

# MSSC 6010/Comp. Probability

**Instructor: Mehdi Maadooliat**

**MSSC 6010**



**Department of Mathematical and Statistical Sciences**

# SYLLABUS - D2L - ONLINE MATERIALS

- **Syllabus:** <http://tinyurl.com/Comp-Prob/course-syllabus.html>
- **Course Website:** <http://tinyurl.com/Comp-Prob>
- **Homework System:** <http://d2l.mu.edu>

Course information > Overview

## Course overview

This is the homepage for MSSC 6010 - Computational Probability by Dr. Mehdi Maadooliat in Fall 2025 at Marquette University. All course materials will be posted on this site.

You can find the course syllabus [here](#) and the course schedule [here](#).

### Class meetings

Meeting	Location	Time
Lecture	CU 208	Tue & Thur 11:00 am - 12:15 pm
Office Hours	T: CU 351 & Th: Stat Help Desk	Tue & Thur 12:15 - 1:30 pm





**Computational Probability**

Dr. Mehdi Maadooliat  
Marquette University  
MSSC 6010 - Fall 2025



## D2L: Marquette University's Learning Management System

Use your CheckMarq username and password to log in. Trouble logging in? [Contact the Help Desk](#).

Username \*

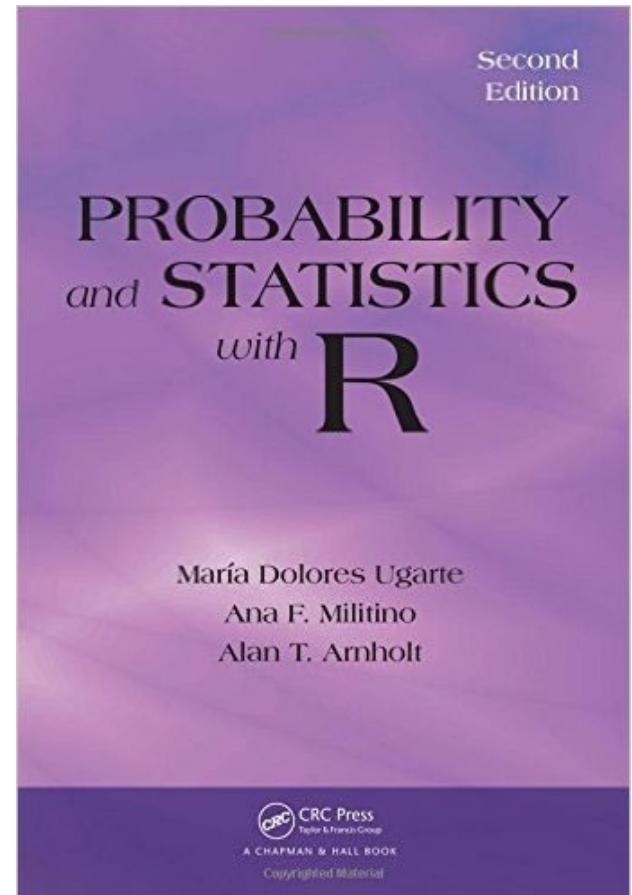
Password \*

**Log In**



# MSSC 6010

- **Text: Probability and Statistics with R, 2<sup>nd</sup> ed.**  
**(ISBN: 9781466504394) Maria Dolores Ugarte, Ana F. Militino, Alan T. Arnholt,**
- **Instructor: Mehdi Maadooliat, Ph.D.**
  - **Office Hours:** T 12:15 - 1:30pm  
**(Cudahy Hall Room 351)**
    - » Th 12:15 - 1:30pm  
**(Statistics Help Desk)**
  - **Office:** Cudahy Hall 351
  - **Email:** [mehdi.maadooliat@marquette.edu](mailto:mehdi.maadooliat@marquette.edu)
- **Homework:** Weekly or biweekly homework assignments will be given.  
11



## MSSC 6010 CONT...

- **Description:** Foundations of probability for modeling random processes and Bayesian approaches, including: counting techniques, probability of events, random variables, distribution functions, probability functions, probability density functions, expectation, moments, moment generating functions, special discrete and continuous distributions, sampling distributions, prior and posterior distributions, Law of Large Numbers, Central Limit Theorem, Bayesian paradigm.
- **Prereq:** Three semesters of mathematics beyond calculus and MATH 4720 or equiv.

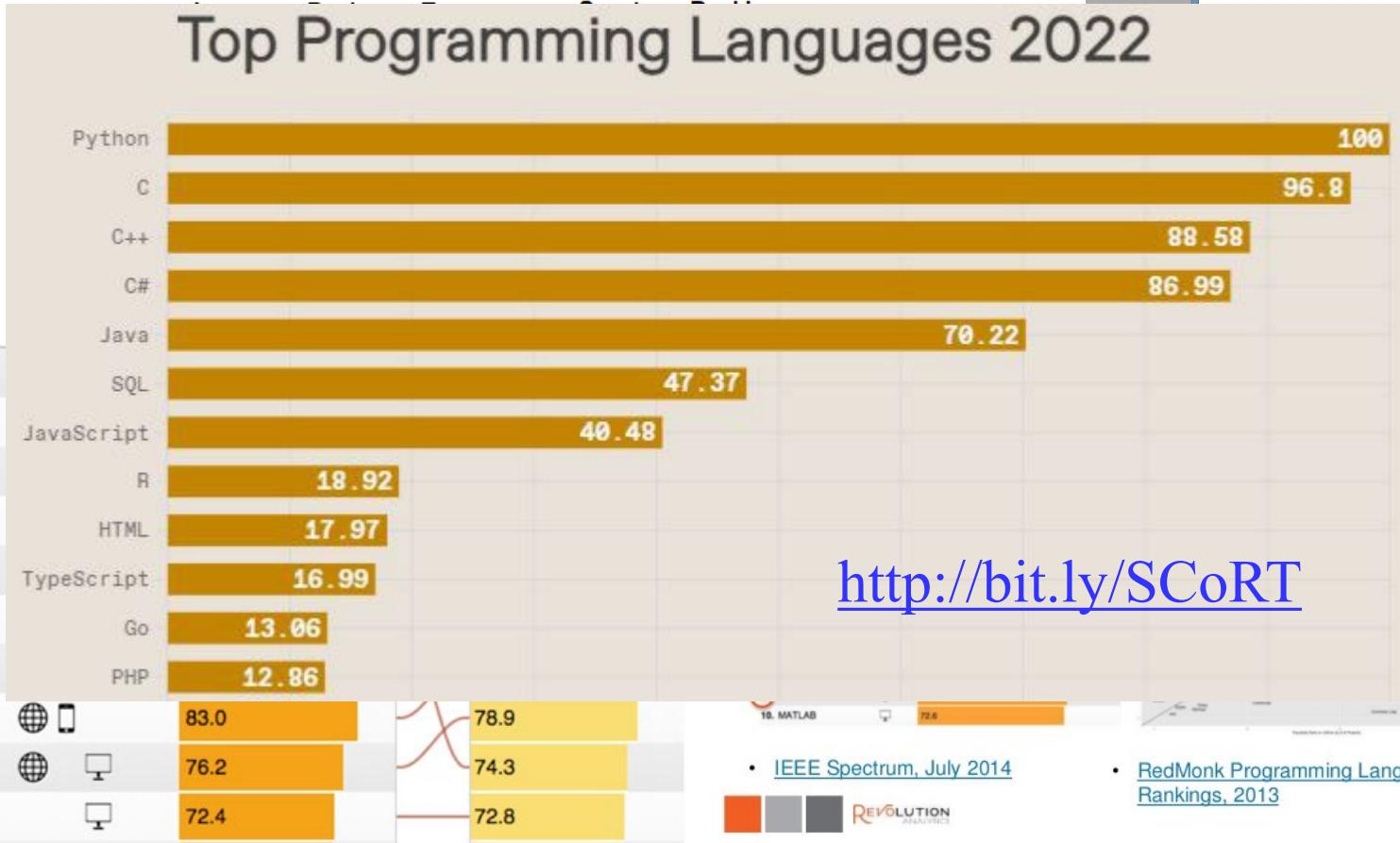
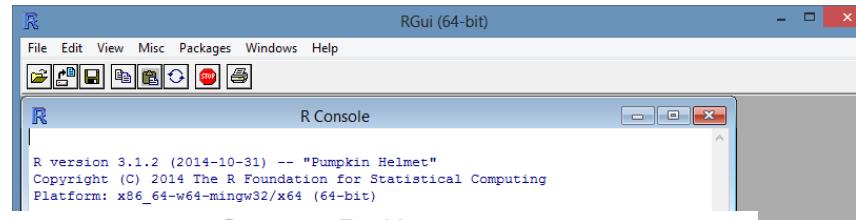
# WHAT DO WE COVER IN THIS COURS?

- “**Probability and Statistics with R, 1st ed.**” Covers:
  - **1. A Brief Introduction to S**
  - **2. Exploring Data**
  - **3. General Probability and Random Variables**
  - **4. Univariate Probability Distributions**
  - **5. Multivariate Probability Distributions**
  - **6. Sampling and Sampling Distributions**
  - **7. Point Estimation**
  - **8. Confidence Intervals**
  - **9. Hypothesis Testing**
  - **10. Nonparametric Methods**
  - **11. Experimental Design**
  - **12. Regression**
- I may cover some of the following topics
  - **Bayesian Statistics**
  - **Transformation of variable**
  - **Hypothesis Testing (LRT)**

# R IS THE STATISTICAL SOFTWARE THAT WE MOSTLY USE IN THIS COURSE

- It is free!!

[IEEE Spectrum, July 2016](#)



[IEEE Spectrum, July 2015](#)

# HISTORY OF R (FROM [WIKIPEDIA](#))

- R is a [programming language](#) and software environment for [statistical computing](#) and graphics. The R language is widely used among [statisticians](#) and [data miners](#) for developing [statistical software](#) and data analysis. Polls, [surveys of data miners](#), and studies of scholarly literature databases show that R's popularity has increased substantially in recent years.
- R is an implementation of the [S programming language](#) combined with [lexical scoping](#) semantics inspired by [Scheme](#). [S](#) was created by [John Chambers](#) while at [Bell Labs](#). There are some important differences, but much of the code written for S runs unaltered.
- R was created by [Ross Ihaka](#) and [Robert Gentleman](#) at the [University of Auckland](#), New Zealand, and is currently developed by the *R Development Core Team*, of which Chambers is a member. R is named partly after the first names of the first two R authors and partly as a play on the name of [S](#).
- R is a [GNU project](#). The [source code](#) for the R software environment is written primarily in [C](#), [Fortran](#), and R. R is freely available under the [GNU General Public License](#), and pre-compiled binary versions are provided for various [operating systems](#). R uses a [command line interface](#); there are also several [graphical front-ends](#) for it.

# MSSC 6010 CONT...

- **Grading:**

- **Weekly homework & class participation**
- **Project**
- **A midterm (in class), tentatively on Oct 16,**
- **plus the final exam (in class or take home)**  
**on Dec 8, 10:30AM - 12:30PM.**
- **No make-ups, If “unavoidable absence” as defined in Arts and Sci. Bulletin, then % added to final %.**
- **Midterm Exams 30%**
- **Project 10%**
- **Homework & Class Participation 30%**
- **Final Exam 30%**

# MSSC 6010 CONT...

- **Questions about Homework, Project and Exams:**
  - **SHOULD be posted in d2l Discussion Board.**
  - **I will NOT answer general emails about Homework and Exams.**

- **Homework:**
  - **Should be submitted as a PDF file:**
    - **How to Combine Images into a PDF file [FREE & EASY + No Software] ([Youtube](#))**
    - **Microsoft Word to PDF in 10 Seconds ([Youtube](#))**
    - **How to: convert Images to PDF in Macbook/iMac ([Youtube](#))**
      - <http://apple.stackexchange.com/questions/11163/how-do-i-combine-two-or-more-images-to-get-a-single-pdf-file>

MATH 1700 102 Mod Elementary Stat

Course Home | Content | Discussions | Dropbox | Quizzes | Classlist | Grades

**Discussions List** Subscriptions Group and Section Restrictions Statistics

New More Actions

Filter by: Unread Unapproved

**Homeworks** Hide Topics for Homeworks

Topic	Threads	Posts
Homework 1	0	0
Homework 2	0	0
Homework 3	0	0
Homework 4	0	0
Homework 5	0	0
Homework 6	0	0
Homework 7	0	0
Homework 8	0	0
Project 1	0	0

# MSSC 6010/Comp. Probability

Instructor: Mehdi Maadooliat

Chapters 2 Review

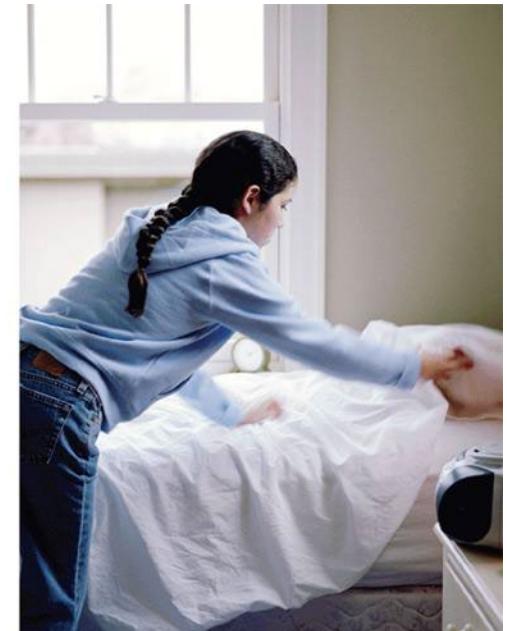


Department of Mathematical and Statistical Sciences

# 1. STATISTICS

## 1.1

### What Is Statistics?



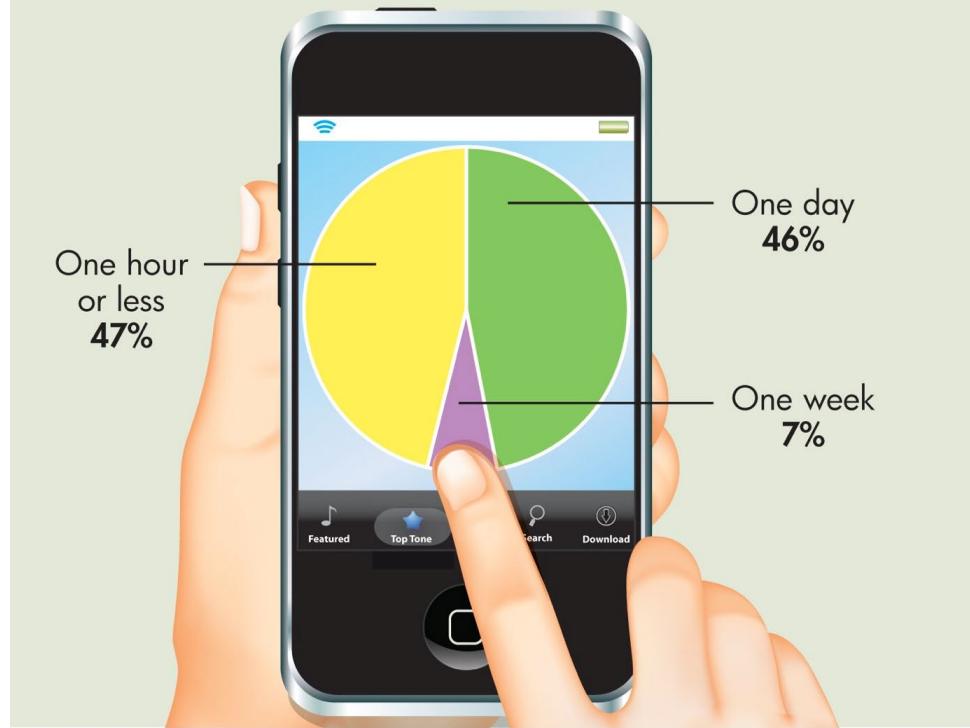
# WHAT IS STATISTICS?

- **Statistics is all around us!**
- **How much time between Internet usage?**

## Fretting Over Messages

### Are you fretting about messages?

How Wi-Fi users responded when asked how long they go before they get "antsy" about checking e-mail, instant messaging and social networking sites:



**Source:** Impulse Research for Qwest Communications online survey of 1,063 adult Wi-Fi users in April 2009.



# WHAT IS STATISTICS?

- **Statistics is all around us!**

## World's Highest Paid Athletes



*Forbes'* list of the highest-paid athletes looks at earnings derived from salaries, bonuses, prize money, endorsements, and licensing income between June 2008 and June 2009 and does not deduct for taxes or agents' fees.

Here are the Top 5:

Rank	Athlete	Sport	Earnings
1	Tiger Woods	Golf	\$110 million
2	Kobe Bryant	Basketball	\$45 million
2	Michael Jordan	Basketball	\$45 million
2	Kimi Raikkonen	Auto Racing	\$45 million
5	David Beckham	Soccer	\$42 million



# WHAT IS STATISTICS?

## Many Teens Use Phones in Class

84% teens have cell phones

16% teens do not

An average of 440 text messages sent per week – 110 of them during class. Works out to be 3 text messages per class period.

