## Chapter 2 Describing, Exploring, and Comparing Data

Slide 1

- 2-1 Overview
- 2-2 Frequency Distributions
- 2-3 Visualizing Data
- 2-4 Measures of Center
- 2-5 Measures of Variation
- 2-6 Measures of Relative Standing
- 2-7 Exploratory Data Analysis



Created by Tom Wegleitner, Centreville, Virginia



## Overview



## Descriptive Statistics

## summarize or describe the important characteristics of a known set of population data

## Inferential Statistics

## use sample data to make inferences (or generalizations) about a population

## Important Characteristics of Data

- 1. Center: A representative or average value that indicates where the middle of the data set is located
- 2. Variation: A measure of the amount that the values vary among themselves
- 3. Distribution: The nature or shape of the distribution of data (such as bell-shaped, uniform, or skewed)
- 4. Outliers: Sample values that lie very far away from the vast majority of other sample values

## 5. Time: Changing characteristics of the data over time

# Section 2-2 Frequency Distributions

Created by Tom Wegleitner, Centreville, Virginia





## **Frequency Distributions**

## Frequency Distribution

lists data values (either individually or by groups of intervals), along with their corresponding frequencies or counts



#### ... from page 37...

#### Table 2-1 Measured Cotinine Levels in Three Groups

**Smoker:** The subjects report tobacco use.

**ETS:** (Environmental Tobacco Smoke) Subjects are nonsmokers who are exposed to environmental tobacco smoke ("secondhand smoke") at home or work.

**NOETS:** (No Environmental Tobacco Smoke) Subjects are nonsmokers who are not exposed to environmental tobacco smoke at home or work. That is, the subjects do not smoke and are not exposed to secondhand smoke.

Smoker:	1	0	131	173	265	210	44	277	32	3
	35	112	477	289	227	103	222	149	313	491
	130	234	164	198	17	253	87	121	266	290
	123	167	250	245	48	86	284	1	208	173
ETS:	384	0	69	19	1	0	178	2	13	1
	4	0	543	17	1	0	51	0	197	3
	0	3	1	45	13	3	1	1	1	0
	0	551	2	1	1	1	0	74	1	241
NOETS:	0	0	0	0	0	0	0	0	0	0
	0	9	0	0	0	0	0	0	244	0
	1	0	0	0	90	1	0	309	0	0
	0	0	0	0	0	0	0	0	0	0



#### ... from page 39...

#### Table 2-2

Frequency Distribution of Cotinine Levels of Smokers

Cotinine	Frequency
0–99	11
100–199	12
200–299	14
300–399	1
400–499	2



## are the smallest numbers that can actually belong to different classes

Cotinine	Frequency
0–99	11
100–199	12
200–299	14
300–399	1
400–499	2

## **Lower Class Limits**



## are the smallest numbers that can actually belong to different classes





## are the largest numbers that can actually belong to different classes





## are the numbers used to separate classes, but without the gaps created by class limits





### number separating classes



## Class Midpoints (p.40)



### midpoints of the classes

Class midpoints can be found by adding the lower class limit to the upper class limit and dividing the sum by two.





### midpoints of the classes







is the difference between two consecutive lower class limits or two consecutive lower class boundaries



**Reasons for Constructing Frequency Distributions** 



- 1. Large data sets can be summarized.
- 2. Can gain some insight into the nature of data.
- 3. Have a basis for constructing graphs.

### **Constructing A Frequency Table (p.40)**

1. Decide on the number of classes (should be between 5 and 20).

Slide 18

2. Calculate (round up).

class width ≈ (highest value) – (lowest value) number of classes

- 3. Starting point: Begin by choosing a lower limit of the first class.
- 4. Using the lower limit of the first class and class width, proceed to list the lower class limits.
- 5. List the lower class limits in a vertical column and proceed to enter the upper class limits.
- 6. Go through the data set putting a tally in the appropriate class for each data value.



