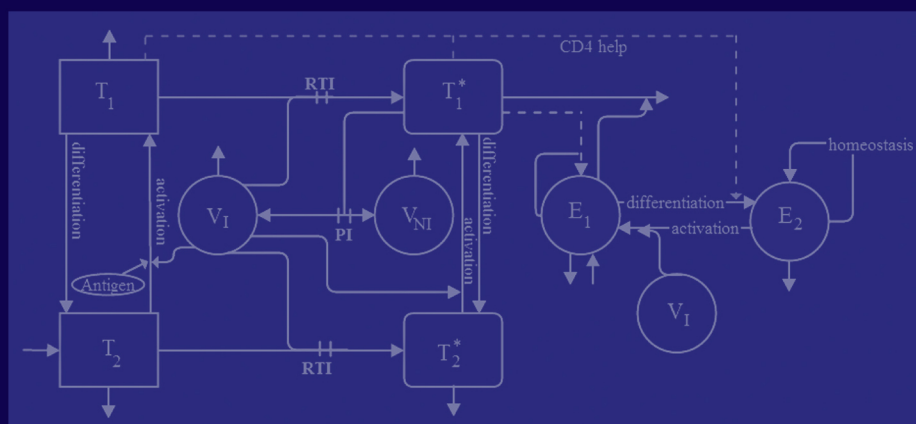


# Modeling and Inverse Problems in the Presence of Uncertainty



H. T. Banks, Shuhua Hu,  
and W. Clayton Thompson

# Modeling and Inverse Problems in the Presence of Uncertainty

# MONOGRAPHS AND RESEARCH NOTES IN MATHEMATICS

## Series Editors

John A. Burns

Thomas J. Tucker

Miklos Bona

Michael Ruzhansky

Chi-Kwong Li

---

## Published Titles

*Iterative Optimization in Inverse Problems*, Charles L. Byrne

*Modeling and Inverse Problems in the Presence of Uncertainty*, H. T. Banks, Shuhua Hu,  
and W. Clayton Thompson

## Forthcoming Titles

*Sinusoids: Theory and Technological Applications*, Prem K. Kythe

*Stochastic Cauchy Problems in Infinite Dimensions: Generalized and Regularized  
Solutions*, Irina V. Melnikova and Alexei Filinkov

*Signal Processing: A Mathematical Approach*, Charles L. Byrne

*Monomial Algebra, Second Edition*, Rafael Villarreal

*Groups, Designs, and Linear Algebra*, Donald L. Kreher

*Geometric Modeling and Mesh Generation from Scanned Images*, Yongjie Zhang

*Difference Equations: Theory, Applications and Advanced Topics, Third Edition*,  
Ronald E. Mickens

*Set Theoretical Aspects of Real Analysis*, Alexander Kharazishvili

*Method of Moments in Electromagnetics, Second Edition*, Walton C. Gibson

*The Separable Galois Theory of Commutative Rings, Second Edition*, Andy R. Magid

*Dictionary of Inequalities, Second Edition*, Peter Bullen

*Actions and Invariants of Algebraic Groups, Second Edition*, Walter Ferrer Santos  
and Alvaro Rittatore

*Practical Guide to Geometric Regulation for Distributed Parameter Systems*,  
Eugenio Aulisa and David S. Gilliam

*Analytical Methods for Kolmogorov Equations, Second Edition*, Luca Lorenzi

*Handbook of the Tutte Polynomial*, Joanna Anthony Ellis-Monaghan and Iain Moffat

*Blow-up Patterns for Higher-Order: Nonlinear Parabolic, Hyperbolic Dispersion and  
Schrödinger Equations*, Victor A. Galaktionov, Enzo L. Mitidieri and Stanislav Pohozaev

*Application of Fuzzy Logic to Social Choice Theory*, John N. Mordeson, Davendar Malik  
and Terry D. Clark

*Microlocal Analysis on  $R^n$  and on NonCompact Manifolds*, Sandro Coriasco

*Cremona Groups and Icosahedron*, Ivan Cheltsov and Constantin Shramov

MONOGRAPHS AND RESEARCH NOTES IN MATHEMATICS

# Modeling and Inverse Problems in the Presence of Uncertainty

H. T. Banks, Shuhua Hu,  
and W. Clayton Thompson

North Carolina State University

Raleigh, USA



CRC Press

Taylor & Francis Group

Boca Raton London New York

---

CRC Press is an imprint of the  
Taylor & Francis Group, an **informa** business  
A CHAPMAN & HALL BOOK

CRC Press  
Taylor & Francis Group  
6000 Broken Sound Parkway NW, Suite 300  
Boca Raton, FL 33487-2742

© 2014 by Taylor & Francis Group, LLC  
CRC Press is an imprint of Taylor & Francis Group, an Informa business

No claim to original U.S. Government works  
Version Date: 20140224

International Standard Book Number-13: 978-1-4822-0643-2 (eBook - PDF)

This book contains information obtained from authentic and highly regarded sources. Reasonable efforts have been made to publish reliable data and information, but the author and publisher cannot assume responsibility for the validity of all materials or the consequences of their use. The authors and publishers have attempted to trace the copyright holders of all material reproduced in this publication and apologize to copyright holders if permission to publish in this form has not been obtained. If any copyright material has not been acknowledged please write and let us know so we may rectify in any future reprint.

Except as permitted under U.S. Copyright Law, no part of this book may be reprinted, reproduced, transmitted, or utilized in any form by any electronic, mechanical, or other means, now known or hereafter invented, including photocopying, microfilming, and recording, or in any information storage or retrieval system, without written permission from the publishers.

For permission to photocopy or use material electronically from this work, please access [www.copyright.com](http://www.copyright.com) (<http://www.copyright.com/>) or contact the Copyright Clearance Center, Inc. (CCC), 222 Rosewood Drive, Danvers, MA 01923, 978-750-8400. CCC is a not-for-profit organization that provides licenses and registration for a variety of users. For organizations that have been granted a photocopy license by the CCC, a separate system of payment has been arranged.

**Trademark Notice:** Product or corporate names may be trademarks or registered trademarks, and are used only for identification and explanation without intent to infringe.

**Visit the Taylor & Francis Web site at**  
**<http://www.taylorandfrancis.com>**

**and the CRC Press Web site at**  
**<http://www.crcpress.com>**

---

# Contents

<b>Preface</b>	<b>xiii</b>
<b>1 Introduction</b>	<b>1</b>
<b>2 Probability and Statistics Overview</b>	<b>3</b>
2.1 Probability and Probability Space . . . . .	3
2.1.1 Joint Probability . . . . .	6
2.1.2 Conditional Probability . . . . .	7
2.2 Random Variables and Their Associated Distribution Functions . . . . .	8
2.2.1 Cumulative Distribution Function . . . . .	9
2.2.2 Probability Mass Function . . . . .	12
2.2.3 Probability Density Function . . . . .	13
2.2.4 Equivalence of Two Random Variables . . . . .	14
2.2.5 Joint Distribution Function and Marginal Distribution Function . . . . .	15
2.2.6 Conditional Distribution Function . . . . .	17
2.2.7 Function of a Random Variable . . . . .	20
2.3 Statistical Averages of Random Variables . . . . .	21
2.3.1 Joint Moments . . . . .	23
2.3.2 Conditional Moments . . . . .	25
2.3.3 Statistical Averages of Random Vectors . . . . .	26
2.3.4 Important Inequalities . . . . .	26
2.4 Characteristic Functions of a Random Variable . . . . .	27
2.5 Special Probability Distributions . . . . .	28
2.5.1 Poisson Distribution . . . . .	29
2.5.2 Uniform Distribution . . . . .	29
2.5.3 Normal Distribution . . . . .	31
2.5.4 Log-Normal Distribution . . . . .	33
2.5.5 Multivariate Normal Distribution . . . . .	35
2.5.6 Exponential Distribution . . . . .	36
2.5.7 Gamma Distribution . . . . .	39
2.5.8 Chi-Square Distribution . . . . .	41
2.5.9 Student's $t$ Distribution . . . . .	42
2.6 Convergence of a Sequence of Random Variables . . . . .	44
2.6.1 Convergence in Distribution . . . . .	44
2.6.2 Convergence in Probability . . . . .	46

2.6.3	Convergence Almost Surely . . . . .	47
2.6.4	Mean Square Convergence . . . . .	49
	References . . . . .	50
<b>3</b>	<b>Mathematical and Statistical Aspects of Inverse Problems</b>	<b>53</b>
3.1	Least Squares Inverse Problem Formulations . . . . .	54
3.1.1	The Mathematical Model . . . . .	54
3.1.2	The Statistical Model . . . . .	54
3.2	Methodology: Ordinary, Weighted and Generalized	
	Least Squares . . . . .	56
3.2.1	Scalar Ordinary Least Squares . . . . .	56
3.2.2	Vector Ordinary Least Squares . . . . .	58
3.2.3	Numerical Implementation of the Vector OLS Procedure . . . . .	60
3.2.4	Weighted Least Squares (WLS) . . . . .	60
3.2.5	Generalized Least Squares Definition and Motivation . . . . .	62
3.2.6	Numerical Implementation of the GLS Procedure . . . . .	64
3.3	Asymptotic Theory: Theoretical Foundations . . . . .	64
3.3.1	Extension to Weighted Least Squares . . . . .	70
3.4	Computation of $\hat{\Sigma}^N$ , Standard Errors and Confidence Intervals	73
3.5	Investigation of Statistical Assumptions . . . . .	80
3.5.1	Residual Plots . . . . .	80
3.5.2	An Example Using Residual Plots: Logistic Growth . . . . .	82
3.5.3	A Second Example Using Residual Plot Analysis: Cell Proliferation . . . . .	87
3.6	Bootstrapping vs. Asymptotic Error Analysis . . . . .	93
3.6.1	Bootstrapping Algorithm: Constant Variance Data . . . . .	94
3.6.2	Bootstrapping Algorithm: Non-Constant Variance Data . . . . .	96
3.6.3	Results of Numerical Simulations . . . . .	97
3.6.3.1	Constant Variance Data with OLS . . . . .	97
3.6.3.2	Non-Constant Variance Data with GLS . . . . .	100
3.6.4	Using Incorrect Assumptions on Errors . . . . .	101
3.6.4.1	Constant Variance Data Using GLS . . . . .	103
3.6.4.2	Non-Constant Variance Data Using OLS . . . . .	103
3.7	The “Corrective” Nature of Bootstrapping Covariance Estimates and Their Effects on Confidence Intervals . . . . .	106
3.8	Some Summary Remarks on Asymptotic Theory vs. Bootstrapping . . . . .	111
	References . . . . .	112
<b>4</b>	<b>Model Selection Criteria</b>	<b>117</b>
4.1	Introduction . . . . .	117
4.1.1	Statistical and Probability Distribution Models . . . . .	117
4.1.2	Risks Involved in the Process of Model Selection . . . . .	119

4.1.3	Model Selection Principle . . . . .	120
4.2	Likelihood Based Model Selection Criteria – Akaike Information Criterion and Its Variations . . . . .	121
4.2.1	Kullback–Leibler Information . . . . .	122
4.2.2	Maximum Likelihood Estimation . . . . .	123
4.2.3	A Large Sample AIC . . . . .	125
4.2.4	A Small Sample AIC . . . . .	126
4.2.4.1	Univariate Observations . . . . .	126
4.2.4.2	Multivariate Observations . . . . .	127
4.2.5	Takeuchi’s Information Criterion . . . . .	128
4.2.6	Remarks on Akaike Information Criterion and Its Variations . . . . .	129
4.2.6.1	Candidate Models . . . . .	130
4.2.6.2	The Selected Best Model . . . . .	130
4.2.6.3	Pitfalls When Using the AIC . . . . .	131
4.3	The AIC under the Framework of Least Squares Estimation . . . . .	132
4.3.1	Independent and Identically Normally Distributed Observations . . . . .	132
4.3.2	Independent Multivariate Normally Distributed Observations . . . . .	134
4.3.2.1	Unequal Number of Observations for Different Observed Components . . . . .	136
4.3.3	Independent Gamma Distributed Observations . . . . .	138
4.3.4	General Remarks . . . . .	140
4.4	Example: CFSE Label Decay . . . . .	142
4.5	Residual Sum of Squares Based Model Selection Criterion . . . . .	146
4.5.1	Ordinary Least Squares . . . . .	146
4.5.2	Application: Cat Brain Diffusion/Convection Problem . . . . .	149
4.5.3	Weighted Least Squares . . . . .	151
4.5.4	Summary Remarks . . . . .	152
	References . . . . .	153

