Chapter 5: Organizing Data

Section 5.1: Introduction	177
Section 5.2: Organizing and Graphing Categorical Data	
Section 5.3: Organizing and Graphing Quantitative Data	192

Section 5.1: Introduction

Objective: Recognize types of data.

The term "statistics" emerged around the 18th century due to the need of governments to collect demographic and economic data. Nowadays, some of the fields that use statistics include finance, economics, actuarial science, biostatistics, etc.

Definition. Statistics *is the art and science of collecting, analyzing, presenting, and interpreting data. It provides tools for predicting and forecasting the use of data through statistical models.*

What are Data?

Data are observations that have been collected. Examples include measurements, income levels, or responses to survey questions. Any set of data contains information about some group of individuals or things. The information is organized in variables.

A variable is a characteristic of an individual or some thing.



Figure 1 Data type classification

Example 1. Eye color



An eye clinic collected the eye color of their patients. Eye color is classified as qualitative data because eye color consists of values that are categories (such as brown, blue, green, hazel, and black).

Example 2. People's Height



The variable "heights of people" takes on an infinite number of possible values. It is classified as quantitative continuous data. Even though a person's height might be 5'8", it is really 5.68241231... feet. It is impossible to measure a length exactly!

Example 3. Household Size



The variable number of people in a household takes on a finite number of values. These possible values are 1, 2, 3, ... which can be plotted on a number line as points with space between each point. Thus, number of people in a household is classified as quantitative discrete data.

Example 4.

The table below contains data organized by variables about a small group of CCBC students. The variables are: Student Name, Gender, No. of classes completed, GPA, Student ID. We would like to identify the data type (variable) presented in each column as either qualitative, quantitative discrete, or quantitative continuous.

Name	Gender	No. of classes completed	GPA	Student ID
Tim	Male	4	3.25	800567230
Joy	Female	9	3.75	800213457
Brad	Male	6	4.00	800120985
Alice	Female	5	3.90	800900521
Michael	Male	10	3.81	800749541

Students' name, as well as students' gender, takes on values that are categories. The data in both of these cases are classified as qualitative.

The variable "number of classes completed" per student is finite. Its possible values are 0, 1, 2, 3, ...; these can be plotted on a number line as points with space between each point. Thus, the data in these cases are classified as quantitative discrete.

GPA can take an infinite number of possible values in the interval 0.0 to 4.0. So, these data are classified as quantitative continuous.

Although the variable Student ID consists of numbers, the numbers are used as labels Thus these data are classified as qualitative.

5.1 Practice

You plan to purchase a car and you have your mind set on a few car characteristics (variables). Determine if these variables are either *quantitative discrete, quantitative continuous*, or *qualitative*.

- 1) Your car budget limits
- 2) Car Make
- 3) Number of cylinders
- 4) Car Model
- 5) The braking distance measured on a scale from 100 ft. to 200 ft

It is Orioles opening season game day. Determine if the following variables are either *quantitative discrete, quantitative continuous,* or *qualitative.*

- 6) The high temperature of the day.
- 7) Birth country of the players.
- 8) Number of home runs during the game.
- 9) Shirt numbers on athletes uniforms.
- 10) The table below provides information of some specifications for five smartphones. Identify the type of data presented in each column as either *qualitative, quantitative discrete,* or *quantitative continuous*.

Rank	Smartphone	Weight (oz)	Touch ID	Battery (talk time in min.)
1	Apple iPhone 6S	5.04	Yes	840
2	Apple iPhone 6S Plus	6.77	Yes	1440
3	Samsung Gallaxy S7	5.37	Yes	1680
4	Sony Xperia Z	5.15	No	840
5	LG G5	5.57	No	1320

11) The table below provides information of some specifications for five breakfast cereals. Identify the type of data presented in each column as either *qualitative*, *quantitative discrete*, or *quantitative continuous*.

Cereal Name	Cereal Brand	Calories (per cup)	Sugar (per	Gluten
			cup)	Free
Special K	Kellogg	110	9g	No
Cheerios	General Mills	100	1g	Yes
Oats and Honey	Cascadian Farm	260	14g	No
Corn Flakes	Millville	100	2g	No
Rice Chex	General Mills	100	2g	Yes

- 12) Below is a small survey about driving history. Identify the type of data being collected with each survey question as either *qualitative*, *quantitative discrete*, or *quantitative continuous*.
 - a. How old are you?
 - b. What is your gender?
 - c. Have you taken a driver's education course?
 - d. What is the fastest speed you have driven an automobile?
 - e. How many speeding tickets have your received in your lifetime?

