

Chapter 1

Statistics Basics

Objectives

What does the word *statistics* bring to mind?

To most people, it suggests numerical facts or data, such as unemployment figures, farm prices, or the number of marriages and divorces.

6.8

3.2

Two common definitions of the word *statistics* are as follows:

1. facts or data, either numerical or nonnumerical, organized and summarized so as to provide useful and accessible information about a particular subject.
2. the science of organizing and summarizing numerical or nonnumerical information.

Definition 1.1

Descriptive Statistics

Descriptive Statistics consists of methods for organizing and summarizing information.

Descriptive statistics includes the construction of graphs, charts, and tables and the calculation of various descriptive measures such as averages, measures of variation, and percentiles.

Example 1.1

The 1948 Baseball Season:

In 1948, the Washington Senators played 153 games, winning 56 and losing 97. They finished seventh in the American League and were led in hitting by Bud Stewart, whose batting average was .279.

The work of baseball statisticians is an illustration of descriptive statistics.

Definition 1.2

Statisticians also analyze data for the purpose of making generalizations and decisions.

For example, a political analyst can use data from a portion of the voting population to predict the political preferences of the entire voting population.

Population and Sample

- **Population:** The collection of all individuals or items under consideration in a statistical study.
- **Sample:** That part of the population from which information is obtained.

Example 1.2

Political polling provides an example of **inferential statistics**.

→ Interviewing everyone of voting age in the United States on their voting preferences would be expensive and unrealistic. Statisticians who want to gauge the sentiment of the entire **population** of U.S.

→ voters can afford to interview only a carefully chosen group of a few thousand voters. This group is called a **sample** of the **population**.

