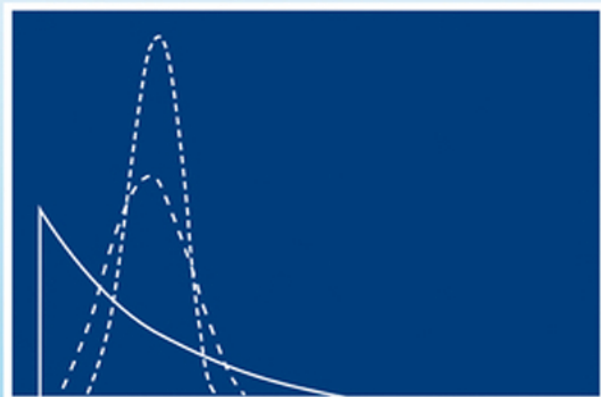
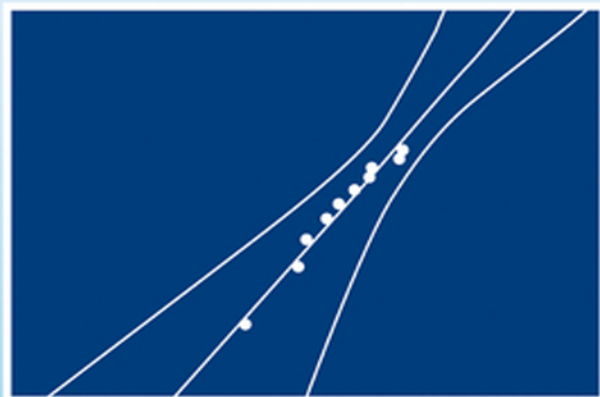


Statistics and Probability with Applications for Engineers and Scientists



Bhisham C. Gupta • Irwin Guttman

***STATISTICS AND PROBABILITY
WITH APPLICATIONS FOR
ENGINEERS AND SCIENTISTS***

***STATISTICS AND PROBABILITY
WITH APPLICATIONS FOR
ENGINEERS AND SCIENTISTS***

Bhisham C. Gupta

*Professor of Statistics
University of Southern Maine
Portland, ME*

Irwin Guttman

*Professor Emeritus of Statistics
SUNY at Buffalo and
University of Toronto, Canada*

WILEY

Copyright © 2013 by John Wiley & Sons, Inc. All rights reserved.

Published by John Wiley & Sons, Inc., Hoboken, New Jersey.
Published simultaneously in Canada.

No part of this publication may be reproduced, stored in a retrieval system, or transmitted in any form or by any means, electronic, mechanical, photocopying, recording, scanning, or otherwise, except as permitted under Section 107 or 108 of the 1976 United States Copyright Act, without either the prior written permission of the Publisher, or authorization through payment of the appropriate per-copy fee to the Copyright Clearance Center, Inc., 222 Rosewood Drive, Danvers, MA 01923, (978) 750-8400, fax (978) 750-4470, or on the web at www.copyright.com. Requests to the Publisher for permission should be addressed to the Permissions Department, John Wiley & Sons, Inc., 111 River Street, Hoboken, NJ 07030, (201) 748-6011, fax (201) 748-6008, or online at <http://www.wiley.com/go/permission>.

Limit of Liability/Disclaimer of Warranty: While the publisher and author have used their best efforts in preparing this book, they make no representations or warranties with respect to the accuracy or completeness of the contents of this book and specifically disclaim any implied warranties of merchantability or fitness for a particular purpose. No warranty may be created or extended by sales representatives or written sales materials. The advice and strategies contained herein may not be suitable for your situation. You should consult with a professional where appropriate. Neither the publisher nor author shall be liable for any loss of profit or any other commercial damages, including but not limited to special, incidental, consequential, or other damages.

For general information on our other products and services or for technical support, please contact our Customer Care Department within the United States at (800) 762-2974, outside the United States at (317) 572-3993 or fax (317) 572-4002.

Wiley also publishes its books in a variety of electronic formats. Some content that appears in print may not be available in electronic formats. For more information about Wiley products, visit our web site at www.wiley.com.

Library of Congress Cataloging-in-Publication Data:

Gupta, Bhisham C., 1942-

Statistics and probability with applications for engineers and scientists / Bhisham C. Gupta,
Department of Mathematics and Statistics, University of Southern Maine, Portland, ME,
Irwin Guttman, Amherst, NY.

pages cm

Includes bibliographical references and index.

ISBN 978-1-118-46404-5 (hardback)

1. Probabilities. 2. Mathematical statistics. I. Guttman, Irwin. II. Title.

QA273.G85 2013

519.5-dc23

2012032883

Printed in the United States of America

10 9 8 7 6 5 4 3 2 1

In the loving memory of my parents, Roshan Lal and Sodhan Devi
-Bhisham

In the loving memory of my parents, Anna and Samuel Guttman
-Irwin

Modern statisticians are familiar with the notion that any finite body of data contains only a limited amount of information on any point under examination; that this limit is set by the nature of the data themselves, and cannot be increased by any amount of ingenuity expended in their statistical examination: that the statistician's task, in fact, is limited to the extraction of the whole of the available information on any particular issue.

R. A. Fisher

Table of Contents

Preface

xvii

Chapter 1 | Introduction

1

- 1.1 Designed Experiment 2
 - 1.1.1 Motivation for the Study 2
 - 1.1.2 Investigation 2
 - 1.1.3 Changing Criteria 2
 - 1.1.4 A Summary of the Various Phases of the Investigation 3
- 1.2 A Survey 5
- 1.3 An Observational Study 6
- 1.4 A Set of Historical Data 6
- 1.5 A Brief Description of What is Covered in This Book 6

PART I

Chapter 2 | Describing Data Graphically and Numerically

11

- 2.1 Getting Started with Statistics 12
 - 2.1.1 What Is Statistics? 12
 - 2.1.2 Population and Sample in a Statistical Study 12
- 2.2 Classification of Various Types of Data 15
 - 2.2.1 Nominal Data 15
 - 2.2.2 Ordinal Data 16
 - 2.2.3 Interval Data 16
 - 2.2.4 Ratio Data 16
- 2.3 Frequency Distribution Tables for Qualitative and Quantitative Data 17
 - 2.3.1 Qualitative Data 17
 - 2.3.2 Quantitative Data 20
- 2.4 Graphical Description of Qualitative and Quantitative Data 25
 - 2.4.1 Dot Plot 25
 - 2.4.2 Pie Chart 25
 - 2.4.3 Bar Chart 27
 - 2.4.4 Histograms 30
 - 2.4.5 Line Graph 35
 - 2.4.6 Stem-and-Leaf Plot 37
- 2.5 Numerical Measures of Quantitative Data 41
 - 2.5.1 Measures of Centrality 42
 - 2.5.2 Measures of Dispersion 46
- 2.6 Numerical Measures of Grouped Data 55
 - 2.6.1 Mean of a Grouped Data 56
 - 2.6.2 Median of a Grouped Data 56
 - 2.6.3 Mode of a Grouped Data 57
 - 2.6.4 Variance of a Grouped Data 57

vii

2.7	Measures of Relative Position	59
2.7.1	Percentiles	59
2.7.2	Quartiles	60
2.7.3	Interquartile Range	60
2.7.4	Coefficient of Variation	61
2.8	Box-Whisker Plot	62
2.8.1	Construction of a Box Plot	62
2.8.2	How to Use the Box Plot	63
2.9	Measures of Association	68
2.10	Case Studies	71
2.11	Using JMP [®]	73
	Review Practice Problems	73

