

PART I

Fundamental Statistical Concepts

CHAPTER 1

Statistics in Engineering and Science

In this chapter we introduce basic statistical concepts and terminology that are fundamental to the use of statistics in experimental work. These concepts include:

- *the role of statistics in engineering and scientific experimentation,*
- *the distinction between samples and populations,*
- *relating sample statistics to populations parameters, and*
- *characterizing deterministic and empirical models.*

The term *scientific* suggests a process of objective investigation that ensures that valid conclusions can be drawn from an experimental study. Scientific investigations are important not only in the academic laboratories of research universities but also in the engineering laboratories of industrial manufacturers. *Quality* and *productivity* are characteristic goals of industrial processes, which are expected to result in goods and services that are highly sought by consumers and that yield profits for the firms that supply them. Recognition is now being given to the necessary link between the scientific study of industrial processes and the quality of the goods produced. The stimulus for this recognition is the intense international competition among firms selling similar products to a limited consumer group.

The setting just described provides one motivation for examining the role of statistics in scientific and engineering investigations. It is no longer satisfactory just to monitor on-line industrial processes and to ensure that products are within desired specification limits. Competition demands that a better product be produced within the limits of economic realities. Better products are initiated in academic and industrial research laboratories, made feasible in

pilot studies and new-product research studies, and checked for adherence to design specifications throughout production. All of these activities require experimentation and the collection of data. The definition of the discipline of statistics in Exhibit 1.1 is used to distinguish the field of statistics from other academic disciplines and is oriented toward the experimental focus of this text. It clearly identifies statistics as a scientific discipline, which demands the same type of rigor and adherence to basic principles as physics or chemistry. The definition also implies that when problem solving involves the collection of data, the science of statistics should be an integral component of the process.

EXHIBIT 1.1

Statistics. Statistics is the science of problem-solving in the presence of variability.

Perhaps the key term in this definition is the last one. The problem-solving process involves a degree of uncertainty through the natural variation of results that occurs in virtually all experimental work.

When the term *statistics* is mentioned, many people think of games of chance as the primary application. In a similar vein, many consider statisticians to be “number librarians,” merely counters of pertinent facts. Both of these views are far too narrow, given the diverse and extensive applications of statistical theory and methodology.

Outcomes of games of chance involve uncertainty, and one relies on probabilities, the primary criteria for statistical decisions, to make choices. Likewise, the determination of environmental standards for automobile emissions, the forces that act on pipes used in drilling oil wells, and the testing of commercial drugs all involve some degree of uncertainty. Uncertainty arises because the level of emissions for an individual automobile, the forces exerted on a pipe in one well, and individual patient reactions to a drug vary with each observation, even if the observations are taken under “controlled” conditions. These types of applications are only a few of many that could be mentioned. Many others are discussed in subsequent chapters of this book.

Figure 1.1 symbolizes the fact that statistics should play a role in every facet of data collection and analysis, from initial problem formulation to the drawing of final conclusions. This figure distinguishes two types of studies: experimental and observational. In experimental studies the variables of interest often can be controlled and fixed at predetermined values for each test run in the experiment. In observational studies many of the variables of interest cannot be controlled, but they can be recorded and analyzed. In this book we emphasize experimental studies, although many of the analytic procedures discussed can be applied to observational studies.

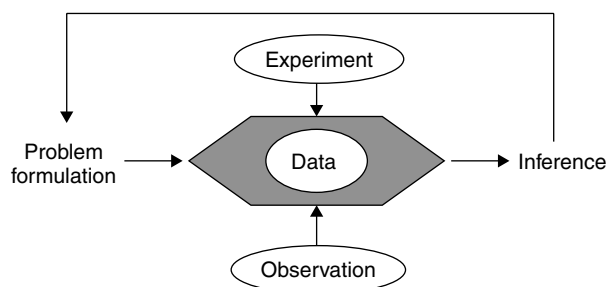


Figure 1.1 Critical stages of statistical input in scientific investigations.

Data are at the center of experimental and observational studies. As will be stressed in Section 1.1, all data are subject to a variety of sources that induce variation in measurements. This variation can occur because of fixed differences among machines, random differences due to changes in ambient conditions, measurement error in instrument readings, or effects due to many other known or unknown influences.

Statistical experimental design will be shown to be effective in eliminating known sources of bias, guarding against unknown sources of bias, ensuring that the experiment provides precise information about the responses of interest, and guaranteeing that excessive experimental resources are not needlessly wasted through the use of an uneconomical design. Likewise, whether one simply wishes to describe the results of an experiment or one wishes to draw inferential conclusions about a process, statistical data-analysis techniques aid in clearly and concisely summarizing salient features of experimental data.

The next section of this chapter discusses the role of statistics in the experimental process, and illustrates how a carefully designed experiment and straightforward statistical graphics can clearly identify major sources of variation in a chemical process. The last three sections of this chapter introduce several concepts that are fundamental to an understanding of statistical inference.

1.1 THE ROLE OF STATISTICS IN EXPERIMENTATION

Statistics is a scientific discipline devoted to the drawing of valid inferences from experimental or observational data. The study of variation, including the construction of experimental designs and the development of models which describe variation, characterizes research activities in the field of statistics. A basic principle that is the cornerstone of the material covered in this book is the following:

All measurements are subject to variation.

The use of the term *measurement* in this statement is not intended to exclude qualitative responses of interest in an experiment, but the main focus of this text is on designs and analyses that are appropriate for quantitative measurements.

In most industrial processes there are numerous sources of possible variation. Frequently studies are conducted to investigate the causes of excessive variation. These studies could focus on a single source or simultaneously examine several sources. Consider, for example, a chemical analysis that involves different specimens of raw materials and that is performed by several operators. Variation could occur because the operators systematically differ in their method of analysis. Variation also could occur because one or more of the operators do not consistently adhere to the analytic procedures, thereby introducing uncontrolled variability to the measurement process. In addition, the specimens sent for analysis could differ on factors other than the ones under examination.

To investigate sources of variability for a chemical analysis similar to the one just described, an experiment was statistically designed and analyzed to ensure that relevant sources of variation could be identified and measured. A test specimen was treated in a combustion-type furnace, and a chemical analysis was performed on it. In the experiment three operators each analyzed two specimens, made three combustion runs on each specimen, and titrated each run in duplicate. The results of the experiment are displayed in Table 1.1 and graphed in Figure 1.2.

Figure 1.2 is an example of a *scatterplot*, a two-dimensional graph of individual data values for pairs of quantitative variables. In Figure 1.2, the abscissa (horizontal axis) is simply the specimen/combustion run index and the ordinate (vertical axis) is the chemical analysis result. Scatterplots can be made for any pair of variables so long as both are quantitative. A scatterplot is constructed by plotting the (x_i, y_i) pairs as indicated in Exhibit 1.2.