### **Relative Frequency Distribution (p.41)**

## relative frequency = $\frac{class frequency}{sum of all frequencies}$

### **Relative Frequency Distribution**

#### ... from page 41...

... from page 39...

Cotinine	Frequency
0–99	11
100–199	12
200–299	14
300–399	1
400–499	2

**Total Frequency = 40** 

Table 2-3 Relative Freque Distribution of Levels in Smo	uency of Cotinine okers	11/40 = 28% 12/40 = 40%
Cotinine	Relative Frequency	etc.
0–99	28%	
100–199	30%	
200–299	35%	
300–399	3%	
400–499	5%	

Slide 20

### **Cumulative Frequency Distribution**



### ... from page 43...

... from page 39...

_

		-
Table 2-4		
Cumulative Freque		
of Cotinine Levels i	n Smokers	
	Cumulative	]
	Cumulative	
Cotinine	Frequency	
Less than 100	11	
Less than 200	23	
Less than 300	37	Cumulative
Less than 400	38	l l l l l l l l l l l l l l l l l l l
Less than 500	40	

## Comparison of Frequency Tables



#### Table 2-2

Frequency Distribution of Cotinine Levels of Smokers

Cotinine	Frequency
0–99	11
100–199	12
200–299	14
300–399	1
400–499	2

#### Table 2-3

Relative Frequency Distribution of Cotinine Levels in Smokers Relative Cotinine Frequency

Countie	riequency
0–99	28%
100–199	30%
200–299	35%
300–399	3%
400–499	5%

#### Table 2-4

Cumulative Frequency Distribution of Cotinine Levels in Smokers

Cotinine	Cumulative Frequency
Less than 100	11
Less than 200	23
Less than 300	37
Less than 400	38
Less than 500	40

## **Recap of Section 2-2**



In this Section we have discussed

- Important characteristics of data
- Frequency distributions
- Procedures for constructing frequency distributions
- Relative frequency distributions
- Cumulative frequency distributions



Created by Tom Wegleitner, Centreville, Virginia



## **Visualizing Data**



# Depict the nature of shape or shape of the data distribution

### Histogram



A bar graph in which the horizontal scale represents the classes of data values and the vertical scale represents the frequencies.

Cotinine	Frequency
0–99	11
100–199	12
200–299	14
300–399	1
400–499	2



#### Figure 2-1 (p.46)

### **Relative Frequency Histogram**



Has the same shape and horizontal scale as a histogram, but the vertical scale is marked with relative frequencies.

Cotinine	Relative Frequency
0–99	28%
100–199	30%
200–299	35%
300–399	3%
400–499	5%



Cotinine Levels of Smokers

Figure 2-2 (p.46)

Chapter 2, Triola, Elementary Statistics, MATH 1342





Figure 2-1

Figure 2-2

## **Frequency Polygon**



## Uses line segments connected to points directly above class midpoint values







#### A line graph that depicts cumulative frequencies



## **Dot Plot**



## Consists of a graph in which each data value is plotted as a point along a scale of values



#### Figure 2-5 (p.48)

## **Stem-and Leaf Plot (p.49)**



#### Represents data by separating each value into two parts: the stem (such as the leftmost digit) and the leaf (such as the rightmost digit)

#### **Stem-and-Leaf Plot**

Stem (tens)	Leaves (units)	
6	449	← Values are 64, 64, 69.
7	011123344445555556666778899	
8	0011122233346899	
9	0024	
10		
11		
12	0	$\leftarrow$ Value is 120.

### **Pareto Chart**



## A bar graph for qualitative data, with the bars arranged in order according to frequencies



## **Pie Chart**



#### A graph depicting qualitative data as slices pf a pie



#### Figure 2-7 (p.51)

## Scatter Diagram (p.51)



#### A plot of paired (x,y) data with a horizontal x-axis and a vertical y-axis



## **Time-Series Graph**



## Data that have been collected at different points in time



### Figure 2-8 (p.52)


Chapter 2, Triola, Elementary Statistics, MATH 1342

**Recap of Section 2-3** 



## In this Section we have discussed graphs that are pictures of distributions.

Keep in mind that the object of this section is not just to construct graphs, but to learn something about the data sets – that is, to understand the nature of their distributions.