Chapter 3 | *Elements of Probability*

83

3.1	Introduction	84
3.2	Random Experiments, Sample Spaces, and Events	84
3.2.1	Random Experiments and Sample Spaces	84
3.2.2	Events	85
3.3	Concepts of Probability	88
3.4	Techniques of Counting Sample Points	93
3.4.1	Tree Diagram	93
3.4.2	Permutations	94
3.4.3	Combinations	95
3.4.4	Arrangements of n Objects Involving Several Kinds of Objects	96
3.5	Conditional Probability	98
3.6	Bayes's Theorem	100
3.7	Introducing Random Variables	104
	Review Practice Problems	105

Chapter 4 | *Discrete Random Variables and Some Important Discrete Probability Distributions*

111

4.1	Graphical Descriptions of Discrete Distributions	112
4.2	Mean and Variance of a Discrete Random Variable	113
4.2.1	Expected Value of Discrete Random Variables and Their Functions	113
4.2.2	The Moment-Generating Function—Expected Value of a Special Function of X	115
4.3	The Discrete Uniform Distribution	117
4.4	The Hypergeometric Distribution	119
4.5	The Bernoulli Distribution	122
4.6	The Binomial Distribution	123
4.7	The Multinomial Distribution	126
4.8	The Poisson Distribution	128
4.8.1	Definition and Properties of the Poisson Distribution	128
4.8.2	Poisson Process	128
4.8.3	Poisson Distribution as a Limiting Form of the Binomial	128
4.9	The Negative Binomial Distribution	132
4.10	Some Derivations and Proofs (Optional)	135
4.11	A Case Study	135
4.12	Using JMP	135
	Review Practice Problems	136

Chapter 5	 	<i>Continuous Random Variables and Some Important Continuous Probability Distributions</i>	143
5.1		Continuous Random Variables	144
5.2		Mean and Variance of Continuous Random Variables	146
5.2.1		Expected Value of Continuous Random Variables and Their Function	146
5.2.2		The Moment-Generating Function–Expected Value of a Special Function of X	149
5.3		Chebychev’s Inequality	151
5.4		The Uniform Distribution	152
5.4.1		Definition and Properties	152
5.4.2		Mean and Standard Deviation of the Uniform Distribution	155
5.5		The Normal Distribution	157
5.5.1		Definition and Properties	157
5.5.2		The Standard Normal Distribution	158
5.5.3		The Moment-Generating Function of the Normal Distribution	164
5.6		Distribution of Linear Combination of Independent Normal Variables	165
5.7		Approximation of the Binomial and Poisson Distribution by the Normal Distribution	169
5.7.1		Approximation of the Binomial Distribution by the Normal Distribution	169
5.7.2		Approximation of the Poisson Distribution by the Normal Distribution	171
5.8		A Test of Normality	171
5.9		Probability Models Commonly Used in Reliability Theory	175
5.9.1		The Lognormal Distribution	176
5.9.2		The Exponential Distribution	180
5.9.3		The Gamma Distribution	184
5.9.4		The Weibull Distribution	187
5.10		A Case Study	191
5.11		Using JMP	192
		Review Practice Problems	192
Chapter 6	 	<i>Distribution of Functions of Random Variables</i>	199
6.1		Introduction	200
6.2		Distribution Functions of Two Random Variables	200
6.2.1		Case of Two Discrete Random Variables	200
6.2.2		Case of Two Continuous Random Variables	202
6.2.3		The Mean Value and Variance of Functions of Two Random Variables	204
6.2.4		Conditional Distributions	206
6.2.5		Correlation between Two Random Variables	208
6.2.6		Bivariate Normal Distribution	211
6.3		Extension to Several Random Variables	214
6.4		The Moment-Generating Function Revisited	214
		Review Practice Problems	218
Chapter 7	 	<i>Sampling Distributions</i>	223
7.1		Random Sampling	224
7.1.1		Random Sampling from an Infinite Population	224
7.1.2		Random Sampling from a Finite Population	225

7.2	The Sampling Distribution of the Mean	228
7.2.1	Normal Sampled Population	228
7.2.2	Nonnormal Sampled Population	228
7.2.3	The Central Limit Theorem	228
7.3	Sampling from a Normal Population	234
7.3.1	The Chi-Square Distribution	234
7.3.2	The Student t -Distribution	240
7.3.3	Snedecor's F -Distribution	244
7.4	Order Statistics	247
7.5	Using JMP	247
	Review Practice Problems	247

Chapter 8 | Estimation of Population Parameters

251

8.1	Introduction	252
8.2	Point Estimators for the Population Mean and Variance	252
8.2.1	Properties of Point Estimators	253
8.2.2	Methods of Finding Point Estimators	256
8.3	Interval Estimators for the Mean μ of a Normal Population	262
8.3.1	σ^2 Known	262
8.3.2	σ^2 Unknown	264
8.3.3	Sample Size Is Large	266
8.4	Interval Estimators for the Difference of Means of Two Normal Populations	272
8.4.1	Variances Are Known	272
8.4.2	Variances Are Unknown	273
8.5	Interval Estimators for the Variance of a Normal Population	280
8.6	Interval Estimator for the Ratio of Variances of Two Normal Populations	284
8.7	Point and Interval Estimators for the Parameters of Binomial Populations	288
8.7.1	One Binomial Population	288
8.7.2	Two Binomial Populations	290
8.8	Determination of Sample Size	294
8.8.1	One Population Mean	294
8.8.2	Difference of Two Population Means	295
8.8.3	One Population Proportion	296
8.8.4	Difference of Two Population Proportions	296
8.9	Some Supplemental Information	298
8.10	A Case Study	298
8.11	Using JMP	299
	Review Practice Problems	299

Chapter 9 | Hypothesis Testing

307

9.1	Introduction	308
9.2	Basic Concepts of Testing a Statistical Hypothesis	308
9.2.1	Hypothesis Formulation	308
9.2.2	Risk Assessment	310
9.3	Tests Concerning the Mean of a Normal Population Having Known Variance	312
9.3.1	Case of a One-Tail (Left-Sided) Test	312
9.3.2	Case of a One-Tail (Right-Sided) Test	316
9.3.3	Case of a Two-Tail Test	317

