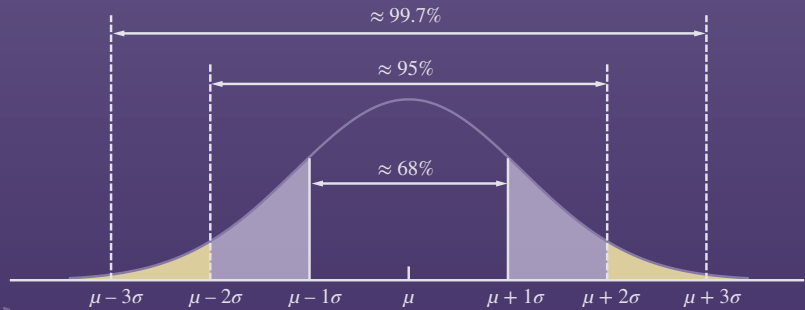


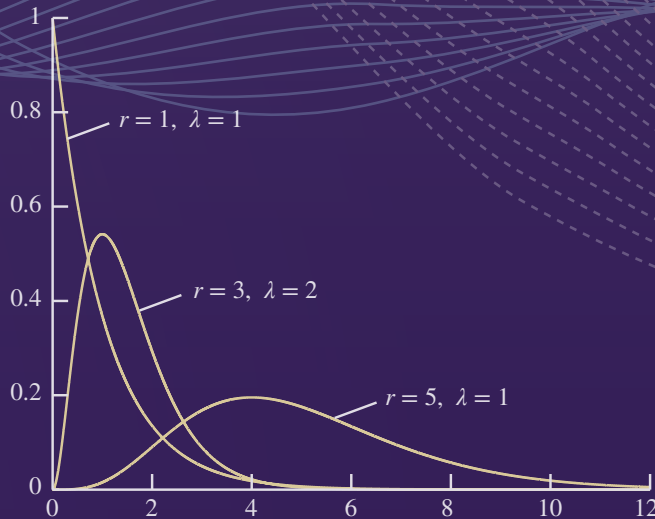
This International Student Edition is for use outside of the U.S.



Second Edition

Principles of Statistics for Engineers and Scientists

William Navidi



**Mc
Graw
Hill**

Principles of Statistics for Engineers and Scientists

Second Edition

William Navidi





PRINCIPLES OF STATISTICS FOR ENGINEERS AND SCIENTISTS

Published by McGraw-Hill Education, 2 Penn Plaza, New York, NY 10121. Copyright © 2021 by McGraw-Hill Education. All rights reserved. Printed in the United States of America. No part of this publication may be reproduced or distributed in any form or by any means, or stored in a database or retrieval system, without the prior written consent of McGraw-Hill Education, including, but not limited to, in any network or other electronic storage or transmission, or broadcast for distance learning.

Some ancillaries, including electronic and print components, may not be available to customers outside the United States.

This book is printed on acid-free paper.

1 2 3 4 5 6 7 8 9 LCR 24 23 22 21 20

ISBN 978-1-260-57073-1

MHID 1-260-57073-8

Cover Image: *McGraw-Hill Education*

All credits appearing on page or at the end of the book are considered to be an extension of the copyright page.

The Internet addresses listed in the text were accurate at the time of publication. The inclusion of a website does not indicate an endorsement by the authors or McGraw-Hill Education, and McGraw-Hill Education does not guarantee the accuracy of the information presented at these sites.

mheducation.com/highered

To Catherine, Sarah, and Thomas

ABOUT THE AUTHOR

William Navidi is Professor of Mathematical and Computer Sciences at the Colorado School of Mines. He received the B.A. degree in mathematics from New College, the M.A. in mathematics from Michigan State University, and the Ph.D. in statistics from the University of California at Berkeley. Professor Navidi has authored more than 80 research papers both in statistical theory and in a wide variety of applications including computer networks, epidemiology, molecular biology, chemical engineering, and geophysics.