Created by Tom Wegleitner, Centreville, Virginia



## Definition



#### Measure of Center

## The value at the center or middle of a data set



## Arithmetic Mean (Mean)

#### the measure of center obtained by adding the values and dividing the total by the number of values

## Notation (p.60)



- $\Sigma$  denotes the addition of a set of values
- *x* is the variable usually used to represent the individual data values
- *n* represents the number of values in a sample
- **N** represents the number of values in a population

## Notation



x is pronounced 'x-bar' and denotes the mean of a set of sample values  $\overline{x} = \underline{\sum x}$ 

 $\mu$  is pronounced 'mu' and denotes the mean of all values in a population

$$\mu = \frac{\sum x}{N}$$

n

## **Definitions (p.61)**



## Median

the middle value when the original data values are arranged in order of increasing (or decreasing) magnitude

**\Rightarrow** often denoted by  $\tilde{x}$  (pronounced 'x-tilde')

is not affected by an extreme value

## **Finding the Median**



If the number of values is odd, the median is the number located in the exact middle of the list

If the number of values is even, the median is found by computing the mean of the two middle numbers





## **Definitions (p.63)**



#### Mode

#### the value that occurs most frequently. The mode is not always unique. A data set may be: Bimodal Multimodal No Mode

#### denoted by M

#### the only measure of central tendency that can be used with nominal data

Examples (p.63)



a. 5.40 1.10 0.42 0.73 0.48 1.10	←Mode is 1.10
b. 27 27 27 55 55 55 88 88 99	Gimodal - 27 & 55
C. 1 2 3 6 7 8 9 10	←No Mode

## **Definitions (p.63)**



## Midrange

#### the value midway between the highest and lowest values in the original data set

#### highest score + lowest score

## Midrange =



### Round-off Rule for Measures of Center (p.65)

## Carry one more decimal place than is present in the original set of values

Mean from a Frequency Distribution (p.65)



# Assume that in each class, all sample values are equal to the class midpoint

#### Mean from a Frequency Distribution



#### use class midpoint of classes for variable x





In some cases, values vary in their degree of importance, so they are weighted accordingly

 $\overline{x} = \frac{\Sigma (w \cdot x)}{\Sigma w}$ 

#### **Best Measure of Center (p.67)**



Table 2-10	Comparison of Mean, Median, Mode, and Midrange						
Measure of Center	Definition	How Common?	Existence	Takes Every Value into Account?	Affected by Extreme Values?	Advantages and Disadvantages	
Mean	$\bar{x} = \frac{\Sigma x}{n}$	most familiar "average"	always exists	yes	yes	used throughout this book; works well with many statistical methods	
Median	middle value	commonly used	always exists	no	no	often a good choice if there are some extreme values	
Mode	most frequent data value	sometimes used	might not exist; may be more than one mode	no	no	appropriate for data at the nominal level	
Midrange	$\frac{\text{high} + \text{low}}{2}$	rarely used	always exists	no	yes	very sensitive to extreme values	

General comments:

- For a data collection that is approximately symmetric with one mode, the mean, median, mode, and midrange tend to be about the same.
- For a data collection that is obviously asymmetric, it would be good to report both the mean and median.
- The mean is relatively *reliable*. That is, when samples are drawn from the same population, the sample means tend to be more consistent than the other measures of center (consistent in the sense that the means of samples drawn from the same population don't vary as much as the other measures of center).

## **Definitions (p.67)**



## Symmetric

Data is symmetric if the left half of its histogram is roughly a mirror image of its right half.