EXHIBIT 1.2 SCATTERPLOTS

1. Construct horizontal and vertical axes that cover the ranges of the two variables.
 2. Plot (x_i, y_i) points for each observation in the data set.
-

Figure 1.2 highlights a major problem with the chemical analysis procedure. There are definite differences in the analytic results of the three operators. Operator 1 exhibits very consistent results for each of the two specimens and each of the three combustion runs. Operator 2 produces analytic results that

TABLE 1.1 Results of an Experiment to Identify Sources of Variation in Chemical Analyses^a

Operator	Specimen	Combustion Run	Chemical Analysis	
			1	2
1	1	1	156	154
		2	151	154
		3	154	160
	2	4	148	150
		5	154	157
		6	147	149
2	3	7	125	125
		8	94	95
		9	98	102
	4	10	118	124
		11	112	117
		12	98	110
3	5	13	184	184
		14	172	186
		15	181	191
	6	16	172	176
		17	181	184
		18	175	177

^a Adapted from Snee, R. D. (1983). "Graphical Analysis of Process Variation Studies," *Journal of Quality Technology*, **15**, 76–88. Copyright, American Society for Quality Control, Inc., Milwaukee, WI. Reprinted by permission.

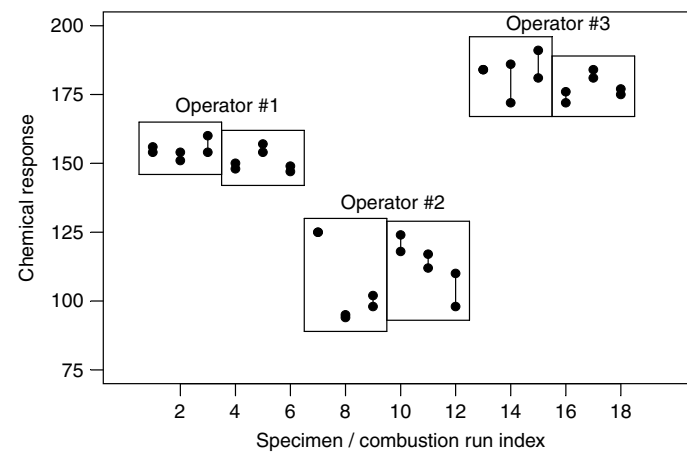


Figure 1.2 Results of a study of variation in a chemical analysis. (Combustion runs are boxed; duplicate analyses are connected by vertical lines.)

are lower on the average than those of the other two operators. Operator 3 shows good consistency between the two specimens, but the repeat analyses of two of the combustion runs on specimen 5 appear to have substantially larger variation than for most of the other repeat analyses in the data set. Operator 2 likewise shows good average consistency for the two specimens, but large variation both for the triplicate combustion runs for each specimen and for at least one of the repeat analyses for the fourth specimen.

Thus, the experimental results indicate that the primary sources of variation in this chemical analysis are the systematic differences (biases) among operators and, in some instances, the (random) inconsistency of the chemical analyses performed by a single operator. In reaching these conclusions statistics played a role in both the design of the experiment and the formal analysis of the results, the foregoing graphical display being one component of the analysis. The quality of this data-collection effort enables straightforward, unambiguous conclusions to be drawn. Such clear-cut inferences are often lacking when data are not collected according to a detailed statistical experimental design.

This example illustrates three general features of the statistical design and analysis of experiments. First, statistical considerations should be included in the project design phase of any experiment. At this stage of a project one should consider the nature of the data to be collected, including what measurements are to be taken, what is known about the likely variation to be encountered, and what factors might influence the variation in the measurements.

Second, a statistical design should be selected that controls, insofar as possible, variation from known sources. The design should allow the estimation of the magnitude of uncontrollable variation and the modeling of relationships between the measurements of interest and factors (sources) believed to influence these measurements.

Uncontrollable variation can arise from many sources. Two general sources of importance to the statistical design of experiments are experimental error and measurement error. Experimental error is introduced whenever test conditions are changed. For example, machine settings are not always exact enough to be fixed at precisely the same value or location when two different test runs call for identical settings. Batches of supposedly identical chemical solutions do not always have exactly the same chemical composition. Measurement errors arise from the inability to obtain exactly the same measurement on two successive test runs when all experimental conditions are unchanged.

Third, a statistical analysis of the experimental results should allow inferences to be drawn on the relationships between the design factors and the measurements. This analysis should be based on both the statistical design

TABLE 1.2 Role of Statistics in Experimentation**Project Planning Phase**

- What is to be measured?
- How large is the likely variation?
- What are the influential factors?

Experimental Design Phase

- Control known sources of variation
- Allow estimation of the size of the uncontrolled variation
- Permit an investigation of suitable models

Statistical Analysis Phase

- Make inferences on design factors
- Guide subsequent designs
- Suggest more appropriate models

and the model used to relate the measurements to the sources of variation. If additional experimentation is necessary or desirable, the analysis should guide the experimenter to an appropriate design and, if needed, a more appropriate model of the measurement process.

Thus, the role of statistics in engineering and scientific experimentation can be described using three basic categories: project planning, experimental design, and data analysis. These three basic steps in the statistical design and analysis of experimental results are depicted in Table 1.2.

1.2 POPULATIONS AND SAMPLES

Experimental data, in the form of a representative sample of observations, enable us to draw inferences about a phenomenon, population, or process of interest. These inferences are obtained by using sample statistics to draw conclusions about postulated models of the underlying data-generating mechanism.

All possible items or units that determine an outcome of a well-defined experiment are collectively called a “population” (see Exhibit 1.3). An item or a unit could be a measurement, or it could be material on which a measurement is taken. For example, in a study of geopressure as an alternative source of electric power, a population of interest might be all geographical locations for which characteristics such as wellhead fluid temperature, pressure, or gas content could be measured. Other examples of populations are:

- all 30-ohm resistors produced by a particular manufacturer under specified manufacturing conditions during a fixed time period;
- all possible fuel-consumption values obtainable with a four-cylinder, 1.7-liter engine using a 10%-methanol, 90%-gasoline fuel blend, tested under controlled conditions on a dynamometer stand;
- all measurements on the fracture strength of one-inch-thick underwater welds on a steel alloy base plate that is located 200 feet deep in a specified salt-water environment; or
- all 1000-lb containers of pelletized, low-density polyethylene produced by a single manufacturing plant under normal operating conditions.

EXHIBIT 1.3

Population. A statistical population consists of all possible items or units possessing one or more common characteristics under specified experimental or observational conditions.

These examples suggest that a population of observations may exist only conceptually, as with the population of fracture-strength measurements. Populations also may represent processes for which the items of interest are not fixed or static; rather, new items are added as the process continues, as in the manufacture of polyethylene.

Populations, as represented by a fixed collection of units or items, are not always germane to an experimental setting. For example, there are no fixed populations in many studies involving chemical mixtures or solutions. Likewise, ongoing production processes do not usually represent fixed populations. The study of physical phenomena such as aging, the effects of drugs, or aircraft engine noise cannot be put in the context of a fixed population of observations. In situations such as these it is a physical process rather than a population that is of interest (see Exhibit 1.4).

EXHIBIT 1.4

Process. A process is a repeatable series of actions that results in an observable characteristic or measurement.

The concepts and analyses discussed in this book relative to samples from populations generally are applicable to processes. For example, one samples both populations and processes in order to draw inferences on models appropriate for each. While one models a fixed population in the former case, one

models a “state of nature” in the latter. A simple random sample may be used to provide observations from which to estimate the population model. A suitably conducted experiment may be used to provide observations from which to estimate the process model. In both situations it is the representative collection of observations and the assumptions made about the data that are important to the modeling procedures.