<b>5</b>	<b>Estimation of Probability Measures Using Aggregate Population Data</b>	<b>157</b>
5.1	Motivation . . . . .	157
5.2	Type I: Individual Dynamics/Aggregate Data Inverse Problems . . . . .	160
5.2.1	Structured Population Models . . . . .	160
5.3	Type II: Aggregate Dynamics/Aggregate Data Inverse Problems . . . . .	163
5.3.1	Probability Measure Dependent Systems—Viscoelasticity . . . . .	163
5.3.2	Probability Measure Dependent Systems—Maxwell’s Equations . . . . .	167
5.4	Aggregate Data and the Prohorov Metric Framework . . . . .	169
5.5	Consistency of the PMF Estimator . . . . .	178



5.6	Further Remarks . . . . .	181
5.7	Non-Parametric Maximum Likelihood Estimation . . . . .	181
5.7.1	Likelihood Formulation . . . . .	182
5.7.2	Computational Techniques . . . . .	184
5.8	Final Remarks . . . . .	187
	References . . . . .	188
<b>6</b>	<b>Optimal Design</b>	<b>195</b>
6.1	Introduction . . . . .	195
6.2	Mathematical and Statistical Models . . . . .	197
6.2.1	Formulation of the Optimal Design Problem . . . . .	198
6.3	Algorithmic Considerations . . . . .	202
6.4	Example: HIV model . . . . .	203
	References . . . . .	206
<b>7</b>	<b>Propagation of Uncertainty in a Continuous Time Dynamical System</b>	<b>209</b>
7.1	Introduction to Stochastic Processes . . . . .	210
7.1.1	Distribution Functions of a Stochastic Process . . . . .	211
7.1.2	Moments, Correlation and Covariance Functions of a Stochastic Process . . . . .	213
7.1.3	Classification of a Stochastic Process . . . . .	215
7.1.3.1	Stationary vs. Non-Stationary Stochastic Processes . . . . .	215
7.1.3.2	Gaussian vs. Non-Gaussian Processes . . . . .	216
7.1.4	Methods of Studying a Stochastic Process . . . . .	217
7.1.4.1	Sample Function Approach . . . . .	217
7.1.4.2	Mean Square Calculus Approach . . . . .	218
7.1.5	Markov Processes . . . . .	220
7.1.5.1	Characterization of a Markov Process . . . . .	222
7.1.5.2	The Chapman–Kolmogorov Equation . . . . .	222
7.1.5.3	An Example of a Markov Process: Wiener Process . . . . .	223
7.1.5.4	An Example of a Markov Process: Diffusion Process . . . . .	226
7.1.5.5	An Example of a Markov Process: Poisson Process . . . . .	226
7.1.5.6	Classification of a Markov Process . . . . .	229
7.1.5.7	Continuous Time Markov Chain . . . . .	230
7.1.6	Martingales . . . . .	233
7.1.6.1	Examples of Sample-Continuous Martingales . . . . .	234
7.1.6.2	The Role of Martingales in the Development of Stochastic Integration Theory . . . . .	234
7.1.7	White Noise vs. Colored Noise . . . . .	236
7.1.7.1	The Power Spectral Density Function . . . . .	236

7.1.7.2	White Noise . . . . .	238
7.1.7.3	Colored Noise . . . . .	240
7.2	Stochastic Differential Equations . . . . .	244
7.2.1	Itô Stochastic Differential Equations . . . . .	245
7.2.1.1	Evolution of the Probability Density Function of $\mathbf{X}(t)$ . . . . .	249
7.2.1.2	Applications of the Fokker–Plank Equation in Population Dynamics . . . . .	252
7.2.2	Stratonovich Stochastic Differential Equations . . . . .	255
7.3	Random Differential Equations . . . . .	257
7.3.1	Differential Equations with Random Initial Conditions . . . . .	258
7.3.1.1	Evolution of the Probability Density Function of $\mathbf{x}(t; \mathbf{X}_0)$ . . . . .	259
7.3.1.2	Applications of Liouville’s Equation in Popu- lation Dynamics . . . . .	261
7.3.2	Differential Equations with Random Model Parameters and Random Initial Conditions . . . . .	262
7.3.2.1	Evolution of the Joint Probability Density Func- tion for $(\mathbf{x}(t; \mathbf{X}_0, \mathbf{Z}), \mathbf{Z})^T$ . . . . .	263
7.3.2.2	Evolution of Conditional Probability Density Function of $\mathbf{x}(t; \mathbf{X}_0, \mathbf{Z})$ Given the Realization $\mathbf{z}$ of $\mathbf{Z}$ . . . . .	264
7.3.2.3	Applications in Population Dynamics . . . . .	265
7.3.3	Differential Equations Driven by Correlated Stochastic Processes . . . . .	268
7.3.3.1	Joint Probability Density Function of the Cou- pled Stochastic Process . . . . .	269
7.3.3.2	The Probability Density Function of $\mathbf{X}(t)$ . . . . .	273
7.4	Relationships between Random and Stochastic Differential Equa- tions . . . . .	276
7.4.1	Markov Operators and Markov Semigroups . . . . .	277
7.4.1.1	Random Differential Equations . . . . .	279
7.4.1.2	Stochastic Differential Equations . . . . .	281
7.4.2	Pointwise Equivalence Results between Stochastic Dif- ferential Equations and Random Differential Equations . . . . .	282
7.4.2.1	Scalar Affine Differential Equations (Class 1) . . . . .	283
7.4.2.2	Scalar Affine Differential Equations (Class 2) . . . . .	285
7.4.2.3	Vector Affine Systems . . . . .	286
7.4.2.4	Non-Linear Differential Equations . . . . .	288
7.4.2.5	Remarks on the Equivalence between the SDE and the RDE . . . . .	293
7.4.2.6	Relationship between the FPPS and GRDPS Population Models . . . . .	295
	References . . . . .	298

<b>8</b>	<b>A Stochastic System and Its Corresponding Deterministic System</b>	<b>309</b>
8.1	Overview of Multivariate Continuous Time Markov Chains . . . . .	310
8.1.1	Exponentially Distributed Holding Times . . . . .	310
8.1.2	Random Time Change Representation . . . . .	311
8.1.2.1	Relationship between the Stochastic Equation and the Martingale Problem . . . . .	312
8.1.2.2	Relationship between the Martingale Problem and Kolmogorov's Forward Equation . . . . .	313
8.2	Simulation Algorithms for Continuous Time Markov Chain Models . . . . .	314
8.2.1	Stochastic Simulation Algorithm . . . . .	314
8.2.1.1	The Direct Method . . . . .	315
8.2.1.2	The First Reaction Method . . . . .	315
8.2.2	The Next Reaction Method . . . . .	316
8.2.2.1	The Original Next Reaction Method . . . . .	317
8.2.2.2	The Modified Next Reaction Method . . . . .	319
8.2.3	Tau-Leaping Methods . . . . .	321
8.2.3.1	An Explicit Tau-Leaping Method . . . . .	321
8.2.3.2	An Implicit Tau-Leaping Method . . . . .	325
8.3	Density Dependent Continuous Time Markov Chains and Kurtz's Limit Theorem . . . . .	327
8.3.1	Kurtz's Limit Theorem . . . . .	328
8.3.2	Implications of Kurtz's Limit Theorem . . . . .	329
8.4	Biological Application: Vancomycin-Resistant Enterococcus Infection in a Hospital Unit . . . . .	331
8.4.1	The Stochastic VRE Model . . . . .	331
8.4.2	The Deterministic VRE Model . . . . .	333
8.4.3	Numerical Results . . . . .	334
8.5	Biological Application: HIV Infection within a Host . . . . .	336
8.5.1	Deterministic HIV Model . . . . .	336
8.5.2	Stochastic HIV Models . . . . .	341
8.5.2.1	The Stochastic HIV Model Based on the Burst Production Mode . . . . .	341
8.5.2.2	The Stochastic HIV Model Based on the Continuous Production Mode . . . . .	343
8.5.3	Numerical Results for the Stochastic HIV Model Based on the Burst Production Mode . . . . .	343
8.5.3.1	Implementation of the Tau-Leaping Algorithms . . . . .	343
8.5.3.2	Comparison of Computational Efficiency of the SSA and the Tau-Leaping Algorithms . . . . .	346
8.5.3.3	Accuracy of the Results Obtained by Tau-Leaping Algorithms . . . . .	348
8.5.3.4	Stochastic Solution vs. Deterministic Solution . . . . .	350

8.5.3.5	Final Remark . . . . .	351
8.6	Application in Agricultural Production Networks . . . . .	352
8.6.1	The Stochastic Pork Production Network Model . . . . .	352
8.6.2	The Deterministic Pork Production Network Model . . . . .	353
8.6.3	Numerical Results . . . . .	354
8.7	Overview of Stochastic Systems with Delays . . . . .	356
8.8	Simulation Algorithms for Stochastic Systems with Fixed Delays . . . . .	358
8.9	Application in the Pork Production Network with a Fixed Delay . . . . .	359
8.9.1	The Stochastic Pork Production Network Model with a Fixed Delay . . . . .	360
8.9.2	The Deterministic Pork Production Network Model with a Fixed Delay . . . . .	360
8.9.3	Comparison of the Stochastic Model with a Fixed Delay and Its Corresponding Deterministic System . . . . .	362
8.10	Simulation Algorithms for Stochastic Systems with Random Delays . . . . .	362
8.11	Application in the Pork Production Network with a Random Delay . . . . .	365
8.11.1	The Corresponding Deterministic System . . . . .	365
8.11.2	Comparison of the Stochastic Model with a Random Delay and Its Corresponding Deterministic System . . . . .	367
8.11.3	The Corresponding Constructed Stochastic System . . . . .	368
8.11.4	Comparison of the Constructed Stochastic System and Its Corresponding Deterministic System . . . . .	370
8.11.5	Comparison of the Stochastic System with a Random Delay and the Constructed Stochastic System . . . . .	370
8.11.5.1	The Effect of Sample Size on the Comparison of These Two Stochastic Systems . . . . .	371
8.11.5.2	The Effect of the Variance of a Random Delay on the Comparison of These Two Stochastic Systems . . . . .	375
8.11.5.3	Summary Remarks . . . . .	376
References	. . . . .	378

<b>Frequently Used Notations and Abbreviations</b>	<b>383</b>
--	------------

<b>Index</b>	<b>387</b>
--------------	------------

This page intentionally left blank

---

# Preface

Writing a research monograph on a “hot topic” such as “uncertainty propagation” is a somewhat daunting undertaking. Nonetheless, we decided to collect our own views, supported by our own research efforts over the past 12–15 years on a number of aspects of this topic, and summarize these for the possible enlightenment they might provide (for us, our students and others). The research results discussed below are thus necessarily filled with a preponderance of references to our own research reports and papers. In numerous references below (given at the conclusion of each chapter), we refer to CRSC-TRXX-YY. This refers to early Technical Report versions of manuscripts which can be found on the Center for Research in Scientific Computation website at North Carolina State University where XX refers to the year, e.g., XX = 03 is 2003, XX = 99 is 1999, while the YY refers to the number of the report in that year. These can be found at and downloaded from <http://www.ncsu.edu/crsc/reports.html> where they are listed by year.

Our presentation here has an intended audience from the community of investigators in applied mathematics interested in deterministic and/or stochastic models and their interactions as well as scientists in biology, medicine, engineering and physics interested in basic modeling and inverse problems, uncertainty in modeling, propagation of uncertainty and statistical modeling.

We owe great thanks to our former and current students, postdocs and colleagues for their patience in enduring lectures, questions, feedback and some proofreading. Special thanks are due (in no particular order) to Zack Kenz, Keri Rehm, Dustin Kapraun, Jared Catenacci, Katie Link, Kris Rinnovatore, Kevin Flores, John Nardini, Karissa Cross and Laura Poag for careful reading of notes and suggested corrections/revisions on subsets of the material for this monograph. However, in a sincere attempt to give credit where it is due, each of the authors firmly insists that any errors in judgment, mathematical content, grammar or typos in the material presented in this monograph are entirely the responsibility of his/her two co-authors!!

We (especially young members of our research group) have been generously supported by research grants and fellowships from US federal funding agencies including AFSOR, DARPA, NIH, NSF, DED, and DOE. For this support and encouragement we are all most grateful.