## Skewed

#### Data is skewed if it is not symmetric and if it extends more to one side than the other.



Chapter 2, Triola, *Elementary Statistics*, MATH 1342

## Recap of Section 2-4 Slide 57

In this section we have discussed:

- Types of Measures of Center Mean Median Mode
- Mean from a frequency distribution
- Weighted means
- Best Measures of Center

#### Skewness



Created by Tom Wegleitner, Centreville, Virginia



## **Measures of Variation**



#### Because this section introduces the concept of variation, this is one of the most important sections in the entire book

## **Definition (p.74)**



# The range of a set of data is the difference between the highest value and the lowest value

### highest \_ lowest value value

## Definition



## The standard deviation of a set of sample values is a measure of variation of values about the mean

#### Sample Standard Deviation Formula





Formula 2-4 (p.75)

## Slide 63 (Shortcut Formula)



Formula 2-5 (p.75)

#### Standard Deviation -Key Points (p.75)



The standard deviation is a measure of variation of all values from the mean

The value of the standard deviation s is usually positive

The value of the standard deviation s can increase dramatically with the inclusion of one or more outliers (data values far away from all others)

The units of the standard deviation s are the same as the units of the original data values



This formula is similar to Formula 2-4, but instead the population mean and population size are used

## **Definition (p.78)**



The variance of a set of values is a measure of variation equal to the square of the standard deviation.

Sample variance: Square of the sample standard deviation s

Population variance: Square of the population standard deviation





### standard deviation squared





#### Round-off Rule for Measures of Variation (p.79)

# Carry one more decimal place than is present in the original set of data.

Round only the final answer, not values in the middle of a calculation.

## **Definition (p.79)**



## The coefficient of variation (or CV) for a set of sample or population data, expressed as a percent, describes the standard deviation relative to the mean

Sample

Population

$$cv = \frac{s}{x} \cdot 100\%$$

$$cv = \frac{\sigma}{\mu} \cdot 100\%$$

### Standard Deviation from a Slide 70 Frequency Distribution

Formula 2-6  
(p.80)  
$$\int \frac{n [\Sigma(f \cdot x^{2})] - [\Sigma(f \cdot x)]^{2}}{n (n - 1)}$$

#### Use the class midpoints as the x values

#### Estimation of Standard Deviation Range Rule of Thumb (p.82)



## For estimating a value of the standard deviation s, Use $S \approx \frac{Range}{4}$

#### Where range = (highest value) – (lowest value)

#### Estimation of Standard Deviation Range Rule of Thumb (p.82)

## For interpreting a known value of the standard deviation s, find rough estimates of the minimum and maximum "usual" values by using:

Minimum "usual" value  $\approx$  (mean) – 2 X (standard deviation)

Maximum "usual" value  $\approx$  (mean) + 2 X (standard deviation)

Compare this definition with the discussion on usual/unusual on p.94.
# **Definition (p.83)**



#### **Empirical (68-95-99.7) Rule**

For data sets having a distribution that is approximately bell shaped, the following properties apply:

# About 68% of all values fall within 1 standard deviation of the mean

# About 95% of all values fall within 2 standard deviations of the mean

About 99.7% of all values fall within 3 standard deviations of the mean

## **The Empirical Rule**





Chapter 2, Triola, Elementary Statistics, MATH 1342

# **Definition (p.85)**



**Chebyshev's Theorem** 

The proportion (or fraction) of any set of data lying within K standard deviations of the mean is always at least  $1-1/K^2$ , where K is any positive number greater than 1.

- For K = 2, at least 3/4 (or 75%) of all values lie within 2 standard deviations of the mean
- For K = 3, at least 8/9 (or 89%) of all values lie within 3 standard deviations of the mean

Note: named after Pafnuty Lvovich Chebyshev (1821-1894), Russian mathematician

### **Rationale for Formula 2-4**



The end of Section 2-5 has a detailed explanation of why Formula 2-4 is employed instead of other possibilities and, specifically, why n - 1 rather than n is used. The student should study it carefully

# **Recap of Section 2-5**



In this section we have looked at:

- Range
- Standard deviation of a sample and population
- Variance of a sample and population
- Coefficient of Variation (CV)
- Standard deviation using a frequency distribution
- Range Rule of Thumb
- Empirical Distribution
- Chebyshev's Theorem