CONTENTS

Preface vii

Chapter 1

Summarizing Univariate Data 1

Introduction 1

- 1.1 Sampling 3
- 1.2 Summary Statistics 11
- 1.3 Graphical Summaries 21

Chapter 2

Summarizing Bivariate Data 37

Introduction 37

- 2.1 The Correlation Coefficient 37
- 2.2 The Least-Squares Line 49
- 2.3 Features and Limitations of the Least-Squares Line 57

Chapter 3

Probability 67

Introduction 67

- 3.1 Basic Ideas 67
- 3.2 Conditional Probability and Independence 75
- 3.3 Random Variables 86
- 3.4 Functions of Random Variables 107

Chapter 4

Commonly Used Distributions 122

Introduction 122

- 4.1 The Binomial Distribution 122
- 4.2 The Poisson Distribution 130
- 4.3 The Normal Distribution 137

4.4 The Lognormal Distribution 148

4.5 The Exponential Distribution 151

4.6 Some Other Continuous Distributions 156

4.7 Probability Plots 162

4.8 The Central Limit Theorem 166

Chapter 5

Point and Interval Estimation for a Single Sample 179

Introduction 179

- 5.1 Point Estimation 180
- 5.2 Large-Sample Confidence Intervals for a Population Mean 183
- 5.3 Confidence Intervals for Proportions 196
- 5.4 Small-Sample Confidence Intervals for a Population Mean 202
- 5.5 Prediction Intervals and Tolerance Intervals 211

Chapter 6

Hypothesis Tests for a Single Sample 219

Introduction 219

- 6.1 Large-Sample Tests for a Population Mean 219
- 6.2 Drawing Conclusions from the Results of Hypothesis Tests 229
- 6.3 Tests for a Population Proportion 237
- 6.4 Small-Sample Tests for a Population Mean 242

- 6.5 The Chi-Square Test 248
- 6.6 Fixed-Level Testing 257
- 6.7 Power 262
- 6.8 Multiple Tests 271

Chapter 7

Inferences for Two Samples 278

Introduction 278

- 7.1 Large-Sample Inferences on the Difference Between Two Population Means 278
- 7.2 Inferences on the Difference Between Two Proportions 287
- 7.3 Small-Sample Inferences on the Difference Between Two Means 295
- 7.4 Inferences Using Paired Data 305
- 7.5 Tests for Variances of Normal Populations 315

Chapter 8

Inference in Linear Models 325

Introduction 325

- 8.1 Inferences Using the Least-Squares Coefficients 326
- 8.2 Checking Assumptions 349
- 8.3 Multiple Regression 360
- 8.4 Model Selection 377

Chapter 9

Factorial Experiments 411

Introduction 411

- 9.1 One-Factor Experiments 411
- 9.2 Pairwise Comparisons in One-Factor Experiments 430
- 9.3 Two-Factor Experiments 436
- 9.4 Randomized Complete Block Designs 456
- 9.5 2^p Factorial Experiments 463

Chapter 10

Statistical Quality Control 492

Introduction 492

- 10.1 Basic Ideas 492
- 10.2 Control Charts for Variables 495
- 10.3 Control Charts for Attributes 514
- 10.4 The CUSUM Chart 519
- 10.5 Process Capability 522

Appendix A: Tables 529

Appendix B: Bibliography 552

Answers to Selected Exercises 555

Index 601

PREFACE

MOTIVATION

This book is based on the author's more comprehensive text *Statistics for Engineers and Scientists*, 5th edition (McGraw-Hill, 2020), which is used for both one- and two-semester courses. The key concepts from that book form the basis for this text, which is designed for a one-semester course. The emphasis is on statistical methods and how they can be applied to problems in science and engineering, rather than on theory. While the fundamental principles of statistics are common to all disciplines, students in science and engineering learn best from examples that present important ideas in realistic settings. Accordingly, the book contains many examples that feature real, contemporary data sets, both to motivate students and to show connections to industry and scientific research. As the text emphasizes applications rather than theory, the mathematical level is appropriately modest. Most of the book will be mathematically accessible to those whose background includes one semester of calculus.

COMPUTER USE

Over the past 40 years, the development of fast and cheap computing has revolutionized statistical practice; indeed, this is one of the main reasons that statistical methods have been penetrating ever more deeply into scientific work. Scientists and engineers today must not only be adept with computer software packages; they must also have the skill to draw conclusions from computer output and to state those conclusions in words. Accordingly, the book contains exercises and examples that involve interpreting, as well as generating, computer output, especially in the chapters on linear models and factorial experiments. Many instructors integrate the use of statistical software into their courses; this book may be used effectively with any package.