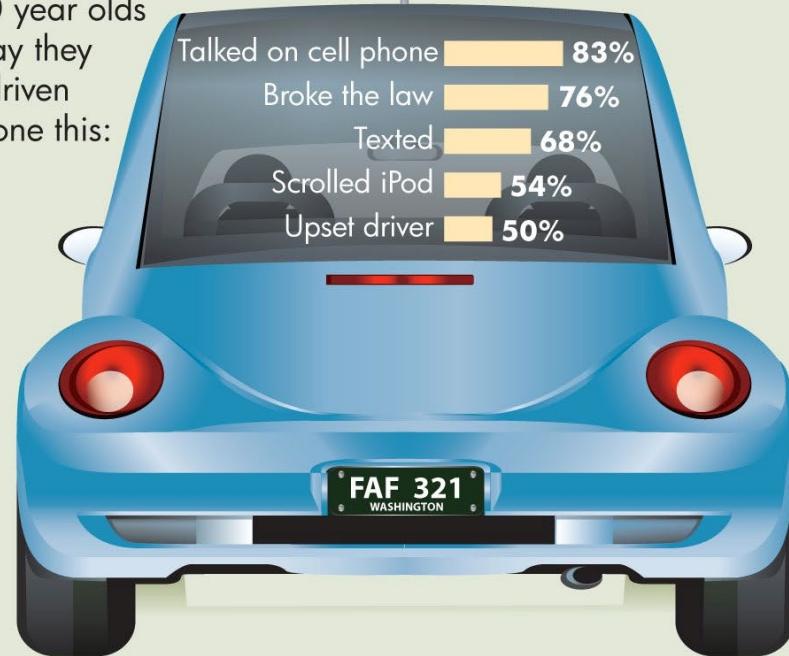
**Source:**

Common Sense Media survey of 1013 teens May/June 2009

## Busy Behind the Wheel

**Most drivers ages 16–20 admit to risky driving habits.**

16–20 year olds who say they have driven and done this:



**Source:** National Organization for Youth Safety, Allstate Foundation online survey of 605 drivers ages 16–20. (6/16/09)

# WHAT IS STATISTICS?

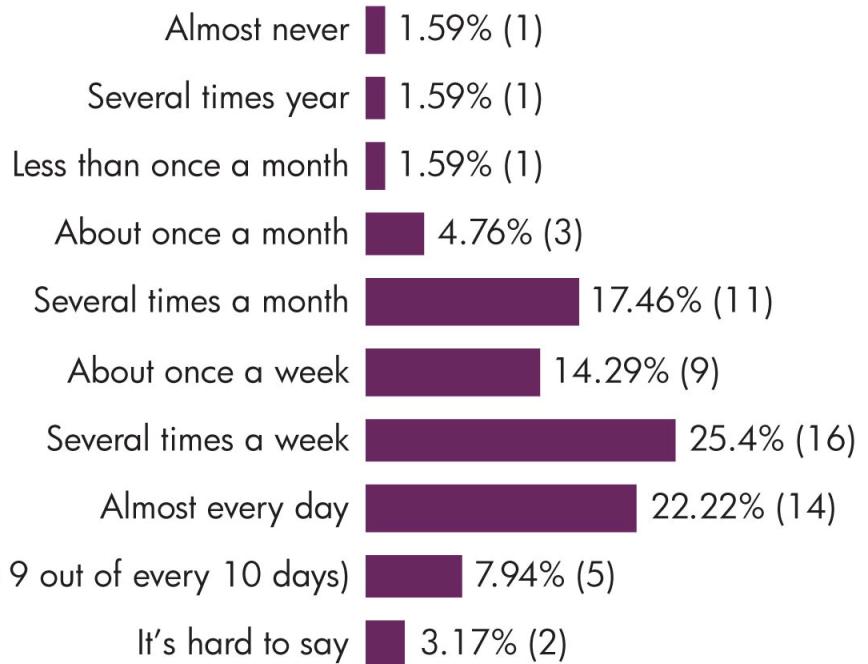
- **Statistics is all around us!**

Postyour.info is a worldwide service where Internet users from arround the world can take part in questionnaires. [<http://postyour.info/>] Below is a graph depicting the combined summary of how users answered one of the posted questions. Results given in Percent (count).

- **How often do you eat fruit?**

## How Often Do You Eat Fruit?

(irrespective of the reasons why)



# WHAT IS STATISTICS?

- **Statistics has become the universal language of the sciences. As potential users of statistics, we need to master both the “science” and the “art” of using statistical methodology correctly.**
- **Careful use of statistical methods will enable us to obtain accurate information from data. These methods include**
  - **(1) carefully defining the situation,**
  - **(2) gathering data,**
  - **(3) accurately summarizing the data, and**
  - **(4) deriving and communicating meaningful conclusions.**

# WHAT IS STATISTICS?

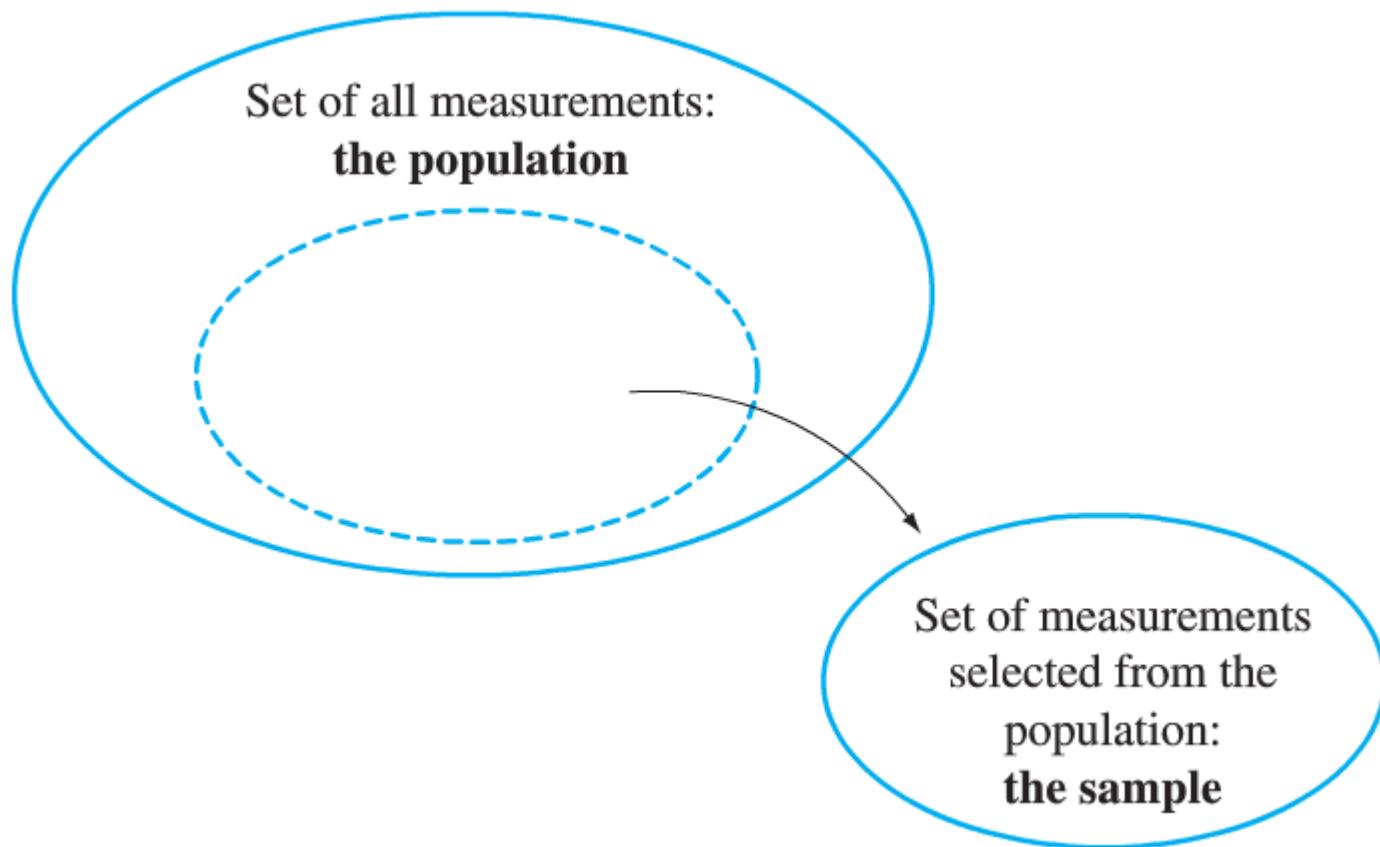
- The field of statistics can be divided into two main branches.
- Descriptive statistics is what most people think of when they hear the word *statistics*. It includes the collection, presentation, and description of sample data.
- The term inferential statistics refers to the technique of interpreting the values resulting from the descriptive techniques and making decisions and drawing conclusions about the population.

# WHAT IS STATISTICS?

- **Statistics** The science of collecting, describing, and interpreting data.
- **Population** A collection, or set, of individuals, objects, or events whose properties are to be analyzed.
  - The set of “all students who have ever attended a U.S. college” is an example of a well-defined population.
- **Sample** A subset of a population.
- **Variable** A characteristic of interest about each individual element of a population or sample.
  - A student’s age at entrance into college, the color of the student’s hair, the student’s height, and the student’s weight are four variables.
- **Data value** The value of the variable associated with one element of a population or sample. This value may be a number, a word, or a symbol.

# POPULATION VS. SAMPLE

- **Population:** The entire group of interest
- **Sample:** A part of the population selected to draw conclusions about the entire population



# WHAT IS STATISTICS?

- **Data** The set of values collected from the variable from each of the elements that belong to the sample.
  - The set of 25 heights collected from 25 students is an example of a set of data.
- **Experiment** A planned activity whose results yield a set of data.
  - An experiment includes the activities for both selecting the elements and obtaining the data values.
- **Parameter** A numerical value summarizing all the data of an entire population.
  - A parameter is a value that describes the entire population. Often a Greek letter is used to symbolize the name of a parameter.
- **Statistic** A numerical value summarizing the sample data.
  - Often letters of the English alphabet is used to symbolize the name of a statistic.

# WHAT IS STATISTICS?

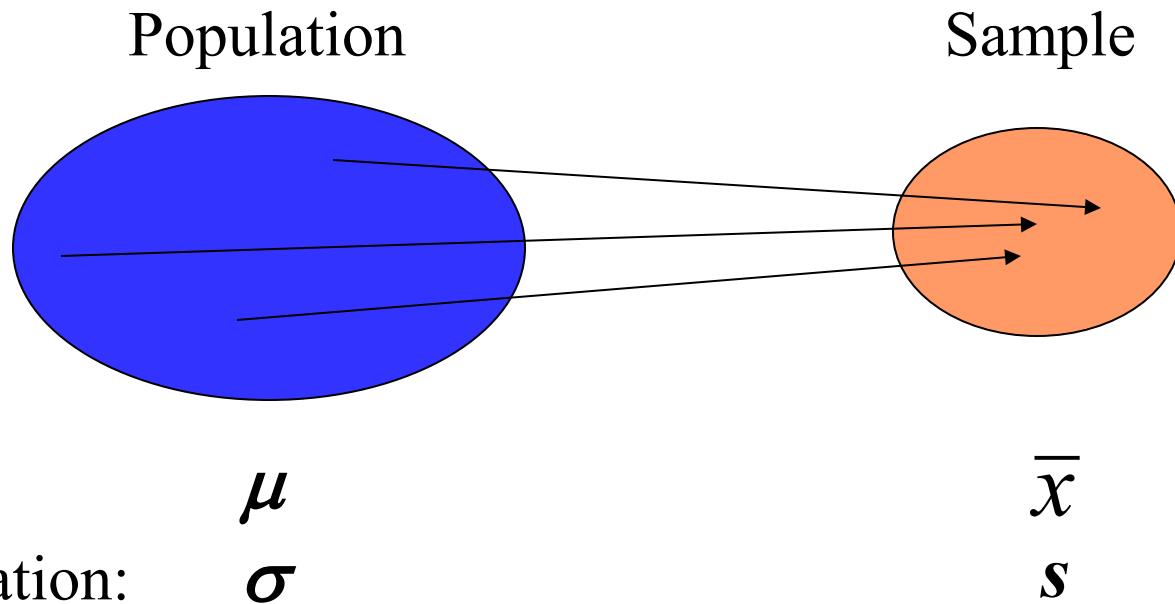
- **Describing a Population**
  - It is common practice to use Greek letters when talking about a population.
  - We call the mean of a population  $\mu$ .
  - We call the standard deviation of a population  $\sigma$  and the variance  $\sigma^2$ .
  - It is important to know that for a given population there is only one true mean and one true standard deviation and variance or one true proportion.
  - There is a special name for these values: **parameters**.

# WHAT IS STATISTICS?

- **Describing a Sample**
  - We call the mean of a sample  $\bar{x}$ .
  - We call the standard deviation of a sample  $s$  and the variance  $s^2$ .
  - There are many different possible samples that could be taken from a given population. For each sample there may be a different mean, standard deviation, variance, or proportion.
  - There is a special name for these values: **statistics**.

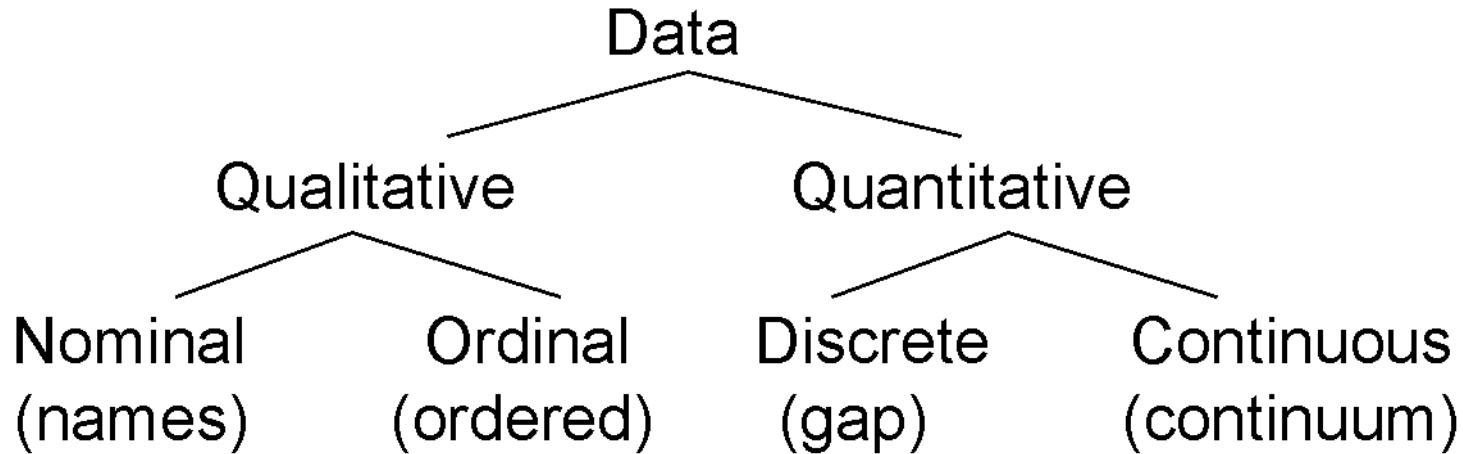
# POPULATION VS SAMPLE

- We use sample statistics to make inference about population parameters



# WHAT IS STATISTICS?

- **Data** The set of values collected from the variable from each of the elements that belong to the sample.



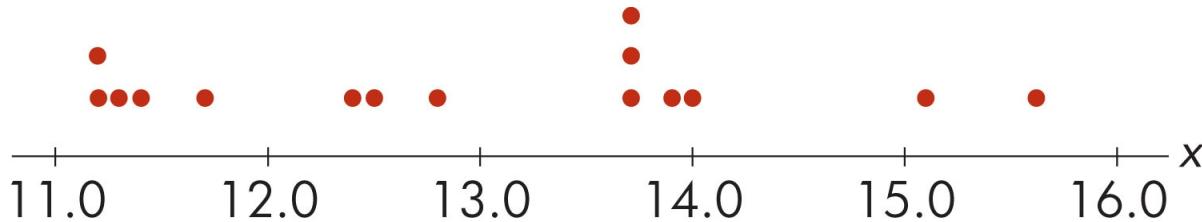
# WHAT IS STATISTICS?

- **Qualitative (Categorical) variable:** A variable that describes or categorizes an element of a population.
  - **Nominal variable:** A qualitative variable that characterizes an element of a population. No ordering. No arithmetic.
  - **Ordinal variable:** A qualitative variable that incorporates an ordered position, or ranking.
- **Quantitative (Numerical) variable:** A variable that quantifies an element of a population.
  - **Discrete variable:** A quantitative variable that can assume a countable number of values. Gap between successive values.
  - **Continuous variable:** A quantitative variable that can assume an uncountable number of values. Continuum of values.

