Chapter 2: Exploring Data with Graphs and Numerical Summaries

Section 2.1: What Are the Types of Data?

Learning Objectives

- 1. Know the definition of variable
- 2. Know the definition and key features of a categorical versus a quantitative variable
- 3. Know the definition of a discrete versus a continuous quantitative variable
- 4. Know the definition of frequency, proportion (relative frequencies), and percentages
- 5. Create Frequency Tables

Learning Objective 1: Variable

- A variable is any characteristic that is recorded for the subjects in a study
- Examples: Marital status, Height, Weight, IQ
- A variable can be classified as either
 - Categorical, or
 - Quantitative (Discrete, Continuous)

Learning Objective 2: Categorical Variable

- A variable can be classified as categorical if each observation belongs to one of a set of categories.
- Examples:
 - Gender (Male or Female)
 - Religious Affiliation (Catholic, Jewish, …)
 - Type of residence (Apt, Condo, ...)
 - Belief in Life After Death (Yes or No)

Learning Objective 2: Quantitative Variable

- A variable is called quantitative if observations on it take numerical values that represent different magnitudes of the variable
 Examples:
 - Age
 - Number of siblings
 - Annual Income

Learning Objective 2: Main Features of Quantitative and Categorical Variables

- For Quantitative variables: key features are the center and spread (variability)
- For Categorical variables: a key feature is the percentage of observations in each of the categories

Learning Objective 3: Discrete Quantitative Variable

- A quantitative variable is discrete if its possible values form a set of separate numbers, such as 0,1,2,3,....
- Discrete variables have a finite number of possible values
- Examples:
 - Number of pets in a household
 - Number of children in a family
 - Number of foreign languages spoken by an individual

Learning Objective 3: Continuous Quantitative Variable

- A quantitative variable is continuous if its possible values form an interval
- Continuous variables have an infinite number of possible values
- Examples:
 - Height/Weight
 - Age
 - Blood pressure

Class Problem #1

- Identify the variable type as either categorical or quantitative
- 1. Number of siblings in a family
- 2. County of residence
- 3. Distance (in miles) of commute to school
- 4. Marital status

Class Problem #2

- Identify each of the following variables as continuous or discrete
- 1. Length of time to take a test
- 2. Number of people waiting in line
- Number of speeding tickets received last year
- 4. Your dog's weight

Learning Objective 4: Proportion & Percentage (Relative Frequencies)

- The proportion of the observations that fall in a certain category is the frequency (count) of observations in that category divided by the total number of observations
 - Frequency of that class
 Sum of all frequencies

The Percentage is the proportion multiplied by 100. Proportions and percentages are also called relative frequencies.

Learning Objective 4: Frequency, Proportion, & Percentage Example

- If 4 students received an "A" out of 40 students, then,
 - 4 is the frequency
 - 0.10 =4/40 is the proportion and relative frequency
 - 10% is the percentage .1*100=10%

Learning Objective 5: Frequency Table

A frequency table is a listing of possible values for a variable, together with the number of observations and/ or relative frequencies for each value

Frequency Table: Daily TV watching			
No. hours	Frequency	Percent	
0-1	232	25.6	
2-3	403	44.5	
4-5	181	20.0	
67	45	5.0	
8 or more	44	4.9	
Total	905	100.0	

Class Problem #3

A stock broker has been following different stocks over the last month and has recorded whether a stock is up, the same, or down in value. The results were

Performance of stock	Up	Same	Down
Count	21	7	12

- 1. What is the variable of interest
- 2. What type of variable is it?
- 3. Add proportions to this frequency table

Chapter 2: Exploring Data with Graphs and Numerical Summaries

Section 2.2: How Can We Describe Data Using Graphical Summaries?