9.4	Tests Concerning the Mean of a Normal Population Having Unknown Variance	324
9.4.1	Case of a Left-Tail Test	324
9.4.2	Case of a Right-Tail Test	326
9.4.3	The Two-Tail Case	326
9.5	Large Sample Theory	330
9.6	Tests Concerning the Difference of Means of Two Populations Having Distributions with Known Variances	332
9.6.1	The Left-Tail Test	332
9.6.2	The Right-Tail Test	333
9.6.3	The Two-Tail Test	334
9.7	Tests Concerning the Difference of Means of Two Populations Having Normal Distributions with Unknown Variances	339
9.7.1	Two Population Variances Are Equal	339
9.7.2	Two Population Variances Are Unequal	342
9.7.3	The Paired t -Test	344
9.8	Testing Population Proportions	349
9.8.1	Test Concerning One Population Proportion	349
9.8.2	Test Concerning the Difference between Two Population Proportions	351
9.9	Tests Concerning the Variance of a Normal Population	355
9.10	Tests Concerning the Ratio of Variances of Two Normal Populations	358
9.11	Testing of Statistical Hypotheses Using Confidence Intervals	362
9.12	Sequential Tests of Hypotheses	367
9.12.1	A One-Tail Sequential Testing Procedure	367
9.12.2	A Two-Tail Sequential Testing Procedure	371
9.13	Case Studies	374
9.14	Using JMP	375
	Review Practice Problems	375

PART II

Chapter 10 | *Elements of Reliability Theory* 389

10.1	The Reliability Function	390
10.1.1	The Hazard Rate Function	391
10.1.2	Employing the Hazard Function	398
10.2	Estimation: Exponential Distribution	399
10.3	Hypothesis Testing: Exponential Distribution	406
10.4	Estimation: Weibull Distribution	407
10.5	Case Studies	414
10.6	Using JMP	416
	Review Practice Problems	416

Chapter 11 | *Statistical Quality Control—Phase I Control Charts* 419

11.1	Basic Concepts of Quality and Its Benefits	420
11.2	What a Process Is and Some Valuable Tools	420
11.2.1	Check Sheet	422
11.2.2	Pareto Chart	422
11.2.3	Cause-and-Effect (Fishbone or Ishikawa) Diagram	425
11.2.4	Defect Concentration Diagram	427
11.3	Common and Assignable Causes	427
11.3.1	Process Evaluation	427

11.3.2	Action on the Process	428
11.3.3	Action on Output	428
11.3.4	Variation	428
11.4	Control Charts	429
11.4.1	Preparation for Use of Control Charts	430
11.4.2	Benefits of a Control Chart	431
11.4.3	Control Limits Versus Specification Limits	433
11.5	Control Charts for Variables	434
11.5.1	Shewhart \bar{X} and R Control Charts	434
11.5.2	Shewhart \bar{X} and R Control Charts When Process Mean μ and Process Standard Deviation σ Are Known	440
11.5.3	Shewhart \bar{X} and S Control Charts	441
11.6	Control Charts for Attributes	448
11.6.1	The p Chart: Control Chart for the Fraction of Nonconforming Units	449
11.6.2	The p Chart: Control Chart for the Fraction Nonconforming with Variable Sample Sizes	454
11.6.3	The np Control Chart: Control Chart for the Number of Nonconforming Units	456
11.6.4	The c Control Chart	458
11.6.5	The u Control Chart	461
11.7	Process Capability	468
11.8	Case Studies	470
11.9	Using JMP	472
	Review Practice Problems	472
Chapter 12	<i>Statistical Quality Control—Phase II Control Charts</i>	479
12.1	Introduction	480
12.2	Basic Concepts of CUSUM Control Chart	480
12.3	Designing a CUSUM Control Chart	483
12.3.1	Two-Sided CUSUM Control Chart Using a Numerical Procedure	484
12.3.2	The Fast Initial Response (FIR) Feature for CUSUM Control Chart	489
12.3.3	The Combined Shewhart–CUSUM Control Chart	492
12.3.4	The CUSUM Control Chart for Controlling Process Variability	493
12.4	The Moving Average (MA) Control Chart	495
12.5	The Exponentially Weighted Moving Average (EWMA) Control Chart	499
12.6	Case Studies	504
12.7	Using JMP	505
	Review Practice Problems	506
Chapter 13	<i>Analysis of Categorical Data</i>	509
13.1	Introduction	509
13.2	The Chi-Square Goodness-of-Fit Test	510
13.3	Contingency Tables	517
13.3.1	The 2×2 Case Parameters Known	517
13.3.2	The 2×2 Case with Unknown Parameters	519
13.3.3	The $r \times s$ Contingency Table	521
13.4	Chi-Square Test for Homogeneity	525
13.5	Comments on the Distribution of the Lack-of-Fit Statistics	528
13.6	Case Studies	529
	Review Practice Problems	531

Chapter 14 | Nonparametric Tests 537

- 14.1 Introduction 537
- 14.2 The Sign Test 538
 - 14.2.1 One-Sample Test 538
 - 14.2.2 The Wilcoxon Signed-Rank Test 541
 - 14.2.3 Two-Sample Test 543
- 14.3 Mann–Whitney (Wilcoxon) W Test for Two Samples 548
- 14.4 Runs Test 551
 - 14.4.1 Runs Above and Below the Median 551
 - 14.4.2 The Wald–Wolfowitz Run Test 553
- 14.5 Spearman Rank Correlation 556
- 14.6 Using JMP 559
- Review Practice Problems 559

Chapter 15 | Simple Linear Regression Analysis 565

- 15.1 Introduction 566
- 15.2 Fitting the Simple Linear Regression Model 567
 - 15.2.1 Simple Linear Regression Model 567
 - 15.2.2 Fitting a Straight Line by Least Squares 569
 - 15.2.3 Sampling Distribution of the Estimators of Regression Coefficients 573
- 15.3 Unbiased Estimator of σ^2 578
- 15.4 Further Inferences Concerning Regression Coefficients (β_0, β_1), $E(Y)$, and Y 580
 - 15.4.1 Confidence Interval for β_1 with Confidence Coefficient $(1 - \alpha)$ 580
 - 15.4.2 Confidence Interval for β_0 with Confidence Coefficient $(1 - \alpha)$ 581
 - 15.4.3 Confidence Interval for $E(Y|X)$ with Confidence Coefficient $(1 - \alpha)$ 582
 - 15.4.4 Prediction Interval for a Future Observation Y with Confidence Coefficient $(1 - \alpha)$ 585
- 15.5 Tests of Hypotheses for β_0 and β_1 590
 - 15.5.1 Test of Hypotheses for β_1 590
 - 15.5.2 Test of Hypotheses for β_0 590
- 15.6 Analysis of Variance Approach to Simple Linear Regression Analysis 596
- 15.7 Residual Analysis 601
- 15.8 Transformations 609
- 15.9 Inference About ρ 615
- 15.10 A Case Study 618
- 15.11 Using JMP 619
- Review Practice Problems 619