H.T. Banks  
Shuhua Hu  
W. Clayton Thompson

This page intentionally left blank

# Chapter 1

---

## Introduction

The terms *uncertainty quantification* and *uncertainty propagation* have become so widely used as to almost have little meaning unless they are further explained. Here we focus primarily on two basic types of problems:

1. Modeling and inverse problems where one assumes that a precise mathematical model without modeling error is available. This is a standard assumption underlying a large segment of what is taught in many modern statistics courses with a frequentist philosophy. More precisely, a mathematical model is given by a dynamical system

$$\frac{d\mathbf{x}}{dt}(t) = \mathbf{g}(t, \mathbf{x}(t), \mathbf{q}) \quad (1.1)$$

$$\mathbf{x}(t_0) = \mathbf{x}_0 \quad (1.2)$$

with **observation process**

$$\mathbf{f}(t; \boldsymbol{\theta}) = \mathcal{C}\mathbf{x}(t; \boldsymbol{\theta}), \quad (1.3)$$

where  $\boldsymbol{\theta} = (\mathbf{q}, \mathbf{x}_0)$ . The mathematical model is an  $n$ -dimensional deterministic system and there is a corresponding “truth” parameter  $\boldsymbol{\theta}_0 = (\mathbf{q}_0, \mathbf{x}_{00})$  so that in the presence of no measurement error the data can be described exactly by the deterministic system at  $\boldsymbol{\theta}_0$ . Thus, uncertainty is present entirely due to some **statistical model** of the form

$$\mathbf{Y}_j = \mathbf{f}(t_j; \boldsymbol{\theta}_0) + \boldsymbol{\mathcal{E}}_j, \quad j = 1, \dots, N, \quad (1.4)$$

where  $\mathbf{f}(t_j; \boldsymbol{\theta}) = \mathcal{C}\mathbf{x}(t_j; \boldsymbol{\theta})$ ,  $j = 1, \dots, N$ , corresponds to the observed part of the solution of the mathematical model (1.1)–(1.2) at the  $j$ th covariate or observation time and  $\boldsymbol{\mathcal{E}}_j$  is some type of (possibly state dependent) measurement error. For example, we consider errors that include those of the form  $\boldsymbol{\mathcal{E}}_j = \mathbf{f}(t_j; \boldsymbol{\theta}_0)^\gamma \circ \bar{\boldsymbol{\mathcal{E}}}_j$  where the operation  $^\gamma \circ$  denotes component-wise exponentiation by  $\gamma$  followed by component-wise multiplication and  $\gamma \geq 0$ .

2. An alternate problem wherein the mathematical modeling *itself* is a major source of uncertainty and this uncertainty usually propagates in time. That is, the mathematical model has major uncertainties in its form and/or its parametrization and/or its initial/boundary data, and this uncertainty is propagated dynamically via some framework as yet to be determined.



Before we begin the inverse problem discussions, we give a brief but useful review of certain basic probability and statistical concepts. After the probability and statistics review we present a chapter summarizing both mathematical and statistical aspects of inverse problem methodology which includes ordinary, weighted and generalized least-squares formulations. We discuss asymptotic theories, bootstrapping and issues related to evaluation of the correctness of the assumed form of statistical models. We follow this with a discussion of methods for evaluating and comparing the validity of appropriateness of a collection of models for describing a given data set, including statistically based model selection and model comparison techniques.

In Chapter 5 we present a summary of recent results on the estimation of probability distributions when they are embedded in complex mathematical models and only aggregate (not individual) data are available. This is followed by a brief chapter on optimal design (what to measure? when and where to measure?) of experiments to be carried out in support of inverse problems for given models.

The last two chapters focus on the uncertainty in model formulation itself (the second item listed above as the focus of this monograph). In Chapter 7 we consider the general problem of evolution of probability density functions in time. This is done in the context of associated processes resulting from stochastic differential equations (SDE), which are driven by white noise, and those resulting from random differential equations (RDE), which are driven by colored noise. We also discuss their respective wide applications in a number of different fields including physics and biology. We also consider the general relationship between SDE and RDE and establish that there are classes of problems for which there is an equivalence between the solutions of the two formulations. This equivalence, which we term pointwise equivalence, is in the sense that the respective probability density functions are the same at each time  $t$ . We show, however, that the stochastic processes resulting from the SDE and its corresponding pointwise equivalent RDE are generally not the same in that they may have different covariance functions.

In a final chapter we consider questions related to the appropriateness of discrete versus continuum models in transitions from small numbers of individuals (particles, populations, molecules, etc.) to large numbers. These investigations are carried out in the context of continuous time Markov chain (CTMC) models and the Kurtz limit theorems for approximations for large number stochastic populations by ordinary differential equations for corresponding mean populations. Algorithms for simulating CTMC models and CTMC models with delays (discrete and random) are explained and simulations are presented for problems arising in specific applications.

The monograph contains illustrative examples throughout, many of them directly related to research projects carried out by our group at North Carolina State University over the past decade.

# Chapter 2

---

## *Probability and Statistics Overview*

The theory of probability and statistics is an essential mathematical tool in the formulation of inverse problems, in the development of subsequent analysis and approaches to statistical hypothesis testing and model selection criteria, and in the study of uncertainty propagation in dynamic systems. Our coverage of these fundamental and important topics is brief and limited in scope. Indeed, we provide in this section a few definitions and basic concepts in the theory of probability and statistics that are essential for the understanding of estimators, confidence intervals, model selection criteria and stochastic processes. For more information on the topics in this chapter, selected references are provided at the end of these as well as subsequent chapters.

---

### 2.1 Probability and Probability Space

The set of all possible outcomes in a statistical experiment is called the *sample space* and is denoted by  $\Omega$ . Each element  $\omega \in \Omega$  is called a *sample point*. A collection of outcomes in which we are interested is called an *event*; that is, an event is a subset of  $\Omega$ . For example, consider the experiment of rolling a six-sided die. In this case, there are six possible outcomes, and the sample space can be represented as

$$\Omega = \{1, 2, 3, 4, 5, 6\}. \quad (2.1)$$

An event  $\mathbb{A}$  might be defined as

$$\mathbb{A} = \{1, 5\}, \quad (2.2)$$

which consists of the outcomes 1 and 5. Note that we say the event  $\mathbb{A}$  occurs if the outcome of the experiment is in the set  $\mathbb{A}$ . Consider another experiment of tossing a coin three times. A sample point in this experiment indicates the result of each toss; for example, HHT indicates that two heads and then a tail were observed. The sample space for this experiment has eight sample points; that is,

$$\Omega = \{\text{HHH}, \text{HHT}, \text{HTH}, \text{THH}, \text{TTH}, \text{THT}, \text{HTT}, \text{TTT}\}. \quad (2.3)$$

An event  $\mathbb{A}$  might be defined as the set of outcomes for which the first toss is a head; that is,

$$\mathbb{A} = \{\text{HHH}, \text{HHT}, \text{HTH}, \text{HTT}\}.$$

Thus, we see that the sample point could be either a numerical value or a character value. Based on the number of sample points contained in the sample space, the sample space can be either finite (as we illustrate above) or infinitely countable (e.g., the number of customers to arrive in a bank) or uncountable (e.g., the sample space for the lifetime of a bulb or the sample space for the reaction time to a certain stimulus).

**Definition 2.1.1** *Let  $\Omega$  be the given sample space. Then the  $\sigma$ -algebra  $\mathcal{F}$  is a collection of subsets of  $\Omega$  with the following properties:*

- (i) *The empty set  $\emptyset$  is an element of  $\mathcal{F}$ ; that is,  $\emptyset \in \mathcal{F}$ .*
- (ii) *(closed under complementation): If  $\mathbb{A} \in \mathcal{F}$ , then  $\mathbb{A}^c \in \mathcal{F}$ , where  $\mathbb{A}^c$  denotes the complement of the event  $\mathbb{A}$ , which consists of all sample points in  $\Omega$  that are not in  $\mathbb{A}$ .*
- (iii) *(closed under countable unions): If  $\mathbb{A}_1, \mathbb{A}_2, \mathbb{A}_3, \dots \in \mathcal{F}$ , then  $\bigcup_j \mathbb{A}_j \in \mathcal{F}$ .*

*The pair  $(\Omega, \mathcal{F})$  is called a measurable space. A subset  $\mathbb{A}$  of  $\Omega$  is said to be a measurable set (or event) if  $\mathbb{A} \in \mathcal{F}$ .*

It is worth noting that a  $\sigma$ -algebra is also called a  $\sigma$ -field in the literature. In addition, one can define many different  $\sigma$ -algebras associated with the sample space  $\Omega$ . The  $\sigma$ -algebra we will mainly consider is the smallest one that contains all of the open sets in the given sample space. In other words, it is the algebra generated by a topological space, whose definition is given as follows.

**Definition 2.1.2** *A topology  $\mathcal{T}$  on a set  $\Omega$  is a collection of subsets of  $\Omega$  having the following properties:*

- $\emptyset \in \mathcal{T}$  and  $\Omega \in \mathcal{T}$ .
- *(closed under finite intersection): If  $\mathbb{U}_i \in \mathcal{T}$ ,  $i = 1, 2, \dots, l$ , with  $l$  being a positive integer, then  $\bigcap_{i=1}^l \mathbb{U}_i \in \mathcal{T}$ .*
- *(closed under arbitrary union): If  $\{\mathbb{U}_\alpha\}$  is an arbitrary collection of members of  $\mathcal{T}$  (finite, countable, uncountable), then  $\bigcup_\alpha \mathbb{U}_\alpha \in \mathcal{T}$ .*

*The pair  $(\Omega, \mathcal{T})$  is called a topological space. A subset  $\mathbb{U}$  of  $\Omega$  is said to be an open set if  $\mathbb{U} \in \mathcal{T}$ . A subset  $\mathbb{A}$  of  $\Omega$  is said to be closed if  $\mathbb{A}^c$  is an open set.*

A  $\sigma$ -algebra generated by all the open sets in a given sample space  $\Omega$  is often called a *Borel algebra* or *Borel  $\sigma$ -algebra* (denoted by  $\mathcal{B}(\Omega)$ ), and the sets in a Borel algebra are called *Borel sets*.

Associated with an event  $A \in \mathcal{F}$  is its probability  $\text{Prob}\{A\}$ , which indicates the likelihood that event  $A$  occurs. For example, in a fair experiment of rolling a die, where one assumes that each possible sample point has probability  $\frac{1}{6}$ , then the event  $A$  as defined by (2.2) has probability  $\text{Prob}\{A\} = \frac{2}{6} = \frac{1}{3}$ . The strict definition for the probability is given as follows.

**Definition 2.1.3** *A probability Prob on the measurable space  $(\Omega, \mathcal{F})$  is a set function  $\text{Prob} : \mathcal{F} \rightarrow [0, 1]$  such that*

- (i)  $\text{Prob}\{\emptyset\} = 0$ ,  $\text{Prob}\{\Omega\} = 1$ .
- (ii) (completely additive): *If  $A_1, A_2, A_3, \dots$ , is a finite or an infinite sequence of disjoint subsets in  $\mathcal{F}$ , then  $\text{Prob}\{\cup_j A_j\} = \sum_j \text{Prob}\{A_j\}$ .*

*The triplet  $(\Omega, \mathcal{F}, \text{Prob})$  is called a probability space.*

It is worth noting that a probability is also called a *probability measure* or a *probability function* (these names will be used interchangeably in this monograph), and a probability measure defined on a Borel algebra is called a *Borel probability measure*.

Using Definition 2.1.3 for probability, a number of immediate consequences can also be derived which have important applications. For example, the probability that an event will occur and that it will not occur always sum to 1. That is,

$$\text{Prob}\{A\} + \text{Prob}\{A^c\} = 1.$$

In addition, if  $A, B \in \mathcal{F}$ , then we have

$$\text{Prob}\{A\} \leq \text{Prob}\{B\} \text{ if } A \subset B.$$

It can also be found that if  $A_j \in \mathcal{F}$ ,  $j = 1, 2, \dots$ , then

$$\text{Prob}\{\cup_{j=1}^{\infty} A_j\} \leq \sum_{j=1}^{\infty} \text{Prob}\{A_j\}.$$

If  $\text{Prob}\{A\} = 1$ , then we say that the event  $A$  occurs “with probability 1” or “*almost surely* (a.s.).” A set  $A \in \mathcal{F}$  is called a *null set* if  $\text{Prob}\{A\} = 0$ . In addition, a probability space  $(\Omega, \mathcal{F}, \text{Prob})$  is said to be *complete* if for any two sets  $A$  and  $B$  the following condition holds: If  $A \subset B$ ,  $B \in \mathcal{F}$  and  $\text{Prob}\{B\} = 0$ , then  $A \in \mathcal{F}$ . It is worth noting that any probability space can be extended into a complete probability space (e.g., see [15, p. 10] and the references therein). Hence, we will assume that all the probability spaces are complete in the remainder of this monograph.

