

Contents

Chapter 1	The Nature of Statistics	1
Chapter 2	Organizing Data	21
Chapter 3	Descriptive Measures	127
Chapter 4	Probability Concepts	215
Chapter 5	Discrete Random Variables	279
Chapter 6	The Normal Distribution	335
Chapter 7	The Sampling Distribution of the Sample Mean	401
Chapter 8	Confidence Intervals for One Population Mean	461
Chapter 9	Hypothesis Tests for One Population Mean	513
Chapter 10	Inferences for Two Population Means	591
Chapter 11	Inferences for Population Standard Deviations	681
Chapter 12	Inferences for Population Proportions	715
Chapter 13	Chi-Square Procedures	749
Chapter 14	Descriptive Methods in Regression and Correlation	801
Chapter 15	Inferential Methods in Regression and Correlation	883
Chapter 16	Analysis of Variance (ANOVA)	953

CHAPTER 1 SOLUTIONS

Exercises 1.1

- 1.1 (a) The *population* is the collection of all individuals or items under consideration in a statistical study.
(b) A *sample* is that part of the population from which information is obtained.
- 1.2 The two major types of statistics are descriptive and inferential statistics. Descriptive statistics consists of methods for organizing and summarizing information. 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.
- 1.3 Descriptive methods are used for organizing and summarizing information and include graphs, charts, tables, averages, measures of variation, and percentiles.
- 1.4 Descriptive statistics are used to organize and summarize information from a sample before conducting an inferential analysis. Preliminary descriptive analysis of a sample may reveal features of the data that lead to the appropriate inferential method.
- 1.5 (a) An *observational study* is a study in which researchers simply observe characteristics and take measurements.
(b) A *designed experiment* is a study in which researchers impose treatments and controls and *then* observe characteristics and take measurements.
- 1.6 Observational studies can reveal only association, whereas designed experiments can help establish causation.
- 1.7 This study is inferential. Data from a sample of Americans are used to make an estimate of (or an inference about) average TV viewing time for all Americans.
- 1.8 This study is descriptive. It is a summary of the average salaries in professional baseball, basketball, and football for 2005 and 2011.
- 1.9 This study is descriptive. It is a summary of information on all homes sold in different cities for the month of September 2012.
- 1.10 This study is inferential. National samples are used to make estimates of (or inferences about) drug use throughout the entire nation.
- 1.11 This study is descriptive. It is a summary of the annual final closing values of the Dow Jones Industrial Average at the end of December for the years 2004-2013.
- 1.12 This study is inferential. Survey results were used to make percentage estimates on which college majors were in demand among U.S firms for all graduating college students.
- 1.13 (a) This study is inferential. It would have been impossible to survey all U.S. adults about their opinions on Darwinism. Therefore, the data must have come from a sample. Then inferences were made about the opinions of all U.S. adults.
(b) The population consists of all U.S. adults. The sample consists only of those U.S. adults who took part in the survey.
- 1.14 (a) The population consists of all U.S. adults. The sample consists of the 1000 U.S. adults who were surveyed.
(b) The percentage of 50% is a descriptive statistic since it describes the opinion of the U.S. adults who were surveyed.
- 1.15 (a) The statement is descriptive since it only tells what was said by the respondents of the survey.

2 Chapter 1

- (b) Then the statement would be inferential since the data has been used to provide an estimate of what all Americans believe.
- 1.16** (a) To change the study to a designed experiment, one would start with a randomly chosen group of men, then randomly divide them into two groups, an experimental group in which all of the men would have vasectomies and a control group in which the men would not have them. This would enable the researcher to make inferences about vasectomies being a cause of prostate cancer.
- (b) This experiment is not feasible, since, in the vasectomy group there would be men who did not want one, and in the control group there would be men who did want one. Since no one can be forced to participate in the study, the study could not be done as planned.
- 1.17** Designed experiment. The researchers did not simply observe the two groups of children, but instead randomly assigned one group to receive the Salk vaccine and the other to get a placebo.
- 1.18** Observational study. The researchers at Harvard University and the National Institute of Aging simply observed the two groups.
- 1.19** Observational study. The researchers simply collected data from the men and women in the study with a questionnaire.
- 1.20** Designed experiment. The researchers did not simply observe the two groups of women, but instead randomly assigned one group to receive aspirin and the other to get a placebo.
- 1.21** Designed experiment. The researchers did not simply observe the three groups of patients, but instead randomly assigned some patients to receive optimal pharmacologic therapy, some to receive optimal pharmacologic therapy and a pacemaker, and some to receive optimal pharmacologic therapy and a pacemaker-defibrillator combination.
- 1.22** Observational studies. The researchers simply collected available information about the starting salaries of new college graduates.
- 1.23** (a) This statement is inferential since it is a statement about all Americans based on a poll. We can be reasonably sure that this is the case since the time and cost of questioning every single American on this issue would be prohibitive. Furthermore, by the time everyone could be questioned, many would have changed their minds.
- (b) To make it clear that this is a descriptive statement, the new statement could be, "Of 1032 American adults surveyed, 73% favored a law that would require every gun sold in the United States to be test-fired first, so law enforcement would have its fingerprint in case it were ever used in a crime." To rephrase it as an inferential statement, use "Based on a sample of 1032 American adults, it is estimated that 73% of American adults favor a law that would require every gun sold in the United States to be test-fired first, so law enforcement would have its fingerprint in case it were ever used in a crime."
- 1.24** Descriptive statistics. The U.S. National Center for Health Statistics collects death certificate information from each state, so the rates shown reflect the causes of all deaths reported on death certificates, not just a sample.
- 1.25** (a) The population consists of all Americans between the ages of 18 and 29.
- (b) The sample consists only of those Americans who took part in the survey.
- (c) The statement in quotes is inferential since it is a statement about all Americans based on a survey.
- (d) "Based on a sample of Americans between the ages of 18 and 29, it is estimated that 59% of Americans oppose medical testing on animals."

- 1.26** (a) The \$5.36 billion lobbying expenditure figure would be a descriptive figure if it was based on the results of all lobbying expenditures during the period from 1998 through 2012.
- (b) The \$5.36 billion lobbying expenditure figure would be an inferential figure if it was an estimate based on the results of a sample of lobbying expenditures during the period from 1998 through 2012.

Exercises 1.2

- 1.27** A census is generally time consuming, costly, frequently impractical, and sometimes impossible.
- 1.28** Sampling and experimentation are two alternative ways to obtain information without conducting a complete census.
- 1.29** The sample should be representative so that it reflects as closely as possible the relevant characteristics of the population under consideration.
- 1.30** There are many possible answers. Surveying people regarding political candidates as they enter or leave an upscale business location, surveying the readers of a particular publication to get information about the population in general, polling college students who live in dormitories to obtain information of interest to all students are all likely to produce samples unrepresentative of the population under consideration.
- 1.31** (a) Probability sampling consists of using a randomizing device such as tossing a coin or consulting a random number table to decide which members of the population will constitute the sample.
- (b) No. It is possible for the randomizing device to randomly produce a sample that is not representative.
- (c) Probability sampling eliminates unintentional selection bias, permits the researcher to control the chance of obtaining a non-representative sample, and guarantees that the techniques of inferential statistics can be applied.
- 1.32** (a) Simple random sampling is a procedure for which each possible sample of a given size is equally likely to be the one obtained.
- (b) A simple random sample is one that was obtained by simple random sampling.
- (c) Random sampling may be done with or without replacement. In sampling with replacement, it is possible for a member of the population to be chosen more than once, i.e., members are eligible for re-selection after they have been chosen once. In sampling without replacement, population members can be selected at most once.
- 1.33** Simple random sampling.
- 1.34** One method would be to place the names of all members of the population under consideration on individual slips of paper, place the slips in a container large enough to allow them to be thoroughly shuffled by shaking or spinning, and then draw out the desired number of slips for the sample while blindfolded. A second method, which is much more practical when the population size is large, is to assign a number to each member of the population, and then use a random number table, random number generating device, or computer program to determine the numbers of those members of the population who are chosen.
- 1.35** The acronym used for simple random sampling without replacement is SRS.
- 1.36** (a) 123, 124, 125, 134, 135, 145, 234, 235, 245, 345

4 Chapter 1

- (b) There are 10 samples, each of size three. Each sample has a one in 10 chance of being selected. Thus, the probability that a sample of three is 1, 3, and 5 is $1/10$.
- (c) Starting in Line 05 and column 20, reading single digit numbers down the column and then up the next column, the first digit that is a one through five is a 5. Ignoring duplicates and skipping digits 6 and above and also skipping zero, the second digit found that is a one through five is a 4. Continuing down column 20 and then up column 21, the third digit found that is a one through five is a 1. Thus the SRS of 1, 4, and 5 is obtained.

1.37 (a) 12, 13, 14, 23, 24, 34

- (b) There are 6 samples, each of size two. Each sample has a one in six chance of being selected. Thus, the probability that a sample of two is 2 and 3 is $1/6$.
- (c) Starting in Line 17 and column 07 (notice there is a column 00), reading single digit numbers down the column and then up the next column, the first digit that is a one through four is a 1. Continue down column 07 and then up column 08. Ignoring duplicates and skipping digits 5 and above and also skipping zero, the second digit found that is a one through four is a 4. Thus the SRS of 1 and 4 is obtained.

1.38 (a) Starting in Line 15 and reading two digits numbers in columns 25 and 26 going down the table, the first two digit number between 01 and 90 is 06. Continuing down the columns and ignoring duplicates and numbers 91-99, the next two numbers are 33 and 61. Then, continuing up columns 27 and 28, the last two numbers selected are 56 and 20. Therefore the SRS of size five consists of observations 06, 33, 61, 56, and 20.