Learning Objectives

- 1. Distribution
- 2. Graphs for categorical data: bar graphs and pie charts
- 3. Graphs for quantitative data: dot plot, stemleaf, and histogram
- 4. Constructing a histogram
- 5. Interpreting a histogram
- 6. Displaying Data over Time: time plots

Learning Objective 1: Distribution

- A graph or frequency table describes a distribution.
- A distribution tells us the possible values a variable takes as well as the occurrence of those values (frequency or relative frequency)

Learning Objective 2: Graphs for Categorical Variables

- Use pie charts and bar graphs to summarize categorical variables
 - Pie Chart: A circle having a "slice of pie" for each category
 - Bar Graph: A graph that displays a vertical bar for each category

Learning Objective 2: Pie Charts

Pie charts:

- used for summarizing a categorical variable
- Drawn as a circle where each category is represented as a "slice of the pie"
- The size of each pie slice is proportional to the percentage of observations falling in that category

Learning Objective 2: Pie Chart Example



Learning Objective 2: Bar Graphs

- Bar graphs are used for summarizing a categorical variable
- Bar Graphs display a vertical bar for each category
- The height of each bar represents either counts ("frequencies") or percentages ("relative frequencies") for that category
- Usually easier to compare categories with a bar graph than with a pie chart

Learning Objective 2: Bar Graph Example

Bar Graphs are called Pareto Charts when the categories are ordered by their frequency, from the tallest bar to the shortest bar

Source	U.S. Percentage	
Coal	51	Pei
Hydropower	6	
Natural gas	16	
Nuclear	21	
Petroleum	3	
Other	3	
Total	100	



Learning Objective 2: Class Exercise

There are 7 students in a class who are either freshman, sophomores, juniors, or seniors.The number of students in this class who are juniors is _____.



Learning Objective 3: Graphs for Quantitative Data

- Dot Plot: shows a dot for each observation placed above its value on a number line
- Stem-and-Leaf Plot: portrays the individual observations
 - Histogram: uses bars to portray the data

Learning Objective 3: Which Graph?

- Dot-plot and stem-and-leaf plot:
 - More useful for small data sets
 - Data values are retained
- Histogram
 - More useful for large data sets
 - Most compact display
 - More flexibility in defining intervals

Learning Objective 3: Dot Plots

- Dot Plots are used for summarizing a quantitative variable
- To construct a dot plot
- 1. Draw a horizontal line



- 2. Label it with the name of the variable
- 3. Mark regular values of the variable on it
- 4. For each observation, place a dot above its value on the number line

Learning Objective 3: Dot plot for Sodium in Cereals

Sodium Data:



Learning Objective 3: Stem-and-leaf plots

- Stem-and-leaf plots are used for summarizing quantitative variables
- Separate each observation into a stem (first part of the number) and a leaf (typically the last digit of the number)
- Write the stems in a vertical column ordered from smallest to largest, including empty stems; draw a vertical line to the right of the stems
- Write each leaf in the row to the right of its stem; order leaves if desired



Observation = 26, in a sample of size 20

Learning Objective 3: Stem-and-Leaf Plot for Sodium in Cereal

	Stems	Leaves
	0	0
	1	
Sodium Data:	2	This data point is the Sugar
0 210	3	Smacks sodium value, 70.
0 210	4 بر	The stem is 7 and the leaf
260 125	5 6	
000 000	7	0
220 290	8	
240 440	9	
210 140	10	
220 200	11	55
220 200	12	
125 170	14	0 These two leaf values are for
	15	0 Frosted Flakes and Wheaties.
250 150	16	Each has 200 mg of sodium
	17	$\begin{array}{c} 00 \\ 0 \end{array}$ per serving. The stem is 20 and each leaf is 0.
170 70	18	
000 000	20	00
230 200	21	00
200 490	22	00
290 100	23	0
	24	
	25	0
	26	0
	21	
	20	00

Learning Objective 4: Histograms



Number of Hours of TV Watching

Learning Objective 4: Steps for Constructing a Histogram

- 1. Divide the range of the data into intervals of equal width
- 2. Count the number of observations in each interval, creating a frequency table
- 3. On the horizontal axis, label the values or the endpoints of the intervals.
- Draw a bar over each value or interval with height equal to its frequency (or percentage), values of which are marked on the vertical axis.
- 5. Label and title appropriately

Learning Objective 4: Histogram for Sodium in Cereals

Sodium

Data:

TABLE 2.4: Frequency Table for Sodium in 20 Breakfast Cereals.