13) You have decided to collect some data on your Netflix account. Identify the type of data you are collecting as either *qualitative*, *quantitative discrete*, or *quantitative continuous*.

- a. The genre of the film.
- b. The length (in minutes) of the film.
- c. The number of family members that have access to the account.

5.1 Answers

- 1) Quantitative continuous
- 2) Qualitative
- 3) Quantitative discrete
- 4) Qualitative
- 5) Quantitative continuous
- 6) Quantitative continuous
- 7) Qualitative
- 8) Quantitative discrete
- 9) Qualitative
- 10) Qualitative data: Rank, Smartphone name, Touch ID. Quantitative continuous: Weight, Battery talk time
- 11) Qualitative data: Cereal Name, Cereal Brand, Gluten Free Quantitative continuous: Calories (per cup), Sugar (per cup)

12)

- a. Quantitative continuous
- b. Qualitative
- c. Qualitative
- d. Quantitative continuous
- e. Quantitative discrete

13)

- a. Qualitative
- b. Quantitative continuous
- c. Quantitative discrete

Section 5.2: Organizing and Graphing Categorical Data

Objective: Create a frequency table.

Data is being collected all the time by businesses, governments, and researchers. The data can range from small to quite large. We need to be able to better understand the nature of the data. Organizing it helps! In this section, we will organize qualitative data. Because qualitative data is classified into categories, we will refer to this data as categorical data.

Definition. A *frequency table* shows how data are divided among several categories (or classes) by listing the categories along with the number (frequency) of data values in each of them.

Example. Coffee Shop

A local coffee shop keeps a list of types of drinks that their customers order each hour. Below is the data from 50 drinks sold an hour before closing on a recent Tuesday.

Coffee	Espresso	Coffee	Tea	Espresso
Tea	Coffee	Tea	Coffee	Espresso
Espresso	Tea	Espresso	Espresso	Espresso
Espresso	Espresso	Coffee	Espresso	Tea
Coffee	Soda	Espresso	Coffee	Coffee
Espresso	Tea	Espresso	Soda	Tea
Coffee	Espresso	Coffee	Tea	Espresso
Coffee	Soda	Coffee	Coffee	Espresso
Soda	Espresso	Tea	Espresso	Coffee
C C				

A frequency table provides a useful way to organize this data. There are four different categories of drinks they sell (Coffee, Espresso, Soda, and Tea). They are in the first column of the table. The second column contains the counts of each type sold that hour.

Drink	Frequency
Coffee	16
Espresso	20
Soda	4
Tea	10

Objective: Compute relative frequencies.

An additional column can be added to the table to give us a better understanding of the data. This column is called a relative frequency. It can be found by dividing the frequency count for a category by the sum of all frequency counts.

Relative frequency for a category	_ frequency of a category
Relative frequency for a category	sum of all frequencies

Drink	Frequency	Relative Frequency
Coffee	16	$\frac{16}{50} = 0.32$
Espresso	20	$\frac{20}{50} = 0.40$
Soda	4	$\frac{4}{50} = 0.08$
Tea	10	$\frac{10}{50} = 0.20$

If we turn the relative frequency into a percent we get a better understanding of the sales that hour.

Percentage for a category	_ frequency of	f a category × 100%
rereentage for a category	sum of all f	requencies

Drink	Frequency	Relative Frequency	Percentage
Coffee	16	$\frac{16}{50} = 0.32$	0.32×100% = 32%
Espresso	20	$\frac{20}{50} = 0.40$	0.40×100% = 40%
Soda	4	$\frac{4}{50} = 0.08$	0.08×100% = 8%
Tea	10	$\frac{10}{50} = 0.20$	0.20×100% = 20%

Initially it was easy to see that Espresso is the most popular drink, but looking at the relative frequency or percentage gives us a better idea of how it relates to the other drink choices. The manager can use information such as this to better stock the shop.