(b) There are many possible answers.

1.39 (a) Starting in Line 10 and reading two digits numbers in columns 10 and 11 going down the table, the first two digit number between 01 and 50 is 43. Continuing down the columns and ignoring duplicates and numbers 51-99, the next two numbers are 45 and 01. Then, continuing up columns 12 and 13, the last three numbers selected are 42, 37, and 47. Therefore the SRS of size six consists of observations 43, 45, 01, 42, 37, and 47.

(b) There are many possible answers.

1.40 The online poll clearly has a built-in non-response bias. Since it was taken over the Memorial Day weekend, most of those who responded were people who stayed at home and had access to their computers. Most people vacationing outdoors over the weekend would not have carried their computers with them and would not have been able to respond.

1.41 Dentists form a high-income group whose incomes are not representative of the incomes of Seattle residents in general.

1.42 (a) The five possible samples of size one are G, L, S, A, and T.

(b) There is no difference between obtaining a SRS of size 1 and selecting one official at random.

(c) The one possible sample of size five is GLSAT.

(d) There is no difference between obtaining a SRS of size 5 and taking a census of the five officials.

1.43 (a) GLS, GLA, GLT, GSA, GST, GAT, LSA, LST, LAT, SAT.

(b) There are 10 samples, each of size three. Each sample has a one in 10 chance of being selected. Thus, the probability that a sample of three officials is the first sample on the list presented in part (a) is $1/10$. The same is true for the second sample and for the tenth sample.

- 1.44** (a) E,M E,A M,L P,L L,A
E,P E,B M,A P,A L,B
E,L M,P M,B P,B A,B
- (b) One procedure for taking a random sample of two representatives from the six is to write the initials of the representatives on six separate pieces of paper, place the six slips of paper into a box, and then, while blindfolded, pick two of the slips of paper. Or, number the representatives 1-6, and use a table of random numbers or a random-number generator to select two different numbers between 1 and 6.
- (c) $1/15$; $1/15$
- 1.45** (a) E,M,P,L E,M,L,B E,P,A,B M,P,A,B
E,M,P,A E,M,A,B E,L,A,B M,L,A,B
E,M,P,B E,P,L,A M,P,L,A P,L,A,B
E,M,L,A E,P,L,B M,P,L,B
- (b) One procedure for taking a random sample of four representatives from the six is to write the initials of the representatives on six separate pieces of paper, place the six slips of paper into a box, and then, while blindfolded, pick four of the slips of paper. Or, number the representatives 1-6, and use a table of random numbers or a random-number generator to select four different numbers between 1 and 6.
- (c) $1/15$; $1/15$
- 1.46** (a) E,M,P E,P,A M,P,L M,A,B
E,M,L E,P,B M,P,A P,L,A
E,M,A E,L,A M,P,B P,L,B
E,M,B E,L,B M,L,A P,A,B
E,P,L E,A,B M,L,B L,A,B
- (b) One procedure for taking a random sample of three representatives from the six is to write the initials of the representatives on six separate pieces of paper, place the six slips of paper into a box, and then, while blindfolded, pick three of the slips of paper. Or, number the representatives 1-6, and use a table of random numbers or a random-number generator to select three different numbers between 1 and 6.
- (c) $1/20$; $1/20$
- 1.47** (a) F,T F,G F,H F,L F,B F,A
T,G T,H T,L T,B T,A G,H
G,L G,B G,A H,L H,B H,A
L,B L,A B,A
- (b) $1/21$; $1/21$
- 1.48** (a) I am using Table I to obtain a list of 20 different random numbers between 1 and 80 as follows.
- I start at the two digit number in line number 5 and column numbers 31-32, which is the number 86. Since I want numbers between 1 and 80 only, I throw out numbers between 81 and 99, inclusive. I also discard the number 00.
- I now go down the table and record the two-digit numbers appearing directly beneath 86.
- After skipping 86, I record 39, 03, skip 97, record 28, 58, 59, skip 81, record 09, 36, skip 81, record 52, skip 94, record 24 and 78.

6 Chapter 1

Now that I've reached the bottom of the table, I move directly rightward to the adjacent column of two-digit numbers and go up.

I skip 84, record 57, 40, skip 89, record 69, 25, skip 95, record 51, 20, 42, 77, skip 89, skip 40(duplicate), record 14, and 34.

I've finished recording the 20 random numbers. In summary, these are

39	03	28	58	59
09	36	52	24	78
57	40	69	25	51
20	42	77	14	34

- (b) We can use Minitab to generate random numbers. Following the instructions in *The Technology Center*, our results are 55, 47, 66, 2, 72, 56, 10, 31, 5, 19, 39, 57, 44, 60, 23, 34, 43, 9, 49, and 62. Your result may be different from ours.

- 1.49** (a) I am using Table I to obtain a list of 10 random numbers between 1 and 500 as follows.

I start at the three digit number in line number 14 and column numbers 10-12, which is the number 452.

I now go down the table and record the three-digit numbers appearing directly beneath 452. Since I want numbers between 1 and 500 only, I throw out numbers between 501 and 999, inclusive. I also discard the number 000.

After 452, I skip 667, 964, 593, 534, and record 016.

Now that I've reached the bottom of the table, I move directly rightward to the adjacent column of three-digit numbers and go up.

I record 343, 242, skip 748, 755, record 428, skip 852, 794, 596, record 378, skip 890, record 163, skip 892, 847, 815, 729, 911, 745, record 182, 293, and 422.

I've finished recording the 10 random numbers. In summary, these are:

452	016	343	242	428
378	163	182	293	422

- (b) We can use Minitab to generate random numbers. Following the instructions in *The Technology Center*, our results are 489, 451, 61, 114, 389, 381, 364, 166, 221, and 437. Your result may be different from ours.

- 1.50** (a) First assign the digits 0 through 9 to the ten cities as listed in the exercise. Select a random starting point in Table I of Appendix A and read in a pre-selected direction until you have encountered 5 different digits. For example, if we start at the top of the fifth column of digits and read down, we encounter the digits 4,1,5,2,5,6. We ignore the second '5'. Thus our sample of five cities consists of Osaka, Tokyo, Miami, San Francisco, and New York. Your answer may be different from this one.

- (b) We can use Minitab to generate random numbers. Following the instructions in *The Technology Center*, our results are 3, 8, 6, 5, 9. Thus our sample of 5 cities is Los Angeles, Manila, New York, Miami, and London. Your result may be different from ours.

- 1.51** (a) First re-assign the elements 93 through 118 as elements 01 to 26.

Select a random starting point in Table I of Appendix A and read in a pre-selected direction until you have encountered 8 different elements.

For example, if we start at the top of the column 10 and read two digit numbers down and then up in the following columns, we encounter

the elements 04, 01, 03, 08, 11, 18, 22, and 15. This corresponds to a sample of the elements Cm, Np, Am, Fm, Lr, Ds, Fl, and Bh. Your answer may be different from this one.

- (b) We can use Minitab to generate random numbers. Following the instructions in The Technology Center, our results are 8, 2, 9, 20, 24, 19, 21, and 13. Thus our sample of 8 elements is Fm, Pu, Md, Cn, Lv, Rg, Uut, and Db. Your result may be different from ours.
- 1.52** (a) One of the biggest reasons for undercoverage in household surveys is that respondents do not correctly indicate all who are living in a household maybe due to deliberate concealment or irregular household structure or living arrangements. The household residents are only partially listed.
- (b) A telephone survey of Americans from a phone book will likely have bias due to undercoverage because many people have unlisted phone numbers and also it is becoming more popular that many people do not even have home phones. This would cause the phone book to be an incomplete list of the population.
- 1.53** (a) One of the dangers of nonresponse is that the individuals who do not respond may have a different observed value than the individuals that do respond causing a nonresponse bias in the estimate. Nonresponse bias may make the measured value too small or too large.
- (b) The lower the response rate, the more likely there is a nonresponse bias in the estimate. Therefore the estimate will either under or over estimate the generalized results to the entire population.
- 1.54** (a) The respondent may wish to please the questioner by answering what is morally or legally right. The respondent might not be willing to admit to the questioner that they smoke marijuana and the measured value of the percentage of people that smoke marijuana would then be underestimated due to response bias.
- (b) Another situation that might be conducive to response bias is perhaps a woman questioning men on their opinion of domestic violence, or an environmentalist questioning people on their recycling habits.
- (c) The wording of a question could lead to response bias. Whether the survey is anonymous or not could lead to response bias. The characteristics of the questioner could lead to response bias. It could also happen if the questioner obviously favors and is pushing for one particular answer.

Exercises 1.3

- 1.55** Systematic random sampling is easier to execute than simple random sampling and usually provides comparable results. The exception is the presence of some kind of cyclical pattern in the listing of the members of the population.
- 1.56** Ideally, in cluster sampling, each cluster should pattern the entire population.
- 1.57** Ideally, in stratified sampling, the members of each stratum should be homogeneous relative to the characteristic under consideration.
- 1.58** Surveys that combine one or more of simple random sampling, systematic random sampling, cluster sampling, and stratified sampling employ what is called multistage sampling.
- 1.59** (a) Answers will vary, but here is the procedure: (1) Divide the population size, 372, by the sample size, 5, and round down to the nearest whole number if necessary; this gives 74. Use a table of random numbers (or a similar device) to select a number between 1 and 74, call it k . (3) List every 74th number, starting with k , until 5 numbers are obtained;

8 Chapter 1

thus, the first number of the required list of 5 numbers is k , the second is $k + 74$, the third is $k + 148$, and so forth.