0 210
 The table summarizes the sodium values using eight intervals and lists the number of observations in each, as well as the proportions and percentages.

220	290	Interval	Frequency	Proportion	Percentage
210	140	0 to 39	I	0.05	5%
220	200	40 to 79	I	0.05	5%
125	170	80 to 119	0	0.00	0%
250	150	120 to 159	4	0.20	20%
200	150	160 to 199	3	0.15	15%
170	70	200 to 239	7	0.35	35%
230	200	240 to 279	2	0.10	10%
290	180	280 to 319	2	0.10	10%

Learning Objective 4: Histogram for Sodium in Cereals



Learning Objective 4: Histogram Example using TI 83+/84



Learning Objective 5: Interpreting Histograms

- Overall pattern consists of center, spread, and shape
 - Assess where a distribution is centered by finding the median (50% of data below median 50% of data above).
 - Assess the spread of a distribution.
 - Shape of a distribution: roughly symmetric, skewed to the right, or skewed to the left

Learning Objective 5: Shape

Symmetric Distributions: if both left and right sides of the histogram are mirror images of each other



Skewed to the left A distribution is skewed to the left if the left tail is longer than the right tail



Skewed to the right A distribution is skewed to the right if the right tail is longer than the left tail ³⁶
Learning Objective 5: Examples of Skewness



Learning Objective 5: Shape: Type of Mound



Learning Objective 5: Shape and Skewness

- Consider a data set containing IQ scores for the general public:
- What shape would you expect a histogram of this data set to have?
- a. Symmetric
- b. Skewed to the left
- c. Skewed to the right
- d. Bimodal

Learning Objective 5: Shape and Skewness

- Consider a data set of the scores of students on a very easy exam in which most score very well but a few score very poorly:
- What shape would you expect a histogram of this data set to have?
- a. Symmetric
- b. Skewed to the left
- c. Skewed to the right
- d. Bimodal

Learning Objective 5: Outlier

An Outlier falls far from the rest of the data



Learning Objective 6: Time Plots

- Used for displaying a time series, a data set collected over time.
- Plots each observation on the vertical scale against the time it was measured on the horizontal scale. Points are usually connected.
- Common patterns in the data over time, known as trends, should be noted.

Learning Objective 6: Time Plots Example

A Time Plot from 1995 – 2001 of the number of people worldwide who use the Internet



Chapter 2: Exploring Data with Graphs and Numerical Summaries

Section 2.3: How Can We Describe the Center of Quantitative Data?

Learning Objectives

- 1. Calculating the mean
- 2. Calculating the median
- 3. Comparing the Mean & Median
- 4. Definition of Resistant
- 5. Know how to identify the mode of a distribution

Learning Objective 1: Mean

The mean is the sum of the observations divided by the number of observations
 It is the center of mass



TI 83 Enter data into L1 STAT; CALC; 1:1-Var Stats; Enter L1;Enter

Learning Objective 1: Calculate Mean

Cereal	Sodium
Frosted Mini Wheats	0
Raisin Bran	210
All Bran	260
Apple Jacks	125
Capt Crunch	220
Cheerios	290
Cinnamon Toast	210
Crackling Oat Bran	140
Crispix	220
rosted Flakes	200
ruit Loops	125
Grape Nuts	170
Honey Nut Cheerios	250
.ife	150
Datmeal Raisin Crisp	170
iugar Smacks	70
ipecial K	230
Wheaties	200
Corn Flakes	290
Honeycomb	180

Learning Objective 2: Median

- The median is the midpoint of the observations when they are ordered from the smallest to the largest (or from the largest to smallest)
- Order observations
- If the number of observations is:
 - Odd, then the median is the middle observation
 - Even, then the median is the average of the two middle observations