CONTENT

Chapter 1 covers sampling and descriptive statistics. The reason that statistical methods work is that samples, when properly drawn, are likely to resemble their populations. Therefore, Chapter 1 begins by describing some ways to draw valid samples. The second part of the chapter discusses descriptive statistics for univariate data.

Chapter 2 presents descriptive statistics for bivariate data. The correlation coefficient and least-squares line are discussed. The discussion emphasizes that linear models are appropriate only when the relationship between the variables is linear, and it describes the effects of outliers and influential points. Placing this chapter early enables instructors to present some coverage of these topics in courses where there is not enough time for a full treatment from an inferential point of view. Alternatively, this chapter may be postponed and covered just before the inferential procedures for linear models in Chapter 8.

Chapter 3 is about probability. The goal here is to present the essential ideas without a lot of mathematical derivations. I have attempted to illustrate each result with an example or two, in a scientific context where possible, to present the intuition behind the result.

Chapter 4 presents many of the probability distribution functions commonly used in practice. Probability plots and the Central Limit Theorem are also covered. Only the normal and binomial distribution are used extensively in the remainder of the text; instructors may choose which of the other distributions to cover.

Chapters 5 and 6 cover one-sample methods for confidence intervals and hypothesis testing, respectively. Point estimation is covered as well, in Chapter 5. The P -value approach to hypothesis testing is emphasized, but fixed-level testing and power calculations are also covered. A discussion of the multiple testing problem is also presented.

Chapter 7 presents two-sample methods for confidence intervals and hypothesis testing. There is often not enough time to cover as many of these methods as one would like; instructors who are pressed for time may choose which of the methods they wish to cover.

Chapter 8 covers inferential methods in linear regression. In practice, scatterplots often exhibit curvature or contain influential points. Therefore, this chapter includes material on checking model assumptions and transforming variables. In the coverage of multiple regression, model selection methods are given particular emphasis, because choosing the variables to include in a model is an essential step in many real-life analyses.

Chapter 9 discusses some commonly used experimental designs and the methods by which their data are analyzed. One-way and two-way analysis of variance methods, along with randomized complete block designs and 2^p factorial designs, are covered fairly extensively.

Chapter 10 presents the topic of statistical quality control, covering control charts, CUSUM charts, and process capability, and concluding with a brief discussion of sixsigma quality.

RECOMMENDED COVERAGE

The book contains enough material for a one-semester course meeting four hours per week. For a three-hour course, it will probably be necessary to make some choices about coverage. One option is to cover the first three chapters, going lightly over the last two sections of Chapter 3, then cover the binomial, Poisson, and normal distributions in Chapter 4, along with the Central Limit Theorem. One can then cover the confidence intervals and hypothesis tests in Chapters 5 and 6, and finish either with the two-sample procedures in Chapter 7 or by covering as much of the material on inferential methods in regression in Chapter 8 as time permits.

For a course that puts more emphasis on regression and factorial experiments, one can go quickly over the power calculations and multiple testing procedures, and cover Chapters 8 and 9 immediately following Chapter 6. Alternatively, one could substitute Chapter 10 on statistical quality control for Chapter 9.

NEW FOR THIS EDITION

The second edition of this book is intended to extend the strengths of the first. Some of the changes are:

- More than 250 new problems have been included.
- Many examples have been updated.
- Material on resistance to outliers has been added to Chapter 1.
- Material on interpreting the slope of the least-squares line has been added to Chapter 2.
- Material on the F -test for variance has been added to Chapter 7.
- The exposition has been improved in a number of places.

ACKNOWLEDGMENTS

I am indebted to many people for contributions at every stage of development. I received many valuable suggestions from my colleagues Gus Greivel, Ashlyn Munson, and Melissa Laeser at the Colorado School of Mines. I am particularly grateful to Jack Miller of The University of Michigan, who found many errors and made many valuable suggestions for improvement.

The staff at McGraw-Hill has been extremely capable and supportive. In particular, I would like to express thanks to Product Developer Tina Bower, Content Project Manager Jeni McAtee and Senior Project Manager Sarita Yadav for their patience and guidance in the preparation of this edition.

William Navidi

Affordability & Outcomes = Academic Freedom!

You deserve choice, flexibility, and control. You know what's best for your students and selecting the course materials that will help them succeed should be in your hands.