**Remark 2.1.1** We remark that we can generalize a probability space to more general measure spaces. Specifically, a measure  $\nu$  on the measurable space  $(\Omega, \mathcal{F})$  is a set function  $\nu : \mathcal{F} \rightarrow [0, \infty]$  such that  $\nu(\emptyset) = 0$  and  $\nu$  is completely additive (that is, the second property of probability in Definition 2.1.3 holds). The triplet  $(\Omega, \mathcal{F}, \nu)$  is called a measure space. If  $\nu(\Omega)$  is finite, then  $\nu$  is said to be a finite measure. In particular,  $\text{Prob}$  is a normalized finite measure with  $\text{Prob}(\Omega) = 1$ . Hence, a probability possesses all the general properties of a finite measure.

**Remark 2.1.2** Let  $(\Omega, \mathcal{F}, \nu)$  be a measure space. If  $\Omega = \bigcup_{j=1}^{\infty} \mathbb{A}_j$  and  $\nu(\mathbb{A}_j)$  is finite for all  $j$ , then we say that  $\nu$  is a  $\sigma$ -finite measure and  $(\Omega, \mathcal{F}, \nu)$  is a  $\sigma$ -finite measure space. Another measure that we will consider in this monograph is the Lebesgue measure, which is defined on the measurable space  $(\mathbb{R}^l, \mathcal{B}(\mathbb{R}^l))$  and is given by

$$\nu((a_1, b_1) \times \cdots \times (a_l, b_l)) = \prod_{j=1}^l (b_j - a_j),$$

that is, the volume of the interval  $(a_1, b_1) \times \cdots \times (a_l, b_l)$ . We thus see that the Lebesgue measure of a countable set of points is zero (that is, a countable set of points is a null set with respect to the Lebesgue measure), and that the Lebesgue measure of a  $k$ -dimensional plane in  $\mathbb{R}^l$  ( $l > k$ ) is also zero. We refer the interested reader to some real analysis textbooks such as [11, 23] for more information on a measure as well as its properties.

### 2.1.1 Joint Probability

Instead of considering a single experiment, let us perform two experiments and consider their outcomes. For example, the two experiments may be two separate tosses of a single die or a single toss of two dice. The sample space in this case consists of 36 pairs  $(k, j)$ , where  $k, j = 1, 2, \dots, 6$ . Note that in a fair dice game, each sample point in the sample space has probability  $\frac{1}{36}$ . We now consider the probability of joint events, such as  $\{k = 2, j = \text{odd}\}$ . We begin by denoting the event of one experiment by  $\mathbb{A}_k$ ,  $k = 1, 2, \dots, l$ , and the event of the second experiment by  $\mathbb{B}_j$ ,  $j = 1, 2, \dots, m$ . The combined experiment has the joint events  $(\mathbb{A}_k, \mathbb{B}_j)$ , where  $k = 1, 2, \dots, l$  and  $j = 1, 2, \dots, m$ .

The joint probability  $\text{Prob}\{\mathbb{A}_k, \mathbb{B}_j\}$ , also denoted by  $\text{Prob}\{\mathbb{A}_k \cap \mathbb{B}_j\}$  (which will be occasionally used in this monograph for notational convenience) or  $\text{Prob}\{\mathbb{A}_k \mathbb{B}_j\}$  in the literature, indicates the likelihood that the events  $\mathbb{A}_k$  and  $\mathbb{B}_j$  occur simultaneously. By Definition 2.1.3, a number of immediate consequences can also be derived for the joint probability. For example,  $\text{Prob}\{\mathbb{A}_k, \mathbb{B}_j\}$  satisfies the condition  $0 \leq \text{Prob}\{\mathbb{A}_k, \mathbb{B}_j\} \leq 1$ . In addition,

if  $\mathbb{B}_j$  for  $j = 1, 2, \dots, m$  are mutually exclusive (i.e.,  $\mathbb{B}_i \cap \mathbb{B}_j = \emptyset, i \neq j$ ) such that  $\bigcup_{j=1}^m \mathbb{B}_j = \Omega$ , then

$$\sum_{j=1}^m \text{Prob}\{\mathbb{A}_k, \mathbb{B}_j\} = \text{Prob}\{\mathbb{A}_k\}. \quad (2.4)$$

Furthermore, if all the outcomes of the two experiments are mutually exclusive such that  $\bigcup_{k=1}^l \mathbb{A}_k = \Omega$  and  $\bigcup_{j=1}^m \mathbb{B}_j = \Omega$ , then  $\sum_{k=1}^l \sum_{j=1}^m \text{Prob}\{\mathbb{A}_k, \mathbb{B}_j\} = 1$ . The generalization of the above concept to more than two experiments follows in a straightforward manner.

### 2.1.2 Conditional Probability

Next, we consider a joint event with probability  $\text{Prob}\{\mathbb{A}, \mathbb{B}\}$ . Assuming that event  $\mathbb{A}$  has occurred and  $\text{Prob}\{\mathbb{A}\} > 0$ , we wish to determine the probability of the event  $\mathbb{B}$ . This is called the *conditional probability* of event  $\mathbb{B}$  given the occurrence of event  $\mathbb{A}$  and is given by

$$\text{Prob}\{\mathbb{B}|\mathbb{A}\} = \frac{\text{Prob}\{\mathbb{A}, \mathbb{B}\}}{\text{Prob}\{\mathbb{A}\}}. \quad (2.5)$$

**Definition 2.1.4** *Two events,  $\mathbb{A}$  and  $\mathbb{B}$ , are said to be statistically independent if and only if*

$$\text{Prob}\{\mathbb{A}, \mathbb{B}\} = \text{Prob}\{\mathbb{A}\}\text{Prob}\{\mathbb{B}\}. \quad (2.6)$$

Statistical independence is often simply called *independence*. By (2.5) and (2.6), we see that if  $\mathbb{A}$  and  $\mathbb{B}$  are independent, then

$$\text{Prob}\{\mathbb{B}|\mathbb{A}\} = \text{Prob}\{\mathbb{B}\}. \quad (2.7)$$

In addition, we observe that if (2.7) holds, then by (2.5) and Definition 2.1.4 we know that  $\mathbb{A}$  and  $\mathbb{B}$  are independent. Thus, (2.7) can also be used as a definition for the independence of two events.

Two very useful relationships for conditional probabilities can be given. If  $\mathbb{A}_k, k = 1, 2, \dots, l$ , are mutually exclusive events such that  $\bigcup_{k=1}^l \mathbb{A}_k = \Omega$  and  $\mathbb{B}$  is an arbitrary event with  $\text{Prob}\{\mathbb{B}\} > 0$ , then by (2.4) and (2.5) we have

$$\text{Prob}\{\mathbb{B}\} = \sum_{j=1}^l \text{Prob}\{\mathbb{A}_j, \mathbb{B}\} = \sum_{j=1}^l \text{Prob}\{\mathbb{B}|\mathbb{A}_j\}\text{Prob}\{\mathbb{A}_j\}, \quad (2.8)$$

and

$$\text{Prob}\{\mathbb{A}_k|\mathbb{B}\} = \frac{\text{Prob}\{\mathbb{A}_k, \mathbb{B}\}}{\text{Prob}\{\mathbb{B}\}} = \frac{\text{Prob}\{\mathbb{B}|\mathbb{A}_k\}\text{Prob}\{\mathbb{A}_k\}}{\sum_{j=1}^l \text{Prob}\{\mathbb{B}|\mathbb{A}_j\}\text{Prob}\{\mathbb{A}_j\}}. \quad (2.9)$$

Equation (2.8) is often called the *law of total probability*. Equation (2.9) is known as Bayes' formula or Bayes' Theorem. Here  $\text{Prob}\{\mathbb{A}_k\}$  is called a *prior* probability of event  $\mathbb{A}_k$ ,  $\text{Prob}\{\mathbb{B}|\mathbb{A}_k\}$  is called the likelihood of  $\mathbb{B}$  given  $\mathbb{A}_k$ , and  $\text{Prob}\{\mathbb{A}_k|\mathbb{B}\}$  is called a *posterior* probability of event  $\mathbb{A}_k$  obtained by using the information gained from  $\mathbb{B}$ .

## 2.2 Random Variables and Their Associated Distribution Functions

In most applications of probability theory, we are not interested in the details associated with each sample point but rather in some numerical description of the outcome of an experiment. For example, in the experiment of tossing a coin three times, we might only be interested in the number of heads obtained in these three tosses. In the language of probability and statistics, the number of heads obtained in these three tosses is called a *random variable*. The values of this particular random variable corresponding to each sample point in (2.3) are given by

$$\begin{array}{cccccccc} \text{HHH} & \text{HHT} & \text{HTH} & \text{THH} & \text{TTH} & \text{THT} & \text{HTT} & \text{TTT} \\ 3 & 2 & 2 & 2 & 1 & 1 & 1 & 0. \end{array}$$

This implies that the range (i.e., the collection of all possible values) of this random variable is  $\{0, 1, 2, 3\}$ . The strict definition of a random variable is as follows.

**Definition 2.2.1** A function  $X : \Omega \rightarrow \mathbb{R}$  is said to be a random variable defined on a measurable space  $(\Omega, \mathcal{F})$  if for any  $x \in \mathbb{R}$  we have  $\{\omega \in \Omega : X(\omega) \leq x\} \in \mathcal{F}$ . Such a function is said to be measurable with respect to  $\mathcal{F}$ . In addition, for any fixed  $\omega \in \Omega$ ,  $X(\omega)$  is called a realization of this random variable.

As is usually done and also for notational convenience, we suppress the dependence of the random variable on  $\omega$  if no confusion occurs; that is, we denote  $X(\cdot)$  by  $X$ . Under this convention, the set  $\{\omega \in \Omega : X(\omega) \leq x\}$  is simply written as  $\{X \leq x\}$ . In addition, a realization of a random variable is simply denoted by its corresponding lower case letter; for example, a realization of random variable  $X$  is denoted by  $x$ . We point out that when we consider

a sequence of random variables in this monograph, we assume that they are defined on the same probability space.

The  $\sigma$ -algebra generated by the random variable  $X$ , often denoted by  $\sigma(X)$ , is given by

$$\sigma(X) = \{\mathbb{A} \subset \Omega \mid \mathbb{A} = X^{-1}(\mathbb{B}), \mathbb{B} \in \mathcal{B}(\mathbb{R})\},$$

where  $\mathcal{B}(\mathbb{R})$  is the Borel algebra on  $\mathbb{R}$ . It is worth noting that this is the smallest  $\sigma$ -algebra with respect to which  $X$  is measurable. This means that if  $X$  is  $\mathcal{F}$ -measurable, then  $\sigma(X) \subset \mathcal{F}$ .

**Remark 2.2.1** *The concept of a random variable can be generalized so that its range can be some complicated space rather than  $\mathbb{R}$ . Let  $(\Omega, \mathcal{F}, \text{Prob})$  be a probability space, and  $(\mathbb{S}, \mathcal{S})$  be a measurable space. A function  $X : \Omega \rightarrow \mathbb{S}$  is said to be a random element if for any  $\mathbb{B} \in \mathcal{S}$  we have  $\{\omega \in \Omega : X(\omega) \in \mathbb{B}\} \in \mathcal{F}$ . Specifically, if  $\mathbb{S} = \mathbb{R}^m$ , then we say that random element  $X$  is an  $m$ -dimensional random vector.*

### 2.2.1 Cumulative Distribution Function

For any random variable, there is an associated function called a cumulative distribution function, which is defined as follows.