Objective: Create a bar graph.

Nothing makes a report look better than a nice graph. In addition to creating frequency tables, an analyst might want to create a graph of categorical data. There are many different types of graphs. **Bar graphs** are probably the most commonly used graphs and they are used to compare things between different groups.

Here is a bar chart of our coffee shop data. The categories are along the horizontal axis and the frequency counts correspond to the height of the bars.



Rules when constructing a bar graph

- 1. The height of each bar represents the frequency or relative frequency for that category.
- 2. The bars should be of the same "width."
- 3. The bars should not overlap.
- 4. Each piece of data should belong to only one category.

Objective: Create a pie chart.

Like bar graphs, pie charts are very common to graph categorical data. **Pie charts** show how the size of the category relates to the whole group. Pie charts are great for showing percentages. Below is the coffee shop pie chart. Notice how the percentages correspond to the size of the pie pieces.



Rules when constructing a pie chart

- 1. Always include the relative frequency or percentage.
- 2. Include labels, either as a legend or directly on pie.

5.2 Practice

1. Twenty-four students answered a survey about pet preferences. Their responses are below.

Cat	Guinea pigs	Guinea pigs	Cat	Rabbit	Dog	Guinea pigs	Dog
Dog	Cat	Dog	Dog	Guinea pigs	Dog	Cat	Rabbit
Dog	Rabbit	Cat	Guinea pigs	Dog	Cat	Rabbit	Dog

- a) Construct a frequency table for this data.
- b) Draw a bar graph.
- c) How many students participated in this survey?
- d) What percent of students like dogs?
- 2. In a software engineering class, the professor asked his students to name their favorite programming language. Their replies are listed in the table below.

Java	Lisp	Perl	Java	Perl	Perl	Perl	C++	Perl	Java
Perl	Java	Java	Perl	Java	Lisp	Java	Perl	Java	Lisp
Perl	C++	C++	Perl	C++	Perl	Java	C++	Perl	C++

- a) Construct a frequency table for this data.
- b) Draw a bar graph.
- c) How many students participated in this survey?
- d) What percent of students like C++?
- 3. The following frequency table represents the number of new HIV/AIDS cases in the US in 2008 according to race/ethnicity. What percent of the new cases were Hispanic/Latino?

Race/Ethnicity	Number of HIV/AIDS Cases
American Indian/Alaskan Native	228
Asian	451
Black/African American	21,443
Hispanic/Latino	7,461
Native Hawaiian/other Pacific Islander	47
White	12,534

4. A school district performed a study to find the main causes leading to its students dropping out of school. Fifty cases were analyzed and a primary cause was assigned to each case. The results for the fifty cases are listed below. What percent of students drop out due to family problems?

Causes to drop out of school	Frequency
Unexcused absences	12
Illness	16
Family problems	14
Other causes	8

5. Relative frequencies allow us to compare groups. Here is the 2008 new HIV/AIDS cases in the US separated by sex.

Males				
Race/Ethnicity	Frequency			
Black	14,247			
White	10,563			
Hispanic	5,906			
Other	565			

Females						
Race/Ethnicity	Frequency					
Black	7,196					
White	1,971					
Hispanic	1,555					
Other	161					

- a) Compute the relative frequencies for each sex.
- b) Write a few sentences explaining the trend of new cases in 2008 using what you learned in part a.
- 6. The table below represents 360 books grouped by their category:

Book Category	Frequency
Science Fiction	150
Romance	125
Adventure	50
Horror	35

Which pie chart best represents this table?



187



7. The bar chart below describes the day of the week workers called in sick for workers at a company. What is the relative frequency for Monday?

8. The bar chart below show the blood groups of O, A, B, and AB of a group of forty randomly selected blood donors. How many donors have a blood group of O?



9. The pie chart below shows the student responses to a survey asking them about their favorite holiday. Use the graph to find the percent of students who answered "4th of July".