- (b) Following part (a) with $k = 10$, the first number of the sample is 10, the second is $10 + 74 = 84$. The remaining three numbers in the sample would be 158, 232, and 306. Thus, the sample of 5 would be 10, 84, 158, 232, and 306.

1.60 (a) Answers will vary, but here is the procedure: (1) Divide the population size, 500, by the sample size, 9, and round down to the nearest whole number if necessary; this gives 55. Use a table of random numbers (or a similar device) to select a number between 1 and 55, call it k . (3) List every 55th number, starting with k , until 9 numbers are obtained; thus, the first number of the required list of 9 numbers is k , the second is $k + 55$, the third is $k + 110$, and so forth.

- (b) Following part (a) with $k = 48$, the first number of the sample is 48, the second is $48 + 55 = 103$. The remaining seven numbers in the sample would be 158, 213, 268, 323, 378, 433, and 488. Thus, the sample of 9 would be 48, 103, 158, 213, 268, 323, 378, 433, and 488.

1.61 (a) Answers will vary, but here is the procedure: (1) The population of size 50 is already divided into five clusters of size 10. (2) Since the required sample size is 20, we will need to take a SRS of 2 clusters. Use a table of random numbers (or a similar device) to select two numbers between 1 and 5. These are the two clusters that are selected. (3) Use all the members of each cluster selected in part (2) as the sample.

- (b) Following part (a) with clusters #1 and #3 selected, we would select all the members in cluster 1, which are 1 - 10, and all the members in cluster 3, which are 21 - 30.

1.62 (a) Answers will vary, but here is the procedure: (1) The population of size 100 is already divided into ten clusters of size 10. (2) Since the required sample size is 30, we will need to take a SRS of 3 clusters. Use a table of random numbers (or a similar device) to select three numbers between 1 and 10. These are the three clusters that are selected. (3) Use all the members of each cluster selected in part (2) as the sample.

- (b) Following part (a) with clusters #2, #6, and #9 selected, we would select all the members in cluster 2 (11-20), all the members in cluster 6 (51-60), and all the members in cluster 9 (81-90). Therefore, our sample would consist of 11-20, 51-60, and 81-90.

1.63 (a) From each strata, we need to obtain a SRS of a size proportional to the size of the stratum. Therefore, since strata #1 is 30% of the population, a SRS equal to 30% of 20, or 6, should be sampled from strata #1. Since strata #2 is 20% of the population, a SRS equal to 20% of 20, or 4, should be sampled from strata #2. Similarly, a SRS of size 8 should be sampled from strata #3 and a SRS of size 2 should be sampled from strata #4. The sample sizes from stratum #1 through #4 are 6, 4, 8, and 2 respectively.

- (b) Answers will vary following the procedure in part (a).

1.64 (a) From each strata, we need to obtain a SRS of a size proportional to the size of the stratum. Therefore, since strata #1 is 40% of the population, a SRS equal to 40% of 10, or 4, should be sampled from strata #1. Since strata #2 is 30% of the population, a SRS equal to 30% of 10, or 3, should be sampled from strata #2. Similarly, a SRS of size 3 should be sampled from strata #3. The sample sizes from stratum #1 through #3 are 4, 3, and 3 respectively.

- (b) Answers will vary following the procedure in part (a).

- 1.65** Stratified Sampling. The entire population is naturally divided into subpopulations, one from each lake, and random sampling is done from each lake. The stratified sampling is not with proportional allocation since that would require knowing how many fish were in each lake.
- 1.66** Stratified Sampling. The entire population is naturally divided into four subpopulations, and random sampling is done from each and then combined into a single sample.
- 1.67** Systematic Random Sampling. Kennedy selected his sample using the fixed periodic interval of every 50th letter, which is the similar to the method presented in procedure 1.1.
- 1.68** Cluster Sampling. The clusters of this sampling design are the 1285 journals. A random sample of 26 clusters was selected and then all articles from the selected journals for a particular year were examined.
- 1.69** Cluster Sampling. The clusters of this sampling design are the 46 schools. A random sample of 10 clusters was selected and then all of the parents of the nonimmunized children at the 10 selected schools were sent a questionnaire.
- 1.70** Systematic Random Sampling. This sampling design follows procedure 1.1. First, dividing the population size of 8493 by 30, they arrived at $k = 283$. Then, the randomly selected starting point was $m = 10$. Then, the sampled stickers were $m = 10$, $m + k = 293$, $m + 2k = 576$, etc.
- 1.71** (a) Answers will vary, but here is the procedure: (1) Divide the population size, 500, by the sample size, 10, and round down to the nearest whole number if necessary; this gives 50. (2) Use a table of random numbers (or a similar device) to select a number between 1 and 50, call it k . (3) List every 50th, starting with k , until 10 numbers are obtained; thus, the first number on the required list of 10 numbers is k , the second is $k+50$, the third is $k+100$, and so forth (e.g., if $k=6$, then the numbers on the list are 6, 56, 106, ...).
- (b) Systematic random sampling is easier.
- (c) The answer depends on the purpose of the sampling. If the purpose of sampling is not related to the size of the sales outside the U.S., systematic sampling will work. However, since the listing is a ranking by amount of sales, if k is low (say 2), then the sample will contain firms that, on the average, have higher sales outside the U.S. than the population as a whole. If the k is high, (say 49) then the sample will contain firms that, on the average, have lower sales than the population as a whole. In either of those cases, the sample would not be representative of the population in regard to the amount of sales outside the U.S.
- 1.72** (a) Answers will vary, but here is the procedure: (1) Divide the population size, 80, by the sample size, 20, and round down to the nearest whole number if necessary; this gives 4. (2) Use a table of random numbers (or a similar device) to select a number between 1 and 4, call it k . (3) List every 4th number, starting with k , until 20 numbers are obtained; thus the first number on the required list of 20 numbers is k , the second is $k+4$, the third is $k+8$, and so forth (e.g., if $k=3$, then the numbers on the list are 3, 7, 11, 15, ...).
- (b) Systematic random sampling is easier.
- (c) No. In Keno, you want every set of 20 balls to have the same chance of being chosen. Systematic sampling would give each of 4 sets of balls [(1, 5, 9, ..., 77), (2, 6, 10, ..., 78), (3, 7, 11, ..., 79) and (4, 8, 12, ..., 80)], a 1/4 chance of occurring, while all of the other possible sets of balls would have no chance of occurring.

10 Chapter 1

- 1.73** (a) Number the suites from 1 to 48, use a table of random numbers to randomly select three of the 48 suites, and take as the sample the 24 dormitory residents living in the three suites obtained.
- (b) Probably not, since friends are more likely to have similar opinions than are strangers.
- (c) There are 384 students in total. Freshmen make up $1/3$ of them. Sophomores make up $7/24$ of them, Juniors $1/4$, and Seniors $1/8$. Multiplying each of these fractions by 24 yields the proportional allocation, which dictates that the number of freshmen, sophomores, juniors, and seniors selected should be, respectively, 8, 7, 6, and 3. Thus a stratified sample of 24 dormitory residents can be obtained as follows: Number the freshmen dormitory residents from 1 to 128 and use a table of random numbers to randomly select 8 of the 128 freshman dormitory residents; number the sophomore dormitory residents from 1 to 112 and use a table of random numbers to randomly select 7 of the 112 sophomore dormitory residents; and so forth.
- 1.74** (a) Each category of "Percent free lunch" should be represented in the sample in the same proportion that it is present in the population of top 100 ranked high schools. Thus 50/100 of the sample of 25 schools should be from the 0 to under 10% free lunch category, 18/100 from the second category, 11/100 from the third, 8/100 from the fourth, and 13/100 from the last. Multiplying each of these fractions by 25 gives us the sample sizes from each category. These sample sizes will not necessarily be integers, so we will need to make some minor adjustments of the results. The first category should have $(50/100)(25) = 12.5$. The second should have $(18/100)(25) = 4.5$. Similarly, the third, fourth, and fifth categories should have 2.75, 2, and 3.25 for their sample sizes. We round the third and fifth sample sizes each to 3. After flipping a coin, we round the first two categories to 12 and 5. Thus the sample sizes for the five Percent free lunch categories should be 12, 5, 3, 2, and 3 respectively. We would now use a random number generator to select 12 out of the 50 in the first category, 5 out of the 18 in the second, 3 out of the 11 in the third, 2 of the 8 in the fourth, and 3 of the 13 in the last category.
- (b) From part (a), two schools would be selected from the strata with a percent free lunch value of 30-under 40.
- 1.75** (a) Answers will vary, but here is the procedure: (1) Divide the population size, 435, by the sample size, 15, and round down to the nearest whole number if necessary; this gives 29. Use a table of random numbers (or a similar device) to select a number between 1 and 29, call it k . (3) List every 29th number, starting with k , until 15 numbers are obtained; thus, the first number of the required list of 15 numbers is k , the second is $k + 29$, the third is $k + 58$, and so forth.
- (b) Following part (a) with $k = 12$, the first number of the sample is 12, the second is $12 + 29 = 41$. The third number selected is $12 + 58 = 70$. The remaining twelve numbers are similarly selected. Thus, the sample of 15 would be 12, 41, 70, 99, 128, 157, 186, 215, 244, 273, 302, 331, 360, 389, and 418.
- 1.76** (a) Each category of "Region" should be represented in the sample in the same proportion that it is present in the population. Thus 43% of the sample of 50 should be volunteers serving in Africa, 21% from Latin America, 15% from Eastern Europe/Central Asia, 10% from Asia, 4% from the Caribbean, 4% from North Africa/Middle East, and 3% from the Pacific Island. Finding each of these proportions of 50 gives us the sample sizes from each category. These sample sizes will not necessarily be integers, so we will need to make some minor adjustments of the results. Volunteers from Africa should have $(0.43)(50) = 21.5$. Volunteers from Latin America should have $(0.21)(50) = 10.5$.

