


Inferential Statistics and Probability a Holistic Approach

Chapter 1

Displaying and Analyzing Data with Graphs


This Course Material by Maurice Geraghty is licensed under a Creative Commons
Attribution-ShareAlike 4.0 International License.
Conditions for use are shown here: <https://creativecommons.org/licenses/by-sa/4.0/>

1

Introduction

- Syllabus– Homework 0
- Projects
- Computer Lab – S44
 - Minitab
- Website
 - <http://nebula2.deanza.edu/~mo>
- Tutor Lab - S43 (S41 for MPS)
 - Drop in or assigned tutors – get form from lab.
 - Group Tutoring
- Other Questions

2

Descriptive Statistics

- Organizing, summarizing and displaying data
 - Graphs
 - Charts
 - Measure of Center
 - Measures of Spread
 - Measures of Relative Standing

3

Problem Solving

- The Role of Probability
- Modeling
- Simulation
- Verification

4

Inferential Statistics

- Population – the set of all measurements of interest to the sample collector
- Sample – a subset of measurements selected from the population
- Inference – A conclusion about the population based on the sample
- Reliability – Measure the strength of the Inference

5

Raw Data – Apple

Monthly Adjusted Stock Price: 12/1998 to 12/2018

115.82	102.97	106.17	75.50	69.86	52.70	41.97	27.42	11.11	25.77	11.04	9.35	4.19	1.39	0.93	1.42	0.97
110.52	115.73	114.39	74.83	76.83	49.73	40.49	26.01	12.06	23.71	11.93	8.82	4.36	1.36	1.01	1.39	1.07
112.96	116.40	103.43	69.93	77.80	52.67	39.16	24.53	14.00	24.72	10.55	7.49	3.41	1.49	1.05	1.14	1.27
112.47	107.44	96.49	63.79	87.18	49.62	36.92	24.12	14.79	19.97	10.02	6.98	2.52	1.35	0.94	1.01	1.68
105.56	109.84	98.16	65.19	86.93	50.07	31.63	21.89	22.06	18.02	8.83	6.10	2.24	1.47	0.96	1.21	3.96
103.12	117.62	91.10	60.15	79.47	50.81	33.47	21.26	20.68	17.14	8.84	5.55	2.10	1.37	0.99	1.22	3.31
94.60	121.63	88.56	52.71	75.99	43.68	32.73	18.53	21.79	15.88	7.45	4.79	2.12	1.24	1.15	1.51	3.41
98.81	126.33	86.17	59.78	75.17	45.26	33.43	17.67	24.56	15.77	7.78	5.17	1.83	1.17	1.52	1.30	2.73
92.20	120.85	79.89	58.47	75.99	45.56	33.97	16.37	22.63	12.99	9.16	4.69	1.68	0.93	1.58	1.66	4.04
107.20	120.15	72.66	58.45	78.01	45.35	30.58	13.68	18.67	12.09	8.16	5.42	1.76	0.92	1.54	1.44	4.42
95.10	124.05	71.24	58.28	70.58	45.96	26.63	11.62	16.27	11.01	8.91	5.84	1.56	0.98	1.41	1.19	3.73
95.22	112.69	67.37	59.80	59.40	44.15	24.99	11.73	17.61	11.16	9.83	5.00	1.47	0.93	1.61	1.41	3.38

6



Crime Rate

- In the last 18 years, has violent crime:
 - Increased?
 - Stayed about the Same?
 - Decreased?

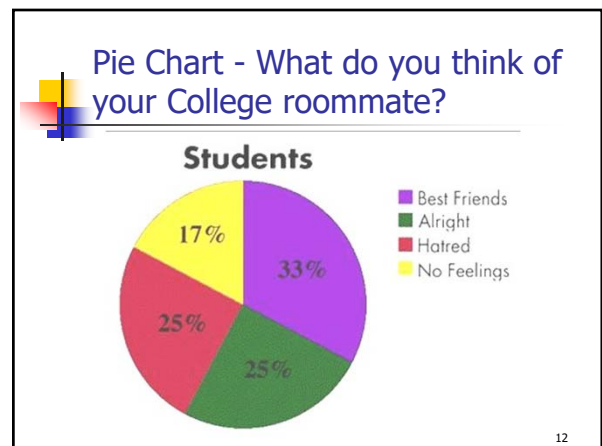
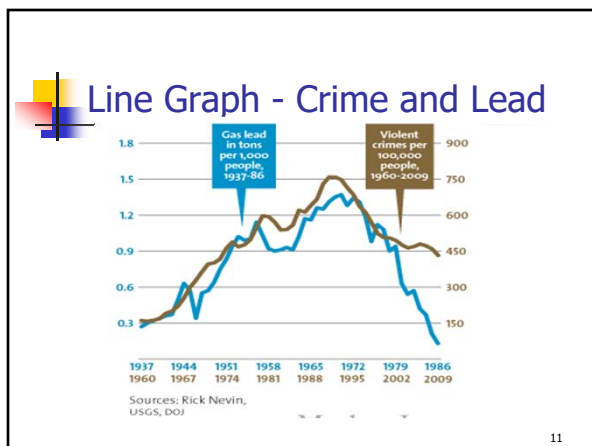
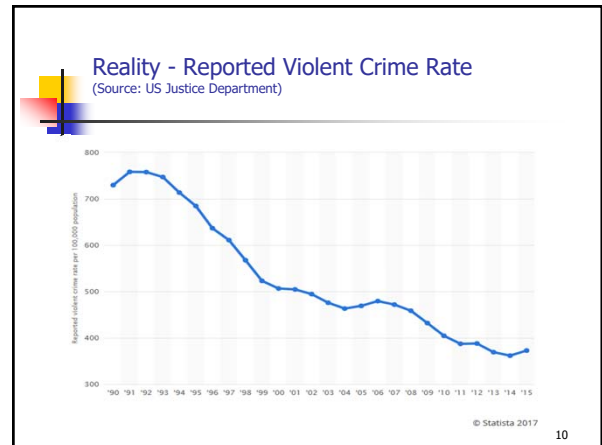
8

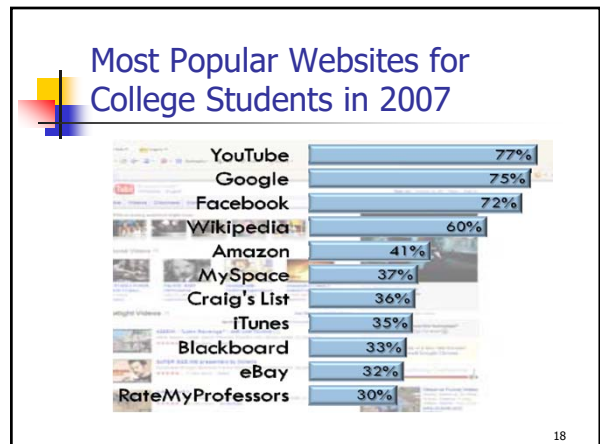
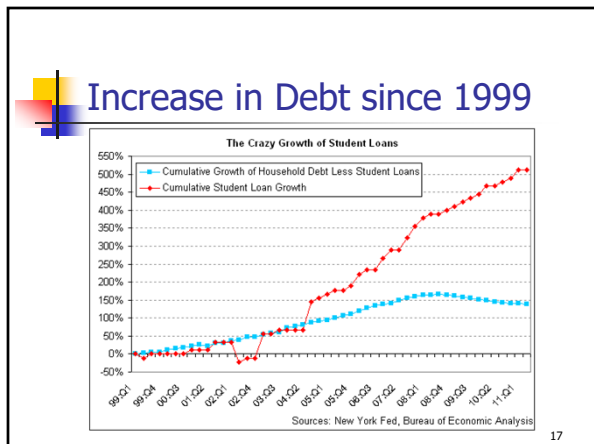
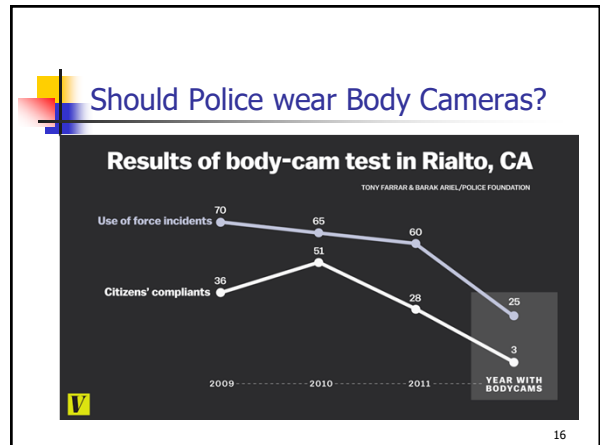
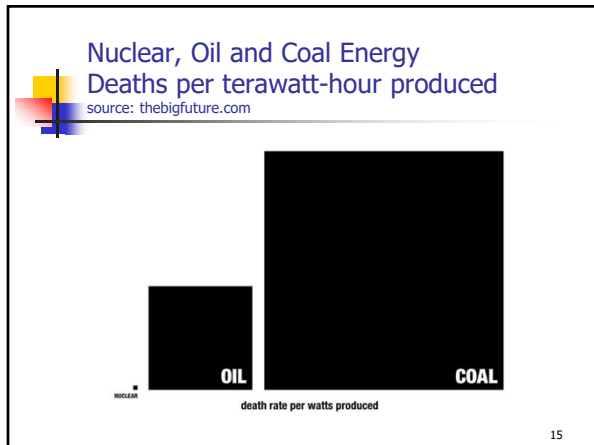
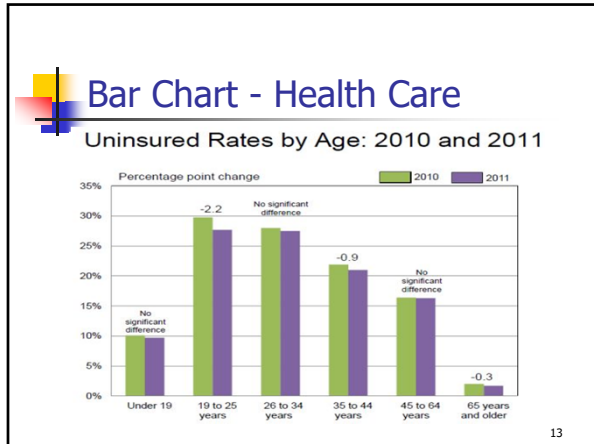
Perception – Gallup Poll

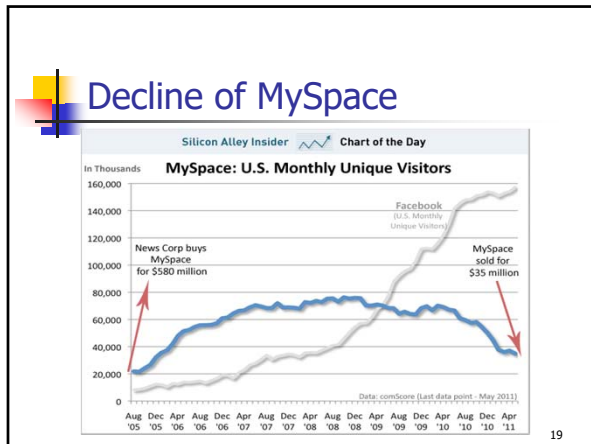
Is there more crime in the U.S. than there was a year ago, or less?

