." It is just this property that has led to calling such a sequence  $\varepsilon = (\varepsilon_n)$  "white noise" by analogy with white light, which consists of different frequencies with the same intensities.

**EXAMPLE 4 (Moving averages)** Starting from the white noise  $\varepsilon = (\varepsilon_n)$  introduced in Example 3, let us form the new sequence

$$\xi_n = \sum_{k=-\infty}^{\infty} a_k \varepsilon_{n-k}, \quad (21)$$

where  $a_k$  are complex numbers such that  $\sum_{k=-\infty}^{\infty} |a_k|^2 < \infty$ . By Parseval's equation,

$$\text{cov}(\xi_{n+m}, \xi_m) = \text{cov}(\xi_n, \xi_0) = \sum_{k=-\infty}^{\infty} a_{n+k} \bar{a}_k,$$

so that  $\xi = (\xi_k)$  is a stationary sequence, which we call the sequence obtained from  $\varepsilon = (\varepsilon_k)$  by a (*two-sided*) moving average.

In the special case when the  $a_k$  of negative index are zero, i.e.

$$\xi_n = \sum_{k=0}^{\infty} a_k \varepsilon_{n-k},$$

the sequence  $\xi = (\xi_n)$  is a *one-sided moving average*. If, in addition,  $a_k = 0$  for  $k > p$ , i.e. if

$$\xi_n = a_0 \varepsilon_n + a_1 \varepsilon_{n-1} + \cdots + a_p \varepsilon_{n-p}, \quad (22)$$

then  $\xi = (\xi_n)$  is a *moving average of order p*.

We can show (Problem 5) that (22) has a covariance function of the form  $R(n) = \int_{-\pi}^{\pi} e^{i\lambda n} f(\lambda) d\lambda$ , where the spectral density is

$$f(\lambda) = \frac{1}{2\pi} |P(e^{-i\lambda})|^2 \quad (23)$$

with

$$P(z) = a_0 + a_1 z + \cdots + a_p z^p.$$

**EXAMPLE 5 (Autoregression).** Again let  $\varepsilon = (\varepsilon_n)$  be white noise. We say that a random sequence  $\xi = (\xi_n)$  is described by an *autoregressive model* of order  $q$  if

$$\xi_n + b_1 \xi_{n-1} + \cdots + b_q \xi_{n-q} = \varepsilon_n. \quad (24)$$

Under what conditions on  $b_1, \dots, b_q$  can we say that (24) has a stationary solution? To find an answer, let us begin with the case  $q = 1$ :

$$\xi_n = \alpha \xi_{n-1} + \varepsilon_n, \quad (25)$$

where  $\alpha = -b_1$ . If  $|\alpha| < 1$ , it is easy to verify that the stationary sequence  $\xi = (\xi_n)$  with

$$\xi_n = \sum_{j=0}^{\infty} \alpha^j \varepsilon_{n-j} \quad (26)$$

is a solution of (25). (The series on the right of (26) converges in mean-square.) Let us now show that, in the class of stationary sequences  $\xi = (\xi_n)$  (with finite second moments) this is the only solution. In fact, we find from (25), by successive iteration, that

$$\xi_n = \alpha \xi_{n-1} + \varepsilon_n = \alpha[\alpha \xi_{n-2} + \varepsilon_{n-1}] + \varepsilon_n = \cdots = \alpha^k \xi_{n-k} + \sum_{j=0}^{k-1} \alpha^j \varepsilon_{n-j}.$$



Hence it follows that

$$\mathbb{E} \left[ \xi_n - \sum_{j=0}^{k-1} \alpha^j \varepsilon_{n-j} \right]^2 = \mathbb{E} [\alpha^k \xi_{n-k}]^2 = \alpha^{2k} \mathbb{E} \xi_{n-k}^2 = \alpha^{2k} \mathbb{E} \xi_0^2 \rightarrow 0, \quad k \rightarrow \infty.$$

Therefore when  $|\alpha| < 1$  a stationary solution of (25) exists and is representable as the one-sided moving average (26).

There is a similar result for every  $q > 1$ : if all the zeros of the polynomial

$$Q(z) = 1 + b_1 z + \dots + b_q z^q \quad (27)$$

lie outside the unit disk, then the autoregression equation (24) has a unique stationary solution, which is representable as a one-sided moving average (Problem 2). Here the covariance function  $R(n)$  can be represented (Problem 5) in the form

$$R(n) = \int_{-\pi}^{\pi} e^{i\lambda n} dF(\lambda), \quad F(\lambda) = \int_{-\pi}^{\lambda} f(v) dv, \quad (28)$$

where

$$f(\lambda) = \frac{1}{2\pi} \cdot \frac{1}{|Q(e^{-i\lambda})|^2}. \quad (29)$$

In the special case  $q = 1$ , we find easily from (25) that  $\mathbb{E} \xi_0 = 0$ ,

$$\mathbb{E} \xi_0^2 = \frac{1}{1 - |\alpha|^2},$$

and

$$R(n) = \frac{\alpha^n}{1 - |\alpha|^2}, \quad n \geq 0$$

(when  $n < 0$  we have  $R(n) = \overline{R(-n)}$ ). Here

$$f(\lambda) = \frac{1}{2\pi} \cdot \frac{1}{|1 - \alpha e^{-i\lambda}|^2}.$$

**EXAMPLE 6.** This example illustrates how autoregression arises in the construction of probabilistic models in hydrology. Consider a body of water; we try to construct a probabilistic model of the deviations of the level of the water from its average value because of variations in the inflow and evaporation from the surface.

If we take a year as the unit of time and let  $H_n$  denote the water level in year  $n$ , we obtain the following *balance equation*:

$$H_{n+1} = H_n - KS(H_n) + \Sigma_{n+1}, \quad (30)$$

where  $\Sigma_{n+1}$  is the inflow in year  $(n+1)$ ,  $S(H)$  is the area of the surface of the water at level  $H$ , and  $K$  is the coefficient of evaporation.

Let  $\xi_n = H_n - \bar{H}$  be the deviation from the mean level (which is obtained from observations over many years) and suppose that  $S(H) = S(\bar{H}) + c(H - \bar{H})$ . Then it follows from the balance equation that  $\xi_n$  satisfies

$$\xi_{n+1} = \alpha \xi_n + \varepsilon_{n+1} \quad (31)$$

with  $\alpha = 1 - cK$ ,  $\varepsilon_n = \Sigma_n - KS(\bar{H})$ . It is natural to assume that the random variables  $\varepsilon_n$  have zero means and are identically distributed. Then, as we showed in Example 5, equation (31) has (for  $|\alpha| < 1$ ) a unique stationary solution, which we think of as the steady-state solution (with respect to time in years) of the oscillations of the level in the body of water.

As an example of practical conclusions that can be drawn from a (theoretical) model (31), we call attention to the possibility of predicting the level for the following year from the results of the observations of the present and preceding years. It turns out (see also Example 2 in §6) that (in the mean-square sense) the optimal linear estimator of  $\xi_{n+1}$  in terms of the values of  $\dots, \xi_{n-1}, \xi_n$  is simply  $\alpha \xi_n$ .

**EXAMPLE 7** (Autoregression and moving average (mixed model)). If we suppose that the right-hand side of (24) contains  $\alpha_0 \varepsilon_n + \alpha_1 \varepsilon_{n-1} + \dots + \alpha_p \varepsilon_{n-p}$  instead of  $\varepsilon_n$ , we obtain a mixed model with autoregression and moving average of order  $(p, q)$ :

$$\xi_n + b_1 \xi_{n-1} + \dots + b_q \xi_{n-q} = a_0 \varepsilon_n + a_1 \varepsilon_{n-1} + \dots + a_p \varepsilon_{n-p}. \quad (32)$$

Under the same hypotheses as in Example 5 on the zeros it will be shown later (Corollary 2 to Theorem 3 of §3) that (32) has the stationary solution  $\xi = (\xi_n)$  for which the covariance function is  $R(n) = \int_{-\pi}^{\pi} e^{i\lambda n} dF(\lambda)$  with  $F(\lambda) = \int_{-\pi}^{\lambda} f(v) dv$ , where

$$f(\lambda) = \frac{1}{2\pi} \cdot \left| \frac{P(e^{-i\lambda})}{Q(e^{-i\lambda})} \right|^2.$$

**3. Theorem (Herglotz).** Let  $R(n)$  be the covariance function of a stationary (wide sense) random sequence with zero mean. Then there is, on

$$([-\pi, \pi), \mathcal{B}([-\pi, \pi))),$$

a finite measure  $F = F(B)$ ,  $B \in \mathcal{B}([-\pi, \pi))$ , such that for every  $n \in \mathbb{Z}$

$$R(n) = \int_{-\pi}^{\pi} e^{i\lambda n} F(d\lambda). \quad (33)$$

**PROOF.** For  $N \geq 1$  and  $\lambda \in [-\pi, \pi]$ , put

$$f_N(\lambda) = \frac{1}{2\pi N} \sum_{k=1}^N \sum_{l=1}^N R(k-l) e^{-ik\lambda} e^{il\lambda}. \quad (34)$$

Since  $R(n)$  is nonnegative definite,  $f_N(\lambda)$  is nonnegative. Since there are  $N - |m|$  pairs  $(k, l)$  for which  $k - l = m$ , we have

$$f_N(\lambda) = \frac{1}{2\pi} \sum_{|m| < N} \left(1 - \frac{|m|}{N}\right) R(m) e^{-im\lambda}. \quad (35)$$

Let

$$F_N(B) = \int_B f_N(\lambda) d\lambda, \quad B \in \mathcal{B}([-\pi, \pi]).$$

Then

$$\int_{-\pi}^{\pi} e^{i\lambda n} F_N(d\lambda) = \int_{-\pi}^{\pi} e^{i\lambda n} f_N(\lambda) d\lambda = \begin{cases} \left(1 - \frac{|n|}{N}\right) R(n), & |n| < N, \\ 0, & |n| \geq N. \end{cases} \quad (36)$$

The measures  $F_N$ ,  $N \geq 1$ , are supported on the interval  $[-\pi, \pi]$  and  $F_N([-\pi, \pi]) = R(0) < \infty$  for all  $N \geq 1$ . Consequently the family of measures  $\{F_N\}$ ,  $N \geq 1$ , is tight, and by Prokhorov's theorem (Theorem 1 of §2 of Chapter III) there are a sequence  $\{N_k\} \subseteq \{N\}$  and a measure  $F$  such that  $F_{N_k} \xrightarrow{w} F$ . (The concepts of tightness, relative compactness, and weak convergence, together with Prokhorov's theorem, can be extended in an obvious way from probability measures to any finite measures.)

It then follows from (36) that

$$\int_{-\pi}^{\pi} e^{i\lambda n} F(d\lambda) = \lim_{N_k \rightarrow \infty} \int_{-\pi}^{\pi} e^{i\lambda n} F_{N_k}(d\lambda) = R(n).$$

The measure  $F$  so constructed is supported on  $[-\pi, \pi]$ . Without changing the integral  $\int_{-\pi}^{\pi} e^{i\lambda n} F(d\lambda)$ , we can redefine  $F$  by transferring the "mass"  $F(\{\pi\})$ , which is concentrated at  $\pi$ , to  $-\pi$ . The resulting new measure (which we again denote by  $F$ ) will be supported on  $[-\pi, \pi)$ .

This completes the proof of the theorem.

**Remark 1.** The measure  $F = F(B)$  involved in (33) is known as the spectral measure, and  $F(\lambda) = F([-\pi, \lambda])$  as the spectral function, of the stationary sequence with covariance function  $R(n)$ .

In Example 2 above the spectral measure was discrete (concentrated at  $\lambda_k$ ,  $k = 0, \pm 1, \dots$ ). In Examples 3–6 the spectral measures were absolutely continuous.

**Remark 2.** The spectral measure  $F$  is uniquely defined by the covariance function. In fact, let  $F_1$  and  $F_2$  be two spectral measures and let

$$\int_{-\pi}^{\pi} e^{i\lambda n} F_1(d\lambda) = \int_{-\pi}^{\pi} e^{i\lambda n} F_2(d\lambda), \quad n \in \mathbb{Z}.$$

Since every bounded continuous function  $g(\lambda)$  can be uniformly approximated on  $[-\pi, \pi]$  by trigonometric polynomials, we have

$$\int_{-\pi}^{\pi} g(\lambda) F_1(d\lambda) = \int_{-\pi}^{\pi} g(\lambda) F_2(d\lambda).$$

It follows (compare the proof in Theorem 2, §12, Chapter II) that  $F_1(B) = F_2(B)$  for all  $B \in \mathcal{B}([-\pi, \pi])$ .

**Remark 3.** If  $\xi = (\xi_n)$  is a stationary sequence of real random variables  $\xi_n$ , then

$$R(n) = \int_{-\pi}^{\pi} \cos \lambda n F(d\lambda).$$

#### 4. PROBLEMS

1. Derive (12) from (11).
2. Show that the autoregression equation (24) has a stationary solution if all the zeros of the polynomial  $Q(z)$  defined by (27) lie outside the unit disk.
3. Prove that the covariance function (28) admits the representation (29) with spectral density given by (30).
4. Show that the sequence  $\xi = (\xi_n)$  of random variables, where

$$\xi_n = \sum_{k=1}^{\infty} (\alpha_k \sin \lambda_k n + \beta_k \cos \lambda_k n)$$

and  $\alpha_k$  and  $\beta_k$  are real random variables, can be represented in the form

$$\xi_n = \sum_{k=-\infty}^{\infty} z_k e^{i\lambda_k n}$$

with  $z_k = \frac{1}{2}(\beta_k - i\alpha_k)$  for  $k \geq 0$  and  $z_k = \bar{z}_{-k}$ ,  $\lambda_k = -\lambda_{-k}$  for  $k < 0$ .

5. Show that the spectral functions of the sequences (22) and (24) have densities given respectively by (23) and (29).
6. Show that if  $\sum |R(n)| < \infty$ , the spectral function  $F(\lambda)$  has density  $f(\lambda)$  given by

$$f(\lambda) = \frac{1}{2\pi} \sum_{n=-\infty}^{\infty} e^{-i\lambda n} R(n).$$

## §2. Orthogonal Stochastic Measures and Stochastic Integrals

1. As we observed in §1, the integral representation of the covariance function and the example of a stationary sequence

$$\xi_n = \sum_{k=-\infty}^{\infty} z_k e^{i\lambda_k n} \quad (1)$$

with pairwise orthogonal random variables  $z_k$ ,  $k \in \mathbb{Z}$ , suggest the possibility of representing an arbitrary stationary sequence as a corresponding integral generalization of (1).

If we put

$$Z(\lambda) = \sum_{\{k: \lambda_k \leq \lambda\}} z_k, \quad (2)$$

we can rewrite (1) in the form

$$\xi_n = \sum_{k=-\infty}^{\infty} e^{i\lambda_k n} \Delta Z(\lambda_k), \quad (3)$$

where  $\Delta Z(\lambda_k) \equiv Z(\lambda_k) - Z(\lambda_k -) = z_k$ .

The right-hand side of (3) reminds us of an approximating sum for an integral  $\int_{-\pi}^{\pi} e^{i\lambda n} dZ(\lambda)$  of Riemann-Stieltjes type. However, in the present case  $Z(\lambda)$  is a random function (it also depends on  $\omega$ ). Hence it is clear that for an integral representation of a general stationary sequence we need to use functions  $Z(\lambda)$  that do not have bounded variation for each  $\omega$ . Consequently the simple interpretation of  $\int_{-\pi}^{\pi} e^{i\lambda n} dZ(\lambda)$  as a Riemann-Stieltjes integral for each  $\omega$  is inapplicable.

2. By analogy with the general ideas of the Lebesgue, Lebesgue-Stieltjes and Riemann-Stieltjes integrals (§6, Chapter II), we begin by defining stochastic measure.

Let  $(\Omega, \mathcal{F}, P)$  be a probability space, and let  $E$  be a subset, with an algebra  $\mathcal{E}_0$  of subsets and the  $\sigma$ -algebra  $\mathcal{E}$  generated by  $\mathcal{E}_0$ .

**Definition 1.** A complex-valued function  $Z(\Delta) = Z(\omega; \Delta)$ , defined for  $\omega \in \Omega$  and  $\Delta \in \mathcal{E}_0$ , is a *finitely additive stochastic measure* if

- (1)  $E|Z(\Delta)|^2 < \infty$  for every  $\Delta \in \mathcal{E}_0$ ;
- (2) for every pair  $\Delta_1$  and  $\Delta_2$  of disjoint sets in  $\mathcal{E}_0$ ,

$$Z(\Delta_1 + \Delta_2) = Z(\Delta_1) + Z(\Delta_2) \quad (P\text{-a.s.}) \quad (4)$$

**Definition 2.** A finitely additive stochastic measure  $Z(\Delta)$  is an *elementary stochastic measure* if, for all disjoint sets  $\Delta_1, \Delta_2, \dots$  of  $\mathcal{E}_0$  such that  $\Delta = \sum_{k=1}^{\infty} \Delta_k \in \mathcal{E}_0$ ,

$$E \left| Z(\Delta) - \sum_{k=1}^n Z(\Delta_k) \right|^2 \rightarrow 0, \quad n \rightarrow \infty. \quad (5)$$

**Remark 1.** In this definition of an elementary stochastic measure on subsets of  $\mathcal{E}_0$ , it is assumed that its values are in the Hilbert space  $H^2 = H^2(\Omega, \mathcal{F}, P)$ , and that countable additivity is understood in the mean-square sense (5). There are other definitions of stochastic measures, without the requirement of the existence of second moments, where countable additivity is defined (for example) in terms of convergence in probability or with probability one.

**Remark 2.** In analogy with nonstochastic measures, one can show that for finitely additive stochastic measures the condition (5) of countable additivity (in the mean-square sense) is equivalent to continuity (in the mean-square sense) at “zero”:

$$E|Z(\Delta_n)|^2 \rightarrow 0, \quad \Delta_n \downarrow \emptyset, \quad \Delta_n \in \mathcal{E}_0. \quad (6)$$

A particularly important class of elementary stochastic measures consists of those that are orthogonal according to the following definition.

**Definition 3.** An elementary stochastic measure  $Z(\Delta)$ ,  $\Delta \in \mathcal{E}_0$ , is *orthogonal* (or a *measure with orthogonal values*) if

$$EZ(\Delta_1)\overline{Z(\Delta_2)} = 0 \quad (7)$$

for every pair of disjoint sets  $\Delta_1$  and  $\Delta_2$  in  $\mathcal{E}_0$ ; or, equivalently, if

$$EZ(\Delta_1)\overline{Z(\Delta_2)} = E|Z(\Delta_1 \cap \Delta_2)|^2 \quad (8)$$

for all  $\Delta_1$  and  $\Delta_2$  in  $\mathcal{E}_0$ .

We write

$$m(\Delta) = E|Z(\Delta)|^2, \quad \Delta \in \mathcal{E}_0. \quad (9)$$

For elementary orthogonal stochastic measures, the set function  $m = m(\Delta)$ ,  $\Delta \in \mathcal{E}_0$ , is, as is easily verified, a finite measure, and consequently by Carathéodory's theorem (§3, Chapter II) it can be extended to  $(E, \mathcal{E})$ . The resulting measure will again be denoted by  $m = m(\Delta)$  and called the *structure function* (of the elementary orthogonal stochastic measure  $Z = Z(\Delta)$ ,  $\Delta \in \mathcal{E}_0$ ).

The following question now arises naturally: since the set function  $m = m(\Delta)$  defined on  $(E, \mathcal{E}_0)$  admits an extension to  $(E, \mathcal{E})$ , where  $\mathcal{E} = \sigma(\mathcal{E}_0)$ , cannot an elementary orthogonal stochastic measure  $Z = Z(\Delta)$ ,  $\Delta \in \mathcal{E}_0$ , be extended to sets  $\Delta$  in  $E$  in such a way that  $E|Z(\Delta)|^2 = m(\Delta)$ ,  $\Delta \in \mathcal{E}$ ?

The answer is affirmative, as follows from the construction given below. This construction, at the same time, leads to the stochastic integral which we need for the integral representation of stationary sequences.

3. Let  $Z = Z(\Delta)$  be an elementary orthogonal stochastic measure,  $\Delta \in \mathcal{E}_0$ , with structure function  $m = m(\Delta)$ ,  $\Delta \in \mathcal{E}$ . For every function

$$f(\lambda) = \sum f_k I_{\Delta_k}, \quad \Delta_k \in \mathcal{E}_0, \quad (10)$$

with only a finite number of different (complex) values, we define the random variable

$$\mathcal{J}(f) = \sum f_k Z(\Delta_k).$$

Let  $L^2 = L^2(E, \mathcal{E}, m)$  be the Hilbert space of complex-valued functions with the scalar product

$$\langle f, g \rangle = \int_E f(\lambda) \bar{g}(\lambda) m(d\lambda)$$

and the norm  $\|f\| = \langle f, f \rangle^{1/2}$ , and let  $H^2 = H^2(\Omega, \mathcal{F}, P)$  be the Hilbert space of complex-valued random variables with the scalar product

$$(\xi, \eta) = E \xi \bar{\eta}$$

and the norm  $\|\xi\| = (\xi, \xi)^{1/2}$ .

Then it is clear that, for every pair of functions  $f$  and  $g$  of the form (10),

$$(\mathcal{J}(f), \mathcal{J}(g)) = \langle f, g \rangle$$

and

$$\|\mathcal{J}(f)\|^2 = \|f\|^2 = \int_E |f(\lambda)|^2 m(d\lambda).$$

Now let  $f \in L^2$  and let  $\{f_n\}$  be functions of type (10) such that  $\|f - f_n\| \rightarrow 0$ ,  $n \rightarrow \infty$  (the existence of such functions follows from Problem 2). Consequently

$$\|\mathcal{J}(f_n) - \mathcal{J}(f_m)\| = \|f_n - f_m\| \rightarrow 0, \quad n, m \rightarrow \infty.$$

Therefore the sequence  $\{\mathcal{J}(f_n)\}$  is fundamental in the mean-square sense and by Theorem 7, §10, Chapter II, there is a random variable (denoted by  $\mathcal{J}(f)$ ) such that  $\mathcal{J}(f) \in H^2$  and  $\|\mathcal{J}(f_n) - \mathcal{J}(f)\| \rightarrow 0$ ,  $n \rightarrow \infty$ .

The random variable  $\mathcal{J}(f)$  constructed in this way is uniquely defined (up to stochastic equivalence) and is independent of the choice of the approximating sequence  $\{f_n\}$ . We call it the *stochastic integral* of  $f \in L^2$  with respect to the elementary orthogonal stochastic measure  $Z$  and denote it by

$$\mathcal{J}(f) = \int_E f(\lambda) Z(d\lambda).$$

We note the following basic properties of the stochastic integral  $\mathcal{J}(f)$ ; these are direct consequences of its construction (Problem 1). Let  $g, f$ , and  $f_n \in L^2$ . Then

$$(\mathcal{J}(f), \mathcal{J}(g)) = \langle f, g \rangle; \quad (11)$$

$$\|\mathcal{J}(f)\| = \|f\|; \quad (12)$$

$$\mathcal{J}(af + bg) = a\mathcal{J}(f) + b\mathcal{J}(g) \quad (\text{P-a.s.}) \quad (13)$$

where  $a$  and  $b$  are constants;

$$\|\mathcal{J}(f_n) - \mathcal{J}(f)\| \rightarrow 0, \quad (14)$$

if  $\|f_n - f\| \rightarrow 0$ ,  $n \rightarrow \infty$ .

4. Let us use the preceding definition of the stochastic integral to *extend* the elementary stochastic measure  $Z(\Delta)$ ,  $\Delta \in \mathcal{E}_0$ , to sets in  $\mathcal{E} = \sigma(\mathcal{E}_0)$ .

Since  $m$  is assumed to be finite, we have  $I_\Delta = I_\Delta(\lambda) \in L^2$  for all  $\Delta \in \mathcal{E}$ . Write  $\tilde{Z}(\Delta) = \mathcal{J}(I_\Delta)$ . It is clear that  $\tilde{Z}(\Delta) = Z(\Delta)$  for  $\Delta \in \mathcal{E}_0$ . It follows from (13) that if  $\Delta_1 \cap \Delta_2 = \emptyset$  for  $\Delta_1$  and  $\Delta_2 \in \mathcal{E}$ , then

$$\tilde{Z}(\Delta_1 + \Delta_2) = \tilde{Z}(\Delta_1) + \tilde{Z}(\Delta_2) \quad (\text{P-a.s.})$$

and it follows from (12) that

$$\mathbb{E}|\tilde{Z}(\Delta)|^2 = m(\Delta), \quad \Delta \in \mathcal{E}.$$

Let us show that the random set function  $\tilde{Z}(\Delta)$ ,  $\Delta \in \mathcal{E}$ , is countably additive in the mean-square sense. In fact, let  $\Delta_k \in \mathcal{E}$  and  $\Delta = \sum_{k=1}^{\infty} \Delta_k$ . Then

$$\tilde{Z}(\Delta) - \sum_{k=1}^n \tilde{Z}(\Delta_k) = \mathcal{J}(g_n),$$

where

$$g_n(\lambda) = I_\Delta(\lambda) - \sum_{k=1}^n I_{\Delta_k}(\lambda) = I_{\Sigma_n}(\lambda), \quad \Sigma_n = \sum_{k=n+1}^{\infty} \Delta_k.$$

But

$$\mathbb{E}|\mathcal{J}(g_n)|^2 = \|g_n\|^2 = m(\Sigma_n) \downarrow 0, \quad n \rightarrow \infty,$$

i.e.

$$\mathbb{E}|\tilde{Z}(\Delta) - \sum_{k=1}^n \tilde{Z}(\Delta_k)|^2 \rightarrow 0, \quad n \rightarrow \infty.$$

It also follows from (11) that

$$\mathbb{E}\tilde{Z}(\Delta_1)\tilde{Z}(\Delta_2) = 0$$

when  $\Delta_1 \cap \Delta_2 = \emptyset$ ,  $\Delta_1, \Delta_2 \in \mathcal{E}$ .

Thus our function  $\tilde{Z}(\Delta)$ , defined on  $\Delta \in \mathcal{E}$ , is countably additive in the mean-square sense and coincides with  $Z(\Delta)$  on the sets  $\Delta \in \mathcal{E}_0$ . We shall call  $\tilde{Z}(\Delta)$ ,  $\Delta \in \mathcal{E}$ , an orthogonal stochastic measure (since it is an extension of the elementary orthogonal stochastic measure  $Z(\Delta)$ ) with respect to the structure function  $m(\Delta)$ ,  $\Delta \in \mathcal{E}$ ; and we call the integral  $\mathcal{J}(f) = \int_E f(\lambda)\tilde{Z}(d\lambda)$ , defined above, a stochastic integral with respect to this measure.

5. We now consider the case  $(E, \mathcal{E}) = (R, \mathcal{B}(R))$ , which is the most important for our purposes. As we know (§3, Chapter II), there is a one-to-one correspondence between finite measures  $m = m(\Delta)$  on  $(R, \mathcal{B}(R))$  and certain (generalized) distribution functions  $G = G(x)$ , with  $m(a, b] = G(b) - G(a)$ .

It turns out that there is something similar for orthogonal stochastic measures. We introduce the following definition.



**Definition 4.** A set of (complex-valued) random variables  $\{Z_\lambda\}$ ,  $\lambda \in R$ , defined on  $(\Omega, \mathcal{F}, P)$ , is a *random process with orthogonal increments* if

- (1)  $E|Z_\lambda|^2 < \infty$ ,  $\lambda \in R$ ;  
 (2) for every  $\lambda \in R$

$$E|Z_\lambda - Z_{\lambda_n}|^2 \rightarrow 0, \quad \lambda_n \downarrow \lambda, \quad \lambda_n \in R;$$

- (3) whenever  $\lambda_1 < \lambda_2 < \lambda_3 < \lambda_4$ ,

$$E(Z_{\lambda_4} - Z_{\lambda_3})(\overline{Z_{\lambda_2} - Z_{\lambda_1}}) = 0.$$

Condition (3) is the condition of orthogonal increments. Condition (1) means that  $Z_\lambda \in H^2$ . Finally, condition (2) is included for technical reasons; it is a requirement of continuity on the right (in the mean-square sense) at each  $\lambda \in R$ .

Let  $Z = Z(\Delta)$  be an orthogonal stochastic measure with respect to the structure function  $m = m(\Delta)$ , of finite measure, with the (generalized) distribution function  $G(\lambda)$ . Let us put

$$Z_\lambda = Z(-\infty, \lambda].$$

Then

$$E|Z_\lambda|^2 = m(-\infty, \lambda] = G(\lambda) < \infty, \quad E|Z_\lambda - Z_{\lambda_n}|^2 = m(\lambda_n, \lambda] \downarrow 0, \quad \lambda_n \downarrow \lambda,$$

and (evidently) 3) is satisfied also. Then this process  $\{Z_\lambda\}$  is called a process with orthogonal increments.

On the other hand, if  $\{Z_\lambda\}$  is such a process with  $E|Z_\lambda|^2 = G(\lambda)$ ,  $G(-\infty) = 0$ ,  $G(+\infty) < \infty$ , we put

$$Z(\Delta) = Z_b - Z_a$$

when  $\Delta = (a, b]$ . Let  $\mathcal{E}_0$  be the algebra of sets

$$\Delta = \sum_{k=1}^n (a_k, b_k] \quad \text{and} \quad Z(\Delta) = \sum_{k=1}^n Z(a_k, b_k].$$

It is clear that

$$E|Z(\Delta)|^2 = m(\Delta),$$

where  $m(\Delta) = \sum_{k=1}^n [G(b_k) - G(a_k)]$  and

$$EZ(\Delta_1)\overline{Z(\Delta_2)} = 0$$

for disjoint intervals  $\Delta_1 = (a_1, b_1]$  and  $\Delta_2 = (a_2, b_2]$ .

Therefore  $Z = Z(\Delta)$ ,  $\Delta \in \mathcal{E}_0$ , is an elementary stochastic measure with orthogonal values. The set function  $m = m(\Delta)$ ,  $\Delta \in \mathcal{E}_0$ , has a unique extension to a measure on  $\mathcal{E} = \mathcal{B}(R)$ , and it follows from the preceding constructions that  $Z = Z(\Delta)$ ,  $\Delta \in \mathcal{E}_0$ , can also be extended to the set  $\Delta \in \mathcal{E}$ , where  $\mathcal{E} = \mathcal{B}(R)$ , and  $E|Z(\Delta)|^2 = m(\Delta)$ ,  $\Delta \in \mathcal{B}(\mathcal{R})$ .

Therefore there is a one-to-one correspondence between processes  $\{Z_\lambda\}$ ,  $\lambda \in R$ , with orthogonal increments and  $E|Z_\lambda|^2 = G(\lambda)$ ,  $G(-\infty) = 0$ ,  $G(+\infty) < \infty$ , and orthogonal stochastic measures  $Z = Z(\Delta)$ ,  $\Delta \in \mathcal{B}(R)$ , with structure functions  $m = m(\Delta)$ . The correspondence is given by

$$Z_\lambda = Z(-\infty, \lambda], \quad G(\lambda) = m(-\infty, \lambda]$$

and

$$Z(a, b] = Z_b - Z_a, \quad m(a, b] = G(b) - G(a).$$

By analogy with the usual notation of the theory of Riemann-Stieltjes integration, the stochastic integral  $\int_R f(\lambda) dZ_\lambda$ , where  $\{Z_\lambda\}$  is a process with orthogonal increments, means the stochastic integral  $\int_R f(\lambda) Z(d\lambda)$  with respect to the corresponding process with an orthogonal stochastic measure.

## 6. PROBLEMS

1. Prove the equivalence of (5) and (6).
2. Let  $f \in L^2$ . Using the results of Chapter II (Theorem 1 of §4, the Corollary to Theorem 3 of §6, and Problem 9 of §3), prove that there is a sequence of functions  $f_n$  of the form (10) such that  $\|f - f_n\| \rightarrow 0$ ,  $n \rightarrow \infty$ .
3. Establish the following properties of an orthogonal stochastic measure  $Z(\Delta)$  with structure function  $m(\Delta)$ :

$$E|Z(\Delta_1) - Z(\Delta_2)|^2 = m(\Delta_1 \triangle \Delta_2),$$

$$Z(\Delta_1 \setminus \Delta_2) = Z(\Delta_1) - Z(\Delta_1 \cap \Delta_2) \quad (\text{P-a.s.}),$$

$$Z(\Delta_1 \triangle \Delta_2) = Z(\Delta_1) + Z(\Delta_2) - 2Z(\Delta_1 \cap \Delta_2) \quad (\text{P-a.s.}).$$

## §3. Spectral Representation of Stationary (Wide Sense) Sequences

1. If  $\xi = (\xi_n)$  is a stationary sequence with  $E\xi_n = 0$ ,  $n \in \mathbb{Z}$ , then by the theorem of §1, there is a finite measure  $F = F(\Delta)$  on  $([-\pi, \pi], \mathcal{B}([-\pi, \pi]))$  such that its covariance function  $R(n) = \text{cov}(\xi_{k+n}, \xi_k)$  admits the spectral representation

$$R(n) = \int_{-\pi}^{\pi} e^{i\lambda n} F(d\lambda). \quad (1)$$

The following result provides the corresponding spectral representation of the sequence  $\xi = (\xi_n)$ ,  $n \in \mathbb{Z}$ , itself.

**Theorem 1.** *There is an orthogonal stochastic measure  $Z = Z(\Delta)$ ,  $\Delta \in \mathcal{B}([-\pi, \pi])$ , such that for every  $n \in \mathbb{Z}$  (P-a.s.)*

$$\xi_n = \int_{-\pi}^{\pi} e^{i\lambda n} Z(d\lambda). \quad (2)$$

Moreover,  $E|Z(\Delta)|^2 = F(\Delta)$ .

The simplest proof is based on properties of Hilbert spaces.

Let  $L^2(F) = L^2(E, \mathcal{E}, F)$  be a Hilbert space of complex functions,  $E = [-\pi, \pi]$ ,  $\mathcal{E} = \mathcal{B}([-\pi, \pi])$ , with the scalar product

$$\langle f, g \rangle = \int_{-\pi}^{\pi} f(\lambda) \bar{g}(\lambda) F(d\lambda), \quad (3)$$

and let  $L_0^2(F)$  be the linear manifold ( $L_0^2(F) \subseteq L^2(F)$ ) spanned by  $e_n = e_n(\lambda)$ ,  $n \in \mathbb{Z}$ , where  $e_n(\lambda) = e^{i\lambda n}$ .

Observe that since  $E = [-\pi, \pi]$  and  $F$  is finite, the closure of  $L_0^2(F)$  coincides (Problem 1) with  $L^2(F)$ :

$$\overline{L_0^2(F)} = L^2(F).$$

Also let  $L_0^2(\xi)$  be the linear manifold spanned by the random variables  $\xi_n$ ,  $n \in \mathbb{Z}$ , and let  $L^2(\xi)$  be its closure in the mean-square sense (with respect to  $P$ ).

We establish a one-to-one correspondence between the elements of  $L_0^2(F)$  and  $L_0^2(\xi)$ , denoted by " $\leftrightarrow$ ", by setting

$$e_n \leftrightarrow \xi_n, \quad n \in \mathbb{Z}, \quad (4)$$

and defining it for elements in general (more precisely, for equivalence classes of elements) by linearity:

$$\sum \alpha_n e_n \leftrightarrow \sum \alpha_n \xi_n \quad (5)$$

(here we suppose that only finitely many of the complex numbers  $\alpha_n$  are different from zero).

Observe that (5) is a consistent definition, in the sense that  $\sum \alpha_n e_n = 0$  almost everywhere with respect to  $F$  if and only if  $\sum \alpha_n \xi_n = 0$  (P-a.s.).

The correspondence " $\leftrightarrow$ " is an *isometry*, i.e. it preserves scalar products. In fact, by (3),

$$\begin{aligned} \langle e_n, e_m \rangle &= \int_{-\pi}^{\pi} e_n(\lambda) \bar{e}_m(\lambda) F(d\lambda) = \int_{-\pi}^{\pi} e^{i\lambda(n-m)} F(d\lambda) = R(n-m) \\ &= E \xi_n \bar{\xi}_m = (\xi_n, \xi_m) \end{aligned}$$

and similarly

$$\langle \sum \alpha_n e_n, \sum \beta_n e_n \rangle = (\sum \alpha_n \xi_n, \sum \beta_n \xi_n). \quad (6)$$

Now let  $\eta \in L^2(\xi)$ . Since  $L^2(\xi) = \bar{L}_0^2(\xi)$ , there is a sequence  $\{\eta_n\}$  such that  $\eta_n \in L_0^2(\xi)$  and  $\|\eta_n - \eta\| \rightarrow 0$ ,  $n \rightarrow \infty$ . Consequently  $\{\eta_n\}$  is a fundamental sequence and therefore so is the sequence  $\{f_n\}$ , where  $f_n \in L_0^2(F)$  and  $f_n \leftrightarrow \eta_n$ . The space  $L^2(F)$  is complete and consequently there is an  $f \in L^2(F)$  such that  $\|f_n - f\| \rightarrow 0$ .

There is an evident converse: if  $f \in L^2(F)$  and  $\|f - f_n\| \rightarrow 0$ ,  $f_n \in L_0^2(F)$ , there is an element  $\eta$  of  $L^2(\xi)$  such that  $\|\eta - \eta_n\| \rightarrow 0$ ,  $\eta_n \in L_0^2(\xi)$  and  $\eta_n \leftrightarrow f_n$ .

Up to now the isometry " $\leftrightarrow$ " has been defined only as between elements of  $L_0^2(\xi)$  and  $L_0^2(F)$ . We extend it by continuity, taking  $f \leftrightarrow \eta$  when  $f$  and  $\eta$  are the elements considered above. It is easily verified that the correspondence obtained in this way is one-to-one (between classes of equivalent random variables and of functions), is linear, and preserves scalar products.

Consider the function  $f(\lambda) = I_\Delta(\lambda)$ , where  $\Delta \in \mathcal{B}([-\pi, \pi])$ , and let  $Z(\Delta)$  be the element of  $L^2(\xi)$  such that  $I_\Delta(\lambda) \leftrightarrow Z(\lambda)$ . It is clear that  $\|I_\Delta(\lambda)\|^2 = F(\Delta)$  and therefore  $E|Z(\Delta)|^2 = F(\Delta)$ . Moreover, if  $\Delta_1 \cap \Delta_2 = \emptyset$ , we have  $EZ(\Delta_1)Z(\Delta_2) = 0$  and  $E|Z(\Delta) - \sum_{k=1}^n Z(\Delta_k)|^2 \rightarrow 0$ ,  $n \rightarrow \infty$ , where  $\Delta = \sum_{k=1}^\infty \Delta_k$ .

Hence the family of elements  $Z(\Delta)$ ,  $\Delta \in \mathcal{B}([-\pi, \pi])$ , form an orthogonal stochastic measure, with respect to which (according to §2) we can define the stochastic integral

$$\mathcal{J}(f) = \int_{-\pi}^{\pi} f(\lambda)Z(d\lambda), \quad f \in L^2(F).$$

Let  $f \in L^2(F)$  and  $\eta \leftrightarrow f$ . Denote the element  $\eta$  by  $\Phi(f)$  (more precisely, select single representatives from the corresponding equivalence classes of random variables or functions). Let us show that (P-a.s.)

$$\mathcal{J}(f) = \Phi(f). \quad (7)$$

In fact, if

$$f(\lambda) = \sum \alpha_k I_{\Delta_k}(\lambda) \quad (8)$$

is a finite linear combination of functions  $I_{\Delta_k}(\lambda)$ ,  $\Delta_k = (a_k, b_k]$ , then, by the very definition of the stochastic integral,  $\mathcal{J}(f) = \sum \alpha_k Z(\Delta_k)$ , which is evidently equal to  $\Phi(f)$ . Therefore (7) is valid for functions of the form (8). But if  $f \in L^2(F)$  and  $\|f_n - f\| \rightarrow 0$ , where  $f_n$  are functions of the form (8), then  $\|\Phi(f_n) - \Phi(f)\| \rightarrow 0$  and  $\|\mathcal{J}(f_n) - \mathcal{J}(f)\| \rightarrow 0$  (by (2.14)). Therefore  $\Phi(f) = \mathcal{J}(f)$  (P-a.s.).

Consider the function  $f(\lambda) = e^{i\lambda n}$ . Then  $\Phi(e^{i\lambda n}) = \xi_n$  by (4), but on the other hand  $\mathcal{J}(e^{i\lambda n}) = \int_{-\pi}^{\pi} e^{i\lambda n} Z(d\lambda)$ . Therefore

$$\xi_n = \int_{-\pi}^{\pi} e^{i\lambda n} Z(d\lambda), \quad n \in \mathbb{Z} \quad (\text{P-a.s.})$$

by (7). This completes the proof of the theorem.

**Corollary 1.** Let  $\xi = (\xi_n)$  be a stationary sequence of real random variables  $\xi_n$ ,  $n \in \mathbb{Z}$ . Then the stochastic measure  $Z = Z(\Delta)$  involved in the spectral representation (2) has the property that

$$Z(\Delta) = \overline{Z(-\Delta)} \quad (9)$$

for every  $\Delta = \mathcal{B}([- \pi, \pi])$ , where  $-\Delta = \{\lambda: -\lambda \in \Delta\}$ .

In fact, let  $f(\lambda) = \sum \alpha_k e^{i\lambda k}$  and  $\eta = \sum \alpha_k \xi_k$  (finite sums). Then  $f \leftrightarrow \eta$  and therefore

$$\bar{\eta} = \sum \bar{\alpha}_k \xi_k \leftrightarrow \sum \bar{\alpha}_k e^{i\lambda k} = \overline{f(-\lambda)}. \quad (10)$$

Since  $\mathcal{J}_\Delta(\lambda) \leftrightarrow Z(\Delta)$ , it follows from (10) that either  $\mathcal{J}_\Delta(-\lambda) \leftrightarrow \bar{Z}(\Delta)$  or  $\mathcal{J}_{-\Delta}(\lambda) \leftrightarrow \bar{Z}(\Delta)$ . On the other hand,  $\mathcal{J}_{-\Delta}(\lambda) \leftrightarrow Z(-\Delta)$ . Therefore  $\bar{Z}(\Delta) = Z(-\Delta)$  (P-a.s.).

**Corollary 2.** Again let  $\xi = (\xi_n)$  be a stationary sequence of real random variables  $\xi_n$  and  $Z(\Delta) = Z_1(\Delta) + iZ_2(\Delta)$ . Then

$$EZ_1(\Delta_1)Z_2(\Delta_2) = 0 \quad (11)$$

for every  $\Delta_1$  and  $\Delta_2$ ; and if  $\Delta_1 \cap \Delta_2 = \emptyset$  then

$$EZ_1(\Delta_1)Z_1(\Delta_2) = 0, \quad EZ_2(\Delta_1)Z_2(\Delta_2) = 0. \quad (12)$$

In fact, since  $Z(\Delta) = \bar{Z}(-\Delta)$ , we have

$$Z_1(-\Delta) = Z_1(\Delta), \quad Z_2(-\Delta) = -Z_2(\Delta). \quad (13)$$

Moreover, since  $EZ(\Delta_1)\bar{Z}(\Delta_2) = E|Z(\Delta_1 \cap \Delta_2)|^2$ , we have  $\text{Im } EZ(\Delta_1)\bar{Z}(\Delta_2) = 0$ , i.e.

$$EZ_1(\Delta_1)Z_2(\Delta_2) + EZ_2(\Delta_1)Z_1(\Delta_2) = 0. \quad (14)$$

If we take the interval  $-\Delta_1$  instead of  $\Delta_1$  we therefore obtain

$$EZ_1(-\Delta_1)Z_2(\Delta_2) + EZ_2(-\Delta_1)Z_1(\Delta_2) = 0,$$

which, by (13), can be transformed into

$$EZ_1(\Delta_1)Z_2(\Delta_2) - EZ_2(\Delta_1)Z_1(\Delta_2) = 0. \quad (15)$$

Then (11) follows from (14) and (15).

On the other hand, if  $\Delta_1 \cap \Delta_2 = \emptyset$  then  $EZ(\Delta_1)\bar{Z}(\Delta_2) = 0$ , whence  $\text{Re } EZ(\Delta_1)\bar{Z}(\Delta_2) = 0$  and  $\text{Re } EZ(-\Delta_1)\bar{Z}(\Delta_2) = 0$ , which, with (13), provides an evident proof of (12).

**Corollary 3.** Let  $\xi = (\xi_n)$  be a Gaussian sequence. Then, for every family  $\Delta_1, \dots, \Delta_k$ , the vector  $(Z_1(\Delta_1), \dots, Z_1(\Delta_k), Z_2(\Delta_1), \dots, Z_2(\Delta_k))$  is normally distributed.

In fact, the linear manifold  $L_0^2(\xi)$  consists of (complex-valued) Gaussian random variables  $\eta$ , i.e. the vector  $(\operatorname{Re} \eta, \operatorname{Im} \eta)$  has a Gaussian distribution. Then, according to Subsection 5, §13, Chapter II, the closure of  $L_0^2(\xi)$  also consists of Gaussian variables. It follows from Corollary 2 that, when  $\xi = (\xi_n)$  is a Gaussian sequence, the real and imaginary parts of  $Z_1$  and  $Z_2$  are independent in the sense that the families of random variables  $(Z_1(\Delta_1), \dots, Z_1(\Delta_k))$  and  $(Z_2(\Delta_1), \dots, Z_2(\Delta_k))$  are independent. It also follows from (12) that when the sets  $\Delta_1, \dots, \Delta_k$  are disjoint, the random variables  $Z_i(\Delta_1), \dots, Z_i(\Delta_k)$  are collectively independent,  $i = 1, 2$ .

**Corollary 4.** *If  $\xi = (\xi_n)$  is a stationary sequence of real random variables, then (P-a.s.)*

$$\xi_n = \int_{-\pi}^{\pi} \cos \lambda n Z_1(d\lambda) + \int_{-\pi}^{\pi} \sin \lambda n Z_2(d\lambda). \quad (16)$$

**Remark.** If  $\{Z_\lambda\}$ ,  $\lambda \in [-\pi, \pi)$ , is a process with orthogonal increments, corresponding to an orthogonal stochastic measure  $Z = Z(\Delta)$ , then in accordance with §2 the spectral representation (2) can also be written in the following form:

$$\xi_n = \int_{-\pi}^{\pi} e^{i\lambda n} dZ_\lambda, \quad n \in \mathbb{Z}. \quad (17)$$

2. Let  $\xi = (\xi_n)$  be a stationary sequence with the spectral representation (2) and let  $\eta \in L^2(\xi)$ . The following theorem describes the structure of such random variables.

**Theorem 2.** *If  $\eta \in L^2(\xi)$ , there is a function  $\varphi \in L^2(F)$  such that (P-a.s.)*

$$\eta = \int_{-\pi}^{\pi} \varphi(\lambda) Z(d\lambda). \quad (18)$$

**PROOF.** If

$$\eta_n = \sum_{|k| \leq n} \alpha_k \xi_k, \quad (19)$$

then by (2)

$$\eta_n = \int_{-\pi}^{\pi} \left( \sum_{|k| \leq n} \alpha_k e^{i\lambda k} \right) Z(d\lambda), \quad (20)$$

i.e. (18) is satisfied with

$$\varphi_n(\lambda) = \sum_{|k| \leq n} \alpha_k e^{i\lambda k}. \quad (21)$$

In the general case, when  $\eta \in L^2(\xi)$ , there are variables  $\eta_n$  of type (19) such that  $\|\eta - \eta_n\| \rightarrow 0, n \rightarrow \infty$ . But then  $\|\varphi_n - \varphi_m\| = \|\eta_n - \eta_m\| \rightarrow 0, n, m \rightarrow \infty$ .

Consequently  $\{\varphi_n\}$  is fundamental in  $L^2(F)$  and therefore there is a function  $\varphi \in L^2(F)$  such that  $\|\varphi - \varphi_n\| \rightarrow 0, n \rightarrow \infty$ .

By property (2.14) we have  $\|\mathcal{J}(\varphi_n) - \mathcal{J}(\varphi)\| \rightarrow 0$ , and since  $\eta_n = \mathcal{J}(\varphi_n)$  we also have  $\eta = \mathcal{J}(\varphi)$  (P-a.s.).

This completes the proof of the theorem.

**Remark.** Let  $H_0(\xi)$  and  $H_0(F)$  be the respective closed linear manifolds spanned by the variables  $\xi_n$  and by the functions  $e_n$  when  $n \leq 0$ . Then if  $\eta \in H_0(\xi)$  there is a function  $\varphi \in H_0(F)$  such that (P-a.s.)  $\eta = \int_{-\infty}^{\infty} \varphi(\lambda) Z(d\lambda)$ .

3. Formula (18) describes the structure of the random variables that are obtained from  $\xi_n, n \in \mathbb{Z}$ , by linear transformations, i.e. in the form of finite sums (19) and their mean-square limits.

A special but important class of such linear transformations are defined by means of what are known as (linear) *filters*. Let us suppose that, at instant  $m$ , a system (filter) receives as input a signal  $x_m$ , and that the output of the system is, at instant  $n$ , the signal  $h(n-m)x_m$ , where  $h = h(s), s \in \mathbb{Z}$ , is a complex valued function called the *impulse response* (of the filter).

Therefore the total signal obtained from the input can be represented in the form

$$y_n = \sum_{m=-\infty}^{\infty} h(n-m)x_m. \quad (22)$$

For physically realizable systems, the values of the input at instant  $n$  are determined only by the "past" values of the signal, i.e. the values  $x_m$  for  $m \leq n$ . It is therefore natural to call a filter with the impulse response  $h(s)$  *physically realizable* if  $h(s) = 0$  for all  $s < 0$ , in other words if

$$y_n = \sum_{m=-\infty}^n h(n-m)x_m = \sum_{m=0}^{\infty} h(m)x_{n-m}. \quad (23)$$

An important *spectral characteristic* of a filter with the impulse response  $h$  is its Fourier transform

$$\varphi(\lambda) = \sum_{m=-\infty}^{\infty} e^{-i\lambda m} h(m), \quad (24)$$

known as the *frequency characteristic* or *transfer function* of the filter.

Let us now take up conditions, about which nothing has been said so far, for the convergence of the series in (22) and (24). Let us suppose that the input is a stationary random sequence  $\xi = (\xi_n), n \in \mathbb{Z}$ , with covariance function  $R(n)$  and spectral decomposition (2). Then if

$$\sum_{k, l=-\infty}^{\infty} h(k) R(k-l) \bar{h}(l) < \infty, \quad (25)$$

the series  $\sum_{m=-\infty}^{\infty} h(n-m)\xi_m$  converges in mean-square and therefore there is a stationary sequence  $\eta = (\eta_n)$  with

$$\eta_n = \sum_{m=-\infty}^{\infty} h(n-m)\xi_m = \sum_{m=-\infty}^{\infty} h(m)\xi_{n-m}. \quad (26)$$

In terms of the spectral measure, (25) is evidently equivalent to saying that  $\varphi(\lambda) \in L^2(F)$ , i.e.

$$\int_{-\pi}^{\pi} |\varphi(\lambda)|^2 F(d\lambda) < \infty. \quad (27)$$

Under (25) or (27), we obtain the spectral representation

$$\eta_n = \int_{-\pi}^{\pi} e^{i\lambda n} \varphi(\lambda) Z(d\lambda). \quad (28)$$

of  $\eta$  from (26) and (2). Consequently the covariance function  $R_{\eta}(n)$  of  $\eta$  is given by the formula

$$R_{\eta}(n) = \int_{-\pi}^{\pi} e^{i\lambda n} |\varphi(\lambda)|^2 F(d\lambda). \quad (29)$$

In particular, if the input to a filter with frequency characteristic  $\varphi = \varphi(\lambda)$  is taken to be white noise  $\varepsilon = (\varepsilon_n)$ , the output will be a stationary sequence (moving average)

$$\eta_n = \sum_{m=-\infty}^{\infty} h(m)\varepsilon_{n-m} \quad (30)$$

with spectral density

$$f_{\eta}(\lambda) = \frac{1}{2\pi} |\varphi(\lambda)|^2.$$

The following theorem shows that there is a sense in which every stationary sequence with a spectral density is obtainable by means of a moving average.

**Theorem 3.** *Let  $\eta = (\eta_n)$  be a stationary sequence with spectral density  $f_{\eta}(\lambda)$ . Then (possibly at the expense of enlarging the original probability space) we can find a sequence  $\varepsilon = (\varepsilon_n)$  representing white noise, and a filter, such that the representation (30) holds.*

**PROOF.** For a given (nonnegative) function  $f_{\eta}(\lambda)$  we can find a function  $\varphi(\lambda)$  such that  $f_{\eta}(\lambda) = (1/2\pi) |\varphi(\lambda)|^2$ . Since  $\int_{-\pi}^{\pi} f_{\eta}(\lambda) d\lambda < \infty$ , we have  $\varphi(\lambda) \in L^2(\mu)$ , where  $\mu$  is Lebesgue measure on  $[-\pi, \pi)$ . Hence  $\varphi$  can be represented as a Fourier series (24) with  $h(m) = (1/2\pi) \int_{-\pi}^{\pi} e^{im\lambda} \varphi(\lambda) d\lambda$ , where convergence is understood in the sense that

$$\int_{-\pi}^{\pi} \left| \varphi(\lambda) - \sum_{|m| \leq n} e^{-i\lambda m} h(m) \right|^2 d\lambda \rightarrow 0, \quad n \rightarrow \infty.$$



Let

$$\eta_n = \int_{-\pi}^{\pi} e^{i\lambda n} Z(d\lambda), \quad n \in \mathbb{Z}.$$

Besides the measure  $Z = Z(\Delta)$  we introduce another independent orthogonal stochastic measure  $\bar{Z} = \bar{Z}(\Delta)$  with  $E|\bar{Z}(a, b)|^2 = (b - a)/2\pi$ . (The possibility of constructing such a measure depends, in general, on having a sufficiently "rich" original probability space.) Let us put

$$\bar{Z}(\Delta) = \int_{\Delta} \varphi^{\oplus}(\lambda) Z(d\lambda) + \int_{\Delta} [1 - \varphi^{\oplus}(\lambda)\varphi(\lambda)] \bar{Z}(d\lambda),$$

where

$$a^{\oplus} = \begin{cases} a^{-1}, & \text{if } a \neq 0, \\ 0, & \text{if } a = 0. \end{cases}$$

The stochastic measure  $\bar{Z} = \bar{Z}(\Delta)$  is a measure with orthogonal values, and for every  $\Delta = (a, b]$  we have

$$E|\bar{Z}(\Delta)|^2 = \frac{1}{2\pi} \int_{\Delta} |\varphi^{\oplus}(\lambda)|^2 |\varphi(\lambda)|^2 d\lambda + \frac{1}{2\pi} \int_{\Delta} |1 - \varphi^{\oplus}(\lambda)\varphi(\lambda)|^2 d\lambda = \frac{|\Delta|}{2\pi},$$

where  $|\Delta| = b - a$ . Therefore the stationary sequence  $\varepsilon = (\varepsilon_n)$ ,  $n \in \mathbb{Z}$ , with

$$\varepsilon_n = \int_{-\pi}^{\pi} e^{i\lambda n} \bar{Z}(d\lambda),$$

is a white noise.

We now observe that

$$\int_{-\pi}^{\pi} e^{i\lambda n} \varphi(\lambda) \bar{Z}(d\lambda) = \int_{-\pi}^{\pi} e^{i\lambda n} Z(d\lambda) = \eta_n \quad (31)$$

and, on the other hand, by property (2.14) (P-a.s.)

$$\begin{aligned} \int_{-\pi}^{\pi} e^{i\lambda n} \varphi(\lambda) \bar{Z}(d\lambda) &= \int_{-\pi}^{\pi} e^{i\lambda n} \left( \sum_{m=-\infty}^{\infty} e^{-i\lambda m} h(m) \right) \bar{Z}(d\lambda) \\ &= \sum_{m=-\infty}^{\infty} h(m) \int_{-\pi}^{\pi} e^{i\lambda(n-m)} \bar{Z}(d\lambda) = \sum_{m=-\infty}^{\infty} h(m) \varepsilon_{n-m}, \end{aligned}$$

which, together with (31), establishes the representation (30).

This completes the proof of the theorem.

**Remark.** If  $f_{\eta}(\lambda) > 0$  (almost everywhere with respect to Lebesgue measure), the introduction of the auxiliary measure  $\bar{Z} = \bar{Z}(\Delta)$  becomes unnecessary (since then  $1 - \varphi^{\oplus}(\lambda)\varphi(\lambda) = 0$  almost everywhere with respect to Lebesgue measure), and the reservation concerning the necessity of extending the original probability space can be omitted.

**Corollary 1.** Let the spectral density  $f_\eta(\lambda) > 0$  (almost everywhere with respect to Lebesgue measure) and

$$f_\eta(\lambda) = \frac{1}{2\pi} |\varphi(\lambda)|^2,$$

where

$$\varphi(\lambda) = \sum_{k=0}^{\infty} e^{-i\lambda k} h(k), \quad \sum_{k=0}^{\infty} |h(k)|^2 < \infty.$$

Then the sequence  $\eta$  admits a representation as a one-sided moving average,

$$\eta_n = \sum_{m=0}^{\infty} h(m) \varepsilon_{n-m}.$$

In particular, let  $P(z) = a_0 + a_1 z + \cdots + a_p z^p$  be a polynomial that has no zeros on  $\{z: |z| = 1\}$ . Then the sequence  $\eta = (\eta_n)$  with spectral density

$$f_\eta(\lambda) = \frac{1}{2\pi} |P(e^{-i\lambda})|^2$$

can be represented in the form

$$\eta_n = a_0 \varepsilon_n + a_1 \varepsilon_{n-1} + \cdots + a_p \varepsilon_{n-p}.$$

**Corollary 2.** Let  $\xi = (\xi_n)$  be a sequence with rational spectral density

$$f_\xi(\lambda) = \frac{1}{2\pi} \left| \frac{P(e^{-i\lambda})}{Q(e^{-i\lambda})} \right|^2, \quad (32)$$

where  $P(z) = a_0 + a_1 z + \cdots + a_p z^p$ ,  $Q(z) = 1 + b_1 z + \cdots + b_q z^q$ .

Let us show that if  $P(z)$  and  $Q(z)$  have no zeros on  $\{z: |z| = 1\}$ , there is a white noise  $\varepsilon = \varepsilon(n)$  such that (P-a.s.)

$$\xi_n + b_1 \xi_{n-1} + \cdots + b_q \xi_{n-q} = a_0 \varepsilon_n + a_1 \varepsilon_{n-1} + \cdots + a_p \varepsilon_{n-p}. \quad (33)$$

Conversely, every stationary sequence  $\xi = (\xi_n)$  that satisfies this equation with some white noise  $\varepsilon = (\varepsilon_n)$  and some polynomial  $Q(z)$  with no zeros on  $\{z: |z| = 1\}$  has a spectral density (32).

In fact, let  $\eta_n = \xi_n + b_1 \xi_{n-1} + \cdots + b_q \xi_{n-q}$ . Then  $f_\eta(\lambda) = (1/2\pi) |P(e^{-i\lambda})|^2$  and the required representation follows from Corollary 1.

On the other hand, if (33) holds and  $F_\xi(\lambda)$  and  $F_\eta(\lambda)$  are the spectral functions of  $\xi$  and  $\eta$ , then

$$F_\eta(\lambda) = \int_{-\pi}^{\lambda} |Q(e^{-i\nu})|^2 dF_\xi(\nu) = \frac{1}{2\pi} \int_{-\pi}^{\lambda} |P(e^{-i\nu})|^2 d\nu.$$

Since  $|Q(e^{-i\nu})|^2 > 0$ , it follows that  $F_\xi(\lambda)$  has a density defined by (32).

4. The following mean-square ergodic theorem can be thought of as an analog of the law of large numbers for stationary (wide sense) random sequences.

**Theorem 4.** Let  $\xi = (\xi_n)$ ,  $n \in \mathbb{Z}$ , be a stationary sequence with  $E\xi_n = 0$ , covariance function (1), and spectral resolution (2). Then

$$\frac{1}{n} \sum_{k=0}^{n-1} \xi_k \xrightarrow{L^2} Z(\{0\}) \quad (34)$$

and

$$\frac{1}{n} \sum_{k=0}^{n-1} R(k) \rightarrow F(\{0\}). \quad (35)$$

**PROOF.** By (2),

$$\frac{1}{n} \sum_{k=0}^{n-1} \xi_k = \int_{-\pi}^{\pi} \frac{1}{n} \sum_{k=0}^{n-1} e^{ik\lambda} Z(d\lambda) = \int_{-\pi}^{\pi} \varphi_n(\lambda) Z(d\lambda),$$

where

$$\varphi_n(\lambda) = \frac{1}{n} \sum_{k=0}^{n-1} e^{ik\lambda}. \quad (36)$$

It is clear that

$$|\varphi_n(\lambda)| \leq 1.$$

Moreover,  $\varphi_n(\lambda) \xrightarrow{L^2(F)} I_{\{0\}}(\lambda)$  and therefore by (2.14)

$$\int_{-\pi}^{\pi} \varphi_n(\lambda) Z(d\lambda) \xrightarrow{L^2} \int_{-\pi}^{\pi} I_{\{0\}}(\lambda) Z(d\lambda) = Z(\{0\}),$$

which establishes (34).

Relation (35) can be proved in a similar way.

This completes the proof of the theorem.

**Corollary.** If the spectral function is continuous at zero, i.e.  $F(\{0\}) = 0$ , then  $Z(\{0\}) = 0$  (P-a.s.) and by (34) and (35),

$$\frac{1}{n} \sum_{k=0}^{n-1} R(k) \rightarrow 0 \Rightarrow \frac{1}{n} \sum_{k=0}^{n-1} \xi_k \xrightarrow{L^2} 0.$$

Since

$$\left| \frac{1}{n} \sum_{k=0}^{n-1} R(k) \right|^2 = \left| E \left( \frac{1}{n} \sum_{k=0}^{n-1} \xi_k \right) \xi_0 \right|^2 \leq E|\xi_0|^2 E \left| \frac{1}{n} \sum_{k=0}^{n-1} \xi_k \right|^2,$$

the converse implication also holds:

$$\frac{1}{n} \sum_{k=0}^{n-1} \xi_k \xrightarrow{L^2} 0 \Rightarrow \frac{1}{n} \sum_{k=0}^{n-1} R(k) \rightarrow 0.$$

Therefore the condition  $(1/n) \sum_{k=0}^{n-1} R(k) \rightarrow 0$  is necessary and sufficient for the convergence (in the mean-square sense) of the arithmetic means  $(1/n) \sum_{k=0}^{n-1} \xi_k$  to zero. It follows that if the original sequences  $\xi = (\xi_n)$  has expectation  $m$  (that is,  $E\xi_0 = m$ ), then

$$\frac{1}{n} \sum_{k=0}^{n-1} R(k) \rightarrow 0 \Leftrightarrow \frac{1}{n} \sum_{k=0}^{n-1} \xi_k \xrightarrow{L^2} m, \quad (37)$$

where  $R(n) = E(\xi_n - E\xi_n)(\overline{\xi_0 - E\xi_0})$ .

Let us also observe that if  $Z(\{0\}) \neq 0$  (P-a.s.) and  $m = 0$ , then  $\xi_n$  "contains a random constant  $\alpha$ ":

$$\xi_n = \alpha + \eta_n,$$

where  $\alpha = Z(\{0\})$ ; and in the spectral representation  $\eta_n = \int_{-\pi}^{\pi} e^{i\lambda n} Z_{\eta}(d\lambda)$  the measure  $Z_{\eta} = Z_{\eta}(\Delta)$  is such that  $Z_{\eta}(\{0\}) = 0$  (P-a.s.). Conclusion (34) means that the arithmetic mean converges in mean-square to precisely this random constant  $\alpha$ .

## 5. PROBLEMS

1. Show that  $\overline{L_0^2}(F) = L^2(F)$  (for the notation see the proof of Theorem 1).
2. Let  $\xi = (\xi_n)$  be a stationary sequence with the property that  $\xi_{n+N} = \xi_n$  for some  $N$  and all  $n$ . Show that the spectral representation of such a sequence reduces to (1.13).
3. Let  $\xi = (\xi_n)$  be a stationary sequence such that  $E\xi_n = 0$  and

$$\frac{1}{N^2} \sum_{k=0}^N \sum_{l=0}^N R(k-l) = \frac{1}{N} \sum_{|k| \leq N-1} R(k) \left[ 1 - \frac{|k|}{N} \right] \leq CN^{-\alpha}$$

for some  $C > 0, \alpha > 0$ . Use the Borel-Cantelli lemma to show that then

$$\frac{1}{N} \sum_{k=0}^N \xi_k \rightarrow 0 \quad (\text{P-a.s.})$$

4. Let the spectral density  $f_{\xi}(\lambda)$  of the sequence  $\xi = (\xi_n)$  be rational,

$$f_{\xi}(\lambda) = \frac{1}{2\pi} \frac{|P_{n-1}(e^{-i\lambda})|}{|Q_n(e^{-i\lambda})|}, \quad (38)$$

where  $P_{n-1}(z) = a_0 + a_1 z + \dots + a_{n-1} z^{n-1}$  and  $Q_n(z) = 1 + b_1 z + \dots + b_n z^n$ , and all the zeros of these polynomials lie outside the unit disk.

Show that there is a white noise  $\varepsilon = (\varepsilon_m)$ ,  $m \in \mathbb{Z}$ , such that the sequence  $(\xi_n)$  is a component of an  $n$ -dimensional sequence  $(\xi_m^1, \xi_m^2, \dots, \xi_m^n)$ ,  $\xi_m^1 = \xi_m$ , that satisfies

the system of equations

$$\begin{aligned}\xi_{m+1}^i &= \xi_m^{i+1} + \beta_i \varepsilon_{m+1}, \quad i = 1, \dots, n-1, \\ \xi_{m+1}^n &= -\sum_{j=0}^{n-1} b_{n-j} \xi_m^{j+1} + \beta_n \varepsilon_{m+1},\end{aligned}\quad (39)$$

where  $\beta_1 = a_0$ ,  $\beta_i = a_{i-1} - \sum_{k=1}^{i-1} \beta_k b_{i-k}$ .

## §4. Statistical Estimation of the Covariance Function and the Spectral Density

1. Problems of the statistical estimation of various characteristics of the probability distributions of random sequences arise in the most diverse branches of science (geophysics, medicine, economics, etc.) The material presented in this section will give the reader an idea of the concepts and methods of estimation, and of the difficulties that are encountered.

To begin with, let  $\xi = (\xi_n)$ ,  $n \in \mathbb{Z}$ , be a sequence, stationary in the wide sense (for simplicity, real) with expectation  $E\xi_n = m$  and covariance  $R(n) = \int_{-\pi}^{\pi} e^{i\lambda n} F(d\lambda)$ .

Let  $x_0, x_1, \dots, x_{N-1}$  be the results of observing the random variables  $\xi_0, \xi_1, \dots, \xi_{N-1}$ . How are we then to construct a "good" estimator of the (unknown) mean value  $m$ ?

Let us put

$$m_N(x) = \frac{1}{N} \sum_{k=0}^{N-1} x_k. \quad (1)$$

Then it follows from the elementary properties of the expectation that this is a "good" estimator of  $m$  in the sense that it is *unbiased* "in the mean over all kinds of data  $x_0, \dots, x_{N-1}$ ", i.e.

$$Em_N(\xi) = E\left(\frac{1}{N} \sum_{k=0}^{N-1} \xi_k\right) = m. \quad (2)$$

In addition, it follows from Theorem 4 of §3 that when  $(1/N) \sum_{k=0}^N R(k) \rightarrow 0$ ,  $N \rightarrow \infty$ , our estimator is *consistent* (in mean-square), i.e.

$$E|m_N(\xi) - m|^2 \rightarrow 0, \quad N \rightarrow \infty. \quad (3)$$

Next we take up the problem of estimating the covariance function  $R(n)$ , the spectral function  $F(\lambda) = F([-\pi, \lambda])$ , and the spectral density  $f(\lambda)$ , all under the assumption that  $m = 0$ .

Since  $R(n) = E\xi_{n+k}\xi_k$ , it is natural to estimate this function on the basis of  $N$  observations  $x_0, x_1, \dots, x_{N-1}$  (when  $0 \leq n < N$ ) by

$$\hat{R}_N(n; x) = \frac{1}{N-n} \sum_{k=0}^{N-n-1} x_{n+k} x_k.$$

It is clear that this estimator is unbiased in the sense that

$$E\hat{R}_N(n; \xi) = R(n), \quad 0 \leq n < N.$$

Let us now consider the question of its consistency. If we replace  $\xi_k$  in (3.37) by  $\xi_{n+k}\xi_k$  and suppose that the sequence  $\xi = (\xi_n)$  under consideration has a fourth moment ( $E\xi_0^4 < \infty$ ), we find that the condition

$$\frac{1}{N} \sum_{k=0}^{N-1} E[\xi_{n+k}\xi_k - R(n)][\xi_n\xi_0 - R(n)] \rightarrow 0, \quad N \rightarrow \infty, \quad (4)$$

is necessary and sufficient for

$$E|\hat{R}_N(n; \xi) - R(n)|^2 \rightarrow 0, \quad N \rightarrow \infty. \quad (5)$$

Let us suppose that the original sequence  $\xi = (\xi_n)$  is Gaussian (with zero mean and covariance  $R(n)$ ). Then by (II.12.51)

$$\begin{aligned} E[\xi_{n+k}\xi_k - R(n)][\xi_n\xi_0 - R(n)] &= E\xi_{n+k}\xi_k\xi_n\xi_0 - R^2(n) \\ &= E\xi_{n+k}\xi_k \cdot E\xi_n\xi_0 + E\xi_{n+k}\xi_n \cdot E\xi_k\xi_0 \\ &\quad + E\xi_{n+k}\xi_0 \cdot E\xi_k\xi_n - R^2(n) \\ &= R^2(k) + R(n+k)R(n-k). \end{aligned}$$

Therefore in the Gaussian case condition (4) is equivalent to

$$\frac{1}{N} \sum_{k=0}^{N-1} [R^2(k) + R(n+k)R(n-k)] \rightarrow 0, \quad N \rightarrow \infty. \quad (6)$$

Since  $|R(n+k)R(n-k)| \leq |R(n+k)|^2 + |R(n-k)|^2$ , the condition

$$\frac{1}{N} \sum_{k=0}^{N-1} R^2(k) \rightarrow 0, \quad N \rightarrow \infty, \quad (7)$$

implies (6). Conversely, if (6) holds for  $n = 0$ , then (7) is satisfied.

We have now established the following theorem.

**Theorem.** Let  $\xi = (\xi_n)$  be a Gaussian stationary sequence with  $E\xi_n = 0$  and covariance function  $R(n)$ . Then (7) is a necessary and sufficient condition that, for every  $n \geq 0$ , the estimator  $\hat{R}_N(n; x)$  is mean-square consistent, (i.e. that (5) is satisfied).

**Remark.** If we use the spectral representation of the covariance function, we obtain

$$\begin{aligned}\frac{1}{N} \sum_{k=0}^{N-1} R^2(k) &= \int_{-\pi}^{\pi} \int_{-\pi}^{\pi} \frac{1}{N} \sum_{k=0}^{N-1} e^{i(\lambda-\nu)k} F(d\lambda) F(d\nu) \\ &= \int_{-\pi}^{\pi} \int_{-\pi}^{\pi} f_N(\lambda, \nu) F(d\lambda) F(d\nu),\end{aligned}$$

where (compare (3.35))

$$f_N(\lambda, \nu) = \begin{cases} 1, & \lambda = \nu, \\ \frac{1 - e^{i(\lambda-\nu)N}}{N[1 - e^{i(\lambda-\nu)}]}, & \lambda \neq \nu. \end{cases}$$

But as  $N \rightarrow \infty$

$$f_N(\lambda, \nu) \rightarrow f(\lambda, \nu) = \begin{cases} 1, & \lambda = \nu, \\ 0, & \lambda \neq \nu. \end{cases}$$

Therefore

$$\begin{aligned}\frac{1}{N} \sum_{k=0}^{N-1} R^2(k) &\rightarrow \int_{-\pi}^{\pi} \int_{-\pi}^{\pi} f(\lambda, \nu) F(d\lambda) F(d\nu) \\ &= \int_{-\pi}^{\pi} F(\{\lambda\}) F(d\lambda) = \sum_{\lambda} F^2(\{\lambda\}),\end{aligned}$$

where the sum over  $\lambda$  contains at most a countable number of terms since the measure  $F$  is finite.

Hence (7) is equivalent to

$$\sum_{\lambda} F^2(\{\lambda\}) = 0, \quad (8)$$

which means that the spectral function  $F(\lambda) = F([-\pi, \lambda])$  is *continuous*.

2. We now turn to the problem of finding estimators for the spectral function  $F(\lambda)$  and the spectral density  $f(\lambda)$  (under the assumption that they exist).

A method that naturally suggests itself for estimating the spectral density follows from the proof of Herglotz' theorem that we gave earlier. Recall that the function

$$f_N(\lambda) = \frac{1}{2\pi} \sum_{|n| \leq N} \left(1 - \frac{|n|}{N}\right) R(n) e^{-i\lambda n} \quad (9)$$

introduced in §1 has the property that the function

$$F_N(\lambda) = \int_{-\pi}^{\lambda} f_N(\nu) d\nu$$

converges on the whole (Chapter III, §1) to the spectral function  $F(\lambda)$ . Therefore if  $F(\lambda)$  has a density  $f(\lambda)$ , we have

$$\int_{-\pi}^{\lambda} f_N(v) dv \rightarrow \int_{-\pi}^{\lambda} f(v) dv \quad (10)$$

for each  $\lambda \in [-\pi, \pi)$ .

Starting from these facts and recalling that an estimator for  $R(n)$  (on the basis of the observations  $x_0, x_1, \dots, x_{N-1}$ ) is  $\hat{R}_N(n; x)$ , we take as an estimator for  $f(\lambda)$  the function

$$\hat{f}_N(\lambda; x) = \frac{1}{2\pi} \sum_{|n| < N} \left(1 - \frac{|n|}{N}\right) \hat{R}_N(n; x) e^{-i\lambda n}, \quad (11)$$

putting  $\hat{R}_N(n; x) = \hat{R}_N(|n|; x)$  for  $|n| < N$ .

The function  $\hat{f}_N(\lambda; x)$  is known as a *periodogram*. It is easily verified that it can also be represented in the following more convenient form:

$$\hat{f}_N(\lambda; x) = \frac{1}{2\pi N} \left| \sum_{n=0}^{N-1} x_n e^{-i\lambda n} \right|^2. \quad (12)$$

Since  $E\hat{R}_N(n; \xi) = R(n)$ ,  $|n| < N$ , we have

$$E\hat{f}_N(\lambda; \xi) = f_N(\lambda).$$

If the spectral function  $F(\lambda)$  has density  $f(\lambda)$ , then, since  $f_N(\lambda)$  can also be written in the form (1.34), we find that

$$\begin{aligned} f_N(\lambda) &= \frac{1}{2\pi N} \sum_{k=0}^{N-1} \sum_{l=0}^{N-1} \int_{-\pi}^{\pi} e^{iv(k-l)} e^{i\lambda(l-k)} f(v) dv \\ &= \int_{-\pi}^{\pi} \frac{1}{2\pi N} \left| \sum_{k=0}^{N-1} e^{i(v-\lambda)k} \right|^2 f(v) dv. \end{aligned}$$

The function

$$\Phi_N(\lambda) = \frac{1}{2\pi N} \left| \sum_{k=0}^{N-1} e^{i\lambda k} \right|^2 = \frac{1}{2\pi N} \left| \frac{\sin \frac{\lambda}{2} N}{\sin \lambda/2} \right|^2$$

is the Fejér kernel. It is known, from the properties of this function, that for almost every  $\lambda$  (with respect to Lebesgue measure)

$$\int_{-\pi}^{\pi} \Phi_N(\lambda - v) f(v) dv \rightarrow f(\lambda). \quad (13)$$

Therefore for almost every  $\lambda \in [-\pi, \pi)$

$$E\hat{f}_N(\lambda; \xi) \rightarrow f(\lambda); \quad (14)$$

in other words, the estimator  $\hat{f}_N(\lambda; x)$  of  $f(\lambda)$  on the basis of  $x_0, x_1, \dots, x_{N-1}$  is *asymptotically unbiased*.