Similarly, the remaining categories should have 7.5, 5, 2, 2, and 1.5 for their sample sizes. After flipping a coin, we round the first two categories either up or down. Thus the sample sizes for the categories should be 21, 11, 7, 5, 2, 2, and 2 respectively. We would now use a random number generator to select the volunteers from each category.

- (b) From part (a), two volunteers would be selected from the strata with volunteers serving in the Caribbean.
- 1.77** (a) This is a poll taken by calling randomly selected U.S. adults. Thus, the sampling design appears to be simple random sampling, although it is possible that a more complex design was used to ensure that various political, religious, educational, or other types of groups were proportionately represented in the sample.
- (b) The sample size for the second question was 78% of 1010 or 788.
- (c) The sample size for the third question was 28% of 788 or 221.
- 1.78** No. In your text, Example 1.10, only 48 different samples are possible. A sample containing students 5,6, and 7 is not possible at all. While the 48 possible samples are equally likely, there are other samples that could be obtained through simple random sampling that are not possible at all in systematic sampling. Thus not all possible samples are equally likely. Nevertheless, if there is no pattern or cycle to the data, this method will tend to give about the same results as simple random sampling.
- 1.79** (a) It is also true for systematic random sampling if the population size divided by the sample size results in an integer for m . The chance for each member to be selected is then still equal to the sample size divided by the population size. For example, suppose the population size is $N=10$ and the sample size is $n=2$. The chance that each member in simple random sampling to be selected is $2/10 = 1/5$. In systematic random sampling for the same example, $m=5$. The possible samples of size two are 1 and 6, 2 and 7, 3 and 8, 4 and 9, and 5 and 10. Therefore, the chance that a member is selected is equal to the chance of one of those five samples being selected, which is the same as simple random sampling of $1/5$.
- (b) It is not true for systematic random sampling if the population size divided by the sample size does not result in an integer for m . For example, suppose the population size is $N=15$ and the sample size is $n=2$. After dividing the population size by the sample size and rounding down to the nearest whole number, we get $m=7$. You would select every 7th member after a random starting place k , between 1 and 7, is determined. If $k=1$, you would select the first and eighth member. If $k=7$, you would select the seventh and fourteenth member. In this situation, the last member (fifteenth) can never be selected. Therefore, the last member of the sample does not have the same chance of being selected as any other member in the population.
- 1.80** Refer to example 1.14. If we approached this problem as a simple random sample each member would have a chance of being selected equal to the sample size divided by the population size: $20/250$, or $2/25$.

If we approached this same example as a stratified sample with proportional allocation, we would select 2 out of 25 households in the upper income group, 14 out of the 175 households in the middle income group, and 4 out of 50 households in the lower income group. Thus the chance that an upper income household is selected is $2/25$. The chance that a middle income household is selected is $14/175 = 2/25$. Finally, the chance that a lower income household is selected is $4/50 = 2/25$. Thus, the chance that each member is selected is the same as a simple random sample.

12 Chapter 1

Exercises 1.4

1.81 (a) Experimental units are the individuals or items on which the experiment is performed.

(b) When the experimental units are humans, we call them subjects.

1.82 The three basic principles of experimental design are 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 differences among the treatments.

1.83 (a) The response variable is the characteristic of the experimental outcome that is to be measured or observed.

(b) A factor is a variable whose effect on the response variable is of interest in the experiment.

(c) The levels are the possible values of the factor.

(d) For a one-factor experiment, the treatments are the levels of the factor. For multifactor experiments, the treatments are the combinations of levels of the factors.

1.84 One type of statistical design is a completely randomized design. In a completely randomized design, all the experimental units are assigned randomly among all the treatments. The second type of statistical design is a randomized block design. In a randomized block design, the experimental units are assigned randomly among all the treatments separately within each block.

1.85 In a one-factor experiment, the number of treatments is equal to the number of levels of the factor. Therefore, there are four treatments.

1.86 In a one-factor experiment, the number of treatments is equal to the number of levels of the factor. Therefore, there are five treatments.

1.87 (a)

		B			
		b_1	b_2	b_3	b_4
A	a_1	a_1b_1	a_1b_2	a_1b_3	a_1b_4
	a_2	a_2b_1	a_2b_2	a_2b_3	a_2b_4
	a_3	a_3b_1	a_3b_2	a_3b_3	a_3b_4

(b) There are twelve combinations of the levels of the factors. Therefore, there are twelve treatments.

(c) Yes, you could have multiplied the number of levels in each factor. There are three levels of factor A and four levels of factor B. Therefore, there are $(3)(4) = 12$ treatments.

1.88 (a)

		B	
		b_1	b_2
A	a_1	a_1b_1	a_1b_2
	a_2	a_2b_1	a_2b_2
	a_3	a_3b_1	a_3b_2
	a_4	a_4b_1	a_4b_2

- (b) There are eight combinations of the levels of the factors. Therefore, there are eight treatments.
- (c) Yes, you could have multiplied the number of levels in each factor. There are four levels of factor A and two levels of factor B. Therefore, there are $(4)(2) = 8$ treatments.
- 1.89** You can multiply the number of levels in each factor. There are m levels in the first factor and n levels in the second factor. Therefore, there are $(m)(n) = m \times n$ treatments.
- 1.90** (a) The treatment group consisted of the 2444 patients who took Prozac.
 (b) The control group consisted of the 1331 patients who received a placebo.
 (c) The treatments were administering Prozac and administering the placebo.
- 1.91** (a) There were three treatments.
 (b) The first group, the one receiving only the pharmacologic therapy, would be considered the control group.
 (c) There were three treatment groups. The first received only pharmacologic therapy, the second received pharmacologic therapy plus a pacemaker, and the third received pharmacologic therapy plus a pacemaker-defibrillator combination.
 (d) The first group (control) contained $1/5$ of the 1520 patients or 304. The other two groups each contained $2/5$ of the 1520 patients or 608.
 (e) Each patient could be randomly assigned a number from 1 to 1520. Any patient assigned a number between 1 and 304 would be assigned to the control group; any patient assigned to the next 608 numbers (305 to 912) would be assigned to receive the pharmacologic therapy plus a pacemaker; and any patient assigned a number between 913 and 1520 would receive pharmacologic therapy plus a pacemaker-defibrillator combination. Each random number would be used only once to ensure that the resulting treatment groups were of the intended sizes.
- 1.92** (a) Experimental units: batches of the product being sold
 (b) Response variable: the number of units of the product sold
 (c) Factors: two factors - display type and pricing scheme
 (d) Levels of each factor: three types of display of the product and three pricing schemes
 (e) Treatments: the nine different combinations of display type and price resulting from testing each of the three pricing schemes with each of the three display types
- 1.93** (a) Experimental units: the drivers
 (b) Response variable: the detection distance, in feet
 (c) Factors: two factors - sign size and sign material

14 Chapter 1

- (d) Levels of each factor: three levels of sign size (small, medium, and large) and three levels of sign material (1, 2, and 3)
- (e) Treatments: the nine different combinations of sign size and sign material resulting from testing each of the three sign sizes with each of the three sign materials

- 1.94**
- (a) Experimental units: fields of oats
 - (b) Response variable: crop yield of the oats per acre
 - (c) Factors: variety of oats and concentration of manure on the fields
 - (d) Levels of each factor: three varieties of oats and four concentrations of manure
 - (e) Treatments: the twelve combinations of oat variety and manure concentration resulting from testing each of the three oat varieties with each of the four concentration levels of the manure
- 1.95**
- (a) Experimental units: female lions
 - (b) Response variable: whether or not the female lion approached the male lion dummy
 - (c) Factors: length and color of the mane on the male lion dummy
 - (d) Levels of each factor: two different mane lengths and two different mane colors
 - (e) Treatments: the four combinations of mane length and color
- 1.96**
- (a) Experimental units: the women in the study
 - (b) Response variable: the color of the shirt chosen
 - (c) Factors: gender and attractiveness of the new acquaintance
 - (d) Levels of each factor: two different genders (male, female) and two different levels of attractiveness (attractive, unattractive)
 - (e) Treatments: the four combinations of gender and attractiveness
- 1.97**
- (a) Experimental units: the children
 - (b) Response variable: IQ score
 - (c) Factor: Whether they were given dexamethasone (control or dexamethasone group)
 - (d) Levels of each factor: two levels of the single factor (control or dexamethasone group)
 - (e) Treatments: the two levels of the single factor
- 1.98**
- (a) This is a completely randomized design since the flashlights were randomly assigned to the different battery brands.
 - (b) This is a randomized block design since the four different battery brands would be randomly assigned within each set of four flashlights from each of the five flashlight brands.
- 1.99**
- (a) This is a randomized block design. The experiment first blocked by gender. All the experimental units are not randomly assigned among all the treatments.
 - (b) The blocks are the two genders (male and female).
- 1.100** Double-blinding guards against bias, both in the evaluations and in the responses. In the Salk vaccine experiment, double-blinding prevented a doctor's evaluation from being influenced by knowing which treatment (vaccine or placebo) a patient received; it also prevented a patient's response to the treatment from being influenced by knowing which treatment he or she received.