	More	Less	Same (incl)	No opinion
	%	%	%	%
2018 Oct 8-9	70	20	6	4
2018 Oct 7-11	70	18	8	4
2014 Oct 12-15	63	21	9	7
2013 Oct 3-6	64	19	9	7
2011 Oct 6-9	68	17	8	8
2010 Oct 7-10	66	17	8	9
2009 Oct 1-4	74	15	6	6
2008 Oct 3-5	67	15	9	9
2007 Oct 4-7	71	14	8	6
2006 Oct 9-12	69	16	8	8
2006 Oct 12-16	67	21	9	3
2004 Oct 11-14	63	28	14	8
2003 Oct 6-8	60	25	11	4
2002 Oct 14-17	62	21	11	6
2001 Oct 11-14	41	43	10	6
2000 Aug 29-Sep 5	47	41	7	5
1998 Oct 23-25	52	35	8	5
1997 Aug 22-25	64	25	6	5
1996 Jul 25-28	71	15	8	6
1993 Oct 13-18	87	4	5	4
1992 Feb 28-Mar 1	89	3	4	4
1990 Sep 10	84	3	7	6

9







19

RATE MY PROFESSORS

Over 6,000 Schools, 1 million professors, 6 million opinions

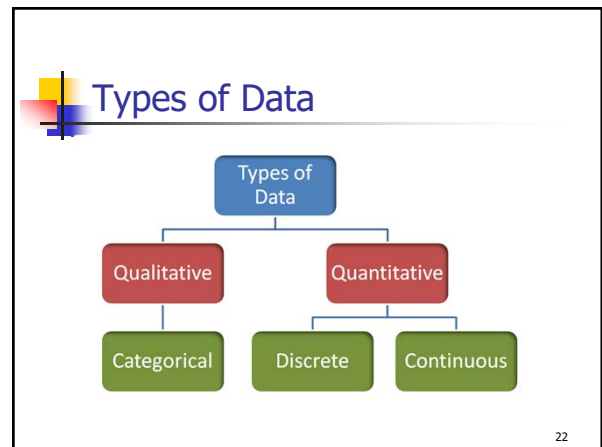
De Anza College

Professor's Name	Department	Total Ratings	Overall Quality	Ease	Hot?
[Profile Icon]	Mandarin	3	4.3	2.0	[Hot Icon]
[Profile Icon]	Mandarin	8	1.6	1.6	[Hot Icon]
[Profile Icon]	Marketing	1	5.0	5.0	[Hot Icon]
[Profile Icon]	Mathematics	66	4.7	4.0	[Hot Icon]
[Profile Icon]	Mathematics	73	1.4	1.7	[Hot Icon]
[Profile Icon]	Mathematics	15	2.7	2.6	[Hot Icon]
[Profile Icon]	Mathematics	41	1.6	2.1	[Hot Icon]

20

- ### Types of Data
- Qualitative
 - Non-numeric
 - Always categorical
 - Quantitative
 - Numeric
 - Categorical numbers are actually qualitative
 - Continuous or discrete

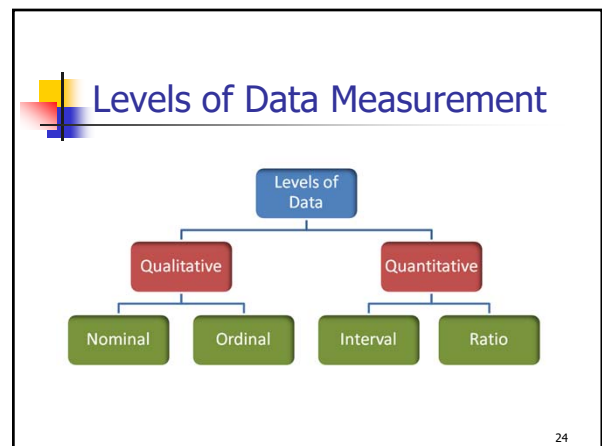
21



22

- ### Levels of Data Measurement
- **Nominal:** Names or labels only
 - Example: What city do you live in?
 - **Ordinal:** Data can be ranked, but no quantifiable difference.
 - Example: Ratings Excellent, Good, Fair, Poor
 - **Interval:** Data can be ranked with quantifiable differences, but no true zero.
 - Example: Temperature
 - **Ratio:** Data can be ranked with quantifiable differences and there is a true zero.
 - Example: Age

23



24

Examples of Data

- Distance from De Anza College
- Number of Grandparents still alive
- Eye Color
- Amount you spend on food each week.
- Number of Facebook "Friends"
- Zip Code
- City you live in.
- Year of Birth
- How to prepare Steak? (rare, medium, well-done)
- Do you drive to De Anza?

25

Graphical Methods

- Qualitative Data
 - Pie Chart
 - Bar Chart
- Quantitative Data
 - Stem and Leaf Chart
 - Histogram
 - Ogive
 - Dot Plot

26

Graphing Categorical Data

A sample of 500 adults (age 18 and over) from Santa Clara County, California were taken from the year 2000 United States Census.

Marital Status	Frequency
Married	270
Widowed	22
Divorced - not remarried	42
Separated	10
Single - never married	156
Total	500

27

Graphing Categorical Data

- n = sample size** - The number of observations in your sample size.
- Frequency** - the number of times a particular value is observed.
- Relative frequency** - The proportion or percentage of times a particular value is observed.
- Relative Frequency = Frequency / n**

28

Graphing Categorical Data

A sample of 500 adults (age 18 and over) from Santa Clara County, California were taken from the year 2000 United States Census.

Marital Status	Frequency	Relative Frequency
Married	270	270/500 = 0.540 or 54.0%
Widowed	22	22/500 = 0.044 or 4.4%
Divorced - not remarried	42	42/500 = 0.084 or 8.4%
Separated	10	10/500 = 0.020 or 2.0%
Single - never married	156	156/500 = 0.312 or 31.2%
Total	500	500/500 = 1.000 or 100.0%

29

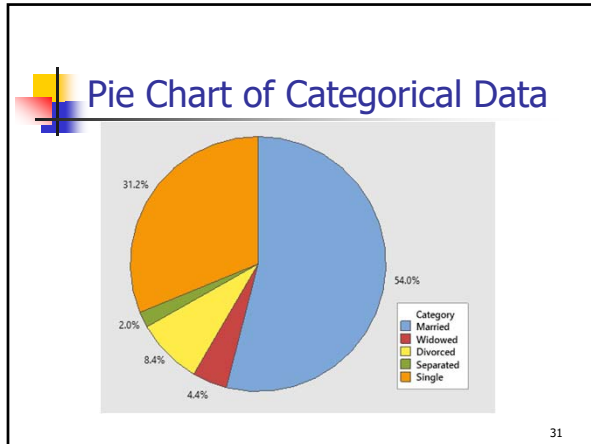
Bar Graph of Categorical Data

Marital Status of 500 Adults in Santa Clara County

Marital Status	Percentage
Married	54
Widowed	4.4
Divorced	8.4
Separated	2
Single	31.2

Percent within all data.

30



Daily Minutes spent on the Internet by 30 students

102	104	85	67	101
71	116	107	99	82
103	97	105	103	95
105	99	86	87	100
109	108	118	87	125
124	112	122	78	92

- ### Describing Numeric Data
- Center?
 - Where is an "average" value
 - Spread?
 - How far are data spread from the center
 - Shape?
 - Symmetric or skewed?
 - Anything Unusual?
 - Outliers, more than 1 peak?

Stem and Leaf Graph

```

6 | 7
7 | 18
8 | 25677
9 | 25799
10 | 01233455789
11 | 268
12 | 245
    
```

Back-to-back Example

- Passenger loading times for two airlines

11, 14, 16, 17,	8, 11, 13, 14,
19, 21, 22, 23,	15, 16, 16, 18,
24, 24, 24, 26,	19, 19, 21, 21,
31, 32, 38, 39	22, 24, 26, 31

Back to Back Example

```

      | 0 |
      | 0 | 8
    14 | 1 | 134
    679 | 1 | 566899
  123444 | 2 | 1124
        | 2 | 6
        | 2 | 6
        | 3 | 1
    89 | 3 |
    
```

Grouping Data

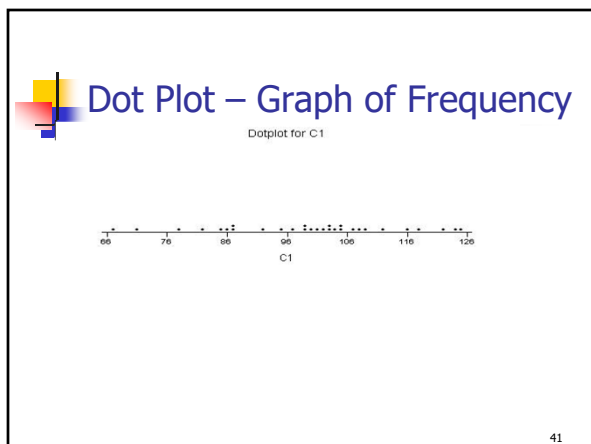
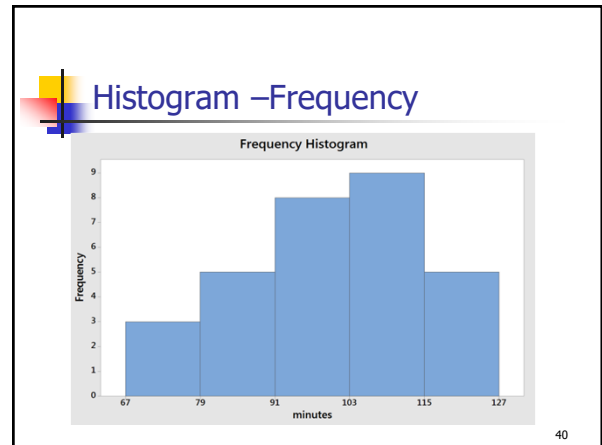
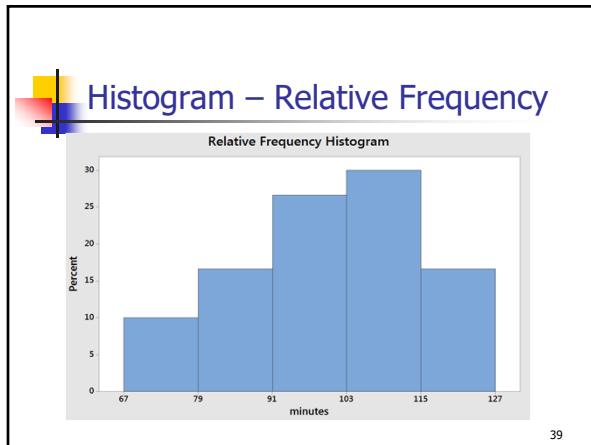
- Choose the number of groups
 - between 5 and 10 is best
- Interval Width = $(\text{Range}+1)/(\text{Number of Groups})$
 - Round **up** to a convenient value
- Start with lowest value and create the groups.
- Example – for 5 categories
Interval Width = $(58+1)/5 = 12$ (rounded up)

37

Grouping Data

Class Interval	Frequency	Relative Frequency
67 to 79	3	0.100 or 10.0%
79 to 91	5	0.167 or 16.7%
91 to 103	8	0.266 or 26.6%
103 to 115	9	0.300 or 30.0%
115 to 127	5	0.167 or 16.7%
Total	30	1.000 or 100%

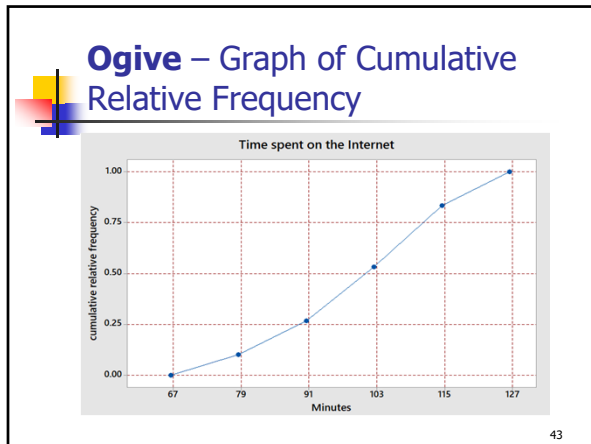
38



Cumulative Relative Frequency


Class Interval	Frequency	Relative Frequency	Cumulative Frequency	Cumulative Relative Frequency
67 to 79	3	0.100 or 10.0%	3	0.100 or 10.0%
79 to 91	5	0.167 or 16.7%	8	0.267 or 26.7%
91 to 103	8	0.266 or 26.6%	16	0.533 or 53.3%
103 to 115	9	0.300 or 30.0%	25	0.833 or 83.3%
115 to 127	5	0.167 or 16.7%	30	1.000 or 100%
Total	30	1.000 or 100%		

42



Inferential Statistics and Probability a Holistic Approach

Chapter 2 Descriptive Statistics



This Course Material by Maurice Geraghty is licensed under a Creative Commons Attribution-ShareAlike 4.0 International License.
Conditions for use are shown here: <https://creativecommons.org/licenses/by-sa/4.0/>

1

Measures of Central Tendency

- Mean
 - Arithmetic Average $\bar{X} = \frac{\sum X_i}{n}$
- Median
 - "Middle" Value after ranking data
 - Not affected by "outliers"
- Mode
 - Most Occurring Value
 - Useful for non-numeric data

2

Example

Anthony's Pizza, a Detroit based company, offers pizza delivery to its customers. A driver for Anthony's Pizza will often make several deliveries on a single delivery run. A sample of 5 delivery runs by a driver showed that the total number of pizzas delivered on each run

2 2 5 9 12

What is the Average?