**Definition 2.2.2** *The cumulative distribution function (CDF) of random variable  $X$  is the function  $P : \mathbb{R} \rightarrow [0, 1]$  defined by*

$$P(x) = \text{Prob}\{X \leq x\}, \quad x \in \mathbb{R}. \quad (2.10)$$

The cumulative distribution function is sometimes simply called the *distribution function*. It has the following properties (inherited from the probability measure):

- (i)  $P$  is a right continuous function of  $x$ ; that is,  $\lim_{\Delta x \rightarrow 0+} P(x + \Delta x) = P(x)$ .
- (ii)  $P$  is a non-decreasing function of  $x$ ; that is, if  $x_1 \leq x_2$ , then  $P(x_1) \leq P(x_2)$ .
- (iii)  $P(-\infty) = 0$ ,  $P(\infty) = 1$ .

The last two properties imply that the cumulative distribution function  $P$  has bounded variation, where the variation of a function is defined as follows.

**Definition 2.2.3** *The  $m$ -variation of a real-valued function  $h$  on the interval  $[\underline{x}, \bar{x}] \subset \mathbb{R}$  is defined as*

$$[h]^{(m)}([\underline{x}, \bar{x}]) = \sup \sum_{j=0}^{l-1} |h(x_{j+1}^l) - h(x_j^l)|^m, \quad (2.11)$$



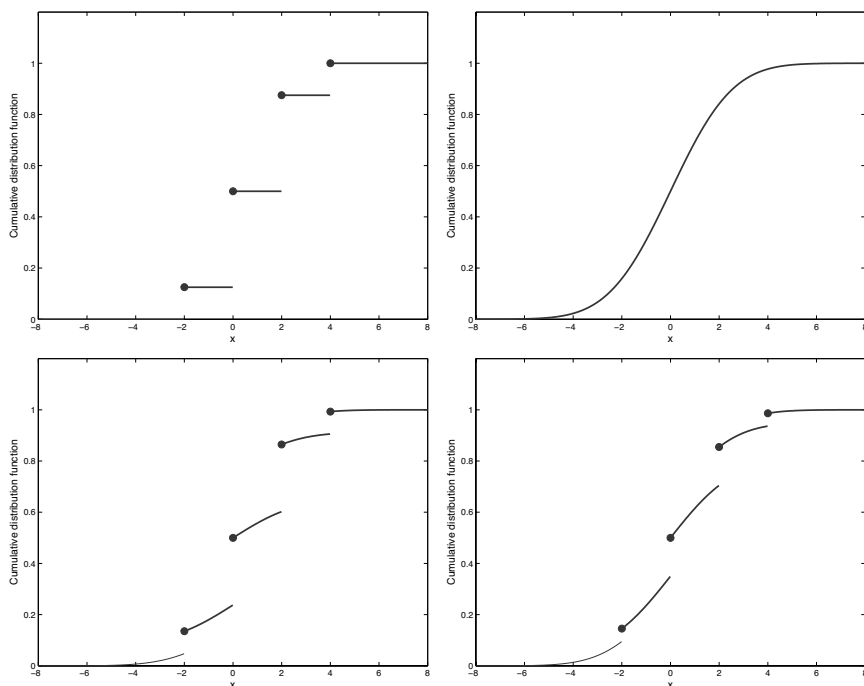
where the supremum is taken over all partitions  $\{x_j^l\}_{j=0}^l$  of  $[\underline{x}, \bar{x}]$ . For the case  $m = 1$  it is simply called variation (it is also called total variation in the literature), and for the case  $m = 2$  it is called quadratic variation. If  $[h]^{(m)}([\underline{x}, \bar{x}])$  is finite, then we say that  $h$  has bounded (or finite)  $m$ -variation on the given interval  $[\underline{x}, \bar{x}]$ .

If  $h$  is a function of  $x \in \mathbb{R}$ , then  $h$  is said to have finite  $m$ -variation if  $[h]^{(m)}([\underline{x}, \bar{x}])$  is finite for any given  $\underline{x}$  and  $\bar{x}$ . In addition,  $h$  is said to have bounded  $m$ -variation if there exists a constant  $c_h$  such that  $[h]^{(m)}([\underline{x}, \bar{x}]) < c_h$  for any  $\underline{x}$  and  $\bar{x}$ , where  $c_h$  is independent of  $\underline{x}$  and  $\bar{x}$ .

The properties of the cumulative distribution function also imply that  $P$  has derivatives almost everywhere with respect to the Lebesgue measure (that is, the set of points at which  $P$  is not differentiable is a null set with respect to the Lebesgue measure), but it should be noted that  $P$  does not have to be equal to the integral of its derivative. These properties also imply that the cumulative distribution function can only have jump discontinuities and it has at most countably many jumps. Thus, we see that the cumulative distribution function could be a step function (illustrated in the upper left panel of Figure 2.1), a continuous function (illustrated in the upper right panel of Figure 2.1) or a function with a mixture of continuous pieces and jumps (illustrated in the bottom plots of Figure 2.1). The last type of cumulative distribution function could result from a convex combination of the first two types of cumulative distribution functions. In fact, this is how we obtained the cumulative distribution functions demonstrated in the bottom plots of Figure 2.1. Specifically, let  $P_j$  be the cumulative distribution function of some random variable,  $j = 1, 2, \dots, m$ , and  $\varpi_j, j = 1, 2, \dots, m$ , be some non-negative numbers such that  $\sum_{j=1}^m \varpi_j = 1$ . Then we easily see that  $\sum_{j=1}^m \varpi_j P_j$  is also a cumulative distribution function. This type of distribution is often called a *mixture distribution*.

**Remark 2.2.2** It is worth noting that there are many different ways found in the literature to define the continuity of a random variable. One is based on the range of the random variable. In this case, a discrete random variable is one with its range consisting of a countable subset in  $\mathbb{R}$  with either a finite or infinite number of elements. For example, the random variable defined in the experiment of tossing a coin three times is a discrete random variable. A continuous random variable is one that takes values in a continuous interval. For example, the random variable defined in the experiment for recording the reaction time to a certain stimulus is a continuous random variable.

Another way to define the continuity of random variables is based on the continuity of the cumulative distribution function for a random variable. Specifically, if the cumulative distribution function is continuous (as illustrated in the upper right panel of Figure 2.1), then the associated random variable is said to be continuous. If the cumulative distribution function is a



**FIGURE 2.1:** Cumulative distribution function: (upper left) a step function  $P_1$ , (upper right) a continuous function  $P_2$ , (lower left)  $0.7P_1 + 0.3P_2$ , (lower right)  $0.4P_1 + 0.6P_2$ .

step function (as illustrated in the upper left panel of Figure 2.1), then the associated random variable is said to be discrete. We see that the discrete random variable defined here is equivalent to that in the first case (as the cumulative distribution function has at most countably many jumps). However, the continuous random variable defined here is more restrictive than that in the first case where the cumulative distribution function of a continuous random variable may be either a continuous function or a function with a mixture of continuous pieces and jumps (as illustrated in the bottom plots of Figure 2.1).

The last way to define the continuity of random variables is based on whether or not a random variable has an associated probability density function (discussed below); specifically, a random variable is said to be continuous if there is a probability density function associated with it. As we shall see below, this definition is even stronger than the second one. In this monograph, we define a random variable to be discrete or continuous based on its range. However, the continuous random variables we will mainly consider in this monograph are those with associated probability density functions.

### 2.2.2 Probability Mass Function

Any discrete random variable has an associated *probability mass function*. Without loss of generality, we assume the range of discrete random variable  $X$  is  $\{x_j\}$ . Then the probability mass function of  $X$  is defined by

$$\Phi(x_j) = \text{Prob}\{X = x_j\}, \quad j = 1, 2, 3, \dots \quad (2.12)$$

Hence, we see that the value of the probability mass function at  $x_j$  is the probability associated with  $x_j$ ,  $j = 1, 2, 3, \dots$ . In addition, by the definition of probability we know that

$$\sum_j \Phi(x_j) = 1.$$

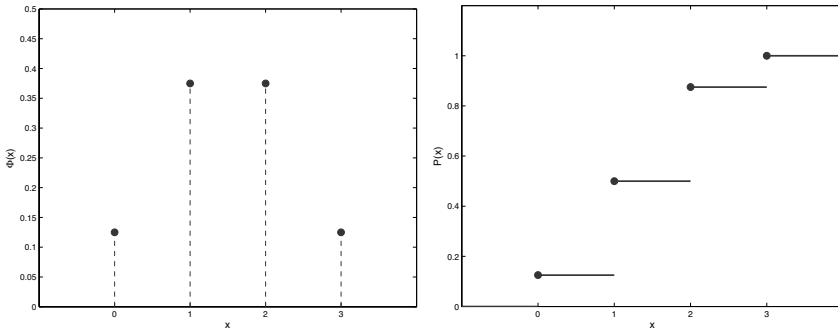
For example, the probability mass function of the discrete random variable defined in the experiment of tossing a coin three times is

$$\Phi(0) = \frac{1}{8}, \quad \Phi(1) = \frac{3}{8}, \quad \Phi(2) = \frac{3}{8}, \quad \Phi(3) = \frac{1}{8}. \quad (2.13)$$

The relationship between the probability mass function and the cumulative distribution function is given by

$$P(x) = \sum_{\{j: x_j \leq x\}} \Phi(x_j), \quad \Phi(x_j) = P(x_j) - \lim_{x \rightarrow x_j -} P(x).$$

In such a case, the cumulative distribution function is a step function, and it is said to be discrete. Figure 2.2 illustrates the probability mass function (2.13) of the discrete random variable defined in the experiment of tossing a coin three times as well as the corresponding cumulative distribution function.



**FIGURE 2.2:** (left) The probability mass function of the discrete random variable defined in the experiment of tossing a coin three times; (right) the corresponding cumulative distribution function.

### 2.2.3 Probability Density Function

If the derivative  $p$  of the cumulative distribution function  $P$  exists for almost all  $x$ , and for all  $x$

$$P(x) = \int_{-\infty}^x p(\xi) d\xi,$$

then  $p$  is called the *probability density function*. It should be noted that a necessary and sufficient condition for a cumulative distribution function to have an associated probability density function is that  $P$  is *absolutely continuous* in the sense that for any positive number  $\epsilon$ , there exists a positive number  $c_l$  such that for any finite collection of disjoint intervals  $(x_j, y_j) \subset \mathbb{R}$  satisfying  $\sum_j |y_j - x_j| < c_l$  then  $\sum_j |P(y_j) - P(x_j)| < \epsilon$ . Thus, we see that

the requirement for absolute continuity is much stronger than that for continuity. This implies that there are cases in which the cumulative distribution function is continuous but not absolutely continuous. An example is the *Cantor distribution* (also called the Cantor–Lebesgue distribution; e.g., see [10, p. 169] or [17, p. 38]), where the cumulative distribution function  $P$  is a constant between the points of the Cantor set (which is an uncountable set in  $\mathbb{R}$  that has Lebesgue measure zero). Such a distribution is often called a *singular continuous distribution* (the derivative of its corresponding cumulative distribution function is zero almost everywhere), and the associated random variable is often termed a *singular continuous random variable*. This type of distribution is rarely encountered in practice. Based on Lebesgue’s decomposition theorem, any cumulative distribution function can be written as a convex combination of a discrete, an absolutely continuous and a singular continuous cumulative distribution function. We refer the interested reader to [2, Section 31] for more information on this topic.

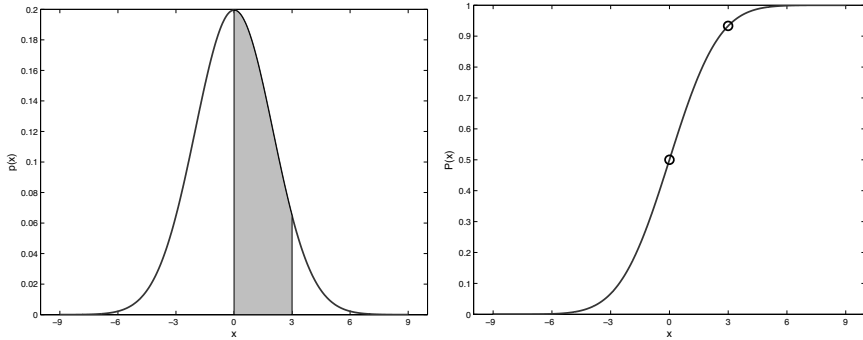
The name “density function” comes from the fact that the probability of the event  $x_1 \leq X \leq x_2$  is given by

$$\begin{aligned} \text{Prob}\{x_1 \leq X \leq x_2\} &= \text{Prob}\{X \leq x_2\} - \text{Prob}\{X \leq x_1\} \\ &= P(x_2) - P(x_1) \\ &= \int_{x_1}^{x_2} p(x) dx. \end{aligned} \tag{2.14}$$

In addition, the probability density function  $p$  satisfies the following properties:

$$p(x) \geq 0, \quad \int_{-\infty}^{\infty} p(x) dx = P(\infty) - P(-\infty) = 1.$$

Figure 2.3 illustrates the probability density function  $p$  (shown in the left panel) of a continuous random variable  $X$  and the corresponding cumulative distribution function  $P$  (shown in the right panel). By (2.14) we know that the shaded area under the probability density function between  $x = 0$  and



**FIGURE 2.3:** (left) Probability density function of a continuous random variable; (right) the corresponding cumulative distribution function.