- 1.101 (a) Simple random sampling corresponds to completely randomized designs since selection is randomly made from the entire population.
- (b) Stratified sampling corresponds to randomized block designs since selection is randomly made from within each strata.

Review Problems for Chapter 1

1. Student exercise.
2. Descriptive statistics are used to display and summarize the data to be used in an inferential study. Preliminary descriptive analysis of a sample often reveals features of the data that lead to the choice or reconsideration of the choice of the appropriate inferential analysis procedure.
3. (a) An *observational study* is a study in which researchers simply observe characteristics and take measurements.
- (b) A *designed experiment* is a study in which researchers impose treatments and controls and *then* observe characteristics and take measurements.
4. A literature search should be made before planning and conducting a study.
5. (a) A representative sample is one that reflects as closely as possible the relevant characteristics of the population under consideration.
- (b) Probability sampling involves the use of a randomizing device such as tossing a coin or die, using a random number table, or using computer software that generates random numbers to determine which members of the population will make up the sample.
- (c) A sample is a simple random sample if all possible samples of a given size are equally likely to be the actual sample selected.
6. (a) This method does not involve probability sampling. No randomizing device is being used and people who do not visit the campus cafeteria have no chance of being included in the sample.
- (b) The dart throwing is a randomizing device that makes all samples of size 20 equally likely. This is probability sampling.
7. (a) Systematic random sampling is done by first dividing the population size by the sample size and rounding the result down to the next integer, say m . Then we select one random number, say k , between 1 and m inclusive. That number will be the first member of the sample. The remaining members of sample will be those numbered $k+m$, $k+2m$, $k+3m$, ... until a sample of size n has been chosen. Systematic sampling will yield results similar to simple random sampling as long as there is nothing systematic about the way the members of the population were assigned their numbers.
- (b) In cluster sampling, clusters of the population (such as blocks, precincts, wards, etc.) are chosen at random from all such possible clusters. Then every member of the population lying within the chosen clusters is sampled. This method of sampling is particularly convenient when members of the population are widely scattered and is most appropriate when the members of each cluster are representative of the entire population. Cluster sampling can save both time and expense in doing the survey, but can yield misleading results if individual clusters are made up of subjects with very similar views on the topic being surveyed.
- (c) In stratified random sampling with proportional allocation, the population is first divided into subpopulations, called strata, and simple random sampling is done within each stratum. Proportional allocation means that the size of the sample from each stratum is proportional to the size of the population in that stratum. This type

16 Chapter 1

of sampling may improve the accuracy of the survey by ensuring that those in each stratum are more proportionately represented than would be the case with cluster sampling or even simple random sampling. Ideally, the members of each stratum should be homogeneous relative to the characteristic under consideration. If they are not homogeneous within each stratum, simple random sampling would work just as well.

8. The three basic principles of experimental design are control, randomization, and replication. Control refers to methods for controlling factors other than those of primary interest. Randomization means randomly dividing the subjects into groups in order to avoid unintentional selection bias in constituting the groups. Replication means using enough experimental units or subjects so that groups resemble each other closely and so that there is a good chance of detecting differences among the treatments when such differences actually exist.
9. Descriptive study. It is a summary of the scores of major league baseball games on August 14, 2013.
10. (a) Descriptive study. It is a summary of the responses from those that participated in the poll.
(b) Inferential statement. It is an implied estimate of the responses of all adults in the U.S.
11. Inferential study. The results of a sample are used to make inferences about the age distribution of all British backpackers in South Africa.
12. (a) Descriptive study. It is a summary of the percentages of Jewish children sampled in Israel and Britain that have peanut allergies.
(b) Observational study. The researchers simply observed the two groups.
13. This is an observational study. To be a designed experiment, the researchers would have to have the ability to assign some children at random to live in persistent poverty during the first 5 years of life or to not suffer any poverty during that period. Clearly that is not possible.
14. This is a designed experiment since the researcher is imposing a treatment and then observing the results.
15. Because Yale is a very expensive school, incomes of parents of Yale students will not be representative of the incomes of all college students' parents.
16. (a)

H,Z,C	H,Z,A	H,Z,J	H,C,A	H,C,J
H,A,J	Z,C,A	Z,C,J	Z,A,J	C,A,J

(b) Since each of the 10 samples of size three is equally likely, there is a 1/10 chance that the sample chosen is the first sample in the list, 1/10 chance that it is the second sample in the list, and 1/10 chance that it is the tenth sample in the list.
(c) (i) Make five slips of paper with each airline on one slip. Draw three slips at random. (ii) Make 10 slips of paper, each having one of the combinations in part (a). Draw one slip at random. (iii) Number the five airlines from 1 to 5. Use a random number table or random number generator to obtain three distinct random numbers between 1 and 5, inclusive.
(d) Your method and result may differ from ours. We rolled a die (ignoring 6's and duplicates) and got 2, 5, 2, 6, 4. Ignoring duplicates and numbers greater than five, our sample consists of Horizon, Jazz, and Alaska Airlines.
17. (a) Table I can be employed to obtain a sample of 15 random numbers between 1 and 100 as follows. First, I pick a random starting point by closing my eyes and putting my finger down on the table.

My finger falls on three digits located at the intersection of a line with three columns. (Notice that the first column of digits is labeled "00" rather than "01".) This is my starting point.

I now go down the table and record all three-digit numbers appearing directly beneath the first three-digit number that are between 001 and 100 inclusive. I throw out numbers between 101 and 999, inclusive. I also discard the number 0000. When the bottom of the column is reached, I move over to the next sequence of three digits and work my way back up the table. Continue in this manner. When 10 distinct three-digit numbers have been recorded, the sample is complete.

- (b) Starting in row 10, columns 7-9, we skip 484, 797, record 082, skip 586, 653, 452, 552, 155, record 008, skip 765, move to the right and record 016, skip 534, 593, 964, 667, 452, 432, 594, 950, 670, record 001, skip 581, 577, 408, 948, 807, 862, 407, record 047, skip 977, move to the right, skip 422 and all of the rest of the numbers in that column, move to the right, skip 732, 192, record 094, skip 615 and all of the rest of the numbers in that column, move to the right, record 097, skip 673, record 074, skip 469, 822, record 052, skip 397, 468, 741, 566, 470, record 076, 098, skip 883, 378, 154, 102, record 003, skip 802, 841, move to the right, skip 243, 198, 411, record 089, skip 701, 305, 638, 654, record 041, skip 753, 790, record 063.

The final list of numbers is 82, 8, 16, 1, 47, 94, 97, 74, 52, 76, 98, 3, 89, 41, 63.

- (c) Using Minitab, our results were the numbers 46, 99, 90, 31, 75, 98, 79, 14, 44, 13, 66, 49, 37, 87, 73, 26, 61, 71, 72, 2. Thus our sample consists of the first 15 numbers 46, 99, 90, 31, 75, 98, 79, 14, 44, 13, 66, 49, 37, 87, 73. Your sample may be different.

18. The statement under the vote is a disclaimer as to the validity of the survey. Since the vote reflects only the responses of volunteers who chose to vote, it cannot be regarded as representative of the public in general, some of which do not use the Internet, nor as representative of Internet users since the sample was not chosen at random from either group.
19. The data in this study were clearly not collected via a controlled experiment in which some participants were forced to do crossword puzzles, practice musical instruments, play board games, or read while others were not allowed to do any of those activities. Therefore, any data relative to these activities and dementia arose as a result of observing whether or not the subjects in the study carried out any of those activities and whether or no they had some form of dementia. Since this would be an observational study, no statement of cause and effect can rightfully be made.
20. The researchers did not impose or manipulate any of the conditions of this study. They didn't decide who had cancer, who didn't have cancer, who had hepatitis B, or who had hepatitis C. This study was an observational study and not a controlled experiment. Observational studies can only reveal an association, not causation. Therefore, the statement in quotes is valid. If the researchers wanted to establish causation, they would need a designed experiment.
21. (a) Answers will vary, but here is the procedure: (1) Divide the population size, 100, by the sample size 15, and round down to the nearest whole number; this gives 6. (2) Use a table of random numbers (or a similar device) to select a number between 1 and 6, call it k . (3) List every 6th number, starting with k , until 15 numbers are obtained; thus the first number on the required list of 15 numbers is k , the second is $k+6$, the third is $k+12$, and so forth (e.g., if $k=4$, then the numbers on the list are 4, 10, 16, ...).
- (b) Yes, unless for some reason there is some kind of trend or a cyclical pattern in the listing of the athletes.