In this sense the estimator  $\hat{f}_N(\lambda; x)$  can be considered to be "good." However, at the individual observed values  $x_0, \dots, x_{N-1}$  the values of the periodogram  $\hat{f}_N(\lambda; x)$  usually turn out to be far from the actual values  $f(\lambda)$ . In fact, let  $\xi = (\xi_n)$  be a stationary sequence of independent Gaussian random variables,  $\xi_n \sim \mathcal{N}(0, 1)$ . Then  $f(\lambda) \equiv 1/2\pi$  and

$$\hat{f}_N(\lambda; \xi) = \frac{1}{2\pi} \left| \frac{1}{\sqrt{N}} \sum_{k=0}^{N-1} \xi_k e^{-i\lambda k} \right|^2.$$

Then at the point  $\lambda = 0$  we have  $\hat{f}_N(0, \xi)$  coinciding in distribution with the square of the Gaussian random variable  $\eta \sim \mathcal{N}(0, 1)$ . Hence, for every  $N$ ,

$$\mathbf{E} |\hat{f}_N(0; \xi) - f(0)|^2 = \frac{1}{4\pi^2} \mathbf{E} |\eta^2 - 1|^2 > 0.$$

Moreover, an easy calculation shows that if  $f(\lambda)$  is the spectral density of a stationary sequence  $\xi = (\xi_n)$  that is constructed as a moving average:

$$\xi_n = \sum_{k=0}^{\infty} a_k \varepsilon_{n-k} \quad (15)$$

with  $\sum_{k=0}^{\infty} |a_k| < \infty$ ,  $\sum_{k=0}^{\infty} |a_k|^2 < \infty$ , where  $\varepsilon = (\varepsilon_n)$  is white noise with  $\mathbf{E} \varepsilon_0^4 < \infty$ , then

$$\lim_{N \rightarrow \infty} \mathbf{E} |\hat{f}_N(\lambda; \xi) - f(\lambda)|^2 = \begin{cases} 2f^2(0), & \lambda = 0, \pm\pi, \\ f^2(\lambda), & \lambda \neq 0, \pm\pi. \end{cases} \quad (16)$$

Hence it is clear that the periodogram cannot be a satisfactory estimator of the spectral density. To improve the situation, one often uses an estimator for  $f(\lambda)$  of the form

$$f_N^W(\lambda; x) = \int_{-\pi}^{\pi} W_N(\lambda - \nu) \hat{f}_N(\nu; x) d\nu, \quad (17)$$

which is obtained from the periodogram  $\hat{f}_N(\lambda; x)$  and a smoothing function  $W_N(\lambda)$ , and which we call a *spectral window*. Natural requirements on  $W_N(\lambda)$  are:

- (a)  $W_N(\lambda)$  has a sharp maximum at  $\lambda = 0$ ;
- (b)  $\int_{-\pi}^{\pi} W_N(\lambda) d\lambda = 1$ ;
- (c)  $\mathbf{P} |\hat{f}_N^W(\lambda; \xi) - f(\lambda)|^2 \rightarrow 0, \quad N \rightarrow \infty, \quad \lambda \in [-\pi, \pi).$

By (14) and (b) the estimators  $\hat{f}_N^W(\lambda; \xi)$  are asymptotically unbiased. Condition (c) is the condition of asymptotic consistency in mean-square, which, as we showed above, is violated for the periodogram. Finally, condition (a) ensures that the required frequency  $\lambda$  is "picked out" from the periodogram.

Let us give some examples of estimators of the form (17).

Bartlett's estimator is based on the spectral window

$$W_N(\lambda) = a_N B(a_N \lambda),$$

where  $a_N \uparrow \infty$ ,  $a_N/N \rightarrow 0$ ,  $N \rightarrow \infty$ , and

$$B(\lambda) = \frac{1}{2\pi} \left| \frac{\sin(\lambda/2)}{\lambda/2} \right|^2.$$

Parzen's estimator takes the spectral window to be

$$W_N(\lambda) = a_N P(a_N \lambda),$$

where  $a_N$  are the same as before and

$$P(\lambda) = \frac{3}{8\pi} \left| \frac{\sin(\lambda/4)}{\lambda/4} \right|^4.$$

Zhurbenko's estimator is constructed from a spectral window of the form

$$W_N(\lambda) = a_N Z(a_N \lambda)$$

with

$$Z(\lambda) = \begin{cases} -\frac{\alpha+1}{2\alpha} |\lambda|^\alpha + \frac{\alpha+1}{2\alpha}, & |\lambda| \leq 1, \\ 0, & |\lambda| > 1, \end{cases}$$

where  $0 < \alpha \leq 2$  and the  $a_N$  are selected in a particular way.

We shall not spend any more time on problems of estimating spectral densities; we merely note that there is an extensive statistical literature dealing with the construction of spectral windows and the comparison of the corresponding estimators  $\hat{f}_N^W(\lambda; x)$ .

**3.** We now consider the problem of estimating the spectral function  $F(\lambda) = F([-\pi, \lambda])$ . We begin by defining

$$F_N(\lambda) = \int_{-\pi}^{\lambda} f_N(v) dv, \quad \hat{F}_N(\lambda; x) = \int_{-\pi}^{\lambda} \hat{f}_N(v; x) dv,$$

where  $\hat{f}_N(v; x)$  is the periodogram constructed with  $(x_0, x_1, \dots, x_{N-1})$ .

It follows from the proof of Herglotz' theorem (§1) that

$$\int_{-\pi}^{\pi} e^{i\lambda n} dF_N(\lambda) \rightarrow \int_{-\pi}^{\pi} e^{i\lambda n} dF(\lambda)$$

for every  $n \in \mathbb{Z}$ . Hence it follows (compare the corollary to Theorem 1, §3, Chapter III) that  $F_N \Rightarrow F$ , i.e.  $F_N(\lambda)$  converges to  $F(\lambda)$  at each point of continuity of  $F(\lambda)$ .

Observe that

$$\int_{-\pi}^{\pi} e^{i\lambda n} d\hat{F}_N(\lambda; \xi) = \hat{R}_N(n; \xi) \left(1 - \frac{|n|}{N}\right)$$

for all  $|n| < N$ . Therefore if we suppose that  $\hat{R}_N(n; \xi)$  converges to  $R(n)$  with probability one as  $N \rightarrow \infty$ , we have

$$\int_{-\pi}^{\pi} e^{i\lambda n} d\hat{F}_N(\lambda; \xi) \rightarrow \int_{-\pi}^{\pi} e^{i\lambda n} dF(\lambda) \quad (\text{P-a.s.})$$

and therefore  $\hat{F}_N(\lambda; \xi) \Rightarrow F(\lambda)$  (P-a.s.).

It is then easy to deduce (if necessary, passing from a sequence to a subsequence) that if  $\hat{R}_N(n; \xi) \rightarrow R(n)$  in probability, then  $\hat{F}_N(\lambda; \xi) \Rightarrow F(\lambda)$  in probability.

#### 4. PROBLEMS

1. In (15) let  $\varepsilon_n \sim \mathcal{N}(0, 1)$ . Show that

$$(N - n)\sqrt{\hat{R}_N(n, \xi)} \rightarrow 2\pi \int_{-\pi}^{\pi} (1 + e^{2in\lambda})f^2(\lambda) d\lambda$$

for every  $n$ , as  $N \rightarrow \infty$ .

2. Establish (16) and the following generalization:

$$\lim_{N \rightarrow \infty} \text{cov}(\hat{f}_N(\lambda; \xi), \hat{f}_N(\nu; \xi)) = \begin{cases} 2f^2(0), & \lambda = \nu = 0, \pm\pi, \\ f^2(\lambda), & \lambda = \nu \neq 0, \pm\pi, \\ 0, & \lambda \neq \pm\nu. \end{cases}$$

### §5. Wold's Expansion

1. In contrast to the representation (3.2) which gives an expansion of a stationary sequence in the *frequency* domain, Wold's expansion operates in the *time* domain. The main point of this expansion is that a stationary sequence  $\xi = (\xi_n)$ ,  $n \in \mathbb{Z}$ , can be represented as the sum of two stationary sequences, one of which is completely predictable (in the sense that its values are completely determined by its "past"), whereas the second does not have this property.

We begin with some definitions. Let  $H_n(\xi) = \bar{L}^2(\xi^n)$  and  $H(\xi) = \bar{L}^2(\xi)$  be closed linear manifolds, spanned respectively by  $\xi^n = (\dots, \xi_{n-1}, \xi_n)$  and  $\xi = (\dots, \xi_{n-1}, \xi_n, \dots)$ . Let

$$S(\xi) = \bigcap_n H_n(\xi).$$

For every  $\eta \in H(\xi)$ , denote by

$$\hat{\pi}_n(\eta) = \hat{E}(\eta | H_n(\xi))$$

the projection of  $\eta$  on the subspace  $H_n(\xi)$  (see §11, Chapter II). We also write

$$\hat{\pi}_{-\infty}(\eta) = \hat{E}(\eta | S(\xi)).$$

Every element  $\eta \in H(\xi)$  can be represented as

$$\eta = \hat{\pi}_{-\infty}(\eta) + (\eta - \hat{\pi}_{-\infty}(\eta)),$$

where  $\eta - \hat{\pi}_{-\infty}(\eta) \perp \hat{\pi}_{-\infty}(\eta)$ . Therefore  $H(\xi)$  is represented as the orthogonal sum

$$H(\xi) = S(\xi) \oplus R(\xi),$$

where  $S(\xi)$  consists of the elements  $\hat{\pi}_{-\infty}(\eta)$  with  $\eta \in H(\xi)$ , and  $R(\xi)$  consists of the elements of the form  $\eta - \hat{\pi}_{-\infty}(\eta)$ .

We shall now assume that  $E\xi_n = 0$  and  $\forall \xi_n > 0$ . Then  $H(\xi)$  is automatically nontrivial (contains elements different from zero).

**Definition 1.** A stationary sequence  $\xi = (\xi_n)$  is *regular* if

$$H(\xi) = R(\xi)$$

and *singular* if

$$H(\xi) = S(\xi).$$

**Remark.** Singular sequences are also called *deterministic* and regular sequences are called *purely* or *completely nondeterministic*. If  $S(\xi)$  is a proper subspace of  $H(\xi)$  we just say that  $\xi$  is *nondeterministic*.

**Theorem 1.** Every stationary (wide sense) random sequence  $\xi$  has a unique decomposition

$$\xi_n = \xi_n^r + \xi_n^s, \quad (1)$$

where  $\xi^r = (\xi_n^r)$  is regular and  $\xi^s = (\xi_n^s)$  is singular. Here  $\xi^r$  and  $\xi^s$  are orthogonal ( $\xi_n^r \perp \xi_m^s$  for all  $n$  and  $m$ ).

**PROOF.** We define

$$\xi_n^s = \hat{E}(\xi_n/S(\xi)), \quad \xi_n^r = \xi_n - \xi_n^s.$$

Since  $\xi_n^r \perp S(\xi)$ , for every  $n$ , we have  $S(\xi^r) \perp S(\xi)$ . On the other hand,  $S(\xi^r) \subseteq S(\xi)$  and therefore  $S(\xi^r)$  is trivial (contains only random sequences that coincide almost surely with zero). Consequently  $\xi^r$  is regular.

Moreover,  $H_n(\xi) \subseteq H_n(\xi^s) \oplus H_n(\xi^r)$  and  $H_n(\xi^s) \subseteq H_n(\xi)$ ,  $H_n(\xi^r) \subseteq H_n(\xi)$ . Therefore  $H_n(\xi) = H_n(\xi^s) \oplus H_n(\xi^r)$  and hence

$$S(\xi) \subseteq H_n(\xi^s) \oplus H_n(\xi^r) \quad (2)$$

for every  $n$ . Since  $\xi_n^r \perp S(\xi)$  it follows from (2) that

$$S(\xi) \subseteq H_n(\xi^s),$$

and therefore  $S(\xi) \subseteq S(\xi^s) \subseteq H(\xi^s)$ . But  $\xi_n^s \subseteq S(\xi)$ ; hence  $H(\xi^s) \subseteq S(\xi)$  and consequently

$$S(\xi) = S(\xi^s) = H(\xi^s),$$

which means that  $\xi^s$  is singular.

The orthogonality of  $\xi^s$  and  $\xi^r$  follows in an obvious way from  $\xi_n^s \in S(\xi)$  and  $\xi_n^r \perp S(\xi)$ .

Let us now show that (1) is unique. Let  $\xi_n = \eta_n^r + \eta_n^s$ , where  $\eta^r$  and  $\eta^s$  are regular and singular orthogonal sequences. Then since  $H_n(\eta^r) = H(\eta^r)$ , we have

$$H_n(\xi) = H_n(\eta^r) \oplus H_n(\eta^s) = H_n(\eta^r) \oplus H(\eta^s),$$

and therefore  $S(\xi) = S(\eta^r) \oplus H(\eta^s)$ . But  $S(\eta^r)$  is trivial, and therefore  $S(\xi) = H(\eta^s)$ .

Since  $\eta_n^s \in H(\eta^s) = S(\xi)$  and  $\eta_n^r \perp H(\eta^s) = S(\xi)$ , we have  $\hat{E}(\xi_n | S(\xi)) = \hat{E}(\eta_n^r + \eta_n^s | S(\xi)) = \eta_n^s$ , i.e.  $\eta_n^s$  coincides with  $\xi_n^s$ ; this establishes the uniqueness of (1).

This completes the proof of the theorem.

**2. Definition 2.** Let  $\xi = (\xi_n)$  be a nondegenerate stationary sequence. A random sequence  $\varepsilon = (\varepsilon_n)$  is an *innovation* sequence (for  $\xi$ ) if

- (a)  $\varepsilon = (\varepsilon_n)$  consists of pairwise orthogonal random variables with  $E\varepsilon_n = 0$ ,  $E|\varepsilon_n|^2 = 1$ ;
- (b)  $H_n(\xi) = H_n(\varepsilon)$  for all  $n \in \mathbb{Z}$ .

**Remark.** The reason for the term "innovation" is that  $\varepsilon_{n+1}$  provides, so to speak, new "information" not contained in  $H_n(\xi)$  (in other words, "innovates" in  $H_n(\xi)$  the information that is needed for forming  $H_{n+1}(\xi)$ ).

The following fundamental theorem establishes a connection between one-sided moving averages (Example 4, §1) and regular sequences.

**Theorem 2.** A necessary and sufficient condition for a nondegenerate sequence  $\xi$  to be regular is that there are an innovation sequence  $\varepsilon = (\varepsilon_n)$  and a sequence  $(a_n)$  of complex numbers,  $n \geq 0$ , with  $\sum_{n=0}^{\infty} |a_n|^2 < \infty$ , such that

$$\xi_n = \sum_{k=0}^{\infty} a_k \varepsilon_{n-k} \quad (\text{P-a.s.}) \quad (3)$$

**PROOF. Necessity.** We represent  $H_n(\xi)$  in the form

$$H_n(\xi) = H_{n-1}(\xi) \oplus B_n.$$

Since  $H_n(\xi)$  is spanned by elements of  $H_{n-1}(\xi)$  and elements of the form  $\beta \cdot \xi_n$ , where  $\beta$  is a complex number, the dimension (dim) of  $B_n$  is either zero or one. But the space  $H_n(\xi)$  cannot coincide with  $H_{n-1}(\xi)$  for any value of  $n$ .

In fact, if  $B_n$  is trivial for some  $n$ , then by stationarity  $B_k$  is trivial for all  $k$ , and therefore  $H(\xi) = S(\xi)$ , contradicting the assumption that  $\xi$  is regular. Thus  $B_n$  has the dimension  $\dim B_n = 1$ . Let  $\eta_n$  be a nonzero element of  $B_n$ . Put

$$\varepsilon_n = \frac{\eta_n}{\|\eta_n\|},$$

where  $\|\eta_n\|^2 = E|\eta_n|^2 > 0$ .

For given  $n$  and  $k \geq 0$ , consider the decomposition

$$H_n(\xi) = H_{n-k}(\xi) \oplus B_{n-k+1} \oplus \cdots \oplus B_n.$$

Then  $\varepsilon_{n-k}, \dots, \varepsilon_n$  is an orthogonal basis in  $B_{n-k+1} \oplus \cdots \oplus B_n$  and

$$\xi_n = \sum_{j=0}^{k-1} a_j \varepsilon_{n-j} + \hat{\pi}_{n-k}(\xi_n), \quad (4)$$

where  $a_j = E \xi_n \bar{\varepsilon}_{n-j}$ .

By Bessel's inequality (II.11.16)

$$\sum_{j=0}^{\infty} |a_j|^2 \leq \|\xi_n\|^2 < \infty.$$

It follows that  $\sum_{j=0}^{\infty} a_j \varepsilon_{n-j}$  converges in mean square, and then, by (4), equation (3) will be established as soon as we show that  $\hat{\pi}_{n-k}(\xi_n) \xrightarrow{L^2} 0$ ,  $k \rightarrow \infty$ .

It is enough to consider the case  $n = 0$ . Since

$$\hat{\pi}_{-k} = \hat{\pi}_0 + \sum_{i=0}^k [\hat{\pi}_{-i} - \hat{\pi}_{-i+1}],$$

and the terms that appear in this sum are orthogonal, we have for every  $k \geq 0$

$$\begin{aligned} \sum_{i=0}^k \|\hat{\pi}_{-i} - \hat{\pi}_{-i+1}\|^2 &= \left\| \sum_{i=0}^k (\hat{\pi}_{-i} - \hat{\pi}_{-i+1}) \right\|^2 \\ &= \|\hat{\pi}_{-k} - \hat{\pi}_0\|^2 \leq 4\|\xi_0\|^2 < \infty. \end{aligned}$$

Therefore the limit  $\lim_{k \rightarrow \infty} \hat{\pi}_{-k}$  exists (in mean square). Now  $\hat{\pi}_{-k} \in H_{-k}(\xi)$  for each  $k$ , and therefore the limit in question must belong to  $\bigcap_{k \geq 0} H_k(\xi) = S(\xi)$ . But, by assumption,  $S(\xi)$  is trivial, and therefore  $\hat{\pi}_{-k} \xrightarrow{L^2} 0$ ,  $k \rightarrow \infty$ .

**Sufficiency.** Let the nondegenerate sequence  $\xi$  have a representation (3), where  $\varepsilon = (\varepsilon_n)$  is an orthonormal system (not necessarily satisfying the condition  $H_n(\xi) = H_n(\varepsilon)$ ,  $n \in \mathbb{Z}$ ). Then  $H_n(\xi) \subseteq H_n(\varepsilon)$  and therefore  $S(\xi) = \bigcap_k H_k(\xi) \subseteq H_n(\varepsilon)$  for every  $n$ . But  $\varepsilon_{n+1} \perp H_n(\varepsilon)$ , and therefore  $\varepsilon_{n+1} \perp S(\xi)$  and at the same time  $\varepsilon = (\varepsilon_n)$  is a basis in  $H(\xi)$ . It follows that  $S(\xi)$  is trivial, and consequently  $\xi$  is regular.

This completes the proof of the theorem.

**Remark.** It follows from the proof that a nondegenerate sequence  $\xi$  is regular if and only if it admits a representation as a one-sided moving average

$$\xi_n = \sum_{k=0}^{\infty} \tilde{a}_k \tilde{\varepsilon}_{n-k}, \quad (5)$$

where  $\tilde{\varepsilon} = \tilde{\varepsilon}_n$  is an orthonormal system which (it is important to emphasize this!) does not necessarily satisfy the condition  $H_n(\xi) = H_n(\tilde{\varepsilon})$ ,  $n \in \mathbb{Z}$ . In this sense the conclusion of Theorem 2 says more, and specifically that for a regular sequence  $\xi$  there exist  $a = (a_n)$  and an orthonormal system  $\varepsilon = (\varepsilon_n)$  such that not only (5), but also (3), is satisfied, with  $H_n(\xi) = H_n(\varepsilon)$ ,  $n \in \mathbb{Z}$ .

The following theorem is an immediate corollary of Theorems 1 and 2.

**Theorem 3 (Wold's Expansion).** *If  $\xi = (\xi_n)$  is a nondegenerate stationary sequence, then*

$$\xi_n = \xi_n^s + \sum_{k=0}^{\infty} a_k \varepsilon_{n-k}, \quad (6)$$

where  $\sum_{k=0}^{\infty} |a_k|^2 < \infty$  and  $\varepsilon = (\varepsilon_n)$  is an innovation sequence (for  $\xi^r$ ).

3. The significance of the concepts introduced here (regular and singular sequences) becomes particularly clear if we consider the following (linear) extrapolation problem, for whose solution the Wold expansion (6) is especially useful.

Let  $H_0(\xi) = \bar{L}^2(\xi^0)$  be the closed linear manifold spanned by the variables  $\xi^0 = (\dots, \xi_{-1}, \xi_0)$ . Consider the problem of constructing an *optimal* (least-squares) *linear estimator*  $\hat{\xi}_n$  of  $\xi_n$  in terms of the "past"  $\xi^0 = (\dots, \xi_{-1}, \xi_0)$ . It follows from §11, Chapter II, that

$$\hat{\xi}_n = \hat{E}(\xi_n | H_0(\xi)). \quad (7)$$

(In the notation of Subsection 1,  $\hat{\xi}_n = \hat{\pi}_0(\xi_n)$ .) Since  $\xi^r$  and  $\xi^s$  are orthogonal and  $H_0(\xi) = H_0(\xi^r) \oplus H_0(\xi^s)$ , we obtain, by using (6),

$$\begin{aligned} \xi_n &= \hat{E}(\xi_n^s + \xi_n^r | H_0(\xi)) = \hat{E}(\xi_n^s | H_0(\xi)) + \hat{E}(\xi_n^r | H_0(\xi)) \\ &= \hat{E}(\xi_n^s | H_0(\xi^r) \oplus H_0(\xi^s)) + \hat{E}(\xi_n^r | H_0(\xi^r) \oplus H_0(\xi^s)) \\ &= \hat{E}(\xi_n^s | H_0(\xi^s)) + \hat{E}(\xi_n^r | H_0(\xi^r)) \\ &= \xi_n^s + \hat{E}\left(\sum_{k=0}^{\infty} a_k \varepsilon_{n-k} | H_0(\xi^r)\right). \end{aligned}$$

In (6), the sequence  $\varepsilon = (\varepsilon_n)$  is an innovation sequence for  $\xi^r = (\xi_n^r)$  and therefore  $H_0(\xi^r) = H_0(\varepsilon)$ . Therefore

$$\hat{\xi}_n = \xi_n^s + \hat{E}\left(\sum_{k=0}^{\infty} a_k \varepsilon_{n-k} | H_0(\varepsilon)\right) = \xi_n^s + \sum_{k=n}^{\infty} a_k \varepsilon_{n-k} \quad (8)$$

and the mean-square error of predicting  $\xi_n$  by  $\xi_0 = (\dots, \xi_{-1}, \xi_0)$  is

$$\sigma_n^2 = E|\xi_n - \hat{\xi}_n|^2 = \sum_{k=0}^{n-1} |a_k|^2. \quad (9)$$

We can draw two important conclusions.

- (a) If  $\xi$  is *singular*, then for every  $n \geq 1$  the error (in the extrapolation)  $\sigma_n^2$  is zero; in other words, we can predict  $\xi_n$  without error from its "past"  $\xi^0 = (\dots, \xi_{-1}, \xi_0)$ .
- (b) If  $\xi$  is *regular*, then  $\sigma_n^2 \leq \sigma_{n+1}^2$  and

$$\lim_{n \rightarrow \infty} \sigma_n^2 = \sum_{k=0}^{\infty} |a_k|^2. \quad (10)$$

Since

$$\sum_{k=0}^{\infty} |a_k|^2 = E|\xi_n|^2,$$

it follows from (10) and (9) that

$$\hat{\xi}_n \xrightarrow{L^2} 0, \quad n \rightarrow \infty;$$

i.e. as  $n$  increases, the prediction of  $\xi_n$  in terms of  $\xi_0 = (\dots, \xi_{-1}, \xi_0)$  becomes trivial (reducing simply to  $E\xi_n = 0$ ).

4. Let us suppose that  $\xi$  is a nondegenerate regular stationary sequence. According to Theorem 2, every such sequence admits a representation as a one-sided moving average

$$\xi_n = \sum_{k=0}^{\infty} a_k \varepsilon_{n-k}, \quad (11)$$

where  $\sum_{k=0}^{\infty} |a_k|^2 < \infty$  and the orthonormal sequence  $\varepsilon = (\varepsilon_n)$  has the important property that

$$H_n(\xi) = H_n(\varepsilon), \quad n \in \mathbb{Z}. \quad (12)$$

The representation (11) means (see Subsection 3, §3) that  $\xi_n$  can be interpreted as the output signal of a physically realizable filter with impulse response  $a = (a_k), k \geq 0$ , when the input is  $\varepsilon = (\varepsilon_n)$ .

Like any sequence of two-sided moving averages, a regular sequence has a spectral density  $f(\lambda)$ . But since a regular sequence admits a representation as a one-sided moving average it is possible to obtain additional information about properties of the spectral density.

In the first place, it is clear that

$$f(\lambda) = \frac{1}{2\pi} |\varphi(\lambda)|^2,$$



where

$$\varphi(\lambda) = \sum_{k=0}^{\infty} e^{-i\lambda k} a_k, \quad \sum_{k=0}^{\infty} |a_k|^2 < \infty. \quad (13)$$

Put

$$\Phi(z) = \sum_{k=0}^{\infty} a_k z^k. \quad (14)$$

This function is analytic in the open domain  $|z| < 1$  and since  $\sum_{k=0}^{\infty} |a_k|^2 < \infty$  it belongs to the *Hardy class*  $H^2$ , the class of functions  $g = g(z)$ , analytic in  $|z| < 1$ , satisfying

$$\sup_{0 \leq r < 1} \frac{1}{2\pi} \int_{-\pi}^{\pi} |g(re^{i\theta})|^2 d\theta < \infty. \quad (15)$$

In fact,

$$\frac{1}{2\pi} \int_{-\pi}^{\pi} |\Phi(re^{i\theta})|^2 d\theta = \sum_{k=0}^{\infty} |a_k|^2 r^{2k}$$

and

$$\sup_{0 \leq r < 1} \sum |a_k|^2 r^{2k} \leq \sum |a_k|^2 < \infty.$$

It is shown in the theory of functions of a complex variable that the boundary function  $\Phi(e^{i\lambda})$ ,  $-\pi \leq \lambda < \pi$ , of  $\Phi \in H^2$ , not identically zero, has the property that

$$\int_{-\pi}^{\pi} \ln |\Phi(e^{-i\lambda})| d\lambda > -\infty. \quad (16)$$

In our case

$$f(\lambda) = \frac{1}{2\pi} |\Phi(e^{-i\lambda})|^2,$$

where  $\Phi \in H^2$ . Therefore

$$\ln f(\lambda) = -\ln 2\pi + 2 \ln |\Phi(e^{-i\lambda})|,$$

and consequently the spectral density  $f(\lambda)$  of a regular process satisfies

$$\int_{-\pi}^{\pi} \ln f(\lambda) d\lambda > -\infty. \quad (17)$$

On the other hand, let the spectral density  $f(\lambda)$  satisfy (17). It again follows from the theory of functions of a complex variable that there is then a function  $\Phi(z) = \sum_{k=0}^{\infty} a_k z^k$  in the Hardy class  $H^2$  such that (almost everywhere with respect to Lebesgue measure)

$$f(\lambda) = \frac{1}{2\pi} |\Phi(e^{-i\lambda})|^2.$$

Therefore if we put  $\varphi(\lambda) = \Phi(e^{-i\lambda})$  we obtain

$$f(\lambda) = \frac{1}{2\pi} |\varphi(\lambda)|^2,$$

where  $\varphi(\lambda)$  is given by (13). Then it follows from the corollary to Theorem 3, §3, that  $\xi$  admits a representation as a one-sided moving average (11), where  $\varepsilon = (\varepsilon_n)$  is an orthonormal sequence. From this and from the Remark on Theorem 2, it follows that  $\xi$  is regular.

Thus we have the following theorem.

**Theorem 4** (Kolmogorov). *Let  $\xi$  be a nondegenerate regular stationary sequence. Then there is a spectral density  $f(\lambda)$  such that*

$$\int_{-\pi}^{\pi} \ln f(\lambda) d\lambda > -\infty. \quad (18)$$

*In particular,  $f(\lambda) > 0$  (almost everywhere with respect to Lebesgue measure).*

*Conversely, if  $\xi$  is a stationary sequence with a spectral density satisfying (18), the sequence is regular.*

## 5. PROBLEMS

1. Show that a stationary sequence with discrete spectrum (piecewise-constant spectral function  $F(\lambda)$ ) is singular.
2. Let  $\sigma_n^2 = E|\xi_n - \hat{\xi}_n|^2$ ,  $\hat{\xi}_n = \hat{E}(\xi_n | H_0(\xi))$ . Show that if  $\sigma_n^2 = 0$  for some  $n \geq 1$ , the sequence is singular; if  $\sigma_n^2 \rightarrow R(0)$  as  $n \rightarrow \infty$ , the sequence is regular.
3. Show that the stationary sequence  $\xi = (\xi_n)$ ,  $\xi_n = e^{in\varphi}$ , where  $\varphi$  is a uniform random variable on  $[0, 2\pi]$ , is regular. Find the estimator  $\hat{\xi}_n$  and the number  $\sigma_n^2$ , and show that the nonlinear estimator

$$\hat{\xi}_n = \left( \frac{\xi_0}{\xi_{-1}} \right)^n$$

provides a correct estimate of  $\xi_n$  by the "past"  $\xi^0 = (\dots, \xi_{-1}, \xi_0)$ , i.e.

$$E|\hat{\xi}_n - \xi_n|^2 = 0, \quad n \geq 1.$$

## §6. Extrapolation, Interpolation and Filtering

**1. Extrapolation.** According to the preceding section, a singular sequence admits an error-free prediction (extrapolation) of  $\xi_n$ ,  $n \geq 1$ , in terms of the "past,"  $\xi^0 = (\dots, \xi_{-1}, \xi_0)$ . Consequently it is reasonable, when considering the problem of extrapolation for arbitrary stationary sequences, to begin with the case of regular sequences.

According to Theorem 2 of §5, every regular sequence  $\xi = (\xi_n)$  admits a representation as a one-sided moving average,

$$\xi_n = \sum_{k=0}^{\infty} a_k \varepsilon_{n-k} \quad (1)$$

with  $\sum_{k=0}^{\infty} |a_k|^2 < \infty$  and some innovation sequence  $\varepsilon = (\varepsilon_n)$ . It follows from §5 that the representation (1) solves the problem of finding the optimal (linear) estimator  $\hat{\xi} = \hat{E}(\xi_n | H_0(\xi))$  since, by (5.8),

$$\hat{\xi}_n = \sum_{k=n}^{\infty} a_k \varepsilon_{n-k} \quad (2)$$

and

$$\sigma_n^2 = E|\xi_n - \hat{\xi}_n|^2 = \sum_{k=0}^{n-1} |a_k|^2. \quad (3)$$

However, this can be considered only as a theoretical solution, for the following reasons.

The sequences that we consider are ordinarily not given to us by means of their representations (1), but by their covariance functions  $R(n)$  or the spectral densities  $f(\lambda)$  (which exist for regular sequences). Hence a solution (2) can only be regarded as satisfactory if the coefficients  $a_k$  are given in terms of  $R(n)$  or of  $f(\lambda)$ , and  $\varepsilon_k$  are given by their values  $\dots, \xi_{k-1}, \xi_k$ .

Without discussing the problem in general, we consider only the special case (of interest in applications) when the spectral density has the form

$$f(\lambda) = \frac{1}{2\pi} |\Phi(e^{-i\lambda})|^2, \quad (4)$$

where  $\Phi(z) = \sum_{k=0}^{\infty} b_k z^k$  has radius of convergence  $r > 1$  and has no zeros in  $|z| \leq 1$ .

Let

$$\xi_n = \int_{-\pi}^{\pi} e^{i\lambda n} Z(d\lambda) \quad (5)$$

be the spectral representation of  $\xi = (\xi_n)$ ,  $n \in \mathbb{Z}$ .

**Theorem 1.** *If the spectral density of  $\xi$  has the density (4), then the optimal (linear) estimator  $\hat{\xi}_n$  of  $\xi_n$  in terms of  $\xi^0 = (\dots, \xi_{-1}, \xi_0)$  is given by*

$$\hat{\xi}_n = \int_{-\pi}^{\pi} \hat{\phi}_n(\lambda) Z(d\lambda), \quad (6)$$

where

$$\hat{\phi}_n(\lambda) = e^{i\lambda n} \frac{\Phi_n(e^{-i\lambda})}{\Phi(e^{-i\lambda})} \quad (7)$$

and

$$\Phi_n(z) = \sum_{k=n}^{\infty} b_k z^k.$$

PROOF. According to the remark on Theorem 2 of §3, every variable  $\xi_n \in H_0(\xi)$  admits a representation in the form

$$\xi_n = \int_{-\pi}^{\pi} \tilde{\varphi}_n(\lambda) Z(d\lambda), \quad \tilde{\varphi}_n \in H_0(F), \quad (8)$$

where  $H_0(F)$  is the closed linear manifold spanned by the functions  $e_n = e^{i\lambda n}$  for  $n \leq 0$  ( $F(\lambda) = \int_{-\pi}^{\lambda} f(v) dv$ ).

Since

$$\begin{aligned} E|\xi_n - \xi_n|^2 &= E \left| \int_{-\pi}^{\pi} (e^{i\lambda n} - \tilde{\varphi}_n(\lambda)) Z(d\lambda) \right|^2 \\ &= \int_{-\pi}^{\pi} |e^{i\lambda n} - \tilde{\varphi}_n(\lambda)|^2 f(\lambda) d\lambda, \end{aligned}$$

the proof that (6) is optimal reduces to proving that

$$\inf_{\tilde{\varphi}_n \in H_0(F)} \int_{-\pi}^{\pi} |e^{i\lambda n} - \tilde{\varphi}_n(\lambda)|^2 f(\lambda) d\lambda = \int_{-\pi}^{\pi} |e^{i\lambda n} - \hat{\varphi}_n(\lambda)|^2 f(\lambda) d\lambda. \quad (9)$$

It follows from Hilbert-space theory (§11, Chapter II) that the optimal function  $\hat{\varphi}_n(\lambda)$  (in the sense of (9)) is determined by the two conditions

- (1)  $\hat{\varphi}_n(\lambda) \in H_0(F)$ ,
  - (2)  $e^{i\lambda n} - \hat{\varphi}_n(\lambda) \perp H_0(F)$ .
- (10)

Since

$$e^{i\lambda n} \Phi_n(e^{-i\lambda}) = e^{i\lambda n} [b_n e^{-i\lambda n} + b_{n+1} e^{-i\lambda(n+1)} + \dots] \in H_0(F)$$

and in a similar way  $1/\Phi(e^{-i\lambda}) \in H_0(F)$ , the function  $\hat{\varphi}_n(\lambda)$  defined in (7) belongs to  $H_0(F)$ . Therefore in proving that  $\hat{\varphi}_n(\lambda)$  is optimal it is sufficient to verify that, for every  $m \geq 0$ ,

$$e^{i\lambda n} - \hat{\varphi}_n(\lambda) \perp e^{i\lambda m},$$

i.e.

$$I_{n,m} \equiv \int_{-\pi}^{\pi} [e^{i\lambda n} - \hat{\varphi}_n(\lambda)] e^{-i\lambda m} f(\lambda) d\lambda = 0, \quad m \geq 0.$$

The following chain of equations shows that this is actually the case:

$$\begin{aligned} I_{n,m} &= \frac{1}{2\pi} \int_{-\pi}^{\pi} e^{i\lambda(n-m)} \left[ 1 - \frac{\Phi_n(e^{-i\lambda})}{\Phi(e^{-i\lambda})} \right] |\Phi(e^{-i\lambda})|^2 d\lambda \\ &= \frac{1}{2\pi} \int_{-\pi}^{\pi} e^{i\lambda(n-m)} [\Phi(e^{-i\lambda}) - \Phi_n(e^{-i\lambda})] \overline{\Phi(e^{-i\lambda})} d\lambda \\ &= \frac{1}{2\pi} \int_{-\pi}^{\pi} e^{i\lambda(n-m)} \left( \sum_{k=0}^{n-1} b_k e^{-i\lambda k} \right) \left( \sum_{l=0}^{\infty} \bar{b}_l e^{i\lambda l} \right) d\lambda \\ &= \frac{1}{2\pi} \int_{-\pi}^{\pi} e^{-i\lambda m} \left( \sum_{k=0}^{n-1} b_k e^{i\lambda(n-k)} \right) \left( \sum_{l=0}^{\infty} \bar{b}_l e^{i\lambda l} \right) d\lambda = 0, \end{aligned}$$

where the last equation follows because, for  $m \geq 0$  and  $r > 1$ ,

$$\int_{-\pi}^{\pi} e^{-i\lambda m} e^{i\lambda r} d\lambda = 0.$$

This completes the proof of the theorem.

**Remark 1.** Expanding  $\hat{\phi}_n(\lambda)$  in a Fourier series, we find that the predicted value  $\hat{\xi}_n$  of  $\xi_n$ ,  $n \geq 1$ , in terms of the past,  $\xi^0 = (\dots, \xi_{-1}, \xi_0)$ , is given by the formula

$$\hat{\xi}_n = C_0 \xi_0 + C_{-1} \xi_{-1} + C_{-2} \xi_{-2} + \dots.$$

**Remark 2.** A typical example of a spectral density represented in the form (4) is the *rational function*

$$f(\lambda) = \frac{1}{2\pi} \left| \frac{P(e^{-i\lambda})}{Q(e^{-i\lambda})} \right|^2,$$

where the polynomials  $P(z) = a_0 + a_1 z + \dots + a_p z^p$  and  $Q(z) = 1 + b_1 z + \dots + b_q z^q$  have no zeros in  $\{z: |z| \leq 1\}$ .

In fact, in this case it is enough to put  $\Phi(z) = P(z)/Q(z)$ . Then  $\Phi(z) = \sum_{k=0}^{\infty} C_k z^k$  and the radius of convergence of this series is greater than one.

Let us illustrate Theorem 1 with two examples.

**EXAMPLE 1.** Let the spectral density be

$$f(\lambda) = \frac{1}{2\pi} (5 + 4 \cos \lambda).$$

The corresponding covariance function  $R(n)$  has the shape of a triangle with

$$R(0) = 5, \quad R(\pm 1) = 2, \quad R(n) = 0 \quad \text{for } |n| \geq 2. \quad (11)$$

Since this spectral density can be represented in the form

$$f(\lambda) = \frac{1}{2\pi} |2 + e^{-i\lambda}|^2,$$

we may apply Theorem 1. We find easily that

$$\hat{\phi}_1(\lambda) = e^{i\lambda} \frac{e^{-i\lambda}}{2 + e^{-i\lambda}}, \quad \hat{\phi}_n(\lambda) = 0 \quad \text{for } n \geq 2. \quad (12)$$

Therefore  $\hat{\xi}_n = 0$  for all  $n \geq 2$ , i.e. the (linear) prediction of  $\xi_n$  in terms of  $\xi^0 = (\dots, \xi_{-1}, \xi_0)$  is trivial, which is not at all surprising if we observe that, by (11), the correlation between  $\xi_n$  and any of  $\xi_0, \xi_{-1}, \dots$  is zero for  $n \geq 2$ .

For  $n = 1$  we find from (6) and (12) that

$$\begin{aligned}\hat{\xi}_1 &= \int_{-\pi}^{\pi} e^{i\lambda} \frac{e^{-i\lambda}}{2 + e^{-i\lambda}} Z(d\lambda) \\ &= \frac{1}{2} \int_{-\pi}^{\pi} \frac{1}{\left(1 + \frac{e^{-i\lambda}}{2}\right)} Z(d\lambda) = \sum_{k=0}^{\infty} \frac{(-1)^k}{2^{k+1}} \int_{-\pi}^{\pi} e^{-ik\lambda} Z(d\lambda) \\ &= \sum_{k=0}^{\infty} \frac{(-1)^k \xi_k}{2^{k+1}} = \frac{1}{2} \xi_0 - \frac{1}{4} \xi_{-1} + \dots\end{aligned}$$

**EXAMPLE 2.** Let the covariance function be

$$R(n) = a^n, \quad |a| < 1.$$

Then (see Example 5 in §1)

$$f(\lambda) = \frac{1}{2\pi} \frac{1 - |a|^2}{|1 - ae^{-i\lambda}|^2},$$

i.e.

$$f(\lambda) = \frac{1}{2\pi} |\Phi(e^{-i\lambda})|^2,$$

where

$$\Phi(z) = \frac{(1 - |a|^2)^{1/2}}{1 - az} = (1 - |a|^2)^{1/2} \sum_{k=0}^{\infty} (az)^k,$$

from which  $\hat{\phi}_n(\lambda) = a^n$  and therefore

$$\hat{\xi}_n = \int_{-\pi}^{\pi} a^n Z(d\lambda) = a^n \xi_0.$$

In other words, in order to predict the value of  $\xi_n$  from the observations  $\xi^0 = (\dots, \xi_{-1}, \xi_0)$  it is sufficient to know only the last observation  $\xi_0$ .

**Remark 3.** It follows from the Wold expansion of the regular sequence  $\xi = (\xi_n)$  with

$$\xi_n = \sum_{k=0}^{\infty} a_k \xi_{n-k} \tag{13}$$

that the spectral density  $f(\lambda)$  admits the representation

$$f(\lambda) = \frac{1}{2\pi} |\Phi(e^{-i\lambda})|^2, \tag{14}$$

where

$$\Phi(z) = \sum_{k=0}^{\infty} a_k z^k. \tag{15}$$

It is evident that the converse also holds, that is, if  $f(\lambda)$  admits the representation (14) with a function  $\Phi(z)$  of the form (15), then the Wold expansion of  $\xi_n$  has the form (13). Therefore the problem of representing the spectral density in the form (14) and the problem of determining the coefficients  $a_k$  in the Wold expansion are equivalent.

The assumptions that  $\Phi(z)$  in Theorem 1 has no zeros for  $|z| \leq 1$  and that  $r > 1$  are in fact not essential. In other words, if the spectral density of a regular sequence is represented in the form (14), then the optimal estimator  $\hat{\xi}_n$  (in the mean square sense) for  $\xi_n$  in terms of  $\xi^0 = (\dots, \xi_{-1}, \xi_0)$  is determined by formulas (6) and (7).

**Remark 4.** Theorem 1 (with the preceding remark) solves the prediction problem for regular sequences. Let us show that in fact the same answer remains valid for arbitrary stationary sequences. More precisely, let

$$\xi_n = \xi_n^s + \xi_n^r, \quad \xi_n = \int_{-\pi}^{\pi} e^{i\lambda n} Z(d\lambda), \quad F(\Delta) = E|Z(\Delta)|^2,$$

and let  $f^r(\lambda) = (1/2\pi)|\Phi(e^{-i\lambda})|^2$  be the spectral density of the regular sequence  $\xi^r = (\xi_n^r)$ . Then  $\hat{\xi}_n$  is determined by (6) and (7).

In fact, let (see Subsection 3, §5)

$$\hat{\xi}_n = \int_{-\pi}^{\pi} \hat{\phi}_n(\lambda) Z(d\lambda), \quad \hat{\xi}_n^r = \int_{-\pi}^{\pi} \hat{\phi}_n^r(\lambda) Z^r(d\lambda),$$

where  $Z^r(\Delta)$  is the orthogonal stochastic measure in the representation of the regular sequence  $\xi^r$ . Then

$$\begin{aligned} E|\xi_n - \hat{\xi}_n|^2 &= \int_{-\pi}^{\pi} |e^{i\lambda n} - \hat{\phi}_n(\lambda)|^2 F(d\lambda) \\ &\geq \int_{-\pi}^{\pi} |e^{i\lambda n} - \hat{\phi}_n(\lambda)|^2 f^r(\lambda) d\lambda \geq \int_{-\pi}^{\pi} |e^{i\lambda n} - \hat{\phi}_n^r(\lambda)|^2 f^r(\lambda) d\lambda \\ &= E|\xi_n^r - \hat{\xi}_n^r|^2. \end{aligned} \quad (16)$$

But  $\xi_n - \hat{\xi}_n = \xi_n^r - \hat{\xi}_n^r$ . Hence  $E|\xi_n - \hat{\xi}_n|^2 = E|\xi_n^r - \hat{\xi}_n^r|^2$ , and it follows from (16) that we may take  $\hat{\phi}_n(\lambda)$  to be  $\hat{\phi}_n^r(\lambda)$ .

**2. Interpolation.** Suppose that  $\xi = (\xi_n)$  is a regular sequence with spectral density  $f(\lambda)$ . The simplest interpolation problem is the problem of constructing the optimal (mean-square) linear estimator from the results of the measurements  $\{\xi_n, n = \pm 1, \pm 2, \dots\}$  omitting  $\xi_0$ .

Let  $H^0(\xi)$  be the closed linear manifold spanned by  $\xi_n, n \neq 0$ . Then according to the results of Theorem 2, §3, every random variable  $\eta \in H^0(\xi)$  can be represented in the form

$$\eta = \int_{-\pi}^{\pi} \varphi(\lambda) Z(d\lambda),$$

where  $\varphi$  belongs to  $H^0(F)$ , the closed linear manifold spanned by the functions  $e^{i\lambda n}$ ,  $n \neq 0$ . The estimator

$$\xi_0 = \int_{-\pi}^{\pi} \tilde{\varphi}(\lambda) Z(d\lambda) \quad (17)$$

will be optimal if and only if

$$\begin{aligned} \inf_{\eta \in H^0(\xi)} E|\xi_0 - \eta|^2 &= \inf_{\varphi \in H^0(F)} \int_{-\pi}^{\pi} |1 - \varphi(\lambda)|^2 F(d\lambda) \\ &= \int_{-\pi}^{\pi} |1 - \tilde{\varphi}(\lambda)|^2 F(d\lambda) = E|\xi_0 - \tilde{\xi}_0|^2. \end{aligned}$$

It follows from the perpendicularity properties of the Hilbert space  $H^0(F)$  that  $\tilde{\varphi}(\lambda)$  is completely determined (compare (10)) by the two conditions

- (1)  $\tilde{\varphi}(\lambda) \in H^0(F)$ ,
  - (2)  $1 - \tilde{\varphi}(\lambda) \perp H^0(F)$ .
- (18)

**Theorem 2** (Kolmogorov). *Let  $\xi = (\xi_n)$  be a regular sequence such that*

$$\int_{-\pi}^{\pi} \frac{d\lambda}{f(\lambda)} < \infty. \quad (19)$$

*Then*

$$\tilde{\varphi}(\lambda) = 1 - \frac{\alpha}{f(\lambda)}, \quad (20)$$

*where*

$$\alpha = \frac{2\pi}{\int_{-\pi}^{\pi} \frac{d\lambda}{f(\lambda)}}, \quad (21)$$

*and the interpolation error  $\delta^2 = E|\xi_0 - \tilde{\xi}_0|^2$  is given by  $\delta^2 = 2\pi \cdot \alpha$ .*

**PROOF.** We shall give the proof only under very stringent hypotheses on the spectral density, specifically that

$$0 < c \leq f(\lambda) \leq C < \infty. \quad (22)$$

It follows from (2) and (18) that

$$\int_{-\pi}^{\pi} [1 - \tilde{\varphi}(\lambda)] e^{in\lambda} f(\lambda) d\lambda = 0 \quad (23)$$

for every  $n \neq 0$ . By (22), the function  $[1 - \tilde{\varphi}(\lambda)]f(\lambda)$  belongs to the Hilbert space  $L^2([-\pi, \pi], \mathcal{B}[-\pi, \pi], \mu)$  with Lebesgue measure  $\mu$ . In this space the functions  $\{e^{in\lambda}/\sqrt{2\pi}, n = 0, \pm 1, \dots\}$  form an orthonormal basis (Problem 7, §11, Chapter II). Hence it follows from (23) that  $[1 - \tilde{\varphi}(\lambda)]f(\lambda)$  is a constant,



which we denote by  $\alpha$ . Thus the second condition in (18) leads to the conclusion that

$$\check{\varphi}(\lambda) = 1 - \frac{\alpha}{f(\lambda)}. \quad (24)$$

Starting from the first condition (18), we now determine  $\alpha$ .

By (22),  $\check{\varphi} \in L^2$  and the condition  $\check{\varphi} \in H^0(F)$  is equivalent to the condition that  $\check{\varphi}$  belongs to the closed (in the  $L^2$  norm) linear manifold spanned by the functions  $e^{i\lambda n}$ ,  $n \neq 0$ . Hence it is clear that the zeroth coefficient in the expansion of  $\check{\varphi}(\lambda)$  must be zero. Therefore

$$0 = \int_{-\pi}^{\pi} \check{\varphi}(\lambda) d\lambda = 2\pi - \alpha \int_{-\pi}^{\pi} \frac{d\lambda}{f(\lambda)}$$

and hence  $\alpha$  is determined by (21).

Finally,

$$\begin{aligned} \delta^2 &= \mathbf{E}|\xi_0 - \check{\xi}_0|^2 = \int_{-\pi}^{\pi} |1 - \check{\varphi}(\lambda)|^2 f(\lambda) d\lambda \\ &= |\alpha|^2 \int_{-\pi}^{\pi} \frac{f(\lambda)}{f^2(\lambda)} d\lambda = \frac{4\pi^2}{\int_{-\pi}^{\pi} \frac{d\lambda}{f(\lambda)}}. \end{aligned}$$

This completes the proof (under condition (22)).

**Corollary.** *If*

$$\check{\varphi}(\lambda) = \sum_{0 < |k| \leq N} c_k e^{i\lambda k},$$

*then*

$$\check{\xi}_0 = \sum_{0 < |k| \leq N} c_k \int_{-\pi}^{\pi} e^{i\lambda k} Z(d\lambda) = \sum_{0 < |k| \leq N} c_k \xi_k.$$

**EXAMPLE 3.** Let  $f(\lambda)$  be the spectral density in Example 2 above. Then an easy calculation shows that

$$\check{\xi}_0 = \int_{-\pi}^{\pi} \frac{a}{1 + |a|^2} [e^{i\lambda} + e^{-i\lambda}] Z(d\lambda) = \frac{a}{1 + |a|^2} [\xi_1 + \xi_{-1}],$$

and the interpolation error is

$$\delta^2 = \frac{1 - |\alpha|^2}{1 + |\alpha|^2}.$$

**3. Filtering.** Let  $(\theta, \xi) = ((\theta_n), (\xi_n))$ ,  $n \in \mathbb{Z}$ , be a *partially observed sequence*, where  $\theta = (\theta_n)$  and  $\xi = (\xi_n)$  are respectively the unobserved and the observed components.

Each of the sequences  $\theta$  and  $\xi$  will be supposed stationary (wide sense) with zero mean; let the spectral densities be

$$\theta_n = \int_{-\pi}^{\pi} e^{i\lambda n} Z_{\theta}(d\lambda), \quad \text{and} \quad \xi_n = \int_{-\pi}^{\pi} e^{i\lambda n} Z_{\xi}(d\lambda).$$

We write

$$F_{\theta}(\Delta) = E|Z_{\theta}(\Delta)|^2, \quad F_{\xi}(\Delta) = E|Z_{\xi}(\Delta)|^2$$

and

$$F_{\theta\xi}(\Delta) = E Z_{\theta}(\Delta) \overline{Z_{\xi}(\Delta)}.$$

In addition, we suppose that  $\theta$  and  $\xi$  are *connected in a stationary way*, i.e. that their covariance functions  $\text{cov}(\theta_n, \xi_m) = E\theta_n \xi_m$  depend only on the differences  $n - m$ . Let  $R_{\theta\xi}(n) = E\theta_n \xi_0$ ; then

$$R_{\theta\xi}(n) = \int_{-\pi}^{\pi} e^{i\lambda n} F_{\theta\xi}(d\lambda).$$

The filtering problem that we shall consider is the construction of the optimal (mean-square) linear estimator  $\hat{\theta}_n$  of  $\theta_n$  in terms of some observation of the sequence  $\xi$ .

The problem is easily solved under the assumption that  $\theta_n$  is to be constructed from *all* the values  $\xi_m$ ,  $m \in \mathbb{Z}$ . In fact, since  $\hat{\theta}_n = \hat{E}(\theta_n | H(\xi))$  there is a function  $\hat{\phi}_n(\lambda)$  such that

$$\hat{\theta}_n = \int_{-\pi}^{\pi} \hat{\phi}_n(\lambda) Z_{\xi}(d\lambda). \quad (25)$$

As in Subsections 1 and 2, the conditions to impose on the optimal  $\hat{\phi}_n(\lambda)$  are that

- (1)  $\hat{\phi}_n(\lambda) \in H(F_{\xi})$ ,
- (2)  $\theta_n - \hat{\theta}_n \perp H(\xi)$ .

From the latter condition we find

$$\int_{-\pi}^{\pi} e^{i\lambda(n-m)} F_{\theta\xi}(d\lambda) - \int_{-\pi}^{\pi} e^{-i\lambda m} \hat{\phi}_n(\lambda) F_{\xi}(d\lambda) = 0 \quad (26)$$

for every  $m \in \mathbb{Z}$ . Therefore if we suppose that  $F_{\theta\xi}(\lambda)$  and  $F_{\xi}(\lambda)$  have densities  $f_{\theta\xi}(\lambda)$  and  $f_{\xi}(\lambda)$ , we find from (26) that

$$\int_{-\pi}^{\pi} e^{i\lambda(n-m)} [f_{\theta\xi}(\lambda) - e^{-i\lambda n} \hat{\phi}_n(\lambda) f_{\xi}(\lambda)] d\lambda = 0.$$

If  $f_{\xi}(\lambda) > 0$  (almost everywhere with respect to Lebesgue measure) we find immediately that

$$\hat{\phi}_n(\lambda) = e^{i\lambda n} \hat{\phi}(\lambda), \quad (27)$$

where

$$\hat{\phi}(\lambda) = f_{\theta\xi}(\lambda) \cdot f_{\xi}^{\oplus}(\lambda)$$

and  $f_{\xi}^{\oplus}(\lambda)$  is the "pseudotransform" of  $f_{\xi}(\lambda)$ , i.e.

$$f_{\xi}^{\oplus}(\lambda) = \begin{cases} f_{\xi}^{-1}(\lambda), & f_{\xi}(\lambda) > 0, \\ 0, & f_{\xi}(\lambda) = 0. \end{cases}$$

Then the filtering error is

$$E|\theta_n - \hat{\theta}_n|^2 = \int_{-\pi}^{\pi} [f_{\theta}(\lambda) - f_{\theta\xi}^2(\lambda)f_{\xi}^{\oplus}(\lambda)] d\lambda. \quad (28)$$

As is easily verified,  $\hat{\phi} \in H(F_{\xi})$ , and consequently the estimator (25), with the function (27), is optimal.

**EXAMPLE 4.** *Detection of a signal in the presence of noise.* Let  $\xi_n = \theta_n + \eta_n$ , where the signal  $\theta = (\theta_n)$  and the noise  $\eta = (\eta_n)$  are uncorrelated sequences with spectral densities  $f_{\theta}(\lambda)$  and  $f_{\eta}(\lambda)$ . Then

$$\hat{\theta}_n = \int_{-\pi}^{\pi} e^{i\lambda n} \hat{\phi}(\lambda) Z_{\xi}(d\lambda),$$

where

$$\hat{\phi}(\lambda) = f_{\theta}(\lambda) [f_{\theta}(\lambda) + f_{\eta}(\lambda)]^{\oplus},$$

and the filtering error is

$$E|\theta_n - \hat{\theta}_n|^2 = \int_{-\pi}^{\pi} [f_{\theta}(\lambda)f_{\eta}(\lambda)] [f_{\theta}(\lambda) + f_{\eta}(\lambda)]^{\oplus} d\lambda.$$

The solution (25) obtained above can now be used to construct an optimal estimator  $\hat{\theta}_{n+m}$  of  $\theta_{n+m}$  as a result of observing  $\xi_k$ ,  $k \leq n$ , where  $m$  is a given element of  $\mathbb{Z}$ . Let us suppose that  $\xi = (\xi_n)$  is regular, with spectral density

$$f(\lambda) = \frac{1}{2\pi} |\Phi(e^{-i\lambda})|^2,$$

where  $\Phi(z) = \sum_{k=0}^{\infty} a_k z^k$ . By the Wold expansion,

$$\xi_n = \sum_{k=0}^{\infty} a_k \varepsilon_{n-k},$$

where  $\varepsilon = (\varepsilon_k)$  is white noise with the spectral resolution

$$\varepsilon_n = \int_{-\pi}^{\pi} e^{i\lambda n} Z_{\varepsilon}(d\lambda).$$

Since

$$\hat{\theta}_{n+m} = \hat{E}[\theta_{n+m} | H_n(\xi)] = \hat{E}[\hat{E}[\theta_{n+m} | H(\xi)] | H_n(\xi)] = \hat{E}[\hat{\theta}_{n+m} | H_n(\xi)]$$

and

$$\hat{\theta}_{n+m} = \int_{-\pi}^{\pi} e^{i\lambda(n+m)} \hat{\phi}(\lambda) \Phi(e^{-i\lambda}) Z_{\varepsilon}(d\lambda) = \sum_{k \leq n+m} \hat{a}_{n+m-k} \varepsilon_k,$$

where

$$\hat{a}_k = \frac{1}{2\pi} \int_{-\pi}^{\pi} e^{i\lambda k} \hat{\phi}(\lambda) \Phi(e^{-i\lambda}) d\lambda, \quad (29)$$

then

$$\tilde{\theta}_{n+m} = \hat{\mathbb{E}} \left[ \sum_{k \leq n+m} \hat{a}_{n+m-k} \varepsilon_k \mid H_n(\xi) \right].$$

But  $H_n(\xi) = H_n(\varepsilon)$  and therefore

$$\begin{aligned} \tilde{\theta}_{n+m} &= \sum_{k \leq n} \hat{a}_{n+m-k} \varepsilon_k = \int_{-\pi}^{\pi} \left[ \sum_{k \leq n} \hat{a}_{n+m-k} e^{i\lambda k} \right] Z_{\varepsilon}(d\lambda) \\ &= \int_{-\pi}^{\pi} e^{i\lambda n} \left[ \sum_{l=0}^{\infty} \hat{a}_{l+m} e^{-i\lambda l} \right] \Phi^{\oplus}(e^{-i\lambda}) Z_{\varepsilon}(d\lambda), \end{aligned}$$

where  $\Phi^{\oplus}$  is the pseudotransform of  $\Phi$ .

We have therefore established the following theorem.

**Theorem 3.** *If the sequence  $\xi = (\xi_n)$  under observation is regular, then the optimal (mean-square) linear estimator  $\tilde{\theta}_{n+m}$  of  $\theta_{n+m}$  in terms of  $\xi_k$ ,  $k \leq n$ , is given by*

$$\tilde{\theta}_{n+m} = \int_{-\pi}^{\pi} e^{i\lambda n} H_m(e^{-i\lambda}) Z_{\xi}(d\lambda), \quad (30)$$

where

$$H_m(e^{-i\lambda}) = \sum_{l=0}^{\infty} \hat{a}_{l+m} e^{-i\lambda l} \Phi^{\oplus}(e^{-i\lambda}) \quad (31)$$

and the coefficients  $a_k$  are defined by (29).

#### 4. PROBLEMS

1. Let  $\xi$  be a nondegenerate regular sequence with spectral density (4). Show that  $\Phi(z)$  has no zeros for  $|z| \leq 1$ .
2. Show that the conclusion of Theorem 1 remains valid even without the hypotheses that  $\Phi(z)$  has radius of convergence  $r > 1$  and that the zeros of  $\Phi(z)$  all lie in  $|z| > 1$ .

3. Show that, for a regular process, the function  $\Phi(z)$  introduced in (4) can be represented in the form

$$\Phi(z) = \sqrt{2\pi} \exp \left\{ \frac{1}{2} c_0 + \sum_{k=1}^{\infty} c_k z^k \right\}, \quad |z| < 1,$$

where

$$c_k = \frac{1}{2\pi} \int_{-\pi}^{\pi} e^{ik\lambda} \ln f(\lambda) d\lambda.$$

Deduce from this formula and (5.9) that the one-step prediction error  $\sigma_1^2 = E|\xi_1 - \xi_1|^2$  is given by the Szegő-Kolmogorov formula

$$\sigma_1^2 = 2\pi \exp \left\{ \frac{1}{2\pi} \int_{-\pi}^{\pi} \ln f(\lambda) d\lambda \right\}.$$

4. Prove Theorem 2 without using (22).  
 5. Let a signal  $\theta$  and a noise  $\eta$ , not correlated with each other, have spectral densities

$$f_{\theta}(\lambda) = \frac{1}{2\pi} \cdot \frac{1}{|1 + b_1 e^{-i\lambda}|^2} \quad \text{and} \quad f_{\eta}(\lambda) = \frac{1}{2\pi} \cdot \frac{1}{|1 + b_2 e^{-i\lambda}|^2}.$$

Using Theorem 3, find an estimator  $\tilde{\theta}_{n+m}$  for  $\theta_{n+m}$  in terms of  $\xi_k$ ,  $k \leq n$ , where  $\xi_k = \theta_k + \eta_k$ . Consider the same problem for the spectral densities

$$f_{\theta}(\lambda) = \frac{1}{2\pi} |2 + e^{-i\lambda}|^2 \quad \text{and} \quad f_{\eta}(\lambda) = \frac{1}{2\pi}.$$

## §7. The Kalman-Bucy Filter and Its Generalizations

1. From a computational point of view, the solution presented above for the problem of filtering out an unobservable component  $\theta$  by means of observations of  $\xi$  is not practical, since, because it is expressed in terms of the spectrum, it has to be carried out by spectral methods. In the method proposed by Kalman and Bucy, the synthesis of the optimal filter is carried out recursively; this makes it possible to do it with a digital computer. There are also other reasons for the wide use of the Kalman-Bucy filter, one being that it still "works" even without the assumption that the sequence  $(\theta, \xi)$  is stationary.

We shall present not only the usual Kalman-Bucy method, but also a generalization in which the recurrent equations determined by  $(\theta, \xi)$  have coefficients that depend on all the data observed in the past.

Thus, let us suppose that  $(\theta, \xi) = ((\theta_n), (\xi_n))$  is a partially observed sequence, and let

$$\theta_n = (\theta_1(n), \dots, \theta_k(n)) \quad \text{and} \quad \xi_n = (\xi_1(n), \dots, \xi_l(n))$$

be governed by the recurrent equations

$$\theta_{n+1} = a_0(n, \xi) + a_1(n, \xi)\theta_n + b_1(n, \xi)\varepsilon_1(n+1) + b_2(n, \xi)\varepsilon_2(n+1),$$

$$\xi_{n+1} = A_0(n, \xi) + A_1(n, \xi)\theta_n + B_1(n, \xi)\varepsilon_1(n+1) + B_2(n, \xi)\varepsilon_2(n+1).$$

(1)

Here

$$\varepsilon_1(n) = (\varepsilon_{11}(n), \dots, \varepsilon_{1k}(n)) \quad \text{and} \quad \varepsilon_2(n) = (\varepsilon_{21}(n), \dots, \varepsilon_{2l}(n))$$

are independent Gaussian vectors with independent components, each of which is normally distributed with parameters 0 and 1;  $a_0(n, \xi) = (a_{01}(n, \xi), \dots, a_{0k}(n, \xi))$  and  $A_0(n, \xi) = (A_{01}(n, \xi), \dots, A_{0l}(n, \xi))$  are vector functions, where the dependence on  $\xi = \{\xi_0, \dots, \xi_n\}$  is determined without looking ahead, i.e. for a given  $n$  the functions  $a_0(n, \xi), \dots, A_0(n, \xi)$  depend only on  $\xi_0, \dots, \xi_n$ ; the matrix functions

$$\begin{aligned} b_1(n, \xi) &= \|b_{ij}^{(1)}(n, \xi)\|, & b_2(n, \xi) &= \|b_{ij}^{(2)}(n, \xi)\|, \\ B_1(n, \xi) &= \|B_{ij}^{(1)}(n, \xi)\|, & B_2(n, \xi) &= \|B_{ij}^{(2)}(n, \xi)\|, \\ a_1(n, \xi) &= \|a_{ij}^{(1)}(n, \xi)\|, & A_1(n, \xi) &= \|A_{ij}^{(1)}(n, \xi)\| \end{aligned}$$

have orders  $k \times k, k \times l, l \times k, l \times l, k \times k, l \times k$ , respectively, and also depend on  $\xi$  without looking ahead. We also suppose that the initial vector  $(\theta_0, \xi_0)$  is independent of the sequences  $\varepsilon_1 = (\varepsilon_1(n))$  and  $\varepsilon_2 = (\varepsilon_2(n))$ .

To simplify the presentation, we shall frequently not indicate the dependence of the coefficients on  $\xi$ .

So that the system (1) will have a solution with finite second moments, we assume that  $E(\|\theta_0\|^2 + \|\xi_0\|^2) < \infty$

$$\left( \|x\|^2 = \sum_{i=1}^k x_i^2, x = (x_1, \dots, x_k) \right), \quad |a_{ij}^{(1)}(n, \xi)| \leq C, \quad |A_{ij}^{(1)}(n, \xi)| \leq C,$$

and if  $g(n, \xi)$  is any of the functions  $a_{0i}, A_{0j}, b_{ij}^{(1)}, b_{ij}^{(2)}, B_{ij}^{(1)}$  or  $B_{ij}^{(2)}$  then  $E|g(n, \xi)|^2 < \infty, n = 0, 1, \dots$ . With these assumptions,  $(\theta, \xi)$  has  $E(\|\theta_n\|^2 + \|\xi_n\|^2) < \infty, n \geq 0$ .

Now let  $\mathcal{F}_n^\xi = \sigma\{\omega: \xi_0, \dots, \xi_n\}$  be the smallest  $\sigma$ -algebra generated by  $\xi_0, \dots, \xi_n$  and

$$m_n = E(\theta_n | \mathcal{F}_n^\xi), \quad \gamma_n = E[(\theta_n - m_n)(\theta_n - m_n)^* | \mathcal{F}_n^\xi].$$

According to Theorem 1, §8, Chapter II,  $m_n = (m_1(n), \dots, m_k(n))$  is an optimal estimator (in the mean square sense) for the vector  $\theta_n = (\theta_1(n), \dots, \theta_k(n))$ , and  $E\gamma_n = E[(\theta_n - m_n)(\theta_n - m_n)^*]$  is the matrix of errors of observation. To determine these matrices for arbitrary sequences  $(\theta, \xi)$  governed by equations (1) is a very difficult problem. However, there is a further supplementary condition on  $(\theta_0, \xi_0)$  that leads to a system of recurrent equations for  $m_n$  and  $\gamma_n$  that still contains the Kalman-Bucy filter. This is the condition that the conditional distribution  $P(\theta_0 \leq a | \xi_0)$  is Gaussian,

$$P(\theta_0 \leq a | \xi_0) = \frac{1}{\sqrt{2\pi\gamma_0}} \int_{-\infty}^a \exp\left\{-\frac{(x - m_0)^2}{2\gamma_0^2}\right\} dx, \quad (2)$$

with parameters  $m_0 = m_0(\xi_0), \gamma_0 = \gamma_0(\xi_0)$ .

To begin with, let us establish an important auxiliary result.

**Lemma 1.** *Under the assumptions made above about the coefficients of (1), together with (2), the sequence  $(\theta, \xi)$  is conditionally Gaussian, i.e. the conditional distribution function*

$$P\{\theta_0 \leq a_0, \dots, \eta_n \leq a_n | \mathcal{F}_n^\xi\}$$

*is (P-a.s.) the distribution function of an  $n$ -dimensional Gaussian vector whose mean and covariance matrix depend on  $(\xi_0, \dots, \xi_n)$ .*

**PROOF.** We prove only the Gaussian character of  $P(\theta_n \leq a | \mathcal{F}_n^\xi)$ ; this is enough to let us obtain equations for  $m_n$  and  $\gamma_n$ .

First we observe that (1) implies that the conditional distribution

$$P(\theta_{n+1} \leq a_1, \xi_{n+1} \leq x | \mathcal{F}_n^\xi, \theta_n = b)$$

is Gaussian with mean-value vector

$$\mathbb{A}_0 + \mathbb{A}_1 b = \begin{pmatrix} a_0 + a_1 b \\ \mathbb{A}_0 + \mathbb{A}_1 b \end{pmatrix}$$

and covariance matrix

$$\mathbb{B} = \begin{pmatrix} b \circ b & b \circ B \\ (b \circ B)^* & B \circ B \end{pmatrix},$$

where  $b \circ b = b_1 b_1^* + b_2 b_2^*$ ,  $b \circ B = b_1 B_1^* + b_2 B_2^*$ ,  $B \circ B = B_1 B_1^* + B_2 B_2^*$ .

Let  $\zeta_n = (\theta_n, \xi_n)$  and  $t = (t_1, \dots, t_{k+1})$ . Then

$$E[\exp(it^* \zeta_{n+1}) | \mathcal{F}_n^\xi, \theta_n] = \exp\{it^*(\mathbb{A}_0(n, \xi) + \mathbb{A}_1(n, \xi)\theta_n) - \frac{1}{2}t^* \mathbb{B}(n, \xi)t\}. \quad (3)$$

Suppose now that the conclusion of the lemma holds for some  $n \geq 0$ . Then

$$E[\exp(it^* \mathbb{A}_1(n, \xi)\theta_n) | \mathcal{F}_n^\xi] = \exp(it^* \mathbb{A}_1(n, \xi)m_n - \frac{1}{2}t^*(\mathbb{A}_1(n, \xi)\gamma_n \mathbb{A}_1^*(n, \xi))t). \quad (4)$$

Let us show that (4) is also valid when  $n$  is replaced by  $n + 1$ .

From (3) and (4), we have

$$E[\exp(it^* \zeta_{n+1}) | \mathcal{F}_n^\xi] = \exp\{it^*(\mathbb{A}_0(n, \xi) + \mathbb{A}_1(n, \xi)m_n) - \frac{1}{2}t^* \mathbb{B}(n, \xi)t - \frac{1}{2}t^*(\mathbb{A}_1(n, \xi)\gamma_n \mathbb{A}_1^*(n, \xi))t\}.$$

Hence the conditional distribution

$$P(\theta_{n+1} \leq a, \xi_{n+1} \leq x | \mathcal{F}_n^\xi) \quad (5)$$

is Gaussian.

As in the proof of the theorem on normal correlation (Theorem 2, §13, Chapter II) we can verify that there is a matrix  $C$  such that the vector

$$\eta = [\theta_{n+1} - E(\theta_{n+1} | \mathcal{F}_n^\xi)] - C[\xi_{n+1} - E(\xi_{n+1} | \mathcal{F}_n^\xi)]$$

has the property that (P-a.s.)

$$E[\eta(\xi_{n+1} - E(\xi_{n+1} | \mathcal{F}_n^\xi))^* | \mathcal{F}_n^\xi] = 0.$$

It follows that the conditionally-Gaussian vectors  $\eta$  and  $\xi_{n+1}$ , considered under the condition  $\mathcal{F}_n^\xi$ , are independent, i.e.

$$P(\eta \in A, \xi_{n+1} \in B | \mathcal{F}_n^\xi) = P(\eta \in A | \mathcal{F}_n^\xi) \cdot P(\xi_{n+1} \in B | \mathcal{F}_n^\xi)$$

for all  $A \in \mathcal{B}(R^k)$ ,  $B \in \mathcal{B}(R^l)$ .

Therefore if  $s = (s_1, \dots, s_n)$  then

$$\begin{aligned} E[\exp(is^* \theta_{n+1}) | \mathcal{F}_n^\xi, \xi_{n+1}] \\ &= E\{\exp(is^* [E(\theta_{n+1} | \mathcal{F}_n^\xi) + \eta + C[\xi_{n+1} - E(\xi_{n+1} | \mathcal{F}_n^\xi)])] | \mathcal{F}_n^\xi, \xi_{n+1}\} \\ &= \exp\{is^* [E(\theta_{n+1} | \mathcal{F}_n^\xi) + C[\xi_{n+1} - E(\xi_{n+1} | \mathcal{F}_n^\xi)]]\} \\ &\quad \times E[\exp(is^* \eta) | \mathcal{F}_n^\xi, \xi_{n+1}] \\ &= \exp\{is^* [E(\theta_{n+1} | \mathcal{F}_n^\xi) + C[\xi_{n+1} - E(\xi_{n+1} | \mathcal{F}_n^\xi)]]\} \\ &\quad \times E(\exp(is^* \eta) | \mathcal{F}_n^\xi). \end{aligned} \quad (6)$$

By (5), the conditional distribution  $P(\eta \leq y | \mathcal{F}_n^\xi)$  is Gaussian. With (6), this shows that the conditional distribution  $P(\theta_{n+1} \leq a | \mathcal{F}_{n+1}^\xi)$  is also Gaussian.

This completes the proof of the lemma.

**Theorem 1.** Let  $(\theta, \xi)$  be a partial observation of a sequence that satisfies the system (1) and condition (2). Then  $(m_n, \gamma_n)$  obey the following recursion relations:

$$\begin{aligned} m_{n+1} &= [a_0 + a_1 m_n] + [b \circ B + a_1 \gamma_n A_1^*] [B \circ B + A_1 \gamma_n A_1^*]^\oplus \\ &\quad \times [\xi_{n+1} - A_0 - A_1 m_n], \end{aligned} \quad (7)$$

$$\begin{aligned} \gamma_{n+1} &= [a_1 \gamma_n A_1^* + b \circ b] - [b \circ B + a_1 \gamma_n A_1^*] [B \circ B + A_1 \gamma_n A_1^*]^\oplus \\ &\quad \times [b \circ B + a_1 \gamma_n A_1^*]^*. \end{aligned} \quad (8)$$

PROOF. From (1),

$$E(\theta_{n+1} | \mathcal{F}_n^\xi) = a_0 + a_1 m_n, \quad E(\xi_{n+1} | \mathcal{F}_n^\xi) = A_0 + A_1 m_n \quad (9)$$

and

$$\begin{aligned} \theta_{n+1} - E(\theta_{n+1} | \mathcal{F}_n^\xi) &= a_1 [\theta_n - m_n] + b_1 \varepsilon_1(n+1) + b_2 \varepsilon_2(n+1), \\ \xi_{n+1} - E(\xi_{n+1} | \mathcal{F}_n^\xi) &= A_1 [\theta_n - m_n] + B_1 \varepsilon_1(n+1) + B_2 \varepsilon_2(n+1). \end{aligned} \quad (10)$$

Let us write

$$\begin{aligned} d_{11} &= \text{cov}(\theta_{n+1}, \theta_{n+1} | \mathcal{F}_n^\xi) \\ &= E\{[\theta_{n+1} - E(\theta_{n+1} | \mathcal{F}_n^\xi)][\theta_{n+1} - E(\theta_{n+1} | \mathcal{F}_n^\xi)]^* | \mathcal{F}_n^\xi\}, \\ d_{12} &= \text{cov}(\theta_{n+1}, \xi_{n+1} | \mathcal{F}_n^\xi) \\ &= E\{[\theta_{n+1} - E(\theta_{n+1} | \mathcal{F}_n^\xi)][\xi_{n+1} - E(\xi_{n+1} | \mathcal{F}_n^\xi)]^* | \mathcal{F}_n^\xi\}, \\ d_{22} &= \text{cov}(\xi_{n+1}, \xi_{n+1} | \mathcal{F}_n^\xi) \\ &= E\{[\xi_{n+1} - E(\xi_{n+1} | \mathcal{F}_n^\xi)][\xi_{n+1} - E(\xi_{n+1} | \mathcal{F}_n^\xi)]^* | \mathcal{F}_n^\xi\}. \end{aligned}$$



Then, by (10),

$$d_{11} = a_1 \gamma_n a_1^* + b \circ b, \quad d_{12} = a_1 \gamma_n A_1^* + b \circ B, \quad d_{22} = A_1 \gamma_n A_1^* + B \circ B. \quad (11)$$

By the theorem on normal correlation (see Theorem 2 and Problem 4, §13, Chapter II),

$$m_{n+1} = E(\theta_{n+1} | \mathcal{F}_n^\xi, \xi_{n+1}) = E(\theta_{n+1} | \mathcal{F}_n^\xi) + d_{12} d_{22}^{\oplus} (\xi_{n+1} - E(\xi_{n+1} | \mathcal{F}_n^\xi))$$

and

$$\gamma_{n+1} = \text{cov}(\theta_{n+1}, \theta_{n+1} | \mathcal{F}_n^\xi, \xi_{n+1}) = d_{11} - d_{12} d_{22}^{\oplus} d_{12}^*.$$

If we then use the expressions from (9) for  $E(\theta_{n+1} | \mathcal{F}_n^\xi)$  and  $E(\xi_{n+1} | \mathcal{F}_n^\xi)$  and those for  $d_{11}$ ,  $d_{12}$ ,  $d_{22}$  from (11), we obtain the required recursion formulas (7) and (8).