$x = 3$  is equal to the probability that the value of  $X$  is between 0 and 3, and is also equal to  $P(3) - P(0)$  (where the points  $(0, P(0))$  and  $(3, P(3))$  are indicated by the circles in the right panel of Figure 2.3).

For a discrete random variable, the corresponding probability density function does not exist in the ordinary sense. But it can be constructed in the general sense (a generalized function sense, sometimes called a “distributional sense”) with the help of Dirac delta “functions” [26]. Specifically, for a discrete random variable  $X$  with range  $\{x_j\}$ , its probability density function can be written as

$$p(x) = \sum_j \Phi(x_j) \delta(x - x_j), \quad (2.15)$$

where  $\delta$  is the Dirac delta “function,” that is,  $\delta(x - x_j) = \begin{cases} 0 & \text{if } x \neq x_j \\ \infty & \text{if } x = x_j \end{cases}$

with the property that  $\int_{-\infty}^{\infty} h(x) \delta(x - x_j) dx = h(x_j)$ . This construction has many advantages; for example, as we shall see later, it can put the definition of moments in the same framework for a discrete random variable and for a continuous random variable with a probability density function.

## 2.2.4 Equivalence of Two Random Variables

One can define the equivalence between two random variables in several senses.

**Definition 2.2.4** *Two random variables are said to be equal in distribution (or identically distributed) if their associated cumulative distribution functions are equal.*

**Definition 2.2.5** Two random variables  $X$  and  $Y$  are said to be equal almost surely if  $\text{Prob}\{X = Y\} = 1$ .

It is worth noting that two random variables that are equal almost surely are equal in distribution. However, the fact that two random variables are equal in distribution does not necessarily imply that they are equal almost surely. For example, in the experiment of tossing a fair coin three times, if we define  $X$  as the number of heads obtained in three tosses and  $Y$  as the number of tails observed in three tosses, then it can easily be seen that  $X$  and  $Y$  have the same probability mass function and hence are equal in distribution. However,  $X(\omega) \neq Y(\omega)$  for any  $\omega \in \Omega$ , where  $\Omega$  is defined in (2.3).

### 2.2.5 Joint Distribution Function and Marginal Distribution Function

The definition of cumulative distribution and probability density functions can be extended from one random variable to two or more random variables. In this case, the cumulative distribution function is often called the *joint cumulative distribution function* or simply the *joint distribution function*, and the probability density function is often called the *joint probability density function*. For example, the joint distribution function of two random variables  $X$  and  $Y$  is defined as

$$P(x, y) = \text{Prob}\{X \leq x, Y \leq y\}. \quad (2.16)$$

Similar to the one-dimensional case, the joint distribution function  $P$  is non-negative, non-decreasing and right continuous with respect to each variable. In addition,

$$\begin{aligned} P(-\infty, -\infty) &= 0, & P(x, -\infty) &= 0, & P(-\infty, y) &= 0, \\ P(\infty, \infty) &= 1, & P(x, \infty) &= P_X(x), & P(\infty, y) &= P_Y(y). \end{aligned} \quad (2.17)$$

Here  $P_X$  and  $P_Y$  are, respectively, the cumulative distribution functions of  $X$  and  $Y$ , and the subscript in  $P_X$  ( $P_Y$ ) is used to emphasize that it is the cumulative distribution function of  $X$  ( $Y$ ). (It should be noted that the subscript in the cumulative distribution function is always suppressed if no confusion occurs. Otherwise, we index a distribution function by the random variable it refers to.) In the context of two or more random variables, each random variable is often called a *marginal variable*, and the cumulative distribution function of each random variable is often called a *marginal distribution function*.

The corresponding joint probability density function of  $X$  and  $Y$ , if it exists, is defined as

$$p(x, y) = \frac{\partial^2}{\partial x \partial y} P(x, y). \quad (2.18)$$

Similar to the one-dimensional case, the joint probability density function of  $X$  and  $Y$  is non-negative and

$$\begin{aligned} \text{Prob}\{x_1 \leq X \leq x_2, y_1 \leq Y \leq y_2\} &= \int_{y_1}^{y_2} \int_{x_1}^{x_2} p(x, y) dx dy, \\ \int_{\mathbb{R}^2} p(x, y) dx dy &= 1. \end{aligned} \quad (2.19)$$

By (2.17) and (2.19), we see that the probability density function of  $X$  can be derived from the joint probability density function of  $X$  and  $Y$ , and is given as

$$p_X(x) = \int_{-\infty}^{\infty} p(x, y) dy, \quad (2.20)$$

where the subscript  $X$  in  $p_X$  is used to emphasize that this is the probability density function of  $X$ . (It should be noted that the subscript in the probability density function is always suppressed if no confusion occurs. Otherwise, we index a probability density function by the random variable it refers to.) Again, in the context of two or more random variables, the probability density function of each random variable is called the *marginal probability density function* (or simply the *marginal density function*). Similarly, the probability density function of  $Y$  can be derived from the joint probability density function of  $X$  and  $Y$ ,

$$p_Y(y) = \int_{-\infty}^{\infty} p(x, y) dx. \quad (2.21)$$

Again the subscript  $Y$  in  $p_Y$  is used to emphasize that it is the probability density function of  $Y$ . Thus, by (2.20) and (2.21) we see that the probability density function of each random variable can be derived from the joint probability density function of  $X$  and  $Y$  by integrating across the other variable. However, the converse is usually not true; that is, the joint probability density function usually cannot be derived based on the associated marginal probability density functions.

In general, the joint cumulative distribution function of  $m$  random variables  $X_1, X_2, \dots, X_m$  is defined as

$$P(x_1, x_2, \dots, x_m) = \text{Prob}\{X_1 \leq x_1, X_2 \leq x_2, \dots, X_m \leq x_m\}, \quad (2.22)$$

and the corresponding joint probability density function of  $X_1, X_2, \dots, X_m$ , if it exists, is defined by

$$p(x_1, x_2, \dots, x_m) = \frac{\partial^m}{\partial x_1 \partial x_2 \dots \partial x_m} P(x_1, x_2, \dots, x_m). \quad (2.23)$$

We can view these  $m$  random variables  $X_1, X_2, \dots, X_m$  as the components of an  $m$ -dimensional random vector (i.e.,  $\mathbf{X} = (X_1, X_2, \dots, X_m)^T$ ). Hence, (2.22) and (2.23) can be respectively viewed as the cumulative distribution

function and probability density function of the random vector  $\mathbf{X}$ . In this case, we write these functions more compactly as  $P(\mathbf{x})$  and  $p(\mathbf{x})$  with  $\mathbf{x} = (x_1, x_2, \dots, x_m)^T$ .

Similar to the two-dimensional case, for any  $j \in \{1, 2, \dots, m\}$ , the probability density function of  $X_j$  can be derived from the joint probability density function of  $X_1, X_2, \dots, X_m$  by integrating with respect to all the other  $m - 1$  variables. For example, the probability density function of  $X_2$  is given by

$$p_{X_2}(x_2) = \int_{\mathbb{R}^{m-1}} p(x_1, x_2, \dots, x_m) dx_1 dx_3 dx_4 \cdots dx_m. \quad (2.24)$$

Similarly, the joint probability density function of any two or more random variables in the set of  $\{X_j\}_{j=1}^m$  can be derived from the joint probability density function of  $X_1, X_2, \dots, X_m$  by integrating with respect to all the rest of the variables. For example, the joint probability density function of  $X_1$  and  $X_2$  is

$$p_{X_1 X_2}(x_1, x_2) = \int_{\mathbb{R}^{m-2}} p(x_1, x_2, \dots, x_m) dx_3 dx_4 \cdots dx_m.$$

Here the subscript  $X_1 X_2$  in  $p_{X_1 X_2}$  is used to emphasize that it is the joint probability density function of  $X_1$  and  $X_2$ . Again, the (joint) probability density function of any subset of random variables in the set of  $\{X_j\}_{j=1}^m$  is often called the *marginal probability density function*.

## 2.2.6 Conditional Distribution Function

The conditional distribution function of  $X$  given  $Y = y$  is defined by

$$\mathcal{P}_{X|Y}(x|y) = \text{Prob}\{X \leq x | Y = y\}. \quad (2.25)$$

Then in analogy to (2.5), we have that

$$\mathcal{P}_{X|Y}(x|y) = \frac{\int_{-\infty}^x p_{XY}(\xi, y) d\xi}{p_Y(y)}. \quad (2.26)$$

Here  $p_{XY}$  is the joint probability density function of  $X$  and  $Y$ , and  $p_Y$  is the probability density function of  $Y$ . The corresponding *conditional probability density function* (or simply *conditional density function*) of  $X$  given  $Y = y$ , if it exists, is

$$\rho_{X|Y}(x|y) = \frac{d}{dx} \mathcal{P}_{X|Y}(x|y) = \frac{p_{XY}(x, y)}{p_Y(y)}. \quad (2.27)$$

By (2.27) we see that

$$p_{XY}(x, y) = \rho_{X|Y}(x|y) p_Y(y). \quad (2.28)$$

**Definition 2.2.6** Two random variables  $X$  and  $Y$  are said to be independent if the conditional distribution of  $X$  given  $Y$  is equal to the unconditional distribution of  $X$ .



Definition 2.2.6 implies that if  $X$  and  $Y$  are independent, then we have  $\rho_{X|Y}(x|y) = p_X(x)$ . Hence, by (2.28), we obtain

$$p_{XY}(x, y) = p_X(x)p_Y(y), \quad (2.29)$$

which is also used as a definition for the independence of two random variables. Thus, we see that if two random variables are independent, then the joint probability density function can be determined by the marginal probability density functions. It should be noted that the definition for the independence of two general (either discrete or continuous) random variables is based on the independence of two events, and is given as follows.

**Definition 2.2.7** *Two random variables  $X$  and  $Y$  are said to be independent if the  $\sigma$ -algebras they generate,  $\sigma(X)$  and  $\sigma(Y)$ , are independent (that is, for any  $\mathbb{A} \in \sigma(X)$  and for any  $\mathbb{B} \in \sigma(Y)$ , events  $\mathbb{A}$  and  $\mathbb{B}$  are independent).*

The above definition implies that if  $X$  and  $Y$  are independent, then their joint distribution function is given by the product of their marginal distribution functions; that is,

$$P_{XY}(x, y) = P_X(x)P_Y(y), \quad (2.30)$$

which is also used as a general definition for the independence of two random variables.

In general, for any positive integer  $k \geq 2$ , the *conditional distribution function* of  $X_k$  given  $X_1 = x_1, \dots, X_{k-1} = x_{k-1}$  is

$$\begin{aligned} & \mathcal{P}_{X_k|X_{k-1}, \dots, X_1}(x_k|x_{k-1}, \dots, x_1) \\ &= \text{Prob}\{X_k \leq x_k | X_{k-1} = x_{k-1}, \dots, X_1 = x_1\} \\ &= \frac{\int_{-\infty}^{x_k} p_{X_1, \dots, X_k}(x_1, \dots, x_{k-1}, \xi) d\xi}{p_{X_1, \dots, X_{k-1}}(x_1, \dots, x_{k-1})}. \end{aligned} \quad (2.31)$$

Here  $p_{X_1, \dots, X_{k-1}}$  is the joint probability density function of  $X_1, \dots, X_{k-1}$ , and  $p_{X_1, \dots, X_k}$  is the joint probability density function of  $X_1, \dots, X_k$ . The corresponding *conditional density function* of  $X_k$  given  $X_1 = x_1, \dots, X_{k-1} = x_{k-1}$  is

$$\begin{aligned} & \rho_{X_k|X_{k-1}, \dots, X_1}(x_k|x_{k-1}, \dots, x_1) \\ &= \frac{d}{dx_k} \mathcal{P}_{X_k|X_{k-1}, \dots, X_1}(x_k|x_{k-1}, \dots, x_1) \\ &= \frac{p_{X_1, \dots, X_k}(x_1, \dots, x_k)}{p_{X_1, \dots, X_{k-1}}(x_1, \dots, x_{k-1})}, \end{aligned} \quad (2.32)$$

which implies that

$$\begin{aligned} & p_{X_1, \dots, X_k}(x_1, \dots, x_k) \\ &= \rho_{X_k|X_{k-1}, \dots, X_1}(x_k|x_{k-1}, \dots, x_1) p_{X_1, \dots, X_{k-1}}(x_1, \dots, x_{k-1}). \end{aligned} \quad (2.33)$$

Hence, by (2.27) and the above equation we see that the joint probability density function  $p_{X_1, \dots, X_m}$  of any  $m$  random variables  $X_1, X_2, \dots, X_m$  can be written as

$$\begin{aligned} p_{X_1, \dots, X_m}(x_1, x_2, \dots, x_m) \\ = p_{X_1}(x_1) \rho_{X_2|X_1}(x_2|x_1) \rho_{X_3|X_2, X_1}(x_3|x_2, x_1) \\ \cdots \rho_{X_m|X_{m-1}, \dots, X_1}(x_m|x_{m-1}, \dots, x_1). \end{aligned} \quad (2.34)$$

The concept of the independence of two random variables can be extended to the case of a sequence of random variables, and it is given as follows.