# TYPES OF VARIABLES

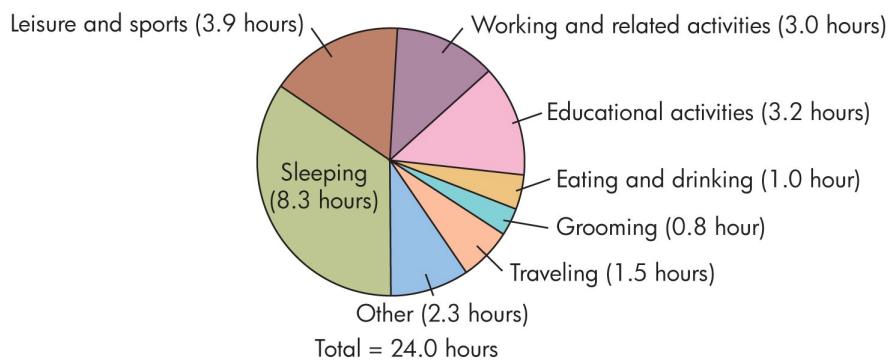
## Examples:

Variable	Numeric		Categorical	
	Discrete	Continuous	Nominal	Ordinal
Weight		X		
Hours Enrolled	X			
Major			X	
Zip Code				X

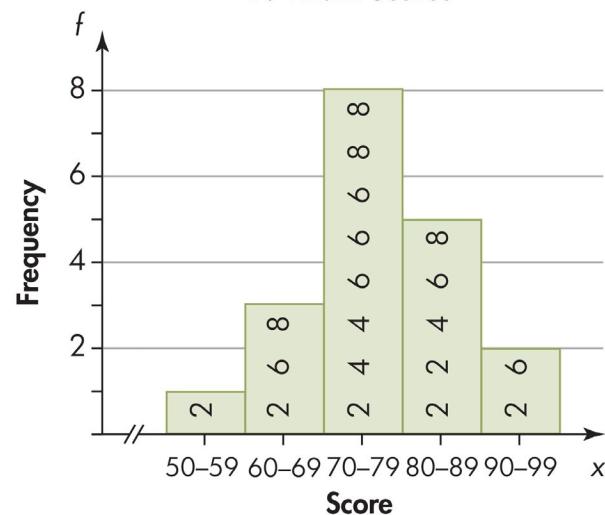


# Graphical Summaries

Time Use on an Average Weekday for Full-time University and College Students



19 Exam Scores



# EXAMPLE 1 - GRAPHING QUALITATIVE DATA

- **Table 2.1 lists the number of cases of each type of operation performed at General Hospital last year.**

Type of Operation	Number of Cases
Thoracic	20
Bones and joints	45
Eye, ear, nose, and throat	58
General	98
Abdominal	115
Urologic	74
Proctologic	65
Neurosurgery	23
<i>Total</i>	498

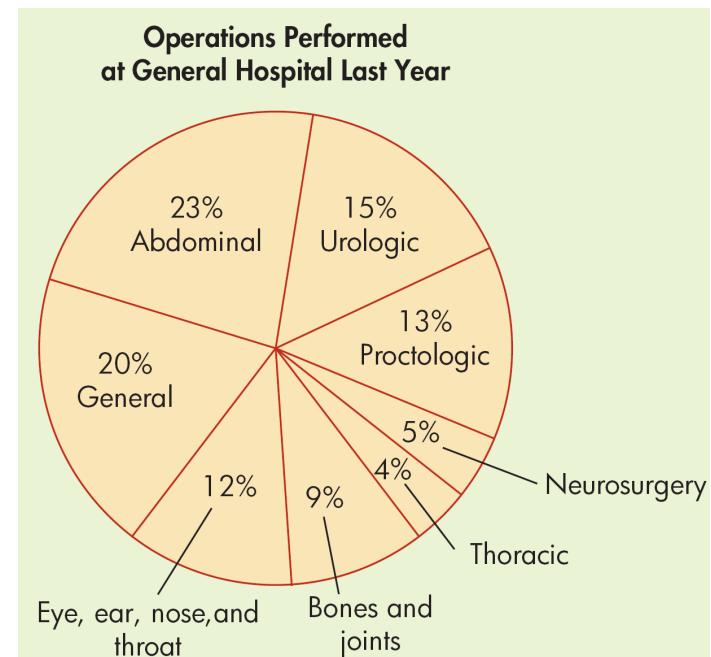
Operations Performed at General Hospital Last Year [TA02-01]

**Table 2.1**

# QUALITATIVE DATA (PIE CHART)

- **Pie charts (circle graphs) and bar graphs** Graphs that are used to summarize qualitative, or attribute, or categorical data.
- **Pie charts (circle graphs) show the amount of data that belong to each category as a proportional part of a circle.**

Type of Operation	Number of Cases
Thoracic	20
Bones and joints	45
Eye, ear, nose, and throat	58
General	98
Abdominal	115
Urologic	74
Proctologic	65
Neurosurgery	23
<i>Total</i>	498

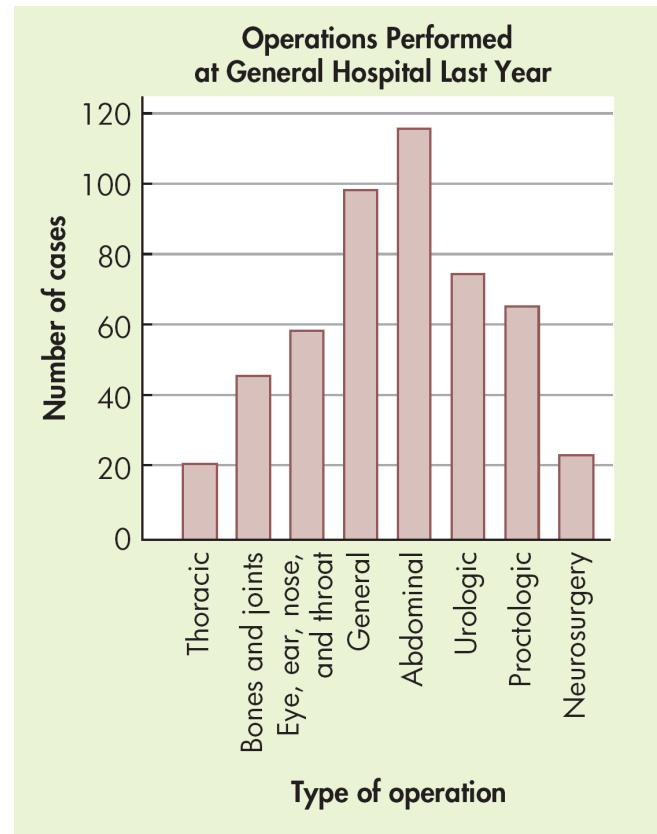


# QUALITATIVE DATA (BAR GRAPH)

- Bar graphs show the amount of data that belong to each category as a proportionally sized rectangular area.**

Type of Operation	Number of Cases
Thoracic	20
Bones and joints	45
Eye, ear, nose, and throat	58
General	98
Abdominal	115
Urologic	74
Proctologic	65
Neurosurgery	23
<i>Total</i>	498

- Bar graphs of attribute data should be drawn with a space between bars of equal width.**



# EXAMPLE 2-

## AUSTRALIAN INSTITUTE OF SPORT DATA

- **Description**

- Data on 102 male and 100 female athletes collected at the Australian Institute of Sport, courtesy of Richard Telford and Ross Cunningham.

- **Source**

- Cook and Weisberg (1994), *An Introduction to Regression Graphics*. John Wiley & Sons, New York.

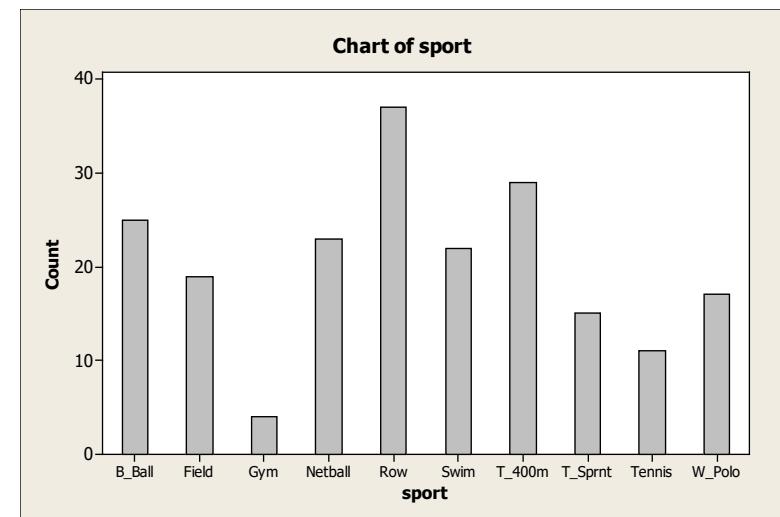
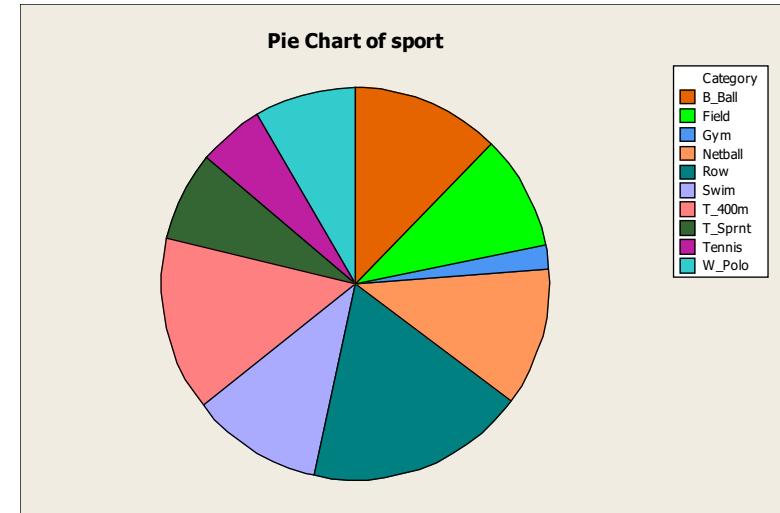
AIS.mjp

Variable	Description
sex	sex
sport	sport
rcc	red cell count
wcc	white cell count
Hc	Hematocrit
Hg	Hemoglobin
Fe	plasma ferritin concentration
bmi	body mass index, weight/(height)
ssf	sum of skin folds
Bfat	body fat percentage
Ibm	lean body mass
Ht	height (cm)
Wt	weight (Kg)

# SUMMARIZING A SINGLE CATEGORICAL VARIABLE

- **Frequency (Count)** - number of times the value occurs in the data
- **Relative frequency (Percent)** - proportion of the data with the value
- **ais.xls (D2L/Content/Datasets)**

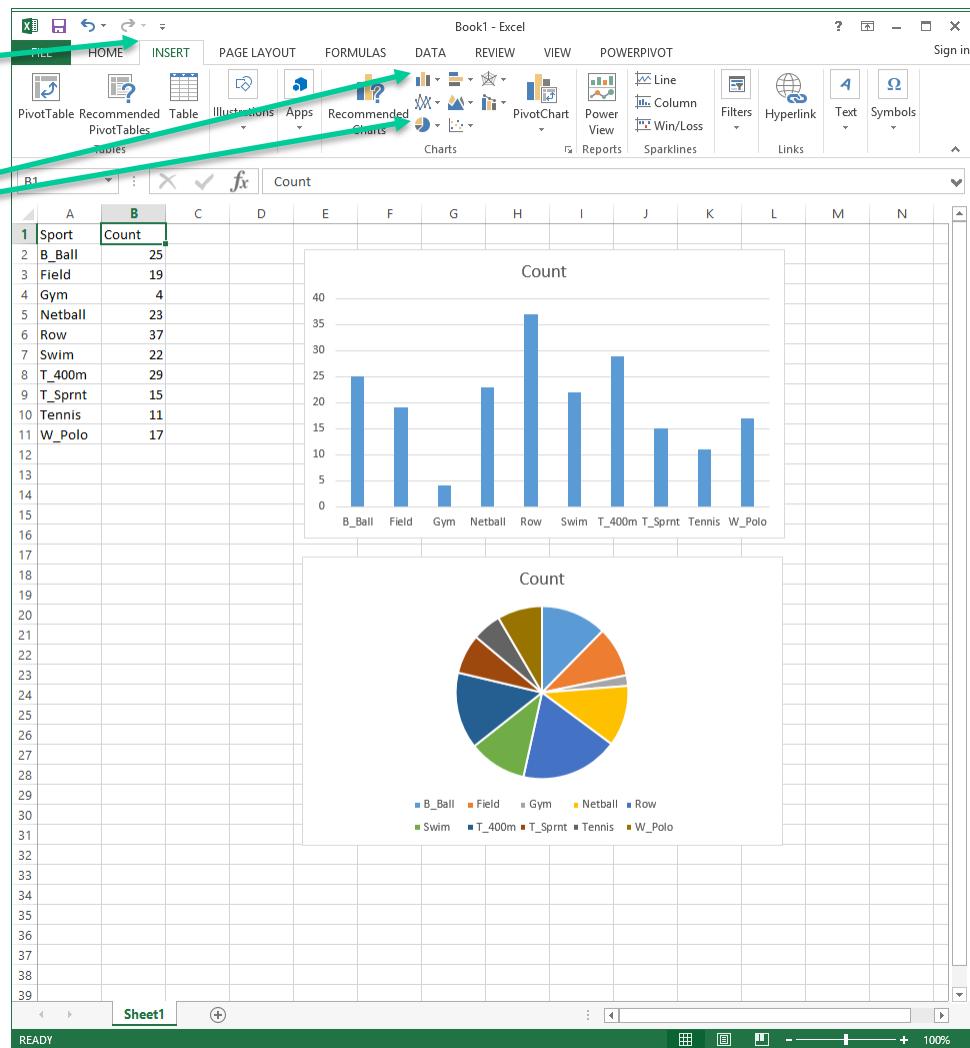
sport	Count	Percent
B_Ball	25	12.38
Field	19	9.41
Gym	4	1.98
Netball	23	11.39
Row	37	18.32
Swim	22	10.89
T_400m	29	14.36
T_Sprnt	15	7.43
Tennis	11	5.45
W_Polo	17	8.42
N=	202	





# HOW TO?

- Enter the Data in Excel:
- Select Insert
- Select Pie or Bar Charts



# GRAPHING QUANTITATIVE DATA

- **Distribution** The pattern of variability displayed by the data of a variable. The distribution displays the frequency of each value of the variable.
- **Dotplot display** Displays the data of a sample by representing each data value with a dot positioned along a scale. The frequency of the values is represented along the other scale.

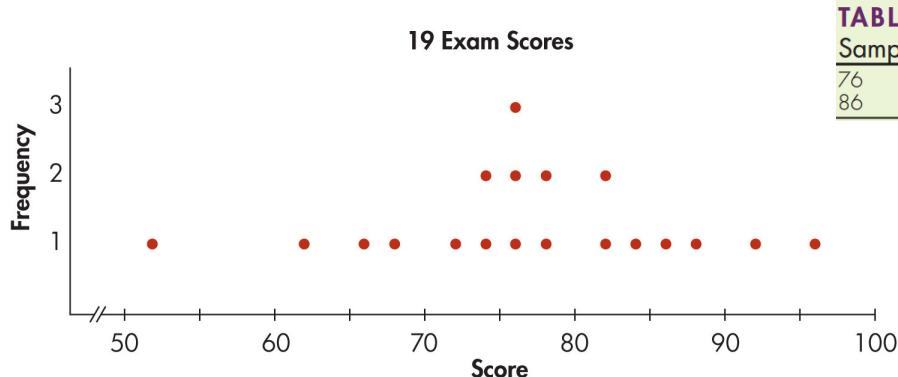


Table 2.2 provides a sample of 19 exam grades randomly selected from a large class.

**TABLE 2.2**  
Sample of 19 Exam Grades [TA02-02]

76	74	82	96	66	76	78	72	52	68
86	84	62	76	78	92	82	74	88	

# GRAPHING QUANTITATIVE DATA

- Stem-and-leaf display** Displays the data of a sample using the actual digits that make up the data values. Each numerical value is divided into two parts: The **leading digit(s)** becomes the **stem**, and the **trailing digit(s)** becomes the **leaf**.

Table 2.2 provides a sample of 19 exam grades randomly selected from a large class.

**TABLE 2.2**

Sample of 19 Exam Grades [TA02-02]

76	74	82	96	66	76	78	72	52	68
86	84	62	76	78	92	82	74	88	

**19 Exam Scores**

5	2
6	6 8 2
7	6 4 6 8 2 6 8 4
8	2 6 4 2 8
9	6 2

**Or**

**19 Exam Scores**

5	2
6	2 6 8
7	2 4 4 6 6 6 8 8
8	2 2 4 6 8
9	2 6



## EXAMPLE 3

# OVERLAPPING DISTRIBUTIONS

- A random sample of 50 college students was selected. Their weights were obtained from their medical records. The resulting data are listed in Table 2.3.**

Student	1	2	3	4	5	6	7	8	9	10
Male/Female	F	M	F	M	M	F	F	M	M	F
Weight	98	150	108	158	162	112	118	167	170	120
Student	11	12	13	14	15	16	17	18	19	20
Male/Female	M	M	M	F	F	M	F	M	M	F
Weight	177	186	191	128	135	195	137	205	190	120
Student	21	22	23	24	25	26	27	28	29	30
Male/Female	M	M	F	M	F	F	M	M	M	M
Weight	188	176	118	168	115	115	162	157	154	148
Student	31	32	33	34	35	36	37	38	39	40
Male/Female	F	M	M	F	M	F	M	F	M	M
Weight	101	143	145	108	155	110	154	116	161	165
Student	41	42	43	44	45	46	47	48	49	50
Male/Female	F	M	F	M	M	F	F	M	M	M
Weight	142	184	120	170	195	132	129	215	176	183

Weights of 50 College Students [TA02-03]

Table 2.3

## EXAMPLE 3

# OVERLAPPING DISTRIBUTIONS

- Notice that the weights range from 98 to 215 pounds. Let's group the weights on stems of 10 units using the hundreds and the tens digits as stems and the units digit as the leaf (see Figure 2.7).
- The leaves have been arranged in numerical order. Close inspection of Figure 2.7 suggests that two overlapping distributions may be involved.

Weights of 50 College Students (lb)	
$N = 50$	Leaf Unit = 1.0
9	8
10	1 8 8
11	0 2 5 5 6 8 8
12	0 0 0 8 9
13	2 5 7
14	2 3 5 8
15	0 4 4 5 7 8
16	1 2 2 5 7 8
17	0 0 6 6 7
18	3 4 6 8
19	0 1 5 5
20	5
21	5

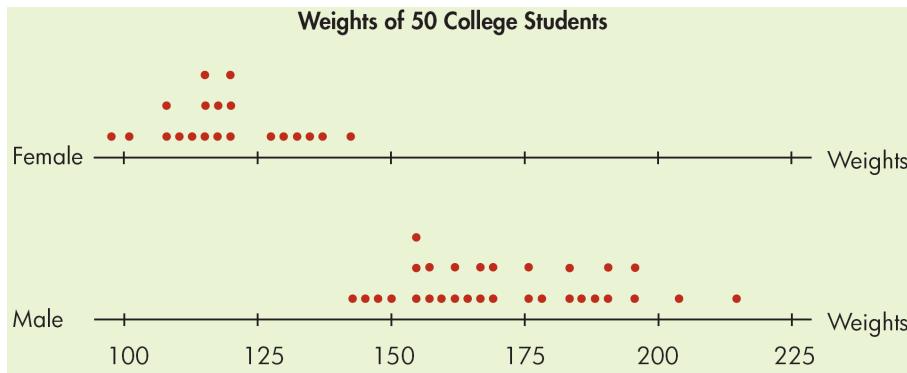
Stem-and-Leaf Display  
**Figure 2.7**

## EXAMPLE 3

# OVERLAPPING DISTRIBUTIONS

cont'd

- That is exactly what we have: a distribution of female weights and a distribution of male weights.
- Figure 2.8 shows a “back-to-back” stem-and-leaf display of this set of data and makes it obvious that two distinct distributions are involved.



Weights of 50 College Students (lb)		
Female		Male
	8	09
1 8 8	10	
0 2 5 5 6 8 8	11	
0 0 0 8 9	12	
2 5 7	13	
2	14	3 5 8
	15	0 4 4 5 7 8
	16	1 2 2 5 7 8
	17	0 0 6 6 7
	18	3 4 6 8
	19	0 1 5 5
	20	5
	21	5

“Back-to-Back” Stem-and-Leaf Display

Figure 2.8

# FREQUENCY DISTRIBUTIONS AND HISTOGRAMS

- **Frequency distribution** A listing, often expressed in chart form, that pairs values of a variable with their frequency.
- Let's use a sample of 50 final exam scores taken from last semester's elementary statistics class.