This completes the proof of the theorem.

**Corollary 1.** *If the coefficients  $a_0(n, \xi)$ ,  $\dots$ ,  $B_2(n, \xi)$  in (1) are independent of  $\xi$  the corresponding method is known as the Kalman-Bucy method, and equations (7) and (8) for  $m_n$  and  $\gamma_n$  describe the Kalman-Bucy filter. It is important to observe that in this case the conditional and unconditional error matrices  $\gamma_n$  agree, i.e.*

$$\gamma_n \equiv E\gamma_n = E[(\theta_n - m_n)(\theta_n - m_n)^*].$$

**Corollary 2.** *Suppose that a partially observed sequence  $(\theta_n, \xi_n)$  has the property that  $\theta_n$  satisfies the first equation (1), and that  $\xi_n$  satisfies the equation*

$$\begin{aligned} \xi_n &= \tilde{A}_0(n-1, \xi) + \tilde{A}_1(n-1, \xi)\theta_n \\ &\quad + \tilde{B}_1(n-1, \xi)\varepsilon_1(n) + \tilde{B}_2(n-1, \xi)\varepsilon_2(n). \end{aligned} \quad (12)$$

Then evidently

$$\begin{aligned} \xi_{n+1} &= \tilde{A}_0(n, \xi) + \tilde{A}_1(n, \xi)[a_0(n, \xi) + a_1(n, \xi)\theta_n \\ &\quad + b_1(n, \xi)\varepsilon_1(n+1) + b_2(n, \xi)\varepsilon_2(n+1)] + \tilde{B}_1(n, \xi)\varepsilon_1(n+1) \\ &\quad + \tilde{B}_2(n, \xi)\varepsilon_2(n+1), \end{aligned}$$

and with the notation

$$\begin{aligned} A_0 &= A_0 + \tilde{A}_1 a_0, & A_1 &= \tilde{A}_1 a_1, \\ B_1 &= A_1 b_1 + \tilde{B}_1, & B_2 &= \tilde{A}_1 b_2 + \tilde{B}_2, \end{aligned}$$

we find that the case under consideration also depends on the model (1), and that  $m_n$  and  $\gamma_n$  satisfy (7) and (8).

2. We now consider a linear model (compare (1))

$$\begin{aligned} \theta_{n+1} &= a_0 + a_1 \theta_n + a_2 \xi_n + b_1 \varepsilon_1(n+1) + b_2 \varepsilon_2(n+1), \\ \xi_{n+1} &= A_0 + A_1 \theta_n + A_2 \xi_n + B_1 \varepsilon_1(n+1) + B_2 \varepsilon_2(n+1), \end{aligned} \quad (13)$$

where the coefficients  $a_0, \dots, B_n$  may depend on  $n$  (but not on  $\xi$ ), and  $\varepsilon_{ij}(n)$  are independent Gaussian random variables with  $E\varepsilon_{ij}(n) = 0$  and  $E\varepsilon_{ij}^2(n) = 1$ .

Let (13) be solved for the initial values  $(\theta_0, \xi_0)$  so that the conditional distribution  $P(\theta_0 \leq a | \xi_0)$  is Gaussian with parameters  $m_0 = E(\theta_0, \xi_0)$  and  $\gamma = \text{cov}(\theta_0, \theta_0 | \xi_0) = E\gamma_0$ . Then, by the theorem on normal correlation and (7) and (8), the optimal estimator  $m_n = E(\theta_n | \mathcal{F}_n^{\xi})$  is a linear function of  $\xi_0, \xi_1, \dots, \xi_n$ .

This remark makes it possible to prove the following important statement about the structure of the optimal linear filter without the assumption that it is Gaussian.

**Theorem 2.** Let  $(\theta, \xi) = (\theta_n, \xi_n)_{n \geq 0}$  be a partially observed sequence that satisfies (13), where  $\varepsilon_{ij}(n)$  are uncorrelated random variables with  $E\varepsilon_{ij}(n) = 0$ ,  $E\varepsilon_{ij}^2(n) = 1$ , and the components of the initial vector  $(\theta_0, \xi_0)$  have finite second moments. Then the optimal linear estimator  $\hat{m}_n = E(\theta_n | \xi_0, \dots, \xi_n)$  satisfies (7) with  $a_0(n, \xi) = a_0(n) + a_2(n)\xi_n$ ,  $A_0(n, \xi) = A_0(n) + A_2(n)\xi_n$ , and the error matrix  $\hat{\gamma}_n = E[(\theta_n - m_n)(\theta_n - m_n)^*]$  satisfies (8) with initial values

$$\begin{aligned}\hat{m}_0 &= \text{cov}(\theta_0, \xi_0) \text{cov}^{\oplus}(\xi_0, \xi_0) \cdot \xi_0, \\ \hat{\gamma}_0 &= \text{cov}(\theta_0, \theta_0) - \text{cov}(\theta_0, \xi_0) \text{cov}^{\oplus}(\xi_0, \xi_0) \text{cov}^*(\theta_0, \xi_0).\end{aligned}\quad (14)$$

For the proof of this lemma, we need the following lemma, which reveals the role of the Gaussian case in determining optimal linear estimators.

**Lemma 2.** Let  $(\alpha, \beta)$  be a two-dimensional random vector with  $E(\alpha^2 + \beta^2) < \infty$ ,  $\alpha(\tilde{\alpha}, \tilde{\beta})$  a two-dimensional Gaussian vector with the same first and second moments as  $(\alpha, \beta)$ , i.e.

$$E\tilde{\alpha}^i = E\alpha^i, \quad E\tilde{\beta}^i = E\beta^i, \quad i = 1, 2; \quad E\tilde{\alpha}\tilde{\beta} = E\alpha\beta.$$

Let  $\lambda(b)$  be a linear function of  $b$  such that

$$\lambda(b) = E(\tilde{\alpha} | \tilde{\beta} = b).$$

Then  $\lambda(\beta)$  is the optimal (in the mean square sense) linear estimator of  $\alpha$  in terms of  $\beta$ , i.e.

$$\hat{E}(\alpha | \beta) = \lambda(\beta).$$

Here  $E\lambda(\beta) = E\alpha$ .

**PROOF.** We first observe that the existence of a linear function  $\lambda(b)$  coinciding with  $E(\tilde{\alpha} | \tilde{\beta} = b)$  follows from the theorem on normal correlation. Moreover, let  $\bar{\lambda}(b)$  be any other linear estimator. Then

$$E[\tilde{\alpha} - \bar{\lambda}(\tilde{\beta})]^2 \geq E[\tilde{\alpha} - \lambda(\tilde{\beta})]^2$$

and since  $\bar{\lambda}(b)$  and  $\lambda(b)$  are linear and the hypotheses of the lemma are satisfied, we have

$$E[\alpha - \bar{\lambda}(\beta)]^2 = E[\tilde{\alpha} - \bar{\lambda}(\tilde{\beta})]^2 \geq E[\tilde{\alpha} - \lambda(\tilde{\beta})]^2 = E[\alpha - \lambda(\beta)]^2,$$

which shows that  $\lambda(\beta)$  is optimal in the class of linear estimators. Finally,

$$E\lambda(\beta) = E\lambda(\tilde{\beta}) = E[E(\tilde{\alpha}|\tilde{\beta})] = E\tilde{\alpha} = E\alpha.$$

This completes the proof of the lemma.

**PROOF OF THEOREM 2.** We consider, besides (13), the system

$$\begin{aligned}\tilde{\theta}_{n+1} &= a_0 + a_1\tilde{\theta}_n + a_2\tilde{\xi}_n + b_1\tilde{\varepsilon}_{11}(n+1) + b_2\tilde{\varepsilon}_{12}(n+1), \\ \tilde{\xi}_{n+1} &= A_0 + A_1\tilde{\theta}_n + A_2\tilde{\xi}_n + B_1\tilde{\varepsilon}_{21}(n+1) + B_2\tilde{\varepsilon}_{22}(n+1),\end{aligned}\quad (15)$$

where  $\tilde{\varepsilon}_{ij}(n)$  are independent Gaussian random variables with  $E\tilde{\varepsilon}_{ij}(n) = 0$  and  $E\tilde{\varepsilon}_{ij}^2(n) = 1$ . Let  $(\tilde{\theta}_0, \tilde{\xi}_0)$  also be a Gaussian vector which has the same first moment and covariance as  $(\theta_0, \xi_0)$  and is independent of  $\tilde{\varepsilon}_{ij}(n)$ . Then since (15) is linear, the vector  $(\tilde{\theta}_0, \dots, \tilde{\theta}_n, \tilde{\xi}_0, \dots, \tilde{\xi}_n)$  is Gaussian and therefore the conclusion of the theorem follows from Lemma 2 (more precisely, from its multidimensional analog) and the theorem on normal covariance.

This completes the proof of the theorem.

### 3. Let us consider some illustrations of Theorems 1 and 2.

**EXAMPLE 1.** Let  $\theta = (\theta_n)$  and  $\eta = (\eta_n)$  be two stationary (wide sense) uncorrelated random sequences with  $E\theta_n = E\eta_n = 0$  and spectral densities

$$f_\theta(\lambda) = \frac{1}{2\pi|1 + b_1e^{-i\lambda}|^2} \quad \text{and} \quad f_\eta(\lambda) = \frac{1}{2\pi} \cdot \frac{1}{|1 + b_2e^{-i\lambda}|^2},$$

where  $|b_1| < 1, |b_2| < 1$ .

We are going to interpret  $\theta$  as a useful signal and  $\eta$  as noise, and suppose that observation produces a sequence  $\xi = (\xi_n)$  with

$$\xi_n = \theta_n + \eta_n.$$

According to Corollary 2 to Theorem 3 of §3 there are (mutually uncorrelated) white noises  $\varepsilon_1 = (\varepsilon_1(n))$  and  $\varepsilon_2 = (\varepsilon_2(n))$  such that

$$\theta_{n+1} + b_1\theta_n = \varepsilon_1(n+1), \quad \eta_{n+1} + b_2\eta_n = \varepsilon_2(n+1).$$

Then

$$\begin{aligned}\xi_{n+1} &= \theta_{n+1} + \eta_{n+1} = -b_1\theta_n - b_2\eta_n + \varepsilon_1(n+1) + \varepsilon_2(n+1) \\ &= -b_2(\theta_n + \eta_n) - \theta_n(b_1 - b_2) + \varepsilon_1(n+1) + \varepsilon_2(n+1) \\ &= -b_2\xi_n - (b_1 - b_2)\theta_n + \varepsilon_1(n+1) + \varepsilon_2(n+1).\end{aligned}$$

Hence  $\theta$  and  $\xi$  satisfy the recursion relations

$$\begin{aligned}\theta_{n+1} &= -b_1\theta_n + \varepsilon_1(n+1), \\ \xi_{n+1} &= -(b_1 - b_2)\theta_n - b_2\xi_n + \varepsilon_1(n+1) + \varepsilon_2(n+1),\end{aligned}\quad (16)$$

and, according to Theorem 2,  $m_n = \hat{E}(\theta_n | \xi_0, \dots, \xi_n)$  and  $\gamma_n = E(\theta_n - m_n)^2$  satisfy the following system of recursion equations for optimal linear filtering:

$$\begin{aligned}m_{n+1} &= -b_1m_n + \frac{b_1(b_1 - b_2)\gamma_n}{2 + (b_1 - b_2)^2\gamma_n} [\xi_{n+1} + (b_1 - b_2)m_n + b_2\xi_n], \\ \gamma_{n+1} &= b_1^2\gamma_n + 1 - \frac{[1 + b_1(b_1 - b_2)\gamma_n]^2}{2 + (b_1 - b_2)^2\gamma_n}.\end{aligned}\quad (17)$$

Let us find the initial conditions under which we should solve this system. Write  $d_{11} = E\theta_n^2$ ,  $d_{12} = E\theta_n\xi_n$ ,  $d_{22} = E\xi_n^2$ . Then we find from (16) that

$$\begin{aligned}d_{11} &= b_1^2d_{11} + 1, \\ d_{12} &= b_1(b_1 - b_2)d_{11} + b_1b_2d_{12} + 1, \\ d_{22} &= (b_1 - b_2)^2d_{11} + b_2^2d_{22} + 2b_2(b_1 - b_2)d_{12} + 2,\end{aligned}$$

from which

$$d_{11} = \frac{1}{1 - b_1^2}, \quad d_{12} = \frac{1}{1 - b_1^2}, \quad d_{22} = \frac{2 - b_1^2 - b_2^2}{(1 - b_1^2)(1 - b_2^2)},$$

which, by (14), leads to the following initial values:

$$\begin{aligned}m_0 &= \frac{d_{12}}{d_{22}} \xi_0 = \frac{1 - b_2^2}{2 - b_1^2 - b_2^2} \xi_0, \\ \gamma_0 &= d_{11} - \frac{d_{12}^2}{d_{22}} = \frac{1}{1 - b_1^2} - \frac{1 - b_2^2}{(1 - b_1^2)(2 - b_1^2 - b_2^2)} = \frac{1}{2 - b_1^2 - b_2^2}.\end{aligned}\quad (18)$$

Thus the optimal (in the least squares sense) linear estimators  $m_n$  for the signal  $\theta_n$  in terms of  $\xi_0, \dots, \xi_n$  and the mean-square error are determined by the system of recurrent equations (17), solved under the initial conditions (18). Observe that the equation for  $\gamma_n$  does not contain any random components, and consequently the number  $\gamma_n$ , which is needed for finding  $m_n$ , can be calculated in advance, before the filtering problem has been solved.

**EXAMPLE 2.** This example is instructive because it shows that the result of Theorem 2 can be applied to find the optimal linear filter in a case where the sequence  $(\theta, \xi)$  is described by a (nonlinear) system which is different from (13).

Let  $\varepsilon_1 = (\varepsilon_1(n))$  and  $\varepsilon_2 = (\varepsilon_2(n))$  be two independent Gaussian sequences of independent random variables with  $E\varepsilon_i(n) = 0$  and  $E\varepsilon_i^2(n) = 1$ ,  $n \geq 1$ . Consider a pair of sequences  $(\theta, \xi) = (\theta_n, \xi_n)$ ,  $n \geq 0$ , with

$$\begin{aligned}\theta_{n+1} &= a\theta_n + (1 + \theta_n)\varepsilon_1(n+1), \\ \xi_{n+1} &= A\theta_n + \varepsilon_2(n+1).\end{aligned}\quad (19)$$

We shall suppose that  $\theta_0$  is independent of  $(\varepsilon_1, \varepsilon_2)$  and that  $\theta_0 \sim \mathcal{N}(m_0, \gamma_0)$ .

The system (19) is *nonlinear*, and Theorem 2 is not immediately applicable. However, if we put

$$\tilde{\varepsilon}_1(n+1) = \frac{1 + \theta_n}{\sqrt{E(1 + \theta_n)^2}} \varepsilon_1(n+1),$$

we can observe that  $E\tilde{\varepsilon}_1(n) = 0$ ,  $E\tilde{\varepsilon}_1(n)\tilde{\varepsilon}_1(m) = 0$ ,  $n \neq m$ ,  $E\tilde{\varepsilon}_1^2(n) = 1$ . Hence we have reduced (19) to a linear system

$$\begin{aligned}\theta_{n+1} &= a_1\theta_n + b_1\tilde{\varepsilon}_1(n+1), \\ \xi_{n+1} &= A_1\theta_n + \varepsilon_2(n+1),\end{aligned}\quad (20)$$

where  $b_1 = \sqrt{E(1 + \theta_n)^2}$ , and  $\{\tilde{\varepsilon}_1(n)\}$  is a sequence of uncorrelated random variables.

Now (20) is a linear system of the same type as (13), and consequently the optimal linear estimator  $\hat{m}_n = \hat{E}(\theta_n | \xi_0, \dots, \xi_n)$  and its error  $\hat{\gamma}_n$  can be determined from (7) and (8) via Theorem 2, applied in the following form in the present case:

$$\begin{aligned}m_{n+1} &= a_1m_n + \frac{a_1A_1\gamma_n}{1 + A_1^2\gamma_n} [\xi_{n+1} - A_1m_n], \\ \gamma_{n+1} &= (a_1^2\gamma_n + b_1^2) - \frac{(a_1A_1\gamma_n)^2}{1 + A_1^2\gamma_n},\end{aligned}$$

where  $b_1 = \sqrt{E(1 + \theta_n)^2}$  must be found from the first equation in (19).

**EXAMPLE 3. Estimators for parameters.** Let  $\theta = (\theta_1, \dots, \theta_k)$  be a Gaussian vector with  $E\theta = m$  and  $\text{cov}(\theta, \theta) = \gamma$ . Suppose that (with known  $m$  and  $v$ ) we want the optimal estimator of  $\theta$  in terms of observations on an  $l$ -dimensional sequence  $\xi = (\xi_n)$ ,  $n \geq 0$ , with

$$\xi_{n+1} = A_0(n, \xi) + A_1(n, \xi)\theta + B_1(n, \xi)\varepsilon_1(n+1), \quad \xi_0 = 0, \quad (21)$$

where  $\varepsilon_1$  is as in (1).

Then from (7) and (8), with  $m_n = E(\theta | \mathcal{F}_n^{\xi})$  and  $\gamma_n$ , we find that

$$\begin{aligned}m_{n+1} &= m_n + \gamma_n A_1^*(n, \xi) [(B_1 B_1^*)(n, \xi) + A_1(n, \xi) \gamma_n A_1^*(n, \xi)]^{\oplus} \\ &\quad \times [\xi_{n+1} - A_0(n, \xi) - A_1(n, \xi) m_n], \\ \gamma_{n+1} &= \gamma_n - \gamma_n A_1^*(n, \xi) [(B_1 B_1^*)(n, \xi) + A_1(n, \xi) \gamma_n A_1^*(n, \xi)]^{\oplus} A_1(n, \xi) \gamma_n\end{aligned}\quad (22)$$

If the matrices  $B_1 B_1^*$  are nonsingular, the solution of (22) is given by

$$\begin{aligned} m_{n+1} &= \left[ E + \gamma \sum_{m=0}^n A_1^*(m, \xi) (B_1 B_1^*)^{-1} (m, \xi) A_1^*(m, \xi) \right]^{-1} \\ &\quad \times \left[ m + \gamma \sum_{m=0}^n A_1^*(m, \xi) (B_1 B_1^*)^{-1} (m, \xi) (\xi_{m+1} - A_0(m, \xi)) \right], \\ \gamma_{n+1} &= \left[ E + \gamma \sum_{m=0}^n A_1^*(m, \xi) (B_1 B_1^*)^{-1} (m, \xi) A_1(m, \xi) \right]^{-1} \gamma, \end{aligned} \quad (23)$$

where  $E$  is a unit matrix.

#### 4. PROBLEMS

1. Show that the vectors  $m_n$  and  $\theta_n - m_n$  in (1) are uncorrelated:

$$E[m_n^*(\theta - m_n)] = 0.$$

2. In (1), let  $\gamma$  and the coefficients other than  $a_0(n, \xi)$  and  $A_0(n, \xi)$  be independent of "chance" (i.e. of  $\xi$ ). Show that then the conditional covariance  $\gamma_n$  is independent of "chance":  $\gamma_n = E\gamma_n$ .

3. Show that the solution of (22) is given by (23).

4. Let  $(\theta, \xi) = (\theta_n, \xi_n)$  be a Gaussian sequence satisfying the following special case of (1):

$$\theta_{n+1} = a\theta_n + b\varepsilon_1(n+1), \quad \xi_{n+1} = A\theta_n + B\varepsilon_2(n+1).$$

Show that if  $A \neq 0$ ,  $b \neq 0$ ,  $B \neq 0$ , the limiting error of filtering,  $\gamma = \lim_{n \rightarrow \infty} \gamma_n$ , exists and is determined as the positive root of the equation

$$\gamma^2 + \left[ \frac{B^2(1 - a^2)}{A^2} - b^2 \right] \gamma - \frac{b^2 B^2}{A^2} = 0.$$

## CHAPTER VII

# Sequences of Random Variables That Form Martingales

### §1. Definitions of Martingales and Related Concepts

1. The study of the dependence of random variables arises in various ways in probability theory. In the theory of stationary (wide sense) random sequences, the basic indicator of dependence is the covariance function, and the inferences made in this theory are determined by the properties of that function. In the theory of Markov chains (§12 of Chapter I; Chapter VIII) the basic dependence is supplied by the transition function, which completely determines the development of the random variables involved in Markov dependence.

In the present chapter (see also §11, Chapter I), we single out a rather wide class of sequences of random variables (martingales and their generalizations) for which dependence can be studied by methods based on a discussion of the properties of conditional expectations.

2. Let  $(\Omega, \mathcal{F}, P)$  be a given probability space, and let  $(\mathcal{F}_n)$  be a family of  $\sigma$ -algebras  $\mathcal{F}_n, n \geq 0$ , such that  $\mathcal{F}_0 \subseteq \mathcal{F}_1 \subseteq \dots \subseteq \mathcal{F}$ .

Let  $X_0, X_1, \dots$  be a sequence of random variables defined on  $(\Omega, \mathcal{F}, P)$ . If, for each  $n \geq 0$ , the variable  $X_n$  is  $\mathcal{F}_n$ -measurable, we say that the set  $X = (X_n, \mathcal{F}_n), n \geq 0$ , or simply  $X = (X_n, \mathcal{F}_n)$ , is a *stochastic sequence*.

If a stochastic sequence  $X = (X_n, \mathcal{F}_n)$  has the property that, for each  $n \geq 1$ , the variable  $X_n$  is  $\mathcal{F}_{n-1}$ -measurable, we write  $X = (X_n, \mathcal{F}_{n-1})$ , taking  $\mathcal{F}_{-1} = \mathcal{F}_0$ , and call  $X$  a *predictable sequence*. We call such a sequence *increasing* if  $X_0 = 0$  and  $X_n \leq X_{n+1}$  (P-a.s.).

**Definition 1.** A stochastic sequence  $X = (X_n, \mathcal{F}_n)$  is a *martingale*, or a *sub-martingale*, if, for all  $n \geq 0$ ,

$$E|X_n| < \infty \quad (1)$$

and, respectively,

$$\mathbf{E}(X_{n+1}|\mathcal{F}_n) = X_n \quad (\text{P-a.s.}) \quad (\text{martingale})$$

or

$$\mathbf{E}(X_{n+1}|\mathcal{F}_n) \geq X_n \quad (\text{P-a.s.}) \quad (\text{submartingale}).$$

(2)

A stochastic sequence  $X = (X_n, \mathcal{F}_n)$  is a *supermartingale* if the sequence  $-X = (-X_n, \mathcal{F}_n)$  is a submartingale.

In the special case when  $\mathcal{F}_n = \mathcal{F}_n^X$ , where  $\mathcal{F}_n^X = \sigma\{\omega: X_0, \dots, X_n\}$ , and the stochastic sequence  $X = (X_n, \mathcal{F}_n)$  is a martingale (or submartingale), we say that the sequence  $(X_n)_{n \geq 0}$  itself is a martingale (or submartingale).

It is easy to deduce from the properties of conditional expectations that (2) is equivalent to the property that, for every  $n \geq 0$  and  $A \in \mathcal{F}_n$ ,

$$\int_A X_{n+1} dP = \int_A X_n dP$$

or

$$\int_A X_{n+1} dP \geq \int_A X_n dP.$$

(3)

EXAMPLE 1. If  $(\xi_n)_{n \geq 0}$  is a sequence of independent random variables with  $\mathbf{E}\xi_n = 0$  and  $X_n = \xi_0 + \dots + \xi_n$ ,  $\mathcal{F}_n = \sigma\{\omega: \xi_0, \dots, \xi_n\}$ , the stochastic sequence  $X = (X_n, \mathcal{F}_n)$  is a martingale.

EXAMPLE 2. If  $(\xi_n)_{n \geq 0}$  is a sequence of independent random variables with  $\mathbf{E}\xi_n = 1$ , the stochastic sequence  $(X_n, \mathcal{F}_n)$  with  $X_n = \prod_{k=0}^n \xi_k$ ,  $\mathcal{F}_n = \sigma\{\omega: \xi_0, \dots, \xi_n\}$  is also a martingale.

EXAMPLE 3. Let  $\xi$  be a random variable with  $\mathbf{E}|\xi| < \infty$  and

$$\mathcal{F}_0 \subseteq \mathcal{F}_1 \subseteq \dots \subseteq \mathcal{F}.$$

Then the sequence  $X = (X_n, \mathcal{F}_n)$  with  $X_n = \mathbf{E}(\xi|\mathcal{F}_n)$  is a martingale.

EXAMPLE 4. If  $(\xi_n)_{n \geq 0}$  is a sequence of nonnegative integrable random variables, the sequence  $(X_n)$  with  $X_n = \xi_0 + \dots + \xi_n$  is a submartingale.

EXAMPLE 5. If  $X = (X_n, \mathcal{F}_n)$  is a martingale and  $g(x)$  is convex downward with  $\mathbf{E}|g(X_n)| < \infty$ ,  $n \geq 0$ , then the stochastic sequence  $(g(X_n), \mathcal{F}_n)$  is a submartingale (as follows from Jensen's inequality).

If  $X = (X_n, \mathcal{F}_n)$  is a submartingale and  $g(x)$  is convex downward and nondecreasing, with  $\mathbf{E}|g(X_n)| < \infty$  for all  $n \geq 0$ , then  $(g(X_n), \mathcal{F}_n)$  is also a submartingale.

Assumption (1) in Definition 1 ensures the existence of the conditional expectations  $\mathbf{E}(X_{n+1}|\mathcal{F}_n)$ ,  $n \geq 0$ . However, these expectations can also exist without the assumption that  $\mathbf{E}|X_{n+1}| < \infty$ . Recall that by §7 of Chapter



II,  $E(X_{n+1}^+|\mathcal{F}_n)$  and  $E(X_{n+1}^-|\mathcal{F}_n)$  are always defined. Let us write  $A = B$  (P-a.s.) when  $P(A \triangle B) = 0$ . Then if

$$\{\omega: E(X_{n+1}^+|\mathcal{F}_n) < \infty\} \cup \{\omega: E(X_{n+1}^-|\mathcal{F}_n) < \infty\} = \Omega \quad (\text{P-a.s.})$$

we say that  $E(X_{n+1}|\mathcal{F}_n)$  is also defined and is given by

$$E(X_{n+1}|\mathcal{F}_n) = E(X_{n+1}^+|\mathcal{F}_n) - E(X_{n+1}^-|\mathcal{F}_n).$$

After this, the following definition is natural.

**Definition 2.** A stochastic sequence  $X = (X_n, \mathcal{F}_n)$  is a *generalized martingale* (or *submartingale*) if the conditional expectations  $E(X_{n+1}|\mathcal{F}_n)$  are defined for every  $n \geq 0$  and (2) is satisfied.

Notice that it follows from this definition that  $E(X_{n+1}^-|\mathcal{F}_n) < \infty$  for a generalized submartingale, and the  $E(|X_{n+1}||\mathcal{F}_n) < \infty$  (P-a.s.) for a generalized martingale.

3. In the following definition we introduce the concept of a Markov time, which plays a very important role in the subsequent theory.

**Definition 3.** A random variable  $\tau = \tau(\omega)$  with values in the set  $\{0, 1, \dots, +\infty\}$  is a *Markov time* (with respect to  $(\mathcal{F}_n)$ ) (or a *random variable independent of the future*) if, for each  $n \geq 0$ ,

$$\{\tau = n\} \in \mathcal{F}_n. \quad (4)$$

When  $P(\tau < \infty) = 1$ , a Markov time  $\tau$  is called a *stopping time*.

Let  $X = (X_n, \mathcal{F}_n)$  be a stochastic sequence and let  $\tau$  be a Markov time (with respect to  $(\mathcal{F}_n)$ ). We write

$$X_\tau = \sum_{n=0}^{\infty} X_n I_{\{\tau \geq n\}}(\omega)$$

(hence  $X_\tau = 0$  on the set  $\{\omega: \tau = \infty\}$ ).

Then for every  $B \in \mathcal{B}(R)$ ,

$$\{\omega: X_\tau \in B\} = \sum_{n=0}^{\infty} \{X_n \in B, \tau \geq n\} \in \mathcal{F},$$

and consequently  $X_\tau$  is a random variable.

**EXAMPLE 6.** Let  $X = (X_n, \mathcal{F}_n)$  be a stochastic sequence and let  $B \in \mathcal{B}(R)$ . Then the time of first hitting the set  $B$ , that is,

$$\tau_B = \inf\{n \geq 0: X_n \in B\}$$

(with  $\tau_B = +\infty$  if  $\{\cdot\} = \emptyset$ ) is a Markov time, since

$$\{\tau_B = n\} = \{X_0 \notin B, \dots, X_{n-1} \notin B, X_n \in B\} \in \mathcal{F}_n$$

for every  $n \geq 0$ .

**EXAMPLE 7.** Let  $X = (X_n, \mathcal{F}_n)$  be a martingale (or submartingale) and  $\tau$  a Markov time (with respect to  $(\mathcal{F}_n)$ ). Then the “stopped” process  $X^\tau = (X_{n \wedge \tau}, \mathcal{F}_n)$  is also a martingale (or submartingale).

In fact, the equation

$$X_{n \wedge \tau} = \sum_{m=0}^{n-1} X_m I_{\{\tau \geq m\}} + X_n I_{\{\tau \geq n\}}$$

implies that the variables  $X_{n \wedge \tau}$  are  $\mathcal{F}_n$ -measurable, are integrable, and satisfy

$$X_{(n+1) \wedge \tau} - X_{n \wedge \tau} = I_{\{\tau > n\}}(X_{n+1} - X_n),$$

whence

$$\mathbb{E}[X_{(n+1) \wedge \tau} - X_{n \wedge \tau} | \mathcal{F}_n] = I_{\{\tau > n\}} \mathbb{E}[X_{n+1} - X_n | \mathcal{F}_n] = 0 \quad (\text{or } \geq 0).$$

Every system  $(\mathcal{F}_n)$  and Markov time  $\tau$  corresponding to it generate a collection of sets

$$\mathcal{F}_\tau = \{A \in \mathcal{F} : A \cap \{\tau = n\} \in \mathcal{F}_n \text{ for all } n \geq 0\}.$$

It is clear that  $\Omega \in \mathcal{F}_\tau$  and  $\mathcal{F}_\tau$  is closed under countable unions. Moreover, if  $A \in \mathcal{F}_\tau$ , then  $\bar{A} \cap \{\tau = n\} = \{\tau = n\} \setminus (A \cap \{\tau = n\}) \in \mathcal{F}_n$  and therefore  $\bar{A} \in \mathcal{F}_\tau$ . Hence it follows that  $\mathcal{F}_\tau$  is a  $\sigma$ -algebra.

If we think of  $\mathcal{F}_n$  as a collection of events observed up to time  $n$  (inclusive), then  $\mathcal{F}_\tau$  can be thought of as a collection of events observed at the “random” time  $\tau$ .

It is easy to show (Problem 3) that the random variables  $\tau$  and  $X_\tau$  are  $\mathcal{F}_\tau$ -measurable.

**4. Definition 4.** A stochastic sequence  $X = (X_n, \mathcal{F}_n)$  is a *local martingale* (or *submartingale*) if there is a (localizing) sequence  $(\tau_k)_{k \geq 1}$  of Markov times such that  $\tau_k \leq \tau_{k+1}$  (P-a.s.),  $\tau_k \uparrow \infty$  (P-a.s.) as  $k \rightarrow \infty$ , and every “stopped” sequence  $X^{\tau_k} = (X_{\tau_k \wedge n} \cdot I_{\{\tau_k > 0\}}, \mathcal{F}_n)$  is a martingale (or submartingale).

In Theorem 1 below, we show that in fact the class of local martingales coincides with the class of generalized martingales. Moreover, every local martingale can be obtained by a “martingale transformation” from a martingale and a predictable sequence.

**Definition 5.** Let  $Y = (Y_n, \mathcal{F}_n)$  be a stochastic sequence and let  $V = (V_n, \mathcal{F}_{n-1})$  be a predictable sequence ( $\mathcal{F}_{-1} = \mathcal{F}_0$ ). The stochastic sequence  $V \cdot Y = ((V \cdot Y)_n, \mathcal{F}_n)$  with

$$(V \cdot Y)_n = V_0 Y_0 + \sum_{i=1}^n V_i \Delta Y_i, \quad (5)$$

where  $\Delta Y_i = Y_i - Y_{i-1}$ , is called the *transform of  $Y$  by  $V$* . If, in addition,  $Y$  is a martingale, we say that  $V \cdot Y$  is a *martingale transform*.

**Theorem 1.** Let  $X = (X_n, \mathcal{F}_n)_{n \geq 0}$  be a stochastic sequence and let  $X_0 = 0$  (P-a.s.). The following conditions are equivalent:

- (a)  $X$  is a local martingale;
- (b)  $X$  is a generalized martingale;
- (c)  $X$  is a martingale transform, i.e., there are a predictable sequence  $V = (V_n, \mathcal{F}_{n-1})$  with  $V_0 = 0$  and a martingale  $Y = (Y_n, \mathcal{F}_n)$  with  $Y_0 = 0$  such that  $X = V \cdot Y$ .

**PROOF.** (a)  $\Rightarrow$  (b). Let  $X$  be a local martingale and let  $(\tau_k)$  be a local sequence of Markov times for  $X$ . Then for every  $m \geq 0$

$$\mathbb{E}[X_{m \wedge \tau_k} | I_{\{\tau_k > 0\}}] < \infty, \quad (6)$$

and therefore

$$\mathbb{E}[X_{(n+1) \wedge \tau_k} | I_{\{\tau_k > n\}}] = \mathbb{E}[X_{n+1} | I_{\{\tau_k > n\}}] < \infty. \quad (7)$$

The random variable  $I_{\{\tau_k > n\}}$  is  $\mathcal{F}_n$ -measurable. Hence it follows from (7) that

$$\mathbb{E}[X_{n+1} | I_{\{\tau_k > n\}}, \mathcal{F}_n] = I_{\{\tau_k > n\}} \mathbb{E}[X_{n+1} | \mathcal{F}_n] < \infty \quad (\text{P-a.s.}).$$

Here  $I_{\{\tau_k > n\}} \rightarrow 1$  (P-a.s.),  $k \rightarrow \infty$ , and therefore

$$\mathbb{E}[X_{n+1} | \mathcal{F}_n] < \infty \quad (\text{P-a.s.}). \quad (8)$$

Under this condition,  $\mathbb{E}[X_{n+1} | \mathcal{F}_n]$  is defined, and it remains only to show that  $\mathbb{E}[X_{n+1} | \mathcal{F}_n] = X_n$  (P-a.s.).

To do this, we need to show that

$$\int_A X_{n+1} d\mathbb{P} = \int_A X_n d\mathbb{P}$$

for  $A \in \mathcal{F}_n$ . By Problem 7, §7, Chapter II, we have  $\mathbb{E}[X_{n+1} | \mathcal{F}_n] < \infty$  (P-a.s.) if and only if the measure  $\int_A X_{n+1} d\mathbb{P}$ ,  $A \in \mathcal{F}_n$ , is  $\sigma$ -finite. Let us show that the measure  $\int_A X_n d\mathbb{P}$ ,  $A \in \mathcal{F}_n$ , is also  $\sigma$ -finite.

Since  $X^{\tau_k}$  is a martingale,  $|X^{\tau_k}| = (|X_{\tau_k \wedge n}| I_{\{\tau_k > 0\}}, \mathcal{F}_n)$  is a submartingale, and therefore (since  $\{\tau_k > n\} \in \mathcal{F}_n$ )

$$\begin{aligned} \int_{A \cap \{\tau_k > n\}} |X_n| dP &= \int_{A \cap \{\tau_k > n\}} |X_{n \wedge \tau_k}| I_{\{\tau_k > 0\}} dP \\ &\leq \int_{A \cap \{\tau_k > n\}} |X_{(n+1) \wedge \tau_k}| I_{\{\tau_k > 0\}} dP = \int_{A \cap \{\tau_k > n\}} |X_{n+1}| dP. \end{aligned}$$

Letting  $k \rightarrow \infty$ , we have

$$\int_A |X_n| dP \leq \int_A |X_{n+1}| dP,$$

from which there follows the required  $\sigma$ -finiteness of the measures  $\int_A |X_n| dP$ ,  $A \in \mathcal{F}_n$ .

Let  $A \in \mathcal{F}_n$  have the property  $\int_A |X_{n+1}| dP < \infty$ . Then, by Lebesgue's theorem on dominated convergence, we may take limits in the relation

$$\int_{A \cap \{\tau_k > n\}} X_n dP = \int_{A \cap \{\tau_k > n\}} X_{n+1} dP,$$

which is valid since  $X$  is a local martingale. Therefore,

$$\int_A X_n dP = \int_A X_{n+1} dP$$

for all  $A \in \mathcal{F}_n$  such that  $\int_A |X_{n+1}| dP < \infty$ . It then follows that the preceding relation also holds for every  $A \in \mathcal{F}_n$ , and therefore,  $E(X_{n+1} | \mathcal{F}_n) = X_n$  (P-a.s.).

b)  $\Rightarrow$  c). Let  $\Delta X_n = X_n - X_{n-1}$ ,  $X_0 = 0$ , and  $V_0 = 0$ ,  $V_n = E[|\Delta X_n| | \mathcal{F}_{n-1}]$ ,  $n \geq 1$ . We set

$$W_n = V_n^{\oplus} \left( = \begin{cases} V_n^{-1}, & V_n \neq 0 \\ 0, & V_n = 0 \end{cases} \right), \quad Y_0 = 0$$

and  $Y_n = \sum_{i=1}^n W_i \Delta X_i$ ,  $n \geq 1$ . It is clear that

$$E[|\Delta Y_n| | \mathcal{F}_{n-1}] \leq 1, \quad E[\Delta Y_n | \mathcal{F}_{n-1}] = 0,$$

and consequently,  $Y = (Y_n, \mathcal{F}_n)$  is a martingale.

Consequently,  $Y = (Y_n, \mathcal{F}_n)$  is a martingale. Moreover,  $X_0 = V_0 \cdot Y_0 = 0$  and  $\Delta(V \cdot Y)_n = \Delta X_n$ . Therefore

$$X = V \cdot Y.$$

(c)  $\Rightarrow$  (a). Let  $X = V \cdot Y$  where  $V$  is a predictable sequence,  $Y$  is a martingale and  $V_0 = Y_0 = 0$ . Put

$$\tau_k = \inf\{n \geq 0: |V_{n+1}| > k\},$$

and suppose that  $\tau_k = \infty$  if the set  $\{\cdot\} = \emptyset$ . Since  $V_{n+1}$  is  $\mathcal{F}_n$ -measurable, the variables  $\tau_k$  are Markov times for every  $k \geq 1$ .

Consider a "stopped" sequence  $X^{\tau_k} = ((V \cdot Y)_{n \wedge \tau_k} I_{\{\tau_k > 0\}}, \mathcal{F}_n)$ . On the set  $\{\tau_k > 0\}$ , the inequality  $|V_{n \wedge \tau_k}| \leq k$  is in effect. Hence it follows that  $E|(V \cdot Y)_{n \wedge \tau_k} I_{\{\tau_k > 0\}}| < \infty$  for every  $n \geq 1$ . In addition, for  $n \geq 1$ ,

$$\begin{aligned} & \mathbf{E}\{[(V \cdot Y)_{(n+1) \wedge \tau_k} - (V \cdot Y)_{n \wedge \tau_k}]I_{\{\tau_k > 0\}} | \mathcal{F}_n\} \\ &= I_{\{\tau_k > 0\}} \cdot V_{(n+1) \wedge \tau_k} \cdot \mathbf{E}\{Y_{(n+1) \wedge \tau_k} - Y_{n \wedge \tau_k} | \mathcal{F}_n\} = 0 \end{aligned}$$

since (see Example 7)  $\mathbf{E}\{Y_{(n+1) \wedge \tau_k} - Y_{n \wedge \tau_k} | \mathcal{F}_n\} = 0$ .

Thus for every  $k \geq 1$  the stochastic sequences  $X^{\tau_k}$  are martingales,  $\tau_k \uparrow \infty$  (P-a.s.), and consequently  $X$  is a local martingale.

This completes the proof of the theorem.

**5. EXAMPLE 8.** Let  $(\eta_n)_{n \geq 1}$  be a sequence of independent identically distributed Bernoulli random variables and let  $\mathbf{P}(\eta_n = 1) = p$ ,  $\mathbf{P}(\eta_n = -1) = q$ ,  $p + q = 1$ . We interpret the event  $\{\eta_n = 1\}$  as success (gain) and  $\{\eta_n = -1\}$  as failure (loss) of a player at the  $n$ th turn. Let us suppose that the player's stake at the  $n$ th turn is  $V_n$ . Then the player's total gain through the  $n$ th turn is

$$X_n = \sum_{i=1}^n V_i \eta_i = X_{n-1} + V_n \eta_n, \quad X_0 = 0.$$

It is quite natural to suppose that the amount  $V_n$  at the  $n$ th turn may depend on the results of the preceding turns, i.e., on  $V_1, \dots, V_{n-1}$  and on  $\eta_1, \dots, \eta_{n-1}$ . In other words, if we put  $F_0 = \{\emptyset, \Omega\}$  and  $F_n = \sigma\{\omega: \eta_1, \dots, \eta_n\}$ , then  $V_n$  is an  $\mathcal{F}_{n-1}$ -measurable random variable, i.e., the sequence  $V = (V_n, \mathcal{F}_{n-1})$  that determines the player's "strategy" is predictable. Putting  $Y_n = \eta_1 + \dots + \eta_n$ , we find that

$$X_n = \sum_{i=1}^n V_i \Delta Y_i,$$

i.e., the sequence  $X = (X_n, \mathcal{F}_n)$  with  $X_0 = 0$  is the transform of  $Y$  by  $V$ .

From the player's point of view, the game in question is *fair* (or *favorable*, or *unfavorable*) if, at every stage, the conditional expectation

$$\mathbf{E}(X_{n+1} - X_n | \mathcal{F}_n) = 0 \text{ (or } \geq 0 \text{ or } \leq 0).$$

Moreover, it is clear that the game is

$$\begin{aligned} & \text{fair if } p = q = \frac{1}{2}, \\ & \text{favorable if } p > q, \\ & \text{unfavorable, if } p < q. \end{aligned}$$

Since  $X = (X_n, \mathcal{F}_n)$  is a

$$\begin{aligned} & \text{martingale if } p = q = \frac{1}{2}, \\ & \text{submartingale if } p > q, \\ & \text{supermartingale if } p < q, \end{aligned}$$

we can say that the assumption that the game is fair (or favorable, or unfavorable) corresponds to the assumption that the sequence  $X$  is a martingale (or submartingale, or supermartingale).

Let us now consider the special class of strategies  $V = (V_n, \mathcal{F}_{n-1})_{n \geq 1}$  with  $V_1 = 1$  and (for  $n > 1$ )

$$V_n = \begin{cases} 2^{n-1} & \text{if } \eta_1 = -1, \dots, \eta_{n-1} = -1, \\ 0 & \text{otherwise.} \end{cases} \quad (9)$$

In such a strategy, a player, having started with a stake  $V_1 = 1$ , doubles the stake after a loss and drops out of the game immediately after a win.

If  $\eta_1 = -1, \dots, \eta_n = -1$ , the total loss to the player after  $n$  turns will be

$$\sum_{i=1}^n 2^{i-1} = 2^n - 1.$$

Therefore if also  $\eta_{n+1} = 1$ , we have

$$X_{n+1} = X_n + V_{n+1} = -(2^n - 1) + 2^n = 1.$$

Let  $\tau = \inf\{n \geq 1: X_n = 1\}$ . If  $p = q = \frac{1}{2}$ , i.e., the game in question is fair, then  $P(\tau = n) = (\frac{1}{2})^n$ ,  $P(\tau < \infty) = 1$ ,  $P(X_\tau = 1) = 1$ , and  $EX_\tau = 1$ . Therefore even for a fair game, by applying the strategy (9), a player can in a finite time (with probability unity) complete the game "successfully," increasing his capital by one unit ( $EX_\tau = 1 > X_0 = 0$ ).

In gambling practice, this system (doubling the stakes after a loss and dropping out of the game after a win) is called a martingale. This is the origin of the mathematical term "martingale."

**Remark.** When  $p = q = \frac{1}{2}$ , the sequence  $X = (X_n, \mathcal{F}_n)$  with  $X_0 = 0$  is a martingale and therefore

$$EX_n = EX_0 = 0 \quad \text{for every } n \geq 1.$$

We may therefore expect that this equation is preserved if the instant  $n$  is replaced by a random instant  $\tau$ . It will appear later (Theorem 1, §2) that  $EX_\tau = EX_0$  in "typical" situations. Violations of this equation (as in the game discussed above) arise in what we may describe as physically unrealizable situations, when either  $\tau$  or  $|X_n|$  takes values that are much too large. (Note that the game discussed above would be physically unrealizable, since it supposes an unbounded time for playing and an unbounded initial capital for the player.)

**6. Definition 6.** A stochastic sequence  $\xi = (\xi_n, \mathcal{F}_n)$  is a *martingale-difference* if  $E|\xi| < \infty$  for all  $n \geq 0$  and

$$E(\xi_{n+1} | \mathcal{F}_n) = 0 \quad (\text{P-a.s.}). \quad (10)$$

The connection between martingales and martingale-differences is clear from Definitions 1 and 6. Thus if  $X = (X_n, \mathcal{F}_n)$  is a martingale, then  $\xi = (\xi_n, \mathcal{F}_n)$  with  $\xi_0 = X_0$  and  $\xi_n = \Delta X_n$ ,  $n \geq 1$ , is a martingale-difference. In turn, if  $\xi = (\xi_n, \mathcal{F}_n)$  is a martingale-difference, then  $X = (X_n, \mathcal{F}_n)$  with  $X_n = \xi_0 + \dots + \xi_n$  is a martingale.

In agreement with this terminology, every sequence  $\xi = (\xi_n)_{n \geq 0}$  of independent integrable random variables with  $E\xi_n = 0$  is a martingale-difference (with  $\mathcal{F}_n = \sigma\{\omega: \xi_0, \xi_1, \dots, \xi_n\}$ ).

7. The following theorem elucidates the structure of submartingales (or supermartingales).

**Theorem 2 (Doob).** *Let  $X = (X_n, \mathcal{F}_n)$  be a submartingale. Then there is a martingale  $m = (m_n, \mathcal{F}_n)$  and a predictable increasing sequence  $A = (A_n, \mathcal{F}_{n-1})$  such that, for every  $n \geq 0$ , Doob's decomposition*

$$X_n = m_n + A_n \quad (\text{P-a.s.}) \quad (11)$$

*holds. A decomposition of this kind is unique.*

PROOF. Let us put  $m_0 = X_0$ ,  $A_0 = 0$  and

$$m_n = m_0 + \sum_{j=0}^{n-1} [X_{j+1} - E(X_{j+1} | \mathcal{F}_j)], \quad (12)$$

$$A_n = \sum_{j=0}^{n-1} [E(X_{j+1} | \mathcal{F}_j) - X_j]. \quad (13)$$

It is evident that  $m$  and  $A$ , defined in this way, have the required properties. In addition, let  $X_n = m'_n + A'_n$ , where  $m' = (m'_n, \mathcal{F}_n)$  is a martingale and  $A' = (A'_n, \mathcal{F}_n)$  is a predictable increasing sequence. Then

$$A'_{n+1} - A'_n = (A_{n+1} - A_n) + (m_{n+1} - m_n) - (m'_{n+1} - m'_n),$$

and if we take conditional expectations on both sides, we find that (P-a.s.)  $A'_{n+1} - A'_n = A_{n+1} - A_n$ . But  $A_0 = A'_0 = 0$ , and therefore  $A_n = A'_n$  and  $m_n = m'_n$  (P-a.s.) for all  $n \geq 0$ .

This completes the proof of the theorem.

It follows from (11) that the sequence  $A = (A_n, \mathcal{F}_{n-1})$  compensates  $X = (X_n, \mathcal{F}_n)$  so that it becomes a martingale. This observation is justified by the following definition.

**Definition 7.** A predictable increasing sequence  $A = (A_n, \mathcal{F}_{n-1})$  appearing in the Doob decomposition (11) is called a *compensator* (of the submartingale  $X$ ).

The Doob decomposition plays a key role in the study of square integrable martingales  $M = (M_n, \mathcal{F}_n)$  i.e., martingales for which  $EM_n^2 < \infty$ ,  $n \geq 0$ ; this depends on the observation that the stochastic sequence  $M^2 = (M_n^2, \mathcal{F}_n)$  is a submartingale. According to Theorem 2 there is a martingale  $m = (m_n, \mathcal{F}_n)$  and a predictable increasing sequence  $\langle M \rangle = (\langle M \rangle_n, \mathcal{F}_{n-1})$  such that

$$M_n^2 = m_n + \langle M \rangle_n. \quad (14)$$

The sequence  $\langle M \rangle$  is called the *predictable quadratic variation* or the *quadratic characteristic* of  $M$  and, in many respects, determines its structure and properties.

It follows from (12) that

$$\langle M \rangle_n = \sum_{j=1}^n E[(\Delta M_j)^2 | \mathcal{F}_{j-1}] \quad (15)$$

and, for all  $l \leq k$ ,

$$E[(M_k - M_l)^2 | \mathcal{F}_l] = E[M_k^2 - M_l^2 | \mathcal{F}_l] = E[\langle M \rangle_k - \langle M \rangle_l | \mathcal{F}_l]. \quad (16)$$

In particular, if  $M_0 = 0$  (P-a.s.) then

$$EM_k^2 = E\langle M \rangle_k. \quad (17)$$

It is useful to observe that if  $M_0 = 0$  and  $M_n = \xi_1 + \cdots + \xi_n$ , where  $(\xi_n)$  is a sequence of independent random variables with  $E\xi_i = 0$  and  $E\xi_i^2 < \infty$ , the quadratic variation

$$\langle M \rangle_n = EM_n^2 = V\xi_1 + \cdots + V\xi_n \quad (18)$$

is not random, and indeed coincides with the variance.

If  $X = (X_n, \mathcal{F}_n)$  and  $Y = (Y_n, \mathcal{F}_n)$  are square integrable martingales, we put

$$\langle X, Y \rangle_n = \frac{1}{4}[\langle X + Y \rangle_n - \langle X - Y \rangle_n]. \quad (19)$$

It is easily verified that  $(X_n Y_n - \langle X, Y \rangle_n, \mathcal{F}_n)$  is a martingale and therefore, for  $l \leq k$ ,

$$E[(X_k - X_l)(Y_k - Y_l) | \mathcal{F}_l] = E[\langle X, Y \rangle_k - \langle X, Y \rangle_l | \mathcal{F}_l]. \quad (20)$$

In the case when  $X_n = \xi_1 + \cdots + \xi_n$ ,  $Y_n = \eta_1 + \cdots + \eta_n$ , where  $(\xi_n)$  and  $(\eta_n)$  are sequences of independent random variables with  $E\xi_i = E\eta_i = 0$  and  $E\xi_i^2 < \infty$ ,  $E\eta_i^2 < \infty$ , the variable  $\langle X, Y \rangle_n$  is given by

$$\langle X, Y \rangle_n = \sum_{i=1}^n \text{cov}(\xi_i, \eta_i).$$

The sequence  $\langle X, Y \rangle = (\langle X, Y \rangle_n, \mathcal{F}_{n-1})$  defined in (19) is often called the *mutual characteristic* of the (square integrable) martingales  $X$  and  $Y$ .

It is easy to show (compare with (15)) that

$$\langle X, Y \rangle_N = \sum_{i=1}^N E[\Delta X_i \Delta Y_i | \mathcal{F}_{i-1}].$$

In the theory of martingales, an important role is also played by the *quadratic covariation*,

$$[X, Y]_n = \sum_{i=1}^n \Delta X_i \Delta Y_i,$$

and the *quadratic variation*,



$$[X]_n = \sum_{i=1}^n (\Delta X_i)^2,$$

which can be defined for all random sequences  $X = (X_n)_{n \geq 1}$  and  $Y = (Y_n)_{n \geq 1}$ .

## 8. PROBLEMS

1. Show that (2) and (3) are equivalent.
2. Let  $\sigma$  and  $\tau$  be Markov times. Show that  $\tau + \sigma$ ,  $\tau \wedge \sigma$ , and  $\tau \vee \sigma$  are also Markov times; and if  $P(\sigma \leq \tau) = 1$ , then  $\mathcal{F}_\sigma \subseteq \mathcal{F}_\tau$ .
3. Show that  $\tau$  and  $X_\tau$  are  $\mathcal{F}_\tau$ -measurable.
4. Let  $Y = (Y_n, \mathcal{F}_n)$  be a martingale (or submartingale), let  $V = (V_n, \mathcal{F}_{n-1})$  be a predictable sequence, and let  $(V \cdot Y)_n$  be integrable random variables,  $n \geq 0$ . Show that  $V \cdot Y$  is a martingale (or submartingale).
5. Let  $\mathcal{F}_1 \subseteq \mathcal{F}_2 \subseteq \dots$  be a nondecreasing family of  $\sigma$ -algebras and  $\xi$  an integrable random variable. Show that  $(X_n)_{n \geq 1}$  with  $X_n = E(\xi | \mathcal{F}_n)$  is a martingale.
6. Let  $\mathcal{G}_1 \supseteq \mathcal{G}_2 \supseteq \dots$  be a nonincreasing family of  $\sigma$ -algebras and let  $\xi$  be an integrable random variable. Show that  $(X_n)_{n \geq 1}$  with  $X_n = E(\xi | \mathcal{G}_n)$  is a *reversed* martingale, i.e.,

$$E(X_n | X_{n+1}, X_{n+2}, \dots) = X_{n+1} \quad (\text{P-a.s.})$$

for every  $n \geq 1$ .

7. Let  $\xi_1, \xi_2, \xi_3, \dots$  be independent random variables,  $P(\xi_i = 0) = P(\xi_i = 2) = \frac{1}{2}$  and  $X_n = \prod_{i=1}^n \xi_i$ . Show that there does not exist an integrable random variable  $\xi$  and a nondecreasing family  $(\mathcal{F}_n)$  of  $\sigma$ -algebras such that  $X_n = E(\xi | \mathcal{F}_n)$ . This example shows that not every martingale  $(X_n)_{n \geq 1}$  can be represented in the form  $(E(\xi | \mathcal{F}_n))_{n \geq 1}$  (compare Example 3, §11, Chapter I).

## §2. Preservation of the Martingale Property Under Time Change at a Random Time

1. If  $X = (X_n, \mathcal{F}_n)_{n \geq 0}$  is a martingale, we have

$$EX_n = EX_0 \tag{1}$$

for every  $n \geq 1$ . Is this property preserved if the time  $n$  is replaced by a Markov time  $\tau$ ? Example 8 of the preceding section shows that, in general, the answer is “no”: there exist a martingale  $X$  and a Markov time  $\tau$  (finite with probability 1) such that

$$EX_\tau \neq EX_0. \tag{2}$$

The following basic theorem describes the “typical” situation, in which, in particular,  $EX_\tau = EX_0$ .

**Theorem 1 (Doob).** Let  $X = (X_n, \mathcal{F}_n)$  be a martingale (or submartingale), and  $\tau_1$  and  $\tau_2$ , stopping times for which

$$E|X_{\tau_i}| < \infty, \quad i = 1, 2, \quad (3)$$

$$\lim_{n \rightarrow \infty} \int_{\{\tau_i > n\}} |X_n| dP = 0, \quad i = 1, 2. \quad (4)$$

Then

$$E(X_{\tau_2} | \mathcal{F}_{\tau_1}) \underset{(\geq)}{=} X_{\tau_1} \quad (\{\tau_2 \geq \tau_1\}; P\text{-a.s.}) \quad (5)$$

If also  $P(\tau_1 \leq \tau_2) = 1$ , then

$$EX_{\tau_2} \underset{(\geq)}{=} EX_{\tau_1}. \quad (6)$$

(Here and in the formulas below, read the upper symbol for martingales and the lower symbol for submartingales.)

PROOF. It is sufficient to show that, for every  $A \in \mathcal{F}_{\tau_1}$ ,

$$\int_{A \cap \{\tau_2 \geq \tau_1\}} X_{\tau_2} dP \underset{(\geq)}{=} \int_{A \cap \{\tau_2 \geq \tau_1\}} X_{\tau_1} dP. \quad (7)$$

For this, in turn, it is sufficient to show that, for every  $n \geq 0$ ,

$$\int_{A \cap \{\tau_2 \geq \tau_1\} \cap \{\tau_1 = n\}} X_{\tau_2} dP \underset{(\geq)}{=} \int_{A \cap \{\tau_2 \geq \tau_1\} \cap \{\tau_1 = n\}} X_{\tau_1} dP,$$

or, what amounts to the same thing,

$$\int_{B \cap \{\tau_2 \geq n\}} X_{\tau_2} dP \underset{(\geq)}{=} \int_{B \cap \{\tau_2 \geq n\}} X_n dP, \quad (8)$$

where  $B = A \cap \{\tau_1 = n\} \in \mathcal{F}_n$ .

We have

$$\begin{aligned} \int_{B \cap \{\tau_2 \geq n\}} X_n dP &= \int_{B \cap \{\tau_2 = n\}} X_n dP + \int_{B \cap \{\tau_2 > n\}} X_n dP \underset{(\leq)}{=} \int_{B \cap \{\tau_2 = n\}} X_n dP \\ &+ \int_{B \cap \{\tau_2 > n\}} E(X_{n+1} | \mathcal{F}_n) dP = \int_{B \cap \{\tau_2 = n\}} X_{\tau_2} dP + \int_{B \cap \{\tau_2 \geq n+1\}} X_{n+1} dP \\ &\underset{(\leq)}{=} \int_{B \cap \{n \leq \tau_2 \leq n+1\}} X_{\tau_2} dP + \int_{B \cap \{\tau_2 \geq n+2\}} X_{n+2} dP \underset{(\leq)}{=} \cdots \\ &\underset{(\leq)}{=} \int_{B \cap \{n \leq \tau_2 \leq m\}} X_{\tau_2} dP + \int_{B \cap \{\tau_2 > m\}} X_m dP, \end{aligned}$$

whence

$$\int_{B \cap \{n \leq \tau_2 \leq m\}} X_{\tau_2} dP \underset{(\geq)}{=} \int_{B \cap \{n \leq \tau_2\}} X_n dP - \int_{B \cap \{m < \tau_2\}} X_m dP$$

and since  $X_m = 2X_m^+ - |X_m|$ , we have, by (4),

$$\begin{aligned} \int_{B \cap \{\tau_2 \geq n\}} X_{\tau_2} dP &\stackrel{(\geq)}{=} \lim_{m \rightarrow \infty} \left[ \int_{B \cap \{n \leq \tau_2\}} X_n dP - \int_{B \cap \{m < \tau_2\}} X_m dP \right] \\ &= \int_{B \cap \{n \leq \tau_2\}} X_n dP - \lim_{m \rightarrow \infty} \int_{B \cap \{m < \tau_2\}} X_m dP = \int_{B \cap \{\tau_2 \geq n\}} X_n dP, \end{aligned}$$

which establishes (8), and hence (5). Finally, (6) follows from (5).

This completes the proof of the theorem.

**Corollary 1.** *If there is a constant  $N$  such that  $P(\tau_1 \leq N) = 1$  and  $P(\tau_2 \leq N) = 1$ , then (3) and (4) are satisfied. Hence if, in addition,  $P(\tau_1 \leq \tau_2) = 1$  and  $X$  is a martingale, then*

$$EX_0 = EX_{\tau_1} = EX_{\tau_2} = EX_N. \quad (9)$$

**Corollary 2.** *If the random variables  $\{X_n\}$  are uniformly integrable (in particular, if  $|X_n| \leq C < \infty$ ,  $n \geq 0$ ), then (3) and (4) are satisfied.*

In fact,  $P(\tau_i > n) \rightarrow 0$ ,  $n \rightarrow \infty$ , and hence (4) follows from Lemma 2, §6, Chapter II. In addition, since the family  $\{X_n\}$  is uniformly integrable, we have (see II.6.(16))

$$\sup E|X_N| < \infty. \quad (10)$$

If  $\tau$  is a stopping time and  $X$  is a submartingale, then by Corollary 1, applied to the bounded time  $\tau_N = \tau \wedge N$ ,

$$EX_0 \leq EX_{\tau_N}.$$

Therefore

$$E|X_{\tau_N}| = 2EX_{\tau_N}^+ - EX_{\tau_N} \leq 2EX_{\tau_N}^+ - EX_0. \quad (11)$$

The sequence  $X^+ = (X_n^+, \mathcal{F}_n)$  is a submartingale (Example 5, §1) and therefore

$$\begin{aligned} EX_{\tau_N}^+ &= \sum_{j=0}^N \int_{\{\tau_N=j\}} X_j^+ dP + \int_{\{\tau > N\}} X_N^+ dP \leq \sum_{j=0}^N \int_{\{\tau_N=j\}} X_N^+ dP \\ &\quad + \int_{\{\tau > N\}} X_N^+ dP = EX_N^+ \leq E|X_N| \leq \sup_N E|X_N|. \end{aligned}$$

From this and (11) we have

$$E|X_{\tau_N}| \leq 3 \sup_N E|X_N|,$$

and hence by Fatou's lemma

$$E|X_\tau| \leq 3 \sup_N E|X_N|.$$

Therefore if we take  $\tau = \tau_i$ ,  $i = 1, 2$ , and use (10), we obtain  $E|X_{\tau_i}| < \infty$ ,  $i = 1, 2$ .

**Remark.** In Example 8 of the preceding section,

$$\int_{\{\tau > n\}} |X_n| dP = (2^n - 1)P\{\tau > n\} = (2^n - 1) \cdot 2^{-n} \rightarrow 1, \quad n \rightarrow \infty,$$

and consequently (4) is violated (for  $\tau_2 = \tau$ ).

2. The following proposition, which we shall deduce from Theorem 1, is often useful in applications.

**Theorem 2.** Let  $X = (X_n)$  be a martingale (or submartingale) and  $\tau$  a stopping time (with respect to  $(\mathcal{F}_n^X)$ , where  $\mathcal{F}_n^X = \sigma\{\omega: X_0, \dots, X_n\}$ ). Suppose that

$$E\tau < \infty,$$

and that for some  $n \geq 0$  and some constant  $C$

$$E\{|X_{n+1} - X_n| \mid \mathcal{F}_n^X\} \leq C \quad (\{\tau \geq n\}; \text{P-a.s.}).$$

Then

$$E|X_{\tau}| < \infty$$

and

$$EX_{\tau} \geq EX_0. \quad (12)$$

We first verify that hypotheses (3) and (4) of Theorem 1 are satisfied with  $\tau_2 = \tau$ .

Let

$$Y_0 = |X_0|, \quad Y_j = |X_j - X_{j-1}|, \quad j \geq 1.$$

Then  $|X_{\tau}| \leq \sum_{j=0}^{\tau} Y_j$  and

$$\begin{aligned} E|X_{\tau}| &\leq E\left(\sum_{j=0}^{\tau} Y_j\right) = \int_{\Omega} \left(\sum_{j=0}^{\tau} Y_j\right) dP = \sum_{n=0}^{\infty} \int_{\{\tau=n\}} \sum_{j=0}^n Y_j dP \\ &= \sum_{n=0}^{\infty} \sum_{j=0}^n \int_{\{\tau=n\}} Y_j dP = \sum_{j=0}^{\infty} \sum_{n=j}^{\infty} \int_{\{\tau=n\}} Y_j dP = \sum_{j=0}^{\infty} \int_{\{\tau \geq j\}} Y_j dP. \end{aligned}$$

The set  $\{\tau \geq j\} = \Omega \setminus \{\tau < j\} \in \mathcal{F}_{j-1}^X$ ,  $j \geq 1$ . Therefore

$$\int_{\{\tau \geq j\}} Y_j dP = \int_{\{\tau \geq j\}} E[Y_j | X_0, \dots, X_{j-1}] dP \leq CP\{\tau \geq j\}$$

for  $j \geq 1$ ; and

$$E|X_\tau| \leq E\left(\sum_{j=0}^{\tau} Y_j\right) \leq E|X_0| + C \sum_{j=1}^{\infty} P\{\tau \geq j\} = E|X_0| + CE\tau < \infty. \quad (13)$$

Moreover, if  $\tau > n$ , then

$$\sum_{j=0}^n Y_j \leq \sum_{j=0}^{\tau} Y_j,$$

and therefore

$$\int_{\{\tau > n\}} |X_n| dP \leq \int_{\{\tau > n\}} \sum_{j=0}^{\tau} Y_j dP.$$

Hence since (by (13))  $E \sum_{j=0}^{\tau} Y_j < \infty$  and  $\{\tau > n\} \downarrow \emptyset, n \rightarrow \infty$ , the dominated convergence theorem yields

$$\lim_{n \rightarrow \infty} \int_{\{\tau > n\}} |X_n| dP \leq \lim_{n \rightarrow \infty} \int_{\{\tau > n\}} \left( \sum_{j=0}^{\tau} Y_j \right) dP = 0.$$

Hence the hypotheses of Theorem 1 are satisfied, and (12) follows as required.

This completes the proof of the theorem.

3. Here we present some applications of the preceding theorems.

**Theorem 3 (Wald's Identities).** Let  $\xi_1, \xi_2, \dots$  be independent identically distributed random variables with  $E|\xi_i| < \infty$  and  $\tau$  a stopping time (with respect to  $\mathcal{F}_n^\xi$ ), where  $\mathcal{F}_n^\xi = \sigma\{\omega: \xi_1, \dots, \xi_n\}$ ,  $\tau \geq 1$ , and  $E\tau < \infty$ . Then

$$E(\xi_1 + \dots + \xi_\tau) = E\xi_1 \cdot E\tau. \quad (14)$$

If also  $E\xi_i^2 < \infty$  then

$$E\{(\xi_1 + \dots + \xi_\tau) - \tau E\xi_1\}^2 = V\xi_1 \cdot E\tau. \quad (15)$$

**PROOF.** It is clear that  $X = (X_n, \mathcal{F}_n^\xi)_{n \geq 1}$  with  $X_n = (\xi_1 + \dots + \xi_n) - nE\xi_1$  is a martingale with

$$\begin{aligned} E[X_{n+1} - X_n | X_1, \dots, X_n] &= E[\xi_{n+1} - E\xi_1 | \xi_1, \dots, \xi_n] \\ &= E[\xi_{n+1} - E\xi_1] = 0. \end{aligned}$$

Therefore  $EX_\tau = EX_0 = 0$ , by Theorem 2, and (14) is established.

Similar considerations applied to the martingale  $Y = (Y_n, \mathcal{F}_n^\xi)$  with  $Y_n = X_n^2 - nV\xi_1$  lead to a proof of (15).

**Corollary.** Let  $\xi_1, \xi_2, \dots$  be independent identically distributed random variables with

$$P(\xi_i = 1) = P(\xi_i = -1) = \frac{1}{2}, S_n = \xi_1 + \dots + \xi_n$$

and  $\tau = \inf\{n \geq 1: S_n = 1\}$ . Then  $P\{\tau < \infty\} = 1$  (see, for example, (I.9.20)) and therefore  $P(S_\tau = 1) = 1$ ,  $ES_\tau = 1$ . Hence it follows from (14) that  $E\tau = \infty$ .

**Theorem 4** (Wald's Fundamental Identity). Let  $\xi_1, \xi_2, \dots$ , be a sequence of independent identically distributed random variables,  $S_n = \xi_1 + \dots + \xi_n$ , and  $n \geq 1$ . Let  $\varphi(t) = Ee^{t\xi_1}$ ,  $t \in R$ , and for some  $t_0 \neq 0$  let  $\varphi(t_0)$  exist and  $\varphi(t_0) \geq 1$ .

If  $\tau$  is a stopping time (with respect to  $(\mathcal{F}_n^\xi)$ ,  $\mathcal{F}_n^\xi = \sigma\{\omega: \xi_1, \dots, \xi_n\}$ ,  $\tau \geq 1$ ), such that  $|S_n| \leq C$  ( $\{\tau \geq n\}$ ; P-a.s.) and  $E\tau < \infty$ , then

$$E\left[\frac{e^{t_0 S_\tau}}{(\varphi(t_0))^\tau}\right] = 1. \quad (16)$$

PROOF. Take

$$Y_n = e^{t_0 S_n} (\varphi(t_0))^{-n}.$$

Then  $Y = (Y_n, \mathcal{F}_n^\xi)_{n \geq 1}$  is a martingale with  $EY_n = 1$  and, on the set  $\{\tau \geq n\}$ ,

$$\begin{aligned} E\{|Y_{n+1} - Y_n| | Y_1, \dots, Y_n\} &= Y_n E\left\{\left|\frac{e^{t_0 \xi_{n+1}}}{\varphi(t_0)} - 1\right| \middle| \xi_1, \dots, \xi_n\right\} \\ &= Y_n \cdot E\{|e^{t_0 \xi_1} \varphi^{-1}(t_0) - 1|\} \leq B < \infty, \end{aligned}$$

where  $B$  is a constant. Therefore Theorem 2 is applicable, and (16) follows since  $EY_1 = 1$ .

This completes the proof.

**EXAMPLE 1.** This example will let us illustrate the use of the preceding examples to find the probabilities of ruin and of mean duration in games (see §9, Chapter I).

Let  $\xi_1, \xi_2, \dots$  be a sequence of independent Bernoulli random variables with  $P(\xi_i = 1) = p$ ,  $P(\xi_i = -1) = q$ ,  $p + q = 1$ ,  $S = \xi_1 + \dots + \xi_n$ , and

$$\tau = \inf\{n \geq 1: S_n = B \text{ or } A\}, \quad (17)$$

where  $(-A)$  and  $B$  are positive integers.

It follows from (I.9.20) that  $P(\tau < \infty) = 1$  and  $E\tau < \infty$ . Then if  $\alpha = P(S_\tau = A)$ ,  $\beta = P(S_\tau = B)$ , we have  $\alpha + \beta = 1$ . If  $p = q = \frac{1}{2}$ , we obtain

$$0 = ES_\tau = \alpha A + \beta B, \quad \text{from (14),}$$

whence

$$\alpha = \frac{B}{B + |A|}, \quad \beta = \frac{|A|}{B + |A|}.$$

Applying (15), we obtain

$$E_\tau = ES_\tau^2 = \alpha A^2 + \beta B^2 = |AB|.$$

However, if  $p \neq q$  we find, by considering the martingale  $((q/p)^{S_n})_{n \geq 1}$ , that

$$\mathbb{E} \left( \frac{q}{p} \right)^{S_\tau} = \mathbb{E} \left( \frac{q}{p} \right)^{S_1} = 1,$$

and therefore

$$\alpha \left( \frac{q}{p} \right)^A + \beta \left( \frac{q}{p} \right)^B = 1.$$

Together with the equation  $\alpha + \beta = 1$  this yields

$$\alpha = \frac{\left( \frac{q}{p} \right)^B - 1}{\left( \frac{q}{p} \right)^B - \left( \frac{q}{p} \right)^{|A|}}, \quad \beta = \frac{1 - \left( \frac{q}{p} \right)^{|A|}}{\left( \frac{q}{p} \right)^B - \left( \frac{q}{p} \right)^{|A|}}. \quad (18)$$

Finally, since  $\mathbb{E} S_\tau = (p - q)\mathbb{E} \tau$ , we find

$$\mathbb{E} \tau = \frac{\mathbb{E} S_\tau}{p - q} = \frac{\alpha A + \beta B}{p - q},$$

where  $\alpha$  and  $\beta$  are defined by (18).

**EXAMPLE 2.** In the example considered above, let  $p = q = \frac{1}{2}$ . Let us show that for every  $\lambda$  in  $0 < \lambda < \pi/(B + |A|)$  and every time  $\tau$  defined in (17),

$$\mathbb{E}(\cos \lambda)^{-\tau} = \frac{\cos \lambda \cdot \frac{B + A}{2}}{\cos \lambda \cdot \frac{B + |A|}{2}}. \quad (19)$$

For this purpose we consider the martingale  $X = (X_n, \mathcal{F}_n^S)_{n \geq 0}$  with

$$X_n = (\cos \lambda)^{-n} \cos \lambda \left( S_n - \frac{B + A}{2} \right) \quad (20)$$

and  $S_0 = 0$ . It is clear that

$$\mathbb{E} X_n = \mathbb{E} X_0 = \cos \lambda \frac{B + A}{2}. \quad (21)$$

Let us show that the family  $\{X_{n \wedge \tau}\}$  is uniformly integrable. For this purpose we observe that, by Corollary 1 to Theorem 1 for  $0 < \lambda < \pi/(B + |A|)$ ,

$$\begin{aligned} \mathbb{E} X_0 &= \mathbb{E} X_{n \wedge \tau} = \mathbb{E} (\cos \lambda)^{-(n \wedge \tau)} \cos \lambda \left( S_{n \wedge \tau} - \frac{B + A}{2} \right) \\ &\geq \mathbb{E} (\cos \lambda)^{-(n \wedge \tau)} \cos \lambda \frac{B - A}{2}. \end{aligned}$$

Therefore, by (21),

$$E(\cos \lambda)^{-(n \wedge \tau)} \leq \frac{\cos \lambda \frac{B+A}{2}}{\cos \lambda \frac{B+|A|}{2}},$$

and consequently by Fatou's lemma,

$$E(\cos \lambda)^{-\tau} \leq \frac{\cos \lambda \frac{B+A}{2}}{\cos \lambda \frac{B+|A|}{2}}. \quad (22)$$

Consequently, by (20),

$$|X_{n \wedge \tau}| \leq (\cos \lambda)^{-\tau}.$$

With (22), this establishes the uniform integrability of the family  $\{X_{n \wedge \tau}\}$ . Then, by Corollary 2 to Theorem 1,

$$\cos \lambda \frac{B+A}{2} = EX_0 = EX_\tau = E(\cos \lambda)^{-\tau} \cos \lambda \frac{B-A}{2},$$

from which the required inequality (19) follows.