18 Chapter 1

22. (a) Each category of "Distance from Plant" should be represented in the sample in the same proportion that it is present in the population of City of Durham's water distribution system. $1310/11707 = 0.112$. Thus, 11.2% of the sample of 80 water samples should be from "Less than 1.5 miles", 27.0% from "1.5 - less than 3.0 miles", 24.1% from "3.0 - less than 4.5 miles", 13.6% from "4.5 - less than 6.0 miles", 11.5% from "6.0 - less than 7.5 miles", and 12.5% from "7.5 miles or greater". Multiplying each of these fractions by 80 gives us the sample sizes from each category. These sample sizes will not necessarily be integers, so we will need to make some minor adjustments of the results. The first category should have $(11.2/100)(80) = 8.96$. The second should have $(27/100)(80) = 21.6$. Similarly, the third, fourth, fifth, and sixth categories should have 19.28, 10.88, 9.2, and 10 for their sample sizes. We round the six sample sizes from the categories to 9, 22, 19, 11, 9, and 10 respectively. We would now randomly select water samples from each region.
23. (a) This is a designed experiment.
(b) The treatment group consists of the 158 patients who took AVONEX. The control group consists of the 143 patients who were given a placebo. The treatments were the AVONEX and the placebo.
24. (a) Experimental units: tomato plants
(b) Response variable: yield of tomatoes
(c) Factor(s): tomato variety and density of plants
(d) Levels of each factor: The four tomato varieties (Harvester, Pusa Early Dwarf, Ife No. 1, and Ibadan Local) would be the levels of variety. The four densities (10,000, 20,000, 30,000, and 40,000 plants/ha) would be the levels of the density.
(e) Treatments: Each treatment would be one of the combinations of a variety planted at a given plant density.
25. (a) Experimental Units: The children
(b) Response variable: Whether or not the child was able to open the bottle
(c) Factors: The container designs
(d) Levels of each factor: Three (types of containers)
(e) Treatments: The container designs
26. This is a completely randomized design. All of the experimental units (batches of doughnuts) were assigned at random to the four treatments (four different fats).
27. (a) This is a completely randomized design since the 24 cars were randomly assigned to the 4 brands of gasoline.
(b) This is a randomized block design. The four different gasoline brands are randomly assigned to the four cars in each of the six car model groups. The blocks are the six groups of four identical cars each.
(c) If the purpose is to learn about the mileage rating of one particular car model with each of the four gasoline brands, then the completely randomized design is appropriate. But if the purpose is to learn about the performance of the gasoline across a variety of cars (and this seems more reasonable), then the randomized block design is more appropriate and will allow the researcher to determine the effect of car model as well as of gasoline type on the mileage obtained.

Case Study: Top Films of All Time

- (a) The population of interest in the AFI survey is the population of film artists, critics, and historians.
- (b) The sample is the 1500 film artists, critics, and historians polled.
- (c) No. The population of all American moviegoers includes many people who are not film artists, critics, nor historians. Furthermore, these members of the film community have very specialized interests and possibly different viewpoints as to what constitutes a great actor or actress than many others in the American movie-going population.
- (d) Descriptive. It merely describes the opinion of those in the sample without trying to draw an inference about the opinions of all moviegoers.
- (e) Inferential. This statement would be an attempt to draw an inference about the opinion of all artists, historians, and critics based on the opinions of those 1500 people who were interviewed.

CHAPTER 2 SOLUTIONS

Exercises 2.1

- 2.1 (a) Hair color, model of car, and brand of popcorn are qualitative variables.
- (b) Number of eggs in a nest, number of cases of flu, and number of employees are discrete, quantitative variables.
- (c) Temperature, weight, and time are quantitative continuous variables.
- 2.2 (a) A qualitative variable is a nonnumerically valued variable. Its possible "values" are descriptive (e.g., color, name, gender).
- (b) A discrete, quantitative variable is one whose possible values can be listed. It is usually obtained by counting rather than by measuring.
- (c) A continuous, quantitative variable is one whose possible values form some interval of numbers. It usually results from measuring.
- 2.3 (a) Qualitative data result from observing and recording values of a qualitative variable, such as, color or shape.
- (b) Discrete, quantitative data are values of a discrete quantitative variable. Values usually result from counting something.
- (c) Continuous, quantitative data are values of a continuous variable. Values are usually the result of measuring something such as temperature that can take on any value in a given interval.
- 2.4 The classification of data is important because it will help you choose the correct statistical method for analyzing the data.
- 2.5 Of qualitative and quantitative (discrete and continuous) types of data, only qualitative yields nonnumerical data.
- 2.6 (a) The first column consists of *quantitative, discrete* data. This column provides ranks of the highest recorded temperature for each continent.
- (b) The second column consists of *qualitative* data since continent names are nonnumerical.
- (c) The fourth column consists of *quantitative, continuous* data. This column provides the highest recorded temperatures for the continents in degrees Fahrenheit.
- (d) The information that Death Valley is in the United States is *qualitative* data since country in which a place is located is nonnumerical.
- 2.7 (a) The first column consists of *quantitative, continuous* data. This column provides the time that the earthquake occurred.
- (b) The second column consists of *quantitative, continuous* data. This column provides the magnitude of each earthquake.
- (c) The third column consists of *quantitative, continuous* data. This column provides the depth of each earthquake in kilometers.
- (d) The fourth column consists of *quantitative, discrete* data. This column provides the number of stations that reported activity on the earthquake.
- (e) The fifth column consists of *qualitative* data since the region of the location of each earthquake is nonnumerical.
- 2.8 (a) The first column consists of *quantitative, discrete* data. This column provides ranks of the top ten IPOs in the United States.

22 Chapter 2

- (b) The second column consists of *qualitative* data since company names are nonnumerical.
- (c) The third column consists of *quantitative, discrete* data. Since money involves discrete units, such as dollars and cents, the data is discrete, although, for all practical purposes, this data might be considered quantitative continuous data.
- (d) The information that Facebook is a social networking business is *qualitative* data since type of business is nonnumerical.

- 2.9**
- (a) The first column consists of *quantitative, discrete* data. This column provides the ranks of the deceased celebrities with the top 10 earnings.
 - (b) The second column consists of *qualitative* data since names are nonnumerical.
 - (c) The third column consists of *quantitative, discrete* data, the earnings of the celebrities. Since money involves discrete units, such as dollars and cents, the data is discrete, although, for all practical purposes, this data might be considered quantitative continuous data.
- 2.10**
- (a) The first column consists of *quantitative, discrete* data. This column provides the ranks of the top 10 universities for 2012-2013.
 - (b) The second column consists of *qualitative* data since names of the institutions are nonnumerical.
 - (c) The third column consists of *quantitative, continuous* data. This column provides the overall score of the top 10 universities for 2012-2013.
- 2.11**
- (a) The first column contains types of products. They are *qualitative* data since they are nonnumerical.
 - (b) The second column contains number of units shipped in the millions. These are whole numbers and are *quantitative, discrete*.
 - (c) The third column contains money values. Technically, these are *quantitative, discrete* data since there are gaps between possible values at the cent level. For all practical purposes, however, these are *quantitative, continuous* data.
- 2.12** Player name, team, and position are nonnumerical and are therefore *qualitative* data. The number of runs batted in, or RBI, are whole numbers and are therefore *quantitative, discrete*. Weight is *quantitative, continuous*.
- 2.13** The first column contains *quantitative, discrete* data in the form of ranks. These are whole numbers. The second and third columns contain *qualitative* data in the form of names. The last column contains the rating of the program which is *quantitative, continuous*.
- 2.14** The first column is *qualitative* since it is nonnumerical. The second and third columns are *quantitative, discrete* since they report the number of grants and applications received. The last column is *quantitative, continuous* since it reports the success rate of the grants.
- 2.15** The first column is *quantitative, discrete* since it is reporting a rank. The second and third columns are *qualitative* since make/model and type are nonnumerical. The last column is *quantitative, continuous* since it is reporting mileage.
- 2.16** Of the eight items presented, only high school class rank involves ordinal data. The rank is ordinal data.

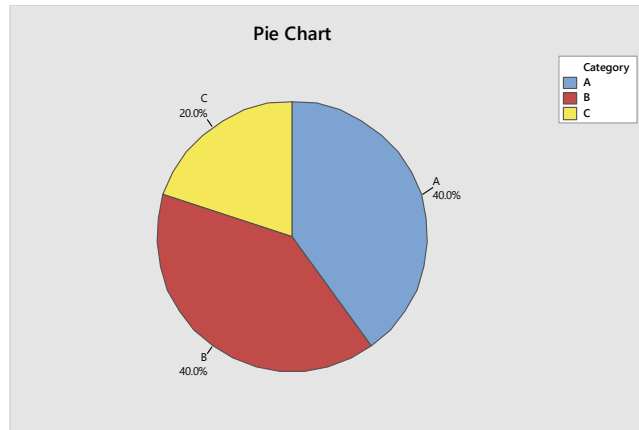
Exercises 2.2

- 2.17** A frequency distribution of qualitative data is a table that lists the distinct values of data and their frequencies. It is useful to organize the data and make it easier to understand.
- 2.18** (a) The frequency of a class is the number of observations in the class, whereas, the relative frequency of a class is the ratio of the class frequency to the total number of observations.
- (b) The percentage of a class is 100 times the relative frequency of the class. Equivalently, the relative frequency of a class is the percentage of the class expressed as a decimal.
- 2.19** (a) True. Having identical frequency distributions implies that the total number of observations and the numbers of observations in each class are identical. Thus, the relative frequencies will also be identical.
- (b) False. Having identical relative frequency distributions means that the ratio of the count in each class to the total is the same for both frequency distributions. However, one distribution may have twice (or some other multiple) the total number of observations as the other. For example, two distributions with counts of 5, 4, 1 and 10, 8, 2 would be different, but would have the same relative frequency distribution.
- (c) If the two data sets have the same number of observations, either a frequency distribution or a relative-frequency distribution is suitable. If, however, the two data sets have different numbers of observations, using relative-frequency distributions is more appropriate because the total of each set of relative frequencies is 1, putting both distributions on the same basis for comparison.
- 2.20** (a)-(b)

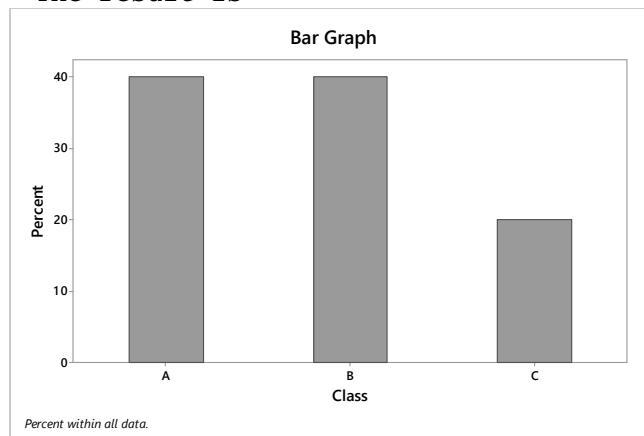
The classes are presented in column 1. The frequency distribution of the classes is presented in column 2. Dividing each frequency by the total number of observations, which is 5, results in each class's relative frequency. The relative frequency distribution is presented in column 3.