Chapter 16 | Multiple Linear Regression Analysis 627

- 16.1 Introduction 628
- 16.2 Multiple Linear Regression Models 628
- 16.3 Estimation of Regression Coefficients 632
 - 16.3.1 Estimation of Regression Coefficients Using Matrix Notation 633
 - 16.3.2 Properties of the Least-Squares Estimators 635
 - 16.3.3 The Analysis of Variance Table 636
 - 16.3.4 More Inferences about Regression Coefficients 639

16.4	Multiple Linear Regression Model Using Quantitative and Qualitative Predictor Variables	646
16.4.1	Single Qualitative Variable with Two Categories	646
16.4.2	Single Qualitative Variable with Three or More Categories	647
16.5	Standardized Regression Coefficients	658
16.5.1	Multicollinearity	660
16.5.2	Consequences of Multicollinearity	661
16.6	Building Regression Type Prediction Models	662
16.6.1	First Variable to Enter into the Model	662
16.7	Residual Analysis and Certain Criteria for Model Selection	665
16.7.1	Residual Analysis	665
16.7.2	Certain Criteria for Model Selection	667
16.8	Logistic Regression	672
16.9	Case Studies	676
16.10	Using JMP	677
	Review Practice Problems	678

Chapter 17 | Analysis of Variance

685

17.1	Introduction	686
17.2	The Design Models	686
17.2.1	Estimable Parameters	686
17.2.2	Estimable Functions	688
17.3	One-Way Experimental Layouts	689
17.3.1	The Model and Its Analysis	689
17.3.2	Confidence Intervals for Treatment Means	695
17.3.3	Multiple Comparisons	700
17.3.4	Determination of Sample Size	706
17.3.5	The Kruskal–Wallis Test for One-Way Layouts (Nonparametric Method)	707
17.4	Randomized Complete Block Designs	710
17.4.1	The Friedman F_r -Test for Randomized Complete Block Design (Nonparametric Method)	718
17.4.2	Experiments with One Missing Observation in an RCB-Design Experiment	719
17.4.3	Experiments with Several Missing Observations in an RCB-Design Experiment	719
17.5	Two-Way Experimental Layouts	722
17.5.1	Two-Way Experimental Layouts with One Observation per Cell	724
17.5.2	Two-Way Experimental Layouts with $r > 1$ Observations per Cell	725
17.5.3	Blocking in Two-Way Experimental Layouts	734
17.5.4	Extending Two-Way Experimental Designs to n -Way Experimental Layouts	734
17.6	Latin Square Designs	736
17.7	Random-Effects and Mixed-Effects Models	742
17.7.1	Random-Effects Model	742
17.7.2	Mixed-Effects Model	744
17.7.3	Nested (Hierarchical) Designs	746
17.8	A Case Study	752
17.9	Using JMP	753
	Review Practice Problems	753

Chapter 18 | *The 2^k Factorial Designs* 765

- 18.1 Introduction 766
- 18.2 The Factorial Designs 766
- 18.3 The 2^k Factorial Design 768
- 18.4 Unreplicated 2^k Factorial Designs 776
- 18.5 Blocking in the 2^k Factorial Design 782
 - 18.5.1 Confounding in the 2^k Factorial Design 783
 - 18.5.2 Yates's Algorithm for the 2^k Factorial Designs 788
- 18.6 The 2^k Fractional Factorial Designs 790
 - 18.6.1 One-half Replicate of a 2^k Factorial Design 790
 - 18.6.2 One-quarter Replicate of a 2^k Factorial Design 795
- 18.7 Case Studies 799
- 18.8 Using JMP 801
- Review Practice Problems 801

Chapter 19 | *Response Surfaces*

This chapter is not included in text, but is available for download via the book's website: www.wiley.com/go/statsforengineers

Appendices 807**Appendix A | *Statistical Tables* 809****Appendix B | *Answers to Selected Problems* 845****Appendix C | *Bibliography* 863*****Index* 867**

Preface

AUDIENCE

This is an introductory textbook in applied statistics and probability for undergraduate students in engineering and the natural sciences. It begins at a level suitable for those with no previous exposure to probability and statistics and carries the reader through to a level of proficiency in various techniques of statistics. This text is divided into two parts: Part I discusses descriptive statistics, concepts of probability, probability distributions, sampling distributions, estimation, and testing of hypotheses, and Part II discusses various topics of applied statistics, including some reliability theory, statistical quality control charts of phase I and phase II, some nonparametric techniques, categorical data analysis, simple and multiple linear regression analysis, design and analysis of variance with emphasis on 2^k factorial designs, and response surface methodology.

This text is suitable for a one- or two-semester undergraduate course sequence. The presentation of material gives instructors a lot of flexibility to pick and choose topics they feel should make up the coverage of material for their courses. However, we feel that in the first course for engineers and science majors, one may cover Chapters 1 and 2, a brief discussion of probability in Chapter 3, selected discrete and continuous distributions from Chapters 4 and 5 with more emphasis on normal distribution, Chapters 7 through 9, and couple of topics from Part II that meet the needs and interests of the particular group of students. For example, some discussion of the material on regression analysis and design of experiments in Chapters 15 and 17 may serve well. A two-semester course may cover the entire book. The only prerequisite is a first course in calculus, which all engineering and science students are required to take.

Because of space considerations, some proofs and derivations, certain advance level topics of interest, including Chapter 19 on response surfaces, are not included in the text, but are available for download via the book's website: www.wiley.com/go/statsforengineers.

MOTIVATION

Students encounter data-analysis problems in many areas of engineering or natural science curricula. Engineers and scientists in their professional lives often encounter situations requiring analysis of data arising from their areas of practice. Very often they have to plan the investigation that generates data (an activity euphemistically called the design of experiments), analyze the data obtained, and interpret the results. Other problems and investigations may pertain to the maintenance of quality of existing products or the development of new products, or to a desired outcome in an investigation of the underlying mechanisms governing a certain process.

Knowing how to “design” a particular investigation to obtain reliable data must be coupled with knowledge of descriptive and inferential statistical tools to analyze properly and interpret such data. The intent of this textbook is to expose the uninitiated to statistical methods that deal with the generation of data for different (but frequently met) types of investigations and to discuss how to analyze and interpret the generated data.

HISTORY

This text has its roots in the three editions of *Introductory Engineering Statistics*, first co-authored by Irwin Guttman and the late, great Samuel Wilks. Professor J. Stuart Hunter (Princeton University), one of the finest expositors in the statistics profession, a noted researcher, and a colleague of Professor Wilks, joined Professor Guttman to produce editions two and three. All editions were published by John Wiley & Sons, with the third edition appearing in 1982.