Learning Objective 2: Median

		1) Sort ol	oservat	tions by size.		
	_	n = hum		Doservations	Order	Data
Order	Data				1	78
1	78				2	91
2	91	2.a) If <i>n</i> is odd , the median is			3	94
3	94	observation (n	+1)/2 C	iown the list	4	98
Δ	98	← <i>n</i> = 9			5	99
- т Б	O	(n+1)/2 = 10/2	2 = 5		6	101
5		Median = 99			7	103
6	101				8	105
7	103		<i>.</i> .		9	114
8	105	2.b) I	† <i>n</i> is e	ven, the median is the	10	121
9	114	mean c	or the ty	wo middle observations	5	
-				<i>n</i> = 10 →		
			Madi	(n+1)/2 = 5.5		
			Iviedia	an = (99+101)/2 = 100		49

Learning Objective 1 &2: Calculate Mean and Median

Cereal	Sodium
Frosted Mini Wheats	0
Raisin Bran	210
All Bran	260
Apple Jacks	125
Capt Crunch	220
Cheerios	290
Cinnamon Toast	210
Crackling Oat Bran	140
Crispix	220
Frosted Flakes	200
Fruit Loops	125
Grape Nuts	170
Honey Nut Cheerios	250
Life	150
Oatmeal Raisin Crisp	170
Sugar Smacks	70
Special K	230
Wheaties	200
Corn Flakes	290
Honeycomb	180

Enter data into L1

STAT; CALC; 1:1-Var Stats; Enter

L1;Enter



1-	Var	Stats	L1

Two screens' worth of statistics are produced - scroll down to read it all.



Leaning Objectives 1 & 2: Find the mean and median

- CO₂ Pollution levels in 8 largest nations measured in metric tons per person:
 2.3 1.1 19.7 9.8 1.8 1.2 0.7 0.2
- a. Mean = 4.6 Median = 1.5
- b. Mean = 4.6 Median = 5.8
- c. Mean = 1.5 Median = 4.6

Learning Objective 3: Comparing the Mean and Median

- The mean and median of a symmetric distribution are close together.
 - For symmetric distributions, the mean is typically preferred because it takes the values of all observations into account

Learning Objective 3: Comparing the Mean and Median

In a skewed distribution, the mean is farther out in the long tail than is the median

For skewed distributions the median is preferred because it is better representative of a typical observation



Learning Objective 4: Resistant Measures

- A numerical summary measure is resistant if extreme observations (outliers) have little, if any, influence on its value
 - The Median is resistant to outliers
 - The Mean is not resistant to outliers

Learning Objective 5: Mode

- Mode
 - Value that occurs most often
 - Highest bar in the histogram
 - The mode is most often used with categorical data

Chapter 2: Exploring Data with Graphs and Numerical Summaries

Section 2.4: How Can We Describe the Spread of Quantitative Data?

Learning Objectives

- 1. Calculate the Range
- 2. Calculate the standard deviation
- 3. Know the properties of the standard deviation
- Know how to interpret the magnitude of s: The Empirical Rule

Learning Objective 1: Range

One way to measure the spread is to calculate the range. The range is the difference between the largest and smallest values in the data set;

Range = max – min

The range is strongly affected by outliers

- Each data value has an associated deviation from the mean, $x \overline{x}$
- A deviation is positive if it falls above the mean and negative if it falls below the mean
- The sum of the deviations is always zero



Gives a measure of variation by summarizing the deviations of each observation from the mean and calculating an adjusted average of these deviations

$$s = \sqrt{\frac{\sum (x - \bar{x})^2}{n - 1}}$$



- Find the mean
- Find the deviation of each value from the mean
- Square the deviations
- Sum the squared deviations
- Divide the sum by n-1

(gives typical <u>squared deviation</u> from mean)

Metabolic rates of 7 men (cal./24hr.): 1792 1666 1362 1614 1460 1867 1439

$$\overline{X} = \frac{1792 + 1666 + 1362 + 1614 + 1460 + 1867 + 1439}{7}$$
$$= \frac{11,200}{7}$$
$$= 1600$$