**Definition 2.2.8** *Random variables  $X_1, X_2, \dots, X_m$  are said to be mutually independent if*

$$p_{X_1, \dots, X_m}(x_1, x_2, \dots, x_m) = \prod_{j=1}^m p_{X_j}(x_j), \quad (2.35)$$

where  $p_{X_1, \dots, X_m}$  is the joint probability density function of  $m$  random variables  $X_1, X_2, \dots, X_m$ , and  $p_{X_j}$  is the probability density function of  $X_j$ ,  $j = 1, 2, \dots, m$ .

In general, the conditional probability density function of a subset of the coordinates of  $(X_1, X_2, \dots, X_m)$  given the values of the remaining coordinates is obtained by dividing the joint probability density function of  $(X_1, X_2, \dots, X_m)$  by the marginal probability density function of the remaining coordinates. For example, the conditional probability density function of  $(X_{k+1}, \dots, X_m)$  given  $X_1 = x_1, \dots, X_k = x_k$  (where  $k$  is a positive integer such that  $1 < k < m$ ) is defined as

$$\rho_{X_{k+1}, \dots, X_m|X_1, \dots, X_k}(x_{k+1}, \dots, x_m | x_1, \dots, x_k) = \frac{p_{X_1, \dots, X_m}(x_1, \dots, x_m)}{p_{X_1, \dots, X_k}(x_1, \dots, x_k)}, \quad (2.36)$$

where  $p_{X_1, \dots, X_k}$  is the joint probability density function of  $X_1, \dots, X_k$ , and  $p_{X_1, \dots, X_m}$  is the joint probability density function of  $X_1, \dots, X_m$ . The concept of mutually independent random variables can also be extended to that of mutually independent random vectors. The definition is given as follows.

**Definition 2.2.9** *Random vectors  $\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_m$  are said to be mutually independent if*

$$p_{\mathbf{X}_1, \dots, \mathbf{X}_m}(\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m) = \prod_{j=1}^m p_{\mathbf{X}_j}(\mathbf{x}_j), \quad (2.37)$$

where  $p_{\mathbf{X}_1, \dots, \mathbf{X}_m}$  is the joint probability density function of  $m$  random vectors  $\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_m$ , and  $p_{\mathbf{X}_j}$  is the probability density function of  $\mathbf{X}_j$ ,  $j = 1, 2, \dots, m$ .

### 2.2.7 Function of a Random Variable

Let  $X$  be a random variable with cumulative distribution function  $P_X$  and probability density function  $p_X$ . Then by Definition 2.2.1 we know that for any measurable function  $\eta : \mathbb{R} \rightarrow \mathbb{R}$ , the composite function  $\eta(X)$  is measurable and indeed is a random variable. Let  $Y = \eta(X)$ . Then the cumulative distribution function of  $Y$  is

$$P_Y(y) = \text{Prob}\{\eta(X) \leq y\}. \quad (2.38)$$

If we assume that  $\eta$  is a monotone function, then by (2.38) we have

$$P_Y(y) = \begin{cases} \text{Prob}\{X \leq \eta^{-1}(y)\} = P_X(\eta^{-1}(y)), & \eta \text{ is increasing} \\ \text{Prob}\{X \geq \eta^{-1}(y)\} = 1 - P_X(\eta^{-1}(y)), & \eta \text{ is decreasing.} \end{cases} \quad (2.39)$$

If we further assume that  $\eta^{-1}$  is differentiable, then differentiating both sides of (2.39) yields the probability density function of  $Y$

$$p_Y(y) = p_X(\eta^{-1}(y)) \left| \frac{d\eta^{-1}(y)}{dy} \right|. \quad (2.40)$$

In general, we consider an  $m$ -dimensional random vector  $\mathbf{X}$  with probability density function  $p_{\mathbf{X}}$ . Let  $\mathbf{Y} = \boldsymbol{\eta}(\mathbf{X})$  with  $\boldsymbol{\eta} = (\eta_1, \eta_2, \dots, \eta_m)^T$  and  $\eta_j : \mathbb{R}^m \rightarrow \mathbb{R}$  be a measurable function for all  $j$ . Assume that  $\boldsymbol{\eta}$  has a unique inverse  $\boldsymbol{\eta}^{-1}$ . Then the probability density function of the  $m$ -dimensional random vector  $\mathbf{Y}$  is given by

$$p_{\mathbf{Y}}(\mathbf{y}) = p_{\mathbf{X}}(\boldsymbol{\eta}^{-1}(\mathbf{y})) |\mathcal{J}|, \quad (2.41)$$

where  $\mathcal{J}$  is the determinant of the Jacobian matrix  $\frac{\partial \mathbf{x}}{\partial \mathbf{y}}$  with its  $(j, k)$ th element being  $\frac{\partial x_j}{\partial y_k}$ .

In the following, whenever we talk about transformation of a random variable or random vector, we always assume that the transformation is measurable so that the resulting function is also a random variable. The following theorem is about the transformation of two independent random variables.

**Theorem 2.2.3** *Let  $X$  and  $Z$  be independent random variables,  $\eta_X$  be a function only of  $x$  and  $\eta_Z$  be a function only of  $z$ . Then the random variables  $U = \eta_X(X)$  and  $V = \eta_Z(Z)$  are independent.*

The above theorem is very important in theory, and it can be generalized as follows.

**Theorem 2.2.4** *Let  $\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_m$  be mutually independent random vectors, and  $\eta_j$  be a function only of  $\mathbf{x}_j$ ,  $j = 1, 2, \dots, m$ . Then random variables  $U_j = \eta_j(\mathbf{X}_j)$ ,  $j = 1, 2, \dots, m$ , are mutually independent.*

### 2.3 Statistical Averages of Random Variables

The concepts of moments of a single random variable and the joint moments between any pair of random variables in a multi-dimensional set of random variables are of particular importance in practice. We begin the discussion of these statistical averages by considering first a single random variable  $X$  and its cumulative distribution function  $P$ . The *expectation* (also called an *expected value* or *mean*) of the random variable  $X$  is defined by

$$\mathbb{E}(X) = \int_{-\infty}^{\infty} x dP(x). \quad (2.42)$$

Here  $\mathbb{E}(\cdot)$  is called the *expectation operator* (or statistical averaging operator), and the integral on the right side is interpreted as a Riemann–Stieltjes integral. We remark that the Riemann–Stieltjes integral is a generalization of the Riemann integral and its definition is given as follows.

**Definition 2.3.1** Let  $\varphi$  and  $h$  be real-valued functions defined on  $[\underline{x}, \bar{x}] \subset \mathbb{R}$ ,  $\{x_j^l\}_{j=0}^l$  be a partition of  $[\underline{x}, \bar{x}]$ ,  $\Delta_l = \max_{0 \leq j \leq l-1} \{x_{j+1}^l - x_j^l\}$ , and  $s_j^l \in [x_j^l, x_{j+1}^l]$  denote intermediate points of the partition. Then  $\varphi$  is said to be Riemann–Stieltjes integrable with respect to  $h$  on  $[\underline{x}, \bar{x}]$  if

$$\lim_{\substack{l \rightarrow \infty \\ \Delta_l \rightarrow 0}} \sum_{j=0}^{l-1} \varphi(s_j^l) [h(x_{j+1}^l) - h(x_j^l)] \quad (2.43)$$

exists and the limit is independent of the choice of the partition and their intermediate points. The limit of (2.43) is called the Riemann–Stieltjes integral of  $\varphi$  with respect to  $h$  on  $[\underline{x}, \bar{x}]$ , and is denoted by  $\int_{\underline{x}}^{\bar{x}} \varphi(x) dh(x)$ ; that is,

$$\int_{\underline{x}}^{\bar{x}} \varphi(x) dh(x) = \lim_{\substack{l \rightarrow \infty \\ \Delta_l \rightarrow 0}} \sum_{j=0}^{l-1} \varphi(s_j^l) [h(x_{j+1}^l) - h(x_j^l)],$$

where  $\varphi$  and  $h$  are called the *integrand* and the *integrator*, respectively.

It is worth noting that the Riemann–Stieltjes integral  $\int_{\underline{x}}^{\bar{x}} \varphi(x) dh(x)$  does not exist for all continuous functions  $\varphi$  on  $[\underline{x}, \bar{x}]$  unless  $h$  has bounded variation. This is why (2.42) can be interpreted as a Riemann–Stieltjes integral (as  $P$  has bounded variation). In general, the Riemann–Stieltjes integral  $\int_{\underline{x}}^{\bar{x}} \varphi(x) dh(x)$  exists if the following conditions are satisfied:

- The functions  $\varphi$  and  $h$  have no discontinuities at the same point  $x \in [\underline{x}, \bar{x}]$ .
- The function  $\varphi$  has bounded  $\kappa_\varphi$ -variation on  $[\underline{x}, \bar{x}]$  and the function  $h$  has bounded  $\kappa_h$ -variation on  $[\underline{x}, \bar{x}]$ , where  $\kappa_\varphi$  and  $\kappa_h$  are some positive constants such that  $\frac{1}{\kappa_\varphi} + \frac{1}{\kappa_h} > 1$ .

We refer the interested reader to [18, Section 2.1], [21], and the references therein for more information on Riemann–Stieltjes integrals.

If we assume that  $P$  is absolutely continuous (that is, the corresponding probability density function  $p$  exists), then we can rewrite (2.42) as

$$\mathbb{E}(X) = \int_{-\infty}^{\infty} xp(x) dx. \quad (2.44)$$

We note from (2.12) and (2.15) that if  $X$  is a discrete random variable with range  $\{x_j\}$ , then the expectation of  $X$  is given by

$$\mathbb{E}(X) = \sum_j x_j \text{Prob}\{X = x_j\}.$$

Thus, we see that with the help of the Dirac delta function one can put the definition of expectation in the same framework for a discrete random variable and for a continuous random variable with a probability density function (as we stated earlier). Since we are mainly interested in discrete random variables and those continuous random variables associated with probability density functions, we will define the statistical average of a random variable in terms of its probability density function in the following presentation.

The expectation of a random variable is also called the *first moment*. In general, the  $k$ th moment of a random variable  $X$  is defined as

$$\mathbb{E}(X^k) = \int_{-\infty}^{\infty} x^k p(x) dx.$$

We can also define the *central moments*, which are the moments of the difference between  $X$  and  $\mathbb{E}(X)$ . For example, the  $k$ th central moment of  $X$  is defined by

$$\mathbb{E}((X - \mathbb{E}(X))^k).$$

Of particular importance is the second central moment, called the *variance* of  $X$ , which is defined as

$$\sigma^2 = \text{Var}(X) = \mathbb{E}((X - \mathbb{E}(X))^2) = \int_{-\infty}^{\infty} (x - \mathbb{E}(X))^2 p(x) dx. \quad (2.45)$$

The square root  $\sigma$  of the variance of  $X$  is called the *standard deviation* of  $X$ . Variance is a measure of the “randomness” of the random variable  $X$ . It is

related to the first and second moments through the relationship

$$\begin{aligned}\text{Var}(X) &= \mathbb{E}((X - \mathbb{E}(X))^2) \\ &= \mathbb{E}(X^2 - 2X\mathbb{E}(X) + (\mathbb{E}(X))^2) \\ &= \mathbb{E}(X^2) - (\mathbb{E}(X))^2.\end{aligned}\tag{2.46}$$

One of the useful concepts in understanding the variation of a random variable  $X$  is the *coefficient of variation* (CV), which is defined as the ratio of the standard deviation to the mean (that is,  $\text{CV} = \sqrt{\text{Var}(X)}/\mathbb{E}(X)$ ). It is the inverse of the so-called *signal-to-noise ratio*. A random variable with  $\text{CV} < 1$  is considered to have low variation, while one with  $\text{CV} > 1$  is considered to have high variation.