That's why providing you with a wide range of options that lower costs and drive better outcomes is our highest priority.



connect®

Students—study more efficiently, retain more, and achieve better outcomes. Instructors—focus on what you love—teaching.



Laptop: McGraw-Hill Education

They'll thank you for it.

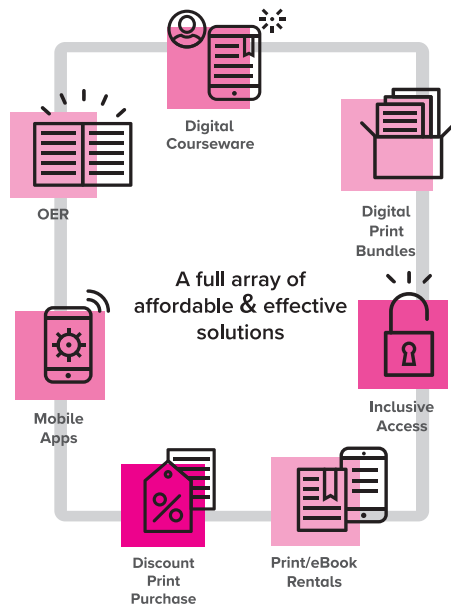
Study resources in Connect help your students be better prepared in less time. You can transform your class time from dull definitions to dynamic discussion. Hear from your peers about the benefits of Connect at www.mheducation.com/highered/connect/smartbook

Make it simple, make it affordable.

Connect makes it easy with seamless integration using any of the major Learning Management Systems—Blackboard®, Canvas, and D2L, among others—to let you organize your course in one convenient location. Give your students access to digital materials at a discount with our inclusive access program. Ask your McGraw-Hill representative for more information.

Learning for everyone.

McGraw-Hill works directly with Accessibility Services Departments and faculty to meet the learning needs of all students. Please contact your Accessibility Services office and ask them to email accessibility@mheducation.com, or visit www.mheducation.com/about/accessibility.html for more information.



Learn more at: www.mheducation.com/realvalue



Rent It

Affordable print and digital rental options through our partnerships with leading textbook distributors including Amazon, Barnes & Noble, Chegg, Follett, and more.



Go Digital

A full and flexible range of affordable digital solutions ranging from Connect, ALEKS, inclusive access, mobile apps, OER and more.



Get Print

Students who purchase digital materials can get a loose-leaf print version at a significantly reduced rate to meet their individual preferences and budget.

Summarizing Univariate Data

Introduction

Advances in science and engineering occur in large part through the collection and analysis of data. Proper analysis of data is challenging, because scientific data are subject to random variation. That is, when scientific measurements are repeated, they come out somewhat differently each time. This poses a problem: How can one draw conclusions from the results of an experiment when those results could have come out differently? To address this question, a knowledge of statistics is essential. The methods of statistics allow scientists and engineers to design valid experiments and to draw reliable conclusions from the data they produce.

While our emphasis in this book is on the applications of statistics to science and engineering, it is worth mentioning that the analysis and interpretation of data are playing an ever-increasing role in all aspects of modern life. For better or worse, huge amounts of data are collected about our opinions and our lifestyles, for purposes ranging from the creation of more effective marketing campaigns to the development of social policies designed to improve our way of life. On almost any given day, newspaper articles are published that purport to explain social or economic trends through the analysis of data. A basic knowledge of statistics is therefore necessary not only to be an effective scientist or engineer, but also to be a well-informed member of society.

The Basic Idea

The basic idea behind all statistical methods of data analysis is to make inferences about a population by studying a relatively small sample chosen from it. As an illustration, consider a machine that makes steel balls for ball bearings used in clutch systems. The specification for the diameter of the balls is 0.65 ± 0.03 cm. During the last hour, the machine has made 2000 balls. The quality engineer wants to know approximately

how many of these balls meet the specification. He does not have time to measure all 2000 balls. So he draws a random sample of 80 balls, measures them, and finds that 72 of them (90%) meet the diameter specification. Now, it is unlikely that the sample of 80 balls represents the population of 2000 perfectly. The proportion of good balls in the population is likely to differ somewhat from the sample proportion of 90%. What the engineer needs to know is just how large that difference is likely to be. For example, is it plausible that the population percentage could be as high as 95%? 98%? As low as 85%? 80%?