APPROACH

In this text we emphasize both descriptive and inferential statistics. We first give details of descriptive statistics and then continue with an elementary discussion of the fundamentals of probability theory underlying many of the statistical techniques discussed in this text. We next cover a wide range of statistical techniques such as statistical estimation, regression methods, statistical quality control (with emphasis on phase I and phase II control charts), and process capability indices, nonparametric methods, elements of reliability theory, and the like. A feature of these discussions is that all statistical concepts are supported by a large number of examples using data encountered in real-life situations. We also illustrate how the statistical packages MINITAB[®] Version 16, Microsoft Excel[®] Version windows 2007, and JMP[®] Version 9, may be used to aid in the analysis of various data sets.

Another feature of this text is the coverage at an adequate and understandable level of the design of experiments. This includes a discussion of randomized block designs, one- and two-way designs, Latin square designs, 2^k factorial designs, response surface designs, among others. As indicated above, all this is illustrated with real-life situations and accompanying data sets, supported by MINITAB, Microsoft Excel, and JMP. We know of no other book in the market that covers all these software packages.

HALLMARK FEATURES

Software Integration

As indicated above, we incorporate MINITAB and Microsoft Excel throughout the text and present JMP at the end of each chapter. Our step-by-step approach to the use of the software packages means no prior knowledge of their use is required. After completing a course that uses this text, students will be able to use these software packages to analyze statistical data in their fields of interest.

Breadth of Coverage

Besides the coverage of many popular statistical techniques, we include discussion of certain aspects of sampling distributions, nonparametric tests, phase II control charts, reliability theory, design of experiments, and response surface methodology. Phase II control charts are discussed in a separate chapter that includes the use of the statistical packages to implement these charts.

Design of experiments and response surface methodology are treated in sufficient breadth and depth to be appropriate for a two-course sequence in engineering statistics that includes probability and the design of experiments.

Real data in examples and homework problems illustrate the importance of statistics and probability as a tool for engineers and scientists in their professional lives. All the data sets with 20 or more data points are available on the website in three formats: MINITAB, Microsoft Excel, and JMP.

Case studies in each chapter further illustrate the importance of statistical techniques in professional practice.

STUDENT RESOURCES

Data sets for all examples and homework exercises from the text are available to students on the website in Minitab, Microsoft Excel, and JMP format. The sample data sets were generated using well-known statistical sampling procedures, ensuring that we are dealing with *random samples*. An inkling of what this may entail is given throughout the text. (See, for example, Section 7.1.2.) The field of sampling is an active topic among research statisticians and practitioners, and references to sampling techniques are widely available in books and journal articles. Some of these references are included in the bibliography section.

Other resources on the book website www.wiley.com/go/statsforengineers available for download include:

Solutions Manual to all odd numbered homework exercises in the text.

Excel worksheets and macros to supplement features that are not available in Excel.

INSTRUCTOR RESOURCES

The following resources are available to adopting instructors on the textbook website www.wiley.com/go/statsforengineers:

Solutions Manual to all homework exercises in the text.

Lecture slides to aid instructors preparing for lectures.

Data sets for all examples and homework exercises from the book, in three formats: Minitab, Microsoft Excel, and JMP.

Excel worksheets and macros to supplement features that are not available in Excel.

Errata

We have thoroughly reviewed the text to make sure it is as error-free as possible. However, any errors discovered will be listed on the textbook website.

If you encounter any errors as you are using the book, please send them directly to the authors (bcgupta@usm.maine.edu), so that the errors can be corrected in a timely manner on the website, and for future editions. We also welcome any suggestions for improvement you may have, and thank you in advance for helping us improve the book for future readers.

ACKNOWLEDGMENTS

We are grateful to the following reviewers and colleagues whose comments and suggestions were invaluable in improving the text:

Zaid Abdo, University of Idaho

Erin Baker, University of Massachusetts

Bob Barnet, University of Wisconsin-Platteville

Mark Gebert, University of Kentucky

Ramesh Gupta, University of Maine

Rameshwar Gupta, University of New Brunswick, Canada

Xiaochun Jiang, North Carolina Agricultural and Technical State University

Dennis Johnston, Baylor University

Gerald Keller, Joseph L. Rotman School of Management, University of Toronto

Kyungduk Ko, Boise State University

Paul Kvam, Georgia Institute of Technology
Thunshun Liao, Louisiana State University
Jye-Chyi Lu, Georgia Institute of Technology
Sumona Mondal, Clarkson University
Janbiao Pan, California Poly State University
Anastassios Perakis, University of Michigan
David Powers, Clarkson University
Ali Touran, Northeastern University
Leigh Williams, Virginia Polytechnic and State University
Tian Zheng, Columbia University
Jingyi Zhu, University of Utah

We thank William Belcher, Darwin Davis, Julie Ellis, Pushpa Gupta, Mohamad Ibourk, James Lucas, Mary McShane-Vaughn, Louis Neveux, and Phil Ramsey who helped find suitable data sets for the case studies. We also thank Laurie McDermott for her help in typing some parts of this manuscript. Special thanks are due to Eric Laflamme for helping write JMP/Excel procedures and creating PowerPoint[®] presentations, George Bernier for helping write Excel work books and macros, and Patricia Miller and Brenda Townsend for editing Power Point Slides and some parts of the manuscript.

We acknowledge Minitab Inc., SAS Institute Inc., and Microsoft for permitting us to print MINITAB, JMP, and Microsoft Excel screen shots in this book.

Portions of the text are reproduced by permission of the American Society for Quality (ASQ), *Applied Statistics for the Six Sigma Green Belt* and *Statistical Quality Control for the Six Sigma Green Belt* by Bhisham C. Gupta and H. Fred Walker (Milwaukee: ASQ Quality Press, 2005, 2007). To order these books, call ASQ at 800-248-1946 or 414-272-8575, or visit <http://www.asq.org/quality-press>.

We would also like to express our thanks and appreciation to the individuals at John Wiley, for their support, confidence, and guidance as we have worked together to develop this project.

The authors would like to gratefully thank their families. Bhisham acknowledges the patience and support of his wife, Swarn; daughters, Anita and Anjali; son, Shiva; sons-in-laws, Prajay and Mark; daughter-in-law, Aditi; and wonderful grandchildren, Priya, Kaviya, and Ayush and Amari. For their patience and support, Irwin is grateful to his wife, Mary; son, Daniel; daughters, Karen and Shaun; wonderful grandchildren, Liam, Teia, and Sebastian; brothers and their better halves, Alvin and Rita, and Stanley and Gloria.

BHISHAM GUPTA
IRWIN GUTTMAN

1

Introduction

Statistics, the discipline, is the study of the scientific method. In pursuing this discipline, Statisticians have developed a set of techniques that are extensively used to solve problems in any field of scientific endeavor, such as in the engineering sciences, biological sciences, and the chemical, pharmaceutical, and social sciences.

This book is concerned with discussing these techniques and their applications for certain experimental situations. It begins at a level suitable for those with no previous exposure to probability and statistics, and carries the reader through to a level of proficiency in various techniques of statistics.