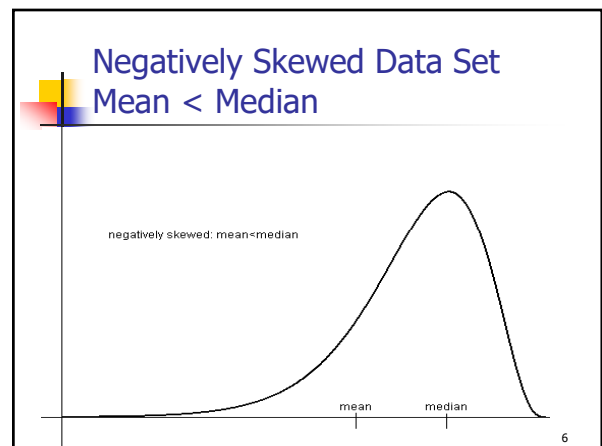
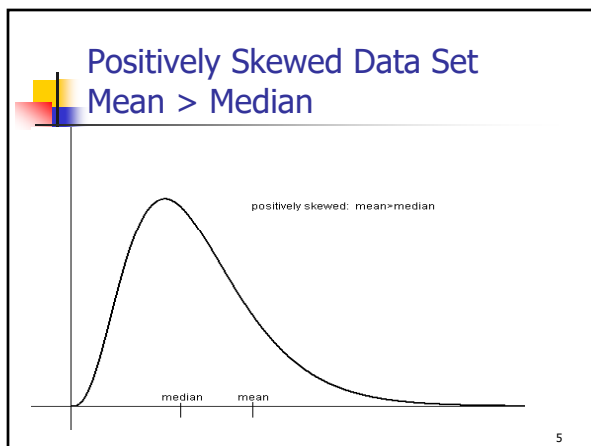
- a) 2
- b) 5
- c) 6

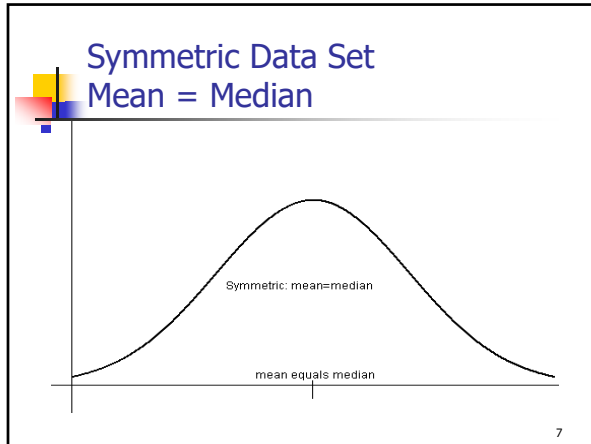
3

Example – 5 Recent Home Sales

- \$500,000
- \$600,000
- \$600,000
- \$700,000
- \$2,600,000

4





- ### Measures of Variability
- Range
 - Variance
 - Standard Deviation
 - Interquartile Range (percentiles)
- 8

Range

$$\text{Max}(X_i) - \text{Min}(X_i)$$

$$125 - 67 = 58$$

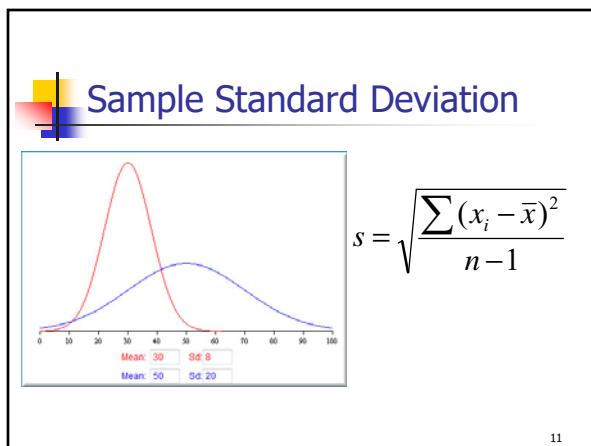
9

Sample Variance

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

$$s^2 = \frac{\sum x_i^2 - (\sum x_i)^2 / n}{n - 1}$$

10



Variance and Standard Deviation

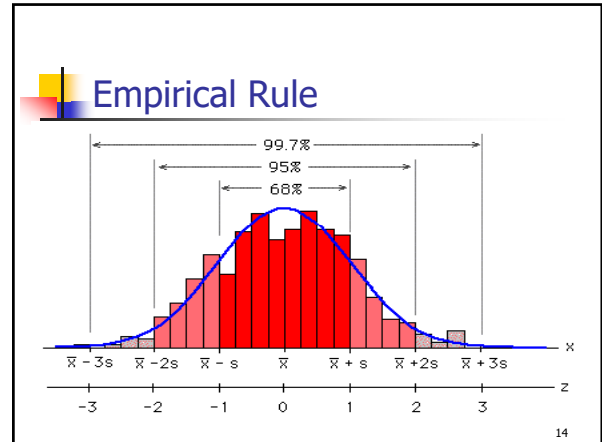
X_i	$X_i - \bar{X}$	$(X_i - \bar{X})^2$	
2	-4	16	$s^2 = \frac{78}{4} = 19.5$
2	-4	16	
5	-1	1	
9	3	9	$s = \sqrt{19.5} \approx 4.42$
<u>12</u>	<u>6</u>	<u>36</u>	
30	0	78	

12

Interpreting the Standard Deviation

- Chebyshev's Rule
 - At least $100 \times (1 - (1/k)^2)\%$ of any data set must be within k standard deviations of the mean.
- Empirical Rule (68-95-99 rule)
 - Bell shaped data
 - 68% within 1 standard deviation of mean
 - 95% within 2 standard deviations of mean
 - 99.7% within 3 standard deviations of mean

13



Measures of Relative Standing

- Z-score
- Percentile
- Quartiles
- Box Plots

15

Z-score

- The number of Standard Deviations from the Mean
- $Z > 0$, X_i is greater than mean
- $Z < 0$, X_i is less than mean

$$Z = \frac{X_i - \bar{X}}{s}$$

16

Percentile Rank

Formula for ungrouped data

- The location is $(n+1)p$ (interpolated or rounded)
- n = sample size
- p = percentile

17

Quartiles

- 25th percentile is 1st quartile
- 50th percentile is median
- 75th percentile is 3rd quartile
- 75th percentile – 25th percentile is called the Interquartile Range which represents the "middle 50%"

18

IQR example

$n+1=31$

$.25 \times 31 = 7.75$ location 8 = **87** ← 1st Quartile

$.75 \times 31 = 23.25$ location 23 = **108** ← 3rd Quartile

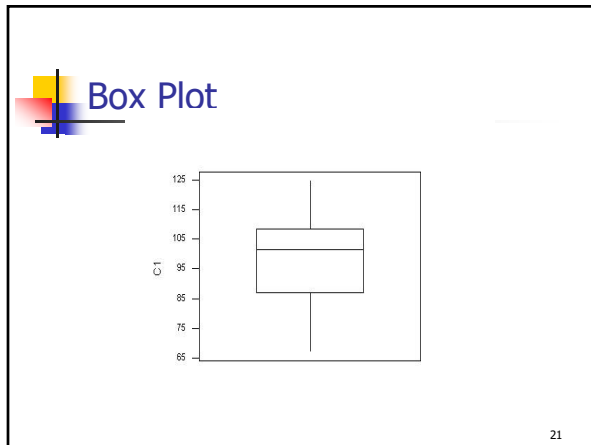
Interquartile Range (IQR) = $108 - 87 = 21$

19

Box Plots

- A **box plot** is a graphical display, based on quartiles, that helps to picture a set of data.
- Five pieces of data are needed to construct a box plot:
 - Minimum Value
 - First Quartile
 - Median
 - Third Quartile
 - Maximum Value.

20



Outliers

- An outlier is data point that is far removed from the other entries in the data set.
- Outliers could be
 - Mistakes made in recording data
 - Data that don't belong in population
 - True rare events

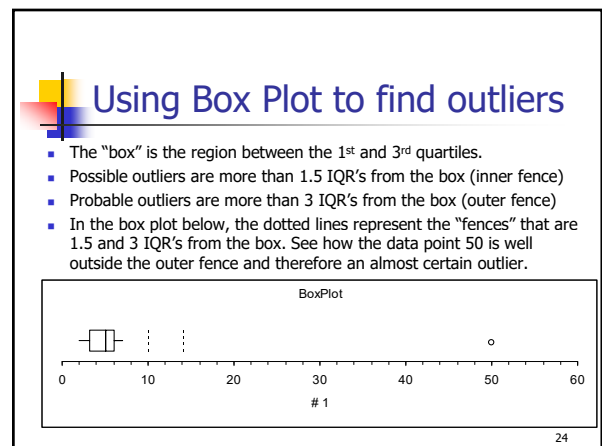
22

Outliers have a dramatic effect on some statistics

- Example quarterly home sales for 10 realtors:

	2	2	3	4	5	5	6	6	7	50
	with outlier					without outlier				
Mean	9.00					4.44				
Median	5.00					5.00				
Std Dev	14.51					1.81				
IQR	3.00					3.50				

23



Using Z-score to detect outliers

- Calculate the mean and standard deviation without the suspected outlier.
- Calculate the Z-score of the suspected outlier.
- If the Z-score is more than 3 or less than -3, that data point is a probable outlier.

$$Z = \frac{50 - 4.4}{1.81} = 25.2$$

25

Outliers – what to do

- Remove or not remove, there is no clear answer.
- For some populations, outliers don't dramatically change the overall statistical analysis. Example: the tallest person in the world will not dramatically change the mean height of 10000 people.
- However, for some populations, a single outlier will have a dramatic effect on statistical analysis (called "**Black Swan**" by Nicholas Taleb) and inferential statistics may be invalid in analyzing these populations. Example: the richest person in the world will dramatically change the mean wealth of 10000 people.

26

Bivariate Data

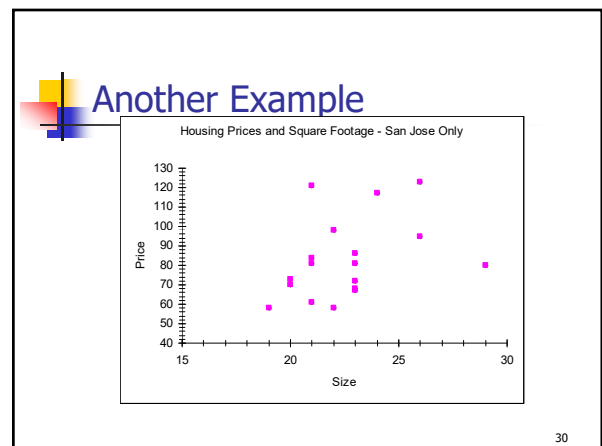
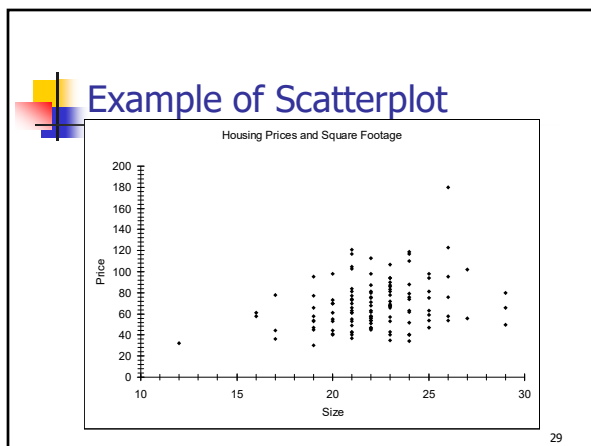
- Ordered numeric pairs (X,Y)
- Both values are numeric
- Paired by a common characteristic
- Graph as Scatterplot