Because of the direct analogies between procedures for populations and for processes, the focus of the discussions in this book could be on either. We shall ordinarily develop concepts and experimental strategies with reference to only one of the two, with the understanding that they should readily be transferrable to the other. In the remainder of this section, we concentrate attention on developing the relationships between samples and populations.

When defining a relevant population (or process) of interest, one must define the exact experimental conditions under which the observations are to be collected. Depending on the experimental conditions, many different populations of observed values could be defined. Thus, while populations may be real or conceptual, they must be explicitly defined with respect to all known sources of variation in order to draw valid statistical inferences.

The items or units that make up a population are usually defined to be the smallest subdivisions of the population for which measurements or observations can take on different values. For the populations defined above, for example, the following definitions represent units of interest. An individual resistor is the natural unit for studying the actual (as opposed to specified) resistance of a brand of resistors. A measurement of fuel consumption from a single test sequence of accelerations and decelerations is the unit for which data are accumulated in a fuel economy study. Individual welds are the appropriate units for investigating fracture strength. A single container of pellets is the unit of interest in the manufacture of polyethylene.

Measurements on a population of units can exhibit many different statistical properties, depending on the characteristic of interest. Thus, it is important to define the fundamental qualities or quantities of interest in an experiment. We term these qualities or quantities *variables* (see Exhibit 1.5).

EXHIBIT 1.5

Variable. A property or characteristic on which information is obtained in an experiment.

An *observation*, as indicated in Exhibit 1.6, refers to the collection of information in an experiment, and an *observed value* refers to an actual measurement or attribute that is the result of an individual observation. We often

use “observation” in both senses; however, the context of its use should make it clear which meaning is implied.

EXHIBIT 1.6

Observation. The collection of information in an experiment, *or* actual values obtained on variables in an experiment.

A delineation of variables into two categories, response variables (see Exhibit 1.7) and factors (see Exhibit 1.8), is an important consideration in the modeling of data. In some instances response variables are defined according to some probability model which is only a function of certain (usually unknown) constants. In other instances the model contains one or more factors in addition to (unknown) constants.

EXHIBIT 1.7

Response Variable. Any outcome or result of an experiment.

EXHIBIT 1.8

Factors. Controllable experimental variables that can influence the observed values of response variables.

The response variable in a resistor study is the actual resistance measured on an individual resistor. In a study of fuel economy one might choose to model the amount of fuel consumed (response variable) as some function of vehicle type, fuel, driver, ambient temperature, and humidity (factors). In the underwater weld study the response variable is the fracture strength. In the manufacture of polyethylene the response variable of interest might be the actual weight of a container of pellets.

Most of the variables just mentioned are quantitative variables, because each observed value can be expressed numerically. There also exist many qualitative or nonnumerical variables that could be used as factors. Among those variables listed above, ones that could be used as qualitative factors include vehicle type, fuel, and driver.

Populations often are too large to be adequately studied in a specified time period or within designated budgetary constraints. This is particularly true when the populations are conceptual, as in most scientific and engineering experiments, when they represent every possible observation that could be

obtained from a manufacturing process under specified conditions, or when the collection of data requires the destruction of the item. If it is not feasible to collect information on every item in a population, inferences on the population can be made by studying a representative subset of the data, a *sample* (see Exhibit 1.9). Figure 1.3 illustrates one of the primary goals of scientific experimentation and observation: induction from a sample to a population or a process.

EXHIBIT 1.9

Sample. A sample is a group of observations taken from a population or a process.

There are many ways to collect samples in experimental work. A *convenience sample* is one that is chosen simply by taking observations that are easily or inexpensively obtained. The key characteristic of a convenience sample is that all other considerations are secondary to the economic or rapid collection of data. For example, small-scale laboratory studies often are necessary prior to the implementation of a manufacturing process. While this type of pilot study is an important strategy in feasibility studies, the results are generally inadequate for inferring characteristics of the full-scale manufacturing process. Sources of variation on the production line may be entirely different from those in the tightly controlled environment of a laboratory.

Similarly, simply entering a warehouse and conveniently selecting a number of units for inspection may result in a sample of units which exhibits less variation than the population of units in the warehouse. From a statistical viewpoint, convenience samples are of dubious value because the population that they represent may have substantially different characteristics than the population of interest.

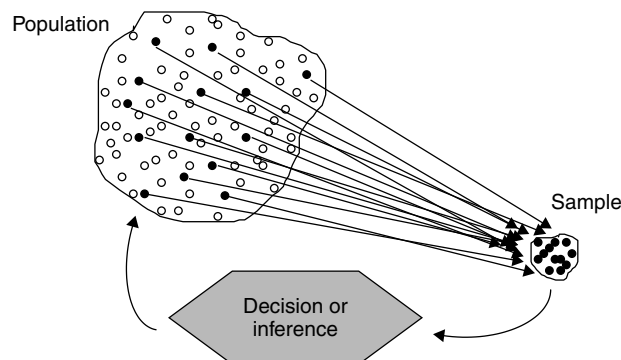


Figure 1.3 Representative samples permit inductive inferences on populations.

Another sampling technique that is frequently used in scientific studies is termed *judgmental sampling*. Here one's experience and professional judgment are used to select representative observations from a population of interest. In the context of a fuel-economy study, an example would be the selection of a particular engine, one fuel blend, and specific laboratory conditions for conducting the study. If the conditions selected for study are not truly representative of typical engine, fuel, and operating conditions, it is difficult to define the relevant population to which the observed fuel-consumption values pertain. A current example of this problem is the E.P.A. fuel-economy ratings posted on automobile stickers. These figures are comparable only under the laboratory conditions under which the estimates are made, not on any typical vehicle, fuel, or operating conditions.

Convenience and judgmental samples are important in exploratory research. The difficulty with these types of samples is that important sources of variation may be held constant or varied over a narrower range than would be the case for the natural occurrence of experimental units from the population of interest. In addition, these sampling schemes may mask the true effects of influential factors. This is an especially acute problem if two or more factors jointly influence a response. Holding one or more of these joint factors constant through convenience or judgmental sampling could lead to erroneous inferences about the effects of the factors on the response.

One of the most important sampling methodologies in experimental work is the *simple random sample*, defined in Exhibit 1.10. In addition to its use in the sampling of observations from a population, simple random sampling has application in the conduct of scientific and engineering experiments. Among the more prominent uses of simple random sampling in experimental work are the selection of experimental units and the randomization of test runs.

EXHIBIT 1.10

Simple Random Sample. In an experimental setting, a simple random sample of size n is obtained when items are selected from a fixed population or a process in such a manner that every group of items of size n has an equal chance of being selected as the sample.

If one wishes to sample 100 resistors from a warehouse, simple random sampling requires that every possible combination of 100 resistors present in the warehouse have an equal chance of being included in the selected sample. Although the requirements of simple random sampling are more stringent than most other sampling techniques, unintentional biases are avoided.

Simple random samples can be obtained in many ways. For example, in the selection of experimental units to be included in an experiment, a common approach is to enumerate or label each item from 1 to N and then use a

table of random numbers to select n of the N units. If a test program is to consist of n test runs, the test runs are sequentially numbered from 1 to n and a random-number table is used to select the run order. Equivalently, one can use random-number generators, which are available on computers. Use of such tables or computer algorithms removes personal bias from the selection of units or the test run order.