Class	Frequency	Relative Frequency
A	2	0.40
B	2	0.40
C	1	0.20
	5	1.00

- (c) We multiply each of the relative frequencies by 360 degrees to obtain the portion of the pie represented by each class. The result using Minitab is



(d) We use the bar chart to show the relative frequency with which each class occurs. The result is

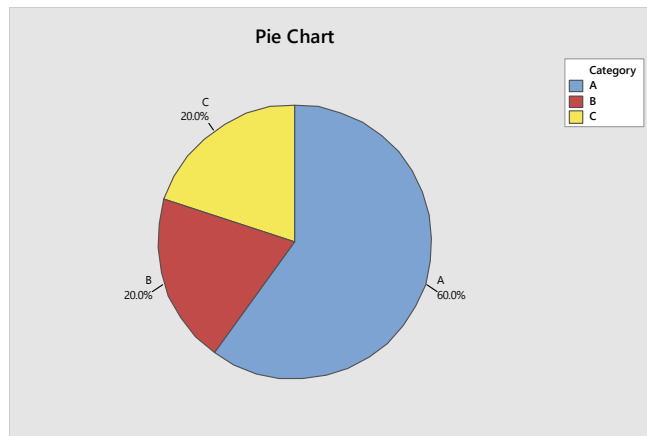


2.21 (a) - (b)

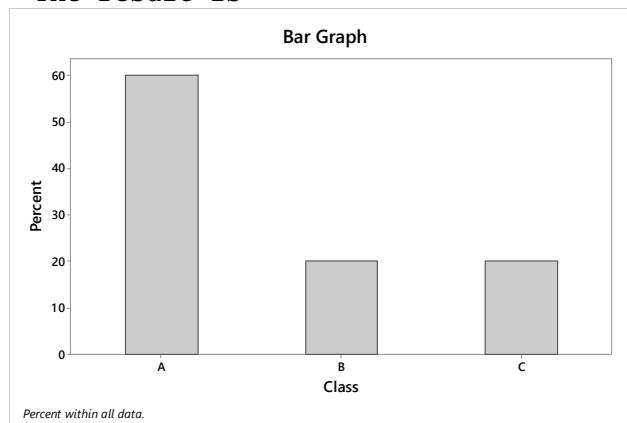
The classes are presented in column 1. The frequency distribution of the classes is presented in column 2. Dividing each frequency by the total number of observations, which is 5, results in each class's relative frequency. The relative frequency distribution is presented in column 3.

Class	Frequency	Relative Frequency
A	3	0.60
B	1	0.20
C	1	0.20
	5	1.00

(c) We multiply each of the relative frequencies by 360 degrees to obtain the portion of the pie represented by each class. The result using Minitab is



(d) We use the bar chart to show the relative frequency with which each class occurs. The result is

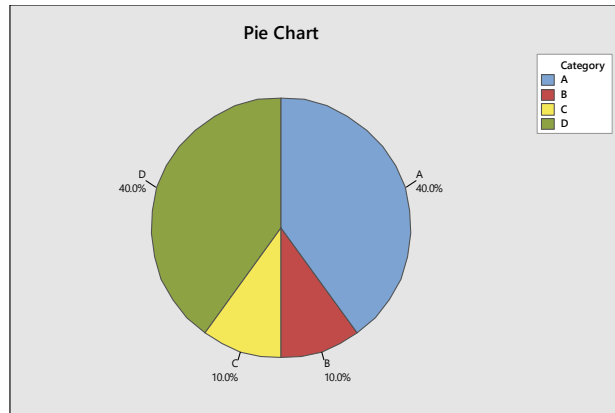


2.22 (a) - (b)

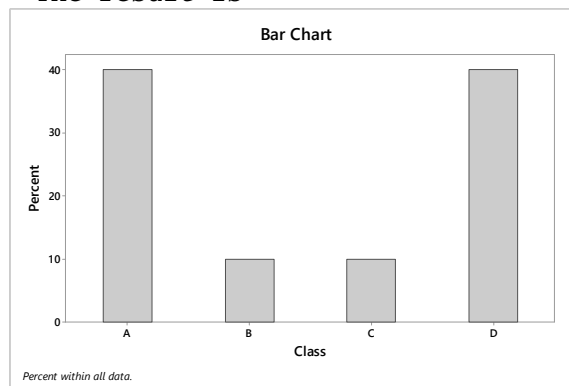
The classes are presented in column 1. The frequency distribution of the classes is presented in column 2. Dividing each frequency by the total number of observations, which is 10, results in each class's relative frequency. The relative frequency distribution is presented in column 3.

Class	Frequency	Relative Frequency
A	4	0.40
B	1	0.10
C	1	0.10
D	4	0.40
	10	1.00

(c) We multiply each of the relative frequencies by 360 degrees to obtain the portion of the pie represented by each class. The result using Minitab is



(d) We use the bar chart to show the relative frequency with which each class occurs. The result is

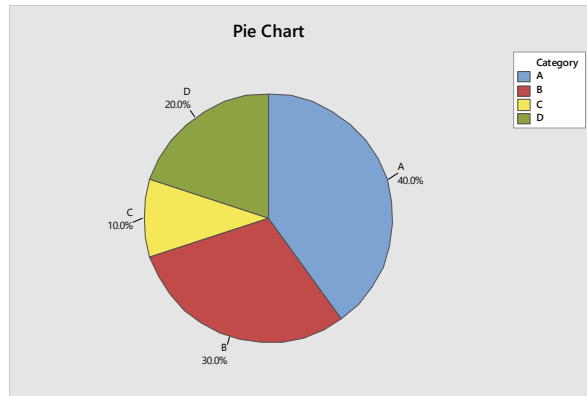


2.23 (a) – (b)

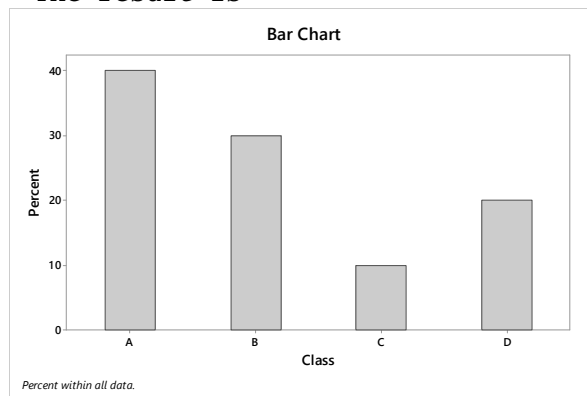
The classes are presented in column 1. The frequency distribution of the classes is presented in column 2. Dividing each frequency by the total number of observations, which is 10, results in each class's relative frequency. The relative frequency distribution is presented in column 3.

Class	Frequency	Relative Frequency
A	4	0.40
B	3	0.30
C	1	0.10
D	2	0.20
	10	1.00

(c) We multiply each of the relative frequencies by 360 degrees to obtain the portion of the pie represented by each class. The result using Minitab is



(d) We use the bar chart to show the relative frequency with which each class occurs. The result is

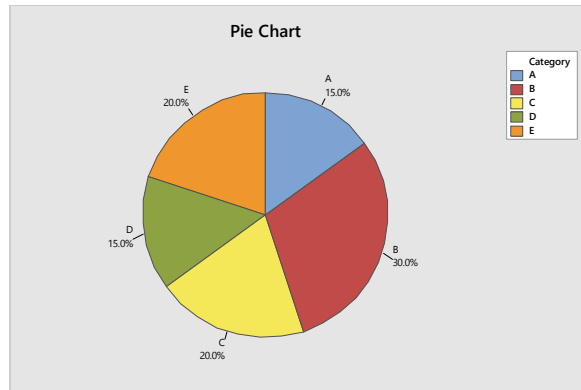


2.24 (a) - (b)

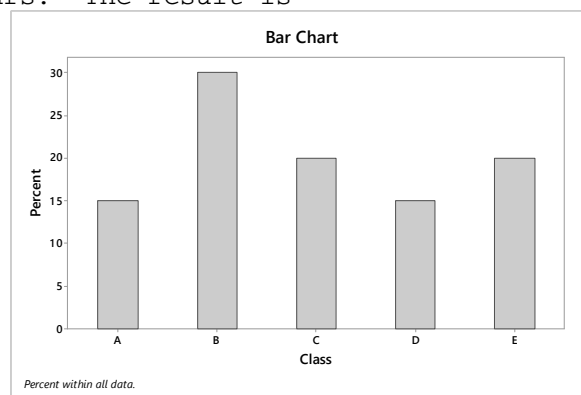
The classes are presented in column 1. The frequency distribution of the classes is presented in column 2. Dividing each frequency by the total number of observations, which is 20, results in each class's relative frequency. The relative frequency distribution is presented in column 3.