#### 4. PROBLEMS

1. Show that Theorem 1 remains valid for submartingales if (4) is replaced by

$$\lim_{n \rightarrow \infty} \int_{\{\tau > n\}} X_n^+ dP = 0, \quad i = 1, 2.$$

2. Let  $X = (X_n, \mathcal{F}_n)_{n \geq 0}$  be a square-integrable martingale,  $\tau$  a stopping time and

$$\lim_{n \rightarrow \infty} \int_{\{\tau > n\}} X_n^2 dP = 0,$$

$$\lim_{n \rightarrow \infty} \int_{\{\tau > n\}} |X_n| dP = 0.$$

Show that then

$$EX_\tau^2 = E\langle X \rangle_\tau \left( = E \sum_{j=0}^{\tau} (\Delta X_j)^2 \right),$$

where  $\Delta X_0 = X_0$ ,  $\Delta X_j = X_j - X_{j-1}$ ,  $j \geq 1$ .

3. Show that

$$E|X_\tau| \leq \lim_{n \rightarrow \infty} E|X_n|$$

for every martingale or nonnegative submartingale  $X = (X_n, \mathcal{F}_n)_{n \geq 0}$  and every stopping time  $\tau$ .



4. Let  $X = (X_n, \mathcal{F}_n)_{n \geq 0}$  be a submartingale such that  $X_n \geq E(\xi | \mathcal{F}_n)$  (P-a.s.),  $n \geq 0$ , where  $E|\xi| < \infty$ . Show that if  $\tau_1$  and  $\tau_2$  are stopping times with  $P(\tau_1 \leq \tau_2) = 1$ , then

$$X_{\tau_1} \geq E(X_{\tau_2} | \mathcal{F}_{\tau_1}) \quad (\text{P-a.s.}).$$

5. Let  $\xi_1, \xi_2, \dots$  be a sequence of independent random variables with  $P(\xi_i = 1) = P(\xi_i = -1) = \frac{1}{2}$ ,  $a$  and  $b$  positive numbers,  $b > a$ ,

$$X_n = a \sum_{k=1}^n I(\xi_k = +1) - b \sum_{k=1}^n I(\xi_k = -1)$$

and

$$\tau = \inf\{n \geq 1: X_n \leq -r\}, \quad r > 0.$$

Show that  $Ee^{\lambda \tau} < \infty$  for  $\lambda \leq \alpha_0$  and  $Ee^{\lambda \tau} = \infty$  for  $\lambda > \alpha_0$ , where

$$\alpha_0 = \frac{b}{a+b} \ln \frac{2b}{a+b} + \frac{a}{a+b} \ln \frac{2a}{a+b}.$$

6. Let  $\xi_1, \xi_2, \dots$  be a sequence of independent random variables with  $E\xi_i = 0, \forall \xi_i = \sigma_i^2$ ,  $S_n = \xi_1 + \dots + \xi_n$ ,  $\mathcal{F}_n^S = \sigma\{\omega: \xi_1, \dots, \xi_n\}$ . Prove the following generalizations of Wald's identities (14) and (15): If  $E \sum_{j=1}^t E|\xi_j| < \infty$  then  $ES_t = 0$ ; if  $E \sum_{j=1}^t E\xi_j^2 < \infty$ , then

$$ES_t^2 = E \sum_{j=1}^t \xi_j^2 = E \sum_{j=1}^t \sigma_j^2. \quad (23)$$

### §3. Fundamental Inequalities

1. Let  $X = (X_n, \mathcal{F}_n)_{n \geq 0}$  be a stochastic sequence,

$$X_n^* = \max_{0 \leq j \leq n} |X_j|, \quad \|X_n\|_p = (E|X_n|^p)^{1/p}, \quad p > 0.$$

In Theorems 1–3 below, we present Doob's fundamental "maximal inequalities for probabilities" and "maximal inequalities in  $L^p$ ," for submartingales, supermartingales and martingales.

**Theorem 1. I.** Let  $X = (X_n, \mathcal{F}_n)_{n \geq 0}$  be a submartingale. Then for all  $\lambda > 0$

$$\lambda P \left\{ \max_{k \leq n} X_k \geq \lambda \right\} \leq E \left[ X_n^+ I \left( \max_{k \leq n} X_k \geq \lambda \right) \right] \leq EX_n^+, \quad (1)$$

$$\lambda P \left\{ \min_{k \leq n} X_k \leq -\lambda \right\} \leq E \left[ X_n I \left( \min_{k \leq n} X_k > -\lambda \right) \right] - EX_0 \leq EX_n^+ - EX_0, \quad (2)$$

$$\lambda P \left\{ \max_{k \leq n} |X_k| \geq \lambda \right\} \leq 3 \max_{k \leq n} E|X_k|. \quad (3)$$

II. Let  $Y = (Y_n, \mathcal{F}_n)_{n \geq 0}$  be a supermartingale. Then for all  $\lambda > 0$

$$\lambda P \left\{ \max_{k \leq n} Y_k \geq \lambda \right\} \leq EY_0 - E \left[ Y_n I \left( \max_{k \leq n} Y_k < \lambda \right) \right] \leq EY_0 + EY_n^-, \quad (4)$$

$$\lambda P \left\{ \min_{k \leq n} Y_k \leq -\lambda \right\} \leq -E \left[ Y_n I \left( \min_{k \leq n} Y_k \leq -\lambda \right) \right] \leq EY_n^-, \quad (5)$$

$$\lambda P \left\{ \max_{k \leq n} |Y_k| \geq \lambda \right\} \leq 3 \max_{k \leq n} E|Y_k|. \quad (6)$$

III. Let  $Y = (Y_n, \mathcal{F}_n)_{n \geq 0}$  be a nonnegative supermartingale. Then for all  $\lambda > 0$

$$\lambda P \left\{ \max_{k \leq n} Y_k \geq \lambda \right\} \leq EY_0, \quad (7)$$

$$\lambda P \left\{ \sup_{k \geq n} Y_k \geq \lambda \right\} \leq EY_n. \quad (8)$$

**Theorem 2.** Let  $X = (X_n, \mathcal{F}_n)_{n \geq 0}$  be a nonnegative submartingale. Then for  $p \geq 1$  we have the following inequalities:

if  $p > 1$ ,

$$\|X_n\|_p \leq \|X_n^*\|_p \leq \frac{p}{p-1} \|X_n\|_p; \quad (9)$$

if  $p = 1$ ,

$$\|X_n\| \leq \|X_n^*\|_1 \leq \frac{e}{e-1} \{1 + \|X_n \ln^+ X_n\|_1\}. \quad (10)$$

**Theorem 3.** Let  $X = (X_n, \mathcal{F}_n)_{n \geq 0}$  be a martingale,  $\lambda > 0$  and  $p \geq 1$ . Then

$$P \left\{ \max_{k \leq n} |X_k| \geq \lambda \right\} \leq \frac{E|X_n|^p}{\lambda^p} \quad (11)$$

and if  $p > 1$

$$\|X_n\|_p \leq \|X_n^*\|_p \leq \frac{p}{p-1} \|X_n\|_p. \quad (12)$$

In particular, if  $p = 2$

$$P \left\{ \max_{k \leq n} |X_k| \geq \lambda \right\} \leq \frac{E|X_n|^2}{\lambda^2}, \quad (13)$$

$$E \left[ \max_{k \leq n} X_k^2 \right] \leq 4EX_n^2. \quad (14)$$

PROOF OF THEOREM 1. Since a submartingale with the opposite sign is a supermartingale, (1)–(3) follow from (4)–(6). Therefore, we consider the case of a supermartingale  $Y = (Y_n, \mathcal{F}_n)_{n \geq 0}$ .

Let us set  $\tau = \inf\{k \leq n: Y_k \geq \lambda\}$  with  $\tau = n$  if  $\max_{k \leq n} Y_k < \lambda$ . Then, by (2.6),

$$\begin{aligned} EY_0 &\geq EY_\tau = E\left[Y_\tau; \max_{k \leq n} Y_k \geq \lambda\right] + E\left[Y_\tau; \max_{k \leq n} Y_k < \lambda\right] \\ &\geq \lambda P\left\{\max_{k \leq n} Y_k \geq \lambda\right\} + E\left[Y_n; \max_{k \leq n} Y_k < \lambda\right], \end{aligned}$$

which proves (4).

Now let us set  $\sigma = \inf\{k \leq n: Y_k \leq -\lambda\}$ , and take  $\sigma = n$  if  $\min_{k \leq n} Y_k > -\lambda$ . Again, by (2.6),

$$\begin{aligned} EY_n &\leq EY_\tau = E\left[Y_\tau; \min_{k \leq n} Y_k \leq -\lambda\right] + E\left[Y_\tau; \min_{k \leq n} Y_k > -\lambda\right] \\ &\leq \lambda P\left\{\min_{k \leq n} Y_k \leq -\lambda\right\} + E\left[Y_n; \min_{k \leq n} Y_k > -\lambda\right]. \end{aligned}$$

Hence,

$$\lambda P\left\{\min_{k \leq n} Y_k \leq -\lambda\right\} \leq -E\left[Y_n; \min_{k \leq n} Y_k \leq -\lambda\right] \leq EY_n^-$$

which proves (5).

To prove (6), we notice that  $Y^- = (-Y)^+$  is a submartingale. Then, by (4) and (1),

$$\begin{aligned} P\left\{\max_{k \leq n} |Y_k| \geq \lambda\right\} &\leq P\left\{\max_{k \leq n} Y_k^+ \geq \lambda\right\} + P\left\{\max_{k \leq n} Y_k^- \geq \lambda\right\} \\ &= P\left\{\max_{k \leq n} Y_k \geq \lambda\right\} + P\left\{\max_{k \leq n} Y_k^- \geq \lambda\right\} \\ &\leq EY_0 + 2EY_n^- \leq 3 \max_{k \leq n} E|Y_k|. \end{aligned}$$

Inequality (7) follows from (4).

To prove (8), we set  $\gamma = \inf\{k \geq n: Y_k \geq \lambda\}$ , taking  $\gamma = \infty$  if  $Y_k < \lambda$  for all  $k \geq n$ . Now let  $n < N < \infty$ . Then, by (2.6),

$$EY_n \geq EY_{\gamma \wedge N} \geq E[Y_{\gamma \wedge N} I(\gamma \leq N)] \geq \lambda P\{\gamma \leq N\},$$

from which, as  $N \rightarrow \infty$ ,

$$EY_n \geq \lambda P\{\gamma < \infty\} = \lambda P\left\{\sup_{k \geq n} Y_k \geq \lambda\right\}.$$

## PROOF OF THEOREM 2.

The first inequalities in (9) and (10) are evident.

To prove the second inequality in (9), we first suppose that

$$\|X_n^*\|_p < \infty, \quad (15)$$

and use the fact that, for every nonnegative random variable  $\xi$  and for  $r > 0$ ,

$$E\xi^r = r \int_0^\infty t^{r-1} P(\xi \geq t) dt. \quad (16)$$

Then we obtain, by (1) and Fubini's theorem, that for  $p > 1$

$$\begin{aligned} E(X_n^*)^p &= p \int_0^\infty t^{p-1} P\{X_n^* \geq t\} dt \leq p \int_0^\infty t^{p-2} \left( \int_{\{X_n^* \geq t\}} X_n dP \right) dt \\ &= p \int_0^\infty t^{p-2} \left[ \int_\Omega X_n I\{X_n^* \geq t\} dP \right] dt \\ &= p \int_\Omega X_n \left[ \int_0^{X_n^*} t^{p-2} dt \right] dP = \frac{p}{p-1} E[X_n (X_n^*)^{p-1}]. \end{aligned} \quad (17)$$

Hence, by Hölder's inequality,

$$E(X_n^*)^p \leq q \|X_n\|_p \cdot \|(X_n^*)^{p-1}\|_q = q \|X_n\|_p [E(X_n^*)^p]^{1/q}, \quad (18)$$

where  $q = p/(p-1)$ .

If (15) is satisfied, we immediately obtain the second inequality in (9) from (18).

However, if (15) is not satisfied, we proceed as follows. In (17), instead of  $X_n^*$  we consider  $(X_n^* \wedge L)$ , where  $L$  is a constant. Then we obtain

$$E(X_n^* \wedge L)^p \leq q E[X_n (X_n^* \wedge L)^{p-1}] \leq q \|X_n\|_p [E(X_n^* \wedge L)^p]^{1/q},$$

from which it follows, by the inequality  $E(X_n^* \wedge L)^p \leq L^p < \infty$ , that

$$E(X_n^* \wedge L)^p \leq q^p E X_n^p = q^p \|X_n\|_p^p$$

and therefore,

$$E(X_n^*)^p = \lim_{L \rightarrow \infty} E(X_n^* \wedge L)^p \leq q^p \|X_n\|_p^p.$$

We now prove the second inequality in (10).

Again applying (1), we obtain

$$\begin{aligned} EX_n^* - 1 &\leq E(X_n^* - 1)^+ = \int_0^\infty P\{X_n^* - 1 \geq t\} dt \\ &\leq \int_0^\infty \frac{1}{1+t} \left[ \int_{\{X_n^* \geq 1+t\}} X_n dP \right] dt = EX_n \int_0^{X_n^*-1} \frac{dt}{1+t} = EX_n \ln X_n^*. \end{aligned}$$

Since, for arbitrary  $a \geq 0$  and  $b > 0$ ,

$$a \ln b \leq a \ln^+ a + b e^{-1}, \quad (19)$$

we have

$$EX_n^* - 1 \leq EX_n \ln X_n^* \leq EX_n \ln^+ X_n + e^{-1} EX_n^*.$$

If  $EX_n^* < \infty$ , we immediately obtain the second inequality (10).

However, if  $EX^* = \infty$ , we proceed, as above, by replacing  $X_n^*$  by  $X_n^* \wedge L$ . This proves the theorem.

The proof of Theorem 3 follows from the remark that  $|X|^p$ ,  $p \geq 1$ , is a nonnegative submartingale (if  $E|X_n|^p < \infty$ ,  $n \geq 0$ ), and from inequalities (1) and (9).

**Corollary of Theorem 3.** Let  $X_n = \xi_0 + \dots + \xi_n$ ,  $n \geq 0$ , where  $(\xi_k)_{k \geq 0}$  is a sequence of independent random variables with  $E\xi_k = 0$  and  $E\xi_k^2 < \infty$ . Then inequality (13) becomes Kolmogorov's inequality (§2, Chapter IV).

2. Let  $X = (X_n, \mathcal{F}_n)$  be a nonnegative submartingale and

$$X_n = M_n + A_n,$$

its Doob decomposition. Then, since  $EM_n = 0$ , it follows from (1) that

$$P\{X_n^* \geq \varepsilon\} \leq \frac{EA_n}{\varepsilon}.$$

Theorem 4, below, shows that this inequality is valid, not only for submartingales, but also for the wider class of sequences that have the property of domination in the following sense.

**Definition.** Let  $X = (X_n, \mathcal{F}_n)$  be a nonnegative stochastic sequence, and  $A = (A_n, \mathcal{F}_{n-1})$  an increasing predictable sequence. We shall say that  $X$  is dominated by the sequence  $A$  if

$$EX_\tau \leq EA_\tau \quad (20)$$

for every stopping time  $\tau$ .

**Theorem 4.** If  $X = (X_n, \mathcal{F}_n)$  is a nonnegative stochastic sequence dominated by an increasing predictable sequence  $A = (A_n, \mathcal{F}_{n-1})$ , then for  $\lambda > 0$ ,  $a > 0$ , and any stopping time  $\tau$ ,

$$P\{X_\tau^* \geq \lambda\} \leq \frac{EA_\tau}{\lambda}, \quad (21)$$

$$P\{X_\tau^* \geq \lambda\} \leq \frac{1}{\lambda} E(A_\tau \wedge a) + P(A_\tau \geq a), \quad (22)$$

$$\|X_\tau^*\|_p \leq \left(\frac{2-p}{1-p}\right)^{1/p} \|A_\tau\|_p, \quad 0 < p < 1. \quad (23)$$

PROOF. We set

$$\sigma_n = \min\{j \leq \tau \wedge n: X_j \geq \lambda\},$$

taking  $\sigma_n = \tau \wedge n$ , if  $\{\cdot\} = \emptyset$ . Then

$$EA_\tau \geq EA_{\sigma_n} \geq EX_{\sigma_n} \geq \int_{\{X_{\sigma_n}^* > \lambda\}} X_{\sigma_n} dP \geq \lambda P\{X_{\tau \wedge n}^* > \lambda\},$$

from which

$$P\{X_{\tau \wedge n}^* > \lambda\} \leq \frac{1}{\lambda} EA_\tau,$$

and we obtain (21) by Fatou's lemma.

For the proof of (22), we introduce the time

$$\gamma = \inf\{j: A_{j+1} \geq a\},$$

setting  $\gamma = \infty$  if  $\{\cdot\} = \emptyset$ . Then

$$\begin{aligned} P\{X_\tau^* \geq \lambda\} &= P\{X_\tau^* \geq \lambda, A_\tau < a\} + P\{X_\tau^* \geq \lambda, A_\tau \geq a\} \\ &\leq P\{I_{\{A_\tau < a\}} X_\tau^* \geq \lambda\} + P\{A_\tau \geq a\} \\ &\leq P\{X_{\tau \wedge \gamma}^* \geq \lambda\} + P\{A_\tau \geq a\} \leq \frac{1}{\lambda} EA_{\tau \wedge \gamma} + P\{A_\gamma \geq a\} \\ &\leq \frac{1}{\lambda} E(A_\tau \wedge a) + P(A_\tau \geq a), \end{aligned}$$

where we used (21) and the inequality  $I_{\{A_\tau < a\}} X_\tau^* \leq X_{\tau \wedge \gamma}^*$ . Finally, by (22),

$$\begin{aligned} \|X_\tau^*\|_p^p &= E(X_\tau^*)^p = \int_0^\infty P\{(X_\tau^*)^p \geq t\} dt = \int_0^\infty P\{X_\tau^* \geq t^{1/p}\} dt \\ &\leq \int_0^\infty t^{-1/p} E[A_\tau \wedge t^{1/p}] dt + \int_0^\infty P\{A_\tau^p \geq t\} dt \\ &= E \int_0^{A_\tau^p} dt + E \int_{A_\tau^p}^\infty (A_\tau t^{-1/p}) dt + EA_\tau^p = \frac{2-p}{1-p} EA_\tau^p. \end{aligned}$$

This completes the proof.

**Remark.** Let us suppose that the hypotheses of Theorem 4 are satisfied, except that the sequence  $A = (A_n, \mathcal{F}_n)_{n \geq 0}$  is not necessarily predictable, but has the property that for some positive constant  $c$

$$P\left\{\sup_{k \geq 1} |\Delta A_k| \leq c\right\} = 1,$$

where  $\Delta A_k = A_k - A_{k-1}$ . Then the following inequality is satisfied (compare (22)):

$$P\{X_\tau^* \geq \lambda\} \leq \frac{1}{\lambda} E[A_\tau \wedge (a + c)] + P\{A_\tau \geq a\}. \quad (24)$$

The proof is analogous to that of (22). We have only to replace the time  $\gamma = \inf\{j: A_{j+1} \geq a\}$  by  $\gamma = \inf\{j: A_j \geq a\}$  and notice that  $A_\gamma \leq a + c$ .

**Corollary.** Let the sequences  $X^k = (X_n^k, \mathcal{F}_n^k)$  and  $A^k = (A_n^k, \mathcal{F}_n^k)$ ,  $n \geq 0$ ,  $k \geq 1$  satisfy the hypotheses of Theorem 4 or the remark. Also, let  $(\tau^k)_{k \geq 1}$  be a sequence of stopping times (with respect to  $\mathcal{F}^k = (\mathcal{F}_n^k)$ ) and  $A_{\tau^k}^k \xrightarrow{P} 0$ . Then  $(X^k)_{\tau^k}^* \xrightarrow{P} 0$ .

3. In this subsection we present (without proofs, but with applications) a number of significant inequalities for martingales. These generalize the inequalities of Khinchin and of Marcinkiewicz and Zygmund for sums of independent random variables.

**Khinchin's Inequalities.** Let  $\xi_1, \xi_2, \dots$  be independent identically distributed Bernoulli random variables with  $P(\xi_i = 1) = P(\xi_i = -1) = \frac{1}{2}$  and let  $(c_n)_{n \geq 1}$  be a sequence of numbers.

Then for every  $p$ ,  $0 < p < \infty$ , there are universal constants  $A_p$  and  $B_p$  (independent of  $(c_n)$ ) such that

$$A_p \left( \sum_{j=1}^n c_j^2 \right)^{1/2} \leq \left\| \sum_{j=1}^n c_j \xi_j \right\|_p \leq B_p \left( \sum_{j=1}^n c_j^2 \right)^{1/2} \quad (25)$$

for every  $n \geq 1$ .

The following result generalizes these inequalities (for  $p \geq 1$ ).

**Marcinkiewicz and Zygmund's Inequalities.** If  $\xi_1, \xi_2, \dots$  is a sequence of independent integrable random variables with  $E\xi_i = 0$ , then for  $p \geq 1$  there are universal constants  $A_p$  and  $B_p$  (independent of  $(\xi_n)$ ) such that

$$A_p \left\| \left( \sum_{j=1}^n \xi_j^2 \right)^{1/2} \right\|_p \leq \left\| \sum_{j=1}^n \xi_j \right\|_p \leq B_p \left\| \left( \sum_{j=1}^n \xi_j^2 \right)^{1/2} \right\|_p \quad (26)$$

for every  $n \geq 1$ .

In (25) and (26) the sequences  $X = (X_n)$  with  $X_n = \sum_{j=1}^n c_j \xi_j$  and  $X_n = \sum_{j=1}^n \xi_j$  are martingales. It is natural to ask whether the inequalities can be extended to arbitrary martingales.

The first result in this direction was obtained by Burkholder.

**Burkholder's Inequalities.** If  $X = (X_n, \mathcal{F}_n)$  is a martingale, then for every  $p > 1$  there are universal constants  $A_p$  and  $B_p$  (independent of  $X$ ) such that

$$A_p \|\sqrt{[X]_n}\|_p \leq \|X_n\|_p \leq B_p \|\sqrt{[X]_n}\|_p, \quad (27)$$

for every  $n \geq 1$ , where  $[X]_n$  is the quadratic variation of  $X_n$ ,

$$[X]_n = \sum_{j=1}^n (\Delta X_j)^2, \quad X_0 = 0. \quad (28)$$

The constants  $A_p$  and  $B_p$  can be taken to have the values

$$A_p = [18p^{3/2}/(p-1)]^{-1}, \quad B_p = 18p^{3/2}/(p-1)^{1/2}.$$

It follows from (17), by using (2), that

$$A_p \|\sqrt{[X]_n}\|_p \leq \|X_n^*\|_p \leq B_p^* \|\sqrt{[X]_n}\|_p, \quad (29)$$

where

$$A_p = [18p^{3/2}/(p-1)]^{-1}, \quad B_p^* = 18p^{5/2}/(p-1)^{3/2}.$$

Burkholder's inequalities (27) hold for  $p > 1$ , whereas the Marcinkiewicz–Zygmund inequalities (26) also hold when  $p = 1$ . What can we say about the validity of (27) for  $p = 1$ ? It turns out that a direct generalization to  $p = 1$  is impossible, as the following example shows.

**EXAMPLE.** Let  $\xi_1, \xi_2, \dots$  be independent Bernoulli random variables with  $P(\xi_i = 1) = P(\xi_i = -1) = \frac{1}{2}$  and let

$$X_n = \sum_{j=1}^{n \wedge \tau} \xi_j,$$

where

$$\tau = \inf \left\{ n \geq 1: \sum_{i=1}^n \xi_i = 1 \right\}.$$

The sequence  $X = (X_n, \mathcal{F}_n^X)$  is a martingale with

$$\|X_n\|_1 = E|X_n| = 2EX_n^+ \rightarrow 2, \quad n \rightarrow \infty.$$

But

$$\|\sqrt{[X]_n}\|_1 = E\sqrt{[X]_n} = E\left(\sum_{j=1}^{n \wedge \tau} 1\right)^{1/2} = E\sqrt{\tau \wedge n} \rightarrow \infty.$$

Consequently the first inequality in (27) fails.

It turns out that when  $p = 1$  we must generalize not (27), but (29) (which is equivalent when  $p > 1$ ).

**Davis's Inequality.** If  $X = (X_n, \mathcal{F}_n)$  is a martingale, there are universal



constants  $A$  and  $B$ ,  $0 < A < B < \infty$ , such that

$$A\|\sqrt{[X]_n}\|_1 \leq \|X_n^*\|_1 \leq B\|\sqrt{[X]_n}\|_1, \quad (30)$$

i.e.,

$$AE\sqrt{\sum_{j=1}^n (\Delta X_j)^2} \leq E\left[\max_{1 \leq j \leq n} |X_n|\right] \leq BE\sqrt{\sum_{j=1}^n (\Delta X_j)^2}.$$

**Corollary 1.** Let  $\xi_1, \xi_2, \dots$  be independent identically distributed random variables;  $S_n = \xi_1 + \dots + \xi_n$ . If  $E|\xi_1| < \infty$  and  $E\xi_1 = 0$ , then according to Wald's inequality (2.14) we have

$$ES_\tau = 0 \quad (31)$$

for every stopping time  $\tau$  (with respect to  $(\mathcal{F}_n^\xi)$ ) for which  $E\tau < \infty$ .

It turns out that (31) is still valid under a weaker hypothesis than  $E\tau < \infty$  if we impose stronger conditions on the random variables. In fact, if

$$E|\xi_1|^r < \infty,$$

where  $1 < r \leq 2$ , the condition  $E\tau^{1/r} < \infty$  is a sufficient condition for  $ES_\tau = 0$ .

For the proof, we put  $\tau_n = \tau \wedge n$ ,  $Y = \sup_n |S_{\tau_n}|$ , and let  $m = [t^r]$  (integral part of  $t^r$ ) for  $t > 0$ . By Corollary 1 to Theorem 1, §2, we have  $ES_{\tau_n} = 0$ . Therefore a sufficient condition for  $ES_\tau = 0$  is (by the dominated convergence theorem) that  $E \sup_n |S_{\tau_n}| < \infty$ .

Using (1) and (27), we obtain

$$\begin{aligned} P(Y \geq t) &= P(\tau \geq t^r, Y \geq t) + P(\tau < t^r, Y \geq t) \\ &\leq P(\tau \geq t^r) + P\left\{\max_{1 \leq j \leq m} |S_{\tau_j}| \geq t\right\} \\ &\leq P(\tau \geq t^r) + t^{-r}E|S_{\tau_m}|^r \\ &\leq P(\tau \geq t^r) + t^{-r}B_r E\left(\sum_{j=1}^{\tau_m} \xi_j^2\right)^{r/2} \\ &\leq P(\tau \geq t^r) + t^{-r}B_r E \sum_{j=1}^{\tau_m} |\xi_j|^r. \end{aligned}$$

Notice that (with  $\mathcal{F}_0^\xi = \{\emptyset, \Omega\}$ )

$$\begin{aligned} E \sum_{j=1}^{\tau_m} |\xi_j|^r &= E \sum_{j=1}^{\infty} I(j \leq \tau_m) |\xi_j|^r \\ &= \sum_{j=1}^{\infty} EE[I(j \leq \tau_m) |\xi_j|^r | \mathcal{F}_{j-1}^\xi] \\ &= E \sum_{j=1}^{\infty} I(j \leq \tau_m) E[|\xi|^r | \mathcal{F}_{j-1}^\xi] = E \sum_{j=1}^{\tau_m} E|\xi_j|^r = \mu_r E\tau_m, \end{aligned}$$

where  $\mu_r = E|\xi_1|^r$ . Consequently

$$\begin{aligned} P(Y \geq t) &\leq P(\tau \geq t^r) + t^{-r} B_r \mu_r E \tau_m \\ &= P(\tau \geq t^r) + B_r \mu_r t^{-r} \left[ m P(\tau \geq t^r) + \int_{\{\tau < t^r\}} \tau dP \right] \\ &\leq (1 + B_r \mu_r) P(\tau \geq t^r) + B_r \mu_r t^{-r} \int_{\{\tau < t^r\}} \tau dP \end{aligned}$$

and therefore

$$\begin{aligned} EY \int_0^\infty P(Y \geq t) dt &\leq (1 + B_r \mu_r) E \tau^{1/r} + B_r \mu_r \int_0^\infty t^{-r} \left[ \int_{\{\tau < t^r\}} \tau dP \right] dt \\ &= (1 + B_r \mu_r) E \tau^{1/r} + B_r \mu_r \int_\Omega \tau \left[ \int_{\tau^{1/r}}^\infty t^{-r} dt \right] dP \\ &= \left( 1 + B_r \mu_r + \frac{B_r \mu_r}{r-1} \right) E \tau^{1/r} < \infty. \end{aligned}$$

**Corollary 2.** Let  $M = (M_n)$  be a martingale with  $E|M_n|^{2r} < \infty$  for some  $r \geq 1$  and such that (with  $M_0 = 0$ )

$$\sum_{n=1}^\infty \frac{E|\Delta M_n|^{2r}}{n^{1+r}} < \infty. \quad (32)$$

Then (compare Theorem 2 of §3, Chapter IV) we have the strong law of large numbers:

$$\frac{M_n}{n} \rightarrow 0 \quad (\text{P-a.s.}), \quad n \rightarrow \infty. \quad (33)$$

When  $r = 1$  the proof follows the same lines as the proof of Theorem 2, §3, Chapter IV. In fact, let

$$m_n = \sum_{k=1}^n \frac{\Delta M_k}{k}.$$

Then

$$\frac{M_n}{n} = \frac{\sum_{k=1}^n \Delta M_k}{n} = \frac{1}{n} \sum_{k=1}^n k \Delta m_k$$

and, by Kronecker's lemma (§3, Chapter IV) a sufficient condition for the limit relation (P-a.s.)

$$\frac{1}{n} \sum_{k=1}^n k \Delta m_k \rightarrow 0, \quad n \rightarrow \infty,$$

is that the limit  $\lim_n m_n$  exists and is finite (P-a.s.) which in turn (Theorems 1 and 4, §10, Chapter II) is true if and only if

$$P\left\{\sup_{k \geq 1} |m_{n+k} - m_n| \geq \varepsilon\right\} \rightarrow 0, \quad n \rightarrow \infty. \quad (34)$$

By (1),

$$\mathbf{P}\left\{\sup_{k \geq 1} |m_{n+k} - m_n| \geq \varepsilon\right\} \leq \varepsilon^{-2} \sum_{k=n}^{\infty} \frac{\mathbf{E}(\Delta M_k)^2}{k^2}.$$

Hence the required result follows from (32) and (34).

Now let  $r > 1$ . Then the statement (33) is equivalent (Theorem 1, §10, Chapter II) to the statement that

$$\varepsilon^{2r} \mathbf{P}\left\{\sup_{j \geq n} \frac{|M_j|}{j} \geq \varepsilon\right\} \rightarrow 0, \quad n \rightarrow \infty \quad (35)$$

for every  $\varepsilon > 0$ . By inequality (53) of Problem 1,

$$\begin{aligned} \varepsilon^{2r} \mathbf{P}\left\{\sup_{j \geq n} \frac{|M_j|}{j} \geq \varepsilon\right\} &= \varepsilon^{2r} \lim_{m \rightarrow \infty} \mathbf{P}\left\{\max_{n \leq j \leq m} \frac{|M_j|^{2r}}{j^{2r}} \geq \varepsilon^{2r}\right\} \\ &\leq \frac{1}{n^{2r}} \mathbf{E}|M_n|^{2r} + \sum_{j \geq n+1} \frac{1}{j^{2r}} \mathbf{E}(|M_j|^{2r} - |M_{j-1}|^{2r}). \end{aligned}$$

It follows from Kronecker's lemma and (32) that

$$\lim_{n \rightarrow \infty} \frac{1}{n^{2r}} \mathbf{E}|M_n|^{2r} = 0.$$

Hence to prove (35) we need only prove that

$$\sum_{j \geq 2} \frac{1}{j^{2r}} \mathbf{E}(|M_j|^{2r} - |M_{j-1}|^{2r}) < \infty. \quad (36)$$

We have

$$\begin{aligned} I_N &= \sum_{j=2}^N \frac{1}{j^{2r}} [\mathbf{E}|M_j|^{2r} - \mathbf{E}|M_{j-1}|^{2r}] \\ &\leq \sum_{j=3}^N \left[ \frac{1}{(j-1)^{2r}} - \frac{1}{j^{2r}} \right] \mathbf{E}|M_{j-1}|^{2r} + \frac{\mathbf{E}|M_N|^{2r}}{N^{2r}}. \end{aligned}$$

By Burkholder's inequality (27) and Hölder's inequality,

$$\mathbf{E}|M_j|^{2r} \leq \mathbf{E}\left[\sum_{i=1}^j (\Delta M_i)^2\right]^r \leq \mathbf{E}j^{r-1} \sum_{i=1}^j |\Delta M_i|^{2r}.$$

Hence

$$\begin{aligned} I_N &\leq \sum_{j=2}^{N-1} \left[ \frac{1}{j^{2r}} - \frac{1}{(j+1)^{2r}} \right] j^{r-1} \sum_{i=1}^j \mathbf{E}|\Delta M_i|^{2r} \\ &\leq C_1 \sum_{j=2}^{N-1} \frac{1}{j^{r+2}} \sum_{i=1}^j \mathbf{E}|\Delta M_i|^{2r} \leq C_2 \sum_{j=2}^N \frac{\mathbf{E}|\Delta M_j|^{2r}}{j^{r+1}} + C_3 \end{aligned}$$

( $C_i$  are constants). By (32), this establishes (36).

4. The sequence of random variables  $\{X_n\}_{n \geq 1}$  has a limit  $\lim X_n$  (finite or infinite) with probability 1, if and only if the number of "oscillations between two arbitrary rational numbers  $a$  and  $b$ ,  $a < b$ " is finite with probability 1. Theorem 5, below, provides an upper bound for the number of "oscillations" for submartingales. In the next section, this will be applied to prove the fundamental result on their convergence.

Let us choose two numbers  $a$  and  $b$ ,  $a < b$ , and define the following times in terms of the stochastic sequence  $X = (X_n, \mathcal{F}_n)$ :

$$\begin{aligned}\tau_0 &= 0, \\ \tau_1 &= \min\{n > 0: X_n \leq a\}, \\ \tau_2 &= \min\{n > \tau_1: X_n \geq b\}, \\ &\dots\dots\dots \\ \tau_{2m-1} &= \min\{n > \tau_{2m-2}: X_n \leq a\}, \\ \tau_{2m} &= \min\{n > \tau_{2m-1}: X_n \geq b\},\end{aligned}$$

taking  $\tau_k = \infty$  if the corresponding set  $\{\cdot\}$  is empty.

In addition, for each  $n \geq 1$  we define the random variables

$$\beta_n(a, b) = \begin{cases} 0, & \text{if } \tau_2 > n, \\ \max\{m: \tau_{2m} \leq n\} & \text{if } \tau_2 \leq n. \end{cases}$$

In words,  $\beta_n(a, b)$  is the number of upcrossings of  $[a, b]$  by the sequence  $X_1, \dots, X_n$ .

**Theorem 5** (Doob). Let  $X = (X_n, \mathcal{F}_n)_{n \geq 1}$  be a submartingale. Then, for every  $n \geq 1$ ,

$$E\beta_n(a, b) \leq \frac{E[X_n - a]^+}{b - a}. \quad (37)$$

PROOF. The number of intersections of  $X = (X_n, \mathcal{F}_n)$  with  $[a, b]$  is equal to the number of intersections of the nonnegative submartingale  $X^+ = ((X_n - a)^+, \mathcal{F}_n)$  with  $[0, b - a]$ . Hence it is sufficient to suppose that  $X$  is nonnegative with  $a = 0$ , and show that

$$E\beta_n(0, b) \leq \frac{EX_n}{b}. \quad (38)$$

Put  $X_0 = 0$ ,  $\mathcal{F}_0 = \{\emptyset, \Omega\}$ , and for  $i = 1, 2, \dots$ , let

$$\varphi_i = \begin{cases} 1 & \text{if } \tau_m < i \leq \tau_{m+1} \text{ for some odd } m, \\ 0 & \text{if } \tau_m < i \leq \tau_{m+1} \text{ for some even } m. \end{cases}$$

It is easily seen that

$$b\beta_n(0, b) \leq \sum_{i=1}^n \varphi_i [X_i - X_{i-1}]$$

and

$$\{\varphi_i = 1\} = \bigcup_{\text{odd } m} [\{\tau_m < i\} \setminus \{\tau_{m+1} < i\}] \in \mathcal{F}_{i-1}.$$

Therefore

$$\begin{aligned}
 bE\beta_n(0, b) &\leq E \sum_{i=1}^n \varphi_i[X_i - X_{i-1}] = \sum_{i=1}^n \int_{\{\varphi_i=1\}} (X_i - X_{i-1}) dP \\
 &= \sum_{i=1}^n \int_{\{\varphi_i=1\}} E(X_i - X_{i-1} | \mathcal{F}_{i-1}) dP \\
 &= \sum_{i=1}^n \int_{\{\varphi_i=1\}} [E(X_i | \mathcal{F}_{i-1}) - X_{i-1}] dP \\
 &\leq \sum_{i=1}^n \int_{\Omega} [E(X_i | \mathcal{F}_{i-1}) - X_{i-1}] dP = EX_n,
 \end{aligned}$$

which establishes (38).

5. In this subsection we discuss some of the simplest inequalities for the probability of large deviations for martingales of integrable square.

Let  $M = (M_n, \mathcal{F}_n)_{n \geq 0}$  be a martingale of integrable square with quadratic variation  $\langle M \rangle = (\langle M \rangle_n, \mathcal{F}_{n-1})$ . If we apply inequality (22) to  $X_n = M_n$ ,  $A_n = \langle M \rangle_n$ , we find that for  $a > 0$  and  $b > 0$

$$\begin{aligned}
 P\left\{\max_{k \leq n} |M_k| \geq an\right\} &= P\left\{\max_{k \leq n} M_k^2 \geq (an)^2\right\} \\
 &\leq \frac{1}{(an)^2} E[\langle M \rangle_n \wedge (bn)] + P\{\langle M \rangle_n \geq an\}. \quad (39)
 \end{aligned}$$

In fact, at least in the case when  $|\Delta M_n| \leq C$  for all  $n$  and  $\omega \in \Omega$ , this inequality can be substantially improved by using the ideas explained in §5, Chapter IV for estimating the probability of large deviations for sums of independent identically distributed random variables.

Let us recall that in §5, Chapter IV, when we introduced the corresponding inequalities, the essential point was to use the property that the sequence

$$(e^{\lambda S_n} / [\varphi(\lambda)]^n, \mathcal{F}_n)_{n \geq 1}, \quad \mathcal{F}_n = \sigma\{\xi_1, \dots, \xi_n\}, \quad (40)$$

formed a nonnegative martingale, to which we could apply the inequality (8). If we now take  $M_n$  instead of  $S_n$ , by analogy with (40), the martingale

$$(e^{\lambda M_n} / \mathcal{E}_n(\lambda), \mathcal{F}_n)_{n \geq 1},$$

will be nonnegative, where

$$\mathcal{E}_n(\lambda) = \prod_{j=1}^n E(e^{\lambda \Delta M_j} | \mathcal{F}_{j-1}) \quad (41)$$

is called the *stochastic exponential*.

This expression is rather complicated. At the same time, in using (8) it is not necessary for the sequence to be a *martingale*. It is enough for it to be a nonnegative *supermartingale*. Here we can arrange this by forming a

sequence  $(Z_n(\lambda), \mathcal{F}_n)$  ((43), below), which sufficiently depends simply on  $M_n$  and  $\langle M \rangle_n$ , and to which we can apply the method used in §5, Chapter IV.

**Lemma 1.** Let  $M = (M_n, \mathcal{F}_n)_{n \geq 0}$  be a square-integrable martingale,  $M_0 = 0$ ,  $\Delta M_0 = 0$ , and  $|\Delta M_n(\omega)| \leq c$  for all  $n$  and  $\omega$ . Let  $\lambda > 0$ ,

$$\psi_c(\lambda) = \begin{cases} \frac{e^{\lambda c} - 1 - \lambda c}{c^2}, & c > 0, \\ \frac{\lambda^2}{2}, & c = 0, \end{cases} \quad (42)$$

and

$$Z_n(\lambda) = e^{\lambda M_n - \psi_c(\lambda) \langle M \rangle_n}. \quad (43)$$

Then for every  $c \geq 0$  the sequence  $Z(\lambda) = (Z_n(\lambda), \mathcal{F}_n)_{n \geq 0}$  is a non-negative supermartingale.

PROOF. For  $|x| \leq c$ ,

$$e^{\lambda x} - 1 - \lambda x = (\lambda x)^2 \sum_{m \geq 2} \frac{(\lambda x)^{m-2}}{m!} \leq (\lambda x)^2 \sum_{m \geq 2} \frac{(\lambda c)^{m-2}}{m!} \leq x^2 \psi_c(\lambda).$$

Using this inequality and the following representation  $(Z_n = Z_n(\lambda))$

$$\Delta Z_n = Z_{n-1} [(e^{\lambda \Delta M_n} - 1) e^{-\Delta \langle M \rangle_n \psi_c(\lambda)} + (e^{-\Delta \langle M \rangle_n \psi_c(\lambda)} - 1)],$$

we find that

$$\begin{aligned} & \mathbb{E}(\Delta Z_n | \mathcal{F}_{n-1}) \\ &= Z_{n-1} [\mathbb{E}(e^{\lambda \Delta M_n} - 1 | \mathcal{F}_{n-1}) e^{-\Delta \langle M \rangle_n \psi_c(\lambda)} + (e^{-\Delta \langle M \rangle_n \psi_c(\lambda)} - 1)] \\ &= Z_{n-1} [\mathbb{E}(e^{\lambda \Delta M_n} - 1 - \lambda \Delta M_n | \mathcal{F}_{n-1}) e^{-\Delta \langle M \rangle_n \psi_c(\lambda)} + (e^{-\Delta \langle M \rangle_n \psi_c(\lambda)} - 1)] \\ &\leq Z_{n-1} [\psi_c(\lambda) \mathbb{E}((\Delta M_n)^2 | \mathcal{F}_{n-1}) e^{-\Delta \langle M \rangle_n \psi_c(\lambda)} + (e^{-\Delta \langle M \rangle_n \psi_c(\lambda)} - 1)] \\ &= Z_{n-1} [\psi_c(\lambda) \Delta \langle M \rangle_n e^{-\Delta \langle M \rangle_n \psi_c(\lambda)} + (e^{-\Delta \langle M \rangle_n \psi_c(\lambda)} - 1)] \leq 0, \end{aligned} \quad (44)$$

where we have also used the fact that, for  $x \geq 0$ ,

$$xe^{-x} + (e^{-x} - 1) \leq 0.$$

We see from (44) that

$$\mathbb{E}(Z_n | \mathcal{F}_{n-1}) \leq Z_{n-1},$$

i.e.,  $Z(\lambda) = (Z_n(\lambda), \mathcal{F}_n)$  is a supermartingale.

This establishes the lemma.

Let the hypotheses of the lemma be satisfied. Then we can always find  $\lambda > 0$  for which, for given  $a > 0$  and  $b > 0$ , we have  $a\lambda - b\psi_c(\lambda) > 0$ . From this, we obtain

$$\begin{aligned}
P\left\{\max_{k \leq n} M_k \geq an\right\} &= P\left\{\max_{k \leq n} e^{\lambda M_k} \geq e^{\lambda an}\right\} \\
&\leq P\left\{\max_{k \leq n} e^{\lambda M_k - \psi_c(\lambda) \langle M \rangle_k} \geq e^{\lambda an - \psi_c(\lambda) \langle M \rangle_n}\right\} \\
&= P\left\{\max_{k \leq n} e^{\lambda M_k - \psi_c(\lambda) \langle M \rangle_k} \geq e^{\lambda an - \psi_c(\lambda) \langle M \rangle_n}, \langle M \rangle_n \leq bn\right\} \\
&\quad + P\left\{\max_{k \leq n} e^{\lambda M_k - \psi_c(\lambda) \langle M \rangle_k} \geq e^{\lambda an - \psi_c(\lambda) \langle M \rangle_n}, \langle M \rangle_n > bn\right\} \\
&\leq P\left\{\max_{k \leq n} e^{\lambda M_k - \psi_c(\lambda) \langle M \rangle_k} \geq e^{\lambda an - \psi_c(\lambda) bn}\right\} \quad (45)
\end{aligned}$$

where the last inequality follows from (7).

Let us write

$$H_c(a, b) = \sup_{\lambda > 0} [a\lambda - b\psi_c(\lambda)].$$

Then it follows from (45) that

$$P\left\{\max_{k \leq n} M_k \geq an\right\} \leq P\{\langle M \rangle_n > bn\} + e^{-nH_c(a, b)}. \quad (46)$$

Passing from  $M$  to  $-M$ , we find that the right-hand side of (46) also provides an upper bound for the probability  $P\{\min_{k \leq n} M_k \leq -an\}$ . Consequently,

$$P\left\{\max_{k \leq n} |M_k| \geq an\right\} \leq 2P\{\langle M \rangle_n > bn\} + 2e^{-nH_c(a, b)}. \quad (47)$$

Thus, we have proved the following theorem.

**Theorem 6.** Let  $M = (M_n, \mathcal{F}_n)$  be a martingale with uniformly bounded steps, i.e.,  $|\Delta M_n| \leq c$  for some constant  $c > 0$  and all  $n$  and  $\omega$ . Then for every  $a > 0$  and  $b > 0$ , we have the inequalities (46) and (47).

**Remark 2.**

$$H_c(a, b) = \frac{1}{c} \left(a + \frac{b}{c}\right) \ln \left(1 + \frac{ac}{b}\right) - \frac{a}{c}. \quad (48)$$

6. Under the hypotheses of Theorem 6, we now consider the question of estimates of probabilities of the type

$$P\left\{\sup_{k \geq n} \frac{M_k}{\langle M \rangle_k} > a\right\},$$

which characterize, in particular, the rapidity of convergence in the strong law of large numbers for martingales (also see Theorem 4 in §5).

Proceeding as in §5, Chapter IV, we find that for every  $a > 0$  there is a  $\lambda > 0$  for which  $a\lambda - \psi_c(\lambda) > 0$ . Then, for every  $b > 0$ ,

$$\begin{aligned} \mathbf{P}\left\{\sup_{k \geq n} \frac{M_k}{\langle M \rangle_k} > a\right\} &\leq \mathbf{P}\left\{\sup_{k \geq n} e^{\lambda M_k - \psi_c(\lambda) \langle M \rangle_k} > e^{[a\lambda - \psi_c(\lambda)] \langle M \rangle_n}\right\} \\ &\leq \mathbf{P}\left\{\sup_{k \geq n} e^{\lambda M_k - \psi_c(\lambda) \langle M \rangle_k} > e^{[a\lambda - \psi_c(\lambda)]bn}\right\} \\ &\quad + \mathbf{P}\{\langle M \rangle_n < bn\} \leq e^{-bn[a\lambda - \psi_c(\lambda)]} + \mathbf{P}\{\langle M \rangle_n < bn\}, \end{aligned} \quad (49)$$

from which

$$\mathbf{P}\left\{\sup_{k \geq n} \frac{M_k}{\langle M \rangle_k} > a\right\} \leq \mathbf{P}\{\langle M \rangle_n < bn\} + e^{-nH_c(ab, b)} \quad (50)$$

$$\mathbf{P}\left\{\sup_{k \geq n} \left|\frac{M_k}{\langle M \rangle_k}\right| > a\right\} \leq 2\mathbf{P}\{\langle M \rangle_n < bn\} + 2e^{-nH_c(ab, b)}. \quad (51)$$

We have therefore proved the following theorem.

**Theorem 7.** *Let the hypotheses of the preceding theorem be satisfied. Then inequalities (50) and (51) are satisfied for all  $a > 0$  and  $b > 0$ .*

**Remark 3.** Comparison of (51) with the estimate (21) in §5, Chapter IV, for the case of a Bernoulli scheme,  $p = 1/2$ ,  $M_n = S_n - (n/2)$ ,  $b = 1/4$ ,  $c = 1/2$ , shows that for small  $\varepsilon > 0$  it leads to the same result

$$\mathbf{P}\left\{\sup_{k \geq n} \left|\frac{M_k}{\langle M \rangle_k}\right| > \varepsilon\right\} = \mathbf{P}\left\{\sup_{k \geq n} \left|\frac{S_k - k/2}{k}\right| > \frac{\varepsilon}{4}\right\} \leq 2e^{-4\varepsilon^2 n}.$$

## 7. PROBLEMS

1. Let  $X = (X_n, \mathcal{F}_n)$  be a nonnegative submartingale and let  $V = (V_n, \mathcal{F}_{n-1})$  be a predictable sequence such that  $0 \leq V_{n+1} \leq V_n \leq C$  (P-a.s.), where  $C$  is a constant. Establish the following generalization of (1):

$$\varepsilon \mathbf{P}\left\{\max_{1 \leq j \leq n} V_j X_j \geq \varepsilon\right\} + \int_{\{\max_{1 \leq j \leq n} V_j X_j < \varepsilon\}} V_n X_n d\mathbf{P} \leq \sum_{j=1}^n \mathbf{E} V_j \Delta X_j. \quad (52)$$

2. Establish *Krickeberg's decomposition*: every martingale  $X = (X_n, \mathcal{F}_n)$  with  $\sup \mathbf{E}|X_n| < \infty$  can be represented as the difference of two nonnegative martingales.
3. Let  $\xi_1, \xi_2, \dots$  be a sequence of independent random variables,  $S_n = \xi_1 + \dots + \xi_n$  and  $S_{m,n} = \sum_{j=m+1}^n \xi_j$ . Establish *Ottaviani's inequality*:

$$\mathbf{P}\left\{\max_{1 \leq j \leq n} |S_j| > 2\varepsilon\right\} \leq \frac{\mathbf{P}\{|S_n| > \varepsilon\}}{\min_{1 \leq j \leq n} \mathbf{P}\{|S_{j,n}| \leq \varepsilon\}}$$



and deduce that

$$\int_0^\infty \mathbf{P}\left\{\max_{1 \leq j \leq n} |S_j| > 2t\right\} dt \leq 2\mathbf{E}|S_n| + 2 \int_{2\mathbf{E}|S_n|}^\infty \mathbf{P}\{|S_n| > t\} dt. \quad (53)$$

4. Let  $\xi_1, \xi_2, \dots$  be a sequence of independent random variables with  $\mathbf{E}\xi_i = 0$ . Use (53) to show that in this case we can strengthen inequality (10) to

$$\mathbf{E}S_n^* \leq 8\mathbf{E}|S_n|.$$

5. Verify formula (16).

6. Establish inequality (19).

7. Let the  $\sigma$ -algebra  $\mathcal{F}_0, \dots, \mathcal{F}_n$  be such that  $\mathcal{F}_0 \subseteq \mathcal{F}_1 \subseteq \dots \subseteq \mathcal{F}_n$  and let the events  $A_k \in \mathcal{F}_k, k = 1, \dots, n$ . Use (22) to establish *Dvoretzky's inequality*: for each  $\varepsilon > 0$ ,

$$\mathbf{P}\left[\bigcup_{k=1}^n A_k \mid \mathcal{F}_0\right] \leq \varepsilon + \mathbf{P}\left[\sum_{k=1}^n \mathbf{P}(A_k \mid \mathcal{F}_{k-1}) > \varepsilon \mid \mathcal{F}_0\right] \quad (\text{P-a.s.}).$$

## §4. General Theorems on the Convergence of Submartingales and Martingales

1. The following result, which is fundamental for all problems about the convergence of submartingales, can be thought of as an analog of the fact that in real analysis a bounded monotonic sequence of numbers has a (finite) limit.

**Theorem 1 (Doob).** *Let  $X = (X_n, \mathcal{F}_n)$  be a submartingale with*

$$\sup_n \mathbf{E}|X_n| < \infty. \quad (1)$$

*Then with probability 1, the limit  $\lim X_n = X_\infty$  exists and  $\mathbf{E}|X_\infty| < \infty$ .*

**PROOF.** Suppose that

$$\mathbf{P}(\overline{\lim} X_n > \underline{\lim} X_n) > 0. \quad (2)$$

Then since

$$\{\overline{\lim} X_n > \underline{\lim} X_n\} = \bigcup_{a < b} \{\overline{\lim} X_n > b > a > \underline{\lim} X_n\}$$

(here  $a$  and  $b$  are rational numbers), there are values  $a$  and  $b$  such that

$$\mathbf{P}\{\overline{\lim} X_n > b > a > \underline{\lim} X_n\} > 0. \quad (3)$$

Let  $\beta_n(a, b)$  be the number of upcrossings of  $(a, b)$  by the sequence  $X_1, \dots, X_n$ , and let  $\beta_\infty(a, b) = \lim_n \beta_n(a, b)$ . By (3.27),

$$E\beta_n(a, b) \leq \frac{E[X_n - a]^+}{b - a} \leq \frac{EX_n^+ + |a|}{b - a}$$

and therefore

$$E\beta_\infty(a, b) = \lim_n E\beta_n(a, b) \leq \frac{\sup_n EX_n^+ + |a|}{b - a} < \infty,$$

which follows from (1) and the remark that

$$\sup_n E|X_n| < \infty \Leftrightarrow \sup_n EX_n^+ < \infty$$

for submartingales (since  $EX_n^+ \leq E|X_n| = 2EX_n^+ - EX_n \leq 2EX_n^+ - EX_1$ ). But the condition  $E\beta_\infty(a, b) < \infty$  contradicts assumption (3). Hence  $\lim X_n = X_\infty$  exists with probability 1, and then by Fatou's lemma

$$E|X_\infty| \leq \sup_n E|X_n| < \infty.$$

This completes the proof of the theorem.

**Corollary 1.** *If  $X$  is a nonpositive submartingale, then with probability 1 the limit  $\lim X_n$  exists and is finite.*

**Corollary 2.** *If  $X = (X_n, \mathcal{F}_n)_{n \geq 1}$  is a nonpositive submartingale, the sequence  $\bar{X} = (X_n, \bar{\mathcal{F}}_n)$  with  $1 \leq n \leq \infty$ ,  $X_\infty = \lim X_n$  and  $\bar{\mathcal{F}}_\infty = \sigma\{\bigcup \mathcal{F}_n\}$  is a (non-positive) submartingale.*

In fact, by Fatou's lemma

$$EX_\infty = E \lim X_n \geq \overline{\lim} EX_n \geq EX_1 > -\infty$$

and (P-a.s.)

$$E(X_\infty | \bar{\mathcal{F}}_m) = E(\lim X_n | \bar{\mathcal{F}}_m) \geq \overline{\lim} E(X_n | \bar{\mathcal{F}}_m) \geq X_m.$$

**Corollary 3.** *If  $X = (X_n, \mathcal{F}_n)$  is a nonnegative martingale, then  $\lim X_n$  exists with probability 1.*

In fact, in that case

$$\sup E|X_n| = \sup EX_n = EX_1 < \infty,$$

and Theorem 1 is applicable.

2. Let  $\xi_1, \xi_2, \dots$  be a sequence of independent random variables with  $P(\xi_i = 0) = P(\xi_i = 2) = \frac{1}{2}$ . Then  $X = (X_n, \mathcal{F}_n^{\xi})$ , with  $X_n = \prod_{i=1}^n \xi_i$  and  $\mathcal{F}_n^{\xi} = \sigma\{\omega: \xi_1, \dots, \xi_n\}$  is a martingale with  $EX_n = 1$  and  $X_n \rightarrow X_\infty \equiv 0$  (P-a.s.). At the same time, it is clear that  $E|X_n - X_\infty| = 1$  and therefore  $X_n \not\rightarrow X_\infty$  in the  $L^1$  sense. Therefore condition (1) does not in general guarantee the convergence of  $X_n$  to  $X_\infty$  in the  $L^1$  sense.

Theorem 2 below shows that if hypothesis (1) is strengthened to uniform integrability of the family  $\{X_n\}$  (from which (1) follows by Subsection 4,

§6, Chapter II), then besides almost sure convergence we also have convergence in  $L^1$ .

**Theorem 2.** Let  $X = \{X_n, \mathcal{F}_n\}$  be a uniformly integrable submartingale (that is, the family  $\{X_n\}$  is uniformly integrable). Then there is a random variable  $X_\infty$  with  $E|X_\infty| < \infty$ , such that as  $n \rightarrow \infty$

$$X_n \rightarrow X_\infty \quad (\text{P-a.s.}), \quad (4)$$

$$X_n \xrightarrow{L^1} X_\infty. \quad (5)$$

Moreover, the sequence  $\bar{X} = (X_n, \mathcal{F}_n)$ ,  $1 \leq n \leq \infty$ , with  $\mathcal{F}_\infty = \sigma(\bigcup \mathcal{F}_n)$ , is also a submartingale.

**PROOF.** Statement (4) follows from Theorem 1, and (5) follows from (4) and Theorem 4, §6, Chapter II.

Moreover, if  $A \in \mathcal{F}_n$  and  $m \geq n$ , then

$$EI_A|X_m - X_\infty| \rightarrow 0, \quad m \rightarrow \infty,$$

and therefore

$$\lim_{m \rightarrow \infty} \int_A X_m dP = \int_A X_\infty dP.$$

The sequence  $\{\int_A X_m dP\}_{m \geq n}$  is nondecreasing and therefore

$$\int_A X_n dP \leq \int_A X_m dP \leq \int_A X_\infty dP,$$

whence  $X_n \leq E(X_\infty | \mathcal{F}_n)$  (P-a.s.) for  $n \geq 1$ .

This completes the proof of the theorem.

**Corollary.** If  $X = (X_n, \mathcal{F}_n)$  is a submartingale and, for some  $p > 1$ ,

$$\sup_n E|X_n|^p < \infty, \quad (6)$$

then there is an integrable random variable  $X_\infty$  for which (4) and (5) are satisfied.

For the proof, it is enough to observe that, by Lemma 3 of §6 of Chapter II, condition (6) guarantees the uniform integrability of the family  $\{X_n\}$ .

**3.** We now present a theorem on the continuity properties of conditional expectations. This was one of the very first results concerning the convergence of martingales.

**Theorem 3 (P. Lévy).** Let  $(\Omega, \mathcal{F}, P)$  be a probability space, and let  $(\mathcal{F}_n)_{n \geq 1}$  be a nondecreasing family of  $\sigma$ -algebras,  $\mathcal{F}_1 \subseteq \mathcal{F}_2 \subseteq \dots \subseteq \mathcal{F}$ . Let  $\xi$  be a random variable with  $E|\xi| < \infty$  and  $\mathcal{F}_\infty = \sigma(\bigcup_n \mathcal{F}_n)$ . Then, both P-a.s. and in the  $L^1$  sense,

$$E(\xi | F_n) \rightarrow E(\xi | F_\infty), \quad n \rightarrow \infty. \quad (7)$$

PROOF. Let  $X_n = E(\xi | \mathcal{F}_n)$ ,  $n \geq 1$ . Then, with  $a > 0$  and  $b > 0$ ,

$$\begin{aligned} \int_{\{|X_i| \geq a\}} |X_i| dP &\leq \int_{\{|X_i| \geq a\}} E(|\xi| | F_i) dP = \int_{\{|X_i| \geq a\}} |\xi| dP \\ &\leq \int_{\{|X_i| \geq a\} \cap \{|\xi| \leq b\}} |\xi| dP + \int_{\{|X_i| \geq a\} \cap \{|\xi| > b\}} |\xi| dP \\ &\leq bP\{|X_i| \geq a\} + \int_{\{|\xi| > b\}} |\xi| dP \\ &\leq \frac{b}{a} E|X_i| + \int_{\{|\xi| > b\}} |\xi| dP \\ &\leq \frac{b}{a} E|\xi| + \int_{\{|\xi| > b\}} |\xi| dP. \end{aligned}$$

Letting  $a \rightarrow \infty$  and then  $b \rightarrow \infty$ , we obtain

$$\lim_{a \rightarrow \infty} \sup_i \int_{\{|X_i| \geq a\}} |X_i| dP = 0,$$

i.e., the family  $\{X_n\}$  is uniformly integrable.

Therefore, by Theorem 2, there is a random variable  $X_\infty$  such that  $X_n = E(\xi | F_n) \rightarrow X_\infty$  ((P-a.s.) and in the  $L^1$  sense). Hence we only have to show that

$$X_\infty = E(\xi | \mathcal{F}_\infty) \quad (\text{P-a.s.}).$$

Let  $m \geq n$  and  $A \in \mathcal{F}_n$ . Then

$$\int_A X_m dP = \int_A X_n dP = \int_A E(\xi | F_n) dP = \int_A \xi dP.$$

Since the family  $\{X_n\}$  is uniformly integrable and since, by Theorem 5, §6, Chapter II, we have  $E|X_m - X_\infty| \rightarrow 0$  as  $m \rightarrow \infty$ , it follows that

$$\int_A X_\infty dP = \int_A \xi dP. \quad (8)$$

This equation is satisfied for all  $A \in \mathcal{F}_n$  and therefore for all  $A \in \bigcup_{n=1}^{\infty} \mathcal{F}_n$ . Since  $E|X_\infty| < \infty$  and  $E|\xi| < \infty$ , the left-hand and right-hand sides of (8) are  $\sigma$ -additive measures; possibly taking negative as well as positive values, but finite and agreeing on the algebra  $\bigcup_{n=1}^{\infty} \mathcal{F}_n$ . Because of the uniqueness of the extension of a  $\sigma$ -additive measure to an algebra over the smallest  $\sigma$ -algebra containing it (Carathéodory's theorem, §3, Chapter II, equation (8) remains valid for sets  $A \in F_\infty = \sigma(\bigcup F_n)$ . Thus,

$$\int_A X_\infty dP = \int_A \xi dP = \int_A E(\xi | \mathcal{F}_\infty) dP, \quad A \in \mathcal{F}_\infty. \quad (9)$$

Since  $X_\infty$  and  $E(\xi | \mathcal{F}_\infty)$  are  $\mathcal{F}_\infty$ -measurable, it follows from Property I of Subsection 2, §6, Chapter II, and from (9), that  $X_\infty = E(\xi | \mathcal{F}_\infty)$  (P-a.s.).

This completes the proof of the theorem.

**Corollary.** A stochastic sequence  $X = (X_n, \mathcal{F}_n)$  is a uniformly integrable martingale if and only if there is a random variable  $\xi$  with  $E|\xi| < \infty$  such that  $X_n = E(\xi | \mathcal{F}_n)$  for all  $n \geq 1$ . Here  $X_n \rightarrow E(\xi | \mathcal{F}_\infty)$  (both P-a.s. and in the  $L^1$  sense) as  $n \rightarrow \infty$ .

In fact, if  $X = (X_n, \mathcal{F}_n)$  is a uniformly integrable martingale, then by Theorem 2 there is an integrable random variable  $X_\infty$  such that  $X_n \rightarrow X_\infty$  (P-a.s. and in the  $L^1$  sense) and  $X_n = E(X_\infty | \mathcal{F}_n)$ . As the random variable  $\xi$  we may take the  $\mathcal{F}_\infty$ -measurable variable  $X_\infty$ .

The converse follows from Theorem 3.

#### 4. We now turn to some applications of these theorems.

**EXAMPLE 1.** The "zero or one" law. Let  $\xi_1, \xi_2, \dots$  be a sequence of independent random variables,  $\mathcal{F}_n^\xi = \sigma\{\omega: \xi_1, \dots, \xi_n\}$  and let  $\mathcal{X}$  be the  $\sigma$ -algebra of the "tail" events. By Theorem 3, we have  $E(I_A | \mathcal{F}_n^\xi) \rightarrow E(I_A | \mathcal{F}_\infty^\xi) = I_A$  (P-a.s.). But  $I_A$  and  $(\xi_1, \dots, \xi_n)$  are independent. Since  $E(I_A | \mathcal{F}_n^\xi) = E I_A$  and therefore  $I_A = E I_A$  (P-a.s.), we find that either  $P(A) = 0$  or  $P(A) = 1$ .

The next two examples illustrate possible applications of the preceding results to convergence theorems in analysis.

**EXAMPLE 2.** If  $f = f(x)$  satisfies a Lipschitz condition on  $[0, 1]$ , it is absolutely continuous and, as is shown in courses in analysis, there is a (Lebesgue) integrable function  $g = g(x)$  such that

$$f(x) - f(0) = \int_0^x g(y) dy. \quad (10)$$

(In this sense,  $g(x)$  is a "derivative" of  $f(x)$ .)

Let us show how this result can be deduced from Theorem 1.

Let  $\Omega = [0, 1]$ ,  $\mathcal{F} = \mathcal{B}([0, 1])$ , and let  $P$  denote Lebesgue measure. Put

$$\xi_n(x) = \sum_{k=1}^{2^n} \frac{k-1}{2^n} I \left\{ \frac{k-1}{2^n} \leq x < \frac{k}{2^n} \right\},$$

$\mathcal{F}_n = \sigma\{x: \xi_1, \dots, \xi_n\} = \sigma\{x: \xi_n\}$ , and

$$X_n = \frac{f(\xi_n + 2^{-n}) - f(\xi_n)}{2^{-n}}.$$

Since for a given  $\xi_n$  the random variable  $\xi_{n+1}$  takes only the values  $\xi_n$  and  $\xi_n + 2^{-(n+1)}$  with conditional probabilities equal to  $\frac{1}{2}$ , we have

$$\begin{aligned} E[X_{n+1} | \mathcal{F}_n] &= E[X_{n+1} | \xi_n] = 2^{n+1} E[f(\xi_{n+1} + 2^{-(n+1)}) - f(\xi_{n+1}) | \xi_n] \\ &= 2^{n+1} \left\{ \frac{1}{2} [f(\xi_n + 2^{-(n+1)}) - f(\xi_n)] + \frac{1}{2} [f(\xi_n + 2^{-n}) - f(\xi_n + 2^{-(n+1)})] \right\} \\ &= 2^n \{f(\xi_n + 2^{-n}) - f(\xi_n)\} = X_n. \end{aligned}$$

It follows that  $X = (X_n, \mathcal{F}_n)$  is a martingale, and it is uniformly integrable since  $|X_n| \leq L$ , where  $L$  is the Lipschitz constant:  $|f(x) - f(y)| \leq L|x - y|$ . Observe that  $\mathcal{F} = \mathcal{B}([0, 1]) = \sigma(\bigcup \mathcal{F}_n)$ . Therefore, by the corollary to Theorem 3, there is an  $\mathcal{F}$ -measurable function  $g = g(x)$  such that  $X_n \rightarrow g$  (P-a.s.) and

$$X_n = E[g | \mathcal{F}_n]. \quad (11)$$

Consider the set  $B = [0, k/2^n]$ . Then by (11)

$$f\left(\frac{k}{2^n}\right) - f(0) = \int_0^{k/2^n} X_n dx = \int_0^{k/2^n} g(x) dx,$$

and since  $n$  and  $k$  are arbitrary, we obtain the required equation (10).

EXAMPLE 3. Let  $\Omega = [0, 1]$ ,  $\mathcal{F} = \mathcal{B}([0, 1])$  and let  $P$  denote Lebesgue measure. Consider the Haar system  $\{H_n(x)\}_{n \geq 1}$ , as defined in Example 3 of §11, Chapter II. Put  $\mathcal{F}_n = \sigma\{x: H_1, \dots, H_n\}$  and observe that  $\sigma(\bigcup \mathcal{F}_n) = \mathcal{F}$ . From the properties of conditional expectations and the structure of the Haar functions, it is easy to deduce that

$$E[f(x) | \mathcal{F}_n] = \sum_{k=1}^n a_k H_k(x) \quad (\text{P-a.s.}), \quad (12)$$

for every Borel function  $f \in L$ , where

$$a_k = (f, H_k) = \int_0^1 f(x) H_k(x) dx.$$

In other words, the conditional expectation  $E[f(x) | \mathcal{F}_n]$  is a partial sum of the Fourier series of  $f(x)$  in the Haar system. Then if we apply Theorem 3 to the martingale we find that, as  $n \rightarrow \infty$ ,

$$\sum_{k=1}^n (f, H_k) H_k(x) \rightarrow f(x) \quad (\text{P-a.s.})$$

and

$$\int_0^1 \left| \sum_{k=1}^n (f, H_k) H_k(x) - f(x) \right| dx \rightarrow 0.$$

EXAMPLE 4. Let  $(\xi_n)_{n \geq 1}$  be a sequence of random variables. By Theorem 2, §10, Chapter II, the P-a.e. convergence of the series  $\sum \xi_n$  implies its convergence in probability and in distribution. It turns out that if the random variables  $\xi_1, \xi_2, \dots$  are independent, the converse is also valid: the convergence in distribution of the series  $\sum \xi_n$  of independent random variables implies its convergence in probability and with probability one.

Let  $S_n = \xi_1 + \dots + \xi_n$ ,  $n \geq 1$  and  $S_n \xrightarrow{d} S$ . Then  $Ee^{itS_n} \rightarrow Ee^{itS}$  for every real number  $t$ . It is clear that there is a  $\delta > 0$  such that  $|Ee^{itS}| > 0$  for all  $|t| < \delta$ . Choose  $t_0$  so that  $|t_0| < \delta$ . Then there is an  $n_0 = n_0(t_0)$  such that  $|Ee^{it_0 S_n}| \geq c > 0$  for all  $n \geq n_0$ , where  $c$  is a constant.

For  $n \geq n_0$ , we form the sequence  $X = (X_n, \mathcal{F}_n)$  with

$$X_n = \frac{e^{it_0 S_n}}{Ee^{it_0 S_{n_0}}}, \quad \mathcal{F}_n = \sigma\{\xi_1, \dots, \xi_n\}.$$

Since  $\xi_1, \xi_2, \dots$  were assumed to be independent, the sequence  $X = (X_n, \mathcal{F}_n)$  is a martingale with

$$\sup_{n \geq n_0} E|X_n| \leq c^{-1} < \infty.$$

Then it follows from Theorem 1 that with probability one the limit  $\lim_n X_n$  exists and is finite. Therefore, the limit  $\lim_{n \rightarrow \infty} e^{itS_n}$  also exists with probability one. Consequently, we can assert that there is a  $\delta > 0$  such that for each  $t$  in the set  $T = \{t: |t| < \delta\}$  the limit  $\lim_n e^{itS_n}$  exists with probability one.

Let  $T \times \Omega = \{(t, \omega): t \in T, \omega \in \Omega\}$ , let  $\bar{\mathcal{B}}(T)$  be the  $\sigma$ -algebra of Lebesgue sets on  $T$  and let  $\lambda$  be Lebesgue measure on  $(T, \bar{\mathcal{B}}(T))$ . Also, let

$$C = \left\{ (t, \omega) \in T \times \Omega: \lim_n e^{itS_n(\omega)} \text{ exists} \right\}.$$

It is clear that  $C \in \bar{\mathcal{B}}(T) \otimes \mathcal{F}$ .

It was shown above that  $P(C_t) = 1$  for every  $t \in T$ , where  $C_t = \{\omega \in \Omega: (t, \omega) \in C\}$  is the section of  $C$  at the point  $t$ . By Fubini's theorem (Theorem 8, §6, Chapter II)

$$\begin{aligned} \int_{T \times \Omega} I_C(t, \omega) d(\lambda \times P) &= \int_T \left( \int_{\Omega} I_C(t, \omega) dP \right) d\lambda \\ &= \int_T P(C_t) d\lambda = \lambda(T) = 2\delta > 0. \end{aligned}$$

On the other hand, again by Fubini's theorem,

$$\lambda(T) = \int_{T \times \Omega} I_C(t, \omega) d(\lambda \times P) = \int_{\Omega} dP \left( \int_T I_C(t, \omega) d\lambda \right) = \int_{\Omega} \lambda(C_{\omega}) dP,$$

where  $C_{\omega} = \{t: (t, \omega) \in C\}$ .

Hence, it follows that there is a set  $\tilde{\Omega}$  with  $P(\tilde{\Omega}) = 1$  such that  $\lambda(C_{\omega}) = \lambda(T) = 2\delta > 0$  for all  $\omega \in \tilde{\Omega}$ .

Consequently, we may say that for every  $\omega \in \tilde{\Omega}$  the limit  $\lim_n e^{itS_n}$  exists for all  $t \in C_{\omega}$ . In addition, the measure of  $C_{\omega}$  is positive. From this and Problem 8, it follows that the limit  $\lim_n S_n(\omega)$  exists and is finite for  $\omega \in \tilde{\Omega}$ . Since  $P(\tilde{\Omega}) = 1$ , the limit  $\lim_n S_n(\omega)$  exists and is finite with probability one.

## 5. PROBLEMS

1. Let  $\{\mathcal{G}_n\}$  be a nonincreasing family of  $\sigma$ -algebras,  $\mathcal{G}_1 \supseteq \mathcal{G}_2 \supseteq \dots \mathcal{G}_{\infty} = \bigcap \mathcal{G}_n$ , and let  $\eta$  be an integrable random variable. Establish the following analog of Theorem 3: as  $n \rightarrow \infty$ ,

$$E(\eta | G_n) \rightarrow E(\eta | G_{\infty}) \quad (\text{P-a.s. and in the } L^1 \text{ sense}).$$

2. Let  $\xi_1, \xi_2, \dots$  be a sequence of independent identically distributed random variables with  $E|\xi_1| < \infty$  and  $E\xi_1 = m$ ; let  $S_n = \xi_1 + \dots + \xi_n$ . Having shown (see Problem 2, §7, Chapter II) that

$$E(\xi_1 | S_n, S_{n+1}, \dots) = E(\xi_1 | S_n) = \frac{S_n}{n} \quad (\text{P-a.s.}),$$

deduce from Problem 1 a stronger form of the law of large numbers: as  $n \rightarrow \infty$ ,

$$\frac{S_n}{n} \rightarrow m \quad (\text{P-a.s. and in the } L^1 \text{ sense}).$$

3. Establish the following result, which combines Lebesgue's dominated convergence theorem and P. Lévy's theorem. Let  $\{\xi_n\}_{n \geq 1}$  be a sequence of random variables such that  $\xi_n \rightarrow \xi$  (P-a.s.),  $|\xi_n| \leq \eta$ ,  $E\eta < \infty$  and  $\{\mathcal{F}_m\}_{m \geq 1}$  is a nondecreasing family of  $\sigma$ -algebras, with  $\mathcal{F}_\infty = \sigma(\bigcup \mathcal{F}_n)$ . Then

$$\lim_{\substack{m \rightarrow \infty \\ n \rightarrow \infty}} E(\xi_n | \mathcal{F}_m) = E(\xi | \mathcal{F}_\infty) \quad (\text{P-a.s.}).$$

4. Establish formula (12).  
 5. Let  $\Omega = [0, 1]$ ,  $\mathcal{F} = \mathcal{B}([0, 1])$ , let  $P$  denote Lebesgue measure, and let  $f = f(x) \in L^1$ . Put

$$f_n(x) = 2^n \int_{k2^{-n}}^{(k+1)2^{-n}} f(y) dy, \quad k2^{-n} \leq x < (k+1)2^{-n}.$$

Show that  $f_n(x) \rightarrow f(x)$  (P-a.s.).

6. Let  $\Omega = [0, 1]$ ,  $\mathcal{F} = \mathcal{B}([0, 1])$ , let  $P$  denote Lebesgue measure and let  $f = f(x) \in L^1$ . Continue this function periodically on  $[0, 2)$  and put

$$f_n(x) = \sum_{i=1}^{2^n} 2^{-n} f(x + i2^{-n}).$$

Show that  $f_n(x) \rightarrow f(x)$  (P-a.s.).