Definition 1.3

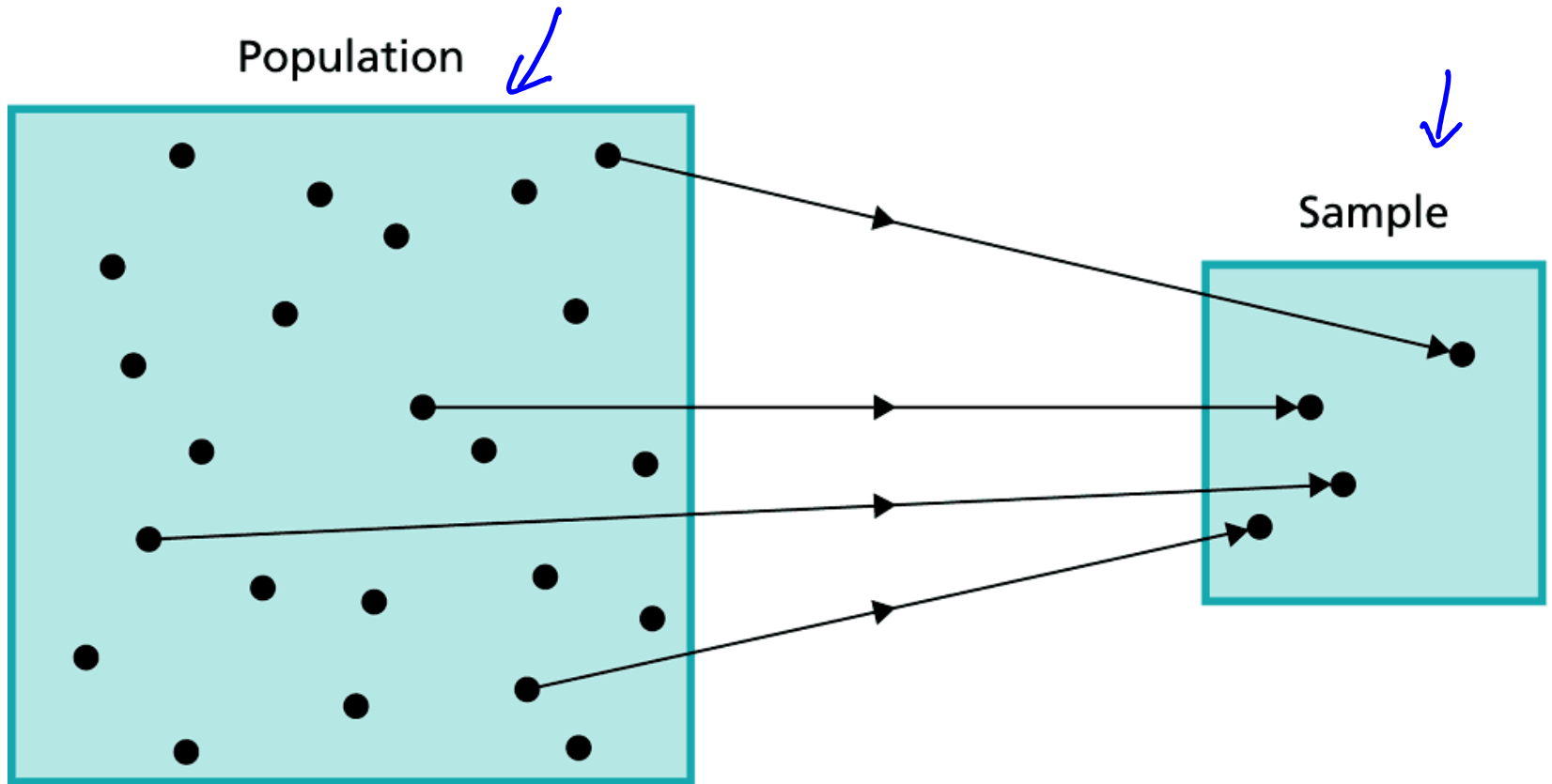
Inferential Statistics

Inferential statistics consists of methods for drawing and measuring the reliability of conclusions about a population based on information obtained from a sample of the population.

→ Statisticians analyze the information obtained from a sample of the voting population to make inferences (draw conclusions) about the preferences of the entire voting population. Inferential statistics provides methods for drawing such conclusions.


Figure 1.1

Relationship between population and sample



Example 1.3 Classifying Statistical Studies

The 1948 Presidential Election Table 1.1 displays the voting results for the 1948 presidential election

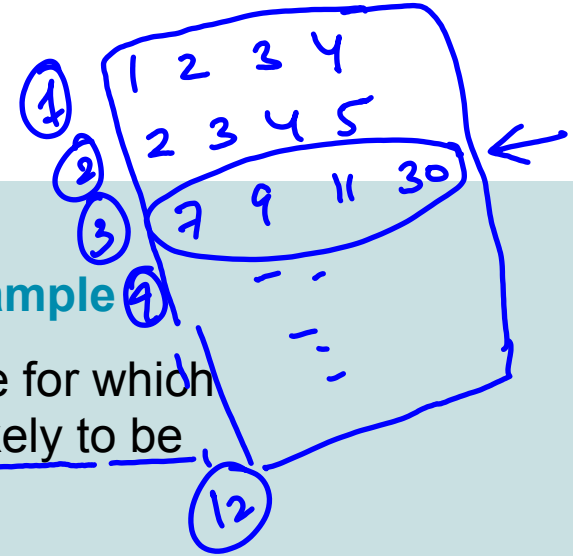


Ticket	Votes	Percentage
Truman–Barkley (<u>Democratic</u>)	<u>24,179,345</u>	49.7 ←
Dewey–Warren (<u>Republican</u>)	<u>21,991,291</u>	45.2 ←
Thurmond–Wright (<u>States Rights</u>)	<u>1,176,125</u>	2.4 ←
Wallace–Taylor (<u>Progressive</u>)	<u>1,157,326</u>	2.4 ←
Thomas–Smith (<u>Socialist</u>)	<u>139,572</u>	0.3 ←

Classification This study is descriptive. It is a summary of the votes cast by U.S. voters in the 1948 presidential election. No inferences are made.

1 2 3 4 - . . 30

Definition 1.4



→ Simple Random Sampling; Simple Random Sample

Simple random sampling: A sampling procedure for which each possible sample of a given size is equally likely to be the one obtained.

→ **Simple random sample:** A sample obtained by simple random sampling.

There are two types of **simple random sampling**. One is simple random sampling **with replacement (SRSWR)**, whereby a member of the population can be selected more than once; the other is **simple random sampling without replacement (SRS)**, whereby a member of the population can be selected at most once.

→ Random-Number Tables

Obtaining a simple random sample by picking slips of paper out of a box is usually impractical, especially when the population is large.

Fortunately, we can use several practical procedures to get simple random samples. One common method involves a **table of random numbers** – a table of randomly chosen digits, as illustrated in Table 1.5.