Simple random samples can be taken *with or without replacement*. Sampling with replacement allows an experimental unit to be selected more than once. One simply obtains n numbers from a random-number table without regard to whether any of the selected numbers occur more than once in the sample. Sampling without replacement prohibits any number from being selected more than once. If a number is sampled more than once, it is discarded after the first selection. In this way n unique numbers are selected. The sequencing of test runs is always performed by sampling without replacement. Ordinarily the selection of experimental units is also performed by sampling without replacement.

Inspection sampling of items from lots in a warehouse is an example for which a complete enumeration of experimental units is possible, at least for those units that are present when the sample is collected. When a population of items is conceptual or an operating production process is being studied, this approach is not feasible. Moreover, while one could conceivably sample at random from a warehouse full of units, the expense suffered through the loss of integrity of bulk lots of product when a single item is selected for inclusion in a sample necessitates alternative sampling schemes.

There are many other types of random sampling schemes besides simple random sampling. *Systematic* random samples are obtained by sampling every k th (e.g., every 5th, 10th, or 100th) unit in the population. *Stratified random samples* are based on subdividing a heterogeneous population into groups, or *strata*, of similar units and selecting simple random samples from each of the strata. *Cluster sampling* is based on subdividing the population into groups, or clusters, of units in such a way that it is convenient to randomly sample the clusters and then either randomly sample or completely enumerate all the observations in each of the sampled clusters. More details on these and other alternatives to simple random sampling are given in the recommended readings at the end of this chapter.

Regardless of which sampling technique is used, the key idea is that the sample should be representative of the population under study. In experimental settings for which the sampling of populations or processes is not germane, the requirement that the data be representative of the phenomenon or the “state of nature” being studied is still pertinent and necessary. Statistics, as a science, seeks to make inferences about a population, process, or phenomenon based on the information contained in a representative sample or collection of observations.

TABLE 1.3 Employee Identification Numbers

1	A11401	41	B09087	81	G07704	121	B04256
2	P04181	42	B00073	82	K20760	122	K05170
3	N00004	43	J08742	83	W00124	123	R07790
4	C03253	44	W13972	84	T00141	124	G15084
5	D07159	45	S00856	85	M25374	125	C16254
6	M00079	46	A00166	86	K03911	126	R20675
7	S15552	47	S01187	87	W01718	127	G06144
8	G01039	48	D00022	88	T04877	128	T12150
9	P00202	49	Z01194	89	M22262	129	R07904
10	R22110	50	M32893	90	C00011	130	M24214
11	D00652	51	K00018	91	W23233	131	D00716
12	M06815	52	H16034	92	K10061	132	M27410
13	C09071	53	F08794	93	K11411	133	J07272
14	S01014	54	S71024	94	B05848	134	L02455
15	D05484	55	G00301	95	L06270	135	D06610
16	D00118	56	B00103	96	K08063	136	M31452
17	M28883	57	B29884	97	P07211	137	L25264
18	G12276	58	G12566	98	F28794	138	M10405
19	M06891	59	P03956	99	L00885	139	D00393
20	B26124	60	B00188	100	M26882	140	B52223
21	D17682	61	J21112	101	M49824	141	M16934
22	B42024	62	J08208	102	R05857	142	M27362
23	K06221	63	S11108	103	L30913	143	B38384
24	C35104	64	M65014	104	B46004	144	H08825
25	M00709	65	M07436	105	R03090	145	S14573
26	P00407	66	H06098	106	H09185	146	B23651
27	P14580	67	S18751	107	J18200	147	S27272
28	P13804	68	W00004	108	W14854	148	G12636
29	P23144	69	M11028	109	S01078	149	R04191
30	D00452	70	L00213	110	G09221	150	D13524
31	B06180	71	J06070	111	M17174	151	G00154
32	B69674	72	B14514	112	L04792	152	B19544
33	H11900	73	H04177	113	S23434	153	V01449
34	M78064	74	B26003	114	T02877	154	F09564
35	L04687	75	B26193	115	K06944	155	L09934
36	F06364	76	H28534	116	E14054	156	A10690
37	G24544	77	B04303	117	F00281	157	N02634
38	T20132	78	S07092	118	H07233	158	W17430
39	D05014	79	H11759	119	K06204	159	R02109
40	R00259	80	L00252	120	K06423	160	C18514

To illustrate the procedures involved in randomly sampling a population, consider the information contained in Table 1.3. The table enumerates a portion of the world-wide sales force of a manufacturer of skin products. The employees are identified in the table by the order of their listing (1–160) and by their employee identification numbers. Such a tabulation might be obtained from a computer printout of personnel records. For the purposes of the study to be described, these 160 individuals form a population that satisfies several criteria set forth in the experimental protocol.

Suppose the purpose of a study involving these employees is to investigate the short-term effects of certain skin products on measurements of skin elasticity. Initial skin measurements are available for the entire population of employees (see Table 1.4). However, the experimental protocol requires that skin measurements be made on a periodic basis, necessitating the transportation of each person in the study to a central measuring laboratory. Because of the expense involved, the researchers would like to limit the participants included in the study to a simple random sample of 25 of the employees listed in Table 1.3.

Because the population of interest has been completely enumerated, one can use a random-number table (e.g., Table A1 of the Appendix) or a computer-generated sequence of random numbers to select 25 numbers between 1 and 160. One such random number sequence is

57, 77, 8, 83, 92, 18, 63, 121, 19, 115, 139, 96, 133,
131, 122, 17, 79, 2, 68, 59, 157, 138, 26, 70, 9.

Corresponding to this sequence of random numbers is the sequence of employee identification numbers that determines which 25 of the 160 employees are to be included in the sample:

B29884, B04303, G01039, W00124, K10061, G12276, S11108,
B04256, M06891, K06944, D00393, K08063, J07272, D00716,
K05170, M28883, H11759, P04181, W00004, P03956, N02634,
M10405, P00407, L00213, P00202.

With periodic measurements taken on only this random sample of employees the researchers wish to draw conclusions about skin elasticity for the population of employees listed in Table 1.3. This statement suggests that a distinction must be made between measured characteristics taken on a population and those taken on a sample. This distinction is made explicit in the next section.