7. Prove that Theorem 1 remains valid for generalized submartingales  $X = (X_n, \mathcal{F}_n)$ , if  $\inf_m \sup_{n \geq m} E(X_n^+ | \mathcal{F}_m) < \infty$  (P-a.s.).  
 8. Let  $a_n$ ,  $n \geq 1$ , be a sequence of real numbers such that for all real numbers  $t$  with  $|t| < \delta$ ,  $\delta > 0$ , the limit  $\lim_n e^{ita_n}$  exists. Prove that then the limit  $\lim a_n$  exists and is finite.

## §5. Sets of Convergence of Submartingales and Martingales

1. Let  $X = (X_n, \mathcal{F}_n)$  be a stochastic sequence. Let us denote by  $\{X_n \rightarrow\}$ , or  $\{-\infty < \lim X_n < \infty\}$ , the set of sample points for which  $\lim X_n$  exists and is finite. Let us also write  $A \subseteq B$  (P-a.s.) if  $P(I_A \leq I_B) = 1$ .

If  $X$  is a submartingale and  $\sup E|X_n| < \infty$  (or, equivalently, if  $\sup EX_n^+ < \infty$ ), then according to Theorem 1 of §4 we have

$$\{X_n \rightarrow\} = \Omega \quad (\text{P-a.s.}).$$

Let us consider the structure of sets  $\{X_n \rightarrow\}$  of convergence for submartingales when the hypothesis  $\sup E|X_n| < \infty$  is not satisfied.

Let  $a > 0$ , and  $\tau_a = \inf\{n \geq 1: X_n > a\}$  with  $\tau_a = \infty$  if  $\{\cdot\} = \emptyset$ .

**Definition.** A stochastic sequence  $X = (X_n, \mathcal{F}_n)$  belongs to class  $C^+$  ( $X \in C^+$ ) if



$$E(\Delta X_{\tau_a})^+ I\{\tau_a < \infty\} < \infty \quad (1)$$

for every  $a > 0$ , where  $\Delta X_n = X_n - X_{n-1}$ ,  $X_0 = 0$ .

It is evident that  $X \in C^+$  if

$$E \sup_n |\Delta X_n| < \infty \quad (2)$$

or, all the more so, if

$$|\Delta X_n| \leq C < \infty \quad (\text{P-a.s.}), \quad (3)$$

for all  $n \geq 1$ .

**Theorem 1.** *If the submartingale  $X \in C^+$  then*

$$\{\sup X_n < \infty\} = \{X_n \rightarrow\} \quad (\text{P-a.s.}). \quad (4)$$

**PROOF.** The inclusion  $\{X_n \rightarrow\} \subseteq \{\sup X_n < \infty\}$  is evident. To establish the inclusion in the opposite direction, we consider the stopped submartingale  $X^{\tau_a} = (X_{\tau_a \wedge n}, \mathcal{F}_n)$ . Then, by (1),

$$\begin{aligned} \sup_n E X_{\tau_a \wedge n}^+ &\leq a + E[X_{\tau_a}^+ \cdot I\{\tau_a < \infty\}] \\ &\leq 2a + E[(\Delta X_{\tau_a})^+ \cdot I\{\tau_a < \infty\}] < \infty, \end{aligned} \quad (5)$$

and therefore by Theorem 1 of §4,

$$\{\tau_a = \infty\} \subseteq \{X_n \rightarrow\} \quad (\text{P-a.s.}).$$

But  $\bigcup_{a>0} \{\tau_a = \infty\} = \{\sup X_n < \infty\}$ ; hence  $\{\sup X_n < \infty\} \subseteq \{X_n \rightarrow\}$  (P-a.s.).

This completes the proof of the theorem.

**Corollary.** *Let  $X$  be a martingale with  $E \sup |\Delta X_n| < \infty$ . Then (P-a.s.)*

$$\{X_n \rightarrow\} \cup \{\underline{\lim} X_n = -\infty, \overline{\lim} X_n = +\infty\} = \Omega. \quad (6)$$

In fact, if we apply Theorem 1 to  $X$  and to  $-X$ , we find that (P-a.s.)

$$\begin{aligned} \{\overline{\lim} X_n < \infty\} &= \{\sup X_n < \infty\} = \{X_n \rightarrow\}, \\ \{\underline{\lim} X_n > -\infty\} &= \{\inf X_n > -\infty\} = \{X_n \rightarrow\}. \end{aligned}$$

Therefore (P-a.s.)

$$\{\overline{\lim} X_n < \infty\} \cup \{\underline{\lim} X_n > -\infty\} = \{X_n \rightarrow\},$$

which establishes (6).

Statement (6) means that, provided that  $E \sup |\Delta X_n| < \infty$ , either almost all trajectories of the martingale  $M$  have finite limits, or all behave very badly, in the sense that  $\overline{\lim} X_n = +\infty$  and  $\underline{\lim} X_n = -\infty$ .

2. If  $\xi_1, \xi_2, \dots$  is a sequence of independent random variables with  $E\xi_i = 0$  and  $|\xi_i| \leq c < \infty$ , then by Theorem 1 of §2, Chapter IV, the series  $\sum \xi_i$

converges (P-a.s.) if and only if  $\sum E\xi_i^2 < \infty$ . The sequence  $X = (X_n, \mathcal{F}_n)$  with  $X_n = \xi_1 + \dots + \xi_n$  and  $\mathcal{F}_n = \sigma\{\omega: \xi_1, \dots, \xi_n\}$  is a square-integrable martingale with  $\langle X \rangle_n = \sum_{i=1}^n E\xi_i^2$ , and the proposition just stated can be interpreted as follows:

$$\{\langle X \rangle_\infty < \infty\} = \{X_n \rightarrow\} = \Omega \quad (\text{P-a.s.}),$$

where  $\langle X \rangle_\infty = \lim_n \langle X \rangle_n$ .

The following proposition generalizes this result to more general martingales and submartingales.

**Theorem 2.** Let  $X = (X_n, \mathcal{F}_n)$  be a submartingale and

$$X_n = m_n + A_n$$

its Doob decomposition.

(a) If  $X$  is a nonnegative submartingale, then (P-a.s.)

$$\{A_\infty < \infty\} \subseteq \{X_n \rightarrow\} \subseteq \{\sup X_n < \infty\}. \quad (7)$$

(b) If  $X \in C^+$  then (P-a.s.)

$$\{X_n \rightarrow\} = \{\sup X_n < \infty\} \subseteq \{A_\infty < \infty\}. \quad (8)$$

(c) If  $X$  is a nonnegative submartingale and  $X \in C^+$ , then (P-a.s.)

$$\{X_n \rightarrow\} = \{\sup X_n < \infty\} = \{A_\infty < \infty\}. \quad (9)$$

**PROOF.** (a) The second inclusion in (7) is obvious. To establish the first inclusion we introduce the times

$$\sigma_a = \inf\{n \geq 1: A_{n+1} > a\}, \quad a > 0,$$

taking  $\sigma_a = +\infty$  if  $\{\cdot\} = \emptyset$ . Then  $A_{\sigma_a} \leq a$  and by Corollary 1 to Theorem 1 of §2, we have

$$EX_{n \wedge \sigma_a} = EA_{n \wedge \sigma_a} \leq a.$$

Let  $Y_n^a = X_{n \wedge \sigma_a}$ . Then  $Y^a = (Y_n^a, \mathcal{F}_n)$  is a submartingale with  $\sup EY_n^a \leq a < \infty$ . Since the martingale is nonnegative, it follows from Theorem 1, §4, that (P-a.s.)

$$\{A_\infty \leq a\} = \{\sigma_a = \infty\} \subseteq \{X_n \rightarrow\}.$$

Therefore (P-a.s.)

$$\{A_\infty < \infty\} = \bigcup_{a>0} \{A_\infty \leq a\} \subseteq \{X_n \rightarrow\}.$$

(b) The first equation follows from Theorem 1. To prove the second, we notice that, in accordance with (5),

$$EA_{\tau_a \wedge n} = EX_{\tau_a \wedge n}^+ \leq EX_{\tau_a \wedge n}^+ \leq 2a + E[(\Delta X_{\tau_a})^+ I\{\tau_a < \infty\}]$$

and therefore

$$EA_{\tau_a} = E \lim_n A_{\tau_a \wedge n} < \infty.$$

Hence  $\{\tau_a = \infty\} \subseteq \{A_\infty < \infty\}$  and we obtain the required conclusion since  $\bigcup_{a>0} \{\tau_a = \infty\} = \{\sup X_n < \infty\}$ .

(c) This is an immediate consequence of (a) and (b).

This completes the proof of the theorem.

**Remark.** The hypothesis that  $X$  is nonnegative can be replaced by the hypothesis  $\sup_n \mathbb{E} X_n^- < \infty$ .

**Corollary 1.** Let  $X_n = \xi_1 + \cdots + \xi_n$ , where  $\xi_i \geq 0$ ,  $\mathbb{E} \xi_i < \infty$ ,  $\xi_i$  are  $\mathcal{F}_i$ -measurable, and  $\mathcal{F}_0 = \{\emptyset, \Omega\}$ . Then (P-a.s.)

$$\left\{ \sum_{n=1}^{\infty} \mathbb{E}(\xi_n | \mathcal{F}_{n-1}) < \infty \right\} \subseteq \{X_n \rightarrow\}, \quad (10)$$

and if, in addition,  $\mathbb{E} \sup_n \xi_n < \infty$  then (P-a.s.)

$$\left\{ \sum_{n=1}^{\infty} \mathbb{E}(\xi_n | \mathcal{F}_{n-1}) < \infty \right\} = \{X_n \rightarrow\}. \quad (11)$$

**Corollary 2 (Borel–Cantelli–Lévy Lemma).** If the events  $B_n \in \mathcal{F}_n$ , then if we put  $\xi_n = I_{B_n}$  in (11), we find that

$$\left\{ \sum_{n=1}^{\infty} \mathbb{P}(B_n | \mathcal{F}_{n-1}) < \infty \right\} = \left\{ \sum_{n=1}^{\infty} I_{B_n} < \infty \right\}. \quad (12)$$

**3. Theorem 3.** Let  $M = (M_n, \mathcal{F}_n)_{n \geq 1}$  be a square-integrable martingale. Then (P-a.s.)

$$\{\langle M \rangle_\infty < \infty\} \subseteq \{M_n \rightarrow\}. \quad (13)$$

If also  $\mathbb{E} \sup |\Delta M_n|^2 < \infty$ , then (P-a.s.)

$$\{\langle M \rangle_\infty < \infty\} = \{M_n \rightarrow\}, \quad (14)$$

where

$$\langle M \rangle_\infty = \sum_{n=1}^{\infty} \mathbb{E}((\Delta M_n)^2 | \mathcal{F}_{n-1}) \quad (15)$$

with  $M_0 = 0$ ,  $\mathcal{F}_0 = \{\emptyset, \Omega\}$ .

**PROOF.** Consider the two submartingales  $M^2 = (M_n^2, \mathcal{F}_n)$  and  $(M + 1)^2 = ((M + 1)^2, \mathcal{F}_n)$ . Let their Doob decompositions be

$$M_n^2 = m'_n + A'_n, \quad (M_n + 1)^2 = m''_n + A''_n.$$

Then  $A'_n$  and  $A''_n$  are the same, since

$$A'_n = \sum_{k=1}^n \mathbb{E}(\Delta M_k^2 | \mathcal{F}_{k-1}) = \sum_{k=1}^n \mathbb{E}((\Delta M_k)^2 | \mathcal{F}_{k-1})$$

and

$$\begin{aligned} A_n'' &= \sum_{k=1}^n \mathbf{E}(\Delta(M_k + 1)^2 | \mathcal{F}_{k-1}) = \sum_{k=1}^n \mathbf{E}(\Delta M_k^2 | \mathcal{F}_{k-1}) \\ &= \sum_{k=1}^n \mathbf{E}((\Delta M_k)^2 | \mathcal{F}_{k-1}). \end{aligned}$$

Hence (7) implies that (P-a.s.)

$$\{\langle M \rangle_\infty < \infty\} = \{A'_\infty < \infty\} \subseteq \{M_n^2 \rightarrow\} \cap \{(M_n + 1)^2 \rightarrow\} = \{M_n \rightarrow\}.$$

Because of (9), equation (14) will be established if we show that the condition  $\mathbf{E} \sup |\Delta M_n|^2 < \infty$  guarantees that  $M^2$  belongs to  $\mathcal{C}^+$ .

Let  $\tau_a = \inf\{n \geq 1: M_n^2 > a\}$ ,  $a > 0$ . Then, on the set  $\{\tau_a < \infty\}$ ,

$$\begin{aligned} |\Delta M_{\tau_a}^2| &= |M_{\tau_a}^2 - M_{\tau_a-1}^2| \leq |M_{\tau_a} - M_{\tau_a-1}|^2 \\ &\quad + 2|M_{\tau_a-1}| \cdot |M_{\tau_a} - M_{\tau_a-1}| \leq (\Delta M_{\tau_a})^2 + 2a^{1/2} |\Delta M_{\tau_a}|, \end{aligned}$$

whence

$$\begin{aligned} \mathbf{E} |\Delta M_{\tau_a}^2| I\{\tau_a < \infty\} &\leq \mathbf{E} (\Delta M_{\tau_a})^2 I\{\tau_a < \infty\} + 2a^{1/2} \sqrt{\mathbf{E} (\Delta M_{\tau_a})^2 I\{\tau_a < \infty\}} \\ &\leq \mathbf{E} \sup |\Delta M_n|^2 + 2a^{1/2} \sqrt{\mathbf{E} \sup |\Delta M_n|^2} < \infty. \end{aligned}$$

This completes the proof of the theorem.

As an illustration of this theorem, we present the following result, which can be considered as a distinctive version of the strong law of large numbers for square-integrable martingales (compare Theorem 2 of §3, Chapter IV and Corollary 2 of Subsection 3 of §3).

**Theorem 4.** *Let  $M = (M_n, \mathcal{F}_n)$  be a square-integrable martingale and let  $A = (A_n, \mathcal{F}_{n-1})$  be a predictable increasing sequence with  $A_1 \geq 1$ ,  $A_\infty = \infty$  (P-a.s.).*

*If (P-a.s.)*

$$\sum_{i=1}^{\infty} \frac{\mathbf{E}[(\Delta M_i)^2 | \mathcal{F}_{i-1}]}{A_i^2} < \infty, \quad (16)$$

*then*

$$M_n/A_n \rightarrow 0, \quad n \rightarrow \infty, \quad (17)$$

*with probability 1.*

*In particular, if  $\langle M \rangle = (M_n, \mathcal{F}_n)$  is the quadratic characteristic of the square-integrable martingale,  $M = (M_n, \mathcal{F}_n)$  and  $\langle M \rangle_\infty = \infty$  (P-a.s.), then with probability 1*

$$\frac{M_n}{\langle M \rangle_n} \rightarrow 0, \quad n \rightarrow \infty. \quad (18)$$

PROOF. Consider the square-integrable martingale  $m = (m_n, \mathcal{F}_n)$  with

$$m_n = \sum_{i=1}^n \frac{\Delta M_i}{A_i}.$$

Then

$$\langle m \rangle_n = \sum_{i=1}^n \frac{\mathbb{E}[(\Delta M_i)^2 | \mathcal{F}_{i-1}]}{A_i^2}. \quad (19)$$

Since

$$\frac{M_n}{A_n} = \frac{\sum_{k=1}^n A_k \Delta m_k}{A_n},$$

we have, by Kronecker's lemma (§3, Chapter IV),  $M_n/A_n \rightarrow 0$  (P-a.s.) if the limit  $\lim_n m_n$  exists (finite) with probability 1. By (13),

$$\{\langle m \rangle_\infty < \infty\} \subseteq \{m_n \rightarrow\}. \quad (20)$$

Therefore it follows from (19) that (16) is a sufficient condition for (17).

If now  $A_n = \langle M \rangle_n$ , then (16) is automatically satisfied (see Problem 6) and consequently we have

$$\frac{M_n}{\langle M \rangle_n} \rightarrow 0 \quad (\text{P-a.s.}).$$

This completes the proof of the theorem.

EXAMPLE. Consider a sequence  $\xi_1, \xi_2, \dots$  of independent random variables with  $\mathbb{E}\xi_i = 0$ ,  $\forall \xi_i = V_i > 0$ , and let the sequence  $X = \{X_n\}_{n \geq 0}$  be defined recursively by

$$X_{n+1} = \theta X_n + \xi_{n+1}, \quad (21)$$

where  $X_0$  is independent of  $\xi_1, \xi_2, \dots$  and  $\theta$  is an unknown parameter,  $-\infty < \theta < \infty$ .

We interpret  $X_n$  as the result of an observation made at time  $n$  and ask for an estimator of the unknown parameter  $\theta$ . As an estimator of  $\theta$  in terms of  $X_0, X_1, \dots, X_n$ , we take

$$\hat{\theta}_n = \frac{\sum_{k=0}^{n-1} \frac{X_k X_{k+1}}{V_{k+1}}}{\sum_{k=0}^{n-1} \frac{X_k^2}{V_{k+1}}}, \quad (22)$$

taking this to be 0 if the denominator is 0. (The number  $\hat{\theta}_n$  is the *least-squares estimator*.)

It is clear from (21) and (22) that

$$\hat{\theta} = \theta + \frac{M_n}{A_n},$$

where

$$M_n = \sum_{k=0}^{n-1} \frac{X_k \xi_{k+1}}{V_{k+1}}, \quad A_n = \langle M \rangle_n = \sum_{k=0}^{n-1} \frac{X_k^2}{V_{k+1}}.$$

Therefore if the true value of the unknown parameter is  $\theta$ , then

$$P(\hat{\theta}_n \rightarrow \theta) = 1, \quad (23)$$

when (P-a.s.)

$$\frac{M_n}{A_n} \rightarrow 0, \quad n \rightarrow \infty. \quad (24)$$

(An estimator  $\theta_n$  with property (23) is said to be *strongly consistent*; compare the notion of consistency in §7, Chapter I.) Let us show that the conditions

$$\sup_n \frac{V_{n+1}}{V_n} < \infty, \quad \sum_{n=1}^{\infty} E\left(\frac{\xi_n^2}{V_n} \wedge 1\right) = \infty \quad (25)$$

are sufficient for (24), and therefore sufficient for (23).

We have

$$\begin{aligned} \sum_{n=1}^{\infty} \left(\frac{\xi_n^2}{V_n} \wedge 1\right) &\leq \sum_{n=1}^{\infty} \frac{\xi_n^2}{V_n} = \sum_{n=1}^{\infty} \frac{(X_n - \theta X_{n-1})^2}{V_n} \\ &\leq 2 \left[ \sum_{n=1}^{\infty} \frac{X_n^2}{V_n} + \theta^2 \sum_{n=1}^{\infty} \frac{X_{n-1}^2}{V_n} \right] \leq 2 \left[ \sup \frac{V_{n+1}}{V_n} + \theta^2 \right] \langle M \rangle_{\infty}. \end{aligned}$$

Therefore

$$\left\{ \sum_{n=1}^{\infty} \left(\frac{\xi_n^2}{V_n} \wedge 1\right) = \infty \right\} \subseteq \{\langle M \rangle_{\infty} = \infty\}.$$

By the three-series theorem (Theorem 3 of §2, Chapter IV) the divergence of  $\sum_{n=1}^{\infty} E((\xi_n^2/V_n) \wedge 1)$  guarantees the divergence (P-a.s.) of  $\sum_{n=1}^{\infty} ((\xi_n^2/V_n) \wedge 1)$ . Therefore  $P\{\langle M \rangle_{\infty} = \infty\} = 1$ . Moreover, if

$$m_n = \sum_{i=1}^n \frac{\Delta M_i}{\langle M \rangle_i},$$

then

$$\langle m \rangle_n = \sum_{i=1}^n \frac{\Delta \langle M \rangle_i}{\langle M \rangle_i^2}$$

and (see Problem 6)  $P(m_{\infty} < \infty) = 1$ . Hence (24) follows directly from Theorem 4.

(In Subsection 5 of the next section we continue the discussion of this example for *Gaussian* variables  $\xi_1, \xi_2, \dots$ )

**Theorem 5.** Let  $X = (X_n, \mathcal{F}_n)$  be a submartingale, and let

$$X_n = m_n + A_n$$

be its Doob decomposition. If  $|\Delta X_n| \leq C$ , then (P-a.s.)

$$\{\langle m \rangle_\infty + A_\infty < \infty\} = \{X_n \rightarrow\}, \quad (26)$$

or equivalently,

$$\left\{ \sum_{n=1}^{\infty} \mathbf{E}[\Delta X_n + (\Delta X_n)^2 | \mathcal{F}_{n-1}] < \infty \right\} = \{X_n \rightarrow\}. \quad (27)$$

**PROOF.** Since

$$A_n = \sum_{k=1}^n \mathbf{E}(\Delta X_k | \mathcal{F}_{k-1}), \quad (28)$$

and

$$m_n = \sum_{k=1}^n [\Delta X_k - \mathbf{E}(\Delta X_k | \mathcal{F}_{k-1})], \quad (29)$$

it follows from the assumption that  $|\Delta X_k| \leq C$  that the martingale  $m = (m_n, \mathcal{F}_n)$  is square-integrable with  $|\Delta m_n| \leq 2C$ . Then by (13)

$$\{\langle m \rangle_\infty + A_\infty < \infty\} \subseteq \{X_n \rightarrow\} \quad (30)$$

and according to (8)

$$\{X_n \rightarrow\} \subseteq \{A_\infty < \infty\}.$$

Therefore, by (14) and (20),

$$\begin{aligned} \{X_n \rightarrow\} &= \{X_n \rightarrow\} \cap \{A_\infty < \infty\} = \{X_n \rightarrow\} \cap \{A_\infty < \infty\} \cap \{m_n \rightarrow\} \\ &= \{X_n \rightarrow\} \cap \{A_\infty < \infty\} \cap \{\langle m \rangle_\infty < \infty\} \\ &= \{X_n \rightarrow\} \cap \{A_\infty + \langle m \rangle_\infty < \infty\} = \{A_\infty + \langle m \rangle_\infty < \infty\}. \end{aligned}$$

Finally, the equivalence of (26) and (27) follows because, by (29),

$$\langle m \rangle_n = \sum \{\mathbf{E}[(\Delta X_k)^2 | \mathcal{F}_{k-1}] - [\mathbf{E}(\Delta X_k | \mathcal{F}_{k-1})]^2\},$$

and the convergence of the series  $\sum_{k=1}^{\infty} \mathbf{E}(\Delta X_k | \mathcal{F}_{k-1})$  of nonnegative terms implies the convergence of  $\sum_{k=1}^{\infty} [\mathbf{E}(\Delta X_k | \mathcal{F}_{k-1})]^2$ . This completes the proof.

4. Kolmogorov's three-series theorem (Theorem 3 of §2, Chapter IV) gives a necessary and sufficient condition for the convergence, with probability 1, of a series  $\sum \xi_n$  of independent random variables. The following theorems, whose proofs are based on Theorems 2 and 3, describe sets of convergence of  $\sum \xi_n$  without the assumption that the random variables  $\xi_1, \xi_2, \dots$  are independent.

**Theorem 6.** Let  $\xi = (\xi_n, \mathcal{F}_n)$ ,  $n \geq 1$ , be a stochastic sequence, let  $\mathcal{F}_0 = \{\emptyset, \Omega\}$ , and let  $c$  be a positive constant. Then the series  $\sum \xi_n$  converges on the set  $A$  of sample points for which the three series

$$\sum \mathbf{P}(|\xi_n| \geq c | \mathcal{F}_{n-1}), \quad \sum \mathbf{E}(\xi_n^c | \mathcal{F}_{n-1}), \quad \sum \mathbf{V}(\xi_n^c | \mathcal{F}_{n-1})$$

converge, where  $\xi_n^c = \xi_n I(|\xi_n| \leq c)$ .

**PROOF.** Let  $X_n = \sum_{k=1}^n \xi_k$ . Since the series  $\sum \mathbf{P}(|\xi_n| \geq c | \mathcal{F}_{n-1})$  converges, by Corollary 2 of Theorem 2, and by the convergence of the series  $\sum \mathbf{E}(\xi_n^c | \mathcal{F}_{n-1})$ , we have

$$\begin{aligned} A \cap \{X_n \rightarrow\} &= A \cap \left\{ \sum_{k=1}^n \xi_k I(|\xi_k| \leq c) \rightarrow \right\} \\ &= A \cap \left\{ \sum_{k=1}^n [\xi_k I(|\xi_k| \leq c) - \mathbf{E}(\xi_k I(|\xi_k| \leq c) | \mathcal{F}_{k-1})] \rightarrow \right\}. \end{aligned} \quad (31)$$

Let  $\eta_k = \xi_k I(|\xi_k| \leq c) - \mathbf{E}(\xi_k I(|\xi_k| \leq c) | \mathcal{F}_{k-1})$  and let  $Y_n = \sum_{k=1}^n \eta_k$ . Then  $Y = (Y_n, \mathcal{F}_n)$  is a square-integrable martingale with  $|\eta_k| \leq 2c$ . By Theorem 3, we have

$$A \subseteq \{\sum \mathbf{V}(\xi_n^c | \mathcal{F}_{n-1}) < \infty\} = \{\langle Y \rangle_\infty < \infty\} = \{Y_n \rightarrow\}. \quad (32)$$

Then it follows from (31) that

$$A \cap \{X_n \rightarrow\} = A,$$

and therefore  $A \subseteq \{X_n \rightarrow\}$ . This completes the proof.

## 5. PROBLEMS

1. Show that if a submartingale  $X = (X_n, \mathcal{F}_n)$  satisfies  $\mathbf{E} \sup_n |X_n| < \infty$ , then it belongs to class  $C^+$ .
2. Show that Theorems 1 and 2 remain valid for generalized submartingales.
3. Show that generalized submartingales satisfy ( $\mathbf{P}$ -a.s.) the inclusion

$$\left\{ \inf_m \sup_{n \geq m} \mathbf{E}(X_n^+ | \mathcal{F}_m) < \infty \right\} \subseteq \{X_n \rightarrow\}.$$

4. Show that the corollary to Theorem 1 remains valid for generalized martingales.
5. Show that every generalized submartingale of class  $C^+$  is a local submartingale.
6. Let  $a_n > 0$ ,  $n \geq 1$ , and let  $b_n = \sum_{k=1}^n a_k$ . Show that

$$\sum_{n=1}^{\infty} \frac{a_n}{b_n^2} < \infty.$$



## §6. Absolute Continuity and Singularity of Probability Distributions

1. Let  $(\Omega, \mathcal{F})$  be a measurable space on which there is defined a family  $(\mathcal{F}_n)_{n \geq 1}$  of  $\sigma$ -algebras such that  $\mathcal{F}_1 \subseteq \mathcal{F}_2 \subseteq \dots \subseteq \mathcal{F}$  and

$$\mathcal{F} = \sigma\left(\bigcup_{n=1}^{\infty} \mathcal{F}_n\right). \quad (1)$$

Let us suppose that two probability measures  $P$  and  $\tilde{P}$  are given on  $(\Omega, \mathcal{F})$ . Let us write

$$P_n = P|_{\mathcal{F}_n}, \quad \tilde{P}_n = \tilde{P}|_{\mathcal{F}_n}$$

for the restrictions of these measures to  $\mathcal{F}_n$ , i.e., let  $P_n$  and  $\tilde{P}_n$  be measures on  $(\Omega, \mathcal{F}_n)$  and for  $B \in \mathcal{F}_n$  let

$$P_n(B) = P(B), \quad \tilde{P}_n(B) = \tilde{P}(B).$$

**Definition 1.** The probability measure  $\tilde{P}$  is *absolutely continuous* with respect to  $P$  (notation,  $\tilde{P} \ll P$ ) if  $\tilde{P}(A) = 0$  whenever  $P(A) = 0$ ,  $A \in \mathcal{F}$ .

When  $\tilde{P} \ll P$  and  $P \ll \tilde{P}$  the measures  $\tilde{P}$  and  $P$  are *equivalent* (notation,  $\tilde{P} \sim P$ ).

The measures  $\tilde{P}$  and  $P$  are *singular* (or *orthogonal*) if there is a set  $A \in \mathcal{F}$  such that  $\tilde{P}(A) = 1$  and  $P(\bar{A}) = 1$  (notation,  $\tilde{P} \perp P$ ).

**Definition 2.** We say that  $\tilde{P}$  is *locally absolutely continuous* with respect to  $P$  (notation,  $\tilde{P} \ll^{\text{loc}} P$ ) if

$$\tilde{P}_n \ll P_n \quad (2)$$

for every  $n \geq 1$ .

The fundamental question that we shall consider in this section is the determination of conditions under which local absolute continuity  $\tilde{P} \ll^{\text{loc}} P$  implies one of the properties  $\tilde{P} \ll P$ ,  $\tilde{P} \sim P$ ,  $\tilde{P} \perp P$ . It will become clear that martingale theory is the mathematical apparatus that lets us give definitive answers to these questions.

Let us then suppose that  $\tilde{P} \ll^{\text{loc}} P$ . We write

$$z_n = \frac{d\tilde{P}_n}{dP_n}$$

the Radon–Nikodým derivative of  $\tilde{P}_n$  with respect to  $P_n$ . It is clear that  $z_n$  is  $\mathcal{F}_n$ -measurable; and if  $A \in \mathcal{F}_n$  then

$$\begin{aligned} \int_A z_{n+1} dP &= \int_A \frac{d\tilde{P}_{n+1}}{dP_{n+1}} dP = \tilde{P}_{n+1}(A) = \tilde{P}_n(A) \\ &= \int_A \frac{d\tilde{P}_n}{dP_n} dP = \int_A z_n dP. \end{aligned}$$

It follows that, with respect to  $P$ , the stochastic sequence  $Z = (z_n, \mathcal{F}_n)_{n \geq 1}$  is a martingale.

Write

$$z_\infty = \overline{\lim} z_n.$$

Since  $Ez_n = 1$ , it follows from Theorem 1, §4, that  $\lim z_n$  exists  $P$ -a.s. and therefore  $P(z_\infty = \lim z_n) = 1$ . (In the course of the proof of Theorem 1 it will be established that  $\lim z_n$  exists also for  $\tilde{P}$ , so that  $\tilde{P}(z_\infty = \lim z_n) = 1$ .)

The key to problems on absolute continuity and singularity is *Lebesgue's decomposition*.

**Theorem 1.** Let  $\tilde{P} \ll^{\text{loc}} P$ . Then for every  $A \in \mathcal{F}$ ,

$$\tilde{P}(A) = \int_A z_\infty dP + \tilde{P}\{A \cap (z_\infty = \infty)\}, \quad (3)$$

and the measures  $\mu(A) = \tilde{P}\{A \cap (z_\infty = \infty)\}$  and  $P(A)$ ,  $A \in \mathcal{F}$ , are singular.

PROOF. Let us notice first that the classical *Lebesgue decomposition* shows that if  $P$  and  $\tilde{P}$  are two measures, there are unique measures  $\lambda$  and  $\mu$  such that  $\tilde{P} = \lambda + \mu$ , where  $\lambda \ll P$  and  $\mu \perp P$ . Conclusion (3) can be thought of as a specialization of this decomposition under the assumption that  $\tilde{P}_n \ll P_n$ ,  $n \geq 1$ .

Let us introduce the probability measures

$$Q = \frac{1}{2}(P + \tilde{P}), \quad Q_n = \frac{1}{2}(P_n + \tilde{P}_n), \quad n \geq 1,$$

and the notation

$$\tilde{z} = \frac{d\tilde{P}}{dQ}, \quad z = \frac{dP}{dQ}, \quad \tilde{z}_n = \frac{d\tilde{P}_n}{dQ_n}, \quad z_n = \frac{dP_n}{dQ_n}.$$

Since  $\tilde{P}(\tilde{z} = 0) = P(z = 0) = 0$ , we have  $Q(\tilde{z} = 0, z = 0) = 0$ . Consequently the product  $\tilde{z} \cdot z^{-1}$  can be defined consistently on the set  $\Omega \setminus \{\tilde{z} = 0, z = 0\}$ ; we define it to be zero on the set  $\{\tilde{z} = 0, z = 0\}$ .

Since  $\tilde{P}_n \ll P_n \ll Q_n$ , we have (see (II.7.36))

$$\frac{d\tilde{P}_n}{dQ_n} = \frac{d\tilde{P}_n}{dP_n} \cdot \frac{dP_n}{dQ_n} \quad (Q\text{-a.s.}) \quad (4)$$

i.e.,

$$\tilde{z}_n = z_n \tilde{z}_n \quad (Q\text{-a.s.}) \quad (5)$$

whence

$$z_n = \tilde{z}_n \cdot \tilde{z}_n^{-1} \quad (Q\text{-a.s.})$$

where, as before, we take  $\tilde{z}_n \cdot \tilde{z}_n^{-1} = 0$  on the set  $\{\tilde{z}_n = 0, z_n = 0\}$ , which is of  $Q$ -measure zero.

Each of the sequences  $(\tilde{z}_n, \mathcal{F}_n)$  and  $(z_n, \mathcal{F}_n)$  is (with respect to  $Q$ ) a uniformly integrable martingale and consequently the limits  $\lim \tilde{z}_n$  and  $\lim z_n$  exist. Moreover (Q-a.s.)

$$\lim \tilde{z}_n = \tilde{z}, \quad \lim z_n = z. \quad (6)$$

From this and the equations  $z_n = \tilde{z}_n z_n^{-1}$  (Q-a.s.) and  $Q(\tilde{z} = 0, z = 0) = 0$ , it follows that (Q-a.s.) the limit  $\lim z_n = z_\infty$  exists and is equal to  $\tilde{z} \cdot z^{-1}$ .

It is clear that  $P \ll Q$  and  $\tilde{P} \ll Q$ . Therefore  $\lim z_n$  exists both with respect to  $P$  and with respect to  $\tilde{P}$ .

Now let

$$\lambda(A) = \int_A z_\infty dP, \quad \mu(A) = \tilde{P}\{A \cap (z_\infty = \infty)\}.$$

To establish (3), we must show that

$$\tilde{P}(A) = \lambda(A) + \mu(A), \quad \lambda \ll P, \quad \mu \perp P.$$

We have

$$\begin{aligned} \tilde{P}(A) &= \int_A \tilde{z} dQ = \int_A \tilde{z} \tilde{z} z dQ + \int_A \tilde{z} [1 - \tilde{z} z] dQ \\ &= \int_A \tilde{z} z dP + \int_A [1 - \tilde{z} z] d\tilde{P} = \int_A z_\infty dP + \tilde{P}\{A \cap (z = 0)\}, \quad (7) \end{aligned}$$

where the last equation follows from

$$P\left\{\tilde{z} = z^{-1}\right\} = 1, \quad \tilde{P}\{z_\infty = \tilde{z} \cdot z^{-1}\} = 1.$$

Furthermore,

$$\begin{aligned} \tilde{P}\{A \cap (z = 0)\} &= \tilde{P}\{A \cap (z = 0) \cap (\tilde{z} > 0)\} \\ &= \tilde{P}\{A \cap (\tilde{z} \cdot z^{-1} = \infty)\} = \tilde{P}\{A \cap (z_\infty = \infty)\}, \end{aligned}$$

which, together with (7), establishes (3).

It is clear from the construction of  $\lambda$  that  $\lambda \ll P$  and that  $P(z_\infty < \infty) = 1$ . But we also have

$$\mu(z_\infty < \infty) = \tilde{P}\{(z_\infty < \infty) \cap (z_\infty = \infty)\} = 0.$$

Consequently the theorem is proved.

The Lebesgue decomposition (3) implies the following useful tests for absolute continuity or singularity for locally absolutely continuous probability measures.

**Theorem 2.** Let  $\tilde{P} \ll^{\text{loc}} P$ , i.e.  $\tilde{P}_n \ll P_n$ ,  $n \geq 1$ . Then

$$\tilde{P} \ll P \Leftrightarrow E z_{\infty} = 1 \Leftrightarrow \tilde{P}(z_{\infty} < \infty) = 1, \quad (8)$$

$$\tilde{P} \perp P \Leftrightarrow E z_{\infty} = 0 \Leftrightarrow \tilde{P}(z_{\infty} = \infty) = 1, \quad (9)$$

where  $E$  denotes averaging with respect to  $P$ .

**PROOF.** Putting  $A = \Omega$  in (3), we find that

$$E z_{\infty} = 1 \Leftrightarrow \tilde{P}(z_{\infty} = \infty) = 0, \quad (10)$$

$$E z_{\infty} = 0 \Leftrightarrow \tilde{P}(z_{\infty} = \infty) = 1. \quad (11)$$

If  $\tilde{P}(z_{\infty} = \infty) = 0$ , it again follows from (3) that  $\tilde{P} \ll P$ .

Conversely, let  $\tilde{P} \ll P$ . Then since  $P(z_{\infty} = \infty) = 0$ , we have  $\tilde{P}(z_{\infty} = \infty) = 0$ .

In addition, if  $\tilde{P} \perp P$  there is a set  $B \in \mathcal{F}$  with  $\tilde{P}(B) = 1$  and  $P(B) = 0$ . Then  $\tilde{P}(B \cap (z_{\infty} = \infty)) = 1$  by (3), and therefore  $\tilde{P}(z_{\infty} = \infty) = 1$ . If, on the other hand,  $\tilde{P}(z_{\infty} = \infty) = 1$  the property  $\tilde{P} \perp P$  is evident, since  $P(z_{\infty} = \infty) = 0$ .

This completes the proof of the theorem.

2. It is clear from Theorem 2 that the tests for absolute continuity or singularity can be expressed either in terms of  $P$  (verify the equation  $E z_{\infty} = 1$  or  $E z_{\infty} = 0$ ), or in terms of  $\tilde{P}$  (verify that  $\tilde{P}(z_{\infty} < \infty) = 1$  or that  $\tilde{P}(z_{\infty} = \infty) = 1$ ).

By Theorem 5 of §6, Chapter II, the condition  $E z_{\infty} = 1$  is equivalent to the uniform integrability (with respect to  $P$ ) of the family  $\{z_n\}_{n \geq 1}$ . This allows us to give simple *sufficient conditions for the absolute continuity*  $\tilde{P} \ll P$ . For example, if

$$\sup_n E[z_n \ln^+ z_n] < \infty \quad (12)$$

or if

$$\sup_n E z_n^{1+\varepsilon} < \infty, \quad \varepsilon > 0, \quad (13)$$

then, by Lemma 3 of §6, Chapter II, the family of random variables  $\{z_n\}_{n \geq 1}$  is uniformly integrable and therefore  $\tilde{P} \ll P$ .

In many cases it is preferable to verify the property of absolute continuity or of singularity by using a test in terms of  $\tilde{P}$ , since then the question is reduced to the investigation of the probability of the "tail" event  $\{z_{\infty} < \infty\}$ , where one can use propositions like the "zero-one" law.

Let us show, by way of illustration, that the "Kakutani dichotomy" can be deduced from Theorem 2.

Let  $(\Omega, \mathcal{F}, P)$  be a probability space, let  $(R^{\infty}, \mathcal{B}_{\infty})$  be a measurable space of sequences  $x = (x_1, x_2, \dots)$  of numbers with  $\mathcal{B}_{\infty} = \mathcal{B}(R^{\infty})$ , and let  $\mathcal{B}_n = \sigma\{x: \{x_1, \dots, x_n\}\}$ . Let  $\xi = (\xi_1, \xi_2, \dots)$  and  $\tilde{\xi} = (\tilde{\xi}_1, \tilde{\xi}_2, \dots)$  be sequences of independent random variables.

Let  $P$  and  $\tilde{P}$  be the probability distributions on  $(R^\infty, \mathcal{B}_\infty)$  for  $\xi$  and  $\tilde{\xi}$ , respectively, i.e.

$$P(B) = P\{\xi \in B\}, \quad \tilde{P}(B) = P\{\tilde{\xi} \in B\}, \quad B \in \mathcal{B}_\infty.$$

Also let

$$P_n = P|_{\mathcal{B}_n}, \quad \tilde{P}_n = \tilde{P}|_{\mathcal{B}_n}$$

be the restrictions of  $P$  and  $\tilde{P}$  to  $\mathcal{B}_n$  and let

$$P_{\xi_n}(A) = P(\xi_n \in A), \quad P_{\tilde{\xi}_n}(A) = P(\tilde{\xi}_n \in A), \quad A \in \mathcal{B}(R^1).$$

**Theorem 3 (Kakutani Dichotomy).** Let  $\xi = (\xi_1, \xi_2, \dots)$  and  $\tilde{\xi} = (\tilde{\xi}_1, \tilde{\xi}_2, \dots)$  be sequences of independent random variables for which

$$P_{\xi_n} \ll P_{\xi_n}, \quad n \geq 1. \quad (14)$$

Then either  $\tilde{P} \ll P$  or  $\tilde{P} \perp P$ .

**PROOF.** Condition (14) is evidently equivalent to  $\tilde{P}_n \ll P_n, n \geq 1$ , i.e.  $\tilde{P} \ll^{\text{loc}} P$ . It is clear that

$$z_n = \frac{d\tilde{P}_n}{dP_n} = q_1(x_1) \cdots q_n(x_n),$$

where

$$q_i(x_i) = \frac{dP_{\tilde{\xi}_i}}{dP_{\xi_i}}(x_i). \quad (15)$$

Consequently

$$\{x: z_\infty < \infty\} = \{x: \ln z_\infty < \infty\} = \left\{x: \sum_{i=1}^{\infty} \ln q_i(x_i) < \infty\right\}.$$

The event  $\{x: \sum_{i=1}^{\infty} \ln q_i(x_i) < \infty\}$  is a tail event. Therefore, by the Kolmogorov zero-one law (Theorem 1 of §1, Chapter IV) the probability  $\tilde{P}\{x: z_\infty < \infty\}$  has only two values (0 or 1), and therefore by Theorem 2 either  $\tilde{P} \perp P$  or  $\tilde{P} \ll P$ .

This completes the proof of the theorem.

**3.** The following theorem provides, in “predictable” terms, a test for absolute continuity or singularity.

**Theorem 4.** Let  $\tilde{P} \ll^{\text{loc}} P$  and let

$$\alpha_n = z_n z_{n-1}^\oplus, \quad n \geq 1,$$

with  $z_0 = 1$ . Then (with  $\mathcal{F}_0 = \{\emptyset, \Omega\}$ )

$$\tilde{P} \ll P \Leftrightarrow \tilde{P}\left\{\sum_{n=1}^{\infty} [1 - E(\sqrt{\alpha_n} | \mathcal{F}_{n-1})] < \infty\right\} = 1, \quad (16)$$

$$\tilde{P} \perp P \Leftrightarrow \tilde{P}\left\{\sum_{n=1}^{\infty} [1 - E(\sqrt{\alpha_n} | \mathcal{F}_{n-1})] = \infty\right\} = 1. \quad (17)$$

PROOF. Since

$$\tilde{\mathbf{P}}_n\{z_n = 0\} = \int_{\{z_n=0\}} z_n d\mathbf{P} = 0,$$

we have (P-a.s.)

$$z_n = \prod_{k=1}^n \alpha_k = \exp\left\{\sum_{k=1}^n \ln \alpha_k\right\}. \quad (18)$$

Putting  $A = \{z_\infty = 0\}$  in (3), we find that  $\tilde{\mathbf{P}}\{z_\infty = 0\} = 0$ . Therefore, by (18), we have ( $\tilde{\mathbf{P}}$ -a.s.)

$$\begin{aligned} \{z_\infty < \infty\} &= \{0 < z_\infty < \infty\} = \{0 < \lim z_n < \infty\} \\ &= \left\{-\infty < \lim_{k \rightarrow \infty} \sum_{k=1}^n \ln \alpha_k < \infty\right\}. \end{aligned} \quad (19)$$

Let us introduce the function

$$u(x) = \begin{cases} x, & |x| \leq 1, \\ \text{sign } x, & |x| > 1. \end{cases}$$

Then

$$\left\{-\infty < \lim_{k \rightarrow \infty} \sum_{k=1}^n \ln \alpha_k < \infty\right\} = \left\{-\infty < \lim_{k \rightarrow \infty} \sum_{k=1}^n u(\ln \alpha_k) < \infty\right\}. \quad (20)$$

Let  $\tilde{\mathbf{E}}$  denote averaging with respect to  $\tilde{\mathbf{P}}$  and let  $\eta$  be an  $\mathcal{F}_n$ -measurable integrable random variable. It follows from the properties of conditional expectations (Problem 4) that

$$z_{n-1} \tilde{\mathbf{E}}(\eta | \mathcal{F}_{n-1}) = \mathbf{E}(\eta z_n | \mathcal{F}_{n-1}) \quad (\mathbf{P}\text{- and } \tilde{\mathbf{P}}\text{-a.s.}), \quad (21)$$

$$\tilde{\mathbf{E}}(\eta | \mathcal{F}_{n-1}) = z_{n-1}^\oplus \mathbf{E}(\eta z_n | \mathcal{F}_{n-1}) \quad (\tilde{\mathbf{P}}\text{-a.s.}). \quad (22)$$

Recalling that  $\alpha_n = z_{n-1}^\oplus z_n$ , we obtain the following useful formula from (22):

$$\tilde{\mathbf{E}}(\eta | \mathcal{F}_{n-1}) = \mathbf{E}(\alpha_n \eta | \mathcal{F}_{n-1}) \quad (\tilde{\mathbf{P}}\text{-a.s.}). \quad (23)$$

From this it follows, in particular, that

$$\mathbf{E}(\alpha_n | \mathcal{F}_{n-1}) = 1 \quad (\tilde{\mathbf{P}}\text{-a.s.}). \quad (24)$$

By (23),

$$\tilde{\mathbf{E}}[u(\ln \alpha_n) | \mathcal{F}_{n-1}] = \mathbf{E}[\alpha_n u(\ln \alpha_n) | \mathcal{F}_{n-1}] \quad (\tilde{\mathbf{P}}\text{-a.s.}).$$

Since  $xu(\ln x) \geq x - 1$  for  $x \geq 0$ , we have, by (24),

$$\tilde{\mathbf{E}}[u(\ln \alpha_n) | \mathcal{F}_{n-1}] \geq 0 \quad (\tilde{\mathbf{P}}\text{-a.s.}).$$

It follows that the stochastic sequence  $X = (X_n, \mathcal{F}_n)$  with

$$X_n = \sum_{k=1}^n u(\ln \alpha_k)$$

is a submartingale with respect to  $\tilde{\mathbf{P}}$ ; and  $|\Delta X_n| = |u(\ln \alpha_n)| \leq 1$ .

Then, by Theorem 5 of §5, we have ( $\tilde{\mathbf{P}}$ -a.s.)

$$\left\{ -\infty < \lim \sum_{k=1}^n u(\ln \alpha_k) < \infty \right\} = \left\{ \sum_{k=1}^{\infty} \tilde{\mathbf{E}}[u(\ln \alpha_k) + u^2(\ln \alpha_k) | \mathcal{F}_{k-1}] < \infty \right\}. \quad (25)$$

Hence we find, by combining (19), (20), (22), and (25), that ( $\mathbf{P}$ -a.s.)

$$\begin{aligned} \{z_{\infty} < \infty\} &= \left\{ \sum_{k=1}^{\infty} \tilde{\mathbf{E}}[u(\ln \alpha_k) + u^2(\ln \alpha_k) | \mathcal{F}_{k-1}] < \infty \right\} \\ &= \left\{ \sum_{k=1}^{\infty} \mathbf{E}[\alpha_k u(\ln \alpha_k) + \alpha_k u^2(\ln \alpha_k) | \mathcal{F}_{k-1}] < \infty \right\} \end{aligned}$$

and consequently, by Theorem 2,

$$\tilde{\mathbf{P}} \ll \mathbf{P} \Leftrightarrow \tilde{\mathbf{P}} \left\{ \sum_{k=1}^{\infty} \mathbf{E}[\alpha_k u(\ln \alpha_k) + \alpha_k u^2(\ln \alpha_k) | \mathcal{F}_{k-1}] < \infty \right\} = 1, \quad (26)$$

$$\tilde{\mathbf{P}} \perp \mathbf{P} \Leftrightarrow \tilde{\mathbf{P}} \left\{ \sum_{k=1}^{\infty} \mathbf{E}[\alpha_k u(\ln \alpha_k) + \alpha_k u^2(\ln \alpha_k) | \mathcal{F}_{k-1}] = \infty \right\} = 1. \quad (27)$$

We now observe that by (24),

$$\mathbf{E}[(1 - \sqrt{\alpha_n})^2 | \mathcal{F}_{n-1}] = 2\mathbf{E}[1 - \sqrt{\alpha_n} | \mathcal{F}_{n-1}] \quad (\mathbf{P}\text{-a.s.})$$

and for  $x \geq 0$  there are constants  $A$  and  $B$  ( $0 < A < B < \infty$ ) such that

$$A(1 - \sqrt{x})^2 \leq xu(\ln x) + xu^2(\ln x) + 1 - x \leq B(1 - \sqrt{x})^2. \quad (28)$$

Hence (16) and (17) follow from (26), (27) and (24), (28).

This completes the proof of the theorem.

**Corollary 1.** *If, for all  $n \geq 1$ , the  $\sigma$ -algebras  $\sigma(\alpha_n)$  and  $\mathcal{F}_{n-1}$  are independent with respect to  $\mathbf{P}$  (or  $\tilde{\mathbf{P}}$ ), and  $\tilde{\mathbf{P}} \stackrel{\text{loc}}{\ll} \mathbf{P}$ , then we have the dichotomy: either  $\tilde{\mathbf{P}} \ll \mathbf{P}$  or  $\tilde{\mathbf{P}} \perp \mathbf{P}$ . Correspondingly,*

$$\tilde{\mathbf{P}} \ll \mathbf{P} \Leftrightarrow \sum_{n=1}^{\infty} [1 - \mathbf{E}\sqrt{\alpha_n}] < \infty,$$

$$\tilde{\mathbf{P}} \perp \mathbf{P} \Leftrightarrow \sum_{n=1}^{\infty} [1 - \mathbf{E}\sqrt{\alpha_n}] = \infty.$$

In particular, in the Kakutani situation (see Theorem 3)  $\alpha_n = q_n$  and

$$\tilde{P} \ll P \Leftrightarrow \sum_{n=1}^{\infty} [1 - E\sqrt{q_n(x_n)}] < \infty,$$

$$\tilde{P} \perp P \Leftrightarrow \sum_{n=1}^{\infty} [1 - E\sqrt{q_n(x_n)}] = \infty.$$

**Corollary 2.** Let  $\tilde{P} \stackrel{\text{loc}}{\ll} P$ . Then

$$\tilde{P} \left\{ \sum_{n=1}^{\infty} E(\alpha_n \ln \alpha_n | \mathcal{F}_{n-1}) < \infty \right\} = 1 \Rightarrow \tilde{P} \ll P.$$

For the proof, it is enough to notice that

$$x \ln x + \frac{3}{2}(1-x) \geq 1 - x^{1/2}, \quad (29)$$

for all  $x \geq 0$ , and apply (16) and (24).

**Corollary 3.** Since the series  $\sum_{n=1}^{\infty} [1 - E(\sqrt{\alpha_n} | \mathcal{F}_{n-1})]$ , which has nonnegative ( $\tilde{P}$ -a.s.) terms, converges or diverges with the series  $\sum |\ln E(\sqrt{\alpha_n} | \mathcal{F}_{n-1})|$ , conclusions (16) and (17) of Theorem 4 can be put in the form

$$\tilde{P} \ll P \Leftrightarrow \tilde{P} \left\{ \sum_{n=1}^{\infty} |\ln E(\sqrt{\alpha_n} | \mathcal{F}_{n-1})| < \infty \right\} = 1, \quad (30)$$

$$\tilde{P} \perp P \Leftrightarrow \tilde{P} \left\{ \sum_{n=1}^{\infty} |\ln E(\sqrt{\alpha_n} | \mathcal{F}_{n-1})| = \infty \right\} = 1. \quad (31)$$

**Corollary 4.** Let there exist constants  $A$  and  $B$  such that  $0 \leq A < 1, B \geq 0$  and

$$P\{1 - A \leq \alpha_n \leq 1 + B\} = 1, \quad n \geq 1.$$

Then if  $\tilde{P} \stackrel{\text{loc}}{\ll} P$  we have

$$\tilde{P} \ll P \Leftrightarrow \tilde{P} \left\{ \sum_{n=1}^{\infty} E[(1 - \alpha_n)^2 | \mathcal{F}_{n-1}] < \infty \right\} = 1,$$

$$\tilde{P} \perp P \Leftrightarrow \tilde{P} \left\{ \sum_{n=1}^{\infty} E[(1 - \alpha_n)^2 | \mathcal{F}_{n-1}] = \infty \right\} = 1.$$

For the proof it is enough to notice that when  $x \in [1 - A, 1 + B]$ , where  $0 \leq A < 1, B \geq 0$ , there are constants  $c$  and  $C$  ( $0 < c < C < \infty$ ) such that

$$c(1 - x)^2 \leq (1 - \sqrt{x})^2 \leq C(1 - x)^2. \quad (32)$$

4. With the notation of Subsection 2, let us suppose that  $\xi = (\xi_1, \xi_2, \dots)$  and  $\tilde{\xi} = (\tilde{\xi}_1, \tilde{\xi}_2, \dots)$  are Gaussian sequences,  $\tilde{P}_n \sim P_n$ ,  $n \geq 1$ . Let us show that, for such sequences, the “Hájek–Feldman dichotomy,” either  $\tilde{P} \sim P$  or  $\tilde{P} \perp P$ , follows from the “predictable” test given above.

By the theorem on normal correlation (Theorem 2 of §13, Chapter II) the conditional expectations  $M(x_n | \mathcal{B}_{n-1})$  and  $\tilde{M}(x_n | \mathcal{B}_{n-1})$ , where  $M$  and  $\tilde{M}$



are averages with respect to  $P$  and  $\tilde{P}$ , respectively, are linear functions of  $x_1, \dots, x_{n-1}$ . We denote these linear functions by  $a_{n-1}(x)$  and  $\tilde{a}_{n-1}(x)$  and put

$$b_{n-1} = (M[x_n - a_{n-1}(x)]^2)^{1/2},$$

$$\tilde{b}_{n-1} = (\tilde{M}[x_n - \tilde{a}_{n-1}(x)]^2)^{1/2}.$$

Again by the theorem on normal correlation, there are sequences  $\varepsilon = (\varepsilon_1, \varepsilon_2, \dots)$  and  $\tilde{\varepsilon} = (\tilde{\varepsilon}_1, \tilde{\varepsilon}_2, \dots)$  of independent Gaussian random variables with zero means and unit variances, such that

$$x_n = a_{n-1}(x) + b_{n-1}\varepsilon_n, \quad (\mathbf{P}\text{-a.s.}),$$

$$x_n = \tilde{a}_{n-1}(x) + \tilde{b}_{n-1}\tilde{\varepsilon}_n \quad (\tilde{\mathbf{P}}\text{-a.s.}). \quad (33)$$

Notice that if  $b_{n-1} = 0$ , or  $\tilde{b}_{n-1} = 0$ , it is generally necessary to extend the probability space in order to construct  $\varepsilon_n$  or  $\tilde{\varepsilon}_n$ . However, if  $b_{n-1} = 0$  the extended vector  $(x_1, \dots, x_n)$  will be contained ( $\mathbf{P}$ -a.s.) in the linear manifold  $x_n = a_{n-1}(x)$ , and since by hypothesis  $\tilde{\mathbf{P}}_n \sim \mathbf{P}_n$ , we have  $b_{n-1} = 0$ ,  $a_{n-1} = \tilde{a}_{n-1}(x)$ , and  $\alpha_n(x) = 1$  ( $\mathbf{P}$ - or  $\tilde{\mathbf{P}}$ -a.s.). Hence we may suppose without loss of generality that  $b_n^2 > 0$ ,  $\tilde{b}_n^2 > 0$  for all  $n \geq 1$ , since otherwise the contribution of the corresponding terms of the sum  $\sum_{n=1}^{\infty} [1 - M\sqrt{\alpha_n}|B_{n-1}|]$  (see (16) and (17)) is zero.

Using the Gaussian hypothesis, we find from (33) that, for  $n \geq 1$ ,

$$\alpha_n = d_{n-1}^{-1} \exp \left\{ -\frac{(x_n - a_{n-1}(x))^2}{2b_{n-1}^2} + \frac{(x_n - \tilde{a}_{n-1}(x))^2}{2\tilde{b}_{n-1}^2} \right\}, \quad (34)$$

where  $d_n = |\tilde{b}_n \cdot \tilde{b}_n^{-1}|$  and

$$a_0(x) = E\xi_1, \quad \tilde{a}_0(x) = E\tilde{\xi}_1,$$

$$b_0^2 = V\xi_1, \quad \tilde{b}_0^2 = V\tilde{\xi}_1.$$

From (34),

$$\ln M(\alpha_n^{1/2} | \mathcal{B}_{n-1}) = \frac{1}{2} \ln \frac{2d_{n-1}}{1 + d_{n-1}^2} - \frac{d_{n-1}^2}{1 + d_{n-1}^2} \left( \frac{a_{n-1}(x) - \tilde{a}_{n-1}(x)}{b_{n-1}} \right)^2.$$

Since  $\ln [2d_{n-1}/(1 + d_{n-1}^2)] \leq 0$ , statement (30) can be written in the form

$$\tilde{P} \ll P \Leftrightarrow \tilde{P} \left\{ \sum_{n=1}^{\infty} \left[ \frac{1}{2} \ln \frac{1 + d_{n-1}^2}{2d_{n-1}} + \frac{d_{n-1}^2}{1 + d_{n-1}^2} \cdot \left( \frac{a_{n-1}(x) - \tilde{a}_{n-1}(x)}{b_{n-1}} \right)^2 \right] < \infty \right\} = 1. \quad (35)$$

The series

$$\sum_{n=1}^{\infty} \ln \frac{1 + d_{n-1}^2}{2d_{n-1}} \quad \text{and} \quad \sum_{n=1}^{\infty} (d_{n-1}^2 - 1)$$

converge or diverge together; hence it follows from (35) that

$$\tilde{P} \ll P \Leftrightarrow \tilde{P} \left\{ \sum_{n=0}^{\infty} \left[ \left( \frac{\tilde{b}_n^2}{b_n^2} - 1 \right)^2 + \frac{\Delta_n^2(x)}{b_n^2} \right] < \infty \right\} = 1, \quad (36)$$

where  $\Delta_n(x) = a_n(x) - \tilde{a}_n(x)$ .

Since  $a_n(x)$  and  $\tilde{a}_n(x)$  are linear, the sequence of random variables  $\{\Delta_n(x)/b_n\}_{n \geq 0}$  is a Gaussian system (with respect to both  $\tilde{P}$  and  $P$ ). As follows from a lemma that will be proved below, such sequences satisfy an analog of the zero-one law:

$$\tilde{P} \left\{ \sum \left( \frac{\Delta_n(x)}{b_n} \right)^2 < \infty \right\} = 1 \Leftrightarrow \sum \tilde{M} \left( \frac{\Delta_n(x)}{b_n} \right)^2 < \infty. \quad (37)$$

Hence it follows from (36) that

$$\tilde{P} \ll P \Leftrightarrow \sum_{n=0}^{\infty} \left[ \tilde{M} \left( \frac{\Delta_n(x)}{b_n} \right)^2 + \left( \frac{\tilde{b}_n^2}{b_n^2} - 1 \right)^2 \right] < \infty$$

and in a similar way

$$\tilde{P} \perp P \Leftrightarrow \sum_{n=0}^{\infty} \left[ \tilde{M} \left( \frac{\Delta_n(x)}{b_n} \right)^2 + \left( \frac{\tilde{b}_n^2}{b_n^2} - 1 \right)^2 \right] = \infty.$$

Then it is clear that if  $\tilde{P}$  and  $P$  are not singular measures, we have  $\tilde{P} \ll P$ . But by hypothesis,  $\tilde{P}_n \sim P_n$ ,  $n \geq 1$ ; hence by symmetry we have  $P \ll \tilde{P}$ . Therefore we have the following theorem.

**Theorem 5** (Hájek–Feldman Dichotomy). *Let  $\xi = (\xi_1, \xi_2, \dots)$  and  $\tilde{\xi} = (\tilde{\xi}_1, \tilde{\xi}_2, \dots)$  be Gaussian sequences whose finite-dimensional distributions are equivalent:  $\tilde{P}_n \sim P_n$ ,  $n \geq 1$ . Then either  $\tilde{P} \sim P$  or  $\tilde{P} \perp P$ . Moreover,*

$$\begin{aligned} \tilde{P} \sim P &\Leftrightarrow \sum_{n=0}^{\infty} \left[ \tilde{M} \left( \frac{\Delta_n(x)}{b_n} \right)^2 + \left( \frac{\tilde{b}_n^2}{b_n^2} - 1 \right)^2 \right] < \infty, \\ \tilde{P} \perp P &\Leftrightarrow \sum_{n=0}^{\infty} \left[ \tilde{M} \left( \frac{\Delta_n(x)}{b_n} \right)^2 + \left( \frac{\tilde{b}_n^2}{b_n^2} - 1 \right)^2 \right] = \infty. \end{aligned} \quad (38)$$

Let us now prove the zero-one law for Gaussian sequences that we need for the proof of Theorem 5.

**Lemma.** *Let  $\beta = (\beta_n)_{n \geq 1}$  be a Gaussian sequence defined on  $(\Omega, \mathcal{F}; P)$ . Then*

$$P \left\{ \sum_{n=1}^{\infty} \beta_n^2 < \infty \right\} = 1 \Leftrightarrow \sum_{n=1}^{\infty} E \beta_n^2 < \infty. \quad (39)$$

**PROOF.** The implication  $\Leftarrow$  follows from Fubini's theorem. To establish the opposite proposition, we first suppose that  $E \beta_n = 0$ ,  $n \geq 1$ . Here it is enough to show that

$$E \sum_{n=1}^{\infty} \beta_n^2 \leq \left[ E \exp \left( - \sum_{n=1}^{\infty} \beta_n^2 \right) \right]^{-2}, \quad (40)$$

since then the condition  $P\{\sum \beta_n^2 < \infty\} = 1$  will imply that the right-hand side of (40) is finite.

Select an  $n \geq 1$ . Then it follows from §§11 and 13, Chapter II, that there are independent Gaussian random variables  $\beta_{k,n}$ ,  $k = 1, \dots, r \leq n$ , with  $E\beta_{k,n} = 0$ , such that

$$\sum_{k=1}^n \beta_k^2 = \sum_{k=1}^r \beta_{k,n}^2.$$

If we write  $E\beta_{k,n}^2 = \lambda_{k,n}$ , we easily see that

$$E \sum_{k=1}^r \beta_{k,n}^2 = \sum_{k=1}^r \lambda_{k,n} \quad (41)$$

and

$$E \exp\left(-\sum_{k=1}^r \beta_{k,n}^2\right) = \prod_{k=1}^r (1 + 2\lambda_{k,n})^{-1/2}. \quad (42)$$

Comparing the right-hand sides of (41) and (42), we obtain

$$E \sum_{k=1}^n \beta_k^2 = E \sum_{k=1}^r \beta_{k,n}^2 \leq \left[ E \exp\left(-\sum_{k=1}^r \beta_{k,n}^2\right) \right]^{-2} = \left[ E \exp\left(-\sum_{k=1}^n \beta_k^2\right) \right]^{-2},$$

from which, by letting  $n \rightarrow \infty$ , we obtain the required inequality (40).

Now suppose that  $E\beta_n \neq 0$ .

Let us consider again the sequence  $\tilde{\beta} = (\tilde{\beta}_n)_{n \geq 1}$  with the same distribution as  $\beta = (\beta_n)_{n \geq 1}$  but independent of it (if necessary, extending the original probability space). If  $P\{\sum_{n=1}^{\infty} \beta_n^2 < \infty\} = 1$ , then  $P\{\sum_{n=1}^{\infty} (\beta_n - \tilde{\beta}_n)^2 < \infty\} = 1$ , and by what we have proved,

$$2 \sum_{n=1}^{\infty} E(\beta_n - E\beta_n)^2 = \sum_{n=1}^{\infty} E(\beta_n - \tilde{\beta}_n)^2 < \infty.$$

Since

$$(E\beta_n)^2 \leq 2\beta_n^2 + 2(\beta_n - E\beta_n)^2,$$

we have  $\sum_{n=1}^{\infty} (E\beta_n)^2 < \infty$  and therefore

$$\sum_{n=1}^{\infty} E\beta_n^2 = \sum_{n=1}^{\infty} (E\beta_n)^2 + \sum_{n=1}^{\infty} E(\beta_n - E\beta_n)^2 < \infty.$$

This completes the proof of the lemma.

**5.** We continue the discussion of the example in Subsection 3 of the preceding section, assuming that  $\xi_0, \xi_1, \dots$  are independent Gaussian random variables with  $E\xi_i = 0, \forall \xi_i = V_i > 0$ .

Again we let

$$X_{n+1} = \theta X_n + \xi_{n+1}$$

for  $n \geq 1$ , where  $X_0 = \xi_0$ , and the unknown parameter  $\theta$  that is to be estimated has values in  $R$ . Let  $\hat{\theta}_n$  be the least-squares estimator (see (5.22)).

**Theorem 6.** *A necessary and sufficient condition for the estimator  $\hat{\theta}_n$ ,  $n \geq 1$ , to be strongly consistent is that*

$$\sum_{n=0}^{\infty} \frac{V_n}{V_{n+1}} = \infty. \quad (43)$$

**PROOF.** *Sufficiency.* Let  $P_\theta$  denote the probability distribution on  $(R^\infty, \mathcal{B}_\infty)$  corresponding to the sequence  $(X_0, X_1, \dots)$  when the true value of the unknown parameter is  $\theta$ . Let  $E_\theta$  denote an average with respect to  $P_\theta$ .

We have already seen that

$$\hat{\theta}_n = \theta + \frac{M_n}{\langle M \rangle_n},$$

where

$$\langle M \rangle_n = \sum_{k=0}^{n-1} \frac{X_k^2}{V_{k+1}}.$$

According to the lemma from the preceding subsection,

$$P_\theta(\langle M \rangle_\infty = \infty) = 1 \Leftrightarrow E_\theta \langle M \rangle_\infty = \infty,$$

i.e.,  $\langle M \rangle_\infty = \infty$  ( $P_\theta$ -a.s.) if and only if

$$\sum_{k=0}^{\infty} \frac{E_\theta X_k^2}{V_{k+1}} = \infty. \quad (44)$$

But

$$E_\theta X_k^2 = \sum_{i=0}^k \theta^{2i} V_{k-i}$$

and

$$\begin{aligned} \sum_{k=0}^{\infty} \frac{E_\theta X_k^2}{V_{k+1}} &= \sum_{k=0}^{\infty} \frac{1}{V_{k+1}} \left( \sum_{i=0}^k \theta^{2i} V_{k-i} \right) \\ &= \sum_{k=0}^{\infty} \theta^{2k} \sum_{i=k}^{\infty} \frac{V_{i-k}}{V_{i+1}} = \sum_{i=0}^{\infty} \frac{V_i}{V_{i+1}} + \sum_{k=1}^{\infty} \theta^{2k} \left( \sum_{i=k}^{\infty} \frac{V_{i-k}}{V_{i+1}} \right). \end{aligned} \quad (45)$$

Hence (44) follows from (43) and therefore, by Theorem 4, the estimator  $\hat{\theta}_n$ ,  $n \geq 1$ , is strongly consistent for every  $\theta$ .

*Necessity.* For all  $\theta \in R$ , let  $P_\theta(\hat{\theta}_n \rightarrow \theta) = 1$ . It follows that if  $\theta_1 \neq \theta_2$ , the measures  $P_{\theta_1}$  and  $P_{\theta_2}$  are singular ( $P_{\theta_1} \perp P_{\theta_2}$ ). In fact, since the sequence  $(X_0, X_1, \dots)$  is Gaussian, by Theorem 5 of §5 the measures  $P_{\theta_1}$  and  $P_{\theta_2}$  are either singular or equivalent. But they cannot be equivalent, since if  $P_{\theta_1} \sim P_{\theta_2}$ ,

but  $P_{\theta_1}(\hat{\theta}_n \rightarrow \theta_1) = 1$ , then also  $P_{\theta_2}(\hat{\theta}_n \rightarrow \theta_1) = 1$ . However, by hypothesis,  $P_{\theta_2}(\hat{\theta}_n \rightarrow \theta_2) = 1$  and  $\theta_2 \neq \theta_1$ . Therefore  $P_{\theta_1} \perp P_{\theta_2}$  for  $\theta_1 \neq \theta_2$ .

According to (5.38),

$$P_{\theta_1} \perp P_{\theta_2} \Leftrightarrow (\theta_1 - \theta_2)^2 \sum_{k=0}^{\infty} E_{\theta_1} \left[ \frac{X_k^2}{V_{k+1}} \right] = \infty$$

for  $\theta_1 \neq \theta_2$ . Taking  $\theta_1 = 0$  and  $\theta_2 \neq 0$ , we obtain from (45) that

$$P_0 \perp P_{\theta_2} \Leftrightarrow \sum_{i=0}^{\infty} \frac{V_i}{V_{i+1}} = \infty,$$

which establishes the necessity of (43).

This completes the proof of the theorem.

## 6. PROBLEMS

1. Prove (6).

2. Let  $\tilde{P}_n \sim P_n$ ,  $n \geq 1$ . Show that

$$\tilde{P} \sim P \Leftrightarrow \tilde{P}\{z_{\infty} < \infty\} = P\{z_{\infty} > 0\} = 1,$$

$$\tilde{P} \perp P \Leftrightarrow \tilde{P}\{z_{\infty} = \infty\} = 1 \quad \text{or} \quad P\{z_{\infty} = 0\} = 1.$$

3. Let  $\tilde{P}_n \ll P_n$ ,  $n \geq 1$ , let  $\tau$  be a stopping time (with respect to  $(\mathcal{F}_n)$ ), and let  $\tilde{P}_{\tau} = \tilde{P}|_{\mathcal{F}_{\tau}}$  and  $P_{\tau} = P|_{\mathcal{F}_{\tau}}$  be the restrictions of  $\tilde{P}$  and  $P$  to the  $\sigma$ -algebra  $\mathcal{F}_{\tau}$ . Show that  $\tilde{P}_{\tau} \ll P_{\tau}$  if and only if  $\{\tau = \infty\} = \{z_{\infty} < \infty\}$  ( $\tilde{P}$ -a.s.). (In particular, if  $\tilde{P}\{\tau < \infty\} = 1$  then  $\tilde{P}_{\tau} \ll P_{\tau}$ .)

4. Prove (21) and (22).

5. Verify (28), (29), and (32).

6. Prove (34).

7. In Subsection 2 let the sequences  $\xi = (\xi_1, \xi_2, \dots)$  and  $\tilde{\xi} = (\tilde{\xi}_1, \tilde{\xi}_2, \dots)$  consist of independent identically distributed random variables. Show that if  $P_{\xi_1} \ll P_{\tilde{\xi}_1}$ , then  $\tilde{P} \ll P$  if and only if the measures  $P_{\xi_1}$  and  $P_{\tilde{\xi}_1}$  coincide. If, however,  $P_{\xi_1} \ll P_{\tilde{\xi}_1}$  and  $P_{\xi_1} \neq P_{\tilde{\xi}_1}$ , then  $\tilde{P} \perp P$ .

## §7. Asymptotics of the Probability of the Outcome of a Random Walk with Curvilinear Boundary

1. Let  $\xi_1, \xi_2, \dots$  be a sequence of independent identically distributed random variables. Let  $S_n = \xi_1 + \dots + \xi_n$ , let  $g = g(n)$  be a "boundary,"  $n \geq 1$ , and let

$$\tau = \inf\{n \geq 1: S_n < g(n)\}$$

be the first time at which the random walk  $(S_n)$  is found below the boundary  $g = g(n)$ . (As usual,  $\tau = \infty$  if  $\{\cdot\} = \emptyset$ .)

It is difficult to discover the exact form of the distribution of the time  $\tau$ . In the present section we find the asymptotic form of the probability  $P(\tau > n)$  as  $n \rightarrow \infty$ , for a wide class of boundaries  $g = g(n)$  and assuming that the  $\xi_i$  are normally distributed. The method of proof is based on the idea of an absolutely continuous change of measure together with a number of the properties of martingales and Markov times that were presented earlier.

**Theorem 1.** Let  $\xi_1, \xi_2, \dots$  be independent identically distributed random variables, with  $\xi_i \sim \mathcal{N}(0, 1)$ . Suppose that  $g = g(n)$  is such that  $g(1) < 0$  and, for  $n \geq 2$ ,

$$0 \leq \Delta g(n+1) \leq \Delta g(n), \quad (1)$$

where  $\Delta g(n) = g(n) - g(n-1)$  and

$$\ln n = o\left(\sum_{k=2}^n [\Delta g(k)]^2\right), \quad n \rightarrow \infty. \quad (2)$$

Then

$$P(\tau > n) = \exp\left\{-\frac{1}{2} \sum_{k=2}^n [\Delta g(k)]^2 (1 + o(1))\right\}, \quad n \rightarrow \infty. \quad (3)$$

Before starting the proof, let us observe that (1) and (2) are satisfied if, for example,

$$g(n) = an^\nu + b, \quad \frac{1}{2} < \nu \leq 1, \quad a + b < 0,$$

or (for sufficiently large  $n$ )

$$g(n) = n^\nu L(n), \quad \frac{1}{2} \leq \nu \leq 1,$$

where  $L(n)$  is a slowly varying function (for example,  $L(n) = C(\ln n)^\beta$  with arbitrary  $\beta$  for  $\frac{1}{2} < \nu < 1$  or with  $\beta > 0$  for  $\nu = \frac{1}{2}$ ).

**2.** We shall need the following two auxiliary propositions for the proof of Theorem 1.

Let us suppose that  $\xi_1, \xi_2, \dots$  is a sequence of independent identically distributed random variables,  $\xi_i \sim \mathcal{N}(0, 1)$ . Let  $\mathcal{F}_0 = \{\emptyset, \Omega\}$ ,  $\mathcal{F}_n = \sigma\{\omega: \xi_1, \dots, \xi_n\}$ , and let  $\alpha = (\alpha_n, \mathcal{F}_{n-1})$  be a predictable sequence with  $P(|\alpha_n| \leq C) = 1$ ,  $n \geq 1$ , where  $C$  is a constant. Form the sequence  $z = (z_n, \mathcal{F}_n)$  with

$$z_n = \exp\left\{\sum_{k=1}^n \alpha_k \xi_k - \frac{1}{2} \sum_{k=1}^n \alpha_k^2\right\}, \quad n \geq 1. \quad (4)$$

It is easily verified that (with respect to  $P$ ) the sequence  $z = (z_n, \mathcal{F}_n)$  is a martingale with  $Ez_n = 1$ ,  $n \geq 1$ .

Choose a value  $n \geq 1$  and introduce a probability measure  $\tilde{P}_n$  on the measurable space  $(\Omega, \mathcal{F}_n)$  by putting

$$\tilde{P}_n(A) = E I(A) z_n, \quad A \in \mathcal{F}_n. \quad (5)$$

**Lemma 1.** With respect to  $\tilde{P}_n$ , the random variables  $\tilde{\xi}_k = \xi_k - \alpha_k$ ,  $1 \leq k \leq n$ , are independent and normally distributed,  $\tilde{\xi}_k \sim \mathcal{N}(0, 1)$ .