# Section 2-6 Measures of Relative Standing

Created by Tom Wegleitner, Centreville, Virginia



**Definition (p.92)** 



# **Z Score** (or standard score)

#### the number of standard deviations that a given value **x** is above or below the mean.



# Sample Population

# $z = \frac{x - \bar{x}}{s} \qquad \qquad z = \frac{x - \mu}{\sigma}$

### **Round to 2 decimal places**

# Interpreting Z Scores Slide 81

#### FIGURE 2-14 (p.94)



Whenever a value is less than the mean, its corresponding z score is negative

Ordinary values: *z* score between –2 and 2 sd Unusual Values: *z* score < -2 or *z* score > 2 sd

# **Definition (p.94)**



- ♦ Q<sub>1</sub> (First Quartile) separates the bottom
  25% of sorted values from the top 75%.
- Q<sub>2</sub> (Second Quartile) same as the median; separates the bottom 50% of sorted values from the top 50%.
- ♦ Q<sub>1</sub> (Third Quartile) separates the bottom
  75% of sorted values from the top 25%.





# $Q_1, Q_2, Q_3$ divides ranked scores into four equal parts $\begin{pmatrix} 25\% & 25\% & 25\% & 25\% \\ (minimum) & Q_1 & Q_2 & Q_3 \end{pmatrix}$ (maximum)

(median)

# **Percentiles (p.95)**



# Just as there are quartiles separating data into four parts, there are 99 percentiles denoted $P_1, P_2, \ldots P_{99}$ , which partition the data into 100 groups.





Percentile of value x =



#### Converting from the *k*th Percentile to the Corresponding Data Value (p.96)

#### Notation

$$L=\frac{k}{100}\boldsymbol{\cdot} n$$

- *n* total number of values in the data set
- *k* percentile being used
- *L* locator that gives the *position* of a value
- $P_k$  kth percentile



#### Converting from the *k*th Percentile to the Corresponding Data Value

The value of the kth percentile is midway between the Lth value and the next value in the sorted set of data. Find  $P_k$  by adding the Lth value and the next value and dividing the total by 2.

#### Figure 2-15 (p.96)







\* Interquartile Range (or IQR):  $Q_3 - Q_1$  $Q_3 - Q_1$ Semi-interquartile Range: 2  $Q_{3} + Q_{1}$ Midguartile:  $\Rightarrow$  10 - 90 Percentile Range:  $P_{10}$  -  $P_{10}$ 

#### **Recap of Section 2-6**



In this section we have discussed:

- ✤ z Scores
- z Scores and unusual values
- Quartiles
- Percentiles
- Converting a percentile to corresponding data values
- Other statistics

# Section 2-7 Exploratory Data Analysis (EDA)

Created by Tom Wegleitner, Centreville, Virginia



# Definition (p.102) Slide 91

Exploratory Data Analysis is the process of using statistical tools (such as graphs, measures of center, and measures of variation) to investigate data sets in order to understand their important characteristics

# **Definition (p.102)**



#### An outlier is a value that is located very far away from almost all the other values





# An outlier can have a dramatic effect on the mean

# An outlier have a dramatic effect on the standard deviation

# An outlier can have a dramatic effect on the scale of the histogram so that the true nature of the distribution is totally obscured

# **Definitions (p.104)**



- For a set of data, the 5-number summary consists of the minimum value; the first quartile Q<sub>1</sub>; the median (or second quartile Q<sub>2</sub>); the third quartile, Q<sub>3</sub>; and the maximum value
- A boxplot ( or box-and-whisker-diagram) is a graph of a data set that consists of a line extending from the minimum value to the maximum value, and a box with lines drawn at the first quartile, Q<sub>1</sub>; the median; and the third quartile, Q<sub>3</sub>

# **Boxplots**





Cotinine Level of Smokers

#### Figure 2-16 (p.105)







#### Figure 2-17 (p.105)



# **Recap of Section 2-7**

In this section we have looked at:

- Exploratory Data Analysis
- Effects of outliers
- 5-number summary and boxplots