TABLE 1.4 Skin Elasticity Measurements

1	31.9	41	36.0	81	36.3	121	33.0
2	33.1	42	28.6	82	36.3	122	37.4
3	33.1	43	38.0	83	41.5	123	33.8
4	38.5	44	39.1	84	33.0	124	35.3
5	39.9	45	39.4	85	36.3	125	37.5
6	36.5	46	30.6	86	36.3	126	31.6
7	34.8	47	34.1	87	30.9	127	33.1
8	38.9	48	40.8	88	32.3	128	38.2
9	40.3	49	35.1	89	39.2	129	31.4
10	33.6	50	34.1	90	35.2	130	35.9
11	36.4	51	36.3	91	35.1	131	37.6
12	34.4	52	35.1	92	33.9	132	35.5
13	35.7	53	35.0	93	42.0	133	34.2
14	33.9	54	39.0	94	35.1	134	34.0
15	36.6	55	34.0	95	34.5	135	31.3
16	36.0	56	35.3	96	35.0	136	32.6
17	30.8	57	36.0	97	35.1	137	34.9
18	31.1	58	34.7	98	35.7	138	35.3
19	37.6	59	39.8	99	36.4	139	35.1
20	35.7	60	35.8	100	39.6	140	35.7
21	29.6	61	35.7	101	35.2	141	32.3
22	37.3	62	39.8	102	37.2	142	38.1
23	31.4	63	36.4	103	33.3	143	36.8
24	31.6	64	36.1	104	33.7	144	38.7
25	34.6	65	37.7	105	37.8	145	40.0
26	34.6	66	32.3	106	34.4	146	35.4
27	33.7	67	35.6	107	36.9	147	34.0
28	30.9	68	38.2	108	31.8	148	34.3
29	34.6	69	39.0	109	35.3	149	32.8
30	37.0	70	34.3	110	38.1	150	30.7
31	35.3	71	40.6	111	34.1	151	34.4
32	36.3	72	37.4	112	35.8	152	34.3
33	31.8	73	37.3	113	33.3	153	35.8
34	38.2	74	36.9	114	33.8	154	37.5
35	34.6	75	29.0	115	36.4	155	34.4
36	36.0	76	39.0	116	36.9	156	35.8
37	40.8	77	33.7	117	35.3	157	31.9
38	39.2	78	32.9	118	37.0	158	36.9
39	33.4	79	33.8	119	33.5	159	34.4
40	34.0	80	36.2	120	40.3	160	30.1

1.3 PARAMETERS AND STATISTICS

Summarization of data can occur in both populations and samples. Parameters, as defined in Exhibit 1.11, are constant population values that summarize the entire collection of observations. Parameters can also be viewed in the context of a stable process or a controlled experiment. In all such settings a parameter is a fixed quantity that represents a characteristic of interest. Some examples are:

- the mean fill level for twelve-ounce cans of a soft drink bottled at one plant,
- the minimum compressive strength of eight-foot-long, residential-grade, oak ceiling supports, and
- the maximum wear on one-half-inch stainless-steel ball bearings subjected to a prescribed wear-testing technique.

EXHIBIT 1.11

Parameters and Statistics. A parameter is a numerical characteristic of a population or a process. A statistic is a numerical characteristic that is computed from a sample of observations.

Parameters often are denoted by Greek letters, such as μ for the mean and σ for the standard deviation (a measure of the variability of the observations in a population), to reinforce the notion that they are (generally unknown) constants. Often population parameters are used to define specification limits or tolerances for a manufactured product. Alternatively they may be used to denote hypothetical values for characteristics of measurements that are to be subjected to scientific or engineering investigations.

In many scientific and engineering contexts the term *parameter* is used as a synonym for variable (as defined in the previous section). The term *parameter* should be reserved for a constant or fixed numerical characteristic of a population and not used for a measured or observed property of interest in an experiment. To emphasize this distinction we will henceforth use Greek letters to represent population parameters and Latin letters to denote variables. Sample statistics, in particular estimates of population parameters, also will generally be denoted by Latin letters.

The term *distribution* (see Exhibit 1.12) is used throughout this text to refer to the possible values of a variable along with some measure of how frequently they occur. In a sample or a population the frequency could be measured by counts or percentages. Often when dealing with populations or processes the frequency is measured in terms of a probability model specifying the likelihood of occurrence of the values.

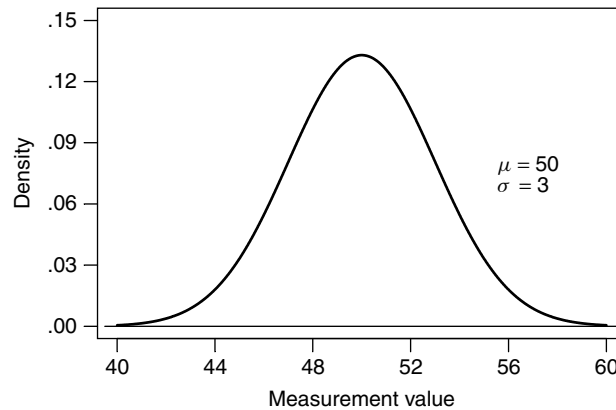


Figure 1.4 Normal distribution of measurement values.

The curve in Figure 1.4 often is used as a probability model, the *normal* distribution, to characterize populations and processes for many types of measurements. The *density* or height of the curve above the axis of measurement values, represents the likelihood of obtaining a value. Probabilities for any range of measurement values can be calculated from the probability model once the model parameters are specified. For this distribution, only the mean and the standard deviation are needed to completely specify the probability model.

EXHIBIT 1.12

Distribution. A tabular, graphical, or theoretical description of the values of a variable using some measure of how frequently they occur in a population, a process, or a sample.

The peak of the curve in Figure 1.4 is located above the measurement value 50, which is the mean μ of the distribution of data values. Because the probability density is highest around the mean, measurement values around the mean are more likely than measurement values greatly distant from it. The standard deviation σ of the distribution in Figure 1.4 is 3. For normal distributions (Section 2.3), approximately 68% of the measurement values lie between $\mu \pm \sigma$ (47 to 53), approximately 95% between $\mu \pm 2\sigma$ (44 to 56), and approximately 99% between $\mu \pm 3\sigma$ (41 to 59). The mean and the standard deviation are very important parameters for the distribution of measurement values for normal distributions such as that of Figure 1.4.

Statistics are sample values that generally are used to estimate population parameters. For example, the average of a sample of observations can be used

to estimate the mean of the population from which the sample was drawn. Similarly, the standard deviation of the sample can be used to estimate the population standard deviation. As we shall see in subsequent chapters, there are often several sample statistics that can be used to estimate a population parameter.

While parameters are fixed constants representing an entire population of data values, statistics are “random” variables and their numerical values depend on which particular observations from the population are included in the sample. One interesting feature about a statistic is that it has its own probability, or *sampling*, distribution: the sample statistic can take on a number of values according to a probability model, which is determined by the probability model for the original population and by the sampling procedure (see Exhibit 1.13). Hence, a statistic has its own probability model as well as its own parameter values, which may be quite different from those of the original population.

EXHIBIT 1.13

Sampling Distribution. A sampling distribution is a theoretical model that describes the probability of obtaining the possible values of a sample statistic.

Histograms are among the most common displays for illustrating the distribution of a set of data. They are especially useful when large numbers of data must be processed. *Histograms* (see Exhibit 1.14) are constructed by dividing the range of the data into several intervals (usually of equal length), counting the number of observations in each interval, and constructing a bar chart of the counts. A by-product of the construction of the histogram is the *frequency distribution*, which is a table of the counts or frequencies for each interval.