60	47	82	95	88	72	67	66	68	98	90	77	86
58	64	95	74	72	88	74	77	39	90	63	68	97
70	64	70	70	58	78	89	44	55	85	82	83	
72	77	72	86	50	94	92	80	91	75	76	78	

Statistics Exam Scores [TA02-06]

Table 2.6

# FORM A FREQUENCY DISTRIBUTION AND HISTOGRAM

60	47	82	95	88	72	67	66	68	98	90	77	86
58	64	95	74	72	88	74	77	39	90	63	68	97
70	64	70	70	58	78	89	44	55	85	82	83	
72	77	72	86	50	94	92	80	91	75	76	78	

1. Identify the high score ( $H = 98$ ) and the low score ( $L = 39$ ), and find the range:
  - $\text{range} = H - L = 98 - 39 = 59$
2. Select a number of classes ( $m = 7$ ) and a class width ( $c = 10$ ) so that the product ( $mc = 70$ ) is a bit larger than the range (range = 59).
3. Pick a starting point. This starting point should be a little smaller than the lowest score,  $L$ .



# FORM A FREQUENCY DISTRIBUTION AND HISTOGRAM

60	47	82	95	88	72	67	66	68	98	90	77	86
58	64	95	74	72	88	74	77	39	90	63	68	97
70	64	70	70	58	78	89	44	55	85	82	83	
72	77	72	86	50	94	92	80	91	75	76	78	

- Let the starting point to be 35. Given class width ( $c = 10$ )

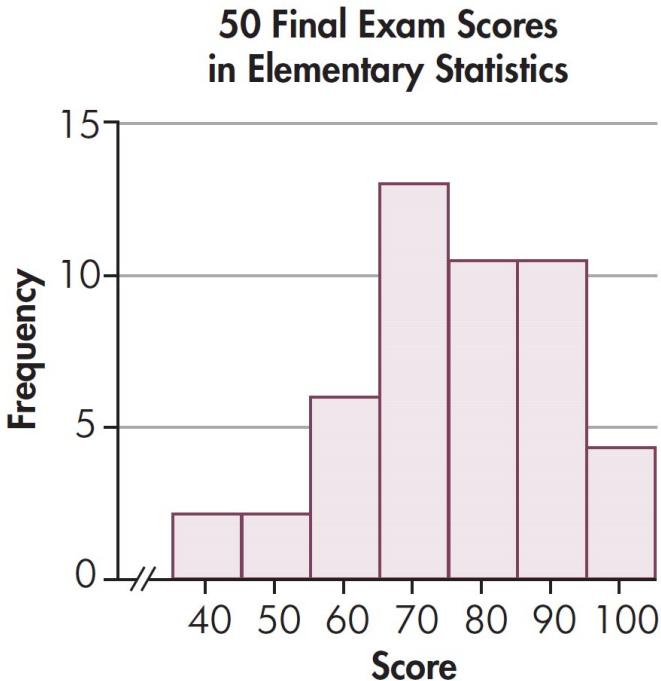
Class Number	Class Tallies	Boundaries	Frequency
1		$35 \leq x < 45$	2
2		$45 \leq x < 55$	2
3		$55 \leq x < 65$	7
4		$65 \leq x < 75$	13
5		$75 \leq x < 85$	11
6		$85 \leq x < 95$	11
7		$95 \leq x \leq 105$	4
			50

Standard Chart for Frequency Distribution

# FORM A FREQUENCY DISTRIBUTION AND HISTOGRAM

Class Number	Class Tallies	Boundaries	Frequency
1		$35 \leq x < 45$	2
2		$45 \leq x < 55$	2
3		$55 \leq x < 65$	7
4		$65 \leq x < 75$	13
5		$75 \leq x < 85$	11
6		$85 \leq x < 95$	11
7		$95 \leq x \leq 105$	4
			50

- Histogram** A bar graph that represents a frequency distribution of a quantitative variable.



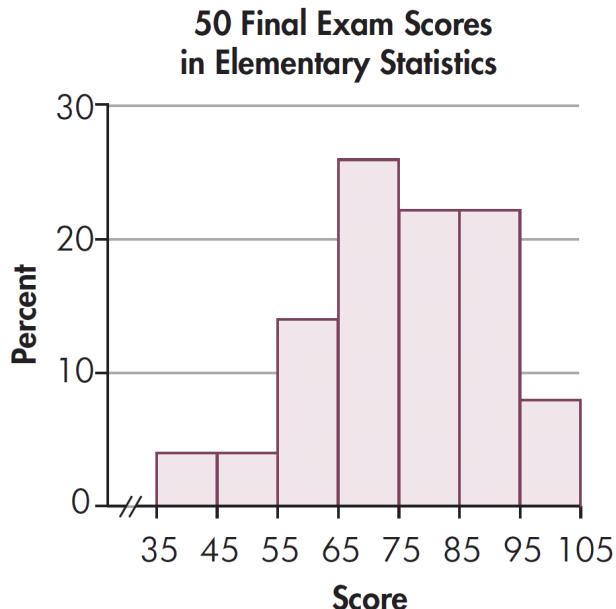


# FORM A FREQUENCY DISTRIBUTION AND HISTOGRAM

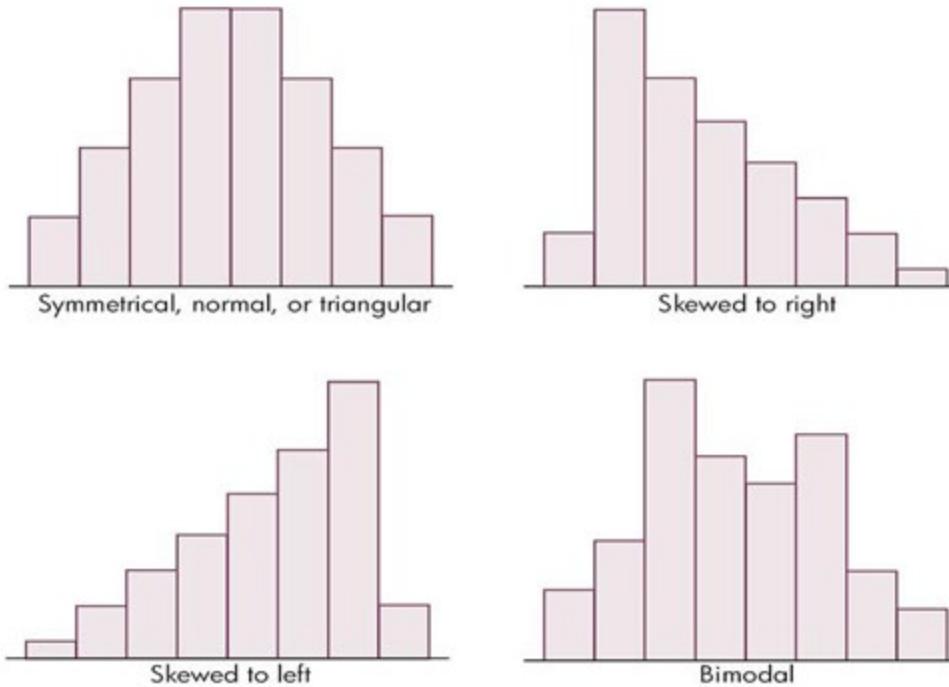
Class Number	Class Tallies	Boundaries	Frequency	Percentage
1		$35 \leq x < 45$	2	4%
2		$45 \leq x < 55$	2	4%
3		$55 \leq x < 65$	7	14%
4		$65 \leq x < 75$	13	26%
5		$75 \leq x < 85$	11	22%
6		$85 \leq x < 95$	11	22%
7		$95 \leq x \leq 105$	4	8%
				100%
				50

Divide all by 50

- **The relative frequency (percentage) is a proportional measure of the frequency for an occurrence.**
- **It is found by dividing the class frequency by the total number of observations.**



# FREQUENCY DISTRIBUTIONS AND HISTOGRAMS



- **Symmetrical** Both sides of this distribution are identical (halves are mirror images).
- **Skewed** One tail is stretched out longer than the other. The direction of skewness is on the side of the longer tail.



# FREQUENCY DISTRIBUTIONS AND HISTOGRAMS

**Bimodal** The two most populous classes are separated by one or more classes. This situation often implies that two populations are being sampled. (See Figure 2.7)

Weights of 50 College Students (lb)	
$N = 50$	Leaf Unit = 1.0
9	8
10	1 8 8
11	0 2 5 5 6 8 8
12	0 0 0 8 9
13	2 5 7
14	2 3 5 8
15	0 4 4 5 7 8
16	1 2 2 5 7 8
17	0 0 6 6 7
18	3 4 6 8
19	0 1 5 5
20	5
21	5

Stem-and-Leaf Display

Figure 2.7



$$\text{midrange} = \frac{\text{low value} + \text{high value}}{2}$$

$$\text{midrange} = \frac{L + H}{2} \quad (2.3)$$

Sample mean:  $\bar{x}$  =  $\frac{\text{sum of all } x}{\text{number of } x}$

$$\bar{x} = \frac{\Sigma x}{n} \quad (2.1)$$

# Numerical Summaries

$$s \text{ squared} = \frac{(\text{sum of } x^2) - \left[ \frac{(\text{sum of } x)^2}{\text{number}} \right]}{\text{number} - 1}$$

$$\text{sample variance: } s^2 = \frac{\Sigma x^2 - \frac{(\Sigma x)^2}{n}}{n - 1} \quad (2.9)$$

$$\text{sample variance: } s \text{ squared} = \frac{\text{sum of (deviations squared)}}{\text{number} - 1}$$

$$s^2 = \frac{\Sigma (x - \bar{x})^2}{n - 1} \quad (2.5)$$

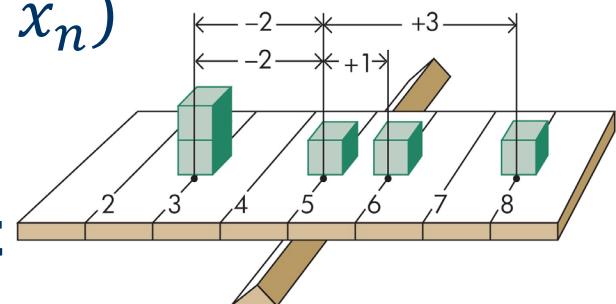
# MEASURES OF CENTRAL TENDENCY

- The measures of central tendency characterize the center of the distribution of data values. The term *average* is often associated with all measures of central tendency.
- Mean (arithmetic mean)** The average with which you are probably most familiar. The sample mean is represented by  $\bar{x}$  (read “*x-bar*” or “sample mean”).

Sample mean:  $\bar{x} = \frac{\text{sum of all } x}{\text{number of } x}$

- $$\bar{x} = \frac{1}{n} \sum_{i=1}^n x_i = \frac{1}{n} (x_1 + x_2 + \cdots + x_n)$$

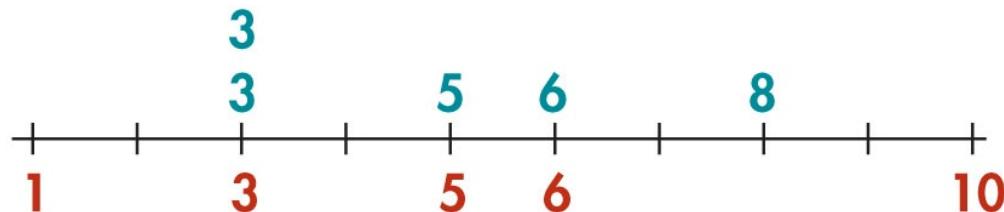
- The center of gravity or balance point**





# MEASURES OF CENTRAL TENDENCY

- **Sample Median:** Middle value when data ordered. 50% above, 50% below. Represented by  $\tilde{x}$  called “*x-tilde*.”
  - Order data from smallest to largest.
  - If n odd,  $\tilde{x}$  = middle value
  - If n even,  $\tilde{x}$  = average of middle two values
- **Sample Mode:** The value that happens most often in sample.
- **Represented by  $\hat{x}$  called “*x-hat*.”**
  - If two or more values in a sample are tied for the highest frequency, we say that there is no mode.



- There are other measures called **measures of dispersion** that characterize the spread or variability in the data.



# MEASURES OF DISPERSION

- **Range** The difference in value between the highest-valued data,  $H$ , and the lowest-valued data,  $L$ :
  - $\text{Range} = \text{high value} - \text{low value} = H - L$
- **Deviation from the mean:** The difference between the data value  $x_i$  and the sample mean  $\bar{x}$ 
  - $i^{\text{th}}$  deviation from the mean  $= x_i - \bar{x}$
- **The sum of the deviations,**  $\sum_{i=1}^n (x_i - \bar{x})$  **is always zero** because the deviations of  $x_i$  values smaller than the mean (which are negative) cancel out those  $x_i$  values larger than the mean (which are positive).



# MEASURES OF DISPERSION

- **Sample Variance:** The mean of the squared deviations using  $n - 1$  as a divisor.

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$$

– where  $x_i$  is the  $i^{th}$  data value,  $\bar{x}$  is the sample mean, and  $n$  is the sample size.

$SS(x)$ : sum of squares for  $x$

- This is equivalent to:

$$s^2 = \frac{1}{n-1} \left\{ \sum_{i=1}^n x_i^2 - \left[ \frac{(\sum_{i=1}^n x_i)^2}{n} \right] \right\}$$

- Therefore

$$s^2 = \frac{1}{n-1} SS(x)$$

# MEASURES OF DISPERSION

- **Sample Standard Deviation:** Square root of the sample variance.
  - Has **same units** data values and sample mean.

$$s = \sqrt{s^2}, \quad \text{where} \quad s^2 = \frac{1}{n-1} SS(x)$$

- Example: Consider a second set of data:  $\{6, 3, 8, 5, 2\}$ . Find the followings:
- Measures of Central Tendency
  - Mean  $\bar{x} = \frac{1}{5}(6 + 3 + 8 + 5 + 2) = 4.8$
  - Median  $\tilde{x} = \text{middle value} = 5$
  - Mode  $\hat{x} = \text{the value with the highest count} \Rightarrow \text{There is no mode}$
- Measures of Dispersion
  - Range  $\text{range} = H - L = 8 - 2 = 6$
  - Sample Variance and Sample Standard Deviation



## EXAMPLE (SAMPLE VARIANCE)

- Consider a second set of data: {6, 3, 8, 5, 2}

$$s^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$$

Step 1. Find $\Sigma x$	Step 2. Find $\bar{x}$	Step 3. Find each $x - \bar{x}$	Step 4. Find $\Sigma(x - \bar{x})^2$	Step 5. Find $s^2$
6	$\bar{x} = \frac{\Sigma x}{n}$	$6 - 4.8 = 1.2$	$(1.2)^2 = 1.44$	$s^2 = \frac{\Sigma(x - \bar{x})^2}{n-1}$
3		$3 - 4.8 = -1.8$	$(-1.8)^2 = 3.24$	
8		$8 - 4.8 = 3.2$	$(3.2)^2 = 10.24$	
5	$\bar{x} = \frac{24}{5}$	$5 - 4.8 = 0.2$	$(0.2)^2 = 0.04$	$s^2 = \frac{22.80}{4}$
2		$2 - 4.8 = -2.8$	$(-2.8)^2 = 7.84$	
$\Sigma x = 24$	$\bar{x} = 4.8$	$\Sigma(x - \bar{x}) = 0$	$\Sigma(x - \bar{x})^2 = 22.80$	$s^2 = 5.7$

- Or  $s^2 = \frac{1}{n-1} \left\{ \sum_{i=1}^n x_i^2 - \frac{1}{n} (\sum_{i=1}^n x_i)^2 \right\}$

$$s = \sqrt{s^2} = \sqrt{5.7}$$

Step 1. Find $\Sigma x$	Step 2. Find $\Sigma x^2$	Step 3. Find $SS(x)$	Step 4. Find $s^2$	Step 5. Find $s$
6	$6^2 = 36$	$SS(x) = \Sigma x^2 - \frac{(\Sigma x)^2}{n}$		$s = \sqrt{s^2}$
3	$3^2 = 9$		$s^2 = \frac{\Sigma x^2 - \frac{(\Sigma x)^2}{n}}{n-1}$	$s = \sqrt{5.7}$
8	$8^2 = 64$	$SS(x) = 138 - \frac{(24)^2}{5}$		$s = 2.4$
5	$5^2 = 25$		$s^2 = \frac{22.8}{4}$	
$\Sigma x = 24$	$\Sigma x^2 = 138$	$SS(x) = 138 - 115.2$	$s^2 = 5.7$	
		$SS(x) = 22.8$		



# EXAMPLE IN MICROSOFT EXCEL

The screenshot shows a Microsoft Excel spreadsheet titled "Book2 - Excel". The data in column A is: 6, 3, 8, 5, 2. The formulas calculated in column I are:

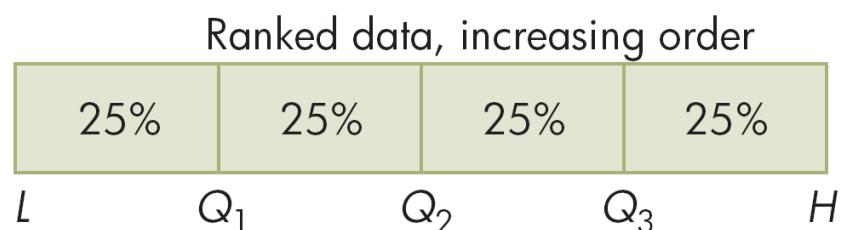
- =AVERAGE(A1:A5) → Value: 4.8
- =MEDIAN(A1:A5) → Value: 5
- =MODE(A1:A5) → Value: #N/A
- =VAR(A1:A5) → Value: 2.387467
- =STDEV(A1:A5) → Value: 5.7

Annotations in red text with arrows point to specific elements:

- "Enter the formula" points to the formula bar with the formula =AVERAGE(A1:A5).
- "Enter the data values" points to the data in column A (A1:A5).
- "=AVERAGE(A1:A5)" points to the value 4.8 in cell I2.
- "=MEDIAN(A1:A5)" points to the value 5 in cell I3.
- "=MODE(A1:A5)" points to the error value #N/A in cell I4.
- "=VAR(A1:A5)" points to the value 2.387467 in cell I5.
- "=STDEV(A1:A5)" points to the value 5.7 in cell I6.