	Observations	Deviations	Squared deviations
	\mathcal{X}_{i}	$x_i - \overline{x}$	$(x_i - \overline{x})^2$
	1792	1792–1600 = 192	$(192)^2 = 36,864$
	1666	1666 –1600 = 66	$(66)^2 = 4,356$
	1362	1362 - 1600 = -238	$(-238)^2 = 56,644$
	1614	1614 - 1600 = 14	$(14)^2 = 196$
	1460	1460 - 1600 = -140	$(-140)^2 = 19,600$
	1867	1867 - 1600 = 267	$(267)^2 = 71,289$
	1439	1439 - 1600 = -161	$(-161)^2 = 25,921$
		sum = 0	sum = 214,870
C	2 _ 214,870 _	- 25 011 67	
2	= =	-33,011.07	- 100 24 colorise
		$S = \sqrt{33.011.0}$	7 = 189.24 calones

Learning Objective 2: Calculate Standard Deviation

Cereal	Sodium
Frosted Mini Wheats	0
Raisin Bran	210
All Bran	260
Apple Jacks	125
Capt Crunch	220
Cheerios	290
Cinnamon Toast	210
Crackling Oat Bran	140
Crispix	220
Frosted Flakes	200
Fruit Loops	125
Grape Nuts	170
Honey Nut Cheerios	250
Life	150
Oatmeal Raisin Crisp	170
Sugar Smacks	70
Special K	230
Wheaties	200
Corn Flakes	290
Honeycomb	180

Enter data into L1

STAT; CALC; 1:1-Var Stats; Enter

L1;Enter



1-Var	Stats	L1

Two screens' worth of statistics are produced - scroll down to read it all.



Learning Objective 3: Properties of the Standard Deviation

- s measures the spread of the data
- s = 0 only when all observations have the same value, otherwise s > 0. As the spread of the data increases, s gets larger.
- s has the same units of measurement as the original observations. The variance=s² has units that are squared
- s is not resistant. Strong skewness or a few outliers can greatly increase s.

Learning Objective 4: Magnitude of s: Empirical Rule

EMPIRICAL RULE

If a distribution of data is bell-shaped, then approximately

- 68% of the observations fall within 1 standard deviation of the mean, that is, between $\overline{x} s$ and $\overline{x} + s$ (denoted $\overline{x} \pm s$).
- 95% of the observations fall within 2 standard deviations of the mean $(\overline{x} \pm 2s)$.
- All or nearly all observations fall within 3 standard deviations of the mean $(\overline{x} \pm 3s)$.



66

Chapter 2: Exploring Data with Graphs and Numerical Summaries

Section 2.5: How Can Measures of Position Describe Spread?

Learning Objectives

- 1. Obtaining quartiles and the 5 number summary
- 2. Calculating interquartile range and detecting potential outliers
- 3. Drawing boxplots
- 4. Comparing Distributions
- 5. Calculating a z-score

Learning Objective 1: Percentile

The pth percentile is a value such that p percent of the observations fall below or at that value



Learning Objective 1: Finding Quartiles

Splits the data into four parts

- Arrange the data in order
- The median is the second quartile, Q₂
- The first quartile, Q₁, is the median of the lower half of the observations
- The third quartile, Q₃, is the median of the upper half of the observations



Learning Objective 1:

Measure of spread: quartiles

Quartiles divide a ranked data set into four equal parts.

The **first quartile**, Q_1 , is the value in the sample that has 25% of the data at or below it and 75% above

The **second quartile** is the same as the **median** of a data set. 50% of the obs are above the median and 50% are below

The **third quartile**, Q_3 , is the value in the sample that has 75% of the data at or below it and 25% above

	0.6	1
	1.2	2
	1.6	3
	1.9	4
	1.5	5
O = first quartile = 2.2	2.1	6
$\mathbf{x}_1 = 1131$ qual the -2.2	2.3	7
	2.3	8
	2.5	9
	2.8	10
	2.9	11
M = machine = 0.4	3.3	12
M = median = 3.4	3.4	13
	3.6	14
	3.7	15
	3.8	16
	3.9	17
	4.1	18
Q_3 = third quartile = 4.3	4.2	19
	4.5	20
	4.7	21
	4.9	22
71	5.3	23
	5.6	24
	6.1	25