Here are some specific questions that the engineer might need to answer on the basis of these sample data:

1. The engineer needs to compute a rough estimate of the likely size of the difference between the sample proportion and the population proportion. How large is a typical difference for this kind of sample?
2. The quality engineer needs to note in a logbook the percentage of acceptable balls manufactured in the last hour. Having observed that 90% of the sample balls were good, he will indicate the percentage of acceptable balls in the population as an interval of the form $90\% \pm x\%$, where x is a number calculated to provide reasonable certainty that the true population percentage is in the interval. How should x be calculated?
3. The engineer wants to be fairly certain that the percentage of good balls is at least 85%; otherwise, he will shut down the process for recalibration. How certain can he be that at least 85% of the 1000 balls are good?

Much of this book is devoted to addressing questions like these. The first of these questions requires the computation of a **standard deviation**, which we will discuss in Chapter 3. The second question requires the construction of a **confidence interval**, which we will learn about in Chapter 5. The third calls for a **hypothesis test**, which we will study in Chapter 6.

The remaining chapters in the book cover other important topics. For example, the engineer in our example may want to know how the amount of carbon in the steel balls is related to their compressive strength. Issues like this can be addressed with the methods of **correlation** and **regression**, which are covered in Chapters 2 and 8. It may also be important to determine how to adjust the manufacturing process with regard to several factors, in order to produce optimal results. This requires the design of **factorial experiments**, which are discussed in Chapter 9. Finally, the engineer will need to develop a plan for monitoring the quality of the product manufactured by the process. Chapter 10 covers the topic of **statistical quality control**, in which statistical methods are used to maintain quality in an industrial setting.

The topics listed here concern methods of drawing conclusions from data. These methods form the field of **inferential statistics**. Before we discuss these topics, we must first learn more about methods of collecting data and of summarizing clearly the basic information they contain. These are the topics of **sampling** and **descriptive statistics**, and they are covered in the rest of this chapter.

1.1 Sampling

As mentioned, statistical methods are based on the idea of analyzing a **sample** drawn from a **population**. For this idea to work, the sample must be chosen in an appropriate way. For example, let us say that we wished to study the heights of students at the Colorado School of Mines by measuring a sample of 100 students. How should we choose the 100 students to measure? Some methods are obviously bad. For example, choosing the students from the rosters of the football and basketball teams would undoubtedly result in a sample that would fail to represent the height distribution of the population of students. You might think that it would be reasonable to use some conveniently obtained sample, for example, all students living in a certain dorm or all students enrolled in engineering statistics. After all, there is no reason to think that the heights of these students would tend to differ from the heights of students in general. Samples like this are not ideal, however, because they can turn out to be misleading in ways that are not anticipated. The best sampling methods involve **random sampling**. There are many different random sampling methods, the most basic of which is **simple random sampling**.

Simple Random Samples

To understand the nature of a simple random sample, think of a lottery. Imagine that 10,000 lottery tickets have been sold and that 5 winners are to be chosen. What is the fairest way to choose the winners? The fairest way is to put the 10,000 tickets in a drum, mix them thoroughly, and then reach in and one by one draw 5 tickets out. These 5 winning tickets are a simple random sample from the population of 10,000 lottery tickets. Each ticket is equally likely to be one of the 5 tickets drawn. More important, each collection of 5 tickets that can be formed from the 10,000 is equally likely to make up the group of 5 that is drawn. It is this idea that forms the basis for the definition of a simple random sample.

Summary

- A **population** is the entire collection of objects or outcomes about which information is sought.
- A **sample** is a subset of a population, containing the objects or outcomes that are actually observed.
- A **simple random sample** of size n is a sample chosen by a method in which each collection of n population items is equally likely to make up the sample, just as in a lottery.

Since a simple random sample is analogous to a lottery, it can often be drawn by the same method now used in many lotteries: with a computer random number generator. Suppose there are N items in the population. One assigns to each item in the population an integer between 1 and N . Then one generates a list of random integers between

1 and N and chooses the corresponding population items to make up the simple random sample.

Example

1.1

A utility company wants to conduct a survey to measure the satisfaction level of its customers in a certain town. There are 10,000 customers in the town, and utility employees want to draw a sample of size 200 to interview over the telephone. They obtain a list of all 10,000 customers, and number them from 1 to 10,000. They use a computer random number generator to generate 200 random integers between 1 and 10,000 and then telephone the customers who correspond to those numbers. Is this a simple random sample?