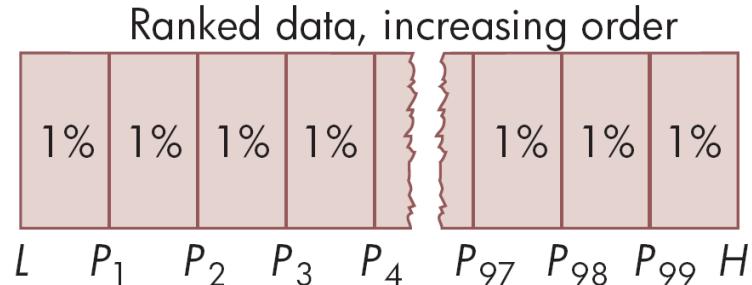
# MEASURES OF POSITION

- **Measures of position** are used to describe the position a specific data value possesses in relation to the rest of the data when in ranked order. *Quartiles* and *percentiles* are two of the most popular measures of position.
- **Quartiles** Values of the variable that divide the ranked data into quarters; each set of data has three quartiles.
  - L = lowest value
  - $Q_1$  = data value where 25% are smaller
  - $Q_2 = \tilde{x}$  =median
  - $Q_3$  = data value where 75% are smaller
  - H = highest value

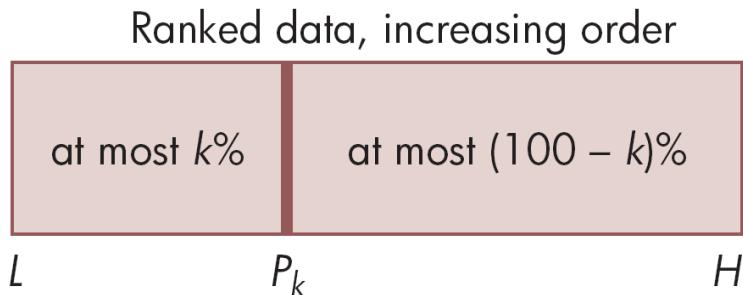


# MEASURES OF POSITION

- **Percentiles** Values of the variable that divide a set of ranked data into 100 equal subsets.



- The  $k^{th}$  percentile,
  - $P_k$  = value where  $k\%$  are smaller



- Quartiles are special percentiles.

# MEASURES OF POSITION

- **The Percentile Process: Finding  $k^{th}$  percentile.**

Step 1 Rank the  $n$  data, lowest to highest



Step 2 Calculate  $\frac{nk}{100}$

An integer **A** results

A number with a fraction results

Step 3  $d(P_k) = \mathbf{A.5}$

$d(P_k) = \mathbf{B}$ , the next larger integer

Step 4

$d(\cdot)$  is called the depth (indicates the location of  $k^{th}$  the percentile)

$P_k$  is halfway between the value of the data in the **A**th position and the value of the data in the **A+1** position.

average of  $A^{\text{th}}$  and  $(A+1)^{\text{th}}$  values

$P_k$  is the value of the data in the **B**th position.

**B**<sup>th</sup> value



## EXAMPLE 12 - FINDING QUARTILES AND PERCENTILES

- Using the sample of 50 elementary statistics final exam scores listed in Table 2.15, find the first quartile,  $Q_1$ ; the 58<sup>th</sup> percentile,  $P_{58}$ .

60	47	82	95	88	72	67	66	68	98	90	77	86
58	64	95	74	72	88	74	77	39	90	63	68	97
70	64	70	70	58	78	89	44	55	85	82	83	
72	77	72	86	50	94	92	80	91	75	76	78	

Raw Scores for Elementary Statistics Exam [TA02-06]

Table 2.15

- Step 1:

- Rank the data from lowest to highest

39	64	72	78	89
44	66	72	80	90
47	67	74	82	90
50	68	74	82	91
55	68	75	83	92
58	70	76	85	94
58	70	77	86	95
60	70	77	86	95
63	72	77	88	97
64	72	78	88	98



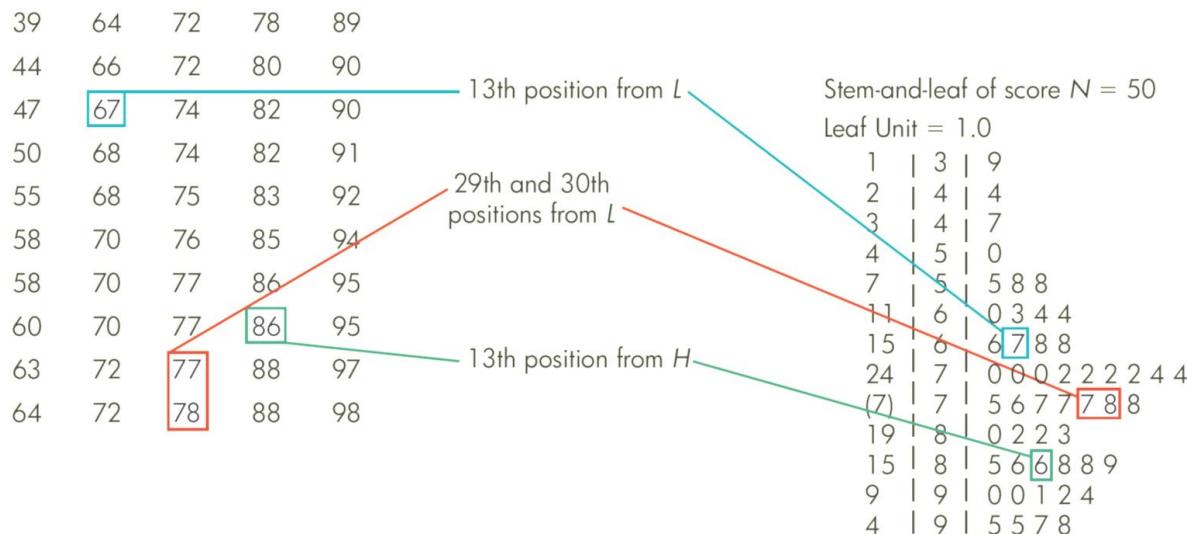
## EXAMPLE 12 - *SOLUTION*

- **Find  $Q_1$ :**

- **Step 2: Find**  $\frac{nk}{100} : \frac{nk}{100} = \frac{(50)(25)}{100} = 12.5$
- **Step 3:  $B$  is the next larger integer, 13.**
- **Step 4: Find  $Q_1$ :  $Q_1$  is the 13<sup>th</sup> value,  $Q_1 = 67$**

- **Find  $P_{58}$ :**

- **Step 2: Find**  $\frac{nk}{100} : \frac{nk}{100} = \frac{(50)(58)}{100} = 29$
- **Step 3: Since  $A = 29$ , an integer, add 0.5 and use 29.5.**
- **Step 4: Find  $P_{58}$ :  $P_{58}$  is the average of 29<sup>th</sup> and 30<sup>th</sup> values,  $P_{58} = 77.5$**



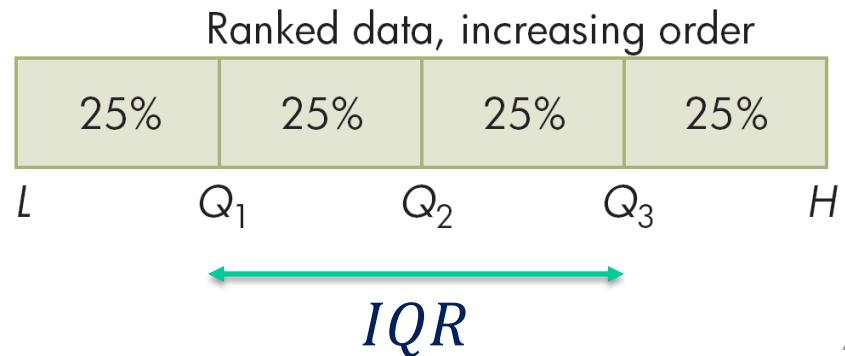
# MEASURES OF POSITION

- **Five Number Summary**

- L = lowest value
- $Q_1$  = data value where 25% are smaller
- $Q_2 = \tilde{x}$  = median
- $Q_3$  = data value where 75% are smaller
- H = highest value

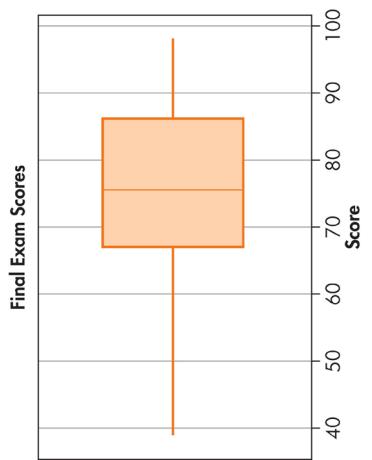
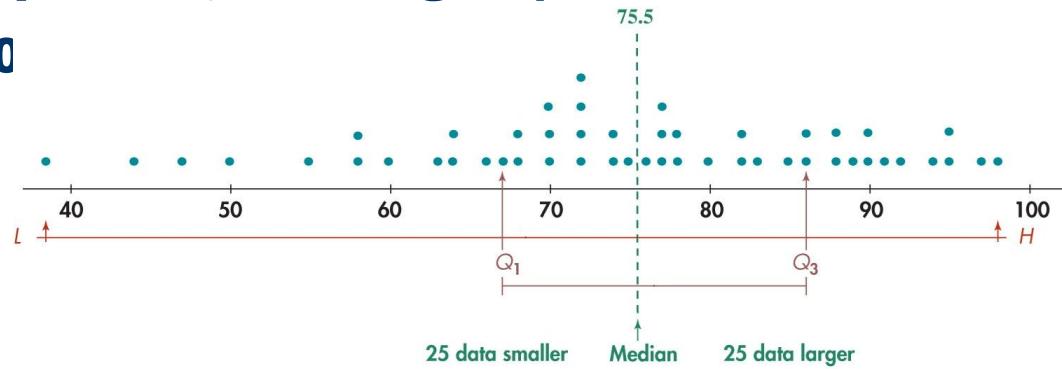
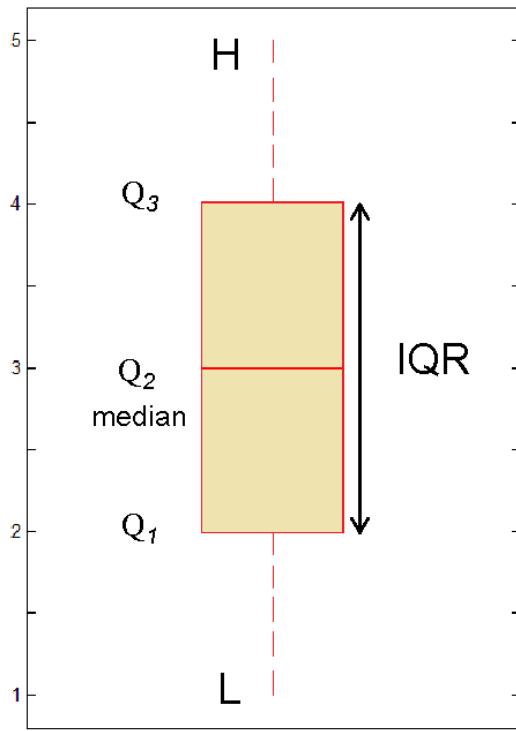
- **Interquartile range: The difference between the first and third quartiles. It is the range of the middle 50% of the data.**

- $IQR = Q_3 - Q_1$

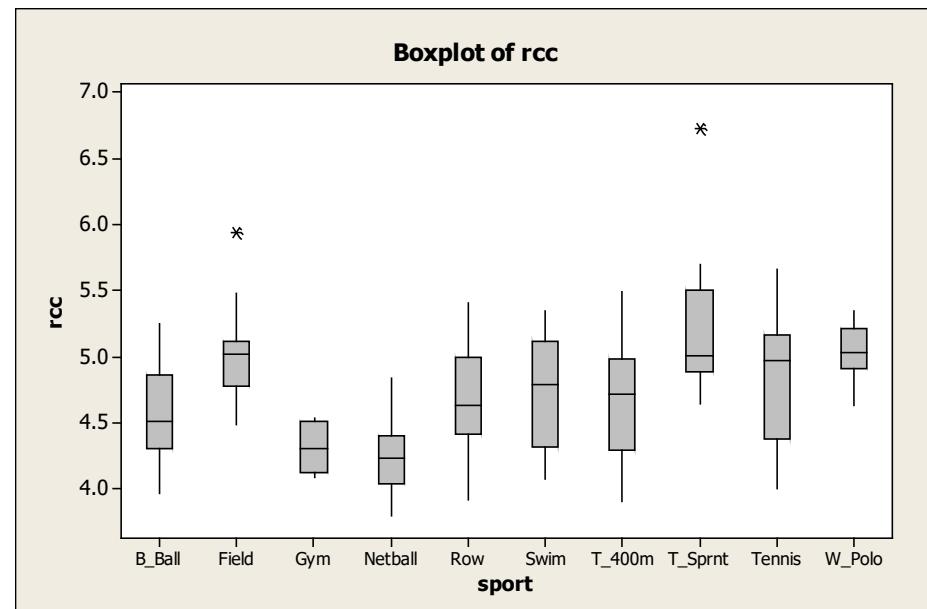
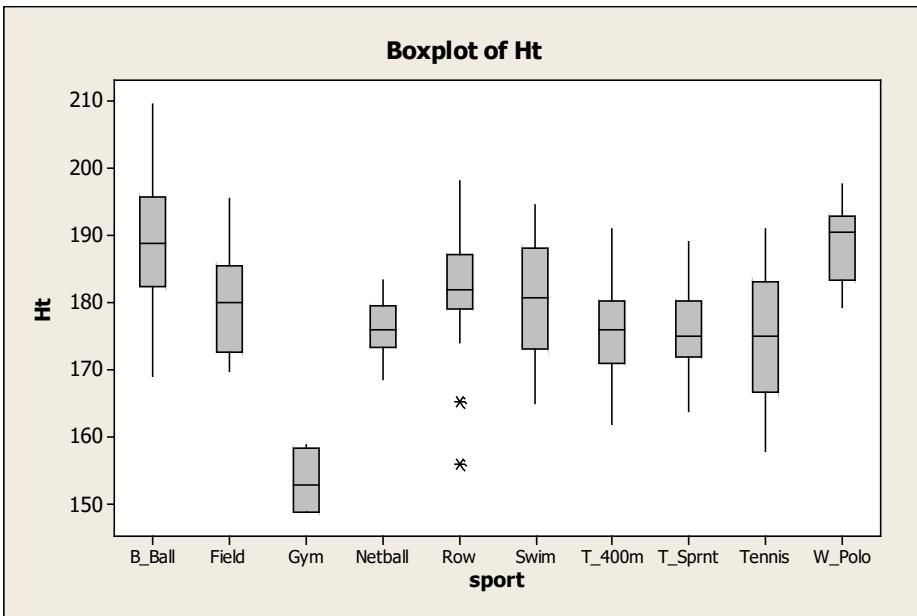


# MEASURES OF POSITION

**Box-and-whiskers display** A graphic representation of the 5-number summary. The five numerical values (smallest, first quartile, median, third quartile, and largest) are located on a scale, either vertical or horizontal.



# SIDE-BY-SIDE BOX PLOT FOR AIS DATA (AUSTRALIAN INSTITUTE OF SPORT)



# MEASURES OF POSITION

- The position of a specific value can also be measured in terms of the mean and standard deviation using the *standard score*, commonly called the *z-score*.
- Standard score, or z-score** The position a particular value of  $x$  has relative to the mean, measured in standard deviations. The z-score is found by the formula

$$z = \frac{\text{value} - \text{mean}}{\text{st.dev.}} = \frac{x - \bar{x}}{s} \quad (2.11)$$

- Example:** Find the standard scores for (a) 92 and (b) 72 with respect to a sample of exam grades that have a mean score of 74.92 and a standard deviation of 14.20.

## EXAMPLE 14 - FINDING Z-SCORES

- **a.**  $x_1 = 92$ ,  $\bar{x} = 74.92$ ,  $s = 14.20$ . Thus,

$$\begin{aligned} z &= \frac{x - \bar{x}}{s} \\ &= \frac{92 - 74.92}{14.20} = \frac{17.08}{14.20} = 1.20. \end{aligned}$$

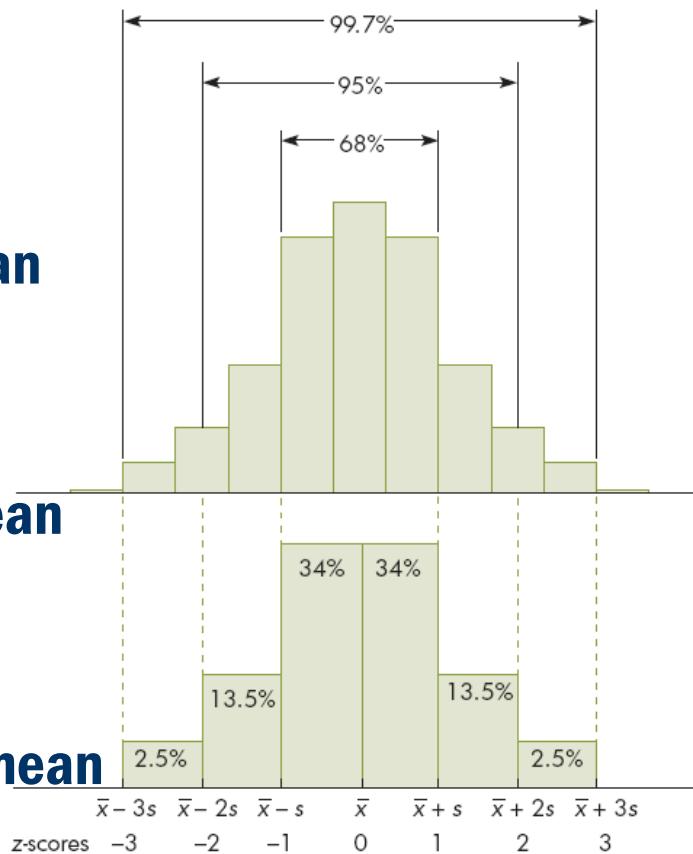
- **b.**  $x_2 = 72$ ,  $\bar{x} = 74.92$ ,  $s = 14.20$ . Thus,

$$\begin{aligned} z &= \frac{x - \bar{x}}{s} \\ &= \frac{72 - 74.92}{14.20} = \frac{-2.92}{14.20} = -0.21. \end{aligned}$$

- This means that the score 92 is approximately 1.2 standard deviations above the mean and
- that the score 72 is approximately one-fifth of a standard deviation below the mean.