5.2 Answers

1. a)	Favorite pets	Frequency
	Rabbits	4
	Cats	6
	Dogs	9
	Guinea pigs	5
	Total	24

b) Frequency



c) 24 students participated in the surveyd) 37.5% of students like dogs

2. a)

Programming language	Frequency
C++	6
Java	9
Lisp	3
Perl	12
Total	30



- c) 30 students participated in the survey
- d) 20% of students like C++
- 3. 17.7%
- 4. 28%
- 5. a)

	Males		Fe	males	
Race/Ethnicity	Relative Frequency		Race/Ethnicity	Relative Frequency	
Black	0.455		Black	0.661	
White	0.338		White	0.181	
Hispanic	0.189		Hispanic	0.143	
Other	0.018		Other	0.015	

- b) When you are comparing new cases of HIV among men and women their relative frequencies for race are nothing alike. With both males and females, the majority of the cases are with Blacks and Whites. Females have two thirds of new cases just among Blacks. The epidemic of HIV is very different according to race for men and women.
- 6. Pie chart (c)
- 7. 0.2556
- 8. 15
- 9. 25%

Section 5.3: Organizing and Graphing Quantitative Data

Objective: Organize quantitative data into frequency tables

An easy way to compile quantitative data would be to make a frequency or relative frequency table as we did with qualitative data.

Commute time to school

The data in the table below are the commute time to school (in minutes) for a group of students attending a particular mathematics course. We would like to create a frequency table.

10	15	17	20	25	20	3	30	17	15	5	15
60	8	25	15	25	22	38	10	14	30	30	18

Since there are no natural categories for this data we must create what are called classes. Classes divide the number line into smaller pieces.

First let's find the minimum commute time (3 minutes). Next, find the maximum commute time (60 minutes). We can start our first class at 3 minutes or back up a few minutes. Let's start the first class at 1 minute and use a class width of 10 minutes. The first class is "1 - 10". The 1 is considered the lower limit of the class and the 10 is considered the upper limit of the class.



Next, we will determine the lower limits of the other classes. Add the class width of 10 to our first lower limit. This will give us 11. If we continue to add 10 the remaining lower limits will be 21, 31, 41, and 51. The upper limits are determined by filling in the numbers that approach the next class' lower limit but do not equal it. For example in the second row we chose 20 because it is the closest whole number less than 21.



The class width of 10 is evident by looking at the differences in consecutive lower class limits. This is also the case with consecutive upper class limits. They are all 10. Once we made all the classes and made sure that the minimum and maximum can be placed in the table, we see that there are 6 classes. Now we must determine how many commute times fall into each of the classes. There are several strategies to tallying up the counts. You can mark them off as you go, or possibly put unique symbols or marks next to the ones that are in the same classes. This will help you avoid classifying a value twice or forgetting a data value altogether. Once all the symbols are there you can count them.

10 *	15 💖	17 🖏	20 💖	25*	20 💖	3 *	30�	17💖	15 💖	5 *	12💖
60∎	8 *	25*	15鬯	25*	22*	38 X	10 *	14 🖗	30�	30�	18💖

	Commute Time	Frequency
*	1 – 10	5
r an	11 - 20	10
÷	21 - 30	7
Ħ	31 - 40	1
X	41 - 50	0
	51 - 60	1

Vocabulary:

Classes	=	range of data values used to group the data.
Lower class limit	=	the smallest value that goes in a class.
Upper class limit	=	the largest value that goes in a class.
Class width	=	the difference between two consecutive lower class limits.

Objective: Create a histogram for quantitative data

We can take our frequency table and create a graph of the information. This graph is called a histogram. A histogram is like a bar graph. The classes are along the horizontal axis and the frequencies are demonstrated with the vertical axis. The bars need to touch in a histogram because we want to imply that the classes are adjacent and represent a continuum on the number line.

Steps to sketch a histogram:

- 1) Draw the horizontal axis with the lower class limits equally spaced along it.
- 2) Draw the vertical axis with the frequencies equally spaced along it.
- 3) Create the bars (rectangles)

The histogram for the "Commute time to school" is shown below.



You can see the height of the bars are exactly as our frequency table (5, 10, 7, 1, 0, and 1). The lower class limits are along the horizontal axis. The first bar gives us the commute times that were in the class 1 - 10. A histogram allows us to visualize the data and see how the data is spread out or possibly similar to one another. Notice there is not a bar for 41 - 50. This is because the frequency was zero for that class. Frequencies of zero are the only

reasons there are gaps in histograms.

