

## KEY TERMS

**Average** also called mean; a number that describes the central tendency of the data

**Blinding** not telling participants which treatment a subject is receiving

**Categorical Variable** variables that take on values that are names or labels

**Cluster Sampling** a method for selecting a random sample and dividing the population into groups (clusters); use simple random sampling to select a set of clusters. Every individual in the chosen clusters is included in the sample.

**Continuous Random Variable** a random variable (RV) whose outcomes are measured; the height of trees in the forest is a continuous RV.

**Control Group** a group in a randomized experiment that receives an inactive treatment but is otherwise managed exactly as the other groups

**Convenience Sampling** a nonrandom method of selecting a sample; this method selects individuals that are easily accessible and may result in biased data.

**Cumulative Relative Frequency** The term applies to an ordered set of observations from smallest to largest. The cumulative relative frequency is the sum of the relative frequencies for all values that are less than or equal to the given value.

**Data** a set of observations (a set of possible outcomes); most data can be put into two groups: **qualitative** (an attribute whose value is indicated by a label) or **quantitative** (an attribute whose value is indicated by a number). Quantitative data can be separated into two subgroups: **discrete** and **continuous**. Data is discrete if it is the result of counting (such as the number of students of a given ethnic group in a class or the number of books on a shelf). Data is continuous if it is the result of measuring (such as distance traveled or weight of luggage)

**Discrete Random Variable** a random variable (RV) whose outcomes are counted

**Double-blinding** the act of blinding both the subjects of an experiment and the researchers who work with the subjects

**Experimental Unit** any individual or object to be measured

**Explanatory Variable** the independent variable in an experiment; the value controlled by researchers

**Frequency** the number of times a value of the data occurs

**Informed Consent** Any human subject in a research study must be cognizant of any risks or costs associated with the study. The subject has the right to know the nature of the treatments included in the study, their potential risks, and their potential benefits. Consent must be given freely by an informed, fit participant.

**Institutional Review Board** a committee tasked with oversight of research programs that involve human subjects

**Lurking Variable** a variable that has an effect on a study even though it is neither an explanatory variable nor a response variable

**Nonsampling Error** an issue that affects the reliability of sampling data other than natural variation; it includes a variety of human errors including poor study design, biased sampling methods, inaccurate information provided by study participants, data entry errors, and poor analysis.

**Numerical Variable** variables that take on values that are indicated by numbers

**Parameter** a number that is used to represent a population characteristic and that generally cannot be determined easily

**Placebo** an inactive treatment that has no real effect on the explanatory variable

**Population** all individuals, objects, or measurements whose properties are being studied

**Probability** a number between zero and one, inclusive, that gives the likelihood that a specific event will occur

**Proportion** the number of successes divided by the total number in the sample

**Qualitative Data** See [Data](#).

**Quantitative Data** See [Data](#).

**Random Assignment** the act of organizing experimental units into treatment groups using random methods

**Random Sampling** a method of selecting a sample that gives every member of the population an equal chance of being selected.

**Relative Frequency** the ratio of the number of times a value of the data occurs in the set of all outcomes to the number of all outcomes to the total number of outcomes

**Representative Sample** a subset of the population that has the same characteristics as the population

**Response Variable** the dependent variable in an experiment; the value that is measured for change at the end of an experiment

**Sample** a subset of the population studied

**Sampling Bias** not all members of the population are equally likely to be selected

**Sampling Error** the natural variation that results from selecting a sample to represent a larger population; this variation decreases as the sample size increases, so selecting larger samples reduces sampling error.

**Sampling with Replacement** Once a member of the population is selected for inclusion in a sample, that member is returned to the population for the selection of the next individual.

**Sampling without Replacement** A member of the population may be chosen for inclusion in a sample only once. If chosen, the member is not returned to the population before the next selection.

**Simple Random Sampling** a straightforward method for selecting a random sample; give each member of the population a number. Use a random number generator to select a set of labels. These randomly selected labels identify the members of your sample.

**Statistic** a numerical characteristic of the sample; a statistic estimates the corresponding population parameter.

**Stratified Sampling** a method for selecting a random sample used to ensure that subgroups of the population are represented adequately; divide the population into groups (strata). Use simple random sampling to identify a proportionate number of individuals from each stratum.