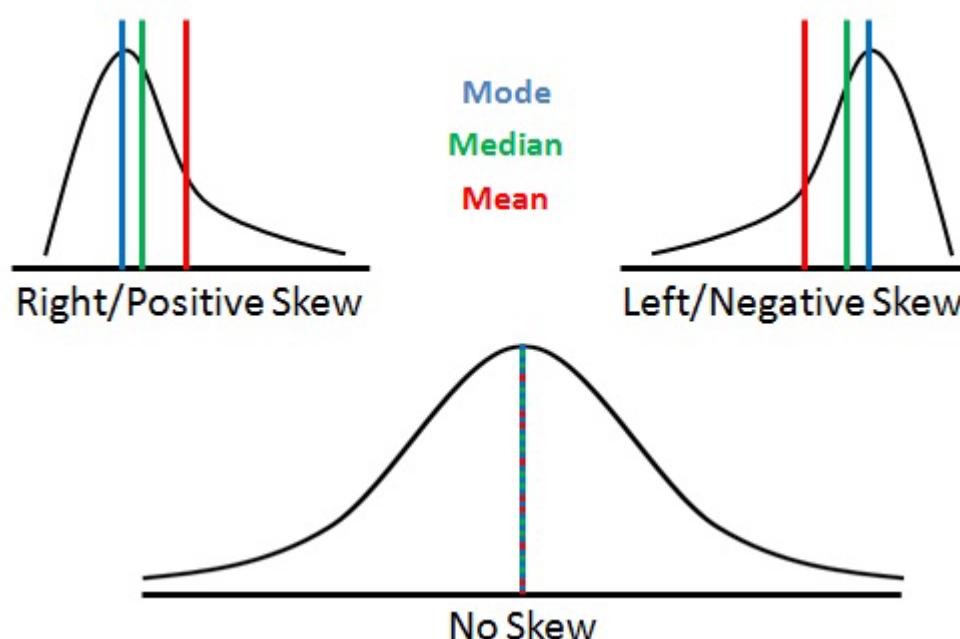
# EMPIRICAL RULE (THE 68-95-99.7 RULE)

- If the distribution is mound-shaped, then
  - Approximately 68% of the data falls within one standard deviation of the mean
  - Approximately 95% of the data falls within two standard deviations of the mean
  - Approximately 99.7% of the data falls within three standard deviations of the mean



# COMPARING MEASURES OF CENTER AND SPREAD

- The sample mean and the sample standard deviation are good measures of center and spread, respectively, for symmetric data
- If the data set is skewed or has outliers, the sample median and the interquartile range are more commonly used
- Mean versus median



# BIVARIATE DATA

- **Bivariate data** The values of two different variables that are obtained from the same population element.
- Each of the two variables may be either *qualitative* or *quantitative*. As a result, three combinations of variable types can form bivariate data:
  1. Both variables are qualitative (categorical).
  2. One variable is qualitative (categorical), and the other is quantitative (numerical).
  3. Both variables are quantitative (numerical).

# TWO QUALITATIVE VARIABLES

- **Qualitative vs Qualitative:**
  - A cross-tabulation or contingency table will be used
- **Example: Thirty students from our college were randomly identified and classified according to two variables: gender (M/F) and major ( liberal arts (LA), business administration (BA), technology(T) )**

Name	Gender	Major	Name	Gender	Major	Name	Gender	Major
Adams	M	LA	Feehey	M	T	McGowan	M	BA
Argento	F	BA	Flanigan	M	LA	Mowers	F	BA
Baker	M	LA	Hodge	F	LA	Ornt	M	T
Bennett	F	LA	Holmes	M	T	Palmer	F	LA
Brand	M	T	Jopson	F	T	Pullen	M	T
Brock	M	BA	Kee	M	BA	Rattan	M	BA
Chun	F	LA	Kleeberg	M	LA	Sherman	F	LA
Crain	M	T	Light	M	BA	Small	F	T
Cross	F	BA	Linton	F	LA	Tate	M	BA
Ellis	F	BA	Lopez	M	T	Yamamoto	M	LA

# EXAMPLE 1 - CONSTRUCTING CROSS-TABULATION TABLES

- We can construct a  $2 \times 3$  table.**

- Given:**

M = male

F = female

LA = liberal arts

BA = business admin

T = technology

Name	Gender	Major	Name	Gender	Major	Name	Gender	Major
Adams	M	LA	Feeaney	M	T	McGowan	M	BA
Argento	F	BA	Flanigan	M	LA	Mowers	F	BA
Baker	M	LA	Hodge	F	LA	Ornt	M	T
Bennett	F	LA	Holmes	M	T	Palmer	F	LA
Brand	M	T	Jopson	F	T	Pullen	M	T
Brock	M	BA	Kee	M	BA	Rattan	M	BA
Chun	F	LA	Kleeberg	M	LA	Sherman	F	LA
Crain	M	T	Light	M	BA	Small	F	T
Cross	F	BA	Linton	F	LA	Tate	M	BA
Ellis	F	BA	Lopez	M	T	Yamamoto	M	LA

- The entry in each cell is found by determining how many students fit into each category. Adams is male (M) and liberal arts (LA) and is classified in the cell in the first row, first column.**

Gender	Major					
	LA	BA	T			
M		(5)		(6)		(7)
F		(6)		(4)		(2)

# EXAMPLE 1 CONT'D

- The resulting  $2 \times 3$  contingency (cross-tabulation) table is:

Gender	Major			Row Total
	LA	BA	T	
M	5	6	7	18
F	6	4	2	12
Col. Total	11	10	9	30

- Percentages Based on the Grand Total (Entire Sample)

Gender	Major			Row Total
	LA	BA	T	
M	17%	20%	23%	60%
F	20%	13%	7%	40%
Col. Total	37%	33%	30%	100%

- Percentages Based on Row Totals
- (Marginal: within Gender)

Gender	Major			Row Total
	LA	BA	T	
M	28%	33%	39%	100%
F	50%	33%	17%	100%
Col. Total	37%	33%	30%	100%

Gender	Major			T
	LA	BA	T	
M		(5)	(6)	(7)
F	(6)	(4)	(2)	

- Percentages Based on Column Totals
- (Marginal: within Major)

Gender	Major			Row Total
	LA	BA	T	
M	45%	60%	78%	60%
F	55%	40%	22%	40%
Col. Total	100%	100%	100%	100%

# QUANTITATIVE VS QUALITATIVE: SIDE-BY-SIDE COMPARISONS

- **Quantitative vs Qualitative:**
- **Example 2: The distance required to stop a 3000-pound automobile on wet pavement was measured to compare the stopping capabilities of three tire tread designs.**

Design A ( $n = 6$ )			Design B ( $n = 6$ )			Design C ( $n = 6$ )		
37	36	38	33	35	38	40	39	40
34	40	32	34	42	34	41	41	43

– **5-Number Summary  
for Each Design**

	Design A	Design B	Design C
High	40	42	43
$Q_3$	38	38	41
Median	36.5	34.5	40.5
$Q_1$	34	34	40
Low	32	33	39

**Mean and Standard Deviation  
for Each Design**

	Design A	Design B	Design C
Mean	36.2	36.0	40.7
Standard deviation	2.9	3.4	1.4

# EXAMPLE 2 - CONSTRUCTING SIDE-BY-SIDE COMPARISONS

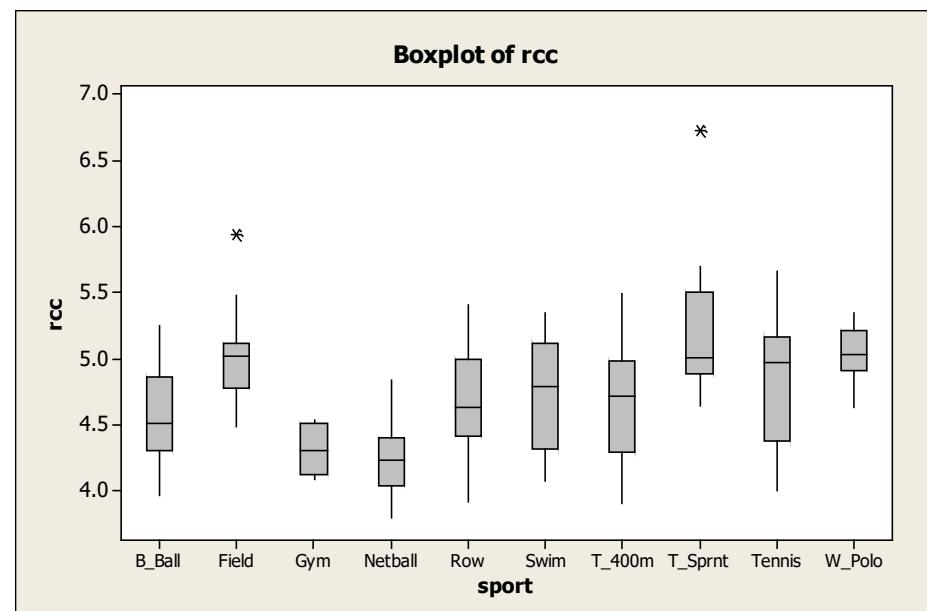
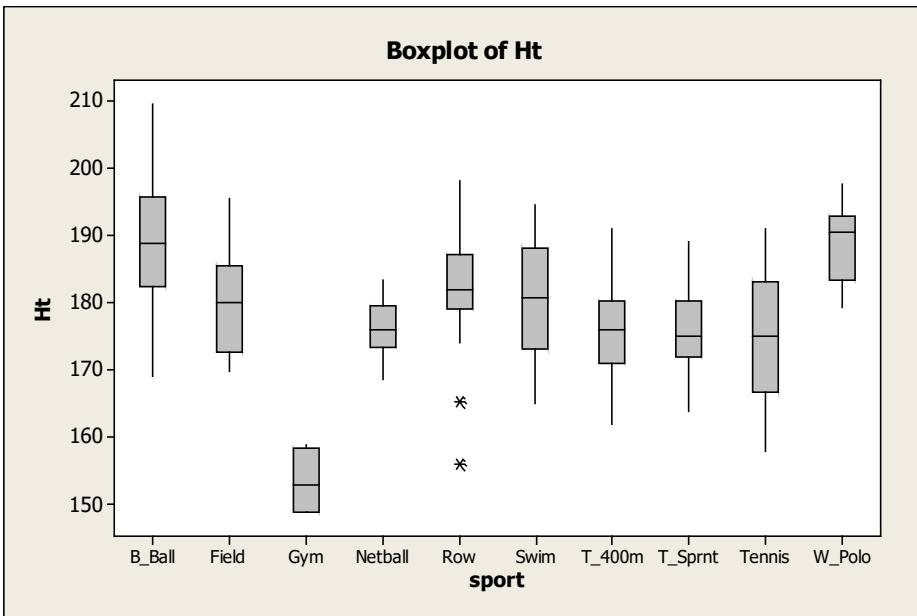
- The distance required to stop a 3000-pound automobile on wet pavement was measured to compare the stopping capabilities of three tire tread designs

Design A ( $n = 6$ )	Design B ( $n = 6$ )	Design C ( $n = 6$ )
37    36    38	33    35    38	40    39    40
34    40    32	34    42    34	41    41    43

- Side by side Box-and-Whiskers



# SIDE-BY-SIDE BOX PLOT FOR AIS DATA (AUSTRALIAN INSTITUTE OF SPORT)



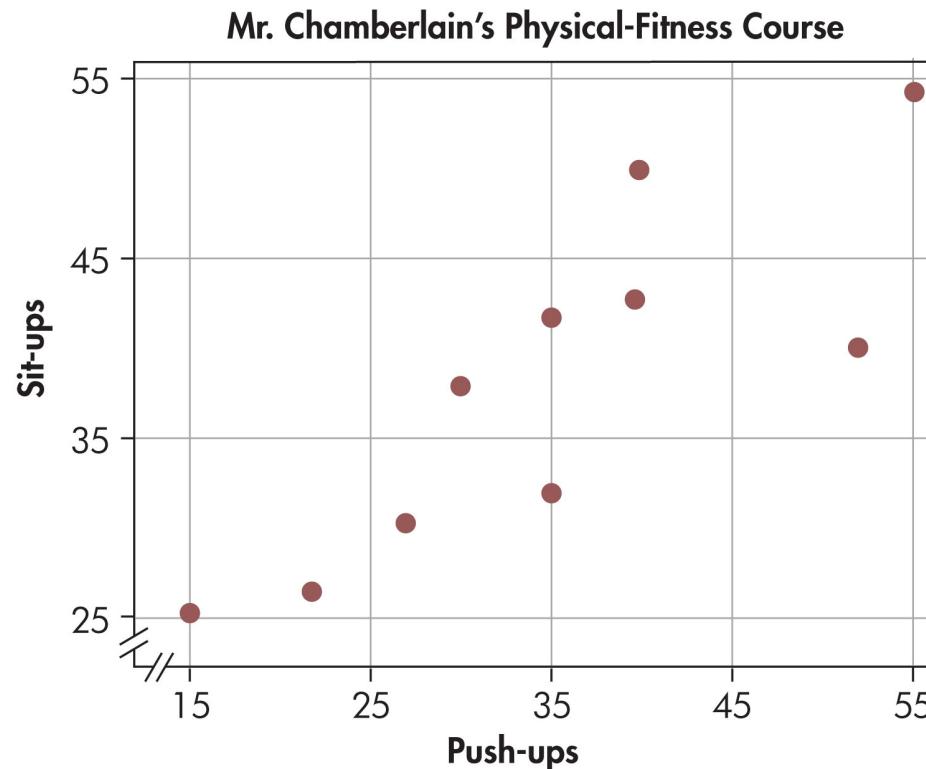
## TWO QUANTITATIVE VARIABLES

- It is customary to express the data mathematically as ordered pairs  $(x, y)$ , where
  - $x$  is the input variable (called the independent variable)
  - $y$  is the output variable (called the dependent variable).
- The data are said to be *ordered* because one value,  $x$ , is always written first.
- They are called *paired* because for each  $x$  value, there is a corresponding  $y$  value from the same source.
- Scatter diagram A plot of all the ordered pairs of bivariate data on a coordinate axis system. The input variable,  $x$ , is plotted on the horizontal axis, and the output variable,  $y$ , is plotted on the vertical axis.

# TWO QUANTITATIVE VARIABLES

- **Example: Push-ups vs Sit-ups**

Student	1	2	3	4	5	6	7	8	9	10
Push-ups, $x$	27	22	15	35	30	52	35	55	40	40
Sit-ups, $y$	30	26	25	42	38	40	32	54	50	43



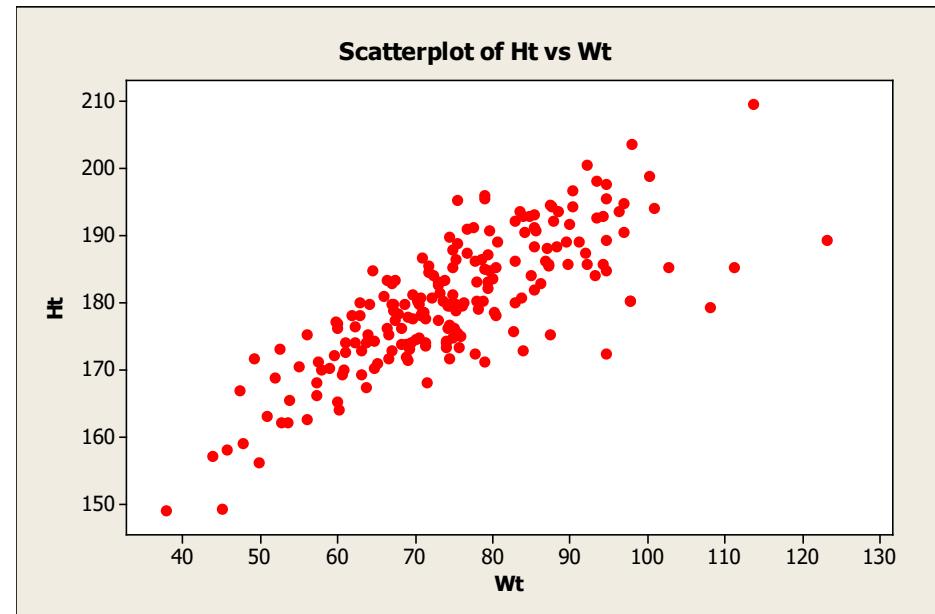
# TWO QUANTITATIVE VARIABLES

- There is not always an explanatory-response (dependent-independent) relationship.