27

Example of Bivariate Data

- Housing Data
 - X = Square Footage
 - Y = Price

28



Correlation Analysis

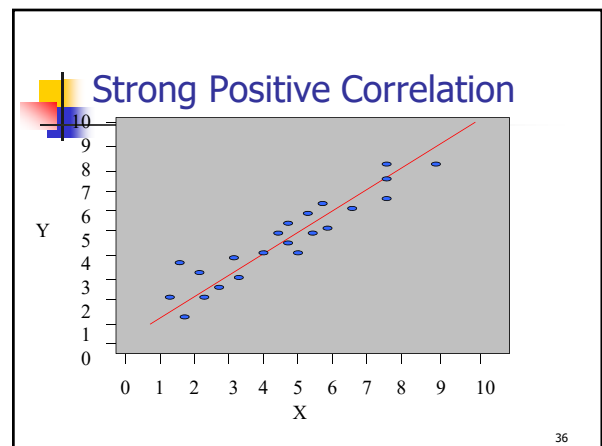
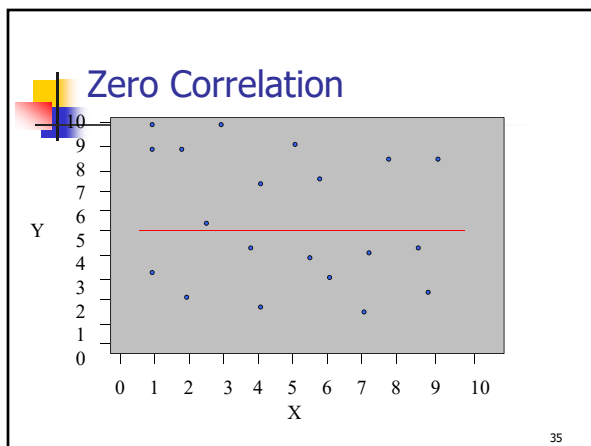
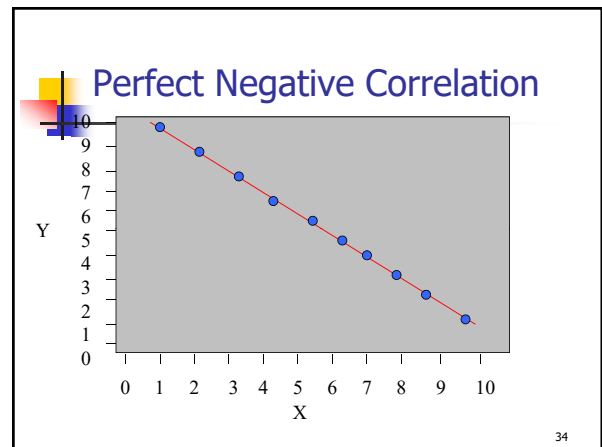
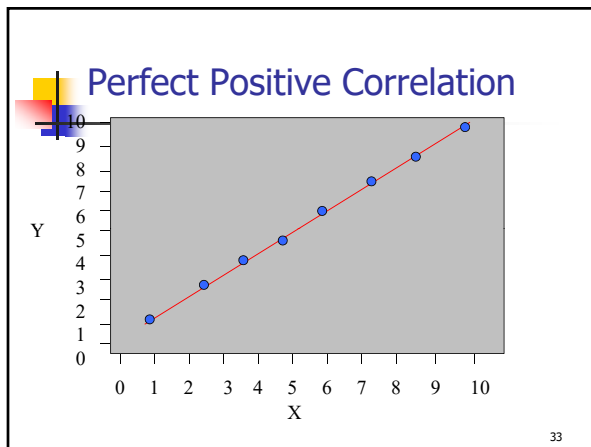
- **Correlation Analysis:** A group of statistical techniques used to measure the strength of the relationship (correlation) between two variables.
- **Scatter Diagram:** A chart that portrays the relationship between the two variables of interest.
- **Dependent Variable:** The variable that is being predicted or estimated. "Effect"
- **Independent Variable:** The variable that provides the basis for estimation. It is the predictor variable. "Cause?" (Maybe!)

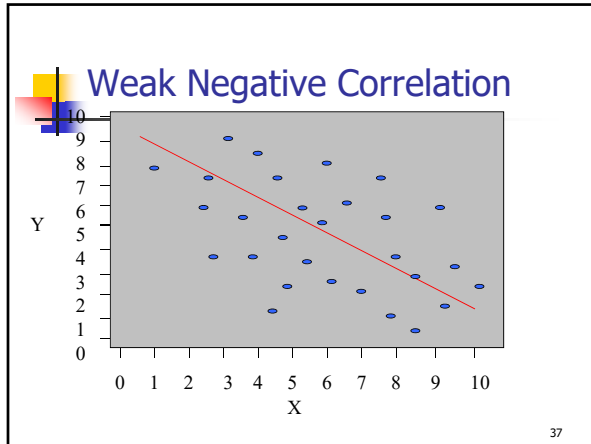
31

The Coefficient of Correlation, r

- **The Coefficient of Correlation (r)** is a measure of the **strength** of the relationship between two variables.
 - It requires interval or ratio-scaled data (variables).
 - It can range from -1 to 1.
 - Values of -1 or 1 indicate perfect and strong correlation.
 - Values close to 0 indicate weak correlation.
 - Negative values indicate an inverse relationship and positive values indicate a direct relationship.

32





- ### Causation
- Correlation does not necessarily imply causation.
 - There are 4 possibilities if X and Y are correlated:
 1. X causes Y
 2. Y causes X
 3. X and Y are caused by something else.
 4. Confounding - The effect of X and Y are hopelessly mixed up with other variables.
- 38

- ### Causation - Examples
- City with more police per capita have more crime per capita.
 - As Ice cream sales go up, shark attacks go up.
 - People with a cold who take a cough medicine feel better after some rest.
- 39

Formula for correlation coefficient r

$$r = \frac{SSXY}{\sqrt{SSX \cdot SSY}}$$

$$SSX = \sum X^2 - \frac{1}{n}(\sum X)^2$$

$$SSY = \sum Y^2 - \frac{1}{n}(\sum Y)^2$$

$$SSXY = \sum XY - \frac{1}{n}(\sum X \cdot \sum Y)$$

40

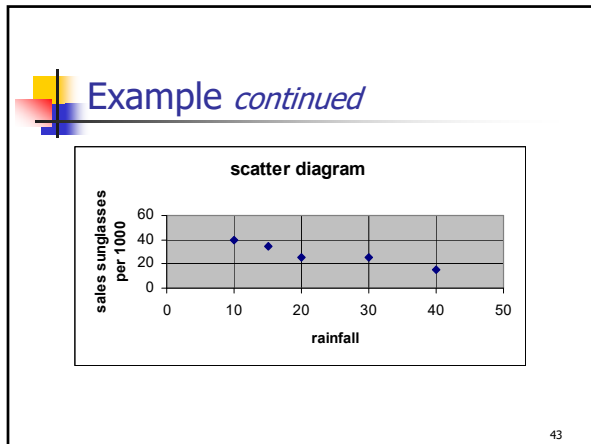
Example

- X = Average Annual Rainfall (Inches)
- Y = Average Sale of Sunglasses/1000
- Make a Scatter Diagram
- Find the correlation coefficient

X	10	15	20	30	40
Y	40	35	25	25	15

41

- ### Example *continued*
- Make a Scatter Diagram
 - Find the correlation coefficient
- 42



Example *continued*

X	Y	X ²	Y ²	XY
10	40	100	1600	400
15	35	225	1225	525
20	25	400	625	500
30	25	900	625	750
40	15	1600	225	600
115	140	3225	4300	2775

- $SSX = 3225 - 115^2/5 = 580$
- $SSY = 4300 - 140^2/5 = 380$
- $SSXY = 2775 - (115)(140)/5 = -445$

44


Example *continued*

$$r = \frac{SSXY}{\sqrt{SSX \cdot SSY}}$$

$$r = \frac{-445}{\sqrt{580 \cdot 380}} = -0.9479$$


- Strong negative correlation

45




Inferential Statistics and Probability a Holistic Approach

Chapter 3 Populations and Sampling



This Course Material by Maurice Geraghty is licensed under a Creative Commons Attribution-ShareAlike 4.0 International License.
Conditions for use are shown here: <https://creativecommons.org/licenses/by-sa/4.0/>


1



Population vs. Sample

- A **population** is the entire group of individuals or objects of interest to us.
- A **sample** is a subset of the population that we can study by collecting or gathering data.
- Quantities that describe populations are called **parameters**.
- Quantities that describe samples are called **statistics**.


2



Example

- A large community college has about 25,000 students. In a study of 85 students from college, it was determined that about 60 of the students have moderate or high math anxiety.
- The **population** is **all** the students at this college.
- The **sample** is the 85 students whose math anxiety was measured.


3



Steps of a Statistical Process

- **Step 1 (Problem)**
Ask a question that can be answered with sample data.
- **Step 2 (Plan)**
Determine what information is needed.
- **Step 3 (Data)**
Collect sample data that is representative of the population.
- **Step 4 (Analysis)**
Summarize, interpret and analyze the sample data.
- **Step 5 (Conclusion)**
State the results and conclusion of the study.


4



Representative Sample

- A **representative sample** has characteristics, behaviors and attitudes similar to the population from which the sample is selected.
- A sample that is not representative is a **biased sample**.

5



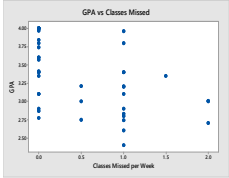
Observational Study

- An **observational study** starts with selecting a representative sample from a population.
- The researcher then takes measurements from the sample, but does not manipulate any of the variables with treatments.
- The goal of an observational study is to interpret and analyze the measured variables, but it is not possible to show a cause and effect relationship.

6

Example of Observational Study

- A group of students at Georgia College conducted a survey asking random students various questions about their scholastic profile.
- One part of their study was to see if there is any correlation between various students' GPA and classes missed.



The scatter plot shows GPA on the y-axis (ranging from 2.0 to 4.0) and Classes Missed per Week on the x-axis (ranging from 0.0 to 2.0). There is a clear downward trend, indicating that as the number of classes missed increases, the GPA tends to decrease.

7

Experiment

- An **experiment** starts with a representative sample from a population.
- The researcher will then randomly break this sample into groups and then apply treatments in order to manipulate a variable of interest.
- The goal of an experiment is to find a cause and effect relationship between a random variable in the population and the variable manipulated by the researcher.
- If an experiment is conducted properly, the researcher can control for confounding or lurking variables and test for a **placebo effect**.

8

Example of Experiment

- Researchers were studying gambling addiction by speed of play using electronic gaming machines.
- 62 participants played a computerized slot machine with either fast, medium, or slow play.
- Gambling speed had no overall effect on either mean bet size, game evaluations or illusion of control, but in the fast machines, at-risk gamblers employed higher bet sizes compared to no-risk gamblers.
- The findings corroborate and elaborate on previous studies and indicate that restrictions on gambling speed may serve as a harm reducing effort for at-risk gamblers.

9

Variables in an Experiment

- **Explanatory Variable:** The variable that is controlled or manipulated by the researcher.
- **Response Variable:** The variable which is being measured and is the focus of the study.
- The researcher tries to answer the question: "Does the explanatory variable (cause) affect the response variable (effect)?"
- In the prior gambling example, the explanatory variable was the speed of the machine, and the response variable was the bet size.

10

Placebos and Blinding


- A **placebo effect** is when participant will respond in a positive way to a treatment with no active ingredients.
- This treatment with no active ingredients is called a **placebo**.
- A **single blind study** is where the participant does not know whether the treatment is real or a placebo.
- A **double blind study** is where neither the administrator of the treatment nor the participant knows whether the treatment is real or a placebo.

11

Example

- An researcher for a pharmaceutical company is conducting research on an experimental drug to reduce the pain from migraine headaches.
- Participants with migraine headaches are randomly split into 3 groups. The first group gets the experimental drug (**Treatment Group**). The second group gets a placebo, a fake drug (**Placebo Group**). The third group gets nothing (**Control Group**).
- The researcher found that pain was reduced for both the treatment group and the placebo group, establishing a placebo effect. The researcher must then compare the amount of pain reduction in the treatment group to the placebo group to determine if the treatment was effective.


12



Probability Sampling Methods

- Properly done, probability or scientific sampling will produce a representative sample.
 - Simple Random Sampling
 - Stratified Sampling
 - Systematic Sampling
 - Cluster Sampling


13



Non-Probability Sampling Methods

- Non-probability sampling methods have immeasurable biases and will usually not produce a representative sample.
 - Convenience Sampling
 - Self-selected Sampling

14




Sources of Bias in Sampling

- **Selection bias** – when the sampling method does not produce a representative sample.
- **Self-selection bias** – when participants who volunteer are not representative of the population.
- **Non-response bias** – when people are intentionally or non-intentionally excluded from participation or choose not to participate in a survey or poll.
- **Response bias** – when the wording of the questions in surveys affect the response.

15

Inferential Statistics and Probability a Holistic Approach

Chapter 4 Probability


This Course Material by Maurice Geraghty is licensed under a Creative Commons Attribution-ShareAlike 4.0 International License. Conditions for use are shown here: <https://creativecommons.org/licenses/by-sa/4.0/>

1

Probability

- Classical probability
 - Based on mathematical formulas
- Empirical probability
 - Based on the relative frequencies of historical data.
- Subjective probability
 - "one-shot" educated guess.