In all scientific areas, whether engineering, biological sciences, medicine, chemical, pharmaceutical, or social sciences, scientists are inevitably confronted with problems that need to be investigated. Consider some examples:

- An engineer wants to determine the role of an electronic component needed to detect the malfunction of the engine of a plane.
- A biologist wants to study various aspects of wildlife, the origin of a disease, or the genetic aspects of a wild animal.
- A medical researcher is interested in determining the cause of a certain type of cancer.
- A manufacturer of lenses wants to study the quality of the finishing on intraocular lenses.
- A chemist is interested in determining the effect of a catalyst in the production of low-density polyethylene.
- A pharmaceutical company is interested in developing a vaccination for swine flu.
- A social scientist is interested in exploring a particular aspect of human society.

In all of the examples above, the first and foremost work is to define clearly the objective of the study and precisely formulate the problem. The next important step is to gather information to help determine what key factors are affecting the problem. Remember that to determine these factors successfully, you should understand not merely statistical methodology but relevant nonstatistical knowledge as well. Once the problem is formulated and the key factors of the problem are identified, the next step is to

collect the data. There are various methods of data collecting. Four basic methods of statistical data collecting are:

- A designed experiment
- A survey
- An observational study
- A set of historical data, that is, data collected by an organization or an individual in an earlier study

1.1 DESIGNED EXPERIMENT

We discuss the concept of a designed experiment with an example, “Development of Screening Facility for Storm Water Overflows” (taken from Box, Hunter, and Hunter, 1978, and used with permission). The example illustrates how a sequence of experiments can enable scientists to gain knowledge of the various *important factors* affecting the problem, and give insight into the objectives of the investigation. It also indicates how unexpected features of the problem can become dominant, and how experimental difficulties can occur so that certain planned experiments cannot be run at all. Most of all, this example shows the importance of common sense in the conduct of any experimental investigation. The reader may rightly conclude from this example that the course of a real investigation, like that of true love, seldom runs smoothly, although the eventual outcome may be satisfactory.

1.1.1 Motivation for the Study

During heavy rainstorms, the total flow coming to a sewage treatment plant may exceed its capacity, making it necessary to bypass the excess flow around the treatment plant, as shown in Figure 1.1.1a. Unfortunately, the storm overflow of untreated sewage causes pollution of the receiving body of water. A possible alternative, sketched in Figure 1.1.1b, is to screen most of the solids out of the overflow in some way and return them to the plant for treatment. Only the less objectionable screened overflow is discharged directly to the river.

To determine whether it was economical to construct and operate such a screening facility, the Federal Water Pollution Control Administration of the Department of the Interior sponsored a research project at the Sullivan Gulch pump station in Portland, Oregon. Usually the flow to the pump station was 20 million gallons per day (mgd), but during a storm, the flow could exceed 50 mgd.

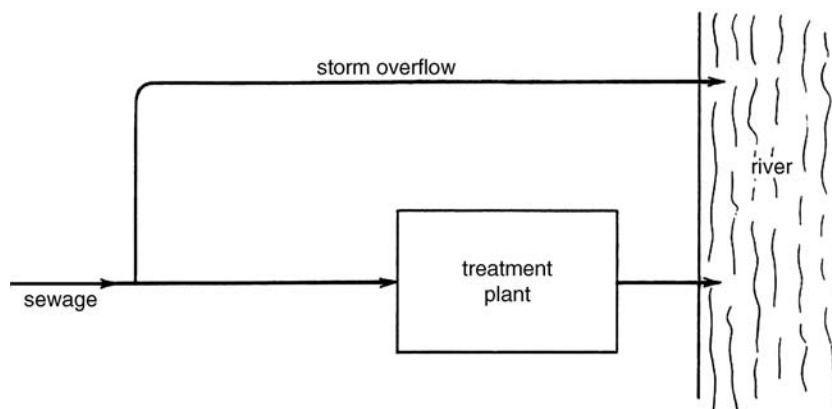
Figure 1.1.2a shows the original version of the experimental screening unit, which could handle approximately 1000 gallons per minute (gpm). Figure 1.1.2a is a perspective view and Figure 1.1.2b is a simplified schematic diagram. A single unit was about seven feet high and seven feet in diameter. The flow of raw sewage struck a rotating collar screen at a velocity of five to 15 feet per second. This speed was a function of the flow rate into the unit and hence a function of the diameter of the influent pipe. Depending on the speed of the rotation of this screen and its fineness, up to 90% of the feed penetrated the collar screen. The rest of the feed dropped to the horizontal screen, which vibrated to remove excess water. The solids concentrate, which passed through neither screen, was sent to the sewage treatment plant. Unfortunately, during operation the screens became clogged with solid matter, not only sewage but also oil, paint, and fish-packing wastes. Backwash sprays were therefore installed for both screens to permit cleaning during operation.

1.1.2 Investigation

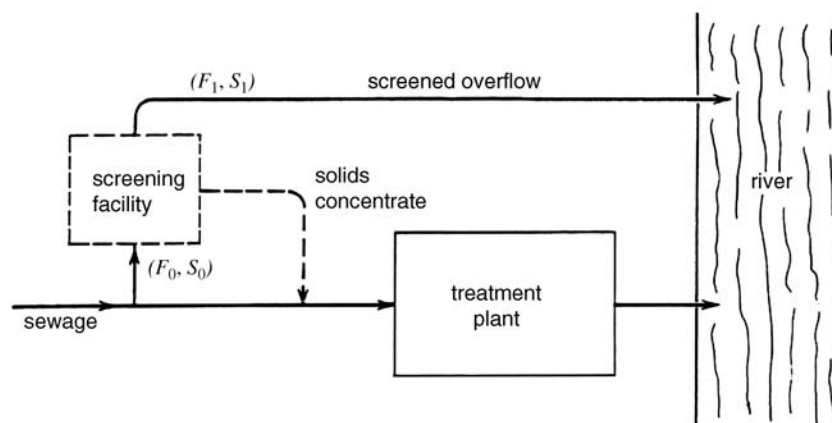
The objective of the investigation was to determine good operating conditions.

1.1.3 Changing Criteria

What are good operating conditions? Initially it was believed they were those resulting in the highest possible removal of solids. In Figure 1.1.1b, settleable solids in the influent are denoted by S_0 and



(a) Standard mode of operation.

(b) Modified mode of operation, with screening facility. F = flow, S = settleable solids.**FIGURE 1.1.1** Operation of the sewage treatment plant.

the settleable solids in the effluent by S_1 . The *percent solids removed* by the screen is therefore $y = 100(S_0 - S_1)/S_0$. Thus, initially, it was believed that good operation meant achieving a high value for y . However, it became evident after the first set of experimental made that the *percentage of the flow retreated* (flow returned to treatment plant), which we denote by z , also had to be taken into account. In Figure 1.1.1b, influent flow to the screens is denoted by F_0 and effluent flow from the screens to the river by F_1 . Thus $z = 100(F_0 - F_1)/F_0$.