- More examples:

- Height and Weight
- Income and Age
- SAT scores on math exam and on verbal exam
- Amount of time spent studying for an exam and exam score

Australian Institute of Sport (AIS.xlsx)

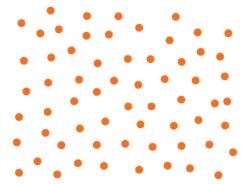


# TWO QUANTITATIVE VARIABLES

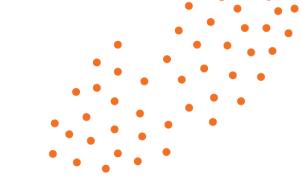
- Why Scatterplots?
- Look for overall pattern and any striking deviations from that pattern.
- Look for outliers, values falling outside the overall pattern of the relationship
- You can describe the overall pattern of a scatterplot by the form, direction, and strength of the relationship.
  - Form: Linear or clusters
  - Direction
    - Two variables are positively associated when above-average values of one tend to accompany above-average values of the other and likewise below-average values also tend to occur together.
    - Two variables are negatively associated when above-average values of one variable accompany below-average values of the other variable, and vice-versa.
  - Strength-how close the points lie to a line

# LINEAR CORRELATION

- **Linear Correlation,  $r$ , is a measure of the strength of a linear relationship between two variables  $x$  and  $y$ .**
- $-1 \leq r \leq 1$       **Will discuss the definition in a minute**



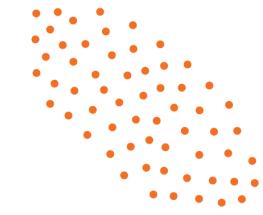
No correlation



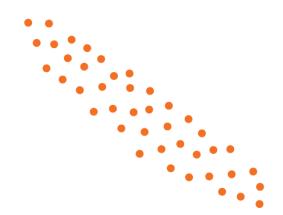
Positive



High positive



Negative



High negative

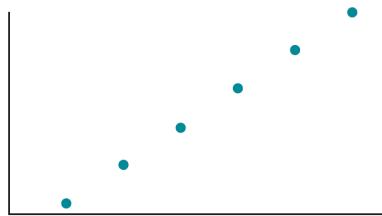
- $r \approx 0$

- $r \approx 0.5$

- $r \approx 0.8$

- $r \approx -0.5$

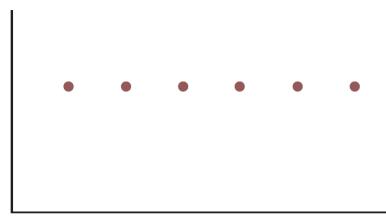
- $r \approx -0.8$



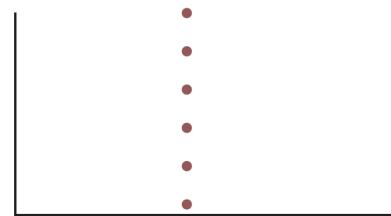
Perfect Positive Correlation



Perfect Negative Correlation



Horizontal—No Correlation

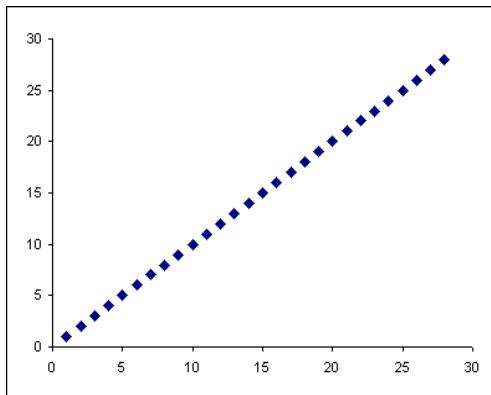


Vertical—No Correlation

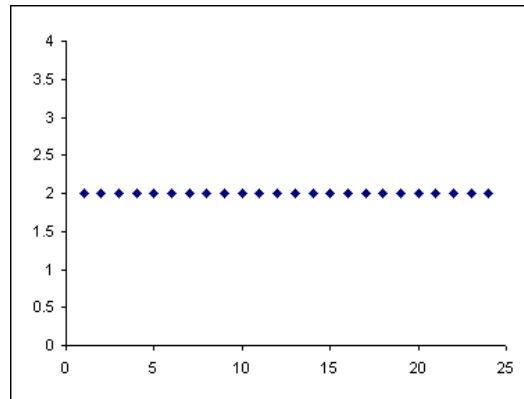
# LINEAR CORRELATION (FORMULA 1)

- $r = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{(n-1)s_x s_y}$ , where
  - $s_x$  and  $s_y$ : are the standard deviations of  $x$ 's and  $y$ 's
  - $s_x = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2}$
  - $s_y = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (y_i - \bar{y})^2}$
- Examples of extreme cases

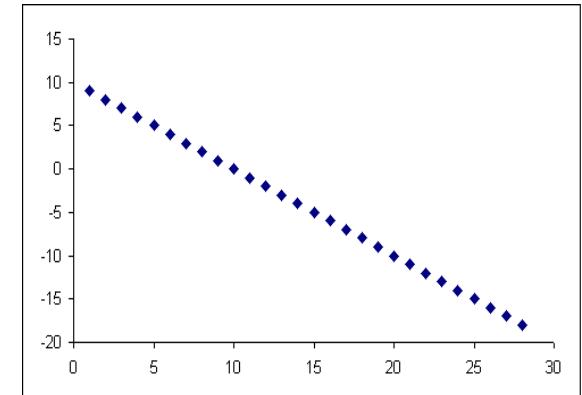
$$r = 1$$



$$r = 0$$



$$r = -1$$



# LINEAR CORRELATION (FORMULA 2)

- $r = \frac{SS(xy)}{\sqrt{SS(x)SS(y)}}$  (\*), where
  - $SS(x) = \sum_{i=1}^n x_i^2 - \frac{1}{n} (\sum_{i=1}^n x_i)^2$  ➤  $SS(x) = \sum_{i=1}^n (x_i - \bar{x})^2$
  - $SS(y) = \sum_{i=1}^n y_i^2 - \frac{1}{n} (\sum_{i=1}^n y_i)^2$  ➤  $SS(y) = \sum_{i=1}^n (y_i - \bar{y})^2$
  - $SS(xy) = \sum_{i=1}^n x_i y_i - \frac{1}{n} (\sum_i x_i)(\sum_i y_i)$  ➤  $SS(xy) = \sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})$
- (\*) is equivalent to the first formula:  $r = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{(n-1)s_x s_y}$
- Example 5: In Mr. Chamberlain's physical-fitness course, several fitness scores were taken. The following sample is the numbers of push-ups and sit-ups done by 10 randomly selected students:

Student	1	2	3	4	5	6	7	8	9	10
Push-ups, x	27	22	15	35	30	52	35	55	40	40
Sit-ups, y	30	26	25	42	38	40	32	54	50	43

## EXAMPLE 5 - *SOLUTION*

- Find the linear correlation coefficient for the push-up/sit-up data.

Student	Push-ups, $x$	$x^2$	Sit-ups, $y$	$y^2$	$xy$
1	27	729	30	900	810
2	22	484	26	676	572
3	15	225	25	625	375
4	35	1,225	42	1,764	1,470
5	30	900	38	1,444	1,140
6	52	2,704	40	1,600	2,080
7	35	1,225	32	1,024	1,120
8	55	3,025	54	2,916	2,970
9	40	1,600	50	2,500	2,000
10	40	1,600	43	1,849	1,720
$\sum x = 351$		$\sum x^2 = 13,717$	$\sum y = 380$	$\sum y^2 = 15,298$	$\sum xy = 14,257$
sum of $x$		sum of $x^2$	sum of $y$	sum of $y^2$	sum of $xy$

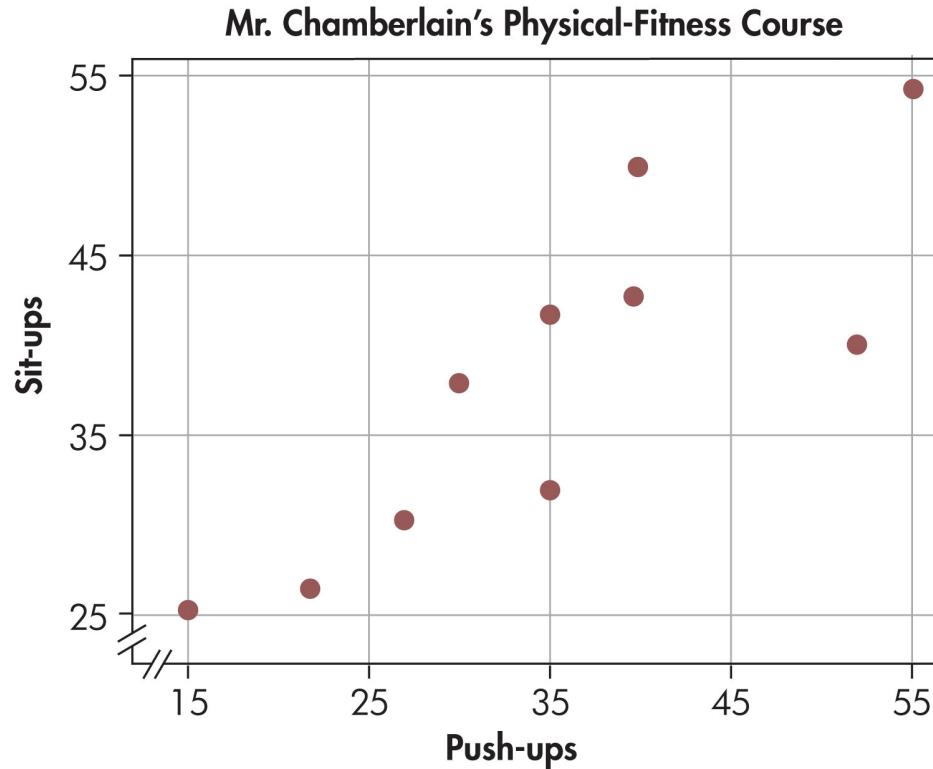
$$SS(x) = \sum x^2 - \frac{(\sum x)^2}{n} = 13,717 - \frac{(351)^2}{10} = 1396.9$$

$$SS(y) = \sum y^2 - \frac{(\sum y)^2}{n} = 15,298 - \frac{(380)^2}{10} = 858.0$$

$$SS(xy) = \sum xy - \frac{\sum x \sum y}{n} = 14,257 - \frac{(351)(380)}{10} = 919.00$$

## EXAMPLE 5 - *SOLUTION*

$$\begin{aligned}
 r &= \frac{SS(xy)}{\sqrt{SS(x)SS(y)}} \\
 &= \frac{919.0}{\sqrt{(1396.9)(858.0)}} \\
 &= 0.8394 = \mathbf{0.84}
 \end{aligned}$$



# RELATIONSHIPS BETWEEN 2 NUMERIC VARIABLES

- **Correlation or  $r$**  : measures the direction and strength of the linear relationship between two numeric variables
- General Properties
  - It must be between -1 and 1, or  $(-1 \leq r \leq 1)$ .
  - If  $r$  is negative, the relationship is negative.
  - If  $r = -1$ , there is a perfect negative linear relationship (extreme case).
  - If  $r$  is positive, the relationship is positive.
  - If  $r = 1$ , there is a perfect positive linear relationship (extreme case).
  - If  $r$  is 0, there is no linear relationship.
  - $r$  measures the strength of the linear relationship.
- Correlation Applet

# CAUSATION AND LURKING VARIABLES

- The **cause-and-effect relationship**: Correlation does not necessarily imply causation. Just because two things are highly related does not mean that one causes the other.
- A perceived relationship between a **dependent(response)** variable and an **independent(explanatory)** variable that has been misestimated due to the failure to account for a confounding factor (lurking variable) is termed a **spurious relationship**
- **Examples of spurious relationship**



# LURKING VARIABLE AND SIMPSON'S PARADOX

- **Lurking variable:** A variable that is not included in a study but has an effect on the variables of the study and makes it appear that those variables are related.
- **Simpson's Paradox:** An association or comparison that holds for all of several groups can reverse direction when a **lurking variable** is present.
- **Example: Kidney stone treatment** (Br Med J (Clin Res Ed) 292 (6524): 879-882)

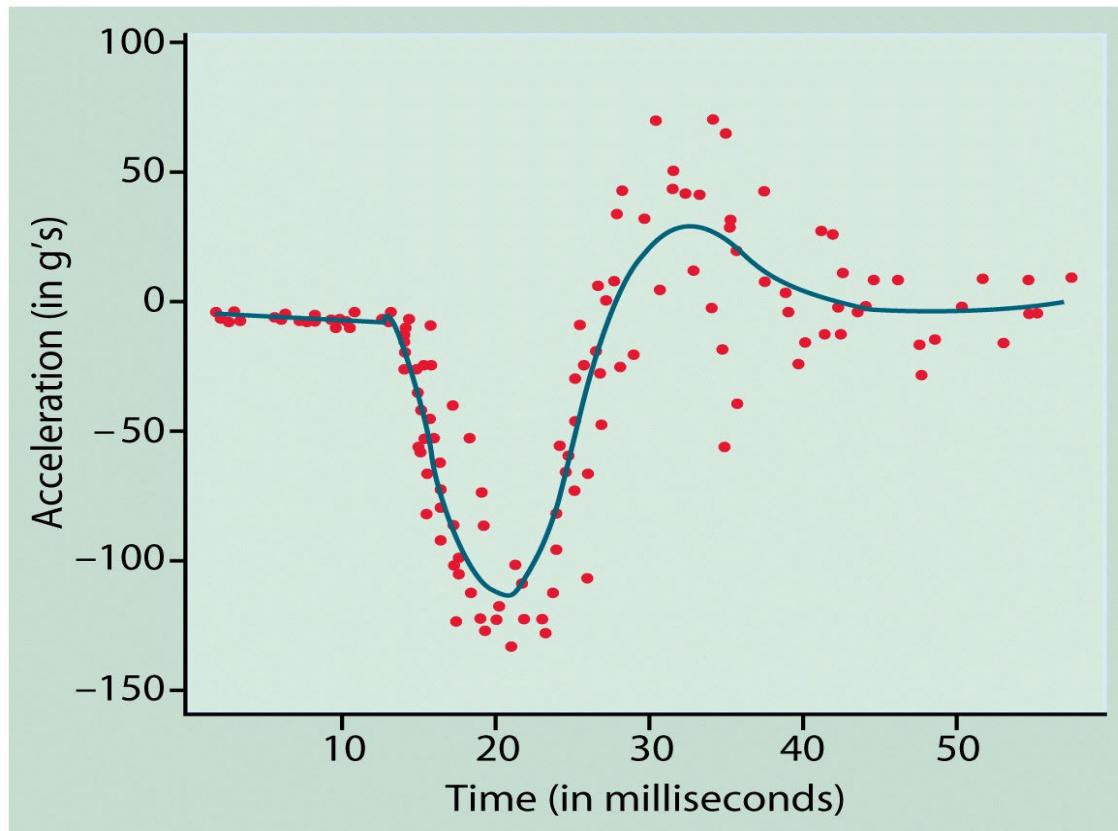
	Treatment A	Treatment B
Small Stones	<i>Group 1</i> <b>93% (81/87)</b>	<i>Group 2</i> <b>87% (234/270)</b>
Large Stones	<i>Group 3</i> <b>73% (192/263)</b>	<i>Group 4</i> <b>69% (55/80)</b>
Both	<b>78% (273/350)</b>	<b>83% (289/350)</b>

- [http://en.wikipedia.org/wiki/Simpson's\\_Paradox](http://en.wikipedia.org/wiki/Simpson's_Paradox)

## RELATIONSHIPS BETWEEN 2 NUMERIC VARIABLES

**It is possible for there to be a strong relationship between two variables and still have  $r \approx 0$ .**

**EX.**



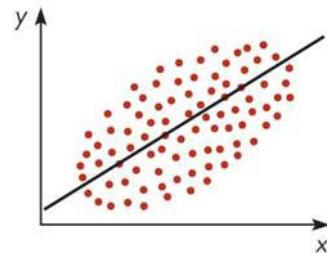
# LINEAR REGRESSION

- Regression analysis finds the equation of the line that best describes the relationship between two variables.
- Here are some examples of various possible relationships, called *models* or prediction equations:
  - Linear (straight-line):  $\hat{y} = b_0 + b_1x$
  - Quadratic:  $\hat{y} = a + bx + cx^2$
  - Exponential:  $\hat{y} = a(b^x)$
  - Logarithmic:  $\hat{y} = a \log_b x$

What we cover  
in this book

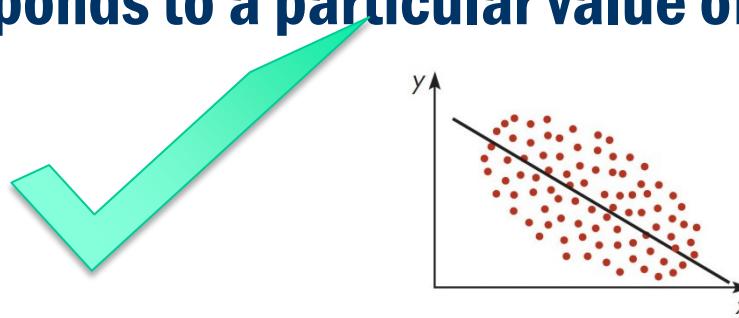
# LINEAR REGRESSION

- Suppose that  $\hat{y} = b_0 + b_1x$  is the equation of a straight line, where  $\hat{y}$  (read “y-hat”) represents the predicted value of  $y$  that corresponds to a particular value of  $x$ .



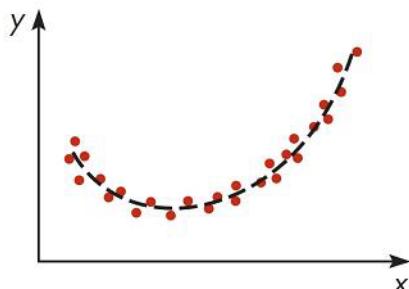
Linear Regression with Positive Slope

Figure 3.17



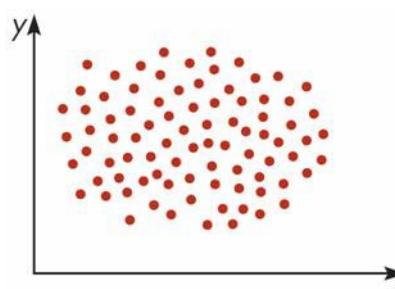
Linear Regression with Negative Slope

Figure 3.18



Curvilinear Regression (Quadratic)

Figure 3.19

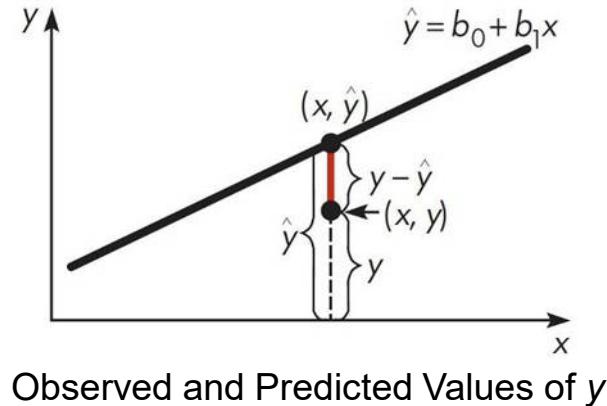


No Relationship

Figure 3.20

# LINEAR REGRESSION

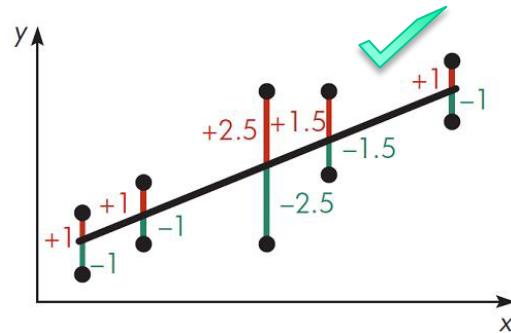
- The least squares criterion requires that we find the constants  $b_0$  and  $b_1$  such that  $\sum_{i=1}^n (y_i - \hat{y})^2$  is as small as possible.