2

Examples of Probability

- What is the probability of rolling a four on a 6-sided die?
- What percentage of De Anza students live in Cupertino?
- What is the chance that the Golden State Warriors will be NBA champions in 2018?

3

Classical Probability

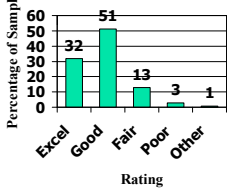
- Event
 - A result of an experiment
- Outcome
 - A result of the experiment that cannot be broken down into smaller events
- Sample Space
 - The set of all possible outcomes
- Probability Event Occurs
 - # of elements in Event / # Elements in Sample Space
- Example – flip two coins, find the probability of exactly 1 head.
 - {HH, HT, TH, TT}
 - $P(1 \text{ head}) = 2/4 = .5$

4

Empirical Probability

- Historical Data
- Relative Frequencies
- Example: What is the chance someone rates their community as good or better?
 - $0.51 + 0.32 = 0.83$

National: Rate Your community

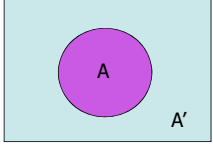


Rating	Percentage of Sample
Excel	32
Good	51
Fair	13
Poor	3
Other	1

5

Rule of Complement

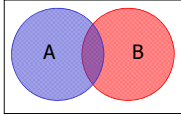
- Complement of an event
- The event does not occur
- A' is the complement of A
- $P(A) + P(A') = 1$
- $P(A) = 1 - P(A')$



6

Additive Rule

- The **UNION** of two events A and B is that either A or B occur (or both). (All colored parts)
- The **INTERSECTION** of two events A and B is that both A and B will occur. (Purple Part only)
- Additive Rule:
 $P(A \text{ or } B) = P(A) + P(B) - P(A \text{ and } B)$



7

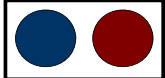
Example

- In a group of students, 40% are taking Math, 20% are taking History.
- 10% of students are taking both Math and History.
- Find the Probability of a Student taking either Math or History or both.
- $P(M \text{ or } H) = 40\% + 20\% - 10\% = 50\%$

8

Mutually Exclusive

- Mutually Exclusive
- Both cannot occur
- If A and B are mutually exclusive, then
 - $P(A \text{ or } B) = P(A) + P(B)$
- Example roll a die
 - A: Roll 2 or less B: Roll 5 or more
 - $P(A)=2/6$ $P(B)=2/6$
 - $P(A \text{ or } B) = P(A) + P(B) = 4/6$



9

Conditional Probability

- The probability of an event occurring GIVEN another event has already occurred.
- $P(A|B) = P(A \text{ and } B) / P(B)$
- Example: Of all cell phone users in the US, 15% have a smart phone with AT&T. 25% of all cell phone users use AT&T. Given a selected cell phone user has AT&T, find the probability the user also has a smart phone.
- A=AT&T subscriber B=Smart Phone User
- $P(A \text{ and } B) = .15$ $P(A)=.25$
- $P(B|A) = .15/.25 = .60$

10

Contingency Tables

- Two data items can be displayed in a contingency table.
- Example: auto accident during year and DUI of driver.

	Accident	No Accident	Total
DUI	70	130	200
Non- DUI	30	770	800
Total	100	900	1000

11

Contingency Tables

	Accident	No Accident	Total
DUI	70	130	200
Non- DUI	30	770	800
Total	100	900	1000

Given the Driver is DUI, find the Probability of an Accident.

A=Accident D=DUI

$P(A \text{ and } D) = .07$ $P(D) = .2$

$P(A|D) = .07/.2 = .35$

12

Marginal, Joint and Conditional Probability

- **Marginal Probability** means the probability of a single event occurring.
- **Joint Probability** means the probability of the union or intersection of multiple events occurring.
- **Conditional Probability** means the probability of an event occurring given that another event has already occurred.

13

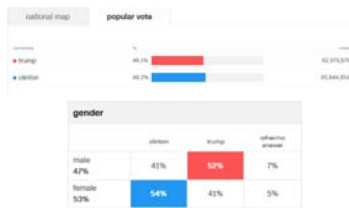
Creating Contingency Tables

- You can create a hypothetical contingency table from reported cross tabulated data.
- First choose a convenient sample size (called a radix) like 10000.
- Then apply the reported marginal probabilities to the radix of one of the variables.
- Then apply the reported conditional probabilities to the total values of one of the other variable.
- Complete the table with arithmetic.

14

Example

Create a two-way table from the cross tabulation of gender from the 2016 election results (from CNN)



15

Example

First select a radix (sample size) of 10000

VOTED FOR	GENDER		Total
	Female	Male	
Trump			
Clinton			
Other			
Total			10000

16

Example

Then apply the marginal probabilities to the radix (53% female, 47% male)

VOTED FOR	GENDER		Total
	Female	Male	
Trump			
Clinton			
Other			
Total	5300	4700	10000

17

Example

Then apply the cross tabulated percentages for each gender. Make sure the numbers add up.

VOTED FOR	GENDER		Total
	Female	Male	
Trump	2173	2444	
Clinton	2862	1927	
Other	265	329	
Total	5300	4700	10000

18

Example

Finally, complete the table using arithmetic.

VOTED FOR	GENDER		Total
	Female	Male	
Trump	2173	2444	4617
Clinton	2862	1927	4789
Other	265	329	594
Total	5300	4700	10000

19

Multiplicative Rule

- $P(A \text{ and } B) = P(A) \times P(B|A)$
- $P(A \text{ and } B) = P(B) \times P(A|B)$
- Example: A box contains 4 green balls and 3 red balls. Two balls are drawn. Find the probability of choosing two red balls.
- A=Red Ball on 1st draw B=Red Ball on 2nd Draw
- $P(A)=3/7$ $P(B|A)=2/6$
- $P(A \text{ and } B) = (3/7)(2/6) = 1/7$

20

Multiplicative Rule – Tree Diagram

$P(A) = \frac{3}{7}$ $P(A') = \frac{4}{7}$
 $P(B|A) = \frac{2}{6}$ $P(B'|A) = \frac{4}{6}$ $P(B|A') = \frac{3}{6}$ $P(B'|A') = \frac{3}{6}$
 $\left(\frac{3}{7}\right)\left(\frac{2}{6}\right) = \frac{1}{7}$ $\left(\frac{3}{7}\right)\left(\frac{4}{6}\right) = \frac{2}{7}$ $\left(\frac{4}{7}\right)\left(\frac{3}{6}\right) = \frac{2}{7}$ $\left(\frac{4}{7}\right)\left(\frac{3}{6}\right) = \frac{2}{7}$

21

Independence

- If A is not dependent on B, then they are **INDEPENDENT** events, and the following statements are true:
 - $P(A|B)=P(A)$
 - $P(B|A)=P(B)$
 - $P(A \text{ and } B) = P(A) \times P(B)$

22

Example

	Accident	No Accident	Total
DUI	70	130	200
Non- DUI	30	770	800
Total	100	900	1000

A: Accident D:DUI Driver

$P(A) = .10$ $P(A|D) = .35 (70/200)$

Therefore A and D are **DEPENDENT** events as $P(A) < P(A|D)$

23

Example

	Accident	No Accident	Total
Domestic Car	60	540	600
Import Car	40	360	400
Total	100	900	1000

A: Accident D:Domestic Car

$P(A) = .10$ $P(A|D) = .10 (60/600)$

Therefore A and D are **INDEPENDENT** events as $P(A) = P(A|D)$

Also $P(A \text{ and } D) = P(A) \times P(D) = (.1)(.6) = .06$

24

Random Sample

- A **random sample** is where each member of the population has an equally likely chance of being chosen, and each member of the sample is **INDEPENDENT** of all other sampled data.

25

Tree Diagram method

- Alternative Method of showing probability
- Example: Flip Three Coins
- Example: A Circuit has three switches. If at least two of the switches function, the Circuit will succeed. Each switch has a 10% failure rate if all are operating, and a 20% failure rate if one switch has already failed. Find the probability the circuit will succeed.

26

Circuit Problem

27

Switching the Conditionality

- Often there are questions where you desire to change the conditionality from one variable to the other variable
- First construct a tree diagram.
- Second, create a Contingency Table using a convenient radix (sample size)
- From the Contingency table it is easy to calculate all conditional probabilities.

28


Example

- 10% of prisoners in a Canadian prison are HIV positive.
- A test will correctly detect HIV 95% of the time, but will incorrectly "detect" HIV in non-infected prisoners 15% of the time (false positive).
- If a randomly selected prisoner tests positive, find the probability the prisoner is HIV+

29

Example

30



Example


	HIV+ A	HIV- A'	Total
Test+ B	950	1350	2300
Test- B'	50	7650	7700
Total	1000	9000	10000

$$P(A | B) = \frac{950}{2300} \approx .413$$

31

Inferential Statistics and Probability a Holistic Approach

Chapter 5 Discrete Random Variables



This Course Material by Maurice Geraghty is licensed under a Creative Commons Attribution-ShareAlike 4.0 International License. Conditions for use are shown here: <https://creativecommons.org/licenses/by-sa/4.0/>

1

Random Variable

- The value of the variable depends on an experiment, observation or measurement.
- The result is not known in advance.
- For the purposes of this class, the variable will be numeric.

2

Random Variables

- Discrete – Data that you Count
 - Defects on an assembly line
 - Reported Sick days
 - RM 7.0 earthquakes on San Andreas Fault
- Continuous – Data that you Measure
 - Temperature
 - Height
 - Time

3

Discrete Random Variable

- List Sample Space
- Assign probabilities $P(x)$ to each event x
- Use "relative frequencies"
- Must follow two rules
 - $P(x) \geq 0$
 - $\sum P(x) = 1$
- $P(x)$ is called a **Probability Distribution Function** or **pdf** for short.

4

Probability Distribution Example

- Students are asked 4 questions and the number of correct answers are determined.
- Assign probabilities to each event.

x	P(x)
0	.1
1	.1
2	.2
3	.4
4	

5

Probability Distribution Example

- Students are asked 4 questions and the number of correct answers are determined.
- Assign probabilities to each event.

x	P(x)
0	.1
1	.1
2	.2
3	.4
4	.2

6

Mean and Variance of Discrete Random Variables

- Population mean μ , is the expected value of x

$$\mu = \sum [(x) P(x)]$$
- Population variance σ^2 , is the expected value of $(x-\mu)^2$

$$\sigma^2 = \sum [(x-\mu)^2 P(x)]$$

7

Example of Mean and Variance

x	P(x)	xP(x)	(x- μ) ² P(x)
0	0.1	0.0	.625
1	0.1	0.1	.225
2	0.2	0.4	.050
3	0.4	1.2	.100
4	0.2	0.8	.450
Total	1.0	2.5=μ	1.450=σ^2

8

Bernoulli Distribution

- Experiment is one trial
- 2 possible outcomes (Success,Failure)
- p=probability of success
- q=probability of failure
- X=number of successes (1 or 0)
- Also known as Indicator Variable

9

Mean and Variance of Bernoulli

x	P(x)	xP(x)	(x- μ) ² P(x)
0	(1-p)	0.0	p ² (1-p)
1	p	p	p(1-p) ²
Total	1.0	p=μ	p(1-p)=σ^2

- $\mu = p$
- $\sigma^2 = p(1-p) = pq$

10

Binomial Distribution

- n identical trials
- Two possible outcomes (success/failure)
- Probability of success in a single trial is p
- Trials are mutually independent
- X is the number of successes
- Note: X is a sum of n independent identically distributed Bernoulli distributions

11

Binomial Distribution

- n independent Bernoulli trials
- Mean and Variance of Binomial Distribution is just sample size times mean and variance of Bernoulli Distribution

$$p(x) = {}_n C_x p^x (1-p)^{n-x}$$

$$\mu = E(X) = np$$

$$\sigma^2 = Var(X) = np(1-p)$$

12

Binomial Examples

- The number of defective parts in a fixed sample.
- The number of adults in a sample who support the war in Iraq.
- The number of correct answers if you guess on a multiple choice test.

13

Binomial Example

- 90% of super duplex globe valves manufactured are good (not defective). A sample of 10 is selected.
- Find the probability of exactly 8 good valves being chosen.
- Find the probability of 9 or more good valves being chosen.
- Find the probability of 8 or less good valves being chosen.