Solution

Yes, this is a simple random sample. Note that it is analogous to a lottery in which each customer has a ticket and 200 tickets are drawn.

Example

1.2

A quality engineer wants to inspect electronic microcircuits in order to obtain information on the proportion that are defective. She decides to draw a sample of 100 circuits from a day's production. Each hour for 5 hours, she takes the 20 most recently produced circuits and tests them. Is this a simple random sample?

Solution

No. Not every subset of 100 circuits is equally likely to make up the sample. To construct a simple random sample, the engineer would need to assign a number to each circuit produced during the day and then generate random numbers to determine which circuits make up the sample.

Samples of Convenience

In some cases, it is difficult or impossible to draw a sample in a truly random way. In these cases, the best one can do is to sample items by some convenient method. For example, imagine that a construction engineer has just received a shipment of 1000 concrete blocks, each weighing approximately 50 pounds. The blocks have been delivered in a large pile. The engineer wishes to investigate the crushing strength of the blocks by measuring the strengths in a sample of 10 blocks. To draw a simple random sample would require removing blocks from the center and bottom of the pile, which might be quite difficult. For this reason, the engineer might construct a sample simply by taking 10 blocks off the top of the pile. A sample like this is called a **sample of convenience**.

Definition

A **sample of convenience** is a sample that is obtained in some convenient way, and not drawn by a well-defined random method.

The big problem with samples of convenience is that they may differ systematically in some way from the population. For this reason samples of convenience should only be

used in situations where it is not feasible to draw a random sample. When it is necessary to take a sample of convenience, it is important to think carefully about all the ways in which the sample might differ systematically from the population. If it is reasonable to believe that no important systematic difference exists, then it may be acceptable to treat the sample of convenience as if it were a simple random sample. With regard to the concrete blocks, if the engineer is confident that the blocks on the top of the pile do not differ systematically in any important way from the rest, then he may treat the sample of convenience as a simple random sample. If, however, it is possible that blocks in different parts of the pile may have been made from different batches of mix or may have different curing times or temperatures, a sample of convenience could give misleading results.

Some people think that a simple random sample is guaranteed to reflect its population perfectly. This is not true. Simple random samples always differ from their populations in some ways, and occasionally they may be substantially different. Two different samples from the same population will differ from each other as well. This phenomenon is known as **sampling variation**. Sampling variation is one of the reasons that scientific experiments produce somewhat different results when repeated, even when the conditions appear to be identical. For example, suppose that a quality inspector draws a simple random sample of 40 bolts from a large shipment, measures the length of each, and finds that 32 of them, or 80%, meet a length specification. Another inspector draws a different sample of 40 bolts and finds that 36 of them, or 90%, meet the specification. By chance, the second inspector got a few more good bolts in her sample. It is likely that neither sample reflects the population perfectly. The proportion of good bolts in the population is likely to be close to 80% or 90%, but it is not likely that it is exactly equal to either value.

Since simple random samples don't reflect their populations perfectly, why is it important that sampling be done at random? The benefit of a simple random sample is that there is no systematic mechanism tending to make the sample unrepresentative. The differences between the sample and its population are due entirely to random variation. Since the mathematical theory of random variation is well understood, we can use mathematical models to study the relationship between simple random samples and their populations. For a sample not chosen at random, there is generally no theory available to describe the mechanisms that caused the sample to differ from its population. Therefore, nonrandom samples are often difficult to analyze reliably.

Tangible and Conceptual Populations

The populations discussed so far have consisted of actual physical objects—the customers of a utility company, the concrete blocks in a pile, the bolts in a shipment. Such populations are called **tangible populations**. Tangible populations are always finite. After an item is sampled, the population size decreases by 1. In principle, one could in some cases return the sampled item to the population, with a chance to sample it again, but this is rarely done in practice.

Engineering data are often produced by measurements made in the course of a scientific experiment, rather than by sampling from a tangible population. To take a simple

example, imagine that an engineer measures the length of a rod five times, being as careful as possible to take the measurements under identical conditions. No matter how carefully the measurements are made, they will differ somewhat from one another, because of variation in the measurement process that cannot be controlled or predicted. It turns out that it is often appropriate to consider data like these to be a simple random sample from a population. The population, in these cases, consists of all the values that might possibly have been observed. Such a population is called a **conceptual population**, since it does not consist of actual objects.