**PROOF.** Let  $\tilde{E}_n$  denote averaging with respect to  $\tilde{P}_n$ . Then for  $\lambda_k \in R$ ,  $1 \leq k \leq n$ ,

$$\begin{aligned}\tilde{E}_n \exp \left\{ i \sum_{k=1}^n \lambda_k \tilde{\xi}_k \right\} &= E \exp \left\{ i \sum_{k=1}^n \lambda_k \tilde{\xi}_k \right\}_{Z_n} \\ &= E \left[ \exp \left\{ i \sum_{k=1}^{n-1} \lambda_k \tilde{\xi}_k \right\}_{Z_{n-1}} \cdot E \left\{ \exp \left( i \lambda_n (\xi_n - \alpha_n) + \alpha_n \xi_n - \frac{\alpha_n^2}{2} \right) \middle| \mathcal{F}_{n-1} \right\} \right] \\ &= E \left[ \exp \left\{ i \sum_{k=1}^{n-1} \lambda_k \tilde{\xi}_k \right\}_{Z_{n-1}} \right] \exp \left\{ -\frac{1}{2} \lambda_n^2 \right\} = \cdots = \exp \left\{ -\frac{1}{2} \sum_{k=1}^n \lambda_k^2 \right\}.\end{aligned}$$

Now the desired conclusion follows from Theorem 4 of §12, Chapter II.

**Lemma 2.** Let  $X = (X_n, \mathcal{F}_n)_{n \geq 1}$  be a square-integrable martingale with mean zero and

$$\sigma = \inf \{n \geq 1: X_n \leq -b\},$$

where  $b$  is a constant,  $b > 0$ . Suppose that

$$P(X_1 < -b) > 0.$$

Then there is a constant  $C > 0$  such that, for all  $n \geq 1$ ,

$$P(\sigma > n) \geq \frac{C}{EX_n^2}. \quad (6)$$

**PROOF.** By Corollary 1 to Theorem VII.2.1 we have  $EX_{\sigma \wedge n} = 0$ , whence

$$-EI(\sigma \leq n)X_\sigma = EI(\sigma > n)X_n. \quad (7)$$

On the set  $\{\sigma \leq n\}$

$$-X_\sigma \geq b > 0.$$

Therefore, for  $n \geq 1$ ,

$$-EI(\sigma \leq n)X_\sigma \geq bP(\sigma \leq n) \geq bP(\sigma = 1) = bP(X_1 < -b) > 0. \quad (8)$$

On the other hand, by the Cauchy-Schwarz inequality,

$$EI(\sigma > n)X_n \leq [P(\sigma > n) \cdot EX_n^2]^{1/2}, \quad (9)$$

which, with (7) and (8), leads to the required inequality with

$$C = (bP(X_1 < -b))^2.$$

**PROOF OF THEOREM 1.** It is enough to show that

$$\lim_{n \rightarrow \infty} \ln P(\tau > n) \bigg/ \sum_{k=2}^n [\Delta g(k)]^2 \geq -\frac{1}{2} \quad (10)$$

and

$$\overline{\lim}_{n \rightarrow \infty} \ln P(\tau > n) \bigg/ \sum_{k=2}^n [\Delta g(k)]^2 \leq -\frac{1}{2}. \quad (11)$$

For this purpose we consider the (nonrandom) sequence  $(\alpha_n)_{n \geq 1}$  with

$$\alpha_1 = 0, \quad \alpha_n = \Delta g(n), \quad n \geq 2,$$

and the probability measure  $(\tilde{P}_n)_{n \geq 1}$  defined by (5). Then by Hölder's inequality

$$\tilde{P}_n(\tau > n) = EI(\tau > n) z_n \leq (P(\tau > n))^{1/q} (E z_n^p)^{1/p}, \quad (12)$$

where  $p > 1$  and  $q = p/(p-1)$ .

The last factor is easily calculated explicitly:

$$(E z_n^p)^{1/p} = \exp \left\{ \frac{p-1}{2} \sum_{k=2}^n [\Delta g(k)]^2 \right\}. \quad (13)$$

Now let us estimate the probability  $\tilde{P}_n(\tau > n)$  that appears on the left-hand side of (12). We have

$$\tilde{P}_n(\tau > n) = \tilde{P}_n(S_k \geq g(k), 1 \leq k \leq n) = \tilde{P}_n(\tilde{S}_k \geq g(1), 1 \leq k \leq n),$$

where  $\tilde{S}_k = \sum_{i=1}^k \tilde{\xi}_i$ ,  $\tilde{\xi}_i = \xi_i - \alpha_i$ . By Lemma 1, the variables are independent and normally distributed,  $\tilde{\xi}_i \sim \mathcal{N}(0, 1)$ , with respect to the measure  $\tilde{P}_n$ . Then by Lemma 2 (applied to  $b = -g(1)$ ,  $P = \tilde{P}_n$ ,  $X_n = \tilde{S}_n$ ) we find that

$$\tilde{P}(\tau > n) \geq \frac{c}{n}, \quad (14)$$

where  $c$  is a constant.

Then it follows from (12)–(14) that, for every  $p > 1$ ,

$$P(\tau > n) \geq C_p \exp \left\{ -\frac{p}{2} \sum_{k=2}^n [\Delta g(k)]^2 - \frac{p}{p-1} \ln n \right\}, \quad (15)$$

where  $C_p$  is a constant. Then (15) implies the lower bound (10) by the hypotheses of the theorem, since  $p > 1$  is arbitrary.

To obtain the upper bound (11), we first observe that since  $z_n > 0$  ( $P$ - and  $\tilde{P}$ -a.s.), we have by (5)

$$P(\tau > n) = \bar{E}_n I(\tau > n) z_n^{-1}, \quad (16)$$

where  $\bar{E}_n$  denotes an average with respect to  $\tilde{P}_n$ .

In the case under consideration,  $\alpha_1 = 0$ ,  $\alpha_n = \Delta g(n)$ ,  $n \geq 2$ , and therefore for  $n \geq 2$

$$z_n^{-1} = \exp \left\{ -\sum_{k=2}^n \Delta g(k) \cdot \xi_k + \frac{1}{2} \sum_{k=2}^n [\Delta g(k)]^2 \right\}.$$



By the formula for summation by parts (see the proof of Lemma 2 of §3, Chapter IV)

$$\sum_{k=2}^n \Delta g(k) \cdot \xi_k = \Delta g(n) \cdot S_n - \sum_{k=2}^n S_{k-1} \Delta(\Delta g(k)).$$

Hence if we recall that by hypothesis  $\Delta g(k) \geq 0$  and  $\Delta(\Delta g(k)) \leq 0$ , we find that, on the set  $\{\tau > n\} = \{S_k \geq g(k), 1 \leq k \leq n\}$ ,

$$\begin{aligned} \sum_{k=2}^n \Delta g(k) \cdot \xi_k &\geq \Delta g(n) \cdot g(n) - \sum_{k=3}^n g(k-1) \Delta(\Delta g(k)) - \xi_1 \Delta g(2) \\ &= \sum_{k=2}^n [\Delta g(k)]^2 + g(1) \Delta g(2) - \xi_1 \Delta g(2). \end{aligned}$$

Thus, by (16),

$$\begin{aligned} P(\tau > n) &\leq \exp \left\{ -\frac{1}{2} \sum_{k=2}^n [\Delta g(k)]^2 - g(1) \Delta g(2) \right\} \tilde{E}_n I(\tau > n) e^{-\xi_1 \Delta g(2)} \\ &\leq \exp \left\{ -\frac{1}{2} \sum_{k=2}^n [\Delta g(k)]^2 \right\} \tilde{E}_n I(\tau > n) e^{-\xi_1 \Delta g(2)}, \end{aligned}$$

where

$$\tilde{E}_n I(\tau > n) e^{-\xi_1 \Delta g(2)} \leq E z_n e^{-\xi_1 \Delta g(2)} = E e^{-\xi_1 \Delta g(2)} < \infty.$$

Therefore

$$P(\tau > n) \leq C \exp \left\{ -\frac{1}{2} \sum_{k=2}^n [\Delta g(k)]^2 \right\},$$

where  $C$  is a positive constant; this establishes the upper bound (11).

This completes the proof of the theorem.

3. The idea of an absolutely continuous change of measure can be used to study similar problems, including the case of a two-sided boundary. We present (without proof) a result in this direction.

**Theorem 2.** Let  $\xi_1, \xi_2, \dots$  be independent identically distributed random variables with  $\xi_i \sim \mathcal{N}(0, 1)$ . Suppose that  $f = f(n)$  is a positive function such that

$$f(n) \rightarrow \infty, \quad n \rightarrow \infty,$$

and

$$\sum_{k=2}^n [\Delta f(k)]^2 = o \left( \sum_{k=1}^n f^{-2}(k) \right), \quad n \rightarrow \infty.$$

Then if

$$\sigma = \inf\{n \geq 1: |S_n| \geq f(n)\},$$

we have

$$P(\sigma > n) = \exp \left\{ -\frac{\pi^2}{8} \sum_{k=1}^n f^{-2}(k)(1 + o(1)) \right\}, \quad n \rightarrow \infty. \quad (17)$$

#### 4. PROBLEMS

1. Show that the sequence defined in (4) is a martingale.
2. Establish (13).
3. Prove (17).

## §8. Central Limit Theorem for Sums of Dependent Random Variables

1. In §4, Chapter III, the central limit theorem for sums  $S_n = \xi_{n1} + \dots + \xi_{nn}$ ,  $n \geq 1$ , of random variables  $\xi_{n1}, \dots, \xi_{nn}$  was established under the assumptions of their independence, finiteness of second moments, and negligibility in the limit of their terms. In the present section, we give up both the assumption of independence and even that of the finiteness of the absolute values of the first-order moments. However, the negligibility in the limit of the terms will be retained.

Thus, we suppose that on the probability space  $(\Omega, \mathcal{F}, P)$  there are given stochastic sequences

$$\xi^n = (\xi_{nk}, \mathcal{F}_k^n), \quad 0 \leq k \leq n, \quad n \geq 1,$$

with  $\xi_{n0} = 0$ ,  $\mathcal{F}_0^n = \{\emptyset, \Omega\}$ ,  $\mathcal{F}_k^n \subseteq \mathcal{F}_{k+1}^n \subseteq \mathcal{F}$ . We set

$$X_t^n = \sum_{k=0}^{[nt]} \xi_{nk}, \quad 0 \leq t \leq 1.$$

**Theorem 1.** For a given  $t$ ,  $0 < t \leq 1$ , let the following conditions be satisfied: for each  $\varepsilon \in (0, 1)$ , as  $n \rightarrow \infty$ ,

- (A)  $\sum_{k=1}^{[nt]} P(|\xi_{nk}| > \varepsilon | \mathcal{F}_{k-1}^n) \xrightarrow{P} 0,$
- (B)  $\sum_{k=1}^{[nt]} E[\xi_{nk} I(|\xi_{nk}| \leq 1) | \mathcal{F}_{k-1}^n] \xrightarrow{P} 0,$
- (C)  $\sum_{k=1}^{[nt]} V[\xi_{nk} I(|\xi_{nk}| \leq \varepsilon) | \mathcal{F}_{k-1}^n] \xrightarrow{P} \sigma_t^2 \geq 0.$

Then

$$X_t^n \xrightarrow{d} \mathcal{N}(0, \sigma_t^2).$$

**Remark 1.** Hypotheses (A) and (B) guarantee that  $X_t^n$  can be represented in the form  $X_t^n = Y_t^n + Z_t^n$  with  $Z_t^n \xrightarrow{P} 0$  and  $Y_t^n = \sum_{k=0}^{[nt]} \eta_{nk}$ , where the sequence  $\eta^n = (\eta_{nk}, \mathcal{F}_k^n)$  is a martingale-difference, and  $E(\eta_{nk} | \mathcal{F}_{k-1}^n) = 0$  with  $|\eta_{nk}| \leq c$ , uniformly for  $1 \leq k \leq n$  and  $n \geq 1$ . Consequently, in the cases under consideration, the proof reduces to proving the central limit theorem for martingale-differences.

In the case when the variables  $\xi_{n1}, \dots, \xi_{nm}$  are independent, conditions (A), (B), and (C), with  $t = 1$ , and  $\sigma^2 = \sigma_1^2$ , become

- (a)  $\sum_{k=1}^n P(|\xi_{nk}| > \varepsilon) \rightarrow 0$ ,  
 (b)  $\sum_{k=1}^n E[\xi_{nk} I(|\xi_{nk}| \leq 1)] \rightarrow 0$ ,  
 (c)  $\sum_{k=1}^n V[\xi_{nk} I(|\xi_{nk}| \leq \varepsilon)] \rightarrow \sigma^2$ .

These are well known; see the book by Gnedenko and Kolmogorov [G5]. Hence we have the following corollary to Theorem 1.

**Corollary.** If  $\xi_{n1}, \dots, \xi_{nm}$  are independent random variables,  $n \geq 1$ , then

$$(a), (b), (c) \Rightarrow X_1^n \xrightarrow{d} \mathcal{N}(0, \sigma^2).$$

**Remark 2.** In hypothesis (C), the case  $\sigma_t^2 = 0$  is not excluded. Hence, in particular, Theorem 1 yields a convergence condition for degenerate distributions ( $X_t^n \xrightarrow{d} 0$ ).

**Remark 3.** The method used to prove Theorem 1 lets us state and prove the following more general proposition.

Let  $0 < t_1 < t_2 < \dots < t_j \leq 1$ ,  $\sigma_{t_1}^2 \leq \sigma_{t_2}^2 \leq \dots \leq \sigma_{t_j}^2$ ,  $\sigma_0^2 = 0$ , and let  $\varepsilon_1, \dots, \varepsilon_j$  be independent Gaussian random variables with zero means and  $E\varepsilon_k^2 = \sigma_{t_k}^2 - \sigma_{t_{k-1}}^2$ . Form the (Gaussian) vectors  $(W_{t_1}, \dots, W_{t_j})$  with  $W_{t_k} = \varepsilon_1 + \dots + \varepsilon_k$ .

Let conditions (A), (B), and (C) be satisfied for  $t = t_1, \dots, t_j$ . Then the joint distribution  $(P_{t_1, \dots, t_j}^n)$  of the random variables  $(X_{t_1}^n, \dots, X_{t_j}^n)$  converges weakly to the Gaussian distribution  $P(t_1, \dots, t_j)$  of the variables  $(W_{t_1}, \dots, W_{t_j})$ :

$$P_{t_1, \dots, t_j}^n \xrightarrow{w} P_{t_1, \dots, t_j}.$$

2. The first assertion of the following theorem shows that condition (A) is equivalent to the condition of negligibility in the limit already introduced in §4, chapter III:

$$(A^*) \quad \max_{1 \leq k \leq [nt]} |\xi_{nk}| \xrightarrow{P} 0.$$

**Theorem 2.**

- (1) Condition (A) is equivalent to (A\*);  
 (2) Assuming (A) or (A\*), condition (C) is equivalent to

$$(C^*) \sum_{k=0}^{[nt]} [\xi_{nk} - E(\xi_{nk} | \mathcal{F}_{k-1}^n)]^2 \xrightarrow{P} \sigma_t^2.$$

**Theorem 3.** For each  $n \geq 1$  let the sequence

$$\xi^n = (\xi_{nk}, \mathcal{F}_k^n), \quad 1 \leq k \leq n,$$

be a square-integrable martingale-difference:

$$E\xi_{nk}^2 < \infty, \quad E(\xi_{nk} | \mathcal{F}_{k-1}^n) = 0.$$

Suppose that the Lindeberg condition is satisfied: for  $\varepsilon > 0$ ,

$$(L) \sum_{k=0}^{[nt]} E[\xi_{nk}^2 I(|\xi_{nk}| > \varepsilon) | \mathcal{F}_{k-1}^n] \xrightarrow{P} 0.$$

Then (C) is equivalent to

$$\langle X^n \rangle_t \xrightarrow{P} \sigma_t^2, \quad (1)$$

where (quadratic variation)

$$\langle X^n \rangle_t = \sum_{k=0}^{[nt]} E(\xi_{nk}^2 | \mathcal{F}_{k-1}^n), \quad (2)$$

and (C\*) is equivalent to

$$[X^n]_t \xrightarrow{P} \sigma_t^2, \quad (3)$$

where (quadratic variation)

$$[X^n]_t = \sum_{k=0}^{[nt]} \xi_{nk}^2. \quad (4)$$

The next theorem is a corollary of Theorems 1–3.

**Theorem 4.** Let the square-integrable martingale-differences  $\xi^n = (\xi_{nk}, \mathcal{F}_k^n)$ ,  $n \geq 1$ , satisfy (for a given  $t$ ,  $0 < t \leq 1$ ) the Lindeberg condition (L). Then

$$\sum_{k=0}^{[nt]} E(\xi_{nk}^2 | \mathcal{F}_{k-1}^n) \xrightarrow{P} \sigma_t^2 \Rightarrow X_t^n \xrightarrow{d} \mathcal{N}(0, \sigma_t^2), \quad (5)$$

$$\sum_{k=0}^{[nt]} \xi_{nk}^2 \xrightarrow{P} \sigma_t^2 \Rightarrow X_t^n \xrightarrow{d} \mathcal{N}(0, \sigma_t^2). \quad (6)$$

### 3. PROOF OF THEOREM 1. Let us represent $X_t^n$ in the form

$$X_t^n = \sum_{k=0}^{[nt]} \xi_{nk} I(|\xi_{nk}| \leq 1) + \sum_{k=0}^{[nt]} \xi_{nk} I(|\xi_{nk}| > 1)$$

$$\begin{aligned}
&= \sum_{k=0}^{[nt]} \mathbb{E}[\xi_{nk} I(|\xi_{nk}| \leq 1) \mid \mathcal{F}_{k-1}^n] + \sum_{k=0}^{[nt]} \xi_{nk} I(|\xi_{nk}| > 1) \\
&\quad + \sum_{k=0}^{[nt]} \{\xi_{nk} I(|\xi_{nk}| \leq 1) - \mathbb{E}[\xi_{nk} I(|\xi_{nk}| \leq 1) \mid \mathcal{F}_{k-1}^n]\}. \quad (7)
\end{aligned}$$

We define

$$\begin{aligned}
B_t^n &= \sum_{k=0}^{[nt]} \mathbb{E}[\xi_{nk} I(|\xi_{nk}| \leq 1) \mid \mathcal{F}_{k-1}^n], \\
\mu_k^n(\Gamma) &= I(\xi_{nk} \in \Gamma), \\
\nu_k^n(\Gamma) &= P(\xi_{nk} \in \Gamma \mid \mathcal{F}_{k-1}^n),
\end{aligned} \quad (8)$$

where  $\Gamma$  is a set from the smallest  $\sigma$ -algebra  $\sigma(R \setminus \{0\})$  and  $P(\xi_{nk} \in \Gamma \mid \mathcal{F}_{k-1}^n)$  is the regular conditional distribution of  $\xi_{nk}$  with respect to  $\mathcal{F}_{k-1}^n$ .

Then (7) can be rewritten in the following form:

$$X_t^n = B_t^n + \sum_{k=0}^{[nt]} \int_{|x| > 1} x d\mu_k^n + \sum_{k=0}^{[nt]} \int_{|x| \leq 1} x d(\mu_k^n - \nu_k^n), \quad (9)$$

which is known as the *canonical* decomposition of  $(X_t^n, \mathcal{F}_t^n)$ . (The integrals are to be understood as Lebesgue–Stieltjes integrals, defined for every sample point.)

According to (B) we have  $B_t^n \xrightarrow{P} 0$ . Let us show that (A) implies

$$\sum_{k=0}^{[nt]} \int_{|x| > 1} |x| d\mu_k^n \xrightarrow{P} 0. \quad (10)$$

We have

$$\sum_{k=0}^{[nt]} \int_{|x| > 1} |x| d\mu_k^n = \sum_{k=0}^{[nt]} |\xi_{nk}| I(|\xi_{nk}| > 1). \quad (11)$$

For every  $\delta \in (0, 1)$ ,

$$\left\{ \sum_{k=0}^{[nt]} |\xi_{nk}| I(|\xi_{nk}| > 1) > \delta \right\} \leq \left\{ \sum_{k=0}^{[nt]} I(|\xi_{nk}| > 1) > \delta \right\}. \quad (12)$$

It is clear that

$$\sum_{k=0}^{[nt]} I(|\xi_{nk}| > 1) = \sum_{k=0}^{[nt]} \int_{|x| > 1} d\mu_k^n (\equiv U_{[nt]}^n).$$

By (A),

$$V_{[nt]}^n \equiv \sum_{k=0}^{[nt]} \int_{|x|>1} dv_k^n \xrightarrow{P} 0, \quad (13)$$

and  $V_k^n$  is  $\mathcal{F}_{k-1}^n$ -measurable.

Then by the corollary to Theorem 2 of §3, Chapter VII,

$$V_{[nt]}^n \xrightarrow{P} 0 \Rightarrow U_{[nt]}^n \xrightarrow{P} 0. \quad (14)$$

(By the same corollary and the inequality  $\Delta U_{[nt]}^n \leq 1$ , we also have the converse implication

$$U_{[nt]}^n \xrightarrow{P} 0 \Rightarrow V_{[nt]}^n \xrightarrow{P} 0, \quad (15)$$

which will be needed in the proof of Theorem 2.)

The required proposition (10) now follows from (11)–(14).

Thus

$$X_t^n = Y_t^n + Z_t^n, \quad (16)$$

where

$$Y_t^n = \sum_{k=0}^{[nt]} \int_{|x| \leq 1} x d(\mu_k^n - \nu_k^n), \quad (17)$$

and

$$Z_t^n = B_t^n + \sum_{k=0}^{[nt]} \int_{|x|>1} x d\mu_k^n \xrightarrow{P} 0. \quad (18)$$

It then follows by Problem 1 that to establish that

$$X_t^n \xrightarrow{d} \mathcal{N}(0, \sigma_t^2)$$

we need only show that

$$Y_t^n \xrightarrow{d} \mathcal{N}(0, \sigma_t^2). \quad (19)$$

Let us represent  $Y_t^n$  in the form

$$Y_t^n = \gamma_{[nt]}^n(\varepsilon) + \Delta_{[nt]}^n(\varepsilon), \quad \varepsilon \in (0, 1],$$

where

$$\gamma_{[nt]}^n(\varepsilon) = \sum_{k=0}^{[nt]} \int_{\varepsilon < |x| \leq 1} x d(\mu_k^n - \nu_k^n), \quad (20)$$

$$\Delta_{[nt]}^n(\varepsilon) = \sum_{k=0}^{[nt]} \int_{|x| \leq \varepsilon} x d(\mu_k^n - \nu_k^n). \quad (21)$$

As in the proof of (14), it is easily verified that, because of (A), we have  $\gamma_{[nt]}^n(\varepsilon) \xrightarrow{P} 0$ ,  $n \rightarrow \infty$ .

The sequence  $\Delta^n(\varepsilon) = (\Delta_k^n(\varepsilon), \mathcal{F}_k^n)$ ,  $1 \leq k \leq n$ , is a square-integrable martingale with quadratic variation

$$\begin{aligned} \langle \Delta^n(\varepsilon) \rangle_k &= \sum_{i=0}^k \left[ \int_{|x| \leq \varepsilon} x^2 dv_i^n - \left( \int_{|x| \leq \varepsilon} x dv_i^n \right)^2 \right] \\ &= \sum_{i=0}^k V[\xi_{ni} I(|\xi_{ni}| \leq \varepsilon) \mid \mathcal{F}_{i-1}^n]. \end{aligned}$$

Because of (C),

$$\langle \Delta^n(\varepsilon) \rangle_{[nt]} \xrightarrow{P} \sigma_t^2.$$

Hence, for every  $\varepsilon \in (0, 1]$ ,

$$\max\{\gamma_{[nt]}^n(\varepsilon), |\langle \Delta^n(\varepsilon) \rangle_{[nt]} - \sigma_t^2|\} \xrightarrow{P} 0.$$

By Problem 2 there is then a sequence of numbers  $\varepsilon_n \downarrow 0$  such that

$$\gamma_{[nt]}^n(\varepsilon_n) \xrightarrow{P} 0, \quad \langle \Delta^n(\varepsilon_n) \rangle_{[nt]} \xrightarrow{P} \sigma_t^2.$$

Therefore, again by Problem 1, it is enough to prove only that

$$M_{[nt]}^n \xrightarrow{d} \mathcal{N}(0, \sigma_t^2), \quad (22)$$

where

$$M_k^n \equiv \Delta_k^n(\varepsilon_n) = \sum_{i=0}^k \int_{|x| \leq \varepsilon_n} x(\mu_i^n - \nu_i^n). \quad (23)$$

For  $\Gamma \in \sigma(R \setminus \{0\})$ , let

$$\tilde{\mu}_k^n(\Gamma) = I(\Delta M_k^n \in \Gamma), \quad \tilde{\nu}_k^n(\Gamma) = P(\Delta M_k^n \in \Gamma \mid \mathcal{F}_{k-1}^n)$$

be a regular conditional probability,  $\Delta M_k^n = M_k^n - M_{k-1}^n$ ,  $k \geq 1$ ,  $M_0^n = 0$ . Then the square-integrable martingale  $M^n = (M_k^n, \mathcal{F}_k^n)$ ,  $1 \leq k \leq n$ , can evidently be written in the form

$$M_k^n = \sum_{i=1}^k \Delta M_i^n = \sum_{i=1}^k \int_{|x| \leq 2\varepsilon_n} x d\tilde{\mu}_i^n.$$

(Notice that  $|\Delta M_i^n| \leq 2\varepsilon_n$  by (23).)

To establish (22) we have to show that, for every real  $\lambda$ ,

$$E \exp\{i\lambda M_{[nt]}^n\} \rightarrow \exp(-\tfrac{1}{2}\lambda^2 \sigma_t^2). \quad (24)$$

Put

$$G_k^n = \sum_{j=1}^k \int_{|x| \leq 2\varepsilon_n} (e^{i\lambda x} - 1) d\tilde{\nu}_j^n$$

and

$$\mathcal{G}_k^n(G^n) = \prod_{j=1}^k (1 + \Delta G_j^n).$$

Observe that

$$\begin{aligned} 1 + \Delta G_k^n &= 1 + \int_{|x| \leq 2\varepsilon_n} (e^{i\lambda x} - 1) d\tilde{v}_k^n = \int_{|x| \leq 2\varepsilon_n} e^{i\lambda x} d\tilde{v}_k^n \\ &= \mathbf{E}[\exp(i\lambda \Delta M_k^n) | \mathcal{F}_{k-1}^n] \end{aligned}$$

and consequently

$$\mathcal{E}_k^n(G^n) = \prod_{j=1}^k (1 + \Delta G_j^n).$$

On the basis of a lemma that will be proved in Subsection 4, (24) will follow if for every real  $\lambda$

$$|\mathcal{E}_{[nt]}^n(G^n)| = \left| \prod_{j=1}^{[nt]} \mathbf{E}[\exp(i\lambda \Delta M_j^n) | \mathcal{F}_{j-1}^n] \right| \geq C(\lambda) > 0 \quad (25)$$

and

$$\mathcal{E}_{[nt]}^n(G^n) \xrightarrow{P} \exp(-\frac{1}{2}\lambda^2\sigma_t^2). \quad (26)$$

To see this we represent  $\mathcal{E}_k^n(G^n)$  in the form

$$\mathcal{E}_k^n(G^n) = \exp(G_k^n) \cdot \prod_{j=1}^k (1 + \Delta G_j^n) \exp(-\Delta G_j^n).$$

(Compare the function  $E_t(A)$  defined by (76) in §6, Chapter II.)

Since

$$\int_{|x| \leq 2\varepsilon_n} x d\tilde{v}_j^n = \mathbf{E}(\Delta M_j^n | \mathcal{F}_{j-1}^n) = 0,$$

we have

$$G_k^n = \sum_{j=1}^k \int_{|x| \leq 2\varepsilon_n} (e^{i\lambda x} - 1 - i\lambda x) d\tilde{v}_j^n. \quad (27)$$

Therefore

$$\begin{aligned} |\Delta G_k^n| &\leq \int_{|x| \leq 2\varepsilon_n} |e^{i\lambda x} - 1 - i\lambda x| d\tilde{v}_k^n \leq \frac{1}{2}\lambda^2 \int_{|x| \leq 2\varepsilon_n} x^2 d\tilde{v}_k^n \\ &\leq \frac{1}{2}\lambda^2 (2\varepsilon_n)^2 \rightarrow 0 \end{aligned} \quad (28)$$

and

$$\sum_{j=1}^k |\Delta G_j^n| \leq \frac{1}{2}\lambda^2 \sum_{j=1}^k \int_{|x| \leq 2\varepsilon_n} x^2 d\tilde{v}_j^n = \frac{1}{2}\lambda^2 \langle M^n \rangle_k. \quad (29)$$

By (C),

$$\langle M^n \rangle_{[nt]} \xrightarrow{P} \sigma_t^2. \quad (30)$$



Suppose first that  $\langle M^n \rangle_k \leq a$  (P-a.s.),  $k \leq [nt]$ , where  $a \geq \sigma_t^2 + 1$ . Then by (28), (29), and Problem 3,

$$\prod_{k=1}^{[nt]} (1 + \Delta G_k^n) \exp(-\Delta G_k^n) \xrightarrow{P} 1, \quad n \rightarrow \infty,$$

and therefore to establish (26) we only have to show that

$$G_{[nt]}^n \rightarrow -\frac{1}{2}\lambda^2\sigma_t^2, \quad (31)$$

i.e., after (27), (29), and (30), that

$$\sum_{k=1}^{[nt]} \int_{|x| \leq 2\varepsilon_n} (e^{i\lambda x} - 1 - i\lambda x + \frac{1}{2}\lambda^2 x^2) d\tilde{v}_k^n \xrightarrow{P} 0. \quad (32)$$

But

$$|e^{i\lambda x} - 1 - i\lambda x + \frac{1}{2}\lambda^2 x^2| \leq \frac{1}{6}|\lambda x|^3$$

and therefore

$$\begin{aligned} \sum_{k=1}^{[nt]} \int_{|x| \leq 2\varepsilon_n} |e^{i\lambda x} - 1 - i\lambda x + \frac{1}{2}\lambda^2 x^2| d\tilde{v}_k^n &\leq \frac{1}{6}|\lambda|^3 (2\varepsilon_n) \sum_{k=1}^{[nt]} \int_{|x| \leq 2\varepsilon_n} x^2 d\tilde{v}_k^n \\ &= \frac{1}{3}\varepsilon_n |\lambda|^3 \langle M_n \rangle_{[nt]} \leq \frac{1}{3}\varepsilon_n |\lambda|^3 a \rightarrow 0, \quad n \rightarrow \infty. \end{aligned}$$

Therefore if  $\langle M^n \rangle_{[nt]} \leq a$  (P-a.s.), (31) is established and consequently (26) is established also.

Let us now verify (25). Since  $|e^{i\lambda x} - 1 - i\lambda x| \leq \frac{1}{2}(\lambda x)^2$ , we find from (28) that, for sufficiently large  $n$ ,

$$\begin{aligned} |\mathcal{E}_k^n(G^n)| &= \left| \prod_{j=1}^k (1 + \Delta G_j^n) \right| \geq \prod_{j=1}^k (1 - \frac{1}{2}\lambda^2 \Delta \langle M^n \rangle_j) \\ &= \exp \left\{ \sum_{j=1}^k \ln(1 - \frac{1}{2}\lambda^2 \Delta \langle M^n \rangle_j) \right\}. \end{aligned}$$

But

$$\ln(1 - \frac{1}{2}\lambda^2 \Delta \langle M^n \rangle_j) \geq -\frac{\frac{1}{2}\lambda^2 \Delta \langle M^n \rangle_j}{1 - \frac{1}{2}\lambda^2 \Delta \langle M^n \rangle_j}$$

and  $\Delta \langle M^n \rangle_j \leq (2\varepsilon_n)^2 \downarrow 0$ ,  $n \rightarrow \infty$ . Therefore there is an  $n_0 = n_0(\lambda)$  such that for all  $n \geq n_0(\lambda)$ ,

$$|\mathcal{E}_k^n(G^n)| \geq \exp\{-\lambda^2 \langle M^n \rangle_k\}$$

and therefore

$$|\mathcal{E}_{[nt]}^n(G^n)| \geq \exp\{-\lambda^2 \langle M^n \rangle_{[nt]}\} \geq e^{-\lambda^2 a}.$$

Hence the theorem is proved under the assumption that  $\langle M^n \rangle_{[nt]} \leq a$  (P-a.s.). To remove this assumption, we proceed as follows.

Let

$$\tau^n = \min\{k \leq [nt]: \langle M^n \rangle_k \geq \sigma_t^2 + 1\},$$

taking  $\tau^n = \infty$  if  $\langle M^n \rangle_{[nt]} \leq \sigma_t^2 + 1$ .

Then for  $\bar{M}_k^n = M_{k \wedge \tau^n}^n$  we have

$$\langle \bar{M}^n \rangle_{[nt]} = \langle M^n \rangle_{[nt] \wedge \tau^n} \leq 1 + \sigma_t^2 + 2\varepsilon_n^2 \leq 1 + \sigma_t^2 + 2\varepsilon_1^2 (=a),$$

and by what has been proved,

$$\mathbf{E} \exp\{i\lambda \bar{M}_{[nt]}^n\} \rightarrow \exp(-\tfrac{1}{2}\lambda^2 \sigma_t^2).$$

But

$$\lim_n |\mathbf{E}\{\exp(i\lambda M_{[nt]}^n) - \exp(i\lambda \bar{M}_{[nt]}^n)\}| \leq 2 \lim_n \mathbf{P}(\tau^n < \infty) = 0.$$

Consequently

$$\begin{aligned} \lim_n \mathbf{E} \exp(i\lambda M_{[nt]}^n) &= \lim_n \mathbf{E}\{\exp(i\lambda M_{[nt]}^n) - \exp(i\lambda \bar{M}_{[nt]}^n)\} \\ &\quad + \lim_n \mathbf{E} \exp(i\lambda \bar{M}_{[nt]}^n) = \exp(-\tfrac{1}{2}\lambda^2 \sigma_t^2). \end{aligned}$$

This completes the proof of Theorem 1.

**Remark.** To prove the statement made in Remark 2 to Theorem 1, we need (using the Cramér–Wold method [B3]) to show that for all real numbers  $\lambda_1, \dots, \lambda_j$

$$\begin{aligned} \mathbf{E} \exp i \left[ \lambda_1 M_{[nt_1]}^n + \sum_{k=2}^j \lambda_k (M_{[nt_k]}^n - M_{[nt_{k-1}]}^n) \right] \\ \rightarrow \exp(-\tfrac{1}{2}\lambda_1^2 \sigma_{t_1}^2) - \sum_{k=2}^j \tfrac{1}{2} \lambda_k^2 (\sigma_{t_k}^2 - \sigma_{t_{k-1}}^2). \end{aligned}$$

The proof of this is similar to the proof of (24), replacing  $(M_k^n, \mathcal{F}_k^n)$  by the square-integrable martingales  $(\hat{M}_k^n, \mathcal{F}_k^n)$ ,

$$\hat{M}_k^n = \sum_{i=1}^k v_i \Delta M_i^n,$$

where  $v_i = \lambda_1$  for  $i \leq [nt_1]$  and  $v_i = \lambda_j$  for  $[nt_{j-1}] < i \leq [nt_j]$ .

4. In this subsection we prove a simple lemma which lets us reduce the verification of (24) to the verification of (25) and (26).

Let  $\eta^n = (\eta_{nk}, \mathcal{F}_k^n)$ ,  $1 \leq k \leq n$ ,  $n \geq 1$ , be stochastic sequences, let

$$Y^n = \sum_{k=1}^n \eta_{nk},$$

let

$$\mathcal{E}^n(\lambda) = \prod_{k=1}^n \mathbf{E}[\exp(i\lambda \eta_{nk}) | \mathcal{F}_{k-1}^n], \quad \lambda \in R,$$

and let  $Y$  be a random variable with

$$\mathcal{E}(\lambda) = \mathbb{E}e^{i\lambda Y}, \quad \lambda \in \mathbb{R}.$$

**Lemma.** *If (for a given  $\lambda$ )  $|\mathcal{E}^n(\lambda)| \geq c(\lambda) > 0$ ,  $n \geq 1$ , a sufficient condition for the limit relation*

$$\mathbb{E}e^{i\lambda Y^n} \rightarrow \mathbb{E}e^{i\lambda Y} \quad (33)$$

is that

$$\mathcal{E}^n(\lambda) \xrightarrow{P} \mathcal{E}(\lambda). \quad (34)$$

**PROOF.** Let

$$m^n(\lambda) = \frac{e^{i\lambda Y^n}}{\mathcal{E}^n(\lambda)}.$$

Then  $|m^n(\lambda)| \leq c^{-1}(\lambda) < \infty$ , and it is easily verified that

$$\mathbb{E}m^n(\lambda) = 1.$$

Hence by (34) and the Lebesgue dominated convergence theorem,

$$\begin{aligned} |\mathbb{E}e^{i\lambda Y^n} - \mathbb{E}e^{i\lambda Y}| &= |\mathbb{E}(e^{i\lambda Y^n} - \mathcal{E}(\lambda))| \\ &= |\mathbb{E}(m^n(\lambda)[\mathcal{E}^n(\lambda) - \mathcal{E}(\lambda)])| \leq c^{-1}(\lambda)\mathbb{E}|\mathcal{E}^n(\lambda) - \mathcal{E}(\lambda)| \rightarrow 0, \quad n \rightarrow \infty. \end{aligned}$$

**Remark.** It follows, from (33) and the hypothesis that  $\mathcal{E}^n(\lambda) \geq c(\lambda) > 0$ , that  $\mathcal{E}(\lambda) \neq 0$ . In fact, the conclusion of the lemma remains valid without the assumption that  $|\mathcal{E}^n(\lambda)| \geq c(\lambda) > 0$ , if restated in the form: if  $\mathcal{E}^n(\lambda) \xrightarrow{P} \mathcal{E}(\lambda)$  and  $\mathcal{E}(\lambda) \neq 0$ , then (33) holds (Problem 5).

**5. PROOF OF THEOREM 2.** (1) Let  $\varepsilon > 0$ ,  $\delta \in (0, \varepsilon)$ , and for simplicity let  $t = 1$ . Since

$$\max_{1 \leq k \leq n} |\xi_{nk}| \leq \varepsilon + \sum_{k=1}^n |\xi_{nk}| I(|\xi_{nk}| > \varepsilon)$$

and

$$\left\{ \sum_{k=1}^n |\xi_{nk}| I(|\xi_{nk}| > \varepsilon) > \delta \right\} \subseteq \left\{ \sum_{k=1}^n I(|\xi_{nk}| > \varepsilon) > \delta \right\},$$

we have

$$\begin{aligned} \mathbb{P}\left\{ \max_{1 \leq k \leq n} |\xi_{nk}| > \varepsilon + \delta \right\} &\leq \mathbb{P}\left\{ \sum_{k=1}^n I(|\xi_{nk}| > \varepsilon) > \delta \right\} \\ &= \mathbb{P}\left\{ \sum_{k=1}^n \int_{|x| > \varepsilon} d\mu_k^n > \delta \right\}. \end{aligned}$$

If (A) is satisfied, i.e.,

$$\mathbb{P}\left\{ \sum_{k=1}^n \int_{|x| > \varepsilon} dv_k^n > \delta \right\} \rightarrow 0$$

then (compare (14)) we also have

$$P\left\{\sum_{k=1}^n \int_{|x|>\varepsilon} d\mu_k^n > \delta\right\} \rightarrow 0.$$

Therefore (A)  $\Rightarrow$  (A\*).

Conversely, let

$$\sigma_n = \min\{k \leq n: |\xi_{nk}| \geq \varepsilon/2\},$$

supposing that  $\sigma_n = \infty$  if  $\max_{1 \leq k \leq n} |\xi_{nk}| < \varepsilon/2$ . By (A\*),  $\lim_n P(\sigma_n < \infty) = 0$ .

Now observe that, for every  $\delta \in (0, 1)$ , the sets

$$\left\{\sum_{k=1}^{n \wedge \sigma_n} I(|\xi_{nk}| \geq \varepsilon/2) > \delta\right\} \quad \text{and} \quad \left\{\max_{1 \leq k \leq n \wedge \sigma_n} |\xi_{nk}| \geq \frac{1}{2}\varepsilon\right\}$$

coincide, and by (A\*)

$$\sum_{k=1}^{n \wedge \sigma_n} I(|\xi_{nk}| \geq \varepsilon/2) = \sum_{k=1}^{n \wedge \sigma_n} \int_{|x| \geq \varepsilon/2} d\mu_k^n \xrightarrow{P} 0.$$

Therefore by (15)

$$\sum_{k=1}^{n \wedge \sigma_n} \int_{|x| \geq \varepsilon} dv_k^n \leq \sum_{k=1}^{n \wedge \sigma_n} \int_{|x| \geq \varepsilon/2} dv_k^n \xrightarrow{P} 0,$$

which, together with the property  $\lim_n P(\sigma_n < \infty) = 0$ , prove that (A\*)  $\Rightarrow$  (A).

(2) Again suppose that  $t = 1$ . Choose an  $\varepsilon \in (0, 1]$  and consider the square-integrable martingales

$$\Delta^n(\delta) = (\Delta_k^n(\delta), \mathcal{F}_k^n) \quad (1 \leq k \leq n),$$

with  $\delta \in (0, \varepsilon]$ . For the given  $\varepsilon \in (0, 1]$ , we have, according to (C),

$$\langle \Delta^n(\varepsilon) \rangle_n \xrightarrow{P} \sigma_1^2.$$

It is then easily deduced from (A) that for every  $\delta \in (0, \varepsilon]$

$$\langle \Delta^n(\delta) \rangle_n \xrightarrow{P} \sigma_1^2. \quad (35)$$

Let us show that it follows from (C\*), (A), and (A\*) that, for every  $\delta \in (0, \varepsilon]$ ,

$$[\Delta^n(\delta)]_n \xrightarrow{P} \sigma_1^2, \quad (36)$$

where

$$[\Delta^n(\delta)]_n = \sum_{k=1}^n [\xi_{nk} I(|\xi_{nk}| \leq \delta) - \int_{|x| \leq \delta} x dv_k^n]^2.$$

In fact, it is easily verified that by (A)

$$[\Delta^n(\delta)]_n - [\Delta^n(1)]_n \xrightarrow{P} 0. \quad (37)$$

But

$$\begin{aligned}
 & \left| \sum_{k=1}^n \left[ \xi_{nk} - \int_{|x| \leq 1} x d\nu_k^n \right]^2 - \sum_{k=1}^n \left[ \xi_{nk} I(|\xi_{nk}| \leq 1) - \int_{|x| \leq 1} x d\nu_k^n \right]^2 \right| \\
 & \leq \sum_{k=1}^n I(|\xi_{nk}| > 1) \left[ (\xi_{nk})^2 + 2|\xi_{nk}| \left| \int_{|x| \leq 1} x d(\mu_k^n - \nu_k^n) \right| \right] \\
 & \leq 5 \sum_{k=1}^n I(|\xi_{nk}| > 1) |\xi_{nk}|^2 \\
 & \leq 5 \max_{1 \leq k \leq n} |\xi_{nk}|^2 \sum_{k=1}^n \int_{|x| > 1} d\mu_k^n \rightarrow 0.
 \end{aligned} \tag{38}$$

Hence (36) follows from (37) and (38).

Consequently to establish the equivalence of (C) and (C\*) it is enough to establish that when (C) is satisfied (for a given  $\varepsilon \in (0, 1]$ ), then (C\*) is also satisfied for every  $a > 0$ :

$$\lim_{\delta \rightarrow 0} \overline{\lim}_n P\{|\Delta^n(\sigma)]_n - \langle \Delta^n(\delta) \rangle_n| > a\} = 0. \tag{39}$$

Let

$$m_k^n(\delta) = [\Delta^n(\delta)]_k - \langle \Delta^n(\delta) \rangle_k, \quad 1 \leq k \leq n.$$

The sequence  $m^n(\delta) = (m_k^n(\delta), \mathcal{F}_k^n)$  is a square-integrable martingale, and  $(m^n(\delta))^2$  is dominated (in the sense of the definition of §3 on p. 467) by the sequences  $[m^n(\delta)]$  and  $\langle m^n(\delta) \rangle$ .

It is clear that

$$\begin{aligned}
 [m^n(\delta)]_n &= \sum_{k=1}^n (\Delta m_k^n(\delta))^2 \leq \max_{1 \leq k \leq n} |\Delta m_k^n(\delta)| \{[\Delta^n(\delta)]_n + \langle \Delta^n(\delta) \rangle_n\} \\
 &\leq 3\delta^2 \{[\Delta^n(\delta)]_n + \langle \Delta^n(\delta) \rangle_n\}.
 \end{aligned} \tag{40}$$

Since  $[\Delta^n(\delta)]$  and  $\langle \Delta^n(\delta) \rangle$  dominate each other, it follows from (40) that  $(m^n(\delta))^2$  is dominated by the sequences  $6\delta^2[\Delta^n(\delta)]$  and  $6\delta^2\langle \Delta^n(\delta) \rangle$ .

Hence if (C) is satisfied, then for sufficiently small  $\delta$  (for example, for  $\delta < \frac{1}{6}b(\sigma_1^2 + 1)$ )

$$\overline{\lim}_n P(6\delta^2\langle \Delta^n(\delta) \rangle_n > b) = 0,$$

and hence, by the corollary to Theorem 2 of §3, we have (39). If (C\*) is also satisfied, then for the same values of  $\delta$ ,

$$\overline{\lim}_n P(6\delta^2[\Delta^n(\delta)]_k > b) = 0. \tag{41}$$

Since  $|\Delta[\Delta^n(\delta)]_k| \leq (2\delta)^2$ , the validity of (39) follows from (41) and another appeal to Theorem 2 of §3.

This completes the proof of Theorem 2.

**6. PROOF OF THEOREM 3.** On account of the Lindeberg condition (L), the equivalence of (C) and (1), and of (C\*) and (3), can be established by direct calculation (Problem 6).

**7. PROOF OF THEOREM 4.** Condition (A) follows from the Lindeberg condition (L). As for condition (B), it is sufficient to observe that when  $\xi^n$  is a martingale-difference, the variables  $B_t^n$  that appear in the canonical decomposition (9) can be represented in the form

$$B_t^n = - \sum_{k=0}^{[nt]} \int_{|x|>1} x \, d\nu_n^k.$$

Therefore  $B_t^n \xrightarrow{P} 0$  by the Lindeberg condition (L).

**8.** The fundamental theorem of the present section, namely Theorem 1, was proved under the hypothesis that the terms that are summed are uniformly asymptotically infinitesimal. It is natural to ask for conditions for the central limit theorem without such a hypothesis. For independent random variables, examples of such theorems are given by Theorem 1 (assuming finite second moments) or Theorem 5 (assuming finite first moments) from §4, Chapter III.

We quote (without proof) an analog of the first of these theorems, applicable only to sequences  $\xi^n = (\xi_{nk}, \mathcal{F}_k^n)$  that are square-integrable martingale differences.

Let  $\mathcal{F}_{nk}(x) = P(\xi_{nk} \leq x | \mathcal{F}_{k-1}^n)$  be a regular distribution function of  $\xi_{nk}$  with respect to  $\mathcal{F}_{k-1}^n$ , and let  $\Delta_{nk} = E(\xi_{nk}^2 | \mathcal{F}_{k-1}^n)$ .

**Theorem 5.** *If a square-integrable martingale-difference  $\xi_n = (\xi_{nk}, \mathcal{F}_k^n)$ ,  $0 \leq k \leq n$ ,  $n \geq 1$ ,  $\xi_{n0} = 0$ , satisfies the condition*

$$\sum_{k=0}^{[nt]} \Delta_{nk} \xrightarrow{P} \sigma_t^2, \quad 0 \leq \sigma_t^2 < \infty,$$

and for every  $\varepsilon > 0$

$$\sum_{k=0}^{[nt]} \int_{|x|>\varepsilon} |x| |\mathcal{F}_{nk}(x) - \Phi\left(\frac{x}{\sqrt{\Delta_{nk}}}\right)| dx \xrightarrow{P} 0,$$

then

$$X_t^n \xrightarrow{d} \mathcal{N}(0, \sigma_t^2).$$

## 9. PROBLEMS

1. Let  $\xi_n = \eta_n + \xi_n$ ,  $n \geq 1$ , where  $\eta_n \xrightarrow{d} \eta$  and  $\xi_n \xrightarrow{d} 0$ . Prove that  $\xi_n \xrightarrow{d} \eta$ .
2. Let  $(\xi_n(\varepsilon))$ ,  $n \geq 1$ ,  $\varepsilon > 0$ , be a family of random variables such that  $\xi_n(\varepsilon) \xrightarrow{P} 0$  for each  $\varepsilon > 0$  as  $n \rightarrow \infty$ . Using, for example, Problem 11 of §10, Chapter II, prove that there is a sequence  $\varepsilon_n \downarrow 0$  such that  $\xi_n(\varepsilon_n) \xrightarrow{P} 0$ .

3. Let  $(\alpha_k^n)$ ,  $1 \leq k \leq n$ ,  $n \geq 1$ , be a complex-valued random variable such that (P-a.s.)

$$\sum_{k=1}^n |\alpha_k^n| \leq C, \quad |\alpha_k^n| \leq a_n \downarrow 0.$$

Show that then (P-a.s.)

$$\lim_n \prod_{k=1}^n (1 + \alpha_k^n) \exp(-\alpha_k^n) = 1.$$

4. Prove the statement made in Remark 2 to Theorem 1.

5. Prove the statement made in Remark 1 to the lemma.

6. Prove Theorem 3.

7. Prove Theorem 5.

## §9. Discrete Version of Itô's Formula

1. In the stochastic analysis of Brownian motion and other related processes (for example, martingales, local martingales, semi-martingales) with continuous time *Itô's change-of-variables formula* plays a key role (see, for example, [J1], [L12]).

This section may be viewed as a prelude to Itô's formula for Brownian motion. In it, we present a discrete (in time) version of Itô's formula and show briefly how a corresponding formula for continuous time could be obtained using a limiting procedure.

2. Let  $X = (X_n)_{0 \leq n \leq N}$  and  $Y = (Y_n)_{0 \leq n \leq N}$  be two sequences of random variables on the probability space  $(\Omega, \mathcal{F}, P)$ ,  $X_0 = Y_0 = 0$  and

$$[X, Y] = ([X, Y]_n)_{0 \leq n \leq N},$$

where

$$[X, Y]_n = \sum_{i=1}^n \Delta X_i \Delta Y_i \quad (1)$$

is the *square covariance* of  $(X_0, X_1, \dots, X_n)$  and  $(Y_0, Y_1, \dots, Y_n)$  (see §1). Also, suppose that  $F = F(x)$  is an absolutely continuous function,

$$F(x) = F(0) + \int_0^x f(y) dy, \quad (2)$$

where  $f = f(y)$ ,  $y \in \mathbb{R}$  is a Borel function such that

$$\int_{|y| \leq c} |f(y)| dy < \infty, \quad c > 0.$$

The change-of-variables formula in which we are interested concerns the

possibility of representing the sequence

$$F(X) = (F(X_n))_{0 \leq n \leq N} \quad (3)$$

in terms of 'natural' functional from the sequence  $X$ . Of course, this requires explanation and will be explained here.

Given the function  $f = f(x)$ , we form the square covariance  $[X, f(X)]$  as follows:

$$\begin{aligned} [X, f(X)]_n &= \sum_{k=1}^n \Delta f(X_k) \Delta X_k \\ &= \sum_{k=1}^n (f(X_k) - f(X_{k-1}))(X_k - X_{k-1}). \end{aligned} \quad (4)$$

We introduce two 'discrete integrals' (cf. Definition 5 of §1):

$$I_n(X, f(X)) = \sum_{k=1}^n f(X_{k-1}) \Delta X_k, \quad (5)$$

$$\tilde{I}_n(X, f(X)) = \sum_{k=1}^n f(X_k) \Delta X_k. \quad (6)$$

Then

$$[X, f(X)]_n = \tilde{I}_n(X, f(X)) - I_n(X, f(X)). \quad (7)$$

For fixed  $N$ , we introduce a new (reversed) sequence  $\tilde{X} = (\tilde{X}_n)_{0 \leq n \leq N}$  with

$$\tilde{X}_n = X_{N-n}. \quad (8)$$

Then clearly,

$$\tilde{I}_N(X, f(X)) = -I_N(\tilde{X}, f(\tilde{X}))$$

and analogously,

$$\tilde{I}_n(X, f(X)) = -\{I_N(\tilde{X}, f(\tilde{X})) - I_{N-n}(\tilde{X}, f(\tilde{X}))\}$$

(we set  $I_0 = \tilde{I}_0 = 0$ ).

Thus,

$$[X, f(X)]_N = -\{I_N(\tilde{X}, f(\tilde{X})) + I_N(\tilde{X}, f(X))\}$$

and for  $0 < n < N$  we have:

$$\begin{aligned} [X, f(X)]_n &= -\{I_N(\tilde{X}, f(\tilde{X})) - I_{N-n}(\tilde{X}, f(\tilde{X}))\} - I_n(X, f(X)) \\ &= -\left\{ \sum_{k=N-n+1}^N f(\tilde{X}_k) \Delta \tilde{X}_k + \sum_{k=1}^n f(X_k) \Delta X_k \right\}. \end{aligned} \quad (9)$$

**Remark 1.** We note that the structures of the right-hand sides of (7) and (9) are different. Equation (7) contains two different forms of "discrete integral." The integral  $I_n(X, f(X))$  is a "forward integral," while  $\tilde{I}_n(X, f(X))$  is a "backward integral." In (9), both integrals are "forward integrals," over two differ-



ent sequences  $X$  and  $\tilde{X}$ .

3. Since for any function  $g = g(x)$

$$g(X_{k-1}) + \frac{1}{2}[g(X_k) - g(X_{k-1})] - \frac{1}{2}[g(X_k) + g(X_{k-1})] = 0,$$

it is clear that

$$\begin{aligned} F(X_n) &= F(X_0) + \sum_{k=1}^n g(X_{k-1})\Delta X_k + \frac{1}{2}[X, g(X)]_n \\ &\quad + \sum_{k=1}^n \left\{ (F(X_k) - F(X_{k-1})) - \frac{g(X_{k-1}) + g(X_k)}{2} \Delta X_k \right\}. \end{aligned} \quad (10)$$

In particular, if  $g(x) = f(x)$ , where  $f(x)$  is the function of (2), then

$$F(X_n) = F(X_0) + I_n(X, f(X)) + \frac{1}{2}[X, f(X)]_n + R_n(X, f(X)), \quad (11)$$

where

$$R_n(X, f(X)) = \sum_{k=1}^n \int_{X_{k-1}}^{X_k} \left[ f(x) - \frac{f(X_{k-1}) + f(X_k)}{2} \right] dx. \quad (12)$$

From analysis, it is well known that if the function  $f''(x)$  is continuous, then the following formula ("trapezoidal rule") holds:

$$\begin{aligned} \int_a^b \left[ f(x) - \frac{f(a) + f(b)}{2} \right] dx &= \int_a^b (x-a)(x-b) \frac{f''(\xi(x))}{2!} dx \\ &= \frac{(b-a)^3}{2} \int_0^1 x(x-1) f''(\xi(a + x(b-a))) dx \\ &= \frac{(b-a)^3}{2} f''(\xi(a + \bar{x}(b-a))) \int_0^1 x(x-1) dx \\ &= -\frac{(b-a)^3}{12} f''(\eta), \end{aligned}$$

where  $\xi(x)$ ,  $\bar{x}$  and  $\eta$  are "intermediate" points in the interval  $[a, b]$ .

Thus, in (11)

$$R_n(X, f(X)) = -\frac{1}{12} \sum_{k=1}^n f''(\eta_k) (\Delta X_k)^3$$

where  $X_{k-1} \leq \eta_k \leq X_k$ , whence

$$|R_n(X, f(X))| \leq \frac{1}{12} \sup f''(\eta) \sum_{k=1}^n |\Delta X_k|^3$$

where the supremum is taken over all  $\eta$  such that

$$\min(X_0, X_1, \dots, X_n) \leq \eta \leq \max(X_0, X_1, \dots, X_n).$$

We shall refer to formula (11) as the *discrete analogue of Itô's formula*. We note that the right-hand side of this formula contains the following three

'natural' ingredients: 'the discrete integral'  $I_n(X, f(X))$ , the square covariance  $[X, f(X)]$ , and the 'residual' term  $R_n(X, f(X))$ .

4. **EXAMPLE 1.** If  $f(x) = a + bx$ , then  $R_n(X, f(X)) = 0$  and formula (11) takes the following form:

$$F(X_n) = F(X_0) + I_n(X, f(X)) + \frac{1}{2}[X, f(X)]_n. \quad (13)$$

**EXAMPLE 2.** Let

$$f(x) = \text{sign } x = \begin{cases} 1, & x > 0 \\ 0, & x = 0 \\ -1, & x < 0. \end{cases}$$

Then  $F(x) = |x|$ .

Let  $X_k = S_k$ , where

$$S_k = \xi_1 + \xi_2 + \cdots + \xi_k,$$

where  $\xi_1, \xi_2, \dots$  are independent Bernoulli random variables taking values  $\pm 1$  with probability  $1/2$ . If we also set  $S_0 = 0$ , we obtain

$$|S_n| = \sum_{k=1}^n \text{sign } S_{k-1} \Delta S_k + N_n, \quad (14)$$

where  $S_0 = 0$  and

$$N_n = \{0 \leq k < n: S_k = 0\}.$$

We note that the sequence of discrete integrals  $(\sum_{k=1}^n \text{sign } S_{k-1} \Delta S_k)_{n \geq 1}$  forms a martingale and therefore,

$$E|S_n| = EN_n.$$

It is not difficult to show that as  $n \rightarrow \infty$ ,

$$E|S_n| \sim \sqrt{\frac{n}{2\pi}}, \quad (15)$$

thus, application of Itô's discrete formula gives the asymptotic behavior of the average number of changes of sign on the path of the random walk

$$(S_k)_{k < n}: EN_n \sim \sqrt{\frac{n}{2\pi}}.$$

5. **Remarks.** Let  $B = (Bt)_{0 \leq t \leq 1}$  be a Brownian movement and  $X_k = B_{k/n}$ ,  $k = 0, 1, \dots, n$ . Then application of formula (11) leads to the following result:

$$F(B_1) = F(B_0) + \sum_{k=1}^n f(B_{(k-1)/n}) \Delta B_{k/n} + \frac{1}{2}[f(B_{\cdot/n}), B_{\cdot/n}]_n + R_n(B_{\cdot/n}, f(B_{\cdot/n})). \quad (16)$$

It is known from the stochastic calculus of Brownian motion that

$$\sum_{k=1}^n |B_{k/n} - B_{(k-1)/n}|^3 \xrightarrow{P} 0 \quad (17)$$

and if  $f = f(x) \in L^2_{\text{loc}}$  (i.e.,  $\int_{|x| \leq k} f^2(x) dx < \infty$  for any  $k > 0$ ), then the limit

$$L^2\text{-}\lim_n \sum_{k=0}^n f(B_{(k-1)/n}) \Delta B_{k/n}, \quad (18)$$

exists and is denoted by

$$\int_0^1 f(B_s) dB_s$$

and is called Itô's *stochastic integral* of  $f(B_s)$  with respect to Brownian motion ([J1], [L12]). In addition, if the function  $f(x)$  has a second derivative and  $|f''(x)| \leq C$ , then from (9), we obtain that  $P\text{-}\lim R_n(B_{\cdot/n}, f(B_{\cdot/n})) = 0$ . Thus, from (8) and (10), it follows that the limit

$$P\text{-}\lim [f(B_{\cdot/n}), B_{\cdot/n}]_n$$

exists and is denoted by

$$[f(B), B]_1,$$

and that the following formula holds (P-a.s.)

$$F(B_1) = F(B_0) + \int_0^1 f(B_s) dB_s + \frac{1}{2} [f(B), B]_1. \quad (19)$$

It can be shown that, for the given smoothness assumptions

$$[f(B), B]_1 = \int_0^1 f'(B_s) ds, \quad (20)$$

which leads to the known *formula of Itô for Brownian motion*:

$$F(B_1) = F(B_0) + \int_0^1 f(B_s) dB_s + \frac{1}{2} \int_0^1 f'(B_s) ds. \quad (21)$$

## 6. PROBLEMS

1. Show that formula (14) is true.
2. Establish that the property (16) is true.
3. Prove formula (15).

## §10. Applications to Calculations of the Probability of Ruin in Insurance

1. The material studied in the present section is a good illustration of the fact that the theory of martingales provides a quick and simple way of calculating the *risk* of an insurance company.

We shall assume that the evolution of the capital  $X = (X_t)_{t \geq 0}$  of a certain insurance company takes place in a probability space  $(\Omega, \mathcal{F}, P)$  as follows.

The initial capital is  $X_0 = u > 0$ . Insurance payments arrive continuously at a constant rate  $c > 0$  (in time  $\Delta t$  the amount arriving is  $c\Delta t$ ) and claims are received at random times  $T_1, T_2, \dots$  ( $0 < T_1 < T_2 < \dots$ ) where the amounts to be paid out at these times are described by a nonnegative random variables  $\xi_1, \xi_2, \dots$ .

Thus, taking into account receipts and claims, the capital  $X_t$  at time  $t > 0$  is determined by the formula

$$X_t = u + ct - S_t, \quad (1)$$

where

$$S_t = \sum_{i \geq 1} \xi_i I(T_i \leq t). \quad (2)$$

We denote

$$T = \inf\{t \geq 0 : X_t \leq 0\}$$

the first time at which the insurance company's capital becomes less than or equal to zero ('time of ruin'). Of course, if  $X_t > 0$  for all  $t \geq 0$ , then the time  $T$  is given to be equal to  $+\infty$ .

One of the main questions relating to the operation of an insurance company is the calculation of the *probability of ruin*,  $P(T < \infty)$ , and the *probability of ruin before time  $t$* ,  $P(T \leq t)$  (inclusively).

**2.** To calculate these probabilities we assume we are in the framework of the classical Cramer–Lundberg model characterized by the following assumptions:

**A** The times  $T_1, T_2, \dots$  at which claims are received are such that the variables ( $T_0 \equiv 0$ )

$$\sigma_i = T_i - T_{i-1}, \quad i \geq 1$$

are independent, identically-distributed random variables having an exponential distribution with density  $\lambda e^{-\lambda t}$ ,  $t \geq 0$  (see Table 2, §3, Chapter II).

**B** The random variables  $\xi_1, \xi_2, \dots$  are independent and identically distributed with distribution function  $F(x) = P(\xi_1 \leq x)$  such that  $F(0) = 0$ ,  $\mu = \int_0^\infty x dF(x) < \infty$ .

**C** The sequences  $(T_1, T_2, \dots)$  and  $(\xi_1, \xi_2, \dots)$  are independent sequences (in the sense of Definition 6, §5, Chapter II).

We denote the process of the number of claims by  $N = (N_t)_{t \geq 0}$ , i.e., set

$$N_t = \sum_{i \geq 1} I(T_i \leq t). \quad (3)$$

It is clear that this process has a piecewise-constant trajectory with jumps by a unit value at times  $T_1, T_2, \dots$  and with value  $N_0 = 0$ .

Since

$$\{T_k > t\} = \{\sigma_1 + \cdots + \sigma_k > t\} = \{N_t < k\},$$

under the assumption A we find that, according to Problem 6, §2, Chapter II,

$$P(N_t < k) = P(\sigma_1 + \cdots + \sigma_k > t) = \sum_{i=0}^{k-1} e^{-\lambda t} \frac{(\lambda t)^i}{i!}.$$

Whence

$$P(N_t = k) = e^{-\lambda t} \frac{(\lambda t)^k}{k!}, \quad k = 0, 1, \dots, \quad (4)$$

i.e., the random variable  $N_t$  has a Poisson distribution (see Table 1. §3, Chapter II) with parameter  $\lambda t$ . Here,  $EN_t = \lambda t$ .

The so-called *Poisson process*  $N = (N_t)_{t \geq 0}$  constructed in this way is (together with the *Brownian motion*; §13, Chapter II) an example of another classical random (stochastic) process with continuous time. Like the Brownian motion, the Poisson process is a process with independent increments (see §13, Chapter II), where, for  $s < t$ , the increments  $N_t - N_s$  have a Poisson distribution with parameter  $\lambda(t - s)$  (these properties are not difficult to derive from assumption A and the explicit construction (3) of the process  $N$ ).

3. From assumption C we find that

$$\begin{aligned} E(X_t - X_s) &= ct - ES_t = ct - E \sum_i \xi_i I(T_i \leq t) \\ &= ct - \sum_i E \xi_i E I(T_i \leq t) = ct - \mu \sum_i P(T_i \leq t) \\ &= ct - \mu \sum_i P(N_t \geq i) = ct - \mu EN_t = t(c - \lambda\mu). \end{aligned}$$

Thus, we see that, in the case under consideration, a natural requirement for an insurance company to operate with a clear profit (i.e.  $E(X_t - X_0) > 0$ ,  $t > 0$ ) is that

$$c > \lambda\mu. \quad (5)$$

In the following analysis, an important role is played by the function

$$h(z) = \int_0^\infty (e^{zx} - 1) dF(x), \quad z \geq 0, \quad (6)$$

which is equal to  $\hat{F}(-z) - 1$ , where

$$\hat{F}(s) = \int_0^\infty e^{-sx} dF(x)$$

is the Laplace-Stieltjes transformation ( $s$  is a complex number).

Denoting

$$g(z) = \lambda h(z) - cz \quad \xi_0 \equiv 0,$$

we find that for  $r > 0$  with  $h(r) < \infty$ ,

$$\begin{aligned} \mathbb{E}e^{-r(X_t - X_0)} &= \mathbb{E}e^{-r(X_t - u)} = e^{-rct} \cdot \mathbb{E}e^{r \sum_{i=0}^{N_t} \xi_i} \\ &= e^{-rct} \sum_{n=0}^{\infty} \mathbb{E}e^{r \sum_{i=0}^{N_t} \xi_i} \mathbb{P}(N_t = n) \\ &= e^{-rct} \sum_{n=0}^{\infty} (1 + h(r))^n \frac{e^{-\lambda t} (\lambda t)^n}{n!} \\ &= e^{-rct} \cdot e^{\lambda th(r)} = e^{t[\lambda h(r) - cr]} = e^{tg(r)}. \end{aligned}$$

Analogously, it can be shown that for any  $s < t$

$$\mathbb{E}e^{-r(X_t - X_s)} = e^{(t-s)g(r)}. \quad (7)$$

Let  $\mathcal{G}_t = \sigma(X_s, s \leq t)$ . Since the process  $X = (X_t)_{t \geq 0}$  is a process with independent increments (Problem 2) (P-a.s.)

$$\mathbb{E}(e^{-r(X_t - X_s)} | \mathcal{G}_s) = \mathbb{E}e^{-r(X_t - X_s)} = e^{(t-s)g(r)},$$

then (P-a.s.)

$$\mathbb{E}(e^{-rX_t - tg(r)} | \mathcal{G}_s) = e^{-rX_s - sg(r)}. \quad (8)$$

Denoting

$$Z = e^{-rX_t - tg(r)}, \quad t \geq 0 \quad (9)$$

we see that property (8) rewritten in the form

$$\mathbb{E}(Z_t | \mathcal{G}_s) = Z_s, \quad s \leq t \quad (10)$$

is a continuous analogue of the *martingale property* (2) of Definition 1 of §1.

By analogy with Definition 3 of §1, we shall say that the random variable  $\tau = r(\omega)$  with values in  $[0, +\infty]$  is a *Markov time* relative to the system of  $\sigma$ -algebras  $(\mathcal{G}_t)_{t \geq 0}$  if for each  $t \geq 0$  the set

$$\{\tau(\omega) \leq t\} \in \mathcal{G}_t.$$

The process  $Z = (Z_t)_{t \geq 0}$  is nonnegative with  $\mathbb{E}Z_t = e^{-ru} < \infty$ . Thus, by analogy with Definition 1 of §1, the process  $Z = (Z_t)_{t \geq 0}$  with *continuous time* is a martingale.

It turns out that for martingales with continuous time, Theorem 1 of §2 remains valid (with self-evident changes to the notation). In particular,

$$\mathbb{E}Z_{t \wedge r} = \mathbb{E}Z_0 \quad (11)$$

for any Markov time  $r$  (see for example [L5], §2, Chapter 3).

By virtue of (9) we find from (11) that for time  $r = T$

$$\begin{aligned} e^{-ru} &= \mathbb{E}e^{-rX_{t \wedge T} - (t \wedge T)g(r)} \\ &\geq \mathbb{E}[e^{-rX_{t \wedge T} - (t \wedge T)g(r)} | T \leq t] \mathbb{P}(T \leq t) \\ &= \mathbb{E}[e^{-rX_T - Tg(r)} | T \leq t] \mathbb{P}(T \leq t) \\ &\geq \mathbb{E}[e^{-Tg(r)} | T \leq t] \mathbb{P}(T \leq t) \geq \min_{0 \leq s \leq t} e^{-sg(r)} \mathbb{P}(T \leq t). \end{aligned}$$

Moreover,

$$P(T \leq t) \leq \frac{e^{-ru}}{\min_{0 \leq s \leq t} e^{-sg(r)}} = e^{-ru} \max_{0 \leq s \leq t} e^{sg(r)}. \quad (12)$$

Let us consider the function

$$g(r) = \lambda h(r) - cr$$

in more detail. Clearly,  $g(0) = 0$ ,  $g'(0) = \lambda\mu - c < 0$  (by virtue of (5)) and  $g''(r) = \lambda h''(r) \geq 0$ . Thus, there exists a unique positive value  $r = R$  with  $g(R) = 0$ .

Noting that for  $r > 0$

$$\begin{aligned} \int_0^\infty e^{rx}(1 - F(x)) dx &= \int_0^\infty \int_x^\infty e^{rx} dF(y) dx \\ &= \int_0^\infty \left( \int_0^y e^{rx} dx \right) dF(y) \\ &= \frac{1}{r} \int_0^\infty (e^{ry} - 1) dF(y) = \frac{1}{r} h(r), \end{aligned}$$

$R$  may be asserted to be the (unique) root of the equation

$$\frac{\lambda}{c} \int_0^\infty e^{rx}(1 - F(x)) dx = 1. \quad (13)$$

Let us set  $r = R$  in (12). Then we obtain, for any  $t > 0$ ,

$$P(T \leq t) \leq e^{-Ru} \quad (14)$$

whence

$$P(T < \infty) \leq e^{-Ru}. \quad (15)$$

Moreover, we prove the following

**Theorem.** Suppose that in the Cramer-Lundberg model assumptions A, B, C and property (5) are satisfied (i.e.,  $\lambda\mu < c$ ). Then the bound of (15) holds for the probability of ruin  $P(T < \infty)$ , where  $R$  is the positive root of equation (13).

**4. Remark.** The above discussion could be greatly simplified from a stochastic point of view if we assumed a geometric distribution for the  $\sigma_i$  ( $P(\sigma_i = k) = q^{k-1}p$ ,  $k = 1, 2, \dots$ ) instead of an exponential distribution. In this case, all the time variables ( $T_i$ ,  $T$ ) would take discrete values, it would not be necessary to call upon the results of the theory of martingales for *continuous* time and the whole study could strictly be carried out based *only* on the 'discrete' (in time) methods from the theory of martingales studied in the present chapter.

However, we have turned our attention to the "continuous" (in time) scheme to illustrate both the method and the usefulness of the general theory of martingales for the case of *continuous time*, based on the given example.

**5. PROBLEMS**

1. Prove that the process  $N = (N_t)_{t \geq 0}$  (under assumption A) is a process with independent increments, where  $N_t - N_s$  has a Poisson distribution with parameter  $\lambda(t - s)$ .
2. Prove that the process  $X = (X_t)_{t \geq 0}$  is a process with independent increments.
3. Consider the problem of determining the probability of ruin  $P(T < \infty)$  assuming that the variables  $\sigma_i$  have a geometric (rather than exponential) distribution ( $P(\sigma_i = k) = q^{k-1}p$ ,  $k = 1, 2, \dots$ ).



## CHAPTER VIII

# Sequences of Random Variables That Form Markov Chains

### §1. Definitions and Basic Properties

1. In Chapter I (§12), for finite probability spaces, we took the basic idea to be that of *Markov dependence* between random variables. We also presented a variety of examples and considered the simplest regularities that are possessed by random variables that are connected by a Markov chain.

In the present chapter we give a general definition of a stochastic sequence of random variables that are connected by Markov dependence, and devote our main attention to the asymptotic properties of Markov chains with countable state spaces.

2. Let  $(\Omega, \mathcal{F}, P)$  be a probability space with a distinguished nondecreasing family  $(\mathcal{F}_n)$  of  $\sigma$ -algebras,  $\mathcal{F}_0 \subseteq \mathcal{F}_1 \subseteq \dots \subseteq \mathcal{F}$ .

**Definition.** A stochastic sequence  $X = (X_n, \mathcal{F}_n)$  is called a *Markov chain* (with respect to the measure  $P$ ) if

$$P\{X_n \in B | \mathcal{F}_m\} = P\{X_n \in B | X_m\} \quad (P\text{-a.s.}) \quad (1)$$

for all  $n \geq m \geq 0$  and all  $B \in \mathcal{B}(R)$ .

Property (1), the *Markov property*, can be stated in a number of ways. For example, it is equivalent to saying that

$$E[g(X_n) | \mathcal{F}_m] = E[g(X_n) | X_m] \quad (P\text{-a.s.}) \quad (2)$$

for every bounded Borel function  $g = g(x)$ .

Property (1) is also equivalent to the statement that, for a given "present"  $X_m$ , the "future"  $F$  and the "past"  $P$  are independent, i.e.

$$P(FP|X_m) = P(F|X_m)P(P|X_m), \quad (3)$$

where  $F \in \sigma\{\omega: X_i, i \geq m\}$ , and  $B \in \mathcal{F}_n, n \leq m$ .

In the special case when

$$\mathcal{F}_n = \mathcal{F}_n^X = \sigma\{\omega: X_0, \dots, X_n\}$$

and the stochastic sequence  $X = (X_n, \mathcal{F}_n^X)$  is a Markov chain, we say that the sequence  $\{X_n\}$  itself is a Markov chain. It is useful to notice that if  $X = \{X_n, \mathcal{F}_n\}$  is a Markov chain, then  $(X_n)$  is also a Markov chain.

**Remark.** It was assumed in the definition that the variables  $X_m$  are real-valued. In a similar way, we can also define Markov chains for the case when  $X_n$  takes values in some measurable space  $(E, \mathcal{E})$ . In this case, if all singletons are measurable, the space is called a *phase space*, and we say that  $X = (X_n, \mathcal{F}_n)$  is a Markov chain with values in the phase space  $(E, \mathcal{E})$ . When  $E$  is finite or countably infinite (and  $\mathcal{E}$  is the  $\sigma$ -algebra of all its subsets) we say that the Markov chain is *discrete*. In turn, a discrete chain with a finite phase space is called a *finite chain*.

The theory of finite Markov chains, as presented in §12, Chapter I, shows that a fundamental role is played by the one-step transition probabilities  $P(X_{n+1} \in B|X_n)$ . By Theorem 3, §7, Chapter II, there are functions  $P_{n+1}(x; B)$ , the *regular conditional probabilities*, which (for given  $x$ ) are measures on  $(R, \mathcal{B}(R))$ , and (for given  $B$ ) are measurable functions of  $x$ , such that

$$P(X_{n+1} \in B|X_n) = P_{n+1}(X_n; B) \quad (\text{P-a.s.}) \quad (4)$$

The functions  $P_n = P_n(x, B)$ ,  $n \geq 0$ , are called *transition functions*, and in the case when they coincide ( $P_1 = P_2 = \dots$ ), the corresponding Markov chain is said to be *homogeneous* (in time).

From now on we shall consider only homogeneous Markov chains, and the transition function  $P_1 = P_1(x, B)$  will be denoted simply by  $P = P(x, B)$ .

Besides the transition function, an important probabilistic property of a Markov chain is the *initial distribution*  $\pi = \pi(B)$ , that is, the probability distribution defined by  $\pi(B) = P(X_0 \in B)$ .