14

Using Technology

<i>X</i>	<i>p(X)</i>	<i>cumulative probability</i>
0	0.00000	0.00000
1	0.00000	0.00000
2	0.00000	0.00000
3	0.00001	0.00001
4	0.00014	0.00015
5	0.00149	0.00163
6	0.01116	0.01280
7	0.05740	0.07019
8	0.19371	0.26390
9	0.38742	0.65132
10	0.34868	1.00000

Use Minitab or Excel to make a table of Binomial Probabilities.

$P(X=8) = .19371$
 $P(X \leq 8) = .26390$
 $P(X \geq 9) = 1 - P(X \leq 8) = .73610$

9.000 expected value
 0.900 variance
 0.949 standard deviation

15

Poisson Distribution

- Occurrences per time period (rate)
- Rate (μ) is constant
- No limit on occurrences over time period

$$P(x) = \frac{e^{-\mu} \mu^x}{x!}$$

$$\mu = \mu$$

$$\sigma = \sqrt{\mu}$$

16

Examples of Poisson

- Text messages in the next hour
- Earthquakes on a fault
- Customers at a restaurant
- Flaws in sheet metal produced
- Lotto winners

Note: A binomial distribution with a large n and small p is approximately Poisson with $\mu \approx np$.

17

Poisson Example

- Earthquakes of Richter magnitude 3 or greater occur on a certain fault at a rate of twice every year.
- Find the probability of at least one earthquake of RM 3 or greater in the next year.

$$P(X > 0) = 1 - P(0)$$

$$= 1 - \frac{e^{-2} 2^0}{0!}$$

$$= 1 - e^{-2} \approx .8647$$

18



Poisson Example (cont)

- Earthquakes of Richter magnitude 3 or greater occur on a certain fault at a rate of twice every year.
- Find the probability of exactly 6 earthquakes of RM 3 or greater in the next 2 years.


$$\mu = 2(2) = 4$$

$$P(X = 6) = \frac{e^{-4}4^6}{6!} \approx .1042$$

19

Inferential Statistics and Probability a Holistic Approach

Chapter 6 Continuous Random Variables



This Course Material by Maurice Geraghty is licensed under a Creative Commons Attribution-ShareAlike 4.0 International License.
Conditions for use are shown here: <https://creativecommons.org/licenses/by-sa/4.0/>

1

Continuous Distributions

- “Uncountable” Number of possibilities
- Probability of a point makes no sense
- Probability is measured over intervals
- Comparable to Relative Frequency Histogram – Find Area under curve.

2

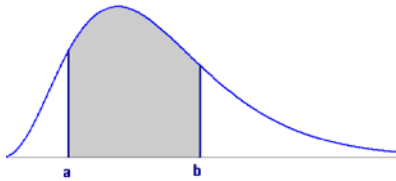
Discrete vs Continuous

<ul style="list-style-type: none"> ■ Countable ■ Discrete Points ■ $p(x)$ is probability distribution function ■ $p(x) \geq 0$ ■ $\sum p(x) = 1$ 	<ul style="list-style-type: none"> ■ Uncountable ■ Continuous Intervals ■ $f(x)$ is probability density function ■ $f(x) \geq 0$ ■ Total Area under curve = 1
--	--

3

Continuous Random Variable

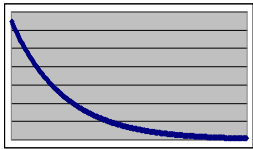
- $f(x)$ is a density function
- $P(X < x)$ is a distribution function.
- $P(a < X < b) =$ area under function between a and b



4

Exponential distribution

- Waiting time
- “Memoryless”
- $f(x) = (1/\mu)e^{-(1/\mu)x}$
- $P(x > a) = e^{-(a/\mu)}$
- $\mu = \mu \quad \sigma^2 = \mu^2$
- $P(x > a + b | x > b) = e^{-(a/\mu)}$



5

Examples of Exponential Distribution

- Time until...
- a circuit will fail
- the next RM 7 Earthquake
- the next customer calls
- An oil refinery accident
- you buy a winning lotto ticket


6

Relationship between Poisson and Exponential Distributions

- If occurrences follow a **Poisson Process** with mean = μ , then the waiting time for the next occurrence has **Exponential** distribution with mean = $1/\mu$.
- Example: If accidents occur at a plant at a constant rate of 3 per month, then the expected waiting time for the next accident is 1/3 month.

7

Exponential Example



The time until a screen is cracked on a smart phone has exponential distribution with $\mu=500$ hours of use.

(a) Find the probability screen will not crack for at least 600 hours.

$$P(x > 600) = e^{-600/500} = e^{-1.2} = .3012$$

(b) Assuming that screen has already lasted 500 hours without cracking, find the chance the display will last an additional 600 hours.

$$P(x > 1100 | x > 500) = P(x > 600) = .3012$$

8

Exponential Example

The time until a screen is cracked on a smart phone has exponential distribution with $\mu=500$ hours of use.

(a) Find the median of the distribution

$$P(x > \text{med}) = e^{-(\text{med})/500} = 0.5$$

$$\text{med} = -500 \ln(.5) = 347$$

p^{th} Percentile = $-\mu \ln(1-p)$

9

Uniform Distribution

- Rectangular distribution
- Example: Random number generator

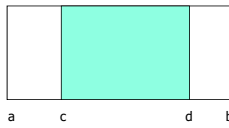
$$f(x) = \frac{1}{b-a} \quad a \leq x \leq b$$

$$\mu = E(X) = \frac{b+a}{2}$$

$$\sigma^2 = \text{Var}(X) = \frac{(b-a)^2}{12}$$

10

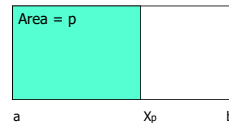
Uniform Distribution - Probability



$$P(c < X < d) = \frac{d-c}{b-a}$$

11

Uniform Distribution - Percentile



Area = p

Formula to find the pth percentile X_p :

$$X_p = a + p(b-a)$$

12

Uniform Example 1

- Find mean, variance, $P(X < 3)$ and 70th percentile for a uniform distribution from 1 to 11.

$$\mu = \frac{1+11}{2} = 6 \quad \sigma^2 = \frac{(11-1)^2}{12} = 8.33$$

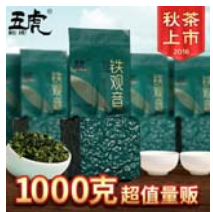
$$P(X < 3) = \frac{3-1}{11-1} = 0.3$$

$$X_{70} = 1 + 0.7(11-1) = 8$$

13

Uniform Example 2

- A tea lover orders 1000 grams of Tie Guan Yin loose leaf when his supply gets to 50 grams.
- The amount of tea currently in stock follows a uniform random variable.
- Determine this model
- Find the mean and variance
- Find the probability of at least 700 grams in stock.
- Find the 80th percentile



14

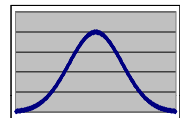
Uniform Example 3

- A bus arrives at a stop every 20 minutes.
 - Find the probability of waiting more than 15 minutes for the bus after arriving randomly at the bus stop.
 - If you have already waited 5 minutes, find the probability of waiting an additional 10 minutes or more. (Hint: recalculate parameters a and b)

15

Normal Distribution

- The normal curve is *bell-shaped*
- The mean, median, and mode of the distribution are equal and located at the peak.
- The normal distribution is *symmetrical* about its mean. Half the area under the curve is above the peak, and the other half is below it.
- The normal probability distribution is *asymptotic* - the curve gets closer and closer to the x-axis but never actually touches it.



$$f(x) = \frac{e^{-\frac{1}{2\sigma^2}(x-\mu)^2}}{\sigma\sqrt{2\pi}}$$

16

The Standard Normal Probability Distribution

- A normal distribution with a mean of 0 and a standard deviation of 1 is called the **standard normal distribution**.
- Z value:** The distance between a selected value, designated x , and the population mean μ , divided by the population standard deviation, σ

$$Z = \frac{X - \mu}{\sigma}$$

17

Areas Under the Normal Curve – Empirical Rule

- About 68 percent of the area under the normal curve is within one standard deviation of the mean. $\mu \pm 1\sigma$
- About 95 percent is within two standard deviations of the mean $\mu \pm 2\sigma$
- 99.7 percent is within three standard deviations of the mean. $\mu \pm 3\sigma$

18

EXAMPLE

- The daily water usage per person in a town is normally distributed with a mean of 20 gallons and a standard deviation of 5 gallons.
- About 68% of the daily water usage per person in New Providence lies between what two values?
- $\mu \pm 1\sigma = 20 \pm 1(5)$. That is, about 68% of the daily water usage will lie between 15 and 25 gallons.

19

Normal Distribution – probability problem procedure

- Given: Interval in terms of X
- Convert to Z by $Z = \frac{X - \mu}{\sigma}$
- Look up probability in table.

20

EXAMPLE

- The daily water usage per person in a town is normally distributed with a mean of 20 gallons and a standard deviation of 5 gallons.
- What is the probability that a person from the town selected at random will use less than 18 gallons per day?
- The associated Z value is $Z = (18 - 20) / 5 = -0.40$.
- Thus, $P(X < 18) = P(Z < -0.40) = .3446$

21

EXAMPLE *continued*

- The daily water usage per person in a town is normally distributed with a mean of 20 gallons and a standard deviation of 5 gallons.
- What proportion of the people uses between 18 and 24 gallons?
- The Z value associated with $x=18$ is $Z = -0.40$ and with $X=24$, $Z = (24 - 20) / 5 = 0.80$.
- Thus, $P(18 < X < 24) = P(-0.40 < Z < 0.80) = .7881 - .3446 = .4435$

22

EXAMPLE *continued*

- The daily water usage per person in a town is normally distributed with a mean of 20 gallons and a standard deviation of 5 gallons.
- What percentage of the population uses more than 26.2 gallons?
- The Z value associated with $X=26.2$, $Z = (26.2 - 20) / 5 = 1.24$.
- Thus $P(X > 26.2) = P(Z > 1.24) = 1 - .8925 = .1075$


23

Normal Distribution – percentile problem procedure

- Given: probability or percentile desired.
- Look up Z value in table that corresponds to probability.
- Convert to X by the formula:

$$X = \mu + Z\sigma$$


24



EXAMPLE

- The daily water usage per person in a town is normally distributed with a mean of 20 gallons and a standard deviation of 5 gallons. A special tax is going to be charged on the top 5% of water users.
- Find the value of daily water usage that generates the special tax
- The Z value associated with 95th percentile = 1.645
- $X = 20 + 5(1.645) = 28.2$ gallons per day


25



EXAMPLE

- Professor Kury has determined that the final averages in his statistics course is normally distributed with a mean of 77.1 and a standard deviation of 11.2.
- He decides to assign his grades for his current course such that the top 15% of the students receive an A.
- What is the lowest average a student can receive to earn an A?
- The top 15% would be the finding the 85th percentile. Find k such that $P(X < k) = .85$.
- The corresponding Z value is 1.04. Thus we have $X = 77.1 + (1.04)(11.2)$, or **$X = 88.75$**

26




EXAMPLE

- The amount of tip the servers in an exclusive restaurant receive per shift is normally distributed with a mean of \$80 and a standard deviation of \$10.
- Shelli feels she has provided poor service if her total tip for the shift is less than \$65.
- What percentage of the time will she feel like she provided poor service?
- Let y be the amount of tip. The Z value associated with $X = 65$ is $Z = (65 - 80) / 10 = -1.5$.
- Thus $P(X < 65) = P(Z < -1.5) = .0668$.