Definition

A simple random sample may consist of values obtained from a process under identical experimental conditions. In this case, the sample comes from a population that consists of all the values that might possibly have been observed. Such a population is called a **conceptual population**.

Example 1.3 involves a conceptual population.

Example

1.3

A geologist weighs a rock several times on a sensitive scale. Each time, the scale gives a slightly different reading. Under what conditions can these readings be thought of as a simple random sample? What is the population?

Solution

If the physical characteristics of the scale remain the same for each weighing, so that the measurements are made under identical conditions, then the readings may be considered to be a simple random sample. The population is conceptual. It consists of all the readings that the scale could in principle produce.

Determining Whether a Sample Is a Simple Random Sample

We saw in Example 1.3 that it is the physical characteristics of the measurement process that determine whether the data are a simple random sample. In general, when deciding whether a set of data may be considered to be a simple random sample, it is necessary to have some understanding of the process that generated the data. Statistical methods can sometimes help, especially when the sample is large, but knowledge of the mechanism that produced the data is more important.

Example

1.4

A new chemical process has been designed that is supposed to produce a higher yield of a certain chemical than does an old process. To study the yield of this process, we run it 50 times and record the 50 yields. Under what conditions might it be reasonable to treat this as a simple random sample? Describe some conditions under which it might not be appropriate to treat this as a simple random sample.

Solution

To answer this, we must first specify the population. The population is conceptual and consists of the set of all yields that will result from this process as many times as it will ever be run. What we have done is to sample the first 50 yields of the process. *If, and only if,* we are confident that the first 50 yields are generated under identical conditions and that they do not differ in any systematic way from the yields of future runs, then we may treat them as a simple random sample.

Be cautious, however. There are many conditions under which the 50 yields could fail to be a simple random sample. For example, with chemical processes, it is sometimes the case that runs with higher yields tend to be followed by runs with lower yields, and vice versa. Sometimes yields tend to increase over time, as process engineers learn from experience how to run the process more efficiently. In these cases, the yields are not being generated under identical conditions and would not be a simple random sample.

Example 1.4 shows once again that a good knowledge of the nature of the process under consideration is important in deciding whether data may be considered to be a simple random sample. Statistical methods can sometimes be used to show that a given data set is *not* a simple random sample. For example, sometimes experimental conditions gradually change over time. A simple but effective method to detect this condition is to plot the observations in the order they were taken. A simple random sample should show no obvious pattern or trend.

Figure 1.1 presents plots of three samples in the order they were taken. The plot in Figure 1.1a shows an oscillatory pattern. The plot in Figure 1.1b shows an increasing trend. Neither of these samples should be treated as a simple random sample. The plot in Figure 1.1c does not appear to show any obvious pattern or trend. It might be appropriate to treat these data as a simple random sample. However, before making that decision, it

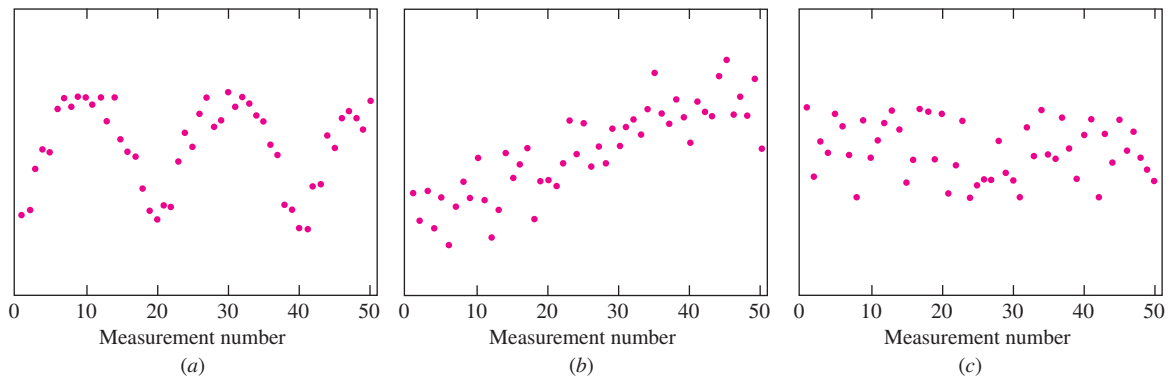


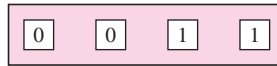
FIGURE 1.1 Three plots of observed values versus the order in which they were made. (a) The values show a definite pattern over time. This is not a simple random sample. (b) The values show a trend over time. This is not a simple random sample. (c) The values do not show a pattern or trend. It may be appropriate to treat these data as a simple random sample.

is still important to think about the process that produced the data, since there may be concerns that don't show up in the plot.