Objective: Create Relative Frequencies

We can also add a column to the frequency table to represent relative frequency. This is similar to what was done for frequency tables for categorical data.

Commute Time	Frequency	Relative Frequency
1 - 10	5	$\frac{5}{24} = 0.21$
11 - 20	10	$\frac{10}{24} = 0.42$
21 - 30	7	$\frac{7}{24} = 0.29$
31 - 40	1	$\frac{1}{24} = 0.04$
41 - 50	0	0
51 - 60	1	$\frac{1}{24} = 0.04$



Relative frequency for a class	_ frequency in that class
Relative frequency for a class	sum of all frequencies

The relative frequencies, like the histogram, allow us to see that the majority of the class had a commute time between 11 and 30 minutes.

Objective: Discuss the shape of histograms

Histograms are valuable tools to display data. There are features that we tend to describe with words. For example, it is helpful to mention how many peaks (humps) are present in the graph. Does the histogram have a single, central peak or several separated peaks? A histogram with one main peak is dubbed **uni**modal; histograms with two peaks are **bi**modal; histograms with three or more peaks are called **multi**modal.

This is an example of a **unimodal** histogram. This histogram displays 90 different rainfalls at a national park. They measured the pH level of each rainfall.



This is an example of a bimodal histogram. These are the heights of 75 singers in a chorus.



Another way we describe histograms is to discuss their symmetry (or lack there of). If you can fold the histogram along a vertical line through the middle and have the edges match pretty closely, the histogram is considered symmetric.



The (usually) thinner ends of a distribution are called the tails. If one tail stretches out farther than the other, the histogram is said to be skewed to the side of the longer tail.

Below the histogram on the left is said to be skewed left, while the histogram on the right is said to be skewed right.



Objective: Create dotplots for quantitative data

Definition. A **dotplot** *is a visual representation of quantitative data and provides a graphical display of the data distribution.*

Steps to construct a dot plot:

- 1) Begin by drawing a number line that reflects the range of values.
- 2) Plot each data point by placing a dot over the appropriate value. For any repeated value stack the points.

The dotplot for the "Commute time to school" is shown below.



Using the dotplot we can answer many questions about our data.

What observation occurred the most?	We have four stacked dots above 15, so we conclude that 15 minutes is the observation that occurred the most.
How many students were surveyed?	Each point represents a reported commute time. By counting all the points on the graph we get the total number of students surveyed. We have 24 dots so the answer is 24 students.
What proportion of students commute to school in 20 minutes or less?	Count the points on the graph that are stacked at the value 20 or below it. 15 students commute in 20 minutes or less. Divide 15 by 24 (total students) to get the proportion. The answer is 15/24=0.625

Objective: Create stem-and-leaf plots for quantitative data

Definition. *A stem-and-leaf plot* is another way to arrange quantitative data. The plot separates each value into two parts: the stem (such as everything of the leftmost digits) and the leaf (such as the rightmost digit).

Steps to construct a stem-and-leaf plot:

1) Draw a vertical bar

- 2) Identify the stems and list them down on the left side of the vertical bar
- 3) List to the right of the vertical bar the leaves corresponding to their stems

Consider the Commute to school data

10	15	17	20	25	20	3	30	17	15	5	15
60	8	25	15	25	22	38	10	14	30	30	18

For this data, the stems are the tens, namely, 0, 1, 2, 3, 4, 5, 6 and the leaves are all the ones. For example, the data value 10 is separated into its stem of 1 and the leaf of 0. The data value 15 is separated into its stem of 1 and the leaf of 5. All the values are separated in the same way and arranged as shown below.

Note: The leaves are arranged in increasing order.

5.3 Practice

1) The frequency distribution below summarizes employee years of service for a particular corporation.

Years of service	Frequency
1 – 5	4
6 - 10	15
11 – 15	20
16 - 20	9
21 - 25	5
26 - 30	2

- a) How many employees participated in the survey?
- b) What is the class width?
- c) Identify the lower class limits.
- d) Identify the upper class limits.
- 2) A business magazine was conducting a study into the amount of travel required for mid-level managers across the US. They surveyed a group of managers and asked them the number of days they spent traveling each year. The frequency distribution below summarizes the results.