27

Inferential Statistics and Probability a Holistic Approach

Chapter 7 Central Limit Theorem

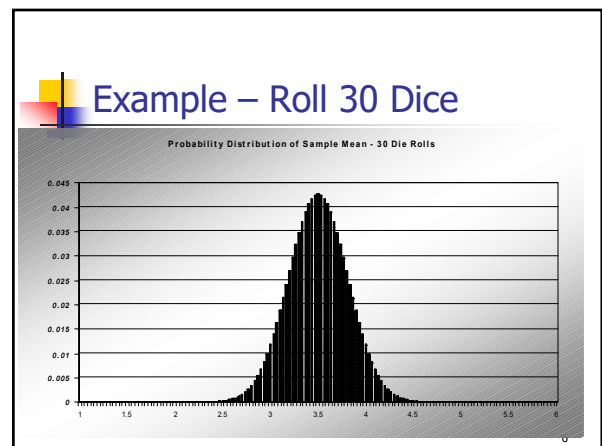
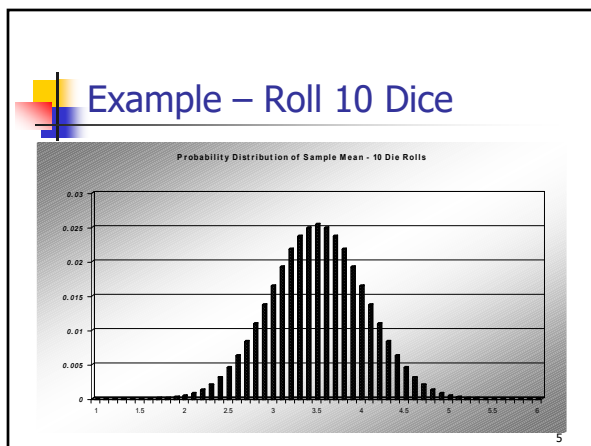
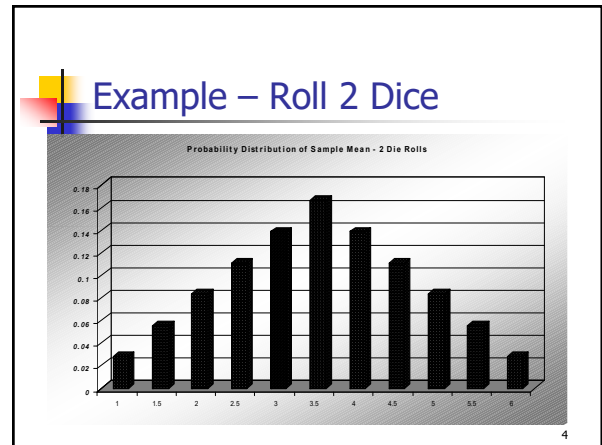
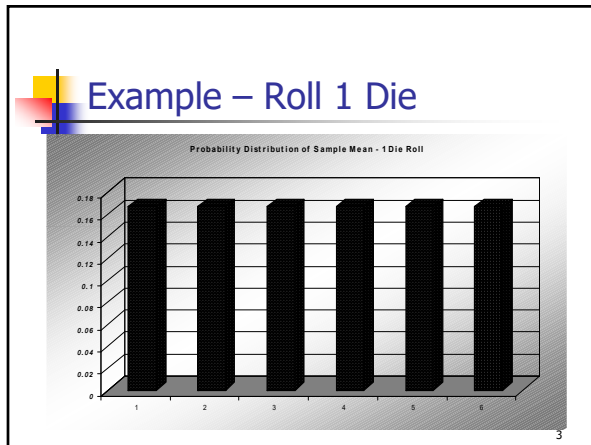

This Course Material by Maurice Geraghty is licensed under a Creative Commons Attribution-ShareAlike 4.0 International License. Conditions for use are shown here: <https://creativecommons.org/licenses/by-sa/4.0/>

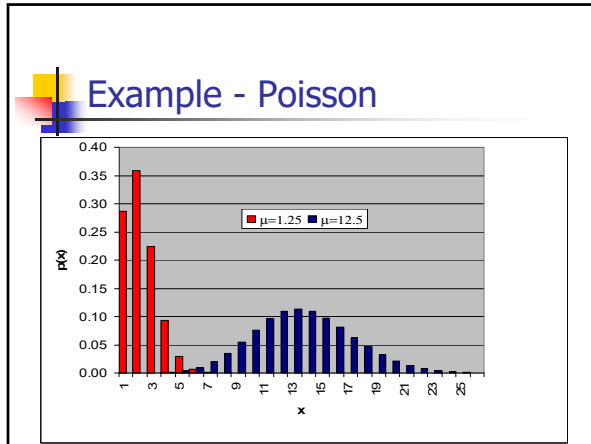
1

Distribution of Sample Mean

- Random Sample: $X_1, X_2, X_3, \dots, X_n$
 - Each X_i is a Random Variable from the same population
 - All X_i 's are Mutually Independent
- \bar{X} is a function of Random Variables, so \bar{X} is itself Random Variable.
- In other words, the Sample Mean can change if the values of the Random Sample change.
- What is the Probability Distribution of \bar{X} ?

2





Central Limit Theorem – Part 1

- IF a Random Sample of **any size** is taken from a population with a **Normal Distribution** with mean = μ and standard deviation = σ

- THEN the distribution of the sample mean has a Normal Distribution with:

$$\mu_{\bar{X}} = \mu \quad \sigma_{\bar{X}} = \frac{\sigma}{\sqrt{n}}$$

Central Limit Theorem – Part 2

- IF a random sample of **sufficiently large size** is taken from a population with **any Distribution** with mean = μ and standard deviation = σ

- THEN the distribution of the sample mean has approximately a Normal Distribution with:

$$\mu_{\bar{X}} = \mu \quad \sigma_{\bar{X}} = \frac{\sigma}{\sqrt{n}}$$

Central Limit Theorem

3 important results for the distribution of \bar{X}

- Mean Stays the same

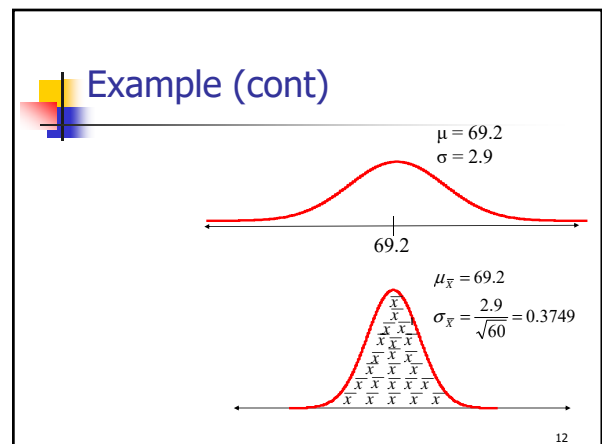
$$\mu_{\bar{X}} = \mu$$
- Standard Deviation Gets Smaller


$$\sigma_{\bar{X}} = \frac{\sigma}{\sqrt{n}}$$
- If n is sufficiently large, \bar{X} has a Normal Distribution

Example

The mean height of American men (ages 20-29) is $\mu = 69.2$ inches. If a random sample of 60 men in this age group is selected, what is the probability the mean height for the sample is greater than 70 inches? Assume $\sigma = 2.9$.

$$P(\bar{X} > 70) = P\left(Z > \frac{(70 - 69.2)}{2.9/\sqrt{60}}\right)$$

$$= P(Z > 2.14) = 0.0162$$




Example – Central Limit Theorem


The waiting time until receiving a text message follows an exponential distribution with an expected waiting time of 1.5 minutes. Find the probability that the mean waiting time for the 50 text messages exceeds 1.6 minutes.

$\mu = 1.5 \quad \sigma = 1.5 \quad n = 50$

Use Normal Distribution ($n > 30$)

$$P(\bar{X} > 1.6) = P\left(Z > \frac{(1.6 - 1.5)}{1.5/\sqrt{50}}\right) = P(Z > 0.47) = 0.3192$$


13



Central Limit Theorem Sample Proportion

- If X is a Random Variable from a Binomial Distribution with parameters n and p, an $np > 10$ and $n(1-p) > 10$, then the following is true for the Sample Proportion, \hat{p} :
 - $\mu_{\hat{p}} = p \quad \sigma_{\hat{p}} = \sqrt{\frac{p(1-p)}{n}}$
 - The Distribution of \hat{p} is approximately Normal.

14



Example

- 45% of all community college students in California receive fee waivers.
- Suppose you randomly sample 1000 community college students to determine the proportion of students with fee waivers in the sample.
- 483 of the sampled students are receiving fee waivers.
- Determine \hat{p} . Is the result unusual?

15

Inferential Statistics and Probability a Holistic Approach

Chapter 8 Point Estimation and Confidence Intervals

This Course Material by Maurice Geraghty is licensed under a Creative Commons Attribution-ShareAlike 4.0 International License. Conditions for use are shown here: <https://creativecommons.org/licenses/by-sa/4.0/>

1

Inference Process

Population

1. From a **Population**, take a sample.

2

Inference Process

Population

Sample

1. From a **Population**, take a sample.
2. Analyze the **Sample**.

3

Inference Process

Population

Sample

1. From a **Population**, take a sample.
2. Analyze the **Sample**.
3. Make an **Inference** about the Population based on the Sample.

4

Inference Process

Population

Sample

1. From a **Population**, take a sample.
2. Analyze the **Sample**.
3. Make an **Inference** about the Population based on the Sample.
4. Measure the **Reliability** of the Inference.

5

Inferential Statistics

- Population Parameters
 - Mean = μ
 - Proportion = p
 - Standard Deviation = σ
- Sample Statistics
 - Mean = \bar{x}
 - Proportion = \hat{p}
 - Standard Deviation = s

6

Inferential Statistics

- Estimation
 - Using sample data to estimate population parameters.
 - Example: Public opinion polls
- Hypothesis Testing
 - Using sample data to make decisions or claims about population
 - Example: A drug effectively treats a disease

7

Estimation of μ

\bar{X} is an unbiased **point estimator** of μ

Example: The number of defective items produced by a machine was recorded for five randomly selected hours during a 40-hour work week. The observed number of defectives were 12, 4, 7, 14, and 10. So the sample mean is 9.4.

Thus a **point estimate** for μ , the hourly mean number of defectives, is 9.4.

8

Confidence Intervals

- An **Interval Estimate** states the range within which a population parameter "probably" lies.
- The interval within which a population parameter is expected to occur is called a **Confidence Interval**.
- The distance from the center of the confidence interval to the endpoint is called the **"Margin of Error"**
- The three confidence intervals that are used extensively are the 90%, 95% and 99%.

9

Confidence Intervals

- A 95% confidence interval means that about 95% of the similarly constructed intervals will contain the parameter being estimated, or 95% of the sample means for a specified sample size will lie within 1.96 standard deviations of the hypothesized population mean.
- For the 99% confidence interval, 99% of the sample means for a specified sample size will lie within 2.58 standard deviations of the hypothesized population mean.
- For the 90% confidence interval, 90% of the sample means for a specified sample size will lie within 1.645 standard deviations of the hypothesized population mean.

10

90%, 95% and 99% Confidence Intervals for μ

- The 90%, 95% and 99% confidence intervals for μ are constructed as follows when $n \geq 30$
- 90% CI for the population mean is given by

$$\bar{X} \pm 1.645 \frac{\sigma}{\sqrt{n}}$$
- 95% CI for the population mean is given by

$$\bar{X} \pm 1.96 \frac{\sigma}{\sqrt{n}}$$
- 99% CI for the population mean is given by

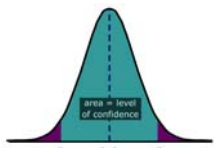
$$\bar{X} \pm 2.58 \frac{\sigma}{\sqrt{n}}$$

11

Constructing General Confidence Intervals for μ

- In general, a confidence interval for the mean is computed by:

$$\bar{X} \pm Z \frac{\sigma}{\sqrt{n}}$$




- This can also be thought of as:

Point Estimator \pm Margin of Error

12

The nature of Confidence Intervals

- The Population mean μ is fixed.
- The confidence interval is centered around the sample mean which is a Random Variable.
- So the Confidence Interval (Random Variable) is like a target trying hit a fixed dart (μ).



13

EXAMPLE

- The Dean wants to estimate the mean number of hours worked per week by students. A sample of 49 students showed a mean of 24 hours with a standard deviation of 4 hours.
- The point estimate is 24 hours (sample mean).
- What is the 95% confidence interval for the average number of hours worked per week by the students?

14

EXAMPLE *continued*

- Using the 95% CI for the population mean, we have
 $24 \pm 1.96(4 / 7) = 22.88 \text{ to } 25.12$
- The endpoints of the confidence interval are the confidence limits. The lower confidence limit is 22.88 and the upper confidence limit is 25.12

15

EXAMPLE *continued*

- Using the 99% CI for the population mean, we have
 $24 \pm 2.58(4 / 7) = 22.53 \text{ to } 25.47$
- Compare to the 95% confidence interval. A higher level of confidence means the confidence interval must be wider.

16

Selecting a Sample Size

- There are 3 factors that determine the size of a sample, none of which has any direct relationship to the size of the population. They are:
 - The degree of confidence selected.
 - The maximum allowable error.
 - The variation of the population.

17

Sample Size for the Mean

- A convenient computational formula for determining n is:

$$n = \left(\frac{Z\sigma}{E} \right)^2$$
- where E is the allowable error (margin of error), Z is the z score associated with the degree of confidence selected, and σ is the sample deviation of the pilot survey.
- σ can be estimated by past data, target sample or range of data.