Table 1.5

Random
numbers

069 ✓
988 X
386 ✓

849 X
578 ✓
404 ✓

Line number	Column number									
	00-09		10-19		20-29		30-39		40-49	
00	15544	80712	97742	21500	97081	42451	50623	56071	28882	28739
01	01011	21285	04729	39986	73150	31548	30168	76189	56996	19210
02	47435	53308	40718	29050	74858	64517	93573	51058	68501	42723
03	91312	75137	86274	59834	69844	19853	06917	17413	44474	86530
04	12775	08768	80791	16298	22934	09630	98862	39746	64623	32768
05	31466	43761	94872	92230	52367	13205	38634	55882	77518	36252
06	09300	43847	40881	51243	97810	18903	53914	31688	06220	40422
07	73582	13810	57784	72454	68997	72229	30340	08844	53924	89630
08	11092	81392	58189	22697	41063	09451	09789	00637	06450	85990
09	93322	98567	00116	35605	66790	52965	62877	21740	56476	49296
10	80134	12484	67089	08674	70753	90959	45842	59844	45214	36505
11	97888	31797	95037	84400	76041	96668	75920	68482	56855	97417
12	92612	27082	59459	69380	98654	20407	88151	56263	27126	63797
13	72744	45586	43279	44218	83638	05422	00995	70217	78925	39097
14	96256	70653	45285	26293	78305	80252	03625	40159	68760	84716
15	07851	47452	66742	83331	54701	06573	98169	37499	67756	68301
16	25594	41552	96475	56151	02089	33748	65289	89956	89559	33687
17	65358	15155	59374	80940	03411	94656	69440	47156	77115	99463
18	09402	31008	53424	21928	02198	61201	02457	87214	59750	51330
19	97424	90765	01634	37328	41243	33564	17884	94747	93650	77668

15 RN in b/w 1 → 728

Random-Number Generators


Nowadays, statisticians prefer statistical software packages or graphing calculators, rather than random-number tables, to obtain simple random samples. The built-in programs for doing so are called **random-number generators**. When using random-number generators, be aware of whether they provide samples with replacement or samples without replacement.

TI-83/84 PLUS

- 1 Press **PRGM**
- 2 Arrow down to SRS and press ENTER twice
- 3 Type 1 for MIN (the smallest possible value) and press ENTER
- 4 Type 728 for MAX (the largest possible value) and press ENTER
- 5 Type 15 for SAMPLE SIZE and press ENTER
- 6 After the program completes, press STAT and then ENTER
- 7 The required sample is in L1 (List 1)



TI-83/84 PLUS

NORMAL FLOAT AUTO REAL RADIAN MP 					
L1	L2	L3	L4	L5	1
86	-----	-----	-----	-----	
176					
622					
272					
435					
442					
57					
13					
10					
507					
511					
L1(1)=86					

Note: Only the first 11 numbers are visible in this output.


Systematic Random Sampling

Procedure 1.1

Systematic Random Sampling

Step 1 Divide the population size by the sample size and round the result down to the nearest whole number, m .

Step 2 Use a random-number table or a similar device to obtain a number, k , between 1 and m .

 **Step 3** Select for the sample those members of the population that are numbered k , $k + m$, $k + 2m$, $k + (n-1)m$

Example: $N=728$, $n=15$ ←

Step 1: $m=728/15 = 48$ (rounded down)

Step 2: $k=22$ (random number between 1 and 48)

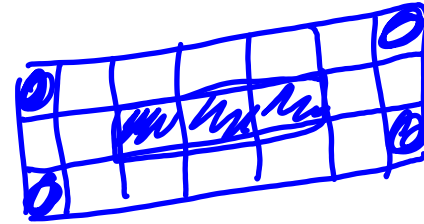
Step 3: Sample numbered: 22, 70, 118, 166, 694.

22, $22+48$, $22+48+48$, ...

$22 + 9 \times 48$ 10th Sample

$694 = 22 + 14 \times 48$

C1 C2 C3 C4 C5 - - - C10
1-30 31-60 61-90 91-120 . . . 271-300



Cluster Sampling

Another sampling method is **cluster sampling**, which is particularly useful when the members of the population are widely scattered geographically.

Procedure 1.2

Cluster Sampling

- Step 1** Divide the population into groups (clusters).
- Step 2** Obtain a simple random sample of the clusters.
- Step 3** Use all the members of the clusters obtained in Step 2 as the sample.