Days Traveling	Frequency
0-6	15
7 – 13	21
14 - 20	27
21 - 27	9
28 - 34	2
35 - 41	1

- a) How many managers participated in the survey?
- b) What is the class width?
- c) Identify the lower class limits.
- d) Identify the upper class limits.
- 3) A teacher at your college selected a group of students from all of his classes and asked them the number of credits each took that semester.

9	6	10	3	9	7	10	9	12	9	6	9
12	9	7	10	9	6	9	7	7	14	9	8

- a) Create a frequency and a relative frequency table to display the data. Let your first lower limit be 3 and the class width be 3 credits.
- b) Create a frequency histogram.
- c) Create a relative frequency histogram.
- d) How does the frequency histogram and relative frequency histogram compare?
- 4) Twenty-four students were asked the number of hours they sleep each night. The results of the survey are listed below.

7	9	5	8	8	11	12	7	10	7	8	9
8	6	4	3	9	10	7	6	12	5	6	11

- a) Create a frequency and a relative frequency table to display the data. Let your first lower limit be 3 and the class width be 2 hours.
- b) Create a frequency histogram.
- c) Create a relative frequency histogram.
- d) How does the frequency histogram and relative frequency histogram compare?
- 5) In a survey, 26 voters were asked their ages. The results are shown below.

43	56	28	63	67	66	52	48	37	56	40	60	62
66	45	21	35	49	32	56	61	53	69	31	48	59

- a) Create a frequency table for the ages. Let your first lower limit be 19 and the class width be 10.
- b) Create a frequency histogram.
- 6) The following histograms describe the number of minutes per day that a group of 50 elementary and 50 middle school students spent on a computer.



a) Which histogram appears to be skewed to the right?

- b) Which histogram appears to be skewed to the left?
- 7) The following histograms represent the pulse rate (beats per minute) for a group of 40 males and 40 females.



- a) Describe the shape of the two histograms.
- b) Compare the histogram of pulse rate of males with females. Is there a notable difference between pulse rates of females and males?
- 8) The stem and leaf plot below shows the number of laps run by each participant in a marathon. Refer to the stem and leaf plot to answer the following questions.

- a) Construct a dotplot.
- b) How many participants ran less than 48 laps?
- c) How many participants ran at least 50 laps?
- d) What observation occured the most?
- 9) The stem and leaf plot below shows the ages of a group of patients who had strokes caused by stress. Refer to the stem and leaf plot to answer the following questions.

- a) Construct a dotplot.
- b) How many patients were at most 36 years old?
- c) How many patients were 72 or older?
- d) What observation occurred the most?

5.3 Answers

1.

- a) 55 employees
- b) Class width: 5
- c) Lower class limits: 1, 6, 11, 16, 21, 26
- d) Upper class limits: 5, 10, 15, 20, 25, 30
- 2.
- a) 75 managers
- b) Class width: 7
- c) Lower class limits: 0, 7, 14, 21, 28, 35
- d) Upper class limits: 6, 13, 20, 27, 34, 41

3.

a)

Number of credits	Frequency	Relative Frequency
3 – 5	1	0.04
6-8	8	0.33
9-11	12	0.50
12 - 14	3	0.13



d) The frequency histogram and relative frequency histogram have the same shape because frequencies and relative frequencies are proportional. The only difference is on the units of the y – axis.

4.

	Hours of Sleep	Frequency	Relative Frequency
a)	3-4	2	0.08
	5-6	5	0.21
	7 - 8	8	0.33
	9 - 10	5	0.21
	11 – 12	4	0.17





5. a) $\begin{array}{c|c} Age & Frequency \\ \hline 19-28 & 2 \\ \hline 29-38 & 4 \\ \hline 39-48 & 5 \\ \hline 49-58 & 6 \\ \hline 59-68 & 8 \\ \hline 69-78 & 1 \\ \hline \end{array}$

b) Frequency



204

6.

- a) Elementary school students' histogram appears to be skewed to the right.
- b) Middle school students' histogram appears to be skewed to the left.
- 7.

a) The histogram of pulse rate of males appears to be unimodal and approximately symmetric. The histogram of pulse rate of females appears to be bimodal and skewed to the right.

b) The pulse rates of males appear to be generally lower than the pulse rates of females.



9.



c) 5 patients were 72 and older

d) 58 occured the most