- The length of this distance represents the value  $(y_i - \hat{y})$  (shown as the red line segment in the Figure. We call it **residual**).
- Note that  $(y_i - \hat{y})$  is **positive** when the point  $(x, y)$  is above the line and **negative** when  $(x, y)$  is below the line

# LINEAR REGRESSION

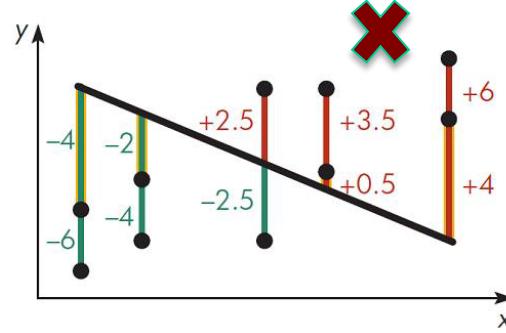
- Figure 3.22 shows a scatter diagram with what appears to be the **line of best fit**, along with 10 individual  $(y_i - \hat{y})$  values. (Positive values are shown in red; negative, in green.)



$\Sigma (y - \hat{y})^2 = (-1)^2 + (+1)^2 + \dots + (+1)^2 = 23.0$

The Line of Best Fit

Figure 3.22



$\Sigma (y - \hat{y})^2 = (-6)^2 + (-4)^2 + \dots + (+6)^2 = 149.0$

Not the Line of Best Fit

Figure 3.23

- Figure 3.23 shows the same data points as Figure 3.22. The 10 individual values  $(y_i - \hat{y})$  are plotted with a line that is definitely not the line of best fit.

# LINEAR REGRESSION

- Our job is to find the one line that will make  $\sum_{i=1}^n (y_i - \hat{y})^2$  the smallest possible value.
- The equation of the line of best fit is determined by its slope ( $b_1$ ) and its  **$y$ -intercept** ( $b_0$ ).
- Slope: 
$$b_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} = \frac{SS(xy)}{SS(x)},$$
 where
  - $SS(xy) = \sum_{i=1}^n x_i y_i - \frac{1}{n} (\sum_i x_i)(\sum_i y_i)$
  - $SS(x) = \sum_{i=1}^n x_i^2 - \frac{1}{n} (\sum_{i=1}^n x_i)^2$
- **$y$ -intercept:** 
$$b_0 = \bar{y} - b_1 \bar{x}$$

# LINEAR REGRESSION

- **Example-5: push-up/sit-up data.**

Student	1	2	3	4	5	6	7	8	9	10
Push-ups, $x$	27	22	15	35	30	52	35	55	40	40
Sit-ups, $y$	30	26	25	42	38	40	32	54	50	43

Student	Push-ups, $x$	$x^2$	Sit-ups, $y$	$y^2$	$xy$
1	27	729	30	900	810
2	22	484	26	676	572
3	15	225	25	625	375
4	35	1,225	42	1,764	1,470
5	30	900	38	1,444	1,140
6	52	2,704	40	1,600	2,080
7	35	1,225	32	1,024	1,120
8	55	3,025	54	2,916	2,970
9	40	1,600	50	2,500	2,000
10	40	1,600	43	1,849	1,720
$\Sigma x = 351$		$\Sigma x^2 = 13,717$	$\Sigma y = 380$	$\Sigma y^2 = 15,298$	$\Sigma xy = 14,257$
sum of $x$		sum of $x^2$	sum of $y$	sum of $y^2$	sum of $xy$

$$SS(x) = \Sigma x^2 - \frac{(\Sigma x)^2}{n} = 13,717 - \frac{(351)^2}{10} = 1396.9$$

$$\Rightarrow \bar{x} = \frac{\sum_{i=1}^n x_i}{n} = \frac{351}{10} = 35.1$$

$$SS(xy) = \Sigma xy - \frac{\Sigma x \Sigma y}{n} = 14,257 - \frac{(351)(380)}{10} = 919.00$$

$$\Rightarrow \bar{y} = \frac{\sum_{i=1}^n y_i}{n} = \frac{380}{10} = 38$$

# LINEAR REGRESSION (EXAMPLE)

Student	1	2	3	4	5	6	7	8	9	10
Push-ups, $x$	27	22	15	35	30	52	35	55	40	40
Sit-ups, $y$	30	26	25	42	38	40	32	54	50	43

- Give  $\bar{x} = 35.1$ ,  $\bar{y} = 38$ ,  $SS(x) = 1396.9$  and  $SS(xy) = 919$
- We want to find the line of best fit,  $\hat{y} = b_0 + b_1x$ , where
- $b_1 = \frac{SS(xy)}{SS(x)} = \frac{919}{1396.9} = 0.6579 = 0.66$
- $b_0 = \bar{y} - b_1\bar{x} = 38 - 0.6579 \times 35.1 = 14.9077 = 14.9$
- Therefore
- $\hat{y} = 14.9 + 0.66x$
- Important Notes:
  - The line goes through  $(\bar{x}, \bar{y})$
  - The slope:  $b_1 = 0.66$
  - The y-intercept:  $b_0 = 14.9$



## EXAMPLE -

### AUSTRALIAN INSTITUTE OF SPORT

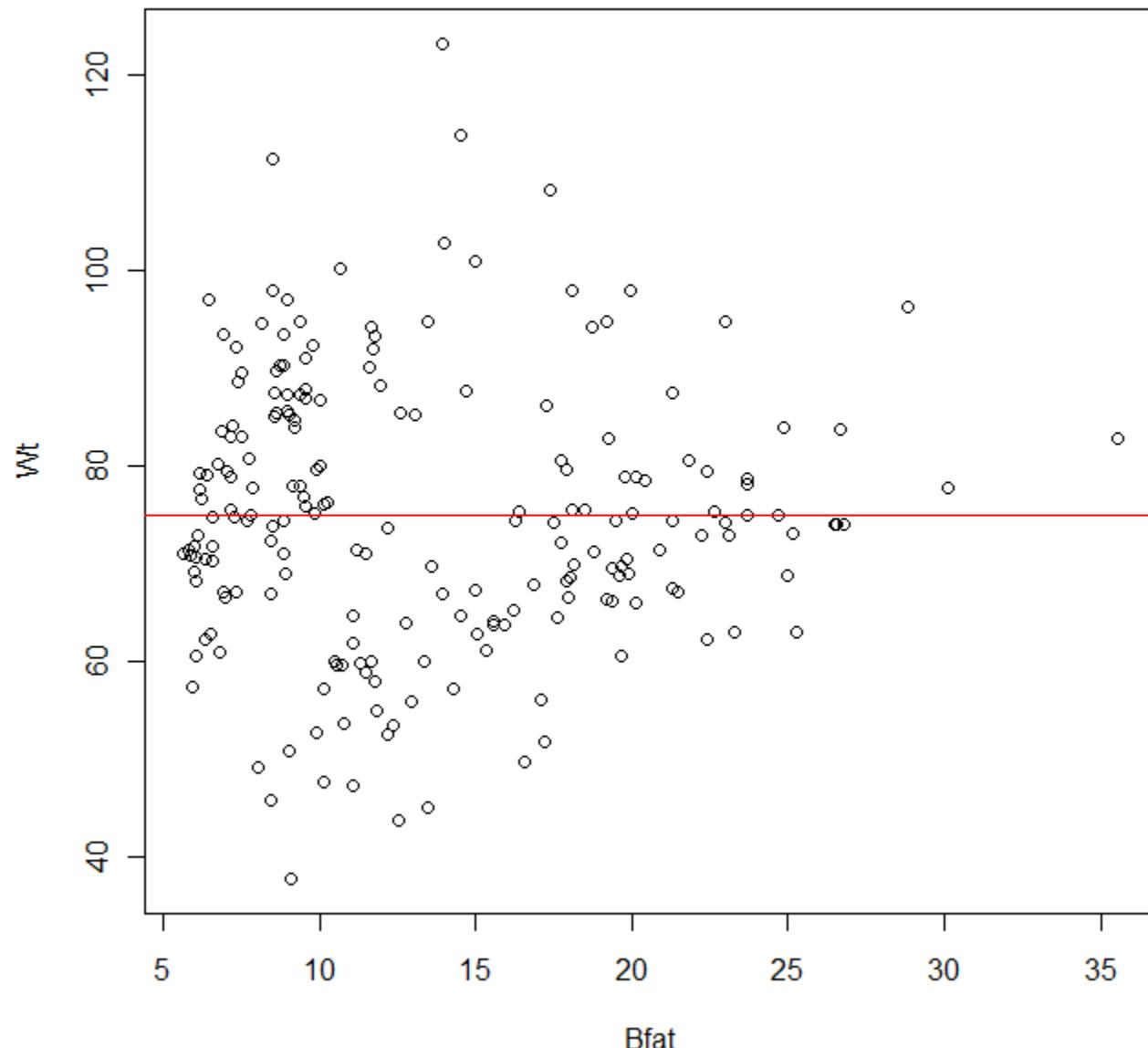
- Data on 102 male and 100 female athletes collected at the Australian Institute of Sport, (courtesy of Richard Telford and Ross Cunningham.)

		<b>Gender</b>	<b>Bfat</b>	<b>Wt</b>
•	<b>1</b>	female	<b>19.75</b>	<b>78.9</b>
•	<b>2</b>	female	<b>21.30</b>	<b>74.4</b>
•	<b>3</b>	female	<b>19.88</b>	<b>69.1</b>
•	<b>4</b>	female	<b>23.66</b>	<b>74.9</b>
•	<b>5</b>	female	<b>17.64</b>	<b>64.6</b>
		:	:	:
•	<b>198</b>	male	<b>11.79</b>	<b>93.2</b>
•	<b>199</b>	male	<b>10.05</b>	<b>80.0</b>
•	<b>200</b>	male	<b>8.51</b>	<b>73.8</b>
•	<b>201</b>	male	<b>11.50</b>	<b>71.1</b>
•	<b>202</b>	male	<b>6.26</b>	<b>76.7</b>



## EXAMPLE CONT'D

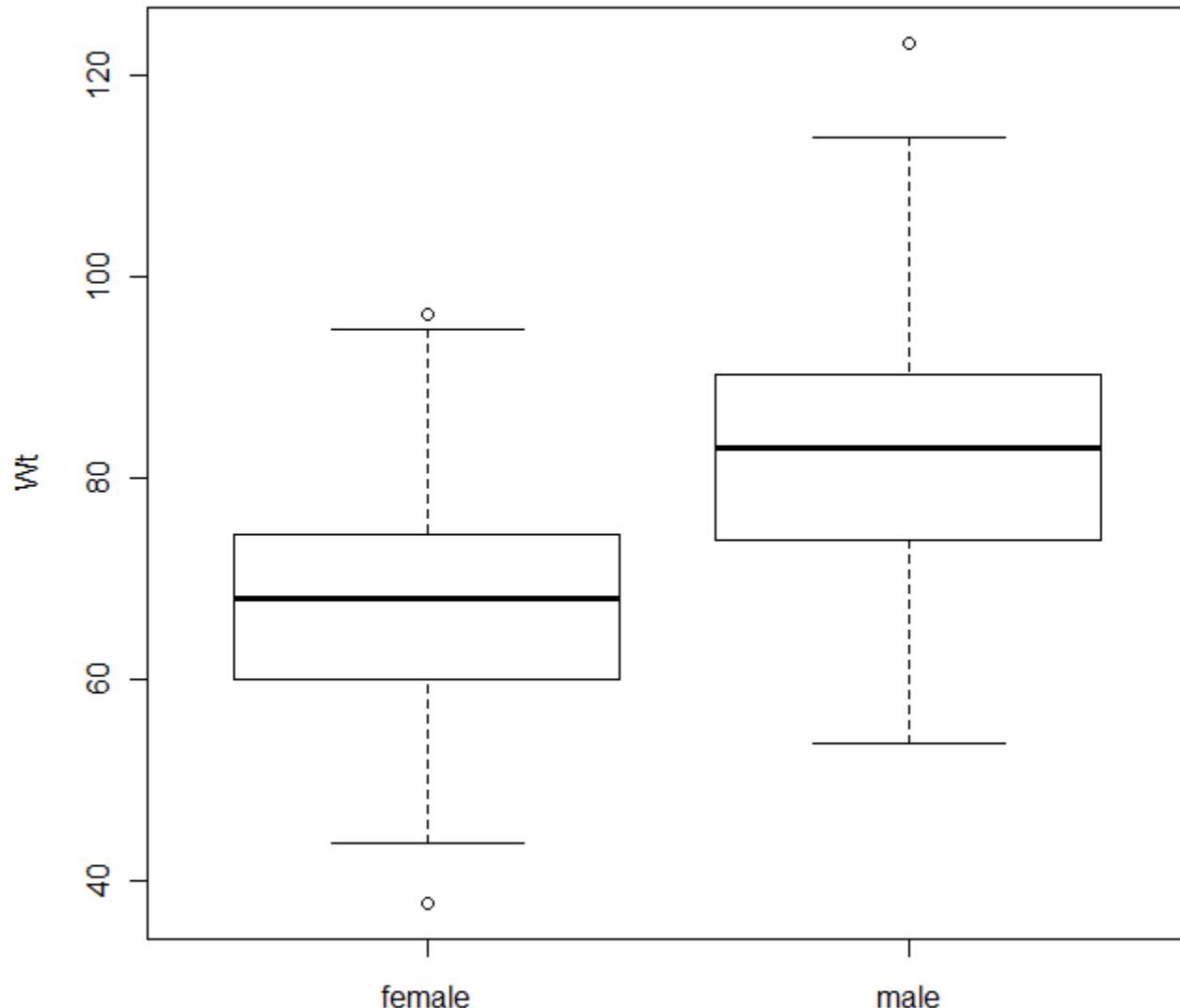
Australian Institute of Sport - Bfat vs Wt





## EXAMPLE CONT'D

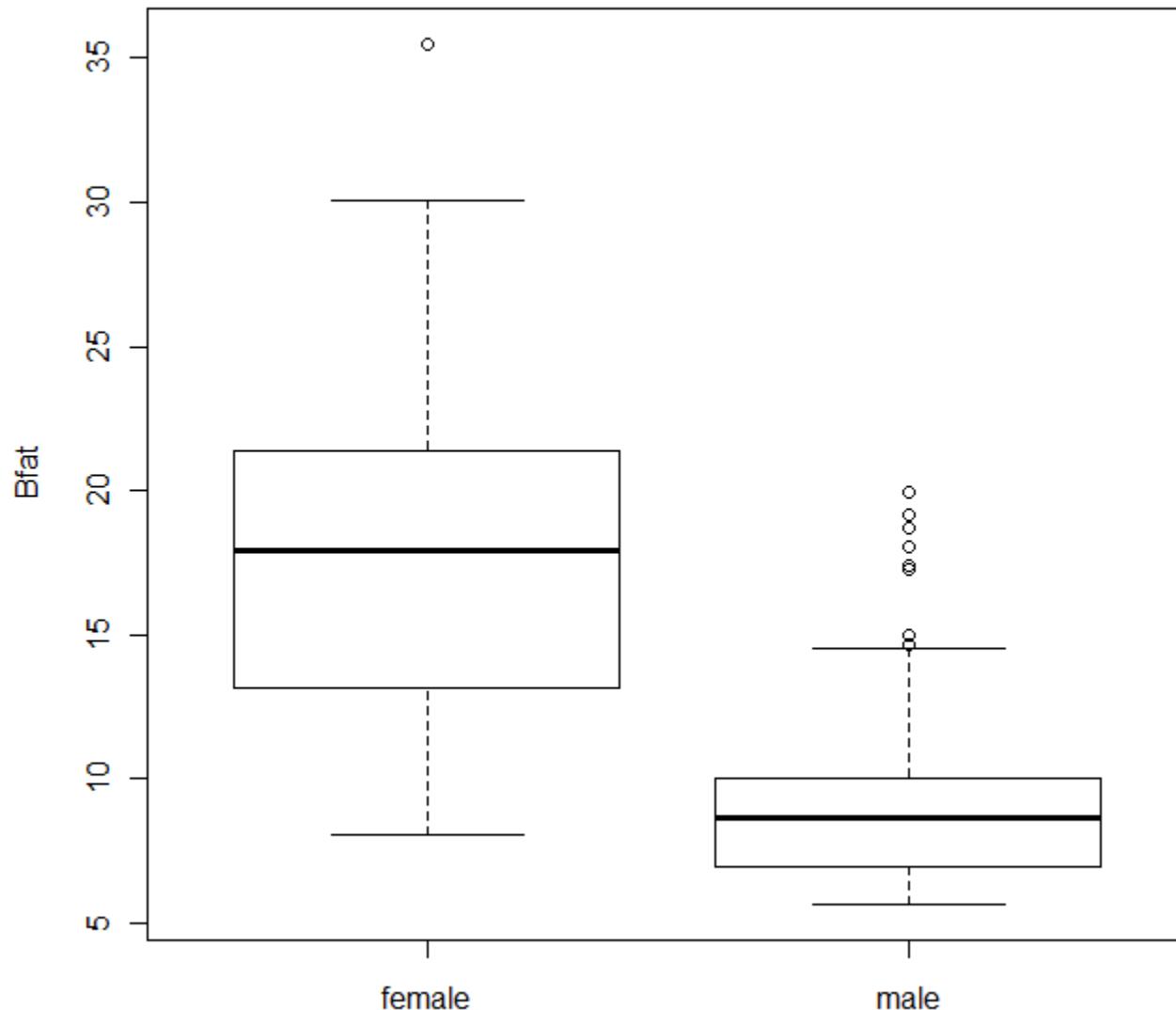
Australian Institute of Sport - Wt vs Gender





## EXAMPLE CONT'D

Australian Institute of Sport - Bfat vs Gender





## EXAMPLE CONT'D

Australian Institute of Sport - Bfat vs Wt