Example: The 300 members of a population have been divided into clusters of equal size 30. Use cluster sampling to obtain a sample of size 60 from the population.

sample = C2, C4
31-60, 91-120

Stratified Sampling

Procedure 1.3

Stratified Random Sampling with Proportional Allocation

Step 1 Divide the population into subpopulations (strata). ←

Step 2 From each stratum, obtain a simple random sample of size proportional to the size of the stratum; that is, the sample size for a stratum equals the total sample size times the stratum size divided by the population size.

Step 3 Use all the members obtained in Step 2 as the sample.

$$\text{s.s. of stratum} = \frac{\text{Total sample size} \times \text{stratum size}}{\text{population size}}$$

Example: The 2000 members of a population have been divided into four strata of sizes 400, 600, 800, and 200. Use stratified sampling with proportional allocation to obtain a sample of size 10 from the population.

Stratum	Size	Numbered	Sample size	Sample
#1	400	1-400	2	166, 264
#2	600	401-1000	3	454, 511, 620
#3	800	1001-1800	4	1246, 1420, 1759, 1793
#4	200	1801-2000	1	1938

+ 2000

$$\text{s.s. of stratum \#1} = \frac{10 \times 400}{2000} = 2$$

$$\text{~ ~ ~ \#4} = \frac{10 \times 200}{2000} = 1$$

Experimental Designs

Other than a census and sampling, another method for obtaining information is experimentation.

Observational Studies and Designed Experiments

Besides classifying statistical studies as either descriptive or inferential, we often need to classify them as either *observational studies* or *designed experiments*.

In an **observational study**, researchers simply observe characteristics and take measurements, as in a sample survey.

In a **designed experiment**, researchers impose treatments and controls and then observe characteristics and take measurements.

Observational studies can reveal only *association*, whereas designed experiments can help establish *causation*.

Vasectomies and Prostate Cancer

One study found 113 cases of prostate cancer among 22,000 men who had a vasectomy. This compares to a rate of 70 cases per 22,000 among men who didn't have a vasectomy." The study shows about a 60% elevated risk of prostate cancer for men who have had a vasectomy, thereby revealing an association between vasectomy and prostate cancer. But does it establish causation: that having a vasectomy causes an increased risk of prostate cancer?

The answer is no, because the study was observational. The researchers simply observed two groups of men, one with vasectomies and the other without. Thus, although an association was established between vasectomy and prostate cancer, the association might be due to other factors (e.g., temperament) that make some men more likely to have vasectomies and also put them at greater risk of prostate cancer.

Folic Acid and Birth Defects

For the study, the doctors enrolled 4753 women prior to conception and divided them randomly into two groups. One group took daily multivitamins containing 0.8 mg of folic acid, whereas the other group received only trace elements (minute amounts of copper, manganese, zinc, and vitamin C). A drastic reduction in the rate of major birth defects occurred among the women who took folic acid: 13 per 1000, as compared to 23 per 1000 for those women who did not take folic acid.

This is a designed experiment and does help establish causation. The researchers did not simply observe two groups of women but, instead, randomly assigned one group to take daily doses of folic acid and the other group to take only trace elements.

Terminology of Experimental Design

Response Variable, Factors, Levels, and Treatments

Response variable: The characteristic of the experimental outcome that is to be measured or observed.

Factor: A variable whose effect on the response variable is of interest in the experiment.

Levels: The possible values of a factor.

Treatment: Each experimental condition. For one-factor experiments, the treatments are the levels of the single factor. For multifactor experiments, each treatment is a combination of levels of the factors.

Principles of Experimental Design

Three basic principles of experimental design: **control**, **randomization**, and **replication**.

Control: Two or more treatments should be compared.

Randomization: The experimental units should be randomly divided into groups to avoid unintentional selection bias in constituting the groups.

Replication: A sufficient number of experimental units should be used to ensure that randomization creates groups that resemble each other closely and to increase the chances of detecting any differences among the treatments.

Principles of Experimental Design

Control: The doctors compared the rate of major birth defects for the women who took folic acid to that for the women who took only trace elements.

Randomization: The women were divided randomly into two groups to avoid unintentional selection bias.

Replication: A large number of women were recruited for the study to make it likely that the two groups created by randomization would be similar and also to increase the chances of detecting any effect due to the folic acid.