EXHIBIT 1.14 FREQUENCY DISTRIBUTIONS AND HISTOGRAMS

1. Construct intervals, ordinarily equally spaced, which cover the range of the data values.
 2. Count the number of observations in each of the intervals. If desirable, form proportions or percentages of counts in each interval.
 3. Clearly label all columns in tables and both axes on histograms, including any units of measurement, and indicate the sample or population size.
 4. For histograms, plot bars whose
 - (a) widths correspond to the measurement intervals, and
 - (b) heights are (proportional to) the counts for each interval (e.g., heights can be counts, proportions, or percentages).
-

Both histograms and the tables of counts that accompany them are sometimes referred to as frequency distributions, because they show how often the data occur in various intervals of the measured variable. The intervals for which counts are made are generally chosen to be equal in width, so that the size (area) of the bar or count is proportional to the number of observations contained in the interval. Selection of the interval width is usually made by simply dividing the range of the data by the number of intervals desired in the histogram or table. Depending on the number of observations, between 8 and 20 intervals are generally selected—the greater the number of observations, the greater the number of intervals.

When the sample size is large, it can be advantageous to construct *relative-frequency* histograms. In these histograms and frequency distributions either the proportions (counts/sample size) or the percentages (proportions $\times 100\%$) of observations in each interval are calculated and graphed, rather than the frequencies themselves. Use of relative frequencies (or percentages) in histograms ensures that the total area under the bars is equal to one (or 100%). This facilitates the comparison of the resultant distribution with that of a theoretical probability distribution, where the total area of the distribution also equals one.

A frequency distribution and histogram for the skin elasticity measurements in Table 1.4 are shown in Table 1.5 and Figure 1.5. The histogram in Figure 1.5 is an example of a relative-frequency histogram. The heights of the bars suggest a shape similar to the form of the normal curve in Figure 1.4. On the basis of these data one might postulate a normal probability model for the skin measurements.

Figure 1.6 shows a normal probability model that has the same mean ($\mu = 35.4$) and standard deviation ($\sigma = 2.65$) as the population of values

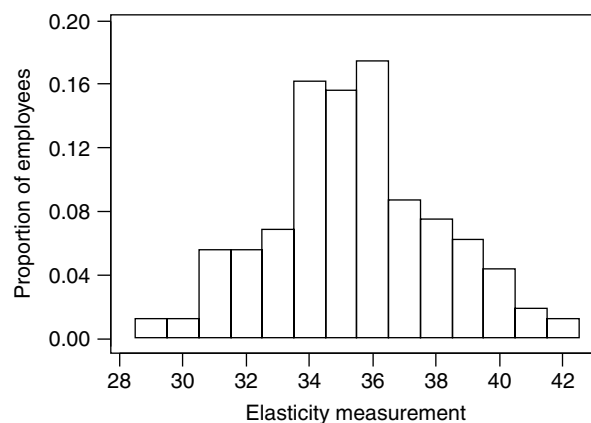
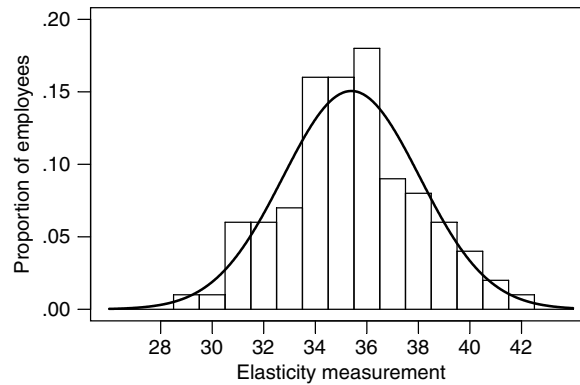


Figure 1.5 Distribution of elasticity measurements ($n = 160$).

TABLE 1.5 Frequency Distribution for Skin Elasticity Data Set

Skin Elasticity*	Interval Midpoint	Frequency	Proportion
28.5–29.5	29	2	0.01
29.5–30.5	30	2	0.01
30.5–31.5	31	9	0.06
31.5–32.5	32	9	0.06
32.5–33.5	33	11	0.07
33.5–34.5	34	26	0.16
34.5–35.5	35	25	0.16
35.5–36.5	36	28	0.18
36.5–37.5	37	14	0.09
37.5–38.5	38	12	0.08
38.5–39.5	39	10	0.06
39.5–40.5	40	7	0.04
40.5–41.5	41	3	0.02
41.5–42.5	42	2	0.01
		160	1.00

*Intervals include lower limits, exclude upper ones.

**Figure 1.6** Normal approximation to elasticity distribution.

in Table 1.4. Observe that the curve for the theoretical normal model provides a good approximation to the actual distribution of the population of measurements, represented by the vertical bars.

One of the features of a normal model is that averages from simple random samples of size n also follow a normal probability model with the same

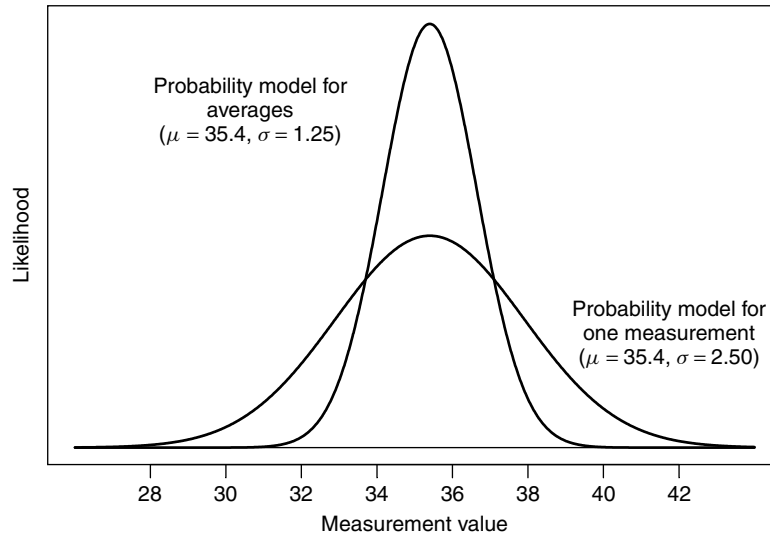


Figure 1.7 Comparison of theoretical normal distributions.

population mean but with a standard deviation that is reduced by a factor of \sqrt{n} from that of the original population. Thus, averages of random samples of size 4 have standard deviations that are half that of the original population. Figure 1.7 shows the relationship between a normal probability model for individual measurements that have a population mean of $\mu = 35.4$ and a standard deviation of $\sigma = 2.5$ and one for the corresponding population of sample averages of size 4. Note that the latter distribution has $\mu = 35.4$ but $\sigma = 2.5/\sqrt{4} = 1.25$. The distribution of the averages is more concentrated around the population mean than is the distribution of individual observations. This indicates that it is much more likely to obtain a sample average that is in a fixed interval around the population mean than it is to obtain a single observation in the same fixed interval.

This discussion is intended to highlight the informational content of population parameters and to shed some light on the model-building processes involved in drawing inferences from sample statistics. The final section in this chapter focuses on one additional issue, which helps to distinguish statistical from mathematical problem solving.

1.4 MATHEMATICAL AND STATISTICAL MODELING

Models and model building are commonplace in the engineering and physical sciences. A research engineer or scientist generally has some basic knowledge about the phenomenon under study and seeks to use this information

to obtain a plausible model of the data-generating process. Experiments are conducted to characterize, confirm, or reject models—in particular, through hypotheses about those models. Models take many shapes and forms, but in general they all seek to characterize one or more response variables, perhaps through relationships with one or more factors.