The set of pairs  $(\pi, P)$ , where  $\pi$  is an initial distribution and  $P$  is a transition function, completely determines the probabilistic properties of  $X$ , since every finite-dimensional distribution can be expressed (Problem 2) in terms of  $\pi$  and  $P$ : for every  $n \geq 0$  and  $A \in \mathcal{B}(R^{n+1})$

$$\begin{aligned} & P\{(X_0, \dots, X_n) \in A\} \\ &= \int_R \pi(dx_0) \int_R P(x_0; dx_1) \cdots \int_R I_A(x_0, \dots, x_n) P(x_{n-1}; dx_n). \end{aligned} \quad (5)$$

We deduce, by a standard limiting process, that for any  $\mathcal{B}(R^{n+1})$ -measurable function  $g(x_0, \dots, x_n)$ , either of constant sign or bounded,

$$\begin{aligned} & \mathbb{E}g(X_0, \dots, X_n) \\ &= \int_R \pi(dx_0) \int_R P(x_0; dx_1) \cdots \int_R g(x_0, \dots, x_n) P(x_{n-1}; dx_n). \end{aligned} \quad (6)$$

3. Let  $P^{(n)} = P^{(n)}(x; B)$  denote a regular variant of the  $n$ -step transition probability:

$$P(X_n \in B | X_0) = P^{(n)}(X_0; B) \quad (\mathbf{P}\text{-a.s.}) \quad (7)$$

It follows at once from the Markov property that for all  $k$  and  $l$ , ( $k, l \geq 1$ ),

$$P^{(k+l)}(X_0; B) = \int_R P^{(k)}(X_0; dy) P^{(l)}(y; B) \quad (\mathbf{P}\text{-a.s.}) \quad (8)$$

It does *not* follow, of course, that for *all*  $x \in R$

$$P^{(k+l)}(x; B) = \int_R P^{(k)}(x; dy) P^{(l)}(y; B). \quad (9)$$

It turns out, however, that regular variants of the transition probabilities *can be chosen* so that (9) will be satisfied for *all*  $x \in R$  (see the discussion in the historical and bibliographical notes, p. 559).

Equation (9) is the *Kolmogorov–Chapman* equation (compare (I.12.13)) and is the starting point for the study of the probabilistic properties of Markov chains.

4. It follows from our discussion that with every Markov chain  $X = (X_n, \mathcal{F}_n)$ , defined on  $(\Omega, \mathcal{F}, \mathbf{P})$  there is associated a set  $(\pi, \mathbf{P})$ . It is natural to ask what properties a set  $(\pi, \mathbf{P})$  must have in order for  $\pi = \pi(B)$  to be a probability distribution on  $(R, \mathcal{B}(R))$  and for  $P = P(x; B)$  to be a function that is measurable in  $x$  for given  $B$ , and a probability measure on  $B$  for every  $x$ , so that  $\pi$  will be the initial distribution, and  $P$  the transition function, for some Markov chain. As we shall now show, no additional hypotheses are required.

In fact, let us take  $(\Omega, \mathcal{F})$  to be the measurable space  $(R^\infty, \mathcal{B}(R^\infty))$ . On the sets  $A \in \mathcal{B}(R^{n+1})$  we define a probability measure by the right-hand side of formula (5). It follows from §9, Chapter II, that a probability measure  $\mathbf{P}$  exists on  $(R^\infty, \mathcal{B}(R^\infty))$  for which

$$\begin{aligned} & \mathbf{P}\{\omega: (x_0, \dots, x_n) \in A\} \\ &= \int_R \pi(dx_0) \int_R P(x_0; dx_1) \cdots \int_R I_A(x_0, \dots, x_n) P(x_{n-1}; dx_n). \end{aligned} \quad (10)$$

Let us show that if we put  $X_n(\omega) = x_n$  for  $\omega = (x_0, x_1, \dots)$ , the sequence  $X = (X_n)_{n \geq 0}$  will constitute a Markov chain (with respect to the measure  $\mathbf{P}$  just constructed).

In fact, if  $B \in \mathcal{B}(R)$  and  $C \in \mathcal{B}(R^{n+1})$ , then

$$\begin{aligned} P\{X_{n+1} \in B, (X_0, \dots, X_n) \in C\} \\ &= \int_R \pi(dx_0) \int_R P(x_0; dx_1) \cdots \int_R I_B(x_{n+1}) I_C(x_0, \dots, x_n) P(x_n; dx_{n+1}) \\ &= \int_R \pi(dx_0) \int_R P(x_0; dx_1) \cdots \int_R P(x_n; B) I_C(x_0, \dots, x_n) P(x_{n-1}; dx_n) \\ &= \int_{\{\omega: (X_0, \dots, X_n) \in C\}} P(X_n; B) dP, \end{aligned}$$

whence (P-a.s.)

$$P\{X_{n+1} \in B | X_0, \dots, X_n\} = P(X_n; B). \quad (11)$$

Similarly we can verify that (P-a.s.)

$$P\{X_{n+1} \in B | X_n\} = P(X_n; B). \quad (12)$$

Equation (1) now follows from (11) and (12). It can be shown in the same way that for every  $k \geq 1$  and  $n \geq 0$ ,

$$P\{X_{n+k} \in B | X_0, \dots, X_n\} = P\{X_{n+k} \in B | X_n\} \quad (\text{P-a.s.}).$$

This implies the homogeneity of Markov chains.

The Markov chain  $X = (X_n)$  that we have constructed is known as the Markov chain generated by  $(\pi, P)$ . To emphasize that the measure  $P$  on  $(R^\infty, \mathcal{B}(R^\infty))$  has precisely the initial distribution  $\pi$ , it is often denoted by  $P_\pi$ .

If  $\pi$  is concentrated at the single point  $x$ , we write  $P_x$  instead of  $P_\pi$ , and the corresponding Markov chain is called the chain *generated by the point  $x$*  (since  $P_x(X_0 = x) = 1$ ).

Consequently, each transition function  $P = P(x, B)$  is in fact connected with the *whole family of probability measures*  $\{P_x, x \in R\}$ , and therefore with the whole family of Markov chains that arise when the sequence  $(X_n)_{n \geq 0}$  is considered with respect to the measures  $P_x, x \in R$ . From now on, we shall use the phrase "Markov chain with given transition function" to mean the family of Markov chains in the sense just described.

We observe that the measures  $P_\pi$  and  $P_x$  constructed from the transition function  $P = P(x, B)$  are *consistent* in the sense that, when  $A \in \mathcal{B}(R^\infty)$ ,

$$P_\pi\{(X_0, X_1, \dots) \in A | X_0 = x\} = P_x\{(X_0, X_1, \dots) \in A\} \quad (\pi\text{-a.s.}) \quad (13)$$

and

$$P_\pi\{(X_0, X_1, \dots) \in A\} = \int_R P_x\{(X_0, X_1, \dots) \in A\} \pi(dx). \quad (14)$$

5. Let us suppose that  $(\Omega, \mathcal{F}) = (R^\infty, \mathcal{B}(R^\infty))$  and that we are considering a sequence  $X = (X_n)$  that is defined coordinate-wise, that is,  $X_n(\omega) = x_n$  for  $\omega = (x_0, x_1, \dots)$ . Also let  $\mathcal{F}_n = \sigma\{\omega: X_0, \dots, X_n\}$ ,  $n \geq 0$ .

Let us define the *shifting operators*  $\theta_n$ ,  $n \geq 0$ , on  $\Omega$  by the equation

$$\theta_n(x_0, x_1, \dots) = (x_n, x_{n+1}, \dots),$$

and let us define, for every random variable  $\eta = \eta(\omega)$ , the random variables  $\theta_n \eta$  by putting

$$(\theta_n \eta)(\omega) = \eta(\theta_n \omega).$$

In this notation, the Markov property of homogeneous chains can (Problem 1) be given the following form: For every  $\mathcal{F}$ -measurable  $\eta = \eta(\omega)$ , every  $n \geq 0$ , and  $B \in \mathcal{B}(R)$ ,

$$P\{\theta_n \eta \in B | \mathcal{F}_n\} = P_{X_n}\{\eta \in B\} \quad (\text{P-a.s.}) \quad (15)$$

This form of the Markov property allows us to give the following important generalization: (15) remains valid if we replace  $n$  by stopping times  $\tau$ .

**Theorem.** Let  $X = (X_n)$  be a homogeneous Markov chain defined on  $(R^\infty, \mathcal{B}(R^\infty), P)$  and let  $\tau$  be a stopping time. Then the following strong Markov property is valid:

$$P\{\theta_\tau \eta \in B | \mathcal{F}_\tau\} = P_{X_\tau}\{\eta \in B\} \quad (\text{P-a.s.}) \quad (16)$$

PROOF. If  $A \in \mathcal{F}_\tau$  then

$$\begin{aligned} P\{\theta_\tau \eta \in B, A\} &= \sum_{n=0}^{\infty} P\{\theta_\tau \eta \in B, A, \tau = n\} \\ &= \sum_{n=0}^{\infty} P\{\theta_n \eta \in B, A, \tau = n\}. \end{aligned} \quad (17)$$

The events  $A \cap \{\tau = n\} \in \mathcal{F}_n$ , and therefore

$$\begin{aligned} P\{\theta_n \eta \in B, A \cap \{\tau = n\}\} &= \int_{A \cap \{\tau = n\}} P\{\theta_n \eta \in B | \mathcal{F}_n\} dP \\ &= \int_{A \cap \{\tau = n\}} P_{X_n}\{\eta \in B\} dP = \int_{A \cap \{\tau = n\}} P_{X_\tau}\{\eta \in B\} dP, \end{aligned}$$

which, with (17), establishes (16).

**Corollary.** If  $\sigma$  is a stopping time such that  $P(\sigma \geq \tau) = 1$  and  $\sigma$  is  $\mathcal{F}_\tau$ -measurable, then

$$P\{X_\sigma \in B, \sigma < \infty | \mathcal{F}_\tau\} = P_{X_\tau}(B) \quad (\{\sigma < \infty\}; \text{P-a.s.}) \quad (18)$$

6. As we said above, we are going to consider only discrete Markov chains (with phase space  $E = \{\dots, i, j, k, \dots\}$ ). To simplify the notation, we shall now denote the transition functions  $P(i; \{j\})$  by  $p_{ij}$  and call them transition

probabilities; an  $n$ -step transition probability from  $i$  to  $j$  will be denoted by  $p_{ij}^{(n)}$ .

Let  $E = \{1, 2, \dots\}$ . The principal questions that we study in §§2–4 are intended to clarify the conditions under which:

- (A) The limits  $\pi_j = \lim p_{ij}^{(n)}$  exist and are independent of  $i$ ;
- (B) The limits  $(\pi_1, \pi_2, \dots)$  form a *probability distribution*, that is,  $\pi_i \geq 0$ ,  $\pi_i = 1$ ;
- (C) The chain is *ergodic*, that is, the limits  $(\pi_1, \pi_2, \dots)$  have the properties  $\pi_i > 0$ ,  $\sum_{i=1}^{\infty} \pi_i = 1$ ;
- (D) There is one and only one *stationary probability distribution*  $\mathbb{Q} = (q_1, q_2, \dots)$ , that is, one such that  $q_i \geq 0$ ,  $\sum_{i=1}^{\infty} q_i = 1$ , and  $q_j = \sum_i q_i p_{ij}$ ,  $j \in E$ .

In the course of answering these questions we shall develop a classification of the states of a Markov chain as they depend on the arithmetic and asymptotic properties of  $p_{ij}^{(n)}$  and  $p_{ii}^{(n)}$ .

## 7. PROBLEMS

1. Prove the equivalence of definitions (1), (2), (3) and (15) of the Markov property.
2. Prove formula (5).
3. Prove equation (18).
4. Let  $(X_n)_{n \geq 0}$  be a Markov chain. Show that the reversed sequence  $(\dots, X_n, X_{n-1}, \dots, X_0)$  is also a Markov chain.

## §2. Classification of the States of a Markov Chain in Terms of Arithmetic Properties of the Transition Probabilities $p_{ij}^{(n)}$

1. We say that a state  $i \in E = \{1, 2, \dots\}$  is *inessential* if, with positive probability, it is possible to escape from it after a finite number of steps, without ever returning to it; that is, there exist  $m$  and  $j$  such that  $p_{ij}^{(m)} > 0$ , but  $p_{ji}^{(n)} = 0$  for all  $n$  and  $j$ .

Let us delete all the inessential states from  $E$ . Then the remaining set of *essential* states has the property that a wandering particle that encounters it can never leave it (Figure 36). As will become clear later, it is essential states that are the most interesting.

Let us now consider the set of essential states. We say that state  $j$  is *accessible* from the point  $i$  ( $i \rightarrow j$ ) if there is an  $m \geq 0$  such that  $p_{ij}^{(m)} > 0$

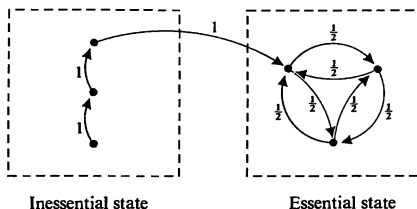


Figure 36

( $p_{ij}^{(0)} = 1$  if  $i = j$ , and 0 if  $i \neq j$ ). States  $i$  and  $j$  *communicate* ( $i \leftrightarrow j$ ) if  $j$  is accessible from  $i$  and  $i$  is accessible from  $j$ .

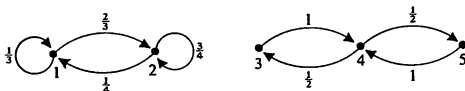
By the definition, the relation " $\leftrightarrow$ " is symmetric and reflexive. It is easy to verify that it is also transitive ( $i \leftrightarrow j, j \leftrightarrow k \Rightarrow i \leftrightarrow k$ ). Consequently the set of essential states separates into a finite or countable number of disjoint sets  $E_1, E_2, \dots$ , each of which consists of communicating sets but with the property that passage between different sets is impossible.

By way of abbreviation, we call the sets  $E_1, E_2, \dots$  *classes* or *indecomposable classes* (of essential communicating sets), and we call a Markov chain *indecomposable* if its states form a single indecomposable class.

As an illustration we consider the chain with matrix

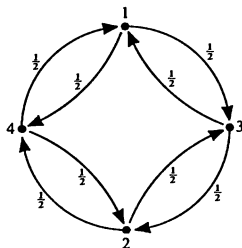
$$\mathbb{P} = \begin{pmatrix} \frac{1}{3} & \frac{2}{3} & 0 & 0 & 0 \\ \frac{1}{4} & \frac{3}{4} & 0 & 0 & 0 \\ \cdots & \cdots & \cdots & \cdots & \cdots \\ 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & \frac{1}{2} & 0 & \frac{1}{2} \\ 0 & 0 & 0 & 1 & 0 \end{pmatrix} = \begin{pmatrix} \mathbb{P}_1 & 0 \\ 0 & \mathbb{P}_2 \end{pmatrix}.$$

The graph of this chain, with set of states  $E = \{1, 2, 3, 4, 5\}$  has the form



It is clear that this chain has two indecomposable classes  $E_1 = \{1, 2\}$ ,  $E_2 = \{3, 4, 5\}$ , and the investigation of their properties reduces to the investigation of the two separate chains whose states are the sets  $E_1$  and  $E_2$ , and whose transition matrices are  $\mathbb{P}_1$  and  $\mathbb{P}_2$ .

Now let us consider any indecomposable class  $E$ , for example the one sketched in Figure 37.

Figure 37. Example of a Markov chain with period  $d = 2$ .

Observe that in this case a return to each state is possible only after an even number of steps; a transition to an adjacent state, after an odd number; the transition matrix has block structure,

$$\mathbb{P} = \begin{pmatrix} 0 & 0 & \frac{1}{2} & \frac{1}{2} \\ 0 & 0 & \frac{1}{2} & \frac{1}{2} \\ \dots & \dots & \dots & \dots \\ \frac{1}{2} & \frac{1}{2} & 0 & 0 \\ \frac{1}{2} & \frac{1}{2} & 0 & 0 \end{pmatrix}.$$

Therefore it is clear that the class  $E = \{1, 2, 3, 4\}$  separates into two subclasses  $C_0 = \{1, 2\}$  and  $C_1 = \{3, 4\}$  with the following *cyclic* property: after one step from  $C_0$  the particle necessarily enters  $C_1$ , and from  $C_1$  it returns to  $C_0$ .

This example suggests a classification of indecomposable classes into *cyclic subclasses*.

2. Let us say that state  $j$  has period  $d = d(j)$  if the following two conditions are satisfied:

- (1)  $p_{jj}^{(n)} > 0$  only for values of  $n$  of the form  $dm$ ;
- (2)  $d$  is the largest number satisfying (1).

In other words,  $d$  is the *greatest common divisor* of the numbers  $n$  for which  $p_{jj}^{(n)} > 0$ . (If  $p_{jj}^{(n)} = 0$  for all  $n \geq 1$ , we put  $d(j) = 0$ .)

Let us show that all states of a single indecomposable class  $E$  have the same period  $d$ , which is therefore naturally called the *period* of the class,  $d = d(E)$ .

Let  $i$  and  $j \in E$ . Then there are numbers  $k$  and  $l$  such that  $p_{ij}^{(k)} > 0$  and  $p_{ji}^{(l)} > 0$ . Consequently  $p_{ii}^{(k+l)} \geq p_{ij}^{(k)} p_{ji}^{(l)} > 0$ , and therefore  $k + l$  is divisible by  $d(i)$ . Suppose that  $n > 0$  and  $n$  is not divisible by  $d(i)$ . Then  $n + k + l$  is also not divisible by  $d(i)$  and consequently  $p_{ii}^{(n+k+l)} = 0$ . But

$$p_{ii}^{(n+k+l)} \geq p_{ij}^{(k)} p_{jj}^{(n)} p_{ji}^{(l)}$$



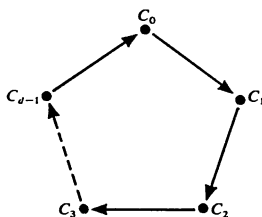


Figure 38. Motion among cyclic subclasses.

and therefore  $p_{jj}^{(n)} = 0$ . It follows that if  $p_{jj}^{(n)} > 0$  we have  $n$  divisible by  $d(i)$ , and therefore  $d(i) \leq d(j)$ . By symmetry,  $d(j) \leq d(i)$ . Consequently  $d(i) = d(j)$ .

If  $d(j) = 1$  ( $d(E) = 1$ ), the state  $j$  (or class  $E$ ) is said to be *aperiodic*.

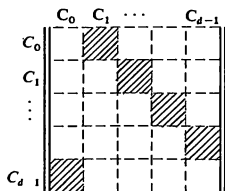
Let  $d = d(E)$  be the period of an indecomposable class  $E$ . The transitions within such a class may be quite freakish, but (as in the preceding example) there is a cyclic character to the transitions from one group of states to another. To show this, let us select a state  $i_0$  and introduce (for  $d \geq 1$ ) the following subclasses:

$$\begin{aligned} C_0 &= \{j \in E: p_{i_0 j}^{(n)} > 0 \Rightarrow n \equiv 0(\text{mod } d)\}; \\ C_1 &= \{j \in E: p_{i_0 j}^{(n)} > 0 \Rightarrow n \equiv 1(\text{mod } d)\}; \\ &\dots\dots\dots \\ C_{d-1} &= \{j \in E: p_{i_0 j}^{(n)} > 0 \Rightarrow n \equiv d-1(\text{mod } d)\}. \end{aligned}$$

Clearly  $E = C_0 + C_1 + \dots + C_{d-1}$ . Let us show that the motion from subclass to subclass is as indicated in Figure 38.

In fact, let state  $i \in C_p$  and  $p_{ij} > 0$ . Let us show that necessarily  $j \in C_{p+1(\text{mod } d)}$ . Let  $n$  be such that  $p_{i_0 i}^{(n)} > 0$ . Then  $n = ad + p$  and therefore  $n \equiv p(\text{mod } d)$  and  $n + 1 \equiv p + 1(\text{mod } d)$ . Hence  $p_{i_0 j}^{(n+1)} > 0$  and  $j \in C_{p+1(\text{mod } d)}$ .

Let us observe that it now follows that the transition matrix  $\mathbb{P}$  of an indecomposable chain has the following block structure:



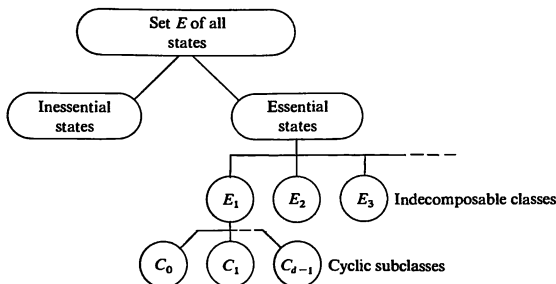


Figure 39. Classification of states of a Markov chain in terms of arithmetic properties of the probabilities  $p_{ij}^{(n)}$ .

Consider a subclass  $C_p$ . If we suppose that a particle is in the set  $C_0$  at the initial time, then at time  $s = p + dt, t = 0, 1, \dots$ , it will be in the subclass  $C_p$ . Consequently, with each subclass  $C_p$  we can connect a new Markov chain with transition matrix  $(p_{ij}^d)_{i,j \in C_p}$ , which is indecomposable and aperiodic. Hence if we take account of the classification that we have outlined (see the summary in Figure 39) we infer that in studying problems on limits of probabilities  $p_{ij}^{(n)}$  we can restrict our attention to *aperiodic indecomposable chains*.

### 3. PROBLEMS

1. Show that the relation " $\leftrightarrow$ " is transitive.
2. For Example 1, §5, show that when  $0 < p < 1$ , all states belong to a single class with period  $d = 2$ .
3. Show that the Markov chains discussed in Examples 4 and 5 of §5 are aperiodic.

## §3. Classification of the States of a Markov Chain in Terms of Asymptotic Properties of the Probabilities $p_{ii}^{(n)}$

1. Let  $\mathbb{P} = \|p_{ij}\|$  be the transition matrix of a Markov chain,

$$f_{ii}^{(k)} = P_i\{X_k = i, X_l \neq i, 1 \leq l \leq k-1\} \quad (1)$$

and for  $i \neq j$

$$f_{ij}^{(k)} = P_i\{X_k = j, X_l \neq j, 1 \leq l \leq k-1\}. \quad (2)$$

For  $X_0 = i$ , these are respectively the *probability of first return to state  $i$  at time  $k$* , and the *probability of first arrival at state  $j$  at time  $k$* .

Using the strong Markov property (1.16), we can show as in (I.12.38) that

$$p_{ij}^{(n)} = \sum_{k=1}^n f_{ij}^{(k)} p_{ij}^{(n-k)}. \quad (3)$$

For each  $i \in E$  we introduce

$$f_{ii} = \sum_{n=1}^{\infty} f_{ii}^{(n)}, \quad (4)$$

which is the probability that a particle that leaves state  $i$  will sooner or later return to that state. In other words,  $f_{ii} = P_i\{\sigma_i < \infty\}$ , where  $\sigma_i = \inf\{n \geq 1: X_n = i\}$  with  $\sigma_i = \infty$  when  $\{\cdot\} = \emptyset$ .

We say that a state  $i$  is *recurrent* if

$$f_{ii} = 1,$$

and *nonrecurrent* if

$$f_{ii} < 1.$$

Every recurrent state can, in turn, be classified according to whether the *average time of return* is finite or infinite.

Let us say that a recurrent state  $i$  is *positive* if

$$\mu_i^{-1} \equiv \left( \sum_{n=1}^{\infty} n f_{ii}^{(n)} \right)^{-1} > 0,$$

and *null* if

$$\mu_i^{-1} \equiv \left( \sum_{n=1}^{\infty} n f_{ii}^{(n)} \right)^{-1} = 0.$$

Thus we obtain the classification of the states of the chain, as displayed in Figure 40.

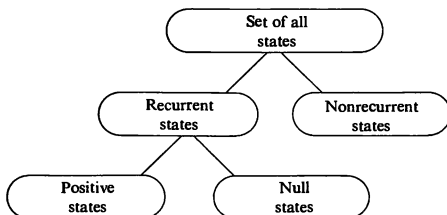


Figure 40. Classification of the states of a Markov chain in terms of the asymptotic properties of the probabilities  $p_{ii}^{(n)}$ .

2. Since the calculation of the functions  $f_{ii}^{(n)}$  can be quite complicated, it is useful to have the following tests for whether a state  $i$  is recurrent or not.

**Lemma 1**

(a) *The state  $i$  is recurrent if and only if*

$$\sum_{n=1}^{\infty} p_{ii}^{(n)} = \infty. \quad (5)$$

(b) *If state  $j$  is recurrent and  $i \leftrightarrow j$  then state  $i$  is also recurrent.*

PROOF. (a) By (3),

$$p_{ii}^{(n)} = \sum_{k=1}^n f_{ii}^{(k)} p_{ii}^{(n-k)},$$

and therefore (with  $p_{ii}^{(0)} = 1$ )

$$\begin{aligned} \sum_{n=1}^{\infty} p_{ii}^{(n)} &= \sum_{n=1}^{\infty} \sum_{k=1}^n f_{ii}^{(k)} p_{ii}^{(n-k)} = \sum_{k=1}^{\infty} f_{ii}^{(k)} \sum_{n=k}^{\infty} p_{ii}^{(n-k)} \\ &= f_{ii} \sum_{n=0}^{\infty} p_{ii}^{(n)} = f_{ii} \left( 1 + \sum_{n=1}^{\infty} p_{ii}^{(n)} \right). \end{aligned}$$

Therefore if  $\sum_{n=1}^{\infty} p_{ii}^{(n)} < \infty$ , we have  $f_{ii} < 1$  and therefore state  $i$  is non-recurrent. Furthermore, let  $\sum_{n=1}^{\infty} p_{ii}^{(n)} = \infty$ . Then

$$\sum_{n=1}^N p_{ii}^{(n)} = \sum_{n=1}^N \sum_{k=1}^n f_{ii}^{(k)} p_{ii}^{(n-k)} = \sum_{k=1}^N f_{ii}^{(k)} \sum_{n=k}^N p_{ii}^{(n-k)} \leq \sum_{k=1}^N f_{ii}^{(k)} \sum_{l=0}^N p_{ii}^{(l)},$$

and therefore

$$f_{ii} = \sum_{k=1}^{\infty} f_{ii}^{(k)} \geq \sum_{k=1}^N f_{ii}^{(k)} \geq \frac{\sum_{n=1}^N p_{ii}^{(n)}}{\sum_{l=0}^N p_{ii}^{(l)}} \rightarrow 1, \quad N \rightarrow \infty.$$

Thus if  $\sum_{n=1}^{\infty} p_{ii}^{(n)} = \infty$  then  $f_{ii} = 1$ , that is, the state  $i$  is recurrent.

(b) Let  $p_{ij}^{(s)} > 0$  and  $p_{ji}^{(t)} > 0$ . Then

$$p_{ii}^{(n+s+t)} \geq p_{ij}^{(s)} p_{ji}^{(t)} p_{ii}^{(n)},$$

and if  $\sum_{n=1}^{\infty} p_{ij}^{(n)} = \infty$ , then also  $\sum_{n=1}^{\infty} p_{ii}^{(n)} = \infty$ , that is, the state  $i$  is recurrent.

3. From (5) it is easy to deduce a first result on the asymptotic behavior of  $p_{ij}^{(n)}$ .

**Lemma 2.** *If state  $j$  is nonrecurrent then*

$$\sum_{n=1}^{\infty} p_{ij}^{(n)} < \infty \quad (6)$$

for every  $i$ , and therefore

$$p_{ij}^{(n)} \rightarrow 0, \quad n \rightarrow \infty. \quad (7)$$

PROOF. By (3) and Lemma 1,

$$\begin{aligned} \sum_{n=1}^{\infty} p_{ij}^{(n)} &= \sum_{n=1}^{\infty} \sum_{k=1}^n f_{ij}^{(k)} p_{jj}^{(n-k)} = \sum_{k=1}^{\infty} f_{ij}^{(k)} \sum_{n=0}^{\infty} p_{jj}^{(n)} \\ &= f_{ij} \sum_{n=0}^{\infty} p_{jj}^{(n)} \leq \sum_{n=0}^{\infty} p_{jj}^{(n)} < \infty. \end{aligned}$$

Here we used the inequality  $f_{ij} = \sum_{k=1}^{\infty} f_{ij}^{(k)} \leq 1$ , which holds because the series represents the probability that a particle starting at  $i$  eventually arrives at  $j$ . This establishes (6) and therefore (7).

Let us now consider recurrent states.

**Lemma 3.** Let  $j$  be a recurrent state with  $d(j) = 1$ .

(a) If  $i$  communicates with  $j$ , then

$$p_{ij}^{(n)} \rightarrow \frac{1}{\mu_j}, \quad n \rightarrow \infty. \quad (8)$$

If in addition  $j$  is a positive state then

$$p_{ij}^{(n)} \rightarrow \frac{1}{\mu_j} > 0, \quad n \rightarrow \infty. \quad (9)$$

If, however,  $j$  is a null state, then

$$p_{ij}^{(n)} \rightarrow 0, \quad n \rightarrow \infty. \quad (10)$$

(b) If  $i$  and  $j$  belong to different classes of communicating states, then

$$p_{ij}^{(n)} \rightarrow \frac{f_{ij}}{\mu_j}, \quad n \rightarrow \infty. \quad (11)$$

The proof of the lemma depends on the following theorem from analysis.

Let  $f_1, f_2, \dots$  be a sequence of nonnegative numbers with  $\sum_{i=1}^{\infty} f_i = 1$ , such that the greatest common divisor of the indices  $j$  for which  $f_j > 0$  is 1. Let  $u_0 = 1$ ,  $u_n = \sum_{k=1}^n f_k u_{n-k}$ ,  $n = 1, 2, \dots$ , and let  $\mu = \sum_{n=1}^{\infty} n f_n$ . Then  $u_n \rightarrow 1/\mu$  as  $n \rightarrow \infty$ . (For a proof, see [F1], §10 of Chapter XIII.)

Taking account of (3), we apply this to  $u_n = p_{jj}^{(n)}$ ,  $f_k = f_{jj}^{(k)}$ . Then we immediately find that

$$p_{jj}^{(n)} \rightarrow \frac{1}{\mu_j},$$

where  $\mu_j = \sum_{n=1}^{\infty} n f_{jj}^{(n)}$ .

Taking  $p_{jj}^{(s)} = 0$  for  $s < 0$ , we can rewrite (3) in the form

$$p_{ij}^{(n)} = \sum_{k=1}^{\infty} f_{ij}^{(k)} p_{jj}^{(n-k)}. \quad (12)$$

By what has been proved, we have  $p_{jj}^{(n-k)} \rightarrow \mu_j^{-1}$ ,  $n \rightarrow \infty$ , for each given  $k$ . Therefore if we suppose that

$$\lim_n \sum_{k=1}^{\infty} f_{ij}^{(k)} p_{jj}^{(n-k)} = \sum_{k=1}^{\infty} f_{ij}^{(k)} \lim_n p_{jj}^{(n-k)}, \quad (13)$$

we immediately obtain

$$p_{ij}^{(n)} \rightarrow \frac{1}{\mu_j} \left( \sum_{k=1}^{\infty} f_{ij}^{(k)} \right) = \frac{1}{\mu_j} f_{ij}, \quad (14)$$

which establishes (11).

Recall that  $f_{ij}$  is the probability that a particle starting from state  $i$  arrives, sooner or later, at state  $j$ . State  $j$  is recurrent, and if  $i$  communicates with  $j$ , it is natural to suppose that  $f_{ij} = 1$ . Let us show that this is indeed the case.

Let  $f'_{ij}$  be the probability that a particle, starting from state  $i$ , visits state  $j$  infinitely often. Clearly  $f_{ij} \geq f'_{ij}$ . Therefore if we show that, for a recurrent state  $j$  and a state  $i$  that communicates with it, the probability  $f'_{ij} = 1$ , we will have established that  $f_{ij} = 1$ .

According to part (b) of Lemma 1, the state  $i$  is also recurrent, and therefore

$$f_{ii} = \sum f_{ii}^{(n)} = 1. \quad (15)$$

Let

$$\sigma_i = \inf\{n \geq 1: X_n = i\}$$

be the first time (for times  $n \geq 1$ ) at which the particle reaches state  $i$ ; take  $\sigma_i = \infty$  if no such time exists.

Then

$$1 = f_{ii} = \sum_{n=1}^{\infty} f_{ii}^{(n)} = \sum_{n=1}^{\infty} P_i(\sigma_i = n) = P_i(\sigma_i < \infty), \quad (16)$$

and consequently to say that state  $i$  is recurrent means that a particle starting at  $i$  will eventually return to the same state (at a random time  $\sigma_i$ ). But after returning to this state the "life" of the particle starts over, so to speak (because of the strong Markov property). Hence it appears that if state  $i$  is recurrent the particle must return to it infinitely often:

$$P_i\{X_n = i \text{ for infinitely many } n\} = 1. \quad (17)$$

Let us now give a formal proof.

Let  $i$  be a state (recurrent or nonrecurrent). Let us show that the probability of return to that state at least  $r$  times is  $(f_{ii})^r$ .

For  $r = 1$  this follows from the definition of  $f_{ii}$ . Suppose that the proposition has been proved for  $r = m - 1$ . Then by using the strong Markov property and (16), we have

$$\begin{aligned}
 P_i(\text{number of returns to } i \text{ is greater than or equal to } m) \\
 &= \sum_{k=1}^{\infty} P_i\left(\sigma_i = k, \text{ and the number of returns to } i \text{ after time } k \right. \\
 &\quad \left. \text{is at least } m - 1 \right) \\
 &= \sum_{k=1}^{\infty} P_i(\sigma_i = k) P_i\left(\text{at least } m - 1 \text{ values} \right. \\
 &\quad \left. \text{of } X_{\sigma_i+1}, X_{\sigma_i+2}, \dots \text{ equal } i \mid \sigma_i = k \right) \\
 &= \sum_{k=1}^{\infty} P_i(\sigma_i = k) P_i(\text{at least } m - 1 \text{ values of } X_1, X_2, \dots \text{ equal } i) \\
 &= \sum_{k=1}^{\infty} f_{ii}^{(k)} (f_{ii})^{m-1} = f_{ii}^m.
 \end{aligned}$$

Hence it follows in particular that formula (17) holds for a recurrent state  $i$ . If the state is nonrecurrent, then

$$P_i\{X_n = i \text{ for infinitely many } n\} = 0. \quad (18)$$

We now turn to the proof that  $f'_{ij} = 1$ . Since the state  $i$  is recurrent, we have by (17) and the strong Markov property

$$\begin{aligned}
 1 &= \sum_{k=1}^{\infty} P_i(\sigma_j = k) + P_i(\sigma_j = \infty) \\
 &= \sum_{k=1}^{\infty} P_i\left(\sigma_j = k, \begin{array}{l} \text{the number of returns to } i \\ \text{after time } k \text{ is infinite} \end{array} \right) + P_i(\sigma_j = \infty) \\
 &= \sum_{k=1}^{\infty} P_i\left(\sigma_j = k, \begin{array}{l} \text{infinitely many values of} \\ X_{\sigma_i+1}, X_{\sigma_i+2}, \dots \text{ equal } i \end{array} \right) + P_i(\sigma_j = \infty) \\
 &= \sum_{k=1}^{\infty} P_i(\sigma_j = k) \cdot P_i\left(\begin{array}{l} \text{infinitely many} \\ \text{values of } X_{\sigma_j+1}, X_{\sigma_j+2}, \dots \\ \text{equal } i \end{array} \mid \sigma_j = k \right) + P_i(\sigma_j = \infty) \\
 &= \sum_{k=1}^{\infty} f_{ij}^{(k)} \cdot P_j\left(\begin{array}{l} \text{infinitely many values} \\ \text{of } X_1, X_2, \dots \text{ equal } i \end{array} \right) + (1 - f_{ij}) \\
 &= \sum_{k=1}^{\infty} f_{ij}^{(k)} f'_{ij} + (1 - f_{ij}) = f'_{ij} f_{ij} + (1 - f_{ij}).
 \end{aligned}$$

Thus

$$1 = f'_{ij} f_{ij} + 1 - f_{ij}$$

and therefore

$$f_{ij} = f'_{ij} \cdot f_{ij}.$$

Since  $i \leftrightarrow j$ , we have  $f_{ij} > 0$ , and consequently  $f'_{ij} = 1$  and  $f_{ij} = 1$ .

Therefore if we assume (13), it follows from (14) and the equation  $f_{ij} = 1$  that, for communicating states  $i$  and  $j$ ,

$$p_{ij}^{(n)} \rightarrow \frac{1}{\mu_j}, \quad n \rightarrow \infty.$$

As for (13), its validity follows from the theorem on dominated convergence together with the remark that

$$p_{jj}^{(n-k)} \rightarrow \frac{1}{\mu_j}, \quad n \rightarrow \infty, \quad \sum_{k=1}^{\infty} f_{ij}^{(k)} = f_{ij} \leq 1.$$

This completes the proof of the lemma.

Next we consider periodic states.

**Lemma 4.** *Let  $j$  be a recurrent state and let  $d(j) > 1$ .*

(a) *If  $i$  and  $j$  belong to the same class (of states), and if  $i$  belongs to the cyclic subclass  $C_r$  and  $j$  to  $C_{r+a}$ , then*

$$p_{ij}^{(nd+a)} \rightarrow \frac{d}{\mu_j}. \quad (19)$$

(b) *With an arbitrary  $i$ ,*

$$p_{ij}^{(nd+a)} \rightarrow \left[ \sum_{r=0}^{\infty} f_{ij}^{(rd+a)} \right] \cdot \frac{d}{\mu_j}, \quad a = 0, 1, \dots, d-1. \quad (20)$$

**PROOF.** (a) First let  $a = 0$ . With respect to the transition matrix  $\mathbb{P}^d$  the state  $j$  is recurrent and aperiodic. Consequently, by (8),

$$p_{ij}^{(nd)} \rightarrow \frac{1}{\sum_{k=1}^{\infty} k f_{jj}^{(kd)}} = \frac{d}{\sum_{k=1}^{\infty} k d f_{jj}^{(kd)}} = \frac{d}{\mu_j}.$$

Suppose that (19) has been proved for  $a = r$ . Then

$$p_{ij}^{(nd+r+1)} = \sum_{k=1}^{\infty} p_{ik} p_{kj}^{(nd+r)} \rightarrow \sum_{k=1}^{\infty} p_{ik} \cdot \frac{d}{\mu_j} = \frac{d}{\mu_j}.$$

(b) Clearly

$$p_{ij}^{(nd+a)} = \sum_{k=1}^{nd+a} f_{ij}^{(k)} p_{jj}^{(nd+a+k)}, \quad a = 0, 1, \dots, d-1.$$

State  $j$  has period  $d$ , and therefore  $p_{jj}^{(nd+a-k)} = 0$ , except when  $k - a$  has the form  $r \cdot d$ . Therefore

$$p_{ij}^{(nd+a)} = \sum_{r=0}^n f_{ij}^{(rd+a)} p_{jj}^{((n-r)d)}$$

and the required result (20) follows from (19).

This completes the proof of the lemma.



Lemmas 2-4 imply, in particular, the following result about limits of  $p_{ij}^{(n)}$ .

**Theorem 1.** *Let a Markov chain be indecomposable (that is, its states form a single class of essential communicating states) and aperiodic.*

*Then:*

(a) *If all states are either null or nonrecurrent, then, for all  $i$  and  $j$ ,*

$$p_{ij}^{(n)} \rightarrow 0, \quad n \rightarrow \infty; \quad (21)$$

(b) *if all states  $j$  are positive, then, for all  $i$ ,*

$$p_{ij}^{(n)} \rightarrow \frac{1}{\mu_j} > 0, \quad n \rightarrow \infty; \quad (22)$$

4. Let us discuss the conclusion of this theorem in the case of a Markov chain with a finite number of states,  $E = \{1, 2, \dots, r\}$ . Let us suppose that the chain is indecomposable and aperiodic. It turns out that then it is automatically increasing and positive:

$$\left( \begin{array}{c} \text{indecomposability} \\ d = 1 \end{array} \right) \Rightarrow \left( \begin{array}{c} \text{indecomposability} \\ \text{recurrence} \\ \text{positivity} \\ d = 1 \end{array} \right) \quad (23)$$

For the proof, we suppose that all states are nonrecurrent. Then by (21) and the finiteness of the set of states of the chain,

$$1 = \lim_n \sum_{j=1}^r p_{ij}^{(n)} = \sum_{j=1}^r \lim_n p_{ij}^{(n)} = 0. \quad (24)$$

The resulting contradiction shows that not all states can be nonrecurrent. Let  $i_0$  be a recurrent state and  $j$  an arbitrary state. Since  $i_0 \leftrightarrow j$ , Lemma 1 shows that  $j$  is also recurrent.

Thus all states of an aperiodic indecomposable chain are recurrent.

Let us now show that all recurrent states are positive.

If we suppose that they are all null states, we again obtain a contradiction with (24). Consequently there is at least one positive state, say  $i_0$ . Let  $i$  be any other state. Since  $i \leftrightarrow i_0$ , there are  $s$  and  $t$  such that  $p_{i_0 i}^{(s)} > 0$  and  $p_{ii_0}^{(t)} > 0$ , and therefore

$$p_{ii}^{(n+s+t)} \geq p_{ii_0}^{(s)} p_{i_0 i_0}^{(n)} p_{i_0 i}^{(t)} \rightarrow p_{ii_0}^{(s)} \frac{1}{\mu_{i_0}} \cdot p_{i_0 i}^{(t)} > 0. \quad (25)$$

Hence there is a positive  $\varepsilon$  such that  $p_{ii}^{(n)} \geq \varepsilon > 0$  for all sufficiently large  $n$ . But  $p_{ii}^{(n)} \rightarrow 1/\mu_i$  and therefore  $\mu_i > 0$ . Consequently (23) is established.

Let  $\pi_j = 1/\mu_j$ . Then  $\pi_j > 0$  by (22) and since

$$1 = \lim_n \sum_{j=1}^r p_{ij}^{(n)} = \sum_{j=1}^r \pi_j,$$

the (aperiodic indecomposable) chain is ergodic. Clearly, for all ergodic finite chains,

$$\text{there is an } n_0 \text{ such that } \min_{i,j} p_{ij}^{(n)} > 0 \text{ for all } n \geq n_0. \quad (26)$$

It was shown in §12 of Chapter I that the converse is also valid: (26) implies ergodicity.

Consequently we have the following implications:

$$\left( \begin{array}{c} \text{indecomposability} \\ d = 1 \end{array} \right) \Leftrightarrow \left( \begin{array}{c} \text{indecomposability} \\ \text{recurrence} \\ \text{positivity} \\ d = 1 \end{array} \right) \Rightarrow \text{ergodicity} \Leftrightarrow (26).$$

However, we can prove more.

**Theorem 2.** *For a finite Markov chain*

$$\left( \begin{array}{c} \text{indecomposability} \\ d = 1 \end{array} \right) \Leftrightarrow \left( \begin{array}{c} \text{indecomposability} \\ \text{recurrence} \\ \text{positivity} \\ d = 1 \end{array} \right) \Leftrightarrow (\text{ergodicity}) \Leftrightarrow (26).$$

**PROOF.** We have only to establish

$$(\text{ergodicity}) \Rightarrow \left( \begin{array}{c} \text{indecomposability} \\ \text{recurrence} \\ \text{positivity} \\ d = 1 \end{array} \right).$$

Indecomposability follows from (26). As for aperiodicity, increasingness, and positivity, they are valid in more general situations (the existence of a limiting distribution is sufficient), as will be shown in Theorem 2, §4.

## 5. PROBLEMS

1. Consider an indecomposable chain with states  $0, 1, 2, \dots$ . A necessary and sufficient condition for it to be nonrecurrent is that the system of equations  $u_j = \sum_i u_i p_{ij}$ ,  $j = 0, 1, \dots$ , has a bounded solution such that  $u_i \neq c$ ,  $i = 0, 1, \dots$ .
2. A sufficient condition for an indecomposable chain with states  $0, 1, \dots$  to be recurrent is that there is a sequence  $(u_0, u_1, \dots)$  with  $u_i \rightarrow \infty$ ,  $i \rightarrow \infty$ , such that  $u_j \geq \sum_i u_i p_{ij}$  for all  $j \neq 0$ .
3. A necessary and sufficient condition for an indecomposable chain with states  $0, 1, \dots$  to be recurrent and positive is that the system of equations  $u_j = \sum_i u_i p_{ij}$ ,  $j = 0, 1, \dots$ , has a solution, not identically zero, such that  $\sum_i |u_i| < \infty$ .

4. Consider a Markov chain with states  $0, 1, \dots$  and transition probabilities

$$p_{00} = r_0, \quad p_{01} = p_0 > 0,$$

$$p_{ij} = \begin{cases} p_i > 0, & j = i + 1, \\ r_i \geq 0, & j = i, \\ q_i > 0, & j = i - 1, \\ 0 & \text{otherwise.} \end{cases}$$

Let  $\rho_0 = 1, \rho_m = (q_1 \dots q_m)/(p_1 \dots p_m)$ . Prove the following propositions.

$$\text{Chain is recurrent} \Leftrightarrow \sum \rho_m = \infty,$$

$$\text{Chain is nonrecurrent} \Leftrightarrow \sum \rho_m < \infty,$$

$$\text{Chain is positive} \Leftrightarrow \sum \frac{1}{p_m \rho_m} < \infty,$$

$$\text{Chain is null} \Leftrightarrow \sum \rho_m = \infty, \sum \frac{1}{p_m \rho_m} = \infty.$$

5. Show that  $f_{ik} \geq f_{ij} f_{jk}$  and  $\sup_n p_{ij}^{(n)} \leq f_{ij} \leq \sum_{n=1}^{\infty} p_{ij}^{(n)}$ .

6. Show that for every Markov chain with countably many states, the limit of  $p_{ij}^{(n)}$  always exists in the *Cesàro* sense:

$$\lim_n \frac{1}{n} \sum_{k=1}^n p_{ij}^{(k)} = \frac{f_{ij}}{\mu_j}.$$

7. Consider a Markov chain  $\xi_0, \xi_1, \dots$  with  $\xi_{k+1} = (\xi_k^+) + \eta_{k+1}$ , where  $\eta_1, \eta_2, \dots$  is a sequence of independent identically distributed random variables with  $P(\eta_k = j) = p_j$ ,  $j = 0, 1, \dots$ . Write the transition matrix and show that if  $p_0 > 0, p_0 + p_1 < 1$ , the chain is recurrent if and only if  $\sum_k k p_k \leq 1$ .

## §4. On the Existence of Limits and of Stationary Distributions

1. We begin with some necessary conditions for the existence of stationary distributions.

**Theorem 1.** Let a Markov chain with countably many states  $E = \{1, 2, \dots\}$  and transition matrix  $\mathbb{P} = \|p_{ij}\|$  be such that the limits

$$\lim_n p_{ij}^{(n)} = \pi_j,$$

exist for all  $i$  and  $j$  and do not depend on  $i$ .

Then

- (a)  $\sum_i \pi_i \leq 1$ ,  $\sum_i \pi_i p_{ij} = \pi_j$ ;  
 (b) either all  $\pi_j = 0$  or  $\sum_j \pi_j = 1$ ;  
 (c) if all  $\pi_j = 0$ , there is no stationary distribution; if  $\sum_i \pi_i = 1$ , then  $\Pi = (\pi_1, \pi_2, \dots)$  is the unique stationary distribution.

PROOF. By Fatou's lemma,

$$\sum_j \pi_j = \sum_j \lim_n p_{ij}^{(n)} \leq \lim_n \sum_j p_{ij}^{(n)} = 1.$$

Moreover,

$$\sum_i \pi_i p_{ij} = \sum_i \left( \lim_n p_{ki}^{(n)} \right) p_{ij} \leq \lim_n \sum_i p_{ki}^{(n)} p_{ij} = \lim_n p_{kj}^{(n+1)} = \pi_j,$$

that is, for each  $j$ ,

$$\sum_i \pi_i p_{ij} \leq \pi_j.$$

Suppose that

$$\sum_i \pi_i p_{ij_0} < \pi_{j_0}$$

for some  $j_0$ . Then

$$\sum_j \pi_j > \sum_j \left( \sum_i \pi_i p_{ij} \right) = \sum_i \pi_i \sum_j p_{ij} = \sum_i \pi_i.$$

This contradiction shows that

$$\sum_i \pi_i p_{ij} = \pi_j \tag{1}$$

for all  $j$ .

It follows from (1) that

$$\sum_i \pi_i p_{ij}^{(n)} = \pi_j.$$

Therefore

$$\pi_j = \lim_n \sum_i \pi_i p_{ij}^{(n)} = \sum_i \pi_i \lim_n p_{ij}^{(n)} = \left( \sum_i \pi_i \right) \pi_j,$$

that is, for all  $j$ ,

$$\pi_j \left( 1 - \sum_i \pi_i \right) = 0,$$

from which (b) follows.

Now let  $\mathbb{Q} = (q_1, q_2, \dots)$  be a stationary distribution. Since  $\sum_i q_i p_{ij}^{(n)} = q_j$  and therefore  $\sum_i q_i \pi_j = q_j$ , that is,  $\pi_j = q_j$  for all  $j$ , this stationary distribution must coincide with  $\Pi = (\pi_1, \pi_2, \dots)$ . Therefore if all  $\pi_j = 0$ , there is no stationary distribution. If, however,  $\sum_j \pi_j = 1$ , then  $\Pi = (\pi_1, \pi_2, \dots)$  is the unique stationary distribution.

This completes the proof of the theorem.

Let us state and prove a fundamental result on the existence of a unique stationary distribution.

**Theorem 2.** *For Markov chains with countably many states, there is a unique stationary distribution if and only if the set of states contains precisely one positive recurrent class (of essential communicating states).*

**PROOF.** Let  $N$  be the number of positive recurrent classes.

Suppose  $N = 0$ . Then all states are either nonrecurrent or are recurrent null states, and by (3.10) and (3.20),  $\lim_n p_{ij}^{(n)} = 0$  for all  $i$  and  $j$ . Consequently, by Theorem 1, there is no stationary distribution.

Let  $N = 1$  and let  $C$  be the unique positive recurrent class. If  $d(C) = 1$  we have, by (3.8),

$$p_{ij}^{(n)} \rightarrow \frac{1}{\mu_j} > 0, \quad i, j \in C.$$

If  $j \notin C$ , then  $j$  is nonrecurrent, and  $p_{ij}^{(n)} \rightarrow 0$  for all  $i$  as  $n \rightarrow \infty$ , by (3.7).

Put

$$q_j = \begin{cases} \frac{1}{\mu_j} > 0, & j \in C, \\ 0, & j \notin C. \end{cases}$$

Then, by Theorem 1, the set  $\mathbb{Q} = (q_1, q_2, \dots)$  is the unique stationary distribution.

Now let  $d = d(C) > 1$ . Let  $C_0, \dots, C_{d-1}$  be the cyclic subclasses. With respect to  $\mathbb{P}^d$ , each subclass  $C_k$  is a recurrent aperiodic class. Then if  $i$  and  $j \in C_k$  we have

$$p_{ij}^{(nd)} \rightarrow \frac{d}{\mu_j} > 0$$

by (3.19). Therefore on each set  $C_k$ , the set  $d/\mu_j, j \in C_k$ , forms (with respect to  $\mathbb{P}^d$ ) the unique stationary distribution. Hence it follows, in particular, that  $\sum_{j \in C_k} (d/\mu_j) = 1$ , that is,  $\sum_{j \in C_k} (1/\mu_j) = 1/d$ .

Let us put

$$q_j = \begin{cases} \frac{1}{\mu_j}, & j \in C = C_0 + \dots + C_{d-1}, \\ 0, & j \notin C, \end{cases}$$

and show that, for the original chain, the set  $\mathbb{Q} = (q_1, q_2, \dots)$  is the unique stationary distribution.

In fact, for  $i \in C$ ,

$$p_{ii}^{(nd)} = \sum_{j \in C} p_{ij}^{(nd-1)} p_{ji}.$$

Then by Fatou's lemma,

$$\frac{d}{\mu_i} = \lim_n p_{ii}^{(nd)} \geq \sum_{j \in C} \lim_n p_{ij}^{(nd-1)} p_{ji} = \sum_{j \in C} \frac{1}{\mu_j} p_{ji}$$

and therefore

$$\frac{1}{\mu_i} \geq \sum_{j \in C} \frac{1}{\mu_j} p_{ji}.$$

But

$$\sum_{i \in C} \frac{1}{\mu_i} = \sum_{k=0}^{d-1} \left( \sum_{i \in C_k} \frac{1}{\mu_i} \right) = \sum_{k=0}^{d-1} \frac{1}{d} = 1.$$

As in Theorem 1, it can now be shown that in fact

$$\frac{1}{\mu_i} = \sum_{j \in C} \frac{1}{\mu_j} p_{ji}.$$

This shows that the set  $\mathbb{Q} = (q_1, q_2, \dots)$  is a stationary distribution, which is unique by Theorem 1.

Now let there be  $N \geq 2$  positive recurrent classes. Denote them by  $C^1, \dots, C^N$ , and let  $\mathbb{Q}^i = (q_1^i, q_2^i, \dots)$  be the stationary distribution corresponding to the class  $C^i$  and constructed according to the formula

$$q_j^i = \begin{cases} \frac{1}{\mu_j} > 0, & j \in C^i, \\ 0, & j \notin C^i. \end{cases}$$

Then, for all nonnegative numbers  $a_1, \dots, a_N$  such that  $a_1 + \dots + a_N = 1$ , the set  $a_1 \mathbb{Q}^1 + \dots + a_N \mathbb{Q}^N$  will also form a stationary distribution, since

$$(a_1 \mathbb{Q}^1 + \dots + a_N \mathbb{Q}^N) \mathbb{P} = a_1 \mathbb{Q}^1 \mathbb{P} + \dots + a_N \mathbb{Q}^N \mathbb{P} = a_1 \mathbb{Q}^1 + \dots + a_N \mathbb{Q}^N.$$

Hence it follows that when  $N \geq 2$  there is a continuum of stationary distributions. Therefore there is a unique stationary distribution only in the case  $N = 1$ .

This completes the proof of the theorem.

2. The following theorem answers the question of when there is a limit distribution for a Markov chain with a countable set of states  $E$ .

**Theorem 3.** *A necessary and sufficient condition for the existence of a limit distribution is that there is, in the set  $E$  of states of the chain, exactly one aperiodic positive recurrent class  $C$  such that  $f_{ij} = 1$  for all  $j \in C$  and  $i \in E$ .*

**PROOF.** *Necessity.* Let  $q_j = \lim p_{ij}^{(n)}$  and let  $\mathbb{Q} = (q_1, q_2, \dots)$  be a distribution ( $q_i \geq 0, \sum_i q_i = 1$ ). Then by Theorem 1 this limit distribution is the unique stationary distribution, and therefore by Theorem 2 there is one and only one recurrent positive class  $C$ . Let us show that this class has period  $d = 1$ . Suppose the contrary, that is, let  $d > 1$ . Let  $C_0, C_1, \dots, C_{d-1}$  be the cyclic subclasses. If  $i \in C_0$  and  $j \in C_1$ , then by (19),  $p_{ij}^{(nd+1)} \rightarrow d/\mu_i$  and  $p_{ij}^{(nd)} = 0$  for all  $n$ . But  $d/\mu_j > 0$ , and therefore  $p_{ij}^{(n)}$  does not have a limit as  $n \rightarrow \infty$ ; this contradicts the hypothesis that  $\lim_n p_{ij}^{(n)}$  exists. Now let  $j \in C$  and  $i \in E$ . Then, by (3.11),  $p_{ij}^{(n)} \rightarrow f_{ij}/\mu_j$ . Consequently  $\pi_j = f_{ij}/\mu_j$ . But  $\pi_j$  is independent of  $i$ . Therefore  $f_{ij} = f_{ji} = 1$ .

*Sufficiency.* By (3.11), (3.10) and (3.7),

$$p_{ij}^{(n)} \rightarrow \begin{cases} \frac{f_{ij}}{\mu_j}, & j \in C, \quad i \in E, \\ 0, & j \notin C, \quad i \in E. \end{cases}$$

Therefore if  $f_{ij} = 1$  for all  $j \in C$  and  $i \in E$ , then  $q_j = \lim_n p_{ij}^{(n)}$  is independent of  $i$ . Class  $C$  is positive and therefore  $q_j > 0$  for  $j \in C$ . Then, by Theorem 1, we have  $\sum_j q_j = 1$  and the set  $\mathbb{Q} = (q_1, q_2, \dots)$  is a limit distribution.

3. Let us summarize the results obtained above on the existence of a limit distribution, the uniqueness of a stationary distribution and ergodicity, for the case of finite chains.

**Theorem 4.** *We have the following implications for finite Markov chains:*

$$\begin{array}{ccc} (\text{ergodicity}) & \stackrel{\{1\}}{\Leftrightarrow} & \left( \begin{array}{l} \text{chain indecomposable,} \\ \text{recurrent, positive,} \\ \text{with } d = 1 \end{array} \right) \\ \Downarrow & & \Downarrow \\ \left( \begin{array}{l} \text{limit distribution} \\ \text{exists} \end{array} \right) & \stackrel{\{2\}}{\Leftrightarrow} & \left( \begin{array}{l} \text{there exists exactly one} \\ \text{recurrent positive class} \\ \text{with } d = 1 \end{array} \right) \\ \Downarrow & & \Downarrow \\ \left( \begin{array}{l} \text{unique stationary} \\ \text{distribution} \end{array} \right) & \stackrel{\{3\}}{\Leftrightarrow} & \left( \begin{array}{l} \text{there exists exactly one} \\ \text{recurrent positive class} \end{array} \right) \end{array}$$

**PROOF.** The "vertical" implications are evident. {1} is established in Theorem 2, §3; {2} in Theorem 3; {3} in Theorem 2.

## 4. PROBLEMS

1. Show that, in Example 1 of §5, neither stationary nor a limit distribution occurs.
2. Discuss the question of stationarity and limit distribution for the Markov chain with transition matrix

$$\mathbb{P} \begin{pmatrix} \frac{1}{2} & 0 & \frac{1}{2} & 0 \\ 0 & 0 & 0 & 1 \\ \frac{1}{4} & \frac{1}{2} & \frac{1}{4} & 0 \\ 0 & \frac{1}{2} & \frac{1}{2} & 0 \end{pmatrix}.$$

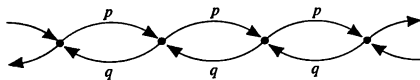
3. Let  $\mathbb{P} = \|p_{ij}\|$  be a finite doubly stochastic matrix, that is,  $\sum_{j=1}^m p_{ij} = 1, j = 1, \dots, m$ . Show that the stationary distribution of the corresponding Markov chain is the vector  $\mathbb{Q} = (1/m, \dots, 1/m)$ .

## §5. Examples

1. We present a number of examples to illustrate the concepts introduced above, and the results on the classification and limit behavior of transition probabilities.

EXAMPLE 1. A *simple random walk* is a Markov chain such that a particle remains in each state with a certain probability, and goes to the next state with a certain probability.

The simple random walk corresponding to the graph



describes the motion of a particle among the states  $E = \{0, \pm 1, \dots\}$  with transitions one unit to the right with probability  $p$  and to the left with probability  $q$ . It is clear that the transition probabilities are

$$p_{ij} = \begin{cases} p, & j = i + 1, \\ q, & j = i - 1, p + q = 1, \\ 0 & \text{otherwise.} \end{cases}$$

If  $p = 0$ , the particle moves deterministically to the left; if  $p = 1$ , to the right. These cases are of little interest since all states are inessential. We therefore assume that  $0 < p < 1$ .

With this assumption, the states of the chain form a single class (of essential communicating states). A particle can return to each state after 2, 4, 6, ... steps. Hence the chain has period  $d = 2$ .



Since, for each  $i \in E$ ,

$$p_{ii}^{(2n)} = C_{2n}^n (pq)^n = \frac{(2n)!}{(n!)^2} (pq)^n,$$

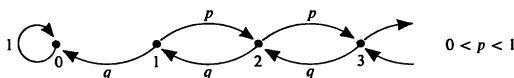
then by Stirling's formula (which says  $n! \sim \sqrt{2\pi n} n^n e^{-n}$ ) we have

$$p_{ii}^{(2n)} \sim \frac{(4pq)^n}{\sqrt{\pi n}}.$$

Therefore  $\sum_n p_{ii}^{(2n)} = \infty$  if  $p = q$ , and  $\sum_n p_{ii}^{(2n)} < \infty$  if  $p \neq q$ . In other words, the chain is recurrent if  $p = q$ , but if  $p \neq q$  it is nonrecurrent. It was shown in §10, Chapter I, that  $f_{ii}^{(2n)} \sim 1/(2\sqrt{\pi n}^{3/2})$ ,  $n \rightarrow \infty$ , if  $p = q = \frac{1}{2}$ . Therefore  $\mu_i = \sum_n (2n) f_{ii}^{(2n)} = \infty$ , that is, all recurrent states are null states. Hence by Theorem 1 of §3,  $p_{ij}^{(n)} \rightarrow 0$  as  $n \rightarrow \infty$  for all  $i$  and  $j$ .

There are no stationary, limit, or ergodic distributions.

EXAMPLE 2. Consider a simple random walk with  $E = \{0, 1, 2, \dots\}$ , where 0 is an absorbing barrier:



State 0 forms a unique positive recurrent class with  $d = 1$ . All other states are nonrecurrent. Therefore, by Theorem 2 of §4, there is a unique stationary distribution

$$\Pi = (\pi_0, \pi_1, \pi_2, \dots)$$

$$\text{with } \pi_0 = 1 \text{ and } \pi_i = 0, i \geq 1.$$

Let us now consider the question of limit distributions. Clearly  $p_{00}^{(n)} = 1$ ,  $p_{ij}^{(n)} \rightarrow 0$ ,  $j \geq 1$ ,  $i \geq 0$ . Let us now show that for  $i \geq 1$  the numbers  $\alpha(i) = \lim_n p_{i0}^{(n)}$  are given by the formulas

$$\alpha(i) = \begin{cases} \left(\frac{q}{p}\right)^i, & p > q, \\ 1, & p \leq q. \end{cases} \quad (1)$$

We begin by observing that since state 0 is absorbing we have  $p_{i0}^{(n)} = \sum_{k \leq n} f_{i0}^{(k)}$  and consequently  $\alpha(i) = f_{i0}$ , that is, the probability  $\alpha(i)$  is the probability that a particle starting from state  $i$  sooner or later reaches the null

state. By the method of §12, Chapter I (see also §2 of Chapter VII) we can obtain the recursion relation

$$\alpha(i) = p\alpha(i+1) + q\alpha(i-1), \quad (2)$$

with  $\alpha(0) = 1$ . The general solution of this equation has the form

$$\alpha(i) = a + b(q/p)^i, \quad (3)$$

and the condition  $\alpha(0) = 1$  imposes the condition  $a + b = 1$ .

If we suppose that  $q > p$ , then since  $\alpha(i)$  is bounded we see at once that  $b = 0$ , and therefore  $\alpha(i) = 1$ . This is quite natural, since when  $q > p$  the particle tends to move toward the null state.

If, on the other hand,  $p > q$  the opposite is true: the particle tends to move to the right, and so it is natural to expect that

$$\alpha(i) \rightarrow 0, \quad i \rightarrow \infty, \quad (4)$$

and consequently  $a = 0$  and

$$\alpha(i) = \left(\frac{q}{p}\right)^i. \quad (5)$$

To establish this equation, we shall not start from (4), but proceed differently.

In addition to the absorbing barrier at 0 we introduce an absorbing barrier at the integral point  $N$ . Let us denote by  $\alpha_N(i)$  the probability that a particle that starts at  $i$  reaches the zero state before reaching  $N$ . Then  $\alpha_N(i)$  satisfies (2) with the boundary conditions

$$\alpha_N(0) = 1, \quad \alpha_N(N) = 0,$$

and, as we have already shown in §9, Chapter I,

$$\alpha_N(i) = \frac{\left(\frac{q}{p}\right)^i - \left(\frac{q}{p}\right)^N}{1 - \left(\frac{q}{p}\right)^N}, \quad 0 \leq i \leq N. \quad (6)$$

Hence

$$\lim_N \alpha_N(i) = \left(\frac{q}{p}\right)^i$$

and consequently to prove (5) we have only to show that

$$\alpha(i) = \lim_N \alpha_N(i). \quad (7)$$

This is intuitively clear. A formal proof can be given as follows.

Let us suppose that the particle starts from a given state  $i$ . Then

$$\alpha(i) = P_i(A), \quad (8)$$

where  $A$  is the event in which there is an  $N$  such that a particle starting from  $i$  reaches the zero state before reaching state  $N$ . If

$$A_N = \{\text{particle reaches } 0 \text{ before } N\},$$

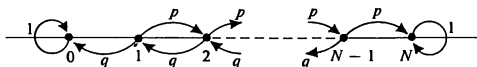
then  $A = \bigcup_{N=i+1}^{\infty} A_N$ . It is clear that  $A_N \subseteq A_{N+1}$  and

$$P_i\left(\bigcup_{N=i+1}^{\infty} A_N\right) = \lim_{N \rightarrow \infty} P_i(A_N). \quad (9)$$

But  $\alpha_N(i) = P_i(A_N)$ , so that (7) follows directly from (8) and (9).

Thus if  $p > q$  the limit  $\lim p_{i0}^{(n)}$  depends on  $i$ , and consequently there is no limit distribution in this case. If, however,  $p \leq q$ , then in all cases  $\lim p_{i0}^{(n)} = 1$  and  $\lim p_{ij}^{(n)} = 0$ ,  $j \geq 1$ . Therefore in this case the limit distribution has the form  $\Pi = (1, 0, 0, \dots)$ .

**EXAMPLE 3.** Consider a simple random walk with absorbing barriers at 0 and  $N$ :

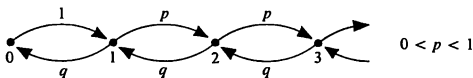


Here there are two positive recurrent classes  $\{0\}$  and  $\{N\}$ . All other states  $\{1, \dots, N-1\}$  are nonrecurrent. It follows from Theorem 1, §3, that there are infinitely many stationary distributions  $\Pi = (\pi_0, \pi_1, \dots, \pi_N)$  with  $\pi_0 = a$ ,  $\pi_N = b$ ,  $\pi_1 = \dots = \pi_{N-1} = 0$ , where  $a \geq 0$ ,  $b \geq 0$ ,  $a + b = 1$ . From Theorem 4, §4 it also follows that there is no limit distribution. This is also a consequence of the equations (Subsection 2, §9, Chapter I)

$$\lim_{n \rightarrow \infty} p_{i0}^{(n)} = \begin{cases} \frac{\left(\frac{q}{p}\right)^i - \left(\frac{q}{p}\right)^N}{1 - \left(\frac{q}{p}\right)^N}, & p \neq q, \\ 1 - \frac{i}{N}, & p = q, \end{cases} \quad (10)$$

$$\lim_n p_{iN}^{(n)} = 1 - \lim_n p_{i0}^{(n)} \quad \text{and} \quad \lim_n p_{ij}^{(n)} = 0, \quad 1 \leq j \leq N-1.$$

**EXAMPLE 4.** Consider a simple random walk with  $E = \{0, 1, \dots\}$  and a reflecting barrier at 0:



It is easy to see that the chain is periodic with period  $d = 2$ . Suppose that  $p > q$  (the moving particle tends to move to the right). Let  $i > 1$ ; to determine the probability  $f_{i1}$  we may use formula (1), from which it follows that

$$f_{i1} = \left(\frac{q}{p}\right)^{i-1} < 1, \quad i > 1.$$

All states of this chain communicate with each other. Therefore if state  $i$  is recurrent, state 1 will also be recurrent. But (see the proof of Lemma 3 in §3) in that case  $f_{i1}$  must be 1. Consequently when  $p > q$  all the states of the chain are nonrecurrent. Therefore  $p_{ij}^{(n)} \rightarrow 0$ ,  $n \rightarrow \infty$  for  $i$  and  $j \in E$ , and there is neither a limit distribution nor a stationary distribution.

Now let  $p \leq q$ . Then, by (1),  $f_{i1} = 1$  for  $i > 1$  and  $f_{11} = q + pf_{21} = 1$ . Hence the chain is recurrent.

Consider the system of equations determining the stationary distribution  $\Pi = (\pi_0, \pi_1, \dots)$ :

$$\begin{aligned}\pi_0 &= \pi_1 q, \\ \pi_1 &= \pi_0 + \pi_2 q, \\ \pi_2 &= \pi_1 p + \pi_3 q, \\ &\vdots\end{aligned}$$

that is,

$$\begin{aligned}\pi_1 &= \pi_1 q + \pi_2 q, \\ \pi_2 &= \pi_2 q + \pi_3 q, \\ &\vdots\end{aligned}$$

whence

$$\pi_j = \left(\frac{p}{q}\right) \pi_{j-1}, \quad j = 2, 3, \dots$$

If  $p = q$  we have  $\pi_1 = \pi_2 = \dots$ , and consequently

$$\pi_0 = \pi_1 = \pi_2 = \dots = 0.$$

In other words, if  $p = q$ , there is no stationary distribution, and therefore no limit distribution. From this and Theorem 3, §4, it follows, in particular, that in this case all states of the chain are null states.

It remains to consider the case  $p < q$ . From the condition  $\sum_{j=0}^{\infty} \pi_j = 1$  we find that

$$\pi_1 \left[ q + 1 + \left(\frac{p}{q}\right) + \left(\frac{p}{q}\right)^2 + \dots \right] = 1,$$

that is

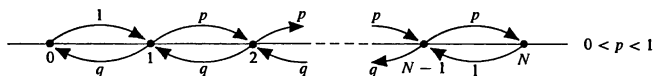
$$\pi_1 = \frac{q-p}{2q}$$

and

$$\pi_j = \frac{q-p}{2q} \cdot \left(\frac{p}{q}\right)^{j-1}; \quad j \geq 2.$$

Therefore the distribution  $\Pi$  is the unique stationary distribution. Hence when  $p < q$  the chain is recurrent and positive (Theorem 2, §4). The distribution  $\Pi$  is also a limit distribution and is ergodic.

EXAMPLE 5. Again consider a simple random walk with *reflecting barriers at 0 and N*:



All the states of the chain are periodic with period  $d = 2$ , recurrent, and positive. According to Theorem 4 of §4, the chain is ergodic. Solving the system  $\pi_j = \sum_{i=0}^N \pi_i p_{ij}$  subject to  $\sum_{i=0}^N \pi_i = 1$ , we obtain the ergodic distribution

$$\pi_i = \frac{\left(\frac{p}{q}\right)^{i-1}}{1 + \sum_{j=1}^{N-1} \left(\frac{p}{q}\right)^{j-1}}, \quad 2 \leq j \leq N-1,$$

and

$$\pi_0 = \pi_1 q, \quad \pi_N = \pi_{N-1} p.$$

2. EXAMPLE 6. It follows from Example 1 that the simple random walk considered there on the integral points of the line is recurrent if  $p = q$ , but nonrecurrent if  $p \neq q$ . Now let us consider simple random walks in the plane and in space, from the point of view of recurrence or nonrecurrence.

For the plane, we suppose that a particle in any state  $(i, j)$  moves up, down, to the right or to the left with probability  $\frac{1}{4}$  (Figure 41).

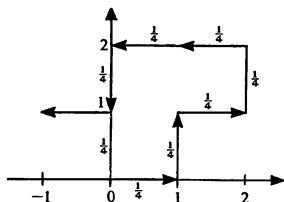


Figure 41. A walk in the plane.

For definiteness, consider the state  $(0, 0)$ . Then the probability  $P_k = p_{(0,0),(0,0)}^{(k)}$  of going from  $(0, 0)$  to  $(0, 0)$  in  $k$  steps is given by

$$P_{2n+1} = 0, \quad n = 0, 1, 2, \dots,$$

$$P_{2n} = \sum_{\{(i,j): i+j=n, 0 \leq i \leq n\}} \frac{(2n)!}{i!i!j!j!} \left(\frac{1}{4}\right)^{2n}, \quad n = 1, 2, \dots$$

Multiplying numerators and denominators by  $(n!)^2$ , we obtain

$$P_{2n} = \left(\frac{1}{4}\right)^{2n} C_{2n}^n \sum_{i=0}^n C_n^i C_n^{n-i} = \left(\frac{1}{4}\right)^{2n} (C_{2n}^n)^2,$$

since

$$\sum_{i=0}^n C_n^i C_n^{n-i} = C_{2n}^n.$$

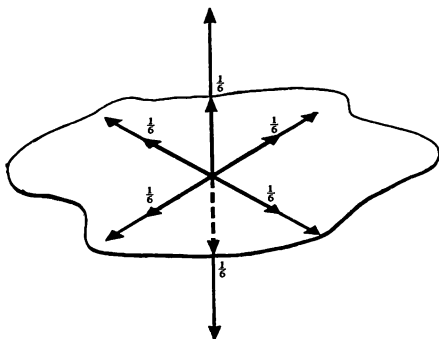
Applying Stirling's formula, we find that

$$P_{2n} \sim \frac{1}{\pi n},$$

and therefore  $\sum P_{2n} = \infty$ . Consequently the state  $(0, 0)$  (likewise any other) is *recurrent*.

It turns out, however, that in *three or more dimensions* the symmetric random walk is *nonrecurrent*. Let us prove this for walks on the integral points  $(i, j, k)$  in space.

Let us suppose that a particle moves from  $(i, j, k)$  by one unit along a coordinate direction, with probability  $\frac{1}{6}$  for each.



Then if  $P_k$  is the probability of going from  $(0, 0, 0)$  to  $(0, 0, 0)$  in  $k$  steps, we have

$$\begin{aligned}
 P_{2n+1} &= 0, \quad n = 0, 1, \dots, \\
 P_{2n} &= \sum_{\{(i, j): 0 \leq i+j \leq n, 0 \leq i \leq n, 0 \leq j \leq n\}} \frac{(2n)!}{(i!)^2 (j!)^2 (n-i-j)!^2} \left(\frac{1}{6}\right)^{2n} \\
 &= \frac{1}{2^{2n}} C_{2n}^n \sum_{\{(i, j): 0 \leq i+j \leq n, 0 \leq i \leq n, 0 \leq j \leq n\}} \left[ \frac{n!}{i! j! (n-i-j)!} \right]^2 \left(\frac{1}{3}\right)^{2n} \\
 &\leq C_n \frac{1}{2^{2n}} C_{2n}^n \frac{1}{3^n} \sum_{\{(i, j): 0 \leq i+j \leq n, 0 \leq i \leq n, 0 \leq j \leq n\}} \frac{n!}{i! j! (n-i-j)!} \left(\frac{1}{3}\right)^n \\
 &= C_n \frac{1}{2^{2n}} C_{2n}^n \frac{1}{3^n}, \quad n = 1, 2, \dots,
 \end{aligned} \tag{11}$$

where

$$C_n = \max_{\{(i, j): 0 \leq i+j \leq n, 0 \leq i \leq n, 0 \leq j \leq n\}} \left[ \frac{n!}{i! j! (n-i-j)!} \right]. \tag{12}$$

Let us show that when  $n$  is large, the max in (12) is attained for  $i \sim n/3$ ,  $j \sim n/3$ . Let  $i_0$  and  $j_0$  be the values at which the max is attained. Then the following inequalities are evident:

$$\begin{aligned}
 \frac{n!}{j_0! (i_0 - 1)! (n - j_0 - i_0 + 1)!} &\leq \frac{n!}{j_0! i_0! (n - j_0 - i_0)!}, \\
 \frac{n!}{j_0! (i_0 + 1)! (n - j_0 - i_0 - 1)!} &\leq \frac{n!}{(j_0 - 1)! i_0! (n - j_0 - i_0 + 1)!} \\
 &\leq \frac{n!}{(j_0 + 1)! i_0! (n - j_0 - i_0 - 1)!},
 \end{aligned}$$

whence

$$\begin{aligned}
 n - i_0 - 1 &\leq 2j_0 \leq n - i_0 + 1, \\
 n - j_0 - 1 &\leq 2i_0 \leq n - j_0 + 1,
 \end{aligned}$$

and therefore we have, for large  $n$ ,  $i_0 \sim n/3$ ,  $j_0 \sim n/3$ , and

$$C_n \sim \frac{n!}{\left[\left(\frac{n}{3}\right)!\right]^3}.$$

By Stirling's formula,

$$C_n \frac{1}{2^{2n}} C_{2n}^n \frac{1}{3^n} \sim \frac{3\sqrt{3}}{2\pi^{3/2} n^{3/2}},$$

and since

$$\sum_{n=1}^{\infty} \frac{3\sqrt{3}}{2\pi^{3/2}n^{3/2}} < \infty,$$

we have  $\sum_n P_{2n} < \infty$ . Consequently the state  $(0, 0, 0)$ , and likewise any other state, is nonrecurrent. A similar result holds for dimensions greater than 3.