**Systematic Sampling** a method for selecting a random sample; list the members of the population. Use simple random sampling to select a starting point in the population. Let  $k = (\text{number of individuals in the population})/(\text{number of individuals needed in the sample})$ . Choose every  $k$ th individual in the list starting with the one that was randomly selected. If necessary, return to the beginning of the population list to complete your sample.

**Treatments** different values or components of the explanatory variable applied in an experiment

**Variable** a characteristic of interest for each person or object in a population

## CHAPTER REVIEW

### 1.1 Definitions of Statistics, Probability, and Key Terms

The mathematical theory of statistics is easier to learn when you know the language. This module presents important terms that will be used throughout the text.

### 1.2 Data, Sampling, and Variation in Data and Sampling

Data are individual items of information that come from a population or sample. Data may be classified as qualitative, quantitative continuous, or quantitative discrete.

Because it is not practical to measure the entire population in a study, researchers use samples to represent the population. A random sample is a representative group from the population chosen by using a method that gives each individual in the population an equal chance of being included in the sample. Random sampling methods include simple random sampling, stratified sampling, cluster sampling, and systematic sampling. Convenience sampling is a nonrandom method of choosing a sample that often produces biased data.

Samples that contain different individuals result in different data. This is true even when the samples are well-chosen and representative of the population. When properly selected, larger samples model the population more closely than smaller samples. There are many different potential problems that can affect the reliability of a sample. Statistical data needs to be critically analyzed, not simply accepted.

### 1.3 Frequency, Frequency Tables, and Levels of Measurement

Some calculations generate numbers that are artificially precise. It is not necessary to report a value to eight decimal places when the measures that generated that value were only accurate to the nearest tenth. Round off your final answer to one more decimal place than was present in the original data. This means that if you have data measured to the nearest tenth of a unit, report the final statistic to the nearest hundredth.

In addition to rounding your answers, you can measure your data using the following four levels of measurement.

- **Nominal scale level:** data that cannot be ordered nor can it be used in calculations
- **Ordinal scale level:** data that can be ordered; the differences cannot be measured
- **Interval scale level:** data with a definite ordering but no starting point; the differences can be measured, but there is no such thing as a ratio.
- **Ratio scale level:** data with a starting point that can be ordered; the differences have meaning and ratios can be calculated.

When organizing data, it is important to know how many times a value appears. How many statistics students study five hours or more for an exam? What percent of families on our block own two pets? Frequency, relative frequency, and cumulative relative frequency are measures that answer questions like these.

### 1.4 Experimental Design and Ethics

A poorly designed study will not produce reliable data. There are certain key components that must be included in every experiment. To eliminate lurking variables, subjects must be assigned randomly to different treatment groups. One of the groups must act as a control group, demonstrating what happens when the active treatment is not applied. Participants in the control group receive a placebo treatment that looks exactly like the active treatments but cannot influence the response variable. To preserve the integrity of the placebo, both researchers and subjects may be blinded. When a study is designed properly, the only difference between treatment groups is the one imposed by the researcher. Therefore, when groups respond differently to different treatments, the difference must be due to the influence of the explanatory variable.

“An ethics problem arises when you are considering an action that benefits you or some cause you support, hurts or reduces benefits to others, and violates some rule.”<sup>[4]</sup> Ethical violations in statistics are not always easy to spot. Professional associations and federal agencies post guidelines for proper conduct. It is important that you learn basic statistical procedures so that you can recognize proper data analysis.

## PRACTICE

### 1.1 Definitions of Statistics, Probability, and Key Terms

Use the following information to answer the next five exercises. Studies are often done by pharmaceutical companies to determine the effectiveness of a treatment program. Suppose that a new AIDS antibody drug is currently under study. It is given to patients once the AIDS symptoms have revealed themselves. Of interest is the average (mean) length of time in

---

4. Andrew Gelman, “Open Data and Open Methods,” Ethics and Statistics, <http://www.stat.columbia.edu/~gelman/research/published/ChanceEthics1.pdf> (accessed May 1, 2013).