Mathematical models, as defined in Exhibit 1.15, have the common trait that the response and predictor variables are assumed to be free of specification error and measurement uncertainty. Mathematical models may be poor descriptors of the physical systems they represent because of this lack of accounting for the various types of errors included in statistical models. Statistical models, as defined in Exhibit 1.16, are approximations to actual physical systems and are subject to specification and measurement errors.

EXHIBIT 1.15

Mathematical Model. A model is termed *mathematical* if it is derived from theoretical or mechanistic considerations that represent exact, error-free assumed relationships among the variables.

EXHIBIT 1.16

Statistical Model. A model is termed *statistical* if it is derived from data that are subject to various types of specification, observation, experimental, and/or measurement errors.

An example of a mathematical model is the well-known fracture mechanics relation:

$$K_{IC} = \gamma S a^{1/2}, \quad (1.1)$$

where K_{IC} is the critical stress intensity factor, S is the fracture strength, a is the size of the flaw that caused the fracture, and γ is a constant relating to the flaw geometry. This formula can be utilized to relate the flaw size of a brittle material to its fracture strength. Its validity is well accepted by mechanical engineers because it is based on the theoretical foundations of fracture mechanics, which have been confirmed through extensive experimental testing.

Empirical studies generally do not operate under the idealized conditions necessary for a model like equation (1.1) to be valid. In fact, it often is not possible to postulate a mathematical model for the mechanism being studied. Even when it is known that a model like equation (1.1) should be valid, experimental error may become a nontrivial problem. In these situations statistical models

are important because they can be used to approximate the response variable over some appropriate range of the other model variables. For example, additive or multiplicative *errors* can be included in the fracture-mechanics model, yielding the statistical models

$$K_{IC} = \gamma Sa^{1/2} + e \quad \text{or} \quad K_{IC} = \gamma Sa^{1/2}e \quad (1.2)$$

where e is the error. Note that the use of “error” in statistical models is not intended to indicate that the model is incorrect, only that unknown sources of uncontrolled variation, often measurement error, are present.

A mathematical model, in practice, can seldom be proved with data. At best, it can be concluded that the experimental data are consistent with a particular hypothesized model. The chosen model might be completely wrong and yet this fact might go unrecognized because of the nature of the experiment; e.g., data collected over a very narrow range of the variables would be consistent with any of a vast number of models. Hence, it is important that proposed mathematical models be sufficiently “strained” by the experimental design so that any substantial discrepancies from the postulated model can be identified.

In many research studies there are mathematical models to guide the investigation. These investigations usually produce statistical models that may be partially based on theoretical considerations but must be validated across wide ranges of the experimental variables. Experimenters must then seek “lawlike relationships” that hold under a variety of conditions rather than try to build separate statistical models for each new data base. In this type of model generalization, one may eventually evolve a “theoretical” model that adequately describes the phenomenon under study.

REFERENCES

Text References

The following books provide excellent case studies. Brief summaries of each are provided below.

- Andrews, D. F. and Herzberg, A. M. (1985). *Data: A Collection of Problems from Many Fields for the Student and Research Worker*, New York: Springer-Verlag, Inc. *This book is a collection of 71 data sets with descriptions of the experiment or the study from which the data were collected. The data sets exemplify the rich variety of problems for which the science of statistics can be used as an integral component in problem-solving.*
- Peck, R., Haugh, L. D., and Goodman, A. (1998). *Statistical Case Studies: A Collaboration Between Academe and Industry*. ASA-SIAM Series on Statistics and Applied Probability. Philadelphia: Society for Industrial and Applied Mathematics.

This collection of 20 case studies is unique in that each is co-authored by at least one academic and at least one industrial partner. Each case study is based on an actual project with specific research goals. A wide variety of data collection methods and data analysis techniques are presented.

- Snee, R. D., Hare, L. B., and Trout, J. R. (1985). *Experiments in Industry: Design, Analysis, and Interpretation of Results*, Milwaukee, WI: American Society for Quality Control.

This collection of eleven case histories focuses on the design, analysis, and interpretation of scientific experiments. The stated objective is "to show scientists and engineers not familiar with this methodology how the statistical approach to experimentation can be used to develop and improve products and processes and to solve problems."

- Tanur, J. M., Mosteller, F., Kruskal, W. H., Link, R. F., Pieters, R. S., and Rising, G. R. (1972). *Statistics: A Guide to the Unknown*, San Francisco: Holden-Day, Inc.

This collection of 44 essays on applications of statistics presents excellent examples of the uses and abuses of statistical methodology. The authors of these essays present applications of statistics and probability in nontechnical expositions which for the most part do not require previous coursework in statistics or probability. Individual essays illustrate the benefits of carefully planned experimental designs as well as numerous examples of statistical analysis of data from designed experiments and observational studies.

The following books contain excellent discussions on the impact of variation on processes and on achieving business objectives:

- Leitnaker, M. G., Sanders, R. D., and Hild, C. (1996). *The Power of Statistical Thinking: Improving Industrial Processes*. New York: Addison Wesley Publishing Company.

- Joner, B. L. (1994). *Fourth Generation Management: The New Business Consciousness*. New York: McGraw-Hill, Inc.

The distinction between populations and samples and between population parameters and sample statistics is stressed in most elementary statistics textbooks. The references at the end of Chapter 2 provide good discussions of these concepts. There are many excellent textbooks on statistical sampling techniques, including:

- Cochran, W. G. (1977). *Sampling Techniques, Third Edition*, New York: John Wiley and Sons, Inc.

- Scheaffer, R. L., Mendenhall, W., and Ott, L. (1996). *Elementary Survey Sampling, Fifth Edition*, North Scituate, MA: Duxbury Press.

The first of these texts is a classic in the statistical literature. The second one is more elementary and a good reference for those not familiar with sampling techniques.

Explicit instructions for constructing frequency distributions, histograms, and scatter-plots can be found in most elementary statistics texts, including:

- Freedman, D., Pisani, R., and Purves, R. (1997). *Statistics Third Edition*, New York: W. W. Norton & Company, Chapters 3 and 7.

- Koopmans, L. (1981). *An Introduction to Contemporary Statistics*, Belmont, CA: Duxbury Press, Chapters 1 and 4.

Ott, L. (1977). *An Introduction to Statistical Methods and Data Analysis*, Belmont, CA: Duxbury Press, Chapters 1 and 6.

Mathematical and statistical modeling is not extensively covered in introductory statistics textbooks. The following text does devote ample space to this important topic:

Box, G. E. P., Hunter, W. G., and Hunter, J. S. (1978). *Statistics for Experimenters*. New York: John Wiley & Sons, Inc., Chapters 9 and 16.