Definition 1.8

Randomized Block Design

In a **randomized block design**, the experimental units are assigned randomly among all the treatments separately within each block.

Although the completely randomized design is commonly used and simple, it is not always the best design. Several alternatives to that design exist. For instance, in a **randomized block design**, experimental units that are similar in ways that are expected to affect the response variable are grouped in **blocks**. Then the random assignment of experimental units to the treatments is made block by block.

Example 1.16 Statistical Designs

Golf Ball Driving Distances

Suppose we want to compare the driving distances for five different brands of golf ball. For 40 golfers, discuss a method of comparison based on

- a. a completely randomized design.
- b. a randomized block design.

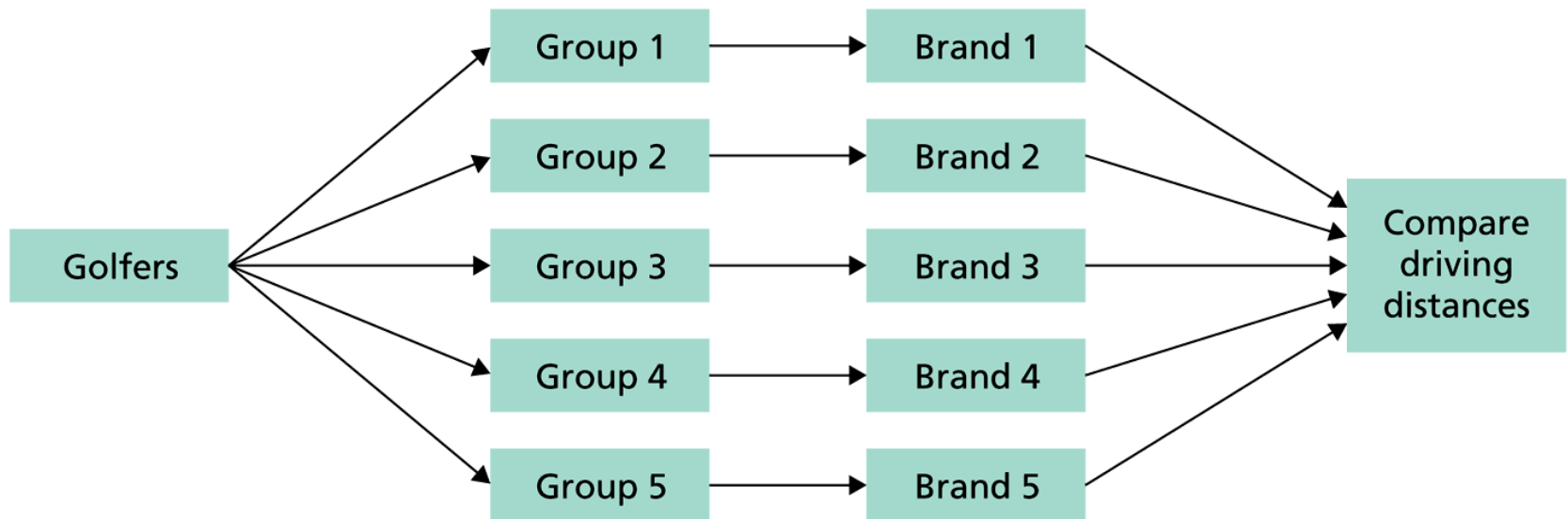
Solution

Here the experimental units are the golfers, the response variable is driving distance, the factor is brand of golf ball, and the levels (and treatments) are the five brands.

- a. For a completely randomized design, we would randomly divide the 40 golfers into five groups of 8 golfers each and then randomly assign each group to drive a different brand of ball, as illustrated in Fig.1.5.