Learning Objective 1 Quartile Example

Find the first and third quartiles

Prices per share of 10 most actively traded stocks on NYSE (rounded to nearest \$)

- 2 4 11 13 14 15 31 32 34 47
- a. $Q_1 = 2$ $Q_3 = 47$
- b. $Q_1 = 12 \quad Q_3 = 31$
- c. $Q_1 = 11$ $Q_3 = 32$ d. $Q_1 = 12$ $Q_3 = 33$
Learning Objective 2: Calculating Interquartile range

- The interquartile range is the distance between the third quartile and first quartile:
- *IQR* = **Q3 Q1**
- IQR gives spread of middle 50% of the data



Learning Objective 2: Criteria for identifying an outlier

An observation is a potential outlier if it falls more than 1.5 x IQR below the first quartile or more than 1.5 x IQR above the third quartile



Learning Objective 3: 5 Number Summary

The five-number summary of a dataset consists of the

- Minimum value
- First Quartile
- Median
- Third Quartile
- Maximum value



Learning Objective 3: Calculate 5 Number Summary

Cereal	Sodium
Frosted Mini Wheats	0
Raisin Bran	210
All Bran	260
Apple Jacks	125
Capt Crunch	220
Cheerios	290
Cinnamon Toast	210
Crackling Oat Bran	140
Crispix	220
Frosted Flakes	200
Fruit Loops	125
Grape Nuts	170
Honey Nut Cheerios	250
Life	150
Oatmeal Raisin Crisp	170
Sugar Smacks	70
Special K	230
Wheaties	200
Corn Flakes	290
Honeycomb	180

Enter data into L1 STAT; CALC; 1:1-Var Stats; Enter L1

Enter

Scroll down to 5 number summary

EDIT DENE TESTS
UB 1-Var Stats
2:2-Var Stats
la:Ned-Ned Mai an Decidentales
5:0usdPag
6:CubicReg
74QuartRe9

Two screens' worth of statistics are produced - scroll down to read it all.

1-Var S	tats
x=185.	5
Σx=371	0
Σx2=78	4650
Sx=71.3	24642189
_σx=69.•	44242219
∔ n=20	

Learning Objective 3: Boxplot

- A box goes from the Q1 to Q3
- A line is drawn inside the box at the median
- A line goes from the lower end of the box to the smallest observation that is not a potential outlier and from the upper end of the box to the largest observation that is not a potential outlier
- The potential outliers are shown separately



Learning Objective 3: Boxplot for Sodium Data

Cereal	Sodium
Frosted Mini Wheats	0
Raisin Bran	210
All Bran	260
Apple Jacks	125
Capt Crunch	220
Cheerios	290
Cinnamon Toast	210
Crackling Oat Bran	140
Crispix	220
Frosted Flakes	200
Fruit Loops	125
Grape Nuts	170
Honey Nut Cheerios	250
Life	150
Oatmeal Raisin Crisp	170
Sugar Smacks	70
Special K	230
Wheaties	200
Corn Flakes	290
Honeycomb	180



Learning Objective 4: Comparing Distributions

Box Plots do not display the shape of the distribution as clearly as histograms, but are useful for making graphical comparisons of two or more distributions



Learning Objective 5: Z-Score

The z-score for an observation is the number of standard deviations that it falls from the mean

 $z = \frac{\text{observation} - \text{mean}}{\text{standard deviation}}$

An observation from a bell-shaped distribution is a potential outlier if its z-score < -3 or > +3

Chapter 2: Exploring Data with Graphs and Numerical Summaries

Section 2.6: How Can Graphical Summaries Be Misused?

Learning Objective 1: Guidelines for Constructing Effective Graphs

- Label both axes and provide proper headings
- To better compare relative size, the vertical axis should start at 0.
- Be cautious in using anything other than bars, lines, or points
- It can be difficult to portray more than one group on a single graph when the variable values differ greatly

Learning Objective 1: Example



FIGURE 2.18: An Example of a Poor Graph. Question: What's misleading about the way the data are presented?

Learning Objective 1: Example