Note that the moments of a function of a random variable,  $Y = \eta(X)$ , can be defined in the same way as above. For example, the  $k$ th moment of  $Y$  is

$$\mathbb{E}(Y^k) = \int_{\Omega_Y} y^k p_Y(y) dy, \tag{2.47}$$

where  $\Omega_Y$  denotes the range of  $Y$ , and  $p_Y$  is the probability density function of  $Y$ . However, due to the relation (2.40) between  $p_Y$  and the probability density function  $p_X$  of  $X$ , we can also calculate the  $k$ th moment of  $Y$  by

$$\mathbb{E}(Y^k) = \mathbb{E}\{\eta^k(X)\} = \int_{-\infty}^{\infty} \eta^k(x) p_X(x) dx. \tag{2.48}$$

For any real numbers  $a$  and  $b$ , by (2.48) we observe that

$$\mathbb{E}(aX + b) = a\mathbb{E}(X) + b, \quad \text{Var}(aX + b) = a^2\text{Var}(X). \tag{2.49}$$

**Remark 2.3.1** *We remark that the infinite sequence of moments is in general not enough to uniquely determine a distribution function. Interested readers can refer to [7, Section 2.3] for a counterexample of two random variables having the same moments but different probability density functions. However, if two random variables have bounded support, then an infinite sequence of moments does uniquely determine the distribution function.*

### 2.3.1 Joint Moments

The definition of moments and central moments can be extended from one random variable to two or more random variables. In this context, the moments are often called joint moments, and central moments are often called joint central moments. For example, the joint moment of two random variables  $X$  and  $Y$  is defined as

$$\mathbb{E}(X^{k_x} Y^{k_y}) = \int_{\mathbb{R}^2} x^{k_x} y^{k_y} p(x, y) dx dy, \tag{2.50}$$

where  $k_x$  and  $k_y$  are positive integers, and  $p$  is the joint probability density function of  $X$  and  $Y$ . The joint central moment of  $X$  and  $Y$  is given by

$$\mathbb{E}((X - \mathbb{E}(X))^{k_x}(Y - \mathbb{E}(Y))^{k_y}). \quad (2.51)$$

However, the joint moment that is most useful in practical applications is the *correlation* of two random variables  $X$  and  $Y$ , defined as

$$\text{Cor}\{X, Y\} = \mathbb{E}(XY) = \int_{\mathbb{R}^2} xyp(x, y)dxdy, \quad (2.52)$$

which implies that the second moment of  $X$  is the correlation of  $X$  with itself. Also of particular importance is the *covariance* of two random variables  $X$  and  $Y$ , defined as

$$\text{Cov}\{X, Y\} = \mathbb{E}((X - \mathbb{E}(X))(Y - \mathbb{E}(Y))) = \mathbb{E}(XY) - \mathbb{E}(X)\mathbb{E}(Y), \quad (2.53)$$

which indicates that the variance of  $X$  is the covariance of  $X$  with itself. It is worth noting that the correlation of two random variables should not be confused with the *correlation coefficient*,  $r$ , which is defined as the covariance of the two random variables divided by the product of their standard deviations.

$$r = \frac{\text{Cov}\{X, Y\}}{\sqrt{\text{Var}(X)\text{Var}(Y)}}.$$

By (2.52) and (2.53) we see that both correlation and covariance are symmetric. That is,

$$\text{Cor}\{X, Y\} = \text{Cor}\{Y, X\}, \quad \text{Cov}\{X, Y\} = \text{Cov}\{Y, X\}.$$

Moreover, they are both linear in each variable. That is, for any real numbers  $a$  and  $b$  we have

$$\begin{aligned} \text{Cor}\{aX + bY, Z\} &= a\text{Cor}\{X, Z\} + b\text{Cor}\{Y, Z\}, \\ \text{Cor}\{Z, aX + bY\} &= a\text{Cor}\{Z, X\} + b\text{Cor}\{Z, Y\}, \\ \text{Cov}\{aX + bY, Z\} &= a\text{Cov}\{X, Z\} + b\text{Cov}\{Y, Z\}, \\ \text{Cov}\{Z, aX + bY\} &= a\text{Cov}\{Z, X\} + b\text{Cov}\{Z, Y\}. \end{aligned}$$

By the above equations, we find that the variance of  $aX + bY$  is given by

$$\text{Var}(aX + bY) = a^2\text{Var}(X) + 2ab\text{Cov}\{X, Y\} + b^2\text{Var}(Y).$$

**Definition 2.3.2** Two random variables  $X$  and  $Y$  are called *uncorrelated* if  $\text{Cov}\{X, Y\} = 0$ .

By the above definition and (2.53) we see that if

$$\mathbb{E}(XY) = \mathbb{E}(X)\mathbb{E}(Y), \quad (2.54)$$

then  $X$  and  $Y$  are uncorrelated. We observe from (2.29) that if  $X$  and  $Y$  are independent, then (2.54) also holds. Therefore, the independence of two random variables implies that these two random variables are uncorrelated. However, in general the converse is not true.

### 2.3.2 Conditional Moments

The  $k$ th *conditional moment* of  $X$  given  $Y = y$  is

$$\mathbb{E}(X^k|Y = y) = \int_{-\infty}^{\infty} x^k \rho_{X|Y}(x|y) dx. \quad (2.55)$$

This implies that if random variables  $X$  and  $Y$  are independent (that is,  $\rho_{X|Y}(x|y) = p_X(x)$ ), then

$$\mathbb{E}(X^k|Y = y) = \mathbb{E}(X^k).$$

The first conditional moment of  $X$  given  $Y = y$  is called the *conditional expectation* of  $X$  given  $Y = y$ .

Observe that  $\mathbb{E}(X|Y = y)$  is a function of  $y$ . Hence,  $\mathbb{E}(X|Y)$  is a random variable, which is called the *conditional expectation* of  $X$  given  $Y$ . By (2.20), (2.27) and (2.55) we find

$$\begin{aligned} \mathbb{E}(\mathbb{E}(X|Y)) &= \int_{-\infty}^{\infty} p_Y(y) \left( \int_{-\infty}^{\infty} x \rho_{X|Y}(x|y) dx \right) dy \\ &= \int_{-\infty}^{\infty} p_Y(y) \left( \int_{-\infty}^{\infty} x \frac{p_{XY}(x, y)}{p_Y(y)} dx \right) dy \\ &= \int_{-\infty}^{\infty} x \left( \int_{-\infty}^{\infty} p_{XY}(x, y) dy \right) dx \\ &= \int_{-\infty}^{\infty} x p_X(x) dx \\ &= \mathbb{E}(X), \end{aligned}$$

which indicates that the expected value of the conditional expectation of  $X$  given  $Y$  is the same as the expected value of  $X$ . This formula is often called the *law of total expectation*. Similarly, we can show that for any positive integer  $k$  we have

$$\mathbb{E}(\mathbb{E}(X^k|Y)) = \mathbb{E}(X^k). \quad (2.56)$$

The *conditional variance* of  $X$  given  $Y$  is defined as

$$\text{Var}(X|Y) = \mathbb{E}((X - \mathbb{E}(X|Y))^2|Y) = \mathbb{E}(X^2|Y) - (\mathbb{E}(X|Y))^2. \quad (2.57)$$

By (2.46), (2.56) and (2.57), we find that the unconditional variance is related to the conditional variance by

$$\text{Var}(X) = \mathbb{E}(\text{Var}(X|Y)) + \text{Var}(\mathbb{E}(X|Y)), \quad (2.58)$$

which indicates that the unconditional variance is equal to the sum of the mean of the conditional variance and the variance of the conditional mean. Equation (2.58) is often called the *law of total variance*, *variance decomposition formula* or *conditional variance formula*.



### 2.3.3 Statistical Averages of Random Vectors

For an  $m$ -dimensional random vector  $\mathbf{X} = (X_1, \dots, X_m)^T$ , its *mean vector* and *covariance matrix* are of particular importance. Specifically, the mean vector of  $\mathbf{X}$  is given by

$$\boldsymbol{\mu} = \mathbb{E}(\mathbf{X}) = (\mathbb{E}(X_1), \dots, \mathbb{E}(X_m))^T,$$

and its *covariance matrix* is defined by

$$\boldsymbol{\Sigma} = \text{Var}(\mathbf{X}) = \mathbb{E}\{(\mathbf{X} - \boldsymbol{\mu})(\mathbf{X} - \boldsymbol{\mu})^T\}. \quad (2.59)$$

Hence, we see that  $\boldsymbol{\Sigma} \in \mathbb{R}^{m \times m}$  is a non-negative definite matrix with its  $(k, j)$ th element being the covariance of random variables  $X_k$  and  $X_j$ :

$$\begin{aligned} \text{Cov}\{X_k, X_j\} &= \mathbb{E}((X_k - \mathbb{E}(X_k))(X_j - \mathbb{E}(X_j))) \\ &= \int_{\mathbb{R}^2} (x_k - \mathbb{E}(X_k))(x_j - \mathbb{E}(X_j)) p_{X_k X_j}(x_k, x_j) dx_k dx_j, \end{aligned}$$

where  $p_{X_k X_j}$  is the joint probability density function of  $X_k$  and  $X_j$ . For any  $\mathcal{A} \in \mathbb{R}^{l \times m}$  and  $\mathbf{a} \in \mathbb{R}^l$ , it can be easily shown that

$$\mathbb{E}(\mathcal{A}\mathbf{X} + \mathbf{a}) = \mathcal{A}\mathbb{E}(\mathbf{X}) + \mathbf{a}, \quad \text{Var}(\mathcal{A}\mathbf{X} + \mathbf{a}) = \mathcal{A}\text{Var}(\mathbf{X})\mathcal{A}^T. \quad (2.60)$$

Similarly, we can extend the covariance of two random variables to the *cross-covariance matrix* between two random vectors. Let  $\mathbf{X} = (X_1, X_2, \dots, X_{m_x})^T$  and  $\mathbf{Y} = (Y_1, Y_2, \dots, Y_{m_y})^T$ , where  $m_x$  and  $m_y$  are positive integers. Then the *cross-covariance matrix* between  $\mathbf{X}$  and  $\mathbf{Y}$  is

$$\begin{aligned} \text{Cov}\{\mathbf{X}, \mathbf{Y}\} &= \mathbb{E}((\mathbf{X} - \mathbb{E}(\mathbf{X}))(\mathbf{Y} - \mathbb{E}(\mathbf{Y}))^T) \\ &= \mathbb{E}(\mathbf{X}\mathbf{Y}^T) - \mathbb{E}(\mathbf{X})(\mathbb{E}(\mathbf{Y}))^T, \end{aligned} \quad (2.61)$$

which implies that  $\text{Var}(\mathbf{X}) = \text{Cov}\{\mathbf{X}, \mathbf{X}\}$ . By (2.60) it can be easily shown that for any matrix  $\mathcal{A} \in \mathbb{R}^{\kappa_x \times m_x}$  and  $\mathcal{B} \in \mathbb{R}^{\kappa_y \times m_y}$  ( $\kappa_x$  and  $\kappa_y$  are positive integers) we have

$$\text{Cov}\{\mathcal{A}\mathbf{X}, \mathcal{B}\mathbf{Y}\} = \mathcal{A}\text{Cov}\{\mathbf{X}, \mathbf{Y}\}\mathcal{B}^T. \quad (2.62)$$

The concept of uncorrelatedness of two random variables can also be extended to two random vectors. Specifically, two random vectors  $\mathbf{X}$  and  $\mathbf{Y}$  are called uncorrelated if  $\text{Cov}\{\mathbf{X}, \mathbf{Y}\} = 0$ .

### 2.3.4 Important Inequalities

In this section, we list a few important inequalities that will be used in this monograph.