Figure 1.5

Completely randomized design for golf ball experiment



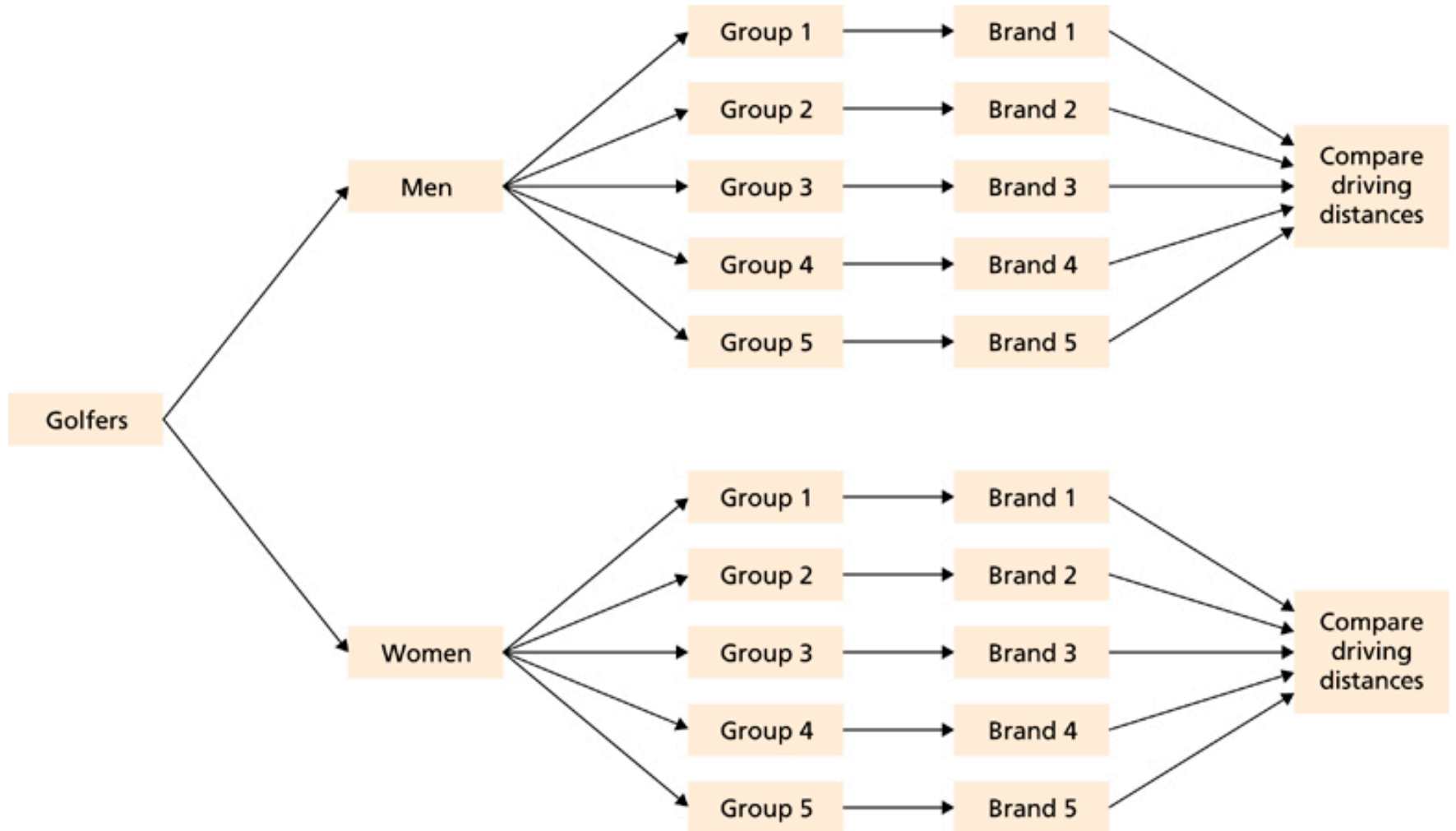
Example 1.16 Statistical Designs

Golf Ball Driving Distances

b. Because driving distance is affected by gender, using a randomized block design that blocks by gender is probably a better approach. We could do so by using 20 men golfers and 20 women golfers. We would randomly divide the 20 men into five groups of 4 men each and then randomly assign each group to drive a different brand of ball, as shown in Fig.1.6. Likewise, we would randomly divide the 20 women into five groups of 4 women each and then randomly assign each group to drive a different brand of ball, as also shown in Fig.1.6.

Figure 1.6

Randomized block design for golf ball experiment



By blocking, we can isolate and remove the variation in driving distances between men and women and thereby make it easier to detect any differences in driving distances among the five brands of golf ball. Additionally, blocking permits us to analyze separately the differences in driving distances among the five brands for men and women.

As illustrated in Example 1.16, blocking can isolate and remove systematic differences among blocks, thereby making any differences among treatments easier to detect. Blocking also makes possible the separate analysis of treatment effects on each block.