Class	Frequency	Relative Frequency
A	3	0.15
B	6	0.30
C	4	0.20
D	3	0.15
E	4	0.20
	20	1.00

(c) We multiply each of the relative frequencies by 360 degrees to obtain the portion of the pie represented by each class. The result using Minitab is



(d) We use the bar chart to show the relative frequency with which each class occurs. The result is

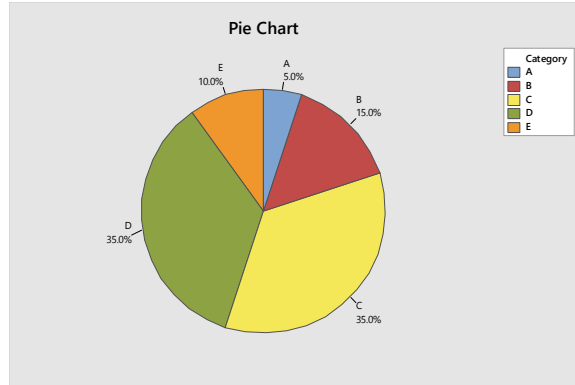


2.25 (a) - (b)

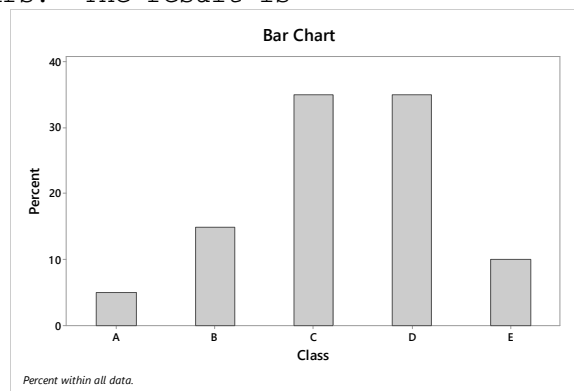
The classes are presented in column 1. The frequency distribution of the classes is presented in column 2. Dividing each frequency by the total number of observations, which is 20, results in each class's relative frequency. The relative frequency distribution is presented in column 3.

Class	Frequency	Relative Frequency
A	1	0.05
B	3	0.15
C	7	0.35
D	7	0.35
E	2	0.10
	20	1.00

(c) We multiply each of the relative frequencies by 360 degrees to obtain the portion of the pie represented by each class. The result using Minitab is



(d) We use the bar chart to show the relative frequency with which each class occurs. The result is

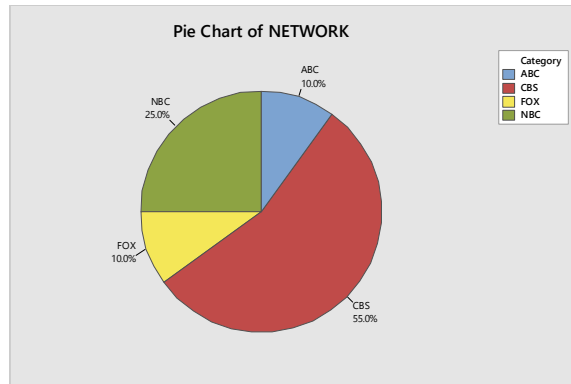


2.26 (a) - (b)

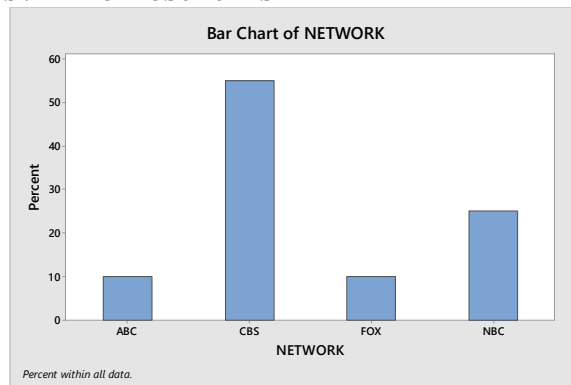
The classes are the networks and are presented in column 1. The frequency distribution of the networks is presented in column 2. Dividing each frequency by the total number of shows, which is 20, results in each class's relative frequency. The relative frequency distribution is presented in column 3.

Network	Frequency	Relative Frequency
ABC	2	0.10
CBS	11	0.55
Fox	2	0.10
NBC	5	0.25
	20	1.00

(c) We multiply each of the relative frequencies by 360 degrees to obtain the portion of the pie represented by each network. The result is



(d) We use the bar chart to show the relative frequency with which each network occurs. The result is

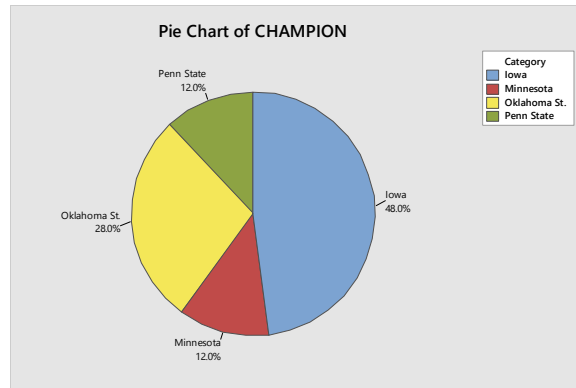


2.27 (a) - (b)

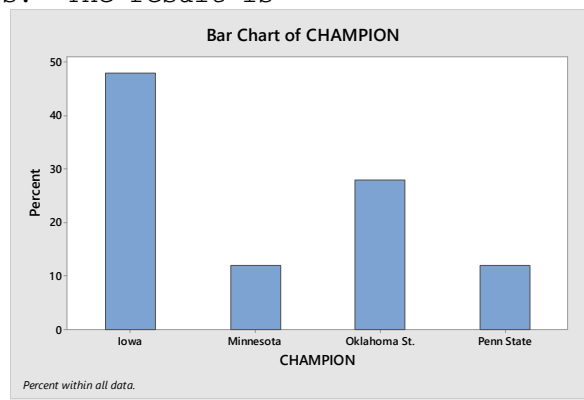
The classes are the NCAA wrestling champions and are presented in column 1. The frequency distribution of the champions is presented in column 2. Dividing each frequency by the total number of champions, which is 25, results in each class's relative frequency. The relative frequency distribution is presented in column 3.

Champion	Frequency	Relative Frequency
Iowa	12	0.48
Penn State	3	0.12
Minnesota	3	0.12
Oklahoma St.	7	0.28
	25	1.00

(b) We multiply each of the relative frequencies by 360 degrees to obtain the portion of the pie represented by each team. The result is



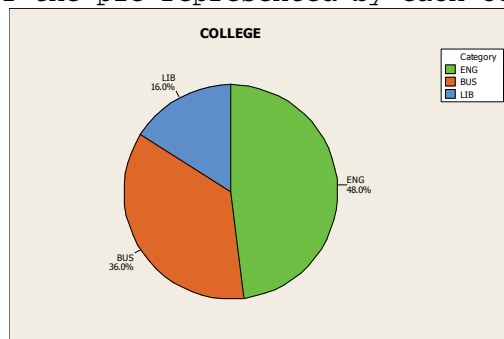
(c) We use the bar chart to show the relative frequency with which each TEAM occurs. The result is



2.28 (a)-(b) The classes are the colleges and are presented in column 1. The frequency distribution of the colleges is presented in column 2. Dividing each frequency by the total number of students in the section of Introduction to Computer Science, which is 25, results in each class's relative frequency. The relative frequency distribution is presented in column 3.

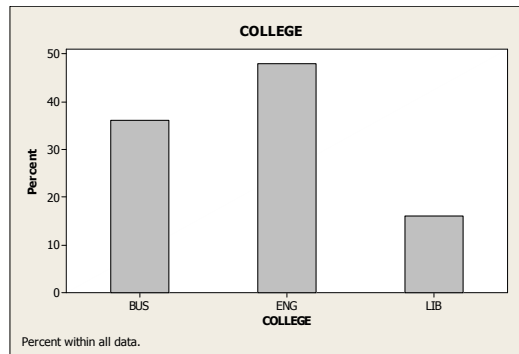
College	Frequency	Relative Frequency
BUS	9	0.36
ENG	12	0.48
LIB	4	0.16
	25	1.00

(c) We multiply each of the relative frequencies by 360 degrees to obtain the portion of the pie represented by each college. The result is



32 Chapter 2

(b) We use the bar chart to show the relative frequency with which each COLLEGE occurs. The result is

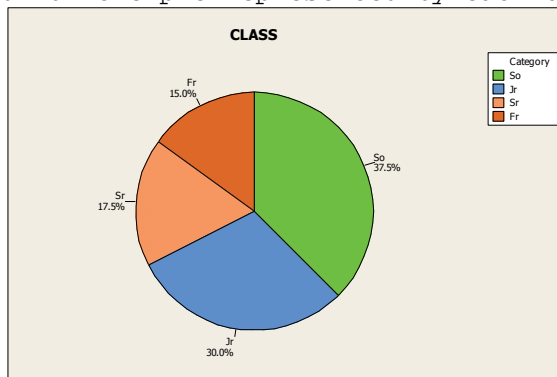


2.29 (a) - (b)

The classes are the class levels and are presented in column 1. The frequency distribution of the class levels is presented in column 2. Dividing each frequency by the total number of students in the introductory statistics class, which is 40, results in each class's relative frequency. The relative frequency distribution is presented in column 3.

Class Level	Frequency	Relative Frequency
Fr	6	0.150
So	15	0.375
Jr	12	0.300
Sr	7	0.175
	40	1.000

(c) We multiply each of the relative frequencies by 360 degrees to obtain the portion of the pie represented by each class level. The result is



(d) We use the bar chart to show the relative frequency with which each CLASS level occurs. The result is