1.1.4 A Summary of the Various Phases of the Investigation

Phase a

In this initial phase an experiment was run in which the roles of three variables were studied: collar screen mesh size (fine, coarse), horizontal screen mesh size (fine, coarse), and flow rate (gallons per minute). At this stage:

1. The experimenters were encouraged by the generally high values achieved for y .
2. Highest values for y were apparently achieved by using a horizontal screen with a coarse mesh and a collar screen with fine mesh.

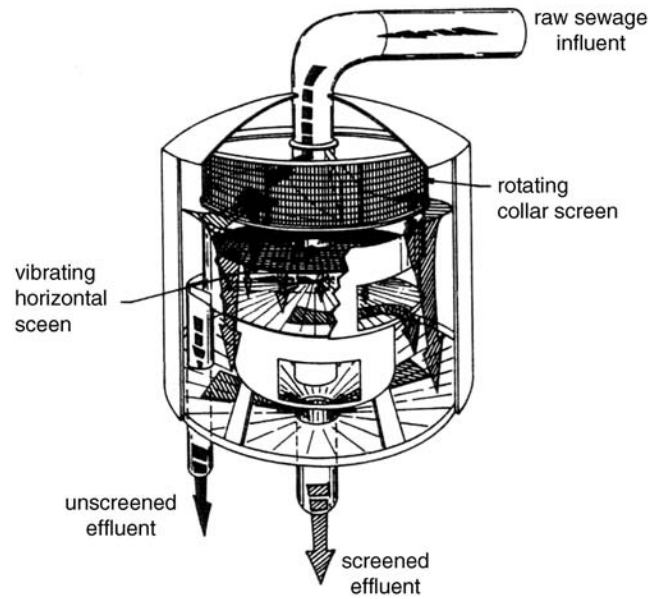


FIGURE 1.1.2a Original version of the screening unit (detailed diagram).

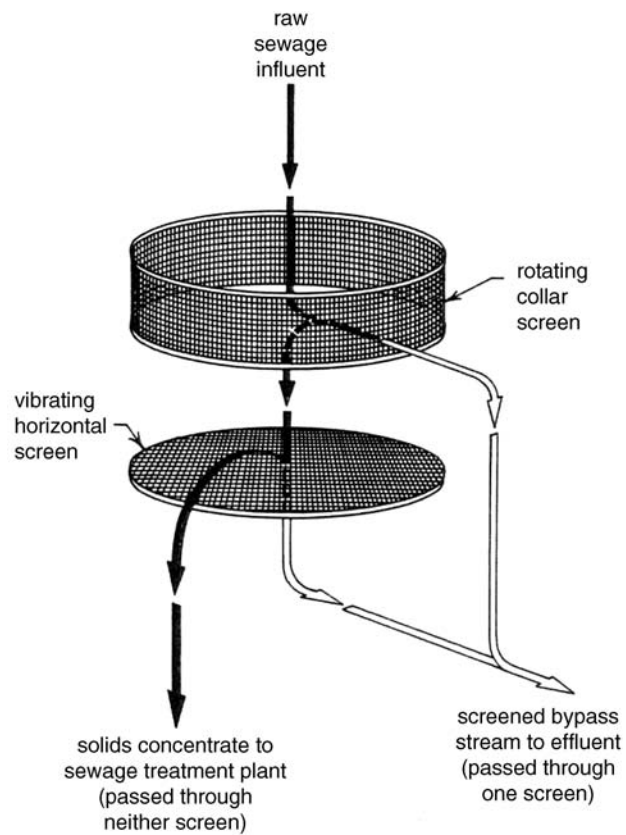


FIGURE 1.1.2b Original version of the screening unit (simplified diagram).

3. Contrary to expectation, flow rate did not show up as an important variable affecting y .
4. Most important, the experiment was unexpectedly dominated by the z values, which measure the flow retreated. These were uniformly very low, with about 0.01% of the flow being returned to the treatment plant and 99.9% leaving the screen for discharge into the river. Although it was desirable that the retreated flow be small, the z values were embarrassingly low. As the experimenters remarked, “[T]he horizontal screen produced a solid concentrate . . . dry enough to shovel . . . This represented a waste of effort of concentrating because the concentrated solids were intended to *flow* from the units.”

Phase b

It was now clear (1) that z as well as y was important and (2) that z was too low. It was conjectured that the matters might be improved by removing the horizontal screen altogether. Another experiment was therefore performed with no horizontal screen. The speed of rotation of the collar screen was introduced as a new variable.

Unfortunately, after only two runs of this experiment this particular phase had to be terminated because of the excessive tearing of the cloth screens. From the scanty results obtained it appeared, however, that with no horizontal screen high solids removal could be achieved with a higher portion of the flow retreated. It was therefore decided to repeat these runs with screens made of stainless steel instead of cloth.

Phase c

A third experiment, using stainless steel collar screens of two mesh sizes, similar to that attempted in phase b, was performed with the same collar screen mesh size, collar screen speed (rpm), and flow rate (gpm) used before.

In this phase with a stainless steel collar screen, high removal rates y were possible for eight sets of conditions for the factors just mentioned. However, these high y values were obtained with retreated flow z at undesirably high values (before, they had been too low). The object was to get reasonably small values for z , but not so small as to make shoveling necessary; values between 5% and 20% were desirable. It was believed that by varying flow rate and speed of rotation of the collar screen, this objective could be achieved without sacrificing solids removal.

Phase d

Again, using a stainless steel collar screen, another experiment, with two factors—collar screen speed (rpm) and flow rate (gpm)—set at two levels each, was run. This time, high values of solids removal were maintained, but unfortunately flow retreated values were even higher than before.

Phase e

It was now conjectured that intermittent back washing could overcome the difficulties. This procedure was now introduced with influent flow rate and collar screen mesh varied.

The results of this experiment lead to a removal efficiency of 89% with a retreated flow of only 8%. This was regarded as a satisfactory and practical solution, and the investigation was terminated at that point.

For detailed analysis of this experiment, the reader should refer to Box, Hunter, and Hunter (1978), p. 354. Of course, these types of experiments and their analyses are discussed in this text (see Chapter 18).

1.2 A SURVEY

The purpose of a sample survey is to make inferences about certain characteristics of a population from which samples are drawn. The inferences to be made for a population usually entails the estimation of population parameters, such as the population total, the mean, or the population proportion of a certain

characteristic of interest. In any sample survey a clear statement of its objective is very important. Without a clear statement about the objectives, it is very easy to miss pertinent information while planning the survey that can cause difficulties at the end of the study.

In any sample survey only relevant information should be collected. Sometimes trying to collect too much information may become very confusing and consequently hinder the determination of the final goal. Moreover, collecting information in sample surveys costs money, so that the interested party must determine which and how much information should be obtained. For example, it is important to describe how much precision in the final results is desired. Too little information may prevent obtaining good estimates with desired precision, while too much information may not be needed and may unnecessarily cost too much money. One way to avoid such problems is to select an appropriate method of sampling the population. In other words, the sample survey needs to be appropriately designed. A brief discussion of such designs is given in Chapter 2. For more details on these designs, the reader may refer to Cochran (1977), Sukhatme et al. (1970), or Schaeffer et al. (2006).