EXERCISES

- 1 The Department of Transportation (DOT) was interested in evaluating the safety performance of motorcycle helmets manufactured in the United States. A total of 264 helmets were obtained from the major U.S. manufacturers and supplied to an independent research testing firm where impact penetration and chin retention tests were performed on the helmets in accordance with DOT standards.
 - (a) What is the population of interest?
 - (b) What is the sample?
 - (c) Is the population finite or infinite?
 - (d) What inferences can be made about the population based on the tested samples?
- 2 List and contrast the characteristics of population parameters and sample statistics.
- 3 A manufacturer of rubber wishes to evaluate certain characteristics of its product. A sample is made from a warehouse containing bales of synthetic rubber. List some of the possible candidate populations from which this sample can be taken.
- 4 It is known that the bales of synthetic rubber described in Exercise 3 are stored on pallets with a total of 15 bales per pallet. What type of sampling methodology is being implemented under the following sample scenarios?
 - (a) Five pallets of bales are randomly chosen; then eight bales of rubber are randomly selected from each pallet.
 - (b) Forty bales are randomly selected from the 4500 bales in the warehouse.
 - (c) All bales are sampled on every fifth pallet in the warehouse.
 - (d) All bales that face the warehouse aisles and can be reached by a forklift truck are selected.
- 5 Recall the normal distribution discussed in Section 1.3. What is the importance of $\mu \pm 3\sigma$?
- 6 The population mean and standard deviation of typical cetane numbers measured on fuels used in compression-ignition engines is known to be

$\mu = 30$ and $\sigma = 5$. Fifteen random samples of these fuels were taken from the relevant fuel population, and the sample means and standard deviations were calculated. This random sampling procedure was repeated (replicated) nine times.

Replicate No.	n	Sample Mean	Sample Standard Deviation
1	15	32.61	4.64
2	15	28.57	6.49
3	15	29.66	4.68
4	15	30.09	5.35
5	15	30.11	6.39
6	15	28.02	4.05
7	15	30.09	5.35
8	15	29.08	3.56
9	15	28.91	4.88

Consider the population of all sample means of size $n = 15$. What proportion of means from this population should be expected to be between the 30 ± 5 limits? How does this sample of nine averages compare with what should be expected?

- 7 A research program was directed toward the design and development of self-restoring traffic-barrier systems capable of containing and redirecting large buses and trucks. Twenty-five tests were conducted in which vehicles were driven into the self-restoring traffic barriers. The range of vehicles used in the study included a 1800-lb car to a 40,000-lb intercity bus. Varying impact angles and vehicle speeds were used, and the car damage, barrier damage, and barrier containment were observed.
 - (a) What is an observation in this study?
 - (b) Which variables are responses?
 - (c) Which variables are factors?
- 8 Space vehicles contain fuel tanks that are subject to liquid sloshing in low-gravity conditions. A theoretical basis for a model of low-gravity sloshing was derived and used to predict the slosh dynamics in a cylindrical tank. Low-gravity simulations were performed in which the experimental results were used to verify a statistical relationship. It was shown in this study that the statistical model closely resembled the theoretical model. What type of errors are associated with the statistical model? Why aren't the statistical and theoretical models exactly the same?
- 9 Use a table of random numbers or a computer-generated random number sequence to draw 20 simple random samples, each of size $n = 10$, from the population of employees listed in Table 1.3. Calculate the average of

the skin elasticity measurements (Table 1.4) for each sample. Graph the distribution of these 20 averages in a form similar to Figure 1.5. Would this graph be sufficient for you to conclude that the sampling distribution of the population of averages is a normal distribution? Why (not)?

- 10 Use the table of random numbers in the Appendix to choose starting points between 1 and 16 for ten systematic random samples of the population of employees listed in Table 1.3. Select every 10th employee. Calculate the average of the skin elasticity measurements for each sample. Does a graph of the distribution of these averages have a similar shape to that of Exercise 9?
- 11 Simple random samples and systematic random samples often result in samples that have very similar characteristics. Give three examples of populations that you would expect to result in similar simple and systematic random samples. Explain why you expect the samples to be similar. Give three examples for which you expect the samples to be different. Explain why you expect them to be different.
- 12 A series of valve closure tests were conducted on a 5-inch-diameter speed-control valve. The valve has a spring-loaded poppet mechanism that allows the valve to remain open until the flow drag on the poppet is great enough to overcome the spring force. The poppet then closes, causing flow through the valve to be greatly reduced. Ten tests were run at different spring locking-nut settings, where the flow rate at which the valve poppet closed was measured in gallons per minute. Produce a scatterplot of these data. What appears to be the effect of nut setting on flow rate?

Nut Setting	Flow Rate	Nut Setting	Flow Rate
0	1250	10	2085
2	1510	12	1503
4	1608	14	2115
6	1650	16	2350
8	1825	18	2411

- 13 A new manufacturing process is being implemented in a factory that produces automobile spark plugs. A random sample of 50 spark plugs is selected each day over a 15-day period. The spark plugs are examined and the number of defective plugs is recorded each day. Plot the following data in a scatterplot with the day number on the horizontal axis. A scatterplot with time on the horizontal axis is often called a *sequence plot*. What does the plot suggest about the new manufacturing process?

Day	No. of Defectives	Day	No. of Defectives
1	3	9	4
2	8	10	3
3	4	11	6
4	5	12	4
5	5	13	4
6	3	14	3
7	4	15	1
8	5		

- 14** Construct two histograms from the solar-energy data in Exercise 3 of Chapter 2. Use the following interval widths and starting lower limits for the first class. What do you conclude about the choices of the interval width for this data set?

	Histogram 1	Histogram 2
Interval width	8	2
Starting lower limit	480	480

- 15** The following data were taken from a study of red-blood-cell counts before and after major surgery. Counts were taken on 23 patients, all of whom were of the same sex (female) and who had the same blood type (O+).

Count			Count		
Patient	Pre-op	Post-op	Patient	Pre-op	Post-op
1	14	0	13	5	6
2	13	26	14	4	0
3	4	2	15	15	3
4	5	4	16	4	2
5	18	8	17	0	3
6	3	1	18	7	0
7	6	0	19	2	0
8	11	3	20	8	13
9	33	23	21	4	24
10	11	2	22	4	6
11	3	2	23	5	0
12	3	2			

- (a) Construct histograms of the pre-op and the post-op blood counts. What distinguishing features, if any, are there in the distributions of the blood counts?
- (b) Make a scatter diagram of the two sets of counts. Is there an apparent relationship between the two sets of counts?
- 16** Satellite sensors can be used to provide estimates of the amounts of certain crops that are grown in agricultural regions of the United States. The following data consist of two sets of estimates of the proportions of each of 33 5×6 -nautical-mile segments of land that are growing corn during one time period during the crop season (the rest of the segment may be growing other crops or consist of roads, lakes, houses, etc.). Use the graphical techniques discussed in this chapter to assess whether these two estimation methods are providing similar information on the proportions of these segments that are growing corn.

Proportion Growing Corn			Proportion Growing Corn		
Segment	Method 1	Method 2	Segment	Method 1	Method 2
1	0.49	0.24	18	0.61	0.33
2	0.63	0.32	19	0.50	0.20
3	0.60	0.51	20	0.62	0.65
4	0.63	0.36	21	0.55	0.51
5	0.45	0.23	22	0.27	0.31
6	0.64	0.26	23	0.65	0.36
7	0.67	0.36	24	0.70	0.33
8	0.66	0.95	25	0.52	0.27
9	0.62	0.56	26	0.60	0.30
10	0.59	0.37	27	0.62	0.38
11	0.60	0.62	28	0.26	0.22
12	0.50	0.31	29	0.46	0.72
13	0.60	0.56	30	0.68	0.76
14	0.90	0.90	31	0.42	0.36
15	0.61	0.32	32	0.68	0.34
16	0.32	0.33	33	0.61	0.28
17	0.63	0.27			