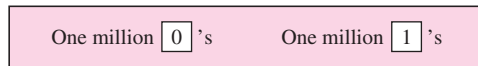
Independence

The items in a sample are said to be **independent** if knowing the values of some of them does not help to predict the values of the others. With a finite, tangible population, the items in a simple random sample are not strictly independent, because as each item is drawn, the population changes. This change can be substantial when the population is small. However, when the population is very large, this change is negligible and the items can be treated as if they were independent.

To illustrate this idea, imagine that we draw a simple random sample of 2 items from the population



For the first draw, the numbers 0 and 1 are equally likely. But the value of the second item is clearly influenced by the first; if the first is 0, the second is more likely to be 1, and vice versa. Thus, the sampled items are dependent. Now assume we draw a sample of size 2 from this population:



Again on the first draw, the numbers 0 and 1 are equally likely. But unlike the previous example, these two values remain almost equally likely the second draw as well, no matter what happens on the first draw. With the large population, the sample items are for all practical purposes independent.

It is reasonable to wonder how large a population must be in order that the items in a simple random sample may be treated as independent. A rule of thumb is that when sampling from a finite population, the items may be treated as independent so long as the sample contains 5% or less of the population.

Interestingly, it is possible to make a population behave as though it were infinitely large, by replacing each item after it is sampled. This method is called **sampling with replacement**. With this method, the population is exactly the same on every draw, and the sampled items are truly independent.

With a conceptual population, we require that the sample items be produced under identical experimental conditions. In particular, then, no sample value may influence the conditions under which the others are produced. Therefore, the items in a simple random sample from a conceptual population may be treated as independent. We may think of a conceptual population as being infinite or, equivalently, that the items are sampled with replacement.

Summary

- The items in a sample are **independent** if knowing the values of some of the items does not help to predict the values of the others.
- Items in a simple random sample may be treated as independent in many cases encountered in practice. The exception occurs when the population is finite and the sample consists of a substantial fraction (more than 5%) of the population.

Other Sampling Methods

In addition to simple random sampling there are other sampling methods that are useful in various situations. In **weighted sampling**, some items are given a greater chance of being selected than others, like a lottery in which some people have more tickets than others. In **stratified random sampling**, the population is divided up into subpopulations, called **strata**, and a simple random sample is drawn from each stratum. In **cluster sampling**, items are drawn from the population in groups, or clusters. Cluster sampling is useful when the population is too large and spread out for simple random sampling to be feasible. For example, many U.S. government agencies use cluster sampling to sample the U.S. population to measure sociological factors such as income and unemployment. A good source of information on sampling methods is Cochran (1977).

Simple random sampling is not the only valid method of sampling. But it is the most fundamental, and we will focus most of our attention on this method. From now on, unless otherwise stated, the terms “sample” and “random sample” will be taken to mean “simple random sample.”

Types of Data

When a numerical quantity designating how much or how many is assigned to each item in a sample, the resulting set of values is called **numerical** or **quantitative**. In some cases, sample items are placed into categories, and category names are assigned to the sample items. Then the data are **categorical** or **qualitative**. Sometimes both quantitative and categorical data are collected in the same experiment. For example, in a loading test of column-to-beam welded connections, data may be collected both on the torque applied at failure and on the location of the failure (weld or beam). The torque is a quantitative variable, and the location is a categorical variable.

Controlled Experiments and Observational Studies

Many scientific experiments are designed to determine the effect of changing one or more factors on the value of a response. For example, suppose that a chemical engineer wants to determine how the concentrations of reagent and catalyst affect the yield of a process. The engineer can run the process several times, changing the concentrations each time, and compare the yields that result. This sort of experiment is called