1.3 AN OBSERVATIONAL STUDY

An observational study is one that does not involve any experimental studies. Consequently observational studies do not control any variables. For example, a realtor wishes to appraise a house value. All the data used for this purpose is observational data. Many psychiatric studies involve observational data.

Frequently, in fitting a regression model (see Chapters 15 and 16), we use observational data. Similarly, in quality control (see Chapters 11 and 12), most of the data used in studying control charts for attributes is observational data. Note that control charts for attributes usually do not provide any cause-and-effect relationships. This is because observational data give us very limited information about cause-and-effect relationships.

As another example, many psychiatric studies involve observational data, and such data do not provide the cause of patient's psychiatric problems. An advantage of observational studies is that they are usually more cost effective than experimental studies. The disadvantage of observational studies is that the data may not be as informative as experimental data.

1.4 A SET OF HISTORICAL DATA

Historical data are not collected by the experimenter. The data are made available to him/her.

Many fields of study such as the many branches of business studies, use historical data. A financial advisor for planning purposes uses sets of historical data. Many investment services provide financial data on a company-by-company basis.

1.5 A BRIEF DESCRIPTION OF WHAT IS COVERED IN THIS BOOK

Data collection is very important since it can greatly influence the final outcome of subsequent data analyses. After collection of the data, it is important to organize, summarize, present the preliminary outcomes, and interpret them. Various types of tables and graphs that summarize the data are presented in Chapter 2. Also in that chapter we give some methods used to determine certain quantities, called *statistics*, which are used to summarize some of the key properties of the data.

The basic principles of probability are necessary to study various probability distributions. We present the basic principles of elementary probability theory in Chapter 3. Probability distributions are fundamental in the development of the various techniques of statistical inference. The concept of random variables is also discussed in Chapter 3.

Chapters 4 and 5 are devoted to some of the important discrete distributions, continuous distributions, and their moment-generating functions. In addition we study in Chapter 5 some special distributions that are used in reliability theory.

In Chapter 6 we study joint distributions of two or more discrete and continuous random variables and their moment-generating functions. Included in Chapter 6 is the study of the bivariate normal distribution.

Chapter 7 is devoted to the probability distributions of some sample statistics, such as the sample mean, sample proportions, and sample variance. In this chapter we also study a fundamental result of probability theory, known as the Central Limit Theorem. This theorem can be used to approximate the probability distribution of the sample mean when the sample size is large. In this chapter we also study some sampling distributions of some sample statistics for the special case in which the population distribution is the so-called normal distribution. In addition we present probability distributions of various "order statistics," such as the largest element in a sample, smallest element in a sample, and sample median.

Chapter 8 discusses the use of sample data for estimating the unknown population parameters of interest, such as the population mean, population variance, and population proportion. In Chapter 8 also discusses the methods of estimating the difference of two population means, the difference of two population proportions, and the ratio of two population variances and standard deviations. Two types of estimators are included, namely point estimators and interval estimators (confidence intervals).

Chapter 9 deals with the important statistical tests of hypotheses that are concerned with the population means, population variance, and population proportion for one and two populations. Methods of testing hypotheses using the confidence intervals studied in Chapter 8 are also presented.

Chapter 10 gives an introduction to the theory of reliability. Methods of estimation and hypothesis testing using the exponential and Weibull distributions are presented.

Chapters 11 and 12 are devoted to control charts for variables and attributes used in phase I and phase II of a process. "Phase I" refers to the initial stage of a new process, and "phase II" refers to a matured process. Control charts are used to determine whether a process involving manufacturing or service is "under statistical control" on the basis of information contained in a sequence of small samples of items of interest.

Chapter 13 is concerned with the chi-square goodness-of-fit test, which is used to test whether a set of sample data support the hypothesis that the sampled population follows some specified probability model. In addition we apply the chi-square goodness-of-fit test for testing hypotheses of independence and homogeneity. These tests involve methods of comparing observed frequencies with those that are expected if a certain hypothesis is true.

Chapter 14 gives a brief look at tests known as "nonparametric tests," which are used when the assumption about the underlying distribution having some specified parametric form cannot be made.

Chapter 15 introduces an important topic of applied statistics: simple linear regression analysis. Linear regression analysis is frequently used by engineers, social scientists, health researchers, and biological scientists. This statistical technique explores the relation between two variables so that one variable can be predicted from the other. In this chapter we discuss the least squares method for estimating the simple regression model, called the fitting of this regression model. Also we discuss how to perform a residual analysis, which is used to check the adequacy of the regression model, and study certain transformations that are used when the model is not adequate.

Chapter 16 extends the results of Chapter 15 to multiple linear regression. Like the simple linear regression model, multiple linear regression analysis is widely used. It provides statistical techniques that explore the relations among more than two variables, so that one variable can be predicted from the use of the other variables. In this chapter we give a discussion of multiple linear regression, including the matrix approach. Finally, a brief discussion of logistic regression is given.

In Chapter 17 we introduce the design and analysis of experiments using one, two, or more factors. Designs for eliminating the effects of one or two nuisance variables along with a method of estimating one or more missing observations are given. We include two nonparametric tests, the Kruskal-Wallis and the Friedman test, for analyzing one-way and randomized complete block designs. Finally, models with fixed effects, mixed effects, and random effects are also discussed.

Chapter 18 introduces a special class of designs, the so-called 2^k factorial designs. These designs are widely used in various industrial and scientific applications. An extensive discussion of unreplicated 2^k factorial designs, blocking of 2^k factorial designs, confounding in the 2^k factorial designs, and Yates's

algorithm for the 2^k factorial designs is also included. We also devote a section to fractional factorial designs, discussing one-half and one-quarter replications of 2^k factorial designs.

In Chapter 19 we introduce the topic of response surface methodology (RSM). First-order and second-order designs used in response surface methodology are discussed. Methods of determining optimum or near optimum points using the “method of steepest ascent” and the analysis of a fitted second-order response surface are also presented. Due to lack of space, this chapter is not included in the text but is available for download on the book website: www.wiley.com/go/statsforengineers.

All chapters are supported by three popular statistical software packages, MINITAB, Microsoft Excel, and JMP. The MINITAB and Microsoft Excel are fully integrated into the text of each chapter, whereas JMP is given in an independent section, which is not included in the text but is available for download on the book website: www.wiley.com/go/statsforengineers. Frequently we use the same examples for the discussion of JMP as are used in the discussion of MINITAB and Microsoft Excel. For the use of each of these software packages, no prior knowledge is assumed, since we give each step, from entering the data to the final analysis of such data under investigation. Finally, a section of case studies is included in almost all the chapters.