Thus we have the following result (Pólya):

**Theorem.** For  $R^1$  and  $R^2$ , the symmetric random walk is recurrent; for  $R^n$ ,  $n \geq 3$ , it is nonrecurrent.

### 3. PROBLEMS

1. Derive the recursion relation (1).
2. Establish (4).
3. Show that in Example 5 all states are aperiodic, recurrent, and positive.
4. Classify the states of a Markov chain with transition matrix

$$\mathbb{P} = \begin{pmatrix} p & q & 0 & 0 \\ 0 & 0 & p & q \\ p & q & 0 & 0 \\ 0 & 0 & p & q \end{pmatrix},$$

where  $p + q = 1$ ,  $p \geq 0$ ,  $q \geq 0$ .





# Historical and Bibliographical Notes

## Introduction

The history of probability theory up to the time of Laplace is described by Todhunter [T1]. The period from Laplace to the end of the nineteenth century is covered by Gnedenko and Sheinin in [K10]. Maistrov [M1] discusses the history of probability theory from the beginning to the thirties of the present century. There is a brief survey in Gnedenko [G4]. For the origin of much of the terminology of the subject see Aleksandrova [A3].

For the basic concepts see Kolmogorov [K8], Gnedenko [G4], Borovkov [B4], Gnedenko and Khinchin [G6], A. M. and I. M. Yaglom [Y1], Prokhorov and Rozanov [P5], Feller [F1, F2], Neyman [N3], Loève [L7], and Doob [D3]. We also mention [M3] which contains a large number of problems on probability theory.

In putting this text together, the author has consulted a wide range of sources. We mention particularly the books by Breiman [B5], Ash [A4, A5], and Ash and Gardner [A6], which (in the author's opinion) contain an excellent selection and presentation of material.

For current work in the field see, for example, *Annals of Probability* (formerly *Annals of Mathematical Statistics*) and *Theory of Probability and its Applications* (translation of *Teoriya Veroyatnostei i ee Primeneniya*).

*Mathematical Reviews* and *Zentralblatt für Mathematik* contain abstracts of current papers on probability and mathematical statistics from all over the world.

For tables for use in computations, see [A1].

## Chapter I

§1. Concerning the construction of probabilistic models see Kolmogorov [K7] and Gnedenko [G4]. For further material on problems of distributing objects among boxes see, e.g., Kolchin, Sevastyanov and Chistyakov [K3].

§2. For other probabilistic models (in particular, the one-dimensional Ising model) that are used in statistical physics, see Işihara [I2].

§3. Bayes's formula and theorem form the basis for the "Bayesian approach" to mathematical statistics. See, for example, De Groot [D1] and Zacks [Z1].

§4. A variety of problems about random variables and their probabilistic description can be found in Meshalkin [M3].

§5. A combinatorial proof of the law of large numbers (originating with James Bernoulli) is given in, for example, Feller [F1]. For the empirical meaning of the law of large numbers see Kolmogorov [K7].

§6. For sharper forms of the local and integrated theorems, and of Poisson's theorem, see Borovkov [B4] and Prokhorov [P3].

§7. The examples of Bernoulli schemes illustrate some of the basic concepts and methods of mathematical statistics. For more detailed discussions see, for example, Cramér [C5] and van der Waerden [W1].

§8. Conditional probability and conditional expectation with respect to a partition will help the reader understand the concepts of conditional probability and conditional expectation with respect to  $\sigma$ -algebras, which will be introduced later.

§9. The ruin problem was considered in essentially the present form by Laplace. See Gnedenko and Sheinin [K10]. Feller [F1] contains extensive material from the same circle of ideas.

§10. Our presentation essentially follows Feller [F1]. The method for proving (10) and (11) is taken from Doherty [D2].

§11. Martingale theory is thoroughly covered in Doob [D3]. A different proof of the ballot theorem is given, for instance, in Feller [F1].

§12. There is extensive material on Markov chains in the books by Feller [F1], Dynkin [D4], Kemeny and Snell [K2], Sarymsakov [S1], and Sirazhdinov [S8]. The theory of branching processes is discussed by Sevastyanov [S3].

## Chapter II

§1. Kolmogorov's axioms are presented in his book [K8].

§2. Further material on algebras and  $\sigma$ -algebras can be found in, for example, Kolmogorov and Fomin [K8], Neveu [N1], Breiman [B5], and Ash [A5].

§3. For a proof of Carathéodory's theorem see Loève [I7] or Halmos [H1].

§§4–5. More material on measurable functions is available in Halmos [H1].

§6. See also Kolmogorov and Fomin [K8], Halmos [H1], and Ash and Gardner [A6]. The Radon–Nikodým theorem is proved in these books. The inequality

$$P(|\xi| \geq \varepsilon) \leq \frac{E\xi^2}{\varepsilon^2}$$

is sometimes called Chebyshev's inequality, and the inequality

$$P(|\xi| \geq \varepsilon) \leq \frac{E|\xi|^r}{\varepsilon^r}, \quad r > 0,$$

is called *Markov's inequality*.

For Pratt's lemma see [P2].

§7. The definitions of conditional probability and conditional expectation with respect to a  $\sigma$ -algebra were given by Kolmogorov [K8]. For additional material see Breiman [B5] and Ash [A5]. The result quoted in the Corollary to Theorem 5 can be found in [M5].

§8. See also Borovkov [B4], Ash [A5], Cramér [C5], and Gnedenko [G4].

§9. Kolmogorov's theorem on the existence of a process with given finite-dimensional distribution is in his book [K8]. For Ionescu–Tulcea's theorem see also Neveu [N1] and Ash [A5]. The proof in the text follows [A5].

§§10–11. See also Kolmogorov and Fomin [K9], Ash [A5], Doob [D3], and Loève [L7].

§12. The theory of characteristic functions is presented in many books. See, for example, Gnedenko [G4], Gnedenko and Kolmogorov [G5], Ramachandran [R1], Lukacs [L8], and Lukacs and Laha [L9]. Our presentation of the connection between moments and semi-invariants follows Leonov and Shiryayev [L4].

§13. See also Ibragimov and Rozanov [I1], Breiman [B5], and Liptser and Shiryayev [L5].

## Chapter III

§1. Detailed investigations of problems on weak convergence of probability measures are given in Gnedenko and Kolmogorov [G5] and Billingsley [B3].

§2. Prokhorov's theorem appears in his paper [P4].

§3. The monograph [G5] by Gnedenko and Kolmogorov studies the limit theorems of probability theory by the method of characteristic functions. See also Billingsley [B3]. Problem 2 includes both Bernoulli's law of large numbers and Poisson's law of large numbers (which assumes that  $\xi_1, \xi_2, \dots$  are independent and take only two values (1 and 0), but in general are differently distributed:  $P(\xi_i = 1) = p_i$ ,  $P(\xi_i = 0) = 1 - p_i$ ,  $i \geq 1$ ).

§4. Here we give the standard proof of the central limit theorem for sums of independent random variables under the Lindeberg condition. Compare [G5] and [P6].

In the first edition, we gave the proof of Theorem 3 here.

§5. Questions of the validity of the central limit theorem without the hypothesis of negligibility in the limit have already attracted the attention of P. Lévy. A detailed account of the current state of the theory of limit theorems in the *nonclassical setting* is contained in Zolotarev [Z4]. The statement and proof of Theorem 1 were given by Rotar [R5].

§6. The presentation uses material from Gnedenko and Kolmogorov [G5], Ash [A5], and Petrov [P1], [P6].

§7. The Lévy-Prokhorov metric was introduced in a well-known work by Prokhorov [P4], to whom the results on metrizability of weak convergence of measures given on metric spaces are also due. Concerning the metric  $\|P - \tilde{P}\|_{BL}^*$ , see Dudley [D6] and Pollard [P7].

§8. Theorem 1 is due to Skorokhod. Useful material on the method of a single probability space may be found in Borovkov [B4] and in Pollard [P7].

§§9–10. A number of books contain a great deal of material touching on these questions: Jacod and Shiryaev [J1], LeCam [L10], Greenwood and Shiryaev [G7], and Liese and Vajda [L11].

§11. Petrov [P6] contains a lot of material on estimates of the rate of convergence in the central limit theorem. The proof given of the theorem of Berry and Esseen is contained in Gnedenko and Kolmogorov [G5].

§12. The proof follows Presman [P8].

## Chapter IV

§1. Kolmogorov's zero-or-one law appears in his book [K8]. For the Hewitt-Savage zero-or-one law see also Borovkov [B4], Breiman [B5], and Ash [A5].

§§2–4. Here the fundamental results were obtained by Kolmogorov and Khinchin (see [K8] and the references given there). See also Petrov [P1] and Stout [S9]. For probabilistic methods in number theory see Kubilius [K11].

§5. Regarding these questions, see Petrov [P6], Borovkov [B4], and Dacunha-Castelle and Duflo [D5].

## Chapter V

§§1–3. Our exposition of the theory of (strict sense) stationary random processes is based on Breiman [B5], Sinai [S7], and Lamperti [I2]. The simple proof of the maximal ergodic theorem was given by Garsia [G1].

## Chapter VI

§1. The books by Rozanov [R4], and Gihman and Skorohod [G2, G3] are devoted to the theory of (wide sense) stationary random processes. Example 6 was frequently presented in Kolmogorov's lectures.

§2. For orthogonal stochastic measures and stochastic integrals see also Doob [D3], Gihman and Skorohod [G3], Rozanov [R4], and Ash and Gardner [A6].

§3. The spectral representation (2) was obtained by Cramér and Loève (see, for example, [I7]). The same representation (in different language) is contained in Kolmogorov [K5]. Also see Doob [D3], Rozanov [R4], and Ash and Gardner [A6].

§4. There is a detailed exposition of problems of statistical estimation of the covariance function and spectral density in Hannan [H2, H3].

§§5–6. See also Rozanov [R4], Lamperti [L2], and Gihman and Skorohod [G2, G3].

§7. The presentation follows Lipster and Shiryaev [L5].

## Chapter VII

§1. Most of the fundamental results of the theory of martingales were obtained by Doob [D3]. Theorem 1 is taken from Meyer [M4]. Also see Meyer and Dellacherie [M5], Lipster and Shiryaev [L5], and Gihman and Skorohod [G3].

§2. Theorem 1 is often called the theorem "on transformation under a system of optional stopping" (Doob [D3]). For the identities (14) and (15) and Wald's fundamental identity see Wald [W2].

§3. Chow and Teicher [C3] contains an illuminating study of the results presented here, including proofs of the inequalities of Khinchin, Marcinkiewicz and Zygmund, Burkholder, and Davis. Theorem 2 was given by Lenglart [L3].

§4. See Doob [D3].

§5. Here we follow Kabanov, Liptser and Shiryaev [K1], Engelbert and Shiryaev [E1], and Neveu [N2]. Theorem 4 and the example were given by Liptser.

§6. This approach to problems of absolute continuity and singularity, and the results given here, can be found in Kabanov, Liptser and Shiryaev [K1]. Theorem 6 was obtained by Kabanov.

§7. Theorems 1 and 2 were given by Novikov [N4]. Lemma 1 is a discrete analog of Girsanov's lemma (see [K1]).

§8. See also Liptser and Shiryaev [L12] and Jacod and Shiryaev [J1], which discuss limit theorems for random processes of a rather general nature (for example, martingales, semi-martingales).

## Chapter VIII

§1. For the basic definitions see Dynkin [D4], Ventzel [V2], Doob [D3], and Gihman and Skorohod [G3]. The existence of regular transition probabilities such that the Kolmogorov–Chapman equation (9) is satisfied for all  $x \in R$  is proved in [N1] (corollary to Proposition V.2.1) and in [G3] (Volume I, Chapter II, §4). Kuznetsov (see Abstracts of the Twelfth European Meeting of Statisticians, Varna, 1979) has established the validity (which is far from trivial) of a similar result for Markov processes with continuous times and values in universal measurable spaces.

§§2–5. Here the presentation follows Kolmogorov [K4], Borovkov [B4], and Ash [A4].

# References<sup>†</sup>

- [A1] M. Abramovitz and I. A. Stegun, editors. *Handbook of Mathematical Functions with Formulas, Graphs and Mathematical Tables*. National Bureau of Standards, Washington, D.C., 1964.
- [A2] P. S. Aleksandrov. *Einführung in die Mengenlehre und die Theorie der reellen Funktionen*. DVW, Berlin, 1956.
- [A3] N. V. Aleksandrova. *Mathematical Terms [Matematicheskie terminy]*. Vysshaya Shkola, Moscow, 1978.
- [A4] R. B. Ash. *Basic Probability Theory*. Wiley, New York, 1970.
- [A5] R. B. Ash. *Real Analysis and Probability*. Academic Press, New York, 1972.
- [A6] R. B. Ash and M. F. Gardner. *Topics in Stochastic Processes*. Academic Press, New York, 1975.
- [B1] S. N. Bernshtein. Chebyshev's work on the theory of probability (in Russian), in *The Scientific Legacy of P. L. Chebyshev [Nauchnoe nasledie P. L. Chebysheva]*, pp. 43–68, Akademiya Nauk SSSR, Moscow–Leningrad, 1945.
- [B2] S. N. Bernshtein. *Theory of Probability [Teoriya veroyatnostei]*, 4th ed. Gostehizdat, Moscow, 1946.
- [B3] P. Billingsley. *Convergence of Probability Measures*. Wiley, New York, 1968.
- [B4] A. A. Borovkov. *Wahrscheinlichkeitstheorie: eine Einführung*, first edition Birkhäuser, Basel–Stuttgart, 1976; *Theory of Probability*, second edition [Teoriya veroyatnostei]. "Nauka," Moscow, 1986.
- [B5] L. Breiman. *Probability*. Addison-Wesley, Reading, MA, 1968.
- [C1] P. L. Chebyshev. *Theory of Probability: Lectures Given in 1879 and 1880 [Teoriya veroyatnostei: Lektsii akad. P. L. Chebysheva chitaniye v 1879, 1880 gg.]*. Edited by A. N. Krylov from notes by A. N. Lyapunov. Moscow–Leningrad, 1936.
- [C2] Y. S. Chow, H. Robbins, and D. Siegmund. *The Theory of Optimal Stopping*. Dover, New York, 1991.

<sup>†</sup> Translator's note: References to translations into Russian have been replaced by references to the original works. Russian references have been replaced by their translations whenever I could locate translations; otherwise they are reproduced (with translated titles). Names of journals are abbreviated according to the forms used in *Mathematical Reviews* (1982 versions).



- [C3] Y. S. Chow and H. Teicher. *Probability Theory: Independence, Interchangeability, Martingales*. Springer-Verlag, New York, 1978.
- [C4] Kai-Lai Chung. *Markov Chains with Stationary Transition Probabilities*. Springer-Verlag, New York, 1967.
- [C5] H. Cramér, *Mathematical Methods of Statistics*. Princeton University Press, Princeton, NJ, 1957.
- [D1] M. H. De Groot. *Optimal Statistical Decisions*. McGraw-Hill, New York, 1970.
- [D2] M. Doherty. An amusing proof in fluctuation theory. *Lecture Notes in Mathematics*, no. 452, 101–104, Springer-Verlag, Berlin, 1975.
- [D3] J. L. Doob. *Stochastic Processes*. Wiley, New York, 1953.
- [D4] E. B. Dynkin. *Markov Processes*. Plenum, New York, 1963.
- [E1] H. J. Engelbert and A. N. Shiryaev. On the sets of convergence and generalized submartingales. *Stochastics* 2 (1979), 155–166.
- [F1] W. Feller. *An Introduction to Probability Theory and Its Applications*, vol. 1, 3rd ed. Wiley, New York, 1968.
- [F2] W. Feller. *An Introduction to Probability Theory and Its Applications*, vol. 2, 2nd ed. Wiley, New York, 1966.
- [G1] A. Garcia. A simple proof of E. Hopf's maximal ergodic theorem. *J. Math. Mech.* 14 (1965), 381–382.
- [G2] I. I. Gihman [Gikhman] and A. V. Skorohod [Skorokhod]. *Introduction to the theory of random processes*, first edition. Saunders, Philadelphia, 1969; second edition [Vvedenie v teoriyu sluchainykh protsessov]. "Nauka", Moscow, 1977.
- [G3] I. I. Gihman and A. V. Skorohod. *Theory of Stochastic Processes*, 3 vols. Springer-Verlag, New York–Berlin, 1974–1979.
- [G4] B. V. Gnedenko. *The theory of probability*. Mir, Moscow, 1988.
- [G5] B. V. Gnedenko and A. N. Kolmogorov. *Limit Distributions for Sums of Independent Random Variables*, revised edition. Addison-Wesley, Reading, 1968.
- [G6] B. V. Gnedenko and A. Ya. Khinchin. *An Elementary Introduction to the Theory of Probability*. Freeman, San Francisco, 1961; ninth edition [Elementarnoe vvedenie v teoriyu veroyatnostei]. "Nauka", Moscow, 1982.
- [H1] P. R. Halmos. *Measure Theory*. Van Nostrand, New York, 1950.
- [H2] E. J. Hannan. *Time Series Analysis*. Methuen, London, 1960.
- [H3] E. J. Hannan. *Multiple Time Series*. New York, Wiley, 1970.
- [I1] I. A. Ibragimov and Yu. V. Linnik. *Independent and Stationary Sequences of Random Variables*. Walters-Noordhoff, Groningen, 1971.
- [I2] I. A. Ibragimov and Yu. A. Rozanov. *Gaussian Random Processes*. Springer-Verlag, New York, 1978.
- [I3] A. Isihara. *Statistical Physics*. Academic Press, New York, 1971.
- [K1] Yu. M. Kabanov, R. Sh. Liptser, and A. N. Shiryaev. On the question of the absolute continuity and singularity of probability measures. *Math. USSR-Sb.* 33 (1977), 203–221.
- [K2] J. Kemeny and L. J. Snell. *Finite Markov Chains*. Van Nostrand, Princeton, 1960.
- [K3] V. F. Kolchin, B. A. Sevastyanov, and V. P. Chistyakov. *Random Allocations*. Halsted, New York, 1978.
- [K4] A. N. Kolmogorov, Markov chains with countably many states. *Byull. Moskov. Univ.* 1 (1937), 1–16 (in Russian).
- [K5] A. N. Kolmogorov. Stationary sequences in Hilbert space. *Byull. Moskov. Univ. Mat.* 2 (1941), 1–40 (in Russian).
- [K6] A. N. Kolmogorov, The contribution of Russian science to the development

- of probability theory. *Uchen. Zap. Moskov. Univ.* 1947, no. 91, 56ff. (in Russian).
- [K7] A. N. Kolmogorov, Probability theory (in Russian), in *Mathematics: Its Contents, Methods, and Value* [*Matematika, ee sodержanie, metody i znachenie*]. Akad. Nauk SSSR, vol. 2, 1956.
- [K8] A. N. Kolmogorov. *Foundations of the Theory of Probability*. Chelsea, New York, 1956; second edition [*Osnovnye poniatiya teorii veroyatnostei*]. "Nauka", Moscow, 1974.
- [K9] A. N. Kolmogorov and S. V. Fomin. *Elements of the Theory of Functions and Functionals Analysis*. Graylok, Rochester, 1957 (vol. 1), 1961 (vol. 2); sixth edition [*Elementy teorii funktsii i funktsional'nogo analiza*]. "Nauka", Moscow, 1989.
- [K10] A. N. Kolmogorov and A. P. Yushkevich, editors. *Mathematics of the Nineteenth Century* [*Matematika XIX veka*]. Nauka, Moscow, 1978.
- [K11] J. Kubilius. *Probabilistic Methods in the Theory of Numbers*. American Mathematical Society, Providence, 1964.
- [L1] J. Lamperti. *Probability*. Benjamin, New York, 1966.
- [L2] J. Lamperti. *Stochastic Processes*. Springer-Verlag, New York, 1977.
- [L3] E. Lengart. Relation de domination entre deux processus. *Ann. Inst. H. Poincaré. Sect. B (N.S.)* 13 (1977), 171–179.
- [L4] V. P. Leonov and A. N. Shiryaev. On a method of calculation of semi-invariants. *Theory Probab. Appl.* 4 (1959), 319–329.
- [L5] R. S. Liptser and A. N. Shiryaev. *Statistics of Random Processes*. Springer-Verlag, New York, 1977.
- [L6] R. Sh. Liptser and A. N. Shiryaev. A functional central limit theorem for semimartingales. *Theory Probab. Appl.* 25 (1980), 667–688.
- [L7] M. Loève. *Probability Theory*. Springer-Verlag, New York, 1977–78.
- [L8] E. Lukacs. *Characteristic Functions*. Hafner, New York, 1960.
- [L9] E. Lukacs and R. G. Laha. *Applications of Characteristic Functions*. Hafner, New York, 1964.
- [M1] D. E. Maistrov. *Probability Theory: A Historical Sketch*. Academic Press, New York, 1974.
- [M2] A. A. Markov. *Calculus of Probabilities* [*Ischislenie veroyatnostei*], 3rd ed. St. Petersburg, 1913.
- [M3] L. D. Meshalkin. *Collection of Problems on Probability Theory* [*Sbornik zadach po teorii veroyatnostei*]. Moscow University Press, 1963.
- [M4] P.-A. Meyer, Martingales and stochastic integrals. I. *Lecture Notes in Mathematics*, no. 284. Springer-Verlag, Berlin, 1972.
- [M5] P.-A. Meyer and C. Dellacherie. Probabilities and potential. *North-Holland Mathematical Studies*, no. 29. Hermann, Paris; North-Holland, Amsterdam, 1978.
- [N1] J. Neveu. *Mathematical Foundations of the Calculus of Probability*. Holden-Day, San Francisco, 1965.
- [N2] J. Neveu. *Discrete Parameter Martingales*. North-Holland, Amsterdam, 1975.
- [N3] J. Neyman. *First Course in Probability and Statistics*. Holt, New York, 1950.
- [N4] A. A. Novikov. On estimates and the asymptotic behavior of the probability of nonintersection of moving boundaries by sums of independent random variables. *Math. USSR-Izv.* 17 (1980), 129–145.
- [P1] V. V. Petrov. *Sums of Independent Random Variables*. Springer-Verlag, Berlin, 1975.
- [P2] J. W. Pratt. On interchanging limits and integrals. *Ann. Math. Stat.* 31 (1960), 74–77.

- [P3] Yu. V. Prohorov [Prokhorov]. Asymptotic behavior of the binomial distribution. *Uspekhi Mat. Nauk* 8, no. 3(55) (1953), 135–142 (in Russian).
- [P4] Yu. V. Prohorov. Convergence of random processes and limit theorems in probability theory. *Theory Probab. Appl.* 1 (1956), 157–214.
- [P5] Yu. V. Prokhorov and Yu. A. Rozanov. *Probability theory*. Springer-Verlag, Berlin–New York, 1969; second edition [*Teoriia veroiatnostei*]. “Nauka”, Moscow, 1973.
- [R1] B. Ramachandran. *Advanced Theory of Characteristic Functions*. Statistical Publishing Society, Calcutta, 1967.
- [R2] A. Rényi. *Probability Theory*, North-Holland, Amsterdam, 1970.
- [R3] V. I. Rotar. An extension of the Lindeberg–Feller theorem. *Math. Notes* 18 (1975), 660–663.
- [R4] Yu. A. Rozanov. *Stationary Random Processes*. Holden-Day, San Francisco, 1967.
- [S1] T. A. Sarymsakov. *Foundations of the Theory of Markov Processes* [*Osnovy teorii protsessov Markova*]. GITTL, Moscow, 1954.
- [S2] V. V. Sazonov. Normal approximation: Some recent advances. *Lecture Notes in Mathematics*, no. 879. Springer-Verlag, Berlin–New York, 1981.
- [S3] B. A. Sevyastyanov [Sewastjanow]. *Verzweigungsprozesse*. Oldenbourg, Munich–Vienna, 1975.
- [S4] A. N. Shiryaev. *Random Processes* [*Sluchainye processy*]. Moscow State University Press, 1972.
- [S5] A. N. Shiryaev. *Probability, Statistics, Random Processes* [*Veroyatnost, statistika, sluchainye protsessy*], vols. I and II. Moscow State University Press, 1973, 1974.
- [S6] A. N. Shiryaev. *Statistical Sequential Analysis* [*Statisticheskii posledovatelnyi analiz*]. Nauka, Moscow, 1976.
- [S7] Ya. G. Sinai. *Introduction to Ergodic Theory*. Princeton Univ. Press, Princeton, 1976.
- [S8] S. H. Sirazhdinov. *Limit Theorems for Stationary Markov Chains* [*Predelnye teoremy dlya odnorodnykh tsepeĭ Markova*]. Akad. Nauk Uzbek. SSR, Tashkent, 1955.
- [S9] W. F. Stout. *Almost Sure Convergence*. Academic Press, New York, 1974.
- [T1] I. Todhunter. *A History of the Mathematical Theory of Probability from the Time of Pascal to that of Laplace*. Macmillan, London, 1865.
- [V1] E. Valkeila. A general Poisson approximation theorem. *Stochastics* 7 (1982), 159–171.
- [V2] A. D. Venttsel. *A Course in the Theory of Stochastic Processes*. McGraw-Hill, New York, 1981.
- [W1] B. L. van der Waerden. *Mathematical Statistics*. Springer-Verlag, Berlin–New York, 1969.
- [W2] A. Wald. *Sequential Analysis*. Wiley, New York, 1947.
- [Y1] A. M. Yaglom and I. M. Yaglom. *Probability and Information*. Reidel, Dordrecht, 1983.
- [Z1] S. Zacks. *The Theory of Statistical Inference*. Wiley, New York, 1971.
- [Z2] V. M. Zolotarev. A generalization of the Lindeberg–Feller theorem. *Theory Probab. Appl.* 12 (1967), 606–618.
- [Z3] V. M. Zolotarev. Théorèmes limites pour les sommes de variables aléatoires indépendantes qui ne sont pas infinitésimales. *C.R. Acad. Sci. Paris. Ser. A–B* 264 (1967), A799–A800.
- [D5] D. Dacunha-Castelle and M. Duflo. *Probabilités et statistiques. 1. Problèmes à temps fixe. 2. Problèmes à temps mobile*. Masson, Paris, 1982; *Probability and Statistics*. Springer-Verlag, New York, 1986 (English translation).

- [D6] R. M. Dudley. Distances of probability measures and random variables. *Ann. Math. Statist.* **39** (1968), 1563–1572.
- [G7] P. E. Greenwood and A. N. Shiryaev. *Contiguity and the Statistical Invariance Principle*. Gordon and Breach, New York, 1985.
- [J1] J. Jacod and A. N. Shiryaev. *Limit Theorems for Stochastic Processes*. Springer-Verlag, Berlin Heidelberg, 1987.
- [K12] V. S. Korolyuk. *Aide-mémoire de théorie des probabilités et de statistique mathématique*. Mir, Moscow, 1983; second edition [*Spravochnik po teorii veroyatnostei i matematicheskoi statistike*]. “Nauka”, Moscow, 1985.
- [L10] L. LeCam. *Asymptotic Methods in Statistical Theory*. Springer-Verlag, New York, 1986.
- [L11] F. Liese and I. Vajda. *Convex Statistical Distances*. Teubner, Leipzig, 1987.
- [L12] R. Sh. Liptser and A. N. Shiryaev. *Theory of Martingales*. Kluwer, Dordrecht, Boston, 1989.
- [P6] V. V. Petrov. *Limit Theorems for Sums of Independent Random Variables* [*Predel'nye teoremy dlya summ nezavisimyh sluchaïnykh velichin*]. Nauka, Moscow, 1987.
- [P7] D. Pollard. *Convergence of Stochastic Processes*. Springer-Verlag, New York, 1984.
- [P8] E. L. Presman. Approximation in variation of the distribution of a sum of independent Bernoulli variables with a Poisson law. *Theory of Probability and Its Applications* **30** (1985), no. 2, 417–422.
- [R6] Yu. A. Rozanov. *The Theory of Probability, Stochastic Processes, and Mathematical Statistics* [*Teoriya veroyatnostei, sluchaïnye protsessy i matematicheskaya statistika*]. Nauka, Moscow, 1985.
- [S10] B. A. Sevastyanov. *A Course in the Theory of Probability and Mathematical Statistics* [*Kurs teorii veroyatnostei i matematicheskoi statistiki*]. Nauka, Moscow, 1982.
- [S11] A. N. Shiryaev. *Probability*. Springer-Verlag, Berlin Heidelberg, 1984 (English translation); *Wahrscheinlichkeit*, 1988 (German translation).
- [Z4] V. M. Zolotarev. *Modern theory of summation of random variables* [*Sovremennaya teoriya summirovaniya nezavisimyh sluchaïnykh velichin*]. Nauka, Moscow, 1986.



# Index of Symbols

- $\cup, \bigcup, \cap, \bigcap$  136, 137  
 $\partial$ , boundary 311  
 $\emptyset$ , empty set 11, 136  
 $\oplus$  447  
 $\otimes$  30, 144  
 $\equiv$ , identity, or definition 151  
 $\sim$ , asymptotic equality 20; or equivalence 298  
 $\Rightarrow, \Leftarrow$ , implications 141, 142  
 $\Rightarrow$ , also used for "convergence in general" 310, 311  
 $\stackrel{d}{=}$  342  
 $\stackrel{f}{=}$  316  
 $\preceq$ , finer than 13  
 $\uparrow, \downarrow$  137  
 $\perp$  265, 524  
 $\xrightarrow{\text{a.s.}}, \xrightarrow{\text{a.e.}}, \xrightarrow{d}, \xrightarrow{P}, \xrightarrow{L^p}$  252;  
 $\xrightarrow{w}$  310  
 $\ll, \ll_{\text{loc}}$  524  
 $\langle X, Y \rangle$  483  
 $\{X_n \rightarrow\}$  515  
 $\overline{A}$ , closure of  $A$  153, 311  
 $\bar{A}$ , complement of  $A$  11, 136  
 $[t_1, \dots, t_n]$  combination, unordered set 7  
 $(t_1, \dots, t_n)$  permutation, ordered set 7  
 $A + B$  union of disjoint sets 11, 136  
 $A \triangle B$  43, 136  
 $\{A_n \text{ i.o.}\} - \limsup A_n$  137  
 $a = \min(a, 0); a^+ = \max(a, 0)$  107  
 $a^\oplus = a^{-1}, a \neq 0; 0, a = 0$  462  
 $a \wedge b = \min(a, b); a \vee b = \max(a, b)$  484  
 $A^\oplus$  307  
 $\mathcal{A}$  13  
 $\mathfrak{A}$ , index set 317  
 $\mathcal{B}(R) = \mathcal{B}(R^1) = \mathcal{B} = \mathcal{B}_1$  143, 144  
 $\mathcal{B}(R^n)$  144, 159  
 $\mathcal{B}_0(R^n)$  146  
 $\mathcal{B}(R^\infty)$  146, 160;  $\mathcal{B}_0(R^\infty)$  147  
 $\mathcal{B}(R^T)$  147, 166  
 $\mathcal{B}(C), \mathcal{B}_0(C)$  150  
 $\mathcal{B}(D)$  150  
 $\mathcal{B}[0, 1]$  154  
 $\mathcal{B}(R) \otimes \dots \otimes \mathcal{B}(R) = \mathcal{B}(R^n)$  144  
 $C$  150  
 $C = C[0, \infty)$  151  
 $(C, \mathcal{B}(C))$  150  
 $C_k^i$  6  
 $C^+$  515  
 $\text{cov}(\xi, \eta)$  41, 234  
 $D, (D, \mathcal{B}(D))$  150  
 $\mathcal{D}$  12, 76, 103, 140, 175  
 $E\xi$  37, 180, 181, 182  
 $E(\xi|D), E(\xi|\mathcal{D})$  78  
 $E(\xi|\mathcal{G})$  213, 215, 226

- $E(\xi|\eta)$  81, 221, 238  
 $E(\xi|\eta_1, \dots, \eta_n)$  81  
 $\hat{E}(\xi|\eta_1, \dots, \eta_k)$  264  
 ess sup 261  
 $(f, g), \langle f, g \rangle$  426  
 $F * G$  241  
 $\mathcal{F}$  133, 138  
 $\overline{\mathcal{F}}^P$  154  
 $\mathcal{F}/\mathcal{E}$  176  
 $\mathcal{F}^*, \mathcal{F}_*, \mathcal{F}_A$  139  
 $H_2$  452  
 $h_n(x), H_n(x)$  268, 271  
 inf 44  
 $\mathcal{I}$  163, 425  
 $I_A, I(A)$  33  
 $i \leftrightarrow j$  570  
 $\int_A \xi dP$  183  
 $\int_\Omega \xi dP$  183  
 $(L-S) \int, (R-S) \int, (L) \int, (R) \int$  183,  
 204, 205  
 $L^2$  262  
 $L^p$  261  
 $\mathcal{L}$  264  
 $\overline{\mathcal{L}}$  267  
 $\lim, \underline{\lim}, \limsup, \liminf, \lim \uparrow,$   
 $\lim \downarrow$  137, 173  
 l.i.m. 253  
 $(M)_n$  7  
 $\mathcal{M}$  140  
 $N(A)$  14, 15  
 $N(\mathcal{A})$  13  
 $N(\Omega)$  5  
 $\mathcal{N}(m, \sigma^2)$  234  
 $\mathcal{O}$  235, 265  
 $p(\omega)$  13, 17, 20, 110  
 $\mathcal{P}$  317  
 $P$  133  
 $P(A)$  10, 134  
 $P(A|D), P(A|\mathcal{D})$  24, 76, 212  
 $P(A|\mathcal{G})$  214  
 $P(A|\xi)$  74, 214  
 $P_x, P_\pi$  178  
 $P_C(F)$  310  
 $\mathbb{P}$  115, 275  
 $\mathbb{P}^k$  116  
 $\|p_{ij}\|$  115  
 $\|p(x, y)\|$  113  
 $\mathcal{P} = \{P_\alpha; \alpha \in \mathfrak{A}\}$  317  
 $\mathbb{Q} = (q_1, q_2, \dots)$  585  
 $R$  143  
 $\bar{R}$  144, 173  
 $R^n$  144, 159  
 $\mathbb{R}$  161, 235, 265  
 $(R, \mathcal{B}) = (R^1; \mathcal{B}) = (R, \mathcal{B}(R))$  143,  
 144, 151  
 $R^n, \mathcal{B}(R^n)$  144, 159  
 $R^\infty, \mathcal{B}(R^\infty)$  146, 162  
 $R^T, (R^T, \mathcal{B}(R^T))$  147, 166;  
 $(R^T, \mathcal{B}(R^T), P)$  247  
 $V\xi$  41, 234  
 $V(\xi|\mathcal{D})$  83  
 $V(\xi|\mathcal{G})$  214  
 $X_n^* = \max_{j \leq n} |X_j|$  492  
 $\|X_n\|_p$  492  
 $Z$  177  
 $\mathbb{Z}$  415  
 $\tilde{Z}$  436  
 $Z(\lambda), Z(\Delta)$  424  
 $\alpha(\mathcal{D})$  12  
 $\delta_{ij}$  268  
 $\Delta$  59, 160, 237, 238, 239, 264  
 $\theta_k \xi$  404, 568  
 $\mu, \mu(A)$  132  
 $\mu_1 \times \mu_2$  198  
 $\xi$  32, 170  
 $\xi^+$  44  
 $\tilde{\xi}$  279  
 $\Pi = \|p_{ij}\|$  115  
 $\boxed{\Pi}$  150  
 $(\prod_{t \leq T} \Omega_t, \boxed{\Pi}_{t \in T} \mathcal{F}_t)$  150  
 $\rho(\xi, \eta)$  41, 234  
 $\tilde{\phi}, \phi$  455  
 $\check{\phi}$  459  
 $\Phi(x)$  61, 66  
 $\chi, \chi^2$  156, 243  
 $\chi_B$  174  
 $\Omega$  5  
 $(\Omega, \mathcal{A}, P)$  18, 29  
 $(\Omega, \mathcal{F}, P)$  138

# Index

Also see the Historical and Bibliographical Notes (pp. 597–602), and the Index of Symbols.

- Absolute continuity with respect to  $P$  195, 524
- Absolute moment 182
- Absolutely continuous
  - distributions 155
  - functions 156
  - measures 155
  - probability measures 195, 524
  - random variables 171
- Absorbing 113, 588
- Accessible state 569
- a.e. 185
- Algebra
  - correspondence with decomposition 82
  - induced by a decomposition 12
  - of events 12, 128
  - of sets 132, 139
  - $\sigma$ - 133, 139
  - smallest 140
  - tail 380
- Almost everywhere 185
- Almost invariant 407
- Almost periodic 417
- Almost surely 185
- Amplitude 418
- Aperiodic state 572
- a posteriori* probability 27
- Appropriate sets, principle of 141
- a priori* probability 27
- Arcsine law 102
- Arithmetic properties 569
  - a.s. 185
  - Asymptotic algebra 380
  - Asymptotic properties 573
  - Asymptotically
    - infinitesimal 337
    - unbiased 443
  - Asymptotics 536
  - Atoms 12
  - Attraction of characteristic functions 298
  - Autoregression, autoregressive 419, 421
  - Average time of return 574
  - Averages, moving 419, 421, 437
  - Axioms 138
- Backward equation 117
- Balance equation 420
- Ballot theorem 107
- Banach space 261
- Barriers 588, 590
- Bartlett's estimator 444
- Basis, orthonormal 267
- Bayes's formula 26
- Bayes's theorem 27, 230
- Bernoulli, James 2
  - distribution 155
  - law of large numbers 49
  - random variable 34, 46
  - scheme 30, 45, 55, 70
- Bernstein, S. N. 4, 307
  - inequality 55



- polynomials 54
- proof of Weierstrass' theorem 54
- Berry-Esseen theorem 63, 374
- Bessel's inequality 264
- Best estimator 42, 69. *Also see* Optimal estimator
- Beta distribution 156
- Bilateral exponential distribution 156
- Binary expansion 131, 394
- Binomial distribution 17, 155
  - negative 155
- Binomial random variable 34
- Birkhoff, G. D. 404
- Birthday problem 15
- Bochner-Khinchin theorem viii, 287, 409
- Borel, E. 4
  - algebra 139
  - function 170
  - rectangle 145
  - sets 143, 147
  - space 229
  - zero-or-one law 380
- Borel-Cantelli lemma 255
- Borel-Cantelli-Lévy theorem 518
- Bose-Einstein 10
- Boundary 536
- Bounded variation 207
- Branching process 115
- Brownian motion 306
- Buffon's needle 224
- Bunyakovskii, V. Ya. 38, 192
- Burkholder's inequality 499
  
- Canonical
  - decomposition 544
  - probability space 247
- Cantelli, F. P. 255, 388
- Cantor, G.
  - diagonal process 319
  - function 157, 158
- Carathéodory's theorem 152
- Carleman's test 296
- Cauchy
  - distribution 156, 344
  - inequality 38
  - sequence 253
- Cauchy Bunyakovskii inequality 38, 192
- Cauchy criterion for
  - almost sure convergence 258
  - convergence in probability 259
  - convergence in mean- $p$  260
- Central limit theorem 4, 322, 326, 348
  - for dependent variables 541
  - nonclassical condition for 337
- Certain event 11, 136
- Cesàro limit 582
- Change of variable in integral 196
- Chapman, D. G. 116, 248, 566
- Characteristic function of
  - distribution 4, 274
  - random vector 275
  - set 33
- Charlier, C. V. L. 269
- Chebyshev, P. L. 3, 321
  - inequality 47, 55, 192
- Chi, chi-squared distributions 156, 243
- Class
  - $C^+$  515
  - convergence-determining 315
  - determining 315
  - monotonic 140
  - of states 570
- Classical
  - method 15
  - models 17
- Classification of states 569, 573
- Closed linear manifold 267
- Coin tossing 1, 5, 17, 33, 83, 131
- Coincidence problem 15
- Collectively independent 36
- Combinations
- Communicating states 570
- Compact
  - relatively 317
  - sequentially 319
- Compensator 482
- Complement 11, 136
- Complete
  - function space 260
  - probability measure 154
  - probability space 154
- Completely nondeterministic 447
- Conditional
  - distribution 227
  - probability 23, 77, 221
    - regular 227
  - with respect to a
    - decomposition 77, 212
    - $\sigma$ -algebra 212, 214
    - random variable 77, 214
  - variance 83, 214
  - Wiener process 307
- Conditional expectation
  - in the wide sense 264, 274
  - with respect to
    - decomposition 78
    - event 220

- set of variables 81
  - $\sigma$ -algebra 214
- Conditionally Gaussian 466
- Confidence interval 74
- Consistency property 163, 246
- Consistent estimator 71, 521, 535
- Construction of a process 245, 246
- Continuity theorem 322
- Continuous at  $\emptyset$  153, 164
- Continuous from above or below 134
- Continuous time 177
- Convergence-determining class 315
- Convergence of
  - martingales and submartingales 508, 515
  - probability measures Chap. III, 308
  - random variables: equivalences 252
  - sequences. *See* Convergence of sequences
  - series 384
- Convergence of sequences
  - almost everywhere 252, 353
  - almost sure 252, 353
  - at points of continuity 253
  - dominated 187
  - in distribution 252, 325, 353
  - in general 310
  - in mean 252
  - in mean of order  $p$  252
  - in mean square 252
  - in measure 252
  - in probability 252, 348
  - monotone 186
  - weak 309, 311
  - with probability 252
- Convolution 241, 377
- Coordinate method 247
- Correlation
  - coefficient 41, 234
  - function 416
  - maximal 244
- Counting measure 233
- Covariance 41, 232, 293
  - function 306, 416
  - matrix 235
- Cramér condition 400
- Cramér–Lundberg model 559
- Cramér transform 401
- Cramér–Wold method 549
- Cumulant 290
- Curve of regression 238
- Curvilinear boundary 536
- Cyclic property 571
- Cyclic subclasses 571
- Cylinder set 146
- Davis's inequality 499
- Decomposition
  - canonical 544
  - countable 140
  - Doob 482
  - Krickeberg 507
  - Lebesgue 525
  - of martingale 507
  - of  $\Omega$  12, 140
  - of probability measure 525
  - of random sequence 447
  - of set 12, 292
  - of submartingale 482
  - trivial 80
- Degenerate
  - distribution 298
  - distribution function 298
  - random variable 298
- Delta function 298
- Delta, Kronecker 268
- De Moivre, A. 2, 49
- De Moivre–Laplace limit theorem 62
- Density
  - Gaussian 66, 156, 161, 238
  - $n$ -dimensional 161
  - of distribution 156
  - of measure with respect to a measure 196
  - of random variable 171
- Dependent random variables 103
  - central limit theorem for 541
- Derivative, Radon–Nikodým 196
- Detection of signal 462
- Determining class 315
- Deterministic 447
  - regularity 1
- Dichotomy
  - Hájek–Feldman 533
  - Kakutani 529
- Difference of sets 11, 136
- Direct product 31, 144, 151
- Dirichlet's function 211
- Discrete
  - measure 155
  - random variable 171
  - time 177
  - uniform density 155
- Discrete version of Ito's formula 554
- Disjoint 136
- Dispersion 41
- Distance in variation 355, 376
- Distribution. *Also see* Distribution function, Probability distribution
  - Bernoulli 155
  - beta 156

- binomial 17, 18, 155
- Cauchy 156, 344
- chi, chi-squared 243
- conditional 212
- discrete (list) 155
- discrete uniform 155
- entropy of 51
- ergodic 118
- exponential 156
- gamma 156
- Gaussian 66, 156, 161, 293
- geometric 155
- hypergeometric 21
- infinitely divisible 341
- initial 112, 565
- invariant 120
- limit 545
- lognormal 240
- multidimensional 160
- multinomial 20
- negative binomial 155
- normal 66, 156
- of process 178
- of sum 36
- Poisson 64, 155
- polynomial 20
- probability 33
- stable 341
- stationary 120
- Student's 155, 244
- $t$ - 155, 244
- two-sided exponential 155
- uniform 156
- with density (list) 156
- Distribution function 34, 35, 152, 171
  - absolutely continuous 156
  - degenerate 288
  - discrete 155
  - finite-dimensional 246
  - generalized 158
  - $n$ -dimensional 160
  - of functions of random variables 36, 239ff.
  - of sum 36, 241
- Distribution of objects in cells 8
- $\mathcal{D}$ -measurable 76
- Dominated convergence 187
- Dominated sequence 496
- Doob, J. L. 482, 485, 492
- Doubling stakes 89, 481
- Doubly stochastic 587
- $d$ -system 142
- Duration of random walk 90
- Dvoretzky's inequality 508
- Efficient estimator 71
- Eigenvalue 130
- Electric circuit 32
- Elementary
  - events 5, 136
  - probability theory Chap. I
  - stochastic measure 424
- Empty set 11, 136
- Entropy 51
- Equivalent measures 524
- Ergodic
  - sequence 407
  - theorems 110, 409, 413
    - maximal 410
    - mean-square 438
  - theory 409
  - transformation 408
- Ergodicity 118, 409, 581
- Errors
  - laws of 298
  - mean-square 43
  - of observation 2, 3
- Esseen's inequality 296
- Essential state 569
- Essential supremum 261
- Estimation 70, 237, 440, 454
- Estimator 42, 70, 237
  - Bartlett's 444
  - best 42, 69
  - consistent 71
  - efficient 71
  - for parameters 472, 520, 535
  - linear 43
  - of spectral quantities 442
  - optimal 70, 237, 303, 454, 461, 463, 469
  - Parzen's 445
  - unbiased 71, 440
  - Zhurbenko's 445
- Events 5, 10, 136
  - certain 11, 136
  - elementary 5
  - impossible 11, 136
  - independent 28, 29
  - mutually exclusive 136
  - symmetric 382
- Existence of limits and stationary
  - distributions 582ff.
- Expectation 37, 182
  - inequalities for 192, 193
  - of function 55
  - of maximum 45
  - of random variable with respect to
    - decomposition 76
    - set of random variables 81

- $\sigma$ -algebra 212
- of sum 38
- Expected value 37
- Exponential distribution 156
- Exponential random variable 156, 244, 245
- Extended random variable 178
- Extension of a measure 150, 163, 249, 427
- Extrapolation 453
- Fair game 480
- Fatou's lemma 187, 211
- Favorable game 89, 480
- $F$ -distribution 156
- Feldman, J. 533
- Feller, W. 597
- Fermat, P. de 2
- Fermi-Dirac 10
- Filter 434, 464
  - physically realizable 451
- Filtering 453, 464
- Finer decomposition 80
- Finite second moment 262
- Finite-dimensional distribution function 246
- Finitely additive 132, 424
- First
  - arrival 129, 574
  - exist 123
  - return 94, 129, 574
- Fisher's information 72
- $\mathcal{F}$ -measurable 170
- Forward equation 117
- Foundations Chap. II, 131
- Fourier transform 276
- Frequencies 418
- Frequency 46
- Frequency characteristic 434
- Fubini's theorem 198
- Fundamental inequalities (for martingales) 492
- Fundamental sequence 253, 258
- Gamma distribution 156, 343
- Garcia, A. viii, 410
- Gauss, C. F. 3
- Gaussian
  - density 66, 156, 161, 236
  - distribution 66, 156, 161, 293
  - measure 268
  - random variables 234, 243, 298
  - random vector 299
  - sequence 306, 413, 439, 441, 466
  - systems 297, 305
- Gauss Markov process 307
- Generalized
  - Bayes theorem 231
  - distribution function 158
  - martingale 476
  - submartingale 476, 523
- Geometric distribution 155
- Gnedenko, B. V. vii, 510, 542
- Gram determinant 265
- Gram-Schmidt process 266
- Haar functions 271, 482
- Hájek-Feldman dichotomy 533
- Hardy class 452
- Harmonics 418
- Hartman, P. 372
- Hellinger integral 3
- Helly's theorem 319
- Herglotz, G. 421
- Hermite polynomials 268
- Hewitt, E. 382
- Hilbert space 262
  - complex 275, 416
  - separable 267
  - unitary 416
- Hinchin. *See* Khinchin
- History 597-602
- Hölder inequality 193
- Huygens, C. 2
- Hydrology 420, 421
- Hypergeometric distribution 21
- Hypotheses 27
- Impossible event 11, 136
- Impulse response 434
- Increasing sequence 137
- Increments
  - independent 306
  - uncorrelated 109, 306
- Indecomposable 580
- Independence 27
  - linear 265, 286
- Independent
  - algebras 28, 29
  - events 28, 29
  - functions 179
  - increments 306
  - random variables 36, 77, 81, 179, 380, 513
- Indicator 33, 43

## Inequalities

- Berry–Esseen 333
  - Bernstein 55
  - Bessel 264
  - Burkholder 499
  - Cauchy–Bunyakovskii 38, 192
  - Cauchy–Schwarz 38
  - Chebyshev 3, 321
  - Davis 499
  - Dvoretzky 508
  - Hölder 193
  - Jensen 192, 233
  - Khinchin 347, 498
  - Kolmogorov 496
  - Lévy 400
  - Lyapunov 193
  - Marcinkiewicz–Zygmund 498
  - Markov 598
  - martingale 492
  - Minkowski 194
  - nonuniform 376
  - Ottaviani 507
  - Rao–Cramér 73
  - Schwarz 38
  - two-dimensional Chebyshev 55
- Inessential state 569
- Infinitely divisible 341
- Infinitely many outcomes 131
- Information 72
- Initial distribution 112, 565
- Innovation sequence 448
- Insurance 558
- Integral
- Lebesgue 180
  - Lebesgue–Stieltjes 183
  - Riemann 183, 205
  - Riemann–Stieltjes 205
  - stochastic 423
- Integral equation 208
- Integral theorem 62
- Integration by
- parts 206
  - substitution 211
- Intensity 418
- Interpolation 453
- Intersection 11, 136
- Introducing probability measures 151
- Invariant set 407, 413
- Inversion formulas 283, 295
- i.o. 137
- Ionescu Tulcea, C. T. vii, 249
- Ising model 23
- Isometry 430
- Iterated logarithm 395

## Ito's formula for

- Brownian motion 558

## Jensen's inequality 192, 233

- Kakutani dichotomy 527, 528
- Kakutani–Hellinger distance 363
- Kalman–Bucy filter 464
- Khinchin, A. Ya 287, 468
- Kolmogorov, A. N. vii, 3, 4, 384, 395, 498, 542
  - axioms 131
  - inequality 384
- Kolmogorov–Chapman equation 116, 248, 566
- Kolmogorov's theorems
  - convergence of series 384
  - existence of process 246
  - extension of measures 167
  - iterated logarithm 395
  - stationary sequences 453, 455
  - strong law of large numbers 366, 389, 391
  - three-series theorem 387
  - two-series theorem 386
  - zero-or-one law 381
- Krickeberg's decomposition 507
- Kronecker, L. 390
  - delta 268
- Kullback information 368

 $\Lambda$  (condition) 338

- Laplace, P. S. 2, 55
- Large deviation 69, 402
- Law of large numbers 45, 49, 325
  - for Markov chains 122
  - for square-integrable martingales 519
- Poisson's 599
- strong 388
- Law of the iterated logarithm 395
- Least squares 3
- Lebesgue, H.
  - decomposition 366, 525
  - derivative 366
  - dominated convergence theorem 187
  - integral 180, 181
    - change of variable in 196
  - measure 154, 159
- Lebesgue–Stieltjes integral 197
- Lebesgue–Stieltjes measure 158, 205
- LeCam, L. 377
- Lévy, P.

- convergence theorem 510
- distance 316
- inequality 400
- Lévy–Khinchin representation 347
- Lévy–Khinchin theorem 344
- Lévy–Prokhorov metric 349
- Likelihood ratio 110
- lim inf, lim sup 137
- Limit theorems 55
- Limits under
  - expectation signs 180
  - integral signs 180
- Lindeberg condition 328
- Lindeberg–Feller theorem 334
- Linear manifold 264
- Linearly independent 265
- Liouville's theorem 406
- Lipschitz condition 512
- Local limit theorem 55, 56
- Local martingale, submartingale 477
- Locally absolutely continuous 524
- Locally bounded variation 206
- Lognormal 240
- Lottery 15, 22
- Lower function 396
- $L^2$ -theory Chap. VI, 415
- Lyapunov, A. M. 3, 322
  - condition 332
  - inequality 193
- Macmillan's theorem 59
- Marcinkiewicz's theorem 288
- Marcinkiewicz–Zygmund inequality 498
- Markov, A. A. viii, 3, 321
  - dependence 564
  - process 248
  - property 112, 127, 564
  - time 476
- Markov chains 110, Chap. VIII, 251, 564
  - classification of 569, 573
  - discrete 565
  - examples 113, 587
  - finite 565
  - homogeneous 113, 565
- Martingale 103, Chap. VII, 474
  - convergence of 508
  - generalized 476
  - inequalities for 492
  - in gambling 480
  - local 477
  - oscillations of 503
  - reversed 484
  - sets of convergence for 515
  - square-integrable 482, 493, 538
  - uniformly integrable 512
- Martingale-difference 481, 543, 559
- Martingale transform 478
- Mathematical expectation 37, 76. *Also see* Expectation
- Mathematical foundations Chap. II, 131
- Matrix
  - covariance 235
  - doubly stochastic 587
  - of transition probabilities 112
  - orthogonal 235, 265
  - stochastic 113
  - transition 112
- Maximal
  - correlation coefficient 244
  - ergodic theorem 410
- Maxwell–Boltzmann 10
- Mean
  - duration 90, 489
  - square 42
  - ergodic theorem 438
  - value 37
  - vector 301
- Measurable
  - function 170
  - random variable 80
  - set 154
  - spaces 133
    - $(C, \mathcal{B}(C))$  150
    - $(D, \mathcal{B}(D))$  150
    - $(\prod \Omega_i, \prod \mathcal{F}(i))$  150
    - $(R, \mathcal{B}(R))$  151
    - $(R^\infty, \mathcal{B}(R^\infty))$  162
    - $(R^n, \mathcal{B}(R^n))$  159
    - $(R^T, \mathcal{B}(R^T))$  166
  - transformation 404
- Measure 133
  - absolutely continuous 155, 524
  - complete 154
  - consistent 567
  - countably additive 133
  - counting 233
  - discrete 155
  - elementary 424
  - extending a 152, 249, 427
  - finite 132
  - finitely additive 132
  - Gaussian 268
  - Lebesgue 154
  - Lebesgue–Stieltjes 158
  - orthogonal 366, 425, 524
  - probability 134
  - restriction of 165

- $\sigma$ -additive 133
- $\sigma$ -finite 133
- signed 196
- singular 158, 366, 524
- stochastic 423
- Wiener 169
- Measure-preserving transformations 404ff.
- Median 44
- Method
  - of characteristic functions 321ff.
  - of moments 4, 321
- Metrically transitive 407
- Minkowski inequality 194
- Mises, R. von 4
- Mixed model 421
- Mixed moment 289
- Mixing 409
- Moivre. *See* De Moivre
- Moment 182
  - absolute 182
  - and semi-invariant 290
  - method 4
  - mixed 289
  - problem 294ff.
- Monotone convergence theorem 186
- Monotonic class 140
- Monte Carlo method 225, 394
- Moving averages 418, 421
- Multinomial distribution 20, 21
- Multiplication formula 26
- Mutual variation 483
  
- Needle (Buffon) 224
- Negative binomial 155
- Noise 418, 435
- Nonclassical hypotheses 328, 337
- Nondeterministic 447
- Nonlinear estimator 453
- Nonnegative definite 235
- Nonrecurrent state 574
- Norm 260
- Normal
  - correlation 303, 307
  - density 66, 161
  - distribution function 62, 66, 156, 161
  - number 394
- Normally distributed 299
- Null state 574
  
- Occurrence of event 136
- Optimal estimator 71, 237, 303, 454, 461, 463, 469
- Optional stopping 601
- Ordered sample 6
- Orthogonal
  - decomposition 265
  - increments 428
  - matrix 235, 265
  - random variables 263
  - stochastic measures 423, 425
  - system 263
- Orthogonalization 266
- Orthonormal 263, 267, 271
- Oscillations of
  - submartingales 503
  - water level 420
- Ottaviani's inequality 507
- Outcome 5, 136
  
- Pairwise independence 29
- P-almost surely, almost everywhere 185
- Parallelogram property 274
- Parseval's equation 268
- Partially observed sequences 460ff.
- Parzen's estimator 445
- Pascal, B. 2
- Path 48, 85, 95
- Pauli exclusion principle 10
- Period of Markov chain 571
- Periodogram 443
- Permutation 7, 382
- Perpendicular 265
- Phase space 112, 565
- Physically realizable 434, 451
- $\pi$  225
- Poincaré recurrence principle 406
- Poisson, D. 3
  - distribution 64, 155
  - law of large numbers 599
  - limit theorem 64, 327
- Poisson-Charlier polynomials 269
- Pólya's theorems
  - characteristic functions 287
  - random walk 595
- Polynomials
  - Bernstein 54
  - Hermite 268
  - Poisson-Charlier 269
- Positive semi-definite 287
- Positive state 574
- Pratt's lemma 211, 599
- Predictable sequence 446, 474
- Predictable quadratic variation characteristic 483
- Preservation of martingale property 484

- Principle of appropriate sets 141
- Probabilistic model 5, 14, 131
  - in the extended sense 133
- Probability 2, 134
  - a posteriori, a priori* 27
  - classical 15
  - conditional 23, 76, 214
  - finitely additive 132
  - measure 131, 151
  - multiplication 26
  - of first arrival or return 574
  - of mean duration 90
  - of ruin 83
  - of success 70
  - total 25, 77, 79
  - transition 566
- Probability distribution 33, 170, 178
  - discrete 155
  - lognormal 240
  - stationary 569
  - table 155, 156
- Probability of ruin in insurance 558
- Probability measure 134, 154, 524
  - absolutely continuous 524
  - complete 154
- Probability space 14, 138
  - canonical 247
  - complete 154
  - universal 252
- Probability of error 361
- Problems on
  - arrangements 8
  - coincidence 15
  - ruin 88
- Process
  - branching 115
  - Brownian motion 306
  - construction of 245ff.
  - Gaussian 306
  - Gauss–Markov 307
  - Markov 248
  - stochastic 4, 177
  - Wiener 306, 307
  - with independent increments 306
- Prohorov, Yu. V. vii, 64, 318
- Projection 265, 273
- Pseudoinverse 307
- Pseudotransform 462
- Purely nondeterministic 447
- Pythagorean property 274
- Quadratic characteristic 483
- Quadratic covariation 483
- Quadratic variation 483
- Queueing theory 114
- Rademacher system 271
- Radon–Nikodym
  - derivative 196
  - theorem 196, 599
- Random
  - elements 176ff.
  - function 177
  - process 177, 306
    - with orthogonal increments 428
  - sequences 4, Chap. V, 404
    - existence of 246, 249
    - orthogonal 447
- Random variables 32ff., 166, 234ff.
  - absolutely continuous 171
  - almost invariant 407
  - complex 177
  - continuous 171
  - degenerate 298
  - discrete 171
  - exponential 156, 244, 245
  - E*-valued 177
  - extended 173
  - Gaussian 234, 243, 298
  - invariant 407
  - normally distributed 234
  - simple 170
  - uncorrelated 234
- Random vectors 35, 177
  - Gaussian 299, 301
- Random walk 18, 83, 381
  - in two and three dimensions 592
  - simple 587
  - symmetric 94, 381
  - with curvilinear boundary 536
- Rao–Cramér inequality 72
- Rapidity of convergence 373, 376, 400, 402
- Realization of a process 178
- Recurrent state 574, 593
- Reflecting barrier 592
- Reflection principle 94, 96
- Regression 238
- Regular
  - conditional distribution 227
  - conditional probability 226
  - stationary sequence 447
- Relatively compact 317
- Reliability 74
- Restriction of a measure 165
- Reversed martingale 105, 403, 484
- Reversed sequence 130



- Riemann integral 204
- Riemann–Stieltjes integral 204
- Ruin 84, 87, 489
- Sample points, space 5
- Sampling
  - with replacement 6
  - without replacement 7, 21, 23
- Savage, L. J. 389
- Scalar product 263
  - Schwarz inequality 38. *Also see*  
Bunyakovskii, Cauchy
- Semicontinuous 313
- Semi-definite 287
- Semi-invariant 290
- Semi-norm 260
- Separable 267
- Sequences
  - almost periodic 417
  - moving average 418
  - of independent random variables 379
  - partially observed 460
  - predictable 446, 474
  - random 176, 404
  - regular 447
  - singular 447
  - stationary (strict sense) 404
  - stationary (wide sense) 416
  - stochastic 474, 483
- Sequential compactness 318
- Series of random variables 384
- Sets of convergence 515
- Shifting operators 568
- $\sigma$ -additive 134
- Sigma algebra 133, 138
  - asymptotic 380
  - generated by  $\xi$  174
  - tail, terminal 380
- Signal, detection of 462
- Significance level 74
- Simple
  - moments 291
  - random variable 32
  - random walk 587
  - semi-invariants 291
- Singular measure 158
- Singular sequence 447
- Singularity of distributions 524
- Skorohod, A. V. 150
- Slowly varying 537
- Spectral
  - characteristic 434
  - density 418
  - rational 437, 456
  - function 422
  - measure 422
  - representation of
    - covariance function 415
    - sequences 429
    - window 444
- Spectrum 418
- Square-integrable martingale 482, 493, 518, 538
- Stable 344
- Standard deviation 41, 234
- State space 112
- States, classification of 234, 569, 573
- Stationary
  - distribution 120, 569, 580
  - Markov chain 110
  - sequence Chap. V, 404; Chap. VI, 415
- Statistical estimation
  - regularity 440
- Statistically independent 28
- Statistics 4, 50
- Stieltjes, T. J. 183, 204
- Stirling's formula 20, 22
- Stochastic
  - exponential 504
  - integral 423, 426
  - matrix 113, 587
  - measure 403, 424
    - extension of 427
  - orthogonal 425
    - with orthogonal values 425, 426
  - process 4, 177
  - sequence 474, 564
- Stochastically independent 42
- Stopped process 477
- Stopping time 84, 105, 476
- Strong law of large numbers 388, 389, 501, 515
- Strong Markov property 127
- Structure function 425
- Student distribution 156, 244
- Submartingales 475
  - convergence of 508
  - generalized 476, 515
  - local 477
  - nonnegative 509
  - nonpositive 509
  - sets of convergence of 515
  - uniformly integrable 510
- Substitution, integration by 211
- Sum of
  - dependent random variables 591
  - events 11, 137

- exponential random variables 245
- Gaussian random variables 243
- independent random variables 328, Chap. IV, 379
- Poisson random variables 244
- sets 11, 136
- Summation by parts 390
- Supermartingale 475
- Symmetric difference  $\Delta$  43, 136
- Symmetric events 382
- Szegő–Kolmogorov formula 464
- Tables
  - continuous densities 156
  - discrete densities 155
  - terms in set theory and probability 136, 137
- Tail 49, 323, 335
  - algebra 380
- Taxi stand 114
- $t$ -distribution 156, 244
- Terminal algebra 380
- Three-series theorem 387
- Tight 318
- Time
  - change (in martingale) 484
  - continuous 177
  - discrete 177
  - domain 177
- Toeplitz, O. 390
- Total probability 25, 77, 79
- Trajectory 178
- Transfer function 434
- Transform 478
- Transformation, measure-preserving 405
- Transition
  - function 565
  - matrix 112
  - probabilities 112, 248, 566
- Trial 30
- Trivial algebra 12
- Tulcea. *See* Ionescu Tulcea
- Two-dimensional Gaussian density 162
- Two-series theorem 386
- Typical
  - path 50, 52
  - realization 50
- Unbiased estimator 71, 440
- Uncorrelated 42, 234
- increments 109
- Unfavorable game 86, 89, 480
- Uniform distribution 155, 156
- Uniformly
  - asymptotically infinitesimal 337
  - continuous 328
  - integrable 188
- Union 11, 136, 137
- Uniqueness of
  - distribution function 282
  - solution of moment problem 295
- Universal probability space 252
- Unordered
  - samples 6
  - sets 166
- Upper function 396
- Variance 41
  - conditional 83
  - of sum 42
- Variation quadratic 483
- Vector
  - Gaussian 238
  - random 35, 177, 238, 301
- Wald's identities 107, 488, 489
- Water level 421
- Weak convergence 309
- Weierstrass approximation theorem
  - for polynomials 54
  - for trigonometric polynomials 282
- White noise 418, 435
- Wiener, N.
  - measure 169
  - process 306, 307
- Window, spectral 444
- Wintner, A. 397
- Wold's
  - expansion 446, 450
  - method 549
- Wolf, R. 225
- Zero-or-one laws 354ff., 379, 512
  - Borel 380
  - for Gaussian sequences 533
  - Hewitt–Savage 382
  - Kolmogorov 381
- Zhurbenko's estimator 445
- Zygmund, A. 498

# Graduate Texts in Mathematics

continued from page ii

- 61 WHITEHEAD. Elements of Homotopy Theory.
- 62 KARGAPOLOV/MERLZIAKOV. Fundamentals of the Theory of Groups.
- 63 BOLLOBAS. Graph Theory.
- 64 EDWARDS. Fourier Series. Vol. I 2nd ed.
- 65 WELLS. Differential Analysis on Complex Manifolds. 2nd ed.
- 66 WATERHOUSE. Introduction to Affine Group Schemes.
- 67 SERRE. Local Fields.
- 68 WEIDMANN. Linear Operators in Hilbert Spaces.
- 69 LANG. Cyclotomic Fields II.
- 70 MASSEY. Singular Homology Theory.
- 71 FARKAS/KRA. Riemann Surfaces. 2nd ed.
- 72 STILLWELL. Classical Topology and Combinatorial Group Theory. 2nd ed.
- 73 HUNGERFORD. Algebra.
- 74 DAVENPORT. Multiplicative Number Theory. 2nd ed.
- 75 HOCHSCHILD. Basic Theory of Algebraic Groups and Lie Algebras.
- 76 ITAKA. Algebraic Geometry.
- 77 HECKE. Lectures on the Theory of Algebraic Numbers.
- 78 BURRIS/SANKAPPANAVAR. A Course in Universal Algebra.
- 79 WALTERS. An Introduction to Ergodic Theory.
- 80 ROBINSON. A Course in the Theory of Groups. 2nd ed.
- 81 FORSTER. Lectures on Riemann Surfaces.
- 82 BOTT/TU. Differential Forms in Algebraic Topology.
- 83 WASHINGTON. Introduction to Cyclotomic Fields. 2nd ed.
- 84 IRELAND/ROSEN. A Classical Introduction to Modern Number Theory. 2nd ed.
- 85 EDWARDS. Fourier Series. Vol. II. 2nd ed.
- 86 VAN LINT. Introduction to Coding Theory. 2nd ed.
- 87 BROWN. Cohomology of Groups.
- 88 PIERCE. Associative Algebras.
- 89 LANG. Introduction to Algebraic and Abelian Functions. 2nd ed.
- 90 BRØNDSTED. An Introduction to Convex Polytopes.
- 91 BEARDON. On the Geometry of Discrete Groups.
- 92 DIESTEL. Sequences and Series in Banach Spaces.
- 93 DUBROVIN/FOMENKO/NOVIKOV. Modern Geometry—Methods and Applications. Part I. 2nd ed.
- 94 WARNER. Foundations of Differentiable Manifolds and Lie Groups.
- 95 SHIRYAEV. Probability. 2nd ed.
- 96 CONWAY. A Course in Functional Analysis. 2nd ed.
- 97 KOBLITZ. Introduction to Elliptic Curves and Modular Forms. 2nd ed.
- 98 BRÖCKER/TOM DIECK. Representations of Compact Lie Groups.
- 99 GROVE/BENSON. Finite Reflection Groups. 2nd ed.
- 100 BERG/CHRISTENSEN/RESSEL. Harmonic Analysis on Semigroups: Theory of Positive Definite and Related Functions.
- 101 EDWARDS. Galois Theory.
- 102 VARADARAJAN. Lie Groups, Lie Algebras and Their Representations.
- 103 LANG. Complex Analysis. 3rd ed.
- 104 DUBROVIN/FOMENKO/NOVIKOV. Modern Geometry—Methods and Applications. Part II.
- 105 LANG.  $SL_2(\mathbb{R})$ .
- 106 SILVERMAN. The Arithmetic of Elliptic Curves.
- 107 OLVER. Applications of Lie Groups to Differential Equations. 2nd ed.
- 108 RANGE. Holomorphic Functions and Integral Representations in Several Complex Variables.
- 109 LEHTO. Univalent Functions and Teichmüller Spaces.
- 110 LANG. Algebraic Number Theory.
- 111 HUSEMÖLLER. Elliptic Curves.
- 112 LANG. Elliptic Functions.
- 113 KARATZAS/SHREVE. Brownian Motion and Stochastic Calculus. 2nd ed.
- 114 KOBLITZ. A Course in Number Theory and Cryptography. 2nd ed.
- 115 BERGER/GOSTIAUX. Differential Geometry: Manifolds, Curves, and Surfaces.
- 116 KELLEY/SRINIVASAN. Measure and Integral. Vol. I.
- 117 SERRE. Algebraic Groups and Class Fields.
- 118 PEDERSEN. Analysis Now.
- 119 ROTMAN. An Introduction to Algebraic Topology.
- 120 ZIEMER. Weakly Differentiable Functions: Sobolev Spaces and Functions of Bounded Variation.
- 121 LANG. Cyclotomic Fields I and II. Combined 2nd ed.
- 122 REMMERT. Theory of Complex Functions. *Readings in Mathematics*

- 123 EBBINGHAUS/HERMES et al. Numbers.  
*Readings in Mathematics*
- 124 DUBROVIN/FOMENKO/NOVIKOV. Modern  
Geometry—Methods and Applications.  
Part III.
- 125 BERENSTEIN/GAY. Complex Variables: An  
Introduction.
- 126 BOREL. Linear Algebraic Groups. 2nd ed.
- 127 MASSEY. A Basic Course in Algebraic  
Topology.
- 128 RAUCH. Partial Differential Equations.
- 129 FULTON/HARRIS. Representation Theory:  
A First Course.  
*Readings in Mathematics*
- 130 DODSON/POSTON. Tensor Geometry.
- 131 LAM. A First Course in Noncommutative  
Rings.
- 132 BEARDON. Iteration of Rational Functions.
- 133 HARRIS. Algebraic Geometry: A First  
Course.
- 134 ROMAN. Coding and Information Theory.
- 135 ROMAN. Advanced Linear Algebra.
- 136 ADKINS/WEINTRAUB. Algebra: An  
Approach via Module Theory.
- 137 AXLER/BOURDON/RAMEY. Harmonic  
Function Theory.
- 138 COHEN. A Course in Computational  
Algebraic Number Theory.
- 139 BREDON. Topology and Geometry.
- 140 AUBIN. Optima and Equilibria. An  
Introduction to Nonlinear Analysis.
- 141 BECKER/WEISFENNING/KREDEL. Gröbner  
Bases. A Computational Approach to  
Commutative Algebra.
- 142 LANG. Real and Functional Analysis.  
3rd ed.
- 143 DOOB. Measure Theory.
- 144 DENNIS/FARB. Noncommutative  
Algebra.
- 145 VICK. Homology Theory. An  
Introduction to Algebraic Topology.  
2nd ed.
- 146 BRIDGES. Computability: A  
Mathematical Sketchbook.
- 147 ROSENBERG. Algebraic  $K$ -Theory  
and Its Applications.
- 148 ROTMAN. An Introduction to the  
Theory of Groups. 4th ed.
- 149 RATCLIFFE. Foundations of  
Hyperbolic Manifolds.
- 150 EISENBUD. Commutative Algebra  
with a View Toward Algebraic  
Geometry.
- 151 SILVERMAN. Advanced Topics in  
the Arithmetic of Elliptic Curves.
- 152 ZIEGLER. Lectures on Polytopes.
- 153 FULTON. Algebraic Topology: A  
First Course.
- 154 BROWN/PEARCY. An Introduction to  
Analysis.
- 155 KASSEL. Quantum Groups.
- 156 KECHRIS. Classical Descriptive Set  
Theory.
- 157 MALLIAVIN. Integration and  
Probability.
- 158 ROMAN. Field Theory.
- 159 CONWAY. Functions of One  
Complex Variable II.
- 160 LANG. Differential and Riemannian  
Manifolds.
- 161 BORWEIN/ERDÉLYI. Polynomials and  
Polynomial Inequalities.
- 162 ALPERIN/BELL. Groups and  
Representations.
- 163 DIXON/MORTIMER. Permutation  
Groups.
- 164 NATHANSON. Additive Number Theory:  
The Classical Bases.
- 165 NATHANSON. Additive Number Theory:  
Inverse Problems and the Geometry of  
Sumsets.
- 166 SHARPE. Differential Geometry: Cartan's  
Generalization of Klein's Erlangen  
Program.
- 167 MORANDI. Field and Galois Theory.
- 168 EWALD. Combinatorial Convexity and  
Algebraic Geometry.
- 169 BHATIA. Matrix Analysis.
- 170 BREDON. Sheaf Theory. 2nd ed.
- 171 PETERSEN. Riemannian Geometry.
- 172 REMMERT. Classical Topics in Complex  
Function Theory.
- 173 DIESTEL. Graph Theory.
- 174 BRIDGES. Foundations of Real and  
Abstract Analysis.
- 175 LICKORISH. An Introduction to Knot  
Theory.
- 176 LEE. Riemannian Manifolds.
- 177 NEWMAN. Analytic Number Theory.
- 178 CLARKE/LEDYAEV/STERN/WOLENSKI.  
Nonsmooth Analysis and Control  
Theory.
- 179 DOUGLAS. Banach Algebra Techniques in  
Operator Theory. 2nd ed.
- 180 SRIVASTAVA. A Course on Borel Sets.
- 181 KRESS. Numerical Analysis.
- 182 WALTER. Ordinary Differential  
Equations.
- 183 MEGGINSON. An Introduction to Banach  
Space Theory.
- 184 BOLLOBAS. Modern Graph Theory.
- 185 COX/LITTLE/O'SHEA. Using Algebraic  
Geometry.
- 186 RAMAKRISHNAN/VALENZA. Fourier  
Analysis on Number Fields.

This book contains a systematic treatment of probability from the ground up, starting with intuitive ideas and gradually developing more sophisticated subjects, such as random walks, martingales, Markov chains, ergodic theory, weak convergence of probability measures, stationary stochastic processes, and the Kalman-Bucy filter. Many examples are discussed in detail, and there are a large number of exercises. The book is accessible to advanced undergraduates and can be used as a text for self-study.

This new edition contains substantial revisions and updated references. The reader will find a deeper study of topics such as the distance between probability measures, metrization of weak convergence, and contiguity of probability measures. Proofs for a number of some important results which were merely stated in the first edition have been added. The author has included new material on the probability of large deviations, on the central limit theorem for sums of dependent random variables, and on a discrete version of Ito's formula.

ISBN 0-387-94549-0