18

EXAMPLE

- A consumer group would like to estimate the mean monthly electric bill for a single family house in July. Based on similar studies the standard deviation is estimated to be \$20.00. A 99% level of confidence is desired, with an accuracy of \$5.00. How large a sample is required?

$$n = [(2.58)(20) / 5]^2 = 106.5024 \approx 107$$

19

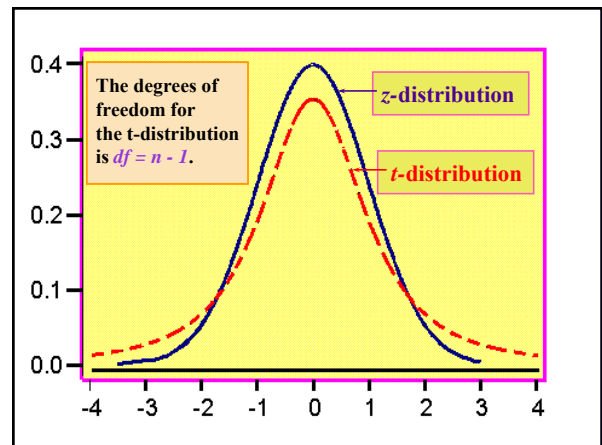
Normal Family of Distributions: Z, t, χ^2 , F

20

Characteristics of Student's *t*-Distribution

- The *t*-distribution has the following properties:
 - It is continuous, bell-shaped, and symmetrical about zero like the z-distribution.
 - There is a **family** of *t*-distributions sharing a mean of zero but having different standard deviations based on **degrees of freedom**.
 - The *t*-distribution is more spread out and flatter at the center than the z-distribution, but approaches the z-distribution as the sample size gets larger.

21



Confidence Interval for μ (σ unknown)

Formula to find a confidence interval using the *t*-distribution for the appropriate level of confidence:

$$\bar{X} \pm t \left(\frac{s}{\sqrt{n}} \right) \quad df = n - 1$$

23

Example – Confidence Interval

- In a random sample of 13 American adults, the mean waste recycled per person per day was 5.3 pounds and the standard deviation was 2.0 pounds.
- Assume the variable is normally distributed and construct a 95% confidence interval for μ .

24

Example- Confidence Interval

level of confidence = 95%
 $df=13-1=12$
 $t=2.18$

$$5.3 \pm 2.18 \frac{2.0}{\sqrt{13}}$$

$$5.3 \pm 1.2 = (4.1, 6.5)$$

25


Confidence Intervals, Population Proportions

- Point estimate for proportion of successes in population is: $\hat{p} = \frac{X}{n}$
- X is the number of successes in a sample of size n.
- Standard deviation of \hat{p} is $\sqrt{\frac{p(1-p)}{n}}$
- Confidence Interval for p:

$$\hat{p} \pm Z \cdot \sqrt{\frac{p(1-p)}{n}} \approx \hat{p} \pm Z \cdot \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$$

26

Population Proportion Example



- In a May 2006 AP/ISPOS Poll, 1000 adults were asked if "Over the next six months, do you expect that increases in the price of gasoline will cause financial hardship for you or your family, or not?"
- 700 of those sampled responded yes!
- Find the **sample proportion** and **margin of error** for this poll. (This means find a 95% confidence interval.)

27

Population Proportion Example

- Sample proportion

$$\hat{p} = \frac{700}{1000} = .70 = 70\%$$

- Margin of Error

$$MOE = 1.96 \sqrt{\frac{.70(1-.70)}{1000}} = .028 = 2.8\%$$

28

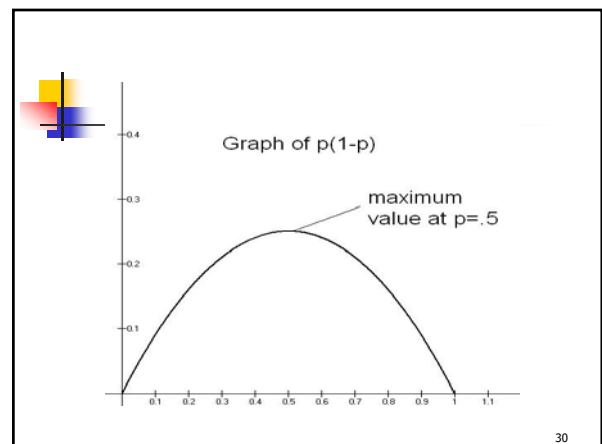
Sample Size for the Proportion

- A convenient computational formula for determining n is:

$$n = (p(1-p)) \left(\frac{Z}{E} \right)^2$$

- where E is the allowable margin of error, Z is the z-score associated with the degree of confidence selected, and p is the population proportion.
- If p is completely unknown, p can be set equal to 1/2 which maximizes the value of (p)(1-p) and guarantees the confidence interval will fall within the margin of error.

29



Example

- In polling, determine the minimum sample size needed to have a margin of error of 3% when p is unknown.

$$n = (.5)(1-.5)\left(\frac{1.96}{.03}\right)^2 = 1068$$

31

Example

- In polling, determine the minimum sample size needed to have a margin of error of 3% when p is known to be close to 1/4.

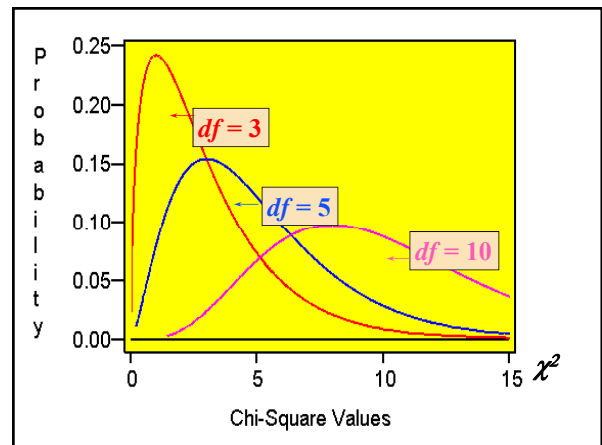
$$n = (.25)(1-.25)\left(\frac{1.96}{.03}\right)^2 = 801$$

32

Characteristics of the Chi-Square Distribution

- The major characteristics of the chi-square distribution are:
 - It is positively skewed
 - It is non-negative
 - It is based on degrees of freedom
 - When the degrees of freedom change, a new distribution is created

33



Inference about Population Variance and Standard Deviation


- s^2 is an unbiased point estimator for σ^2
- s is a point estimator for σ
- Interval estimates and hypothesis testing for both σ^2 and σ require a new distribution – the χ^2 (Chi-square)

35

Distribution of s^2

- $\frac{(n-1)s^2}{\sigma^2}$ has a chi-square distribution
- $n-1$ is degrees of freedom
- s^2 is sample variance
- σ^2 is population variance

36




Confidence interval for σ^2

- Confidence is **NOT** symmetric since chi-square distribution is not symmetric. You must find separate left and right bounds.
- We can construct a confidence interval for σ^2

$$\left(\frac{(n-1)s^2}{\chi_L^2}, \frac{(n-1)s^2}{\chi_R^2} \right)$$

- Take square root of both endpoints to get confidence interval for σ , the population standard deviation.


37



Example

- In performance measurement of investments, standard deviation is a measure of volatility or risk.
- Twenty monthly returns from a mutual fund show an average monthly return of 1% and a sample standard deviation of 5%
- Find a 95% confidence interval for the monthly standard deviation of the mutual fund.

38



Example (cont)


- df = n-1 = 19
- 95% CI for σ

$$\left(\sqrt{\frac{(19)5^2}{32.8523}}, \sqrt{\frac{(19)5^2}{8.90655}} \right) = (3.8, 7.3)$$

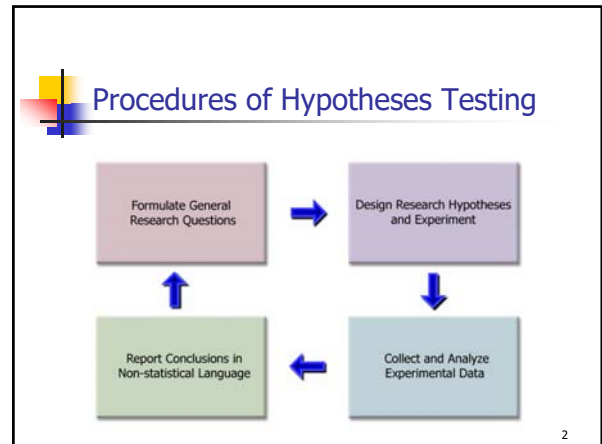
39

Inferential Statistics and Probability a Holistic Approach

Chapter 9 One Population Hypothesis Testing


This Course Material by Maurice Geraghty is licensed under a Creative Commons
Attribution-ShareAlike 4.0 International License.
Conditions for use are shown here: <https://creativecommons.org/licenses/by-sa/4.0/>

1



Hypotheses Testing – Procedure 1

Formulate General Research Questions

3

General Research Question

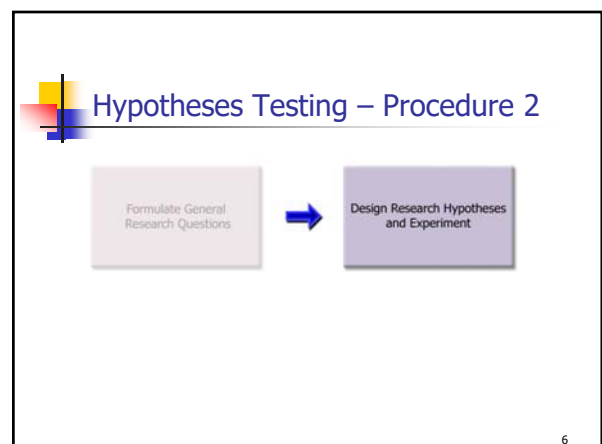
- Decide on a topic or phenomena that you want to research.
- Formulate general research questions based on the topic.
- Example:
 - Topic: Health Care Reform
 - Some General Questions:
 - Would a Single Payer Plan be less expensive than Private Insurance?
 - Do HMOs provide the same quality care as PPOs?
 - Would the public support mandated health coverage?

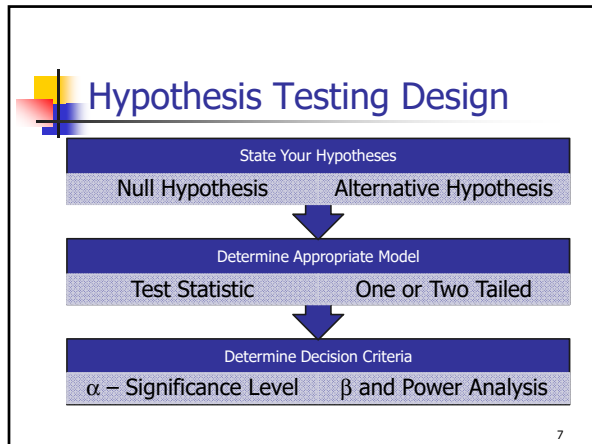
4

EXAMPLE – General Question

- A food company has a policy that the stated contents of a product match the actual results.
- A General Question might be “Does the stated net weight of a food product match (on average) the actual weight?”
- The quality control statistician could then decide to test various food products for accuracy.

5





What is a Hypothesis?

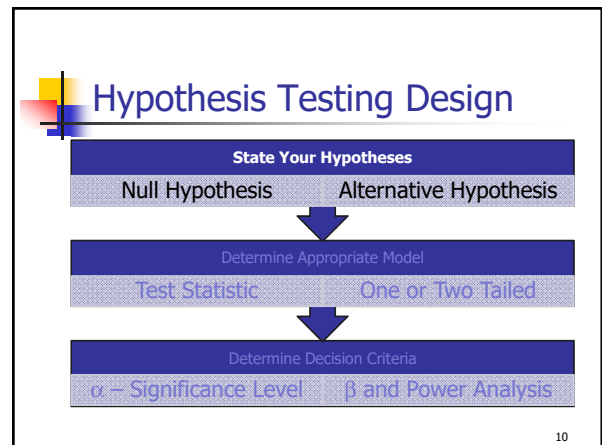
- **Hypothesis:** A statement about the value of a population parameter developed for the purpose of testing.
- Examples of hypotheses made about a population parameter are:
 - The mean monthly income for programmers is \$9,000.
 - At least twenty percent of all juvenile offenders are caught and sentenced to prison.
 - The standard deviation for an investment portfolio is no more than 10 percent per month.

8

What is Hypothesis Testing?

- **Hypothesis testing:** A procedure, based on sample evidence and probability theory, used to determine whether the hypothesis is a reasonable statement and should not be rejected, or is unreasonable and should be rejected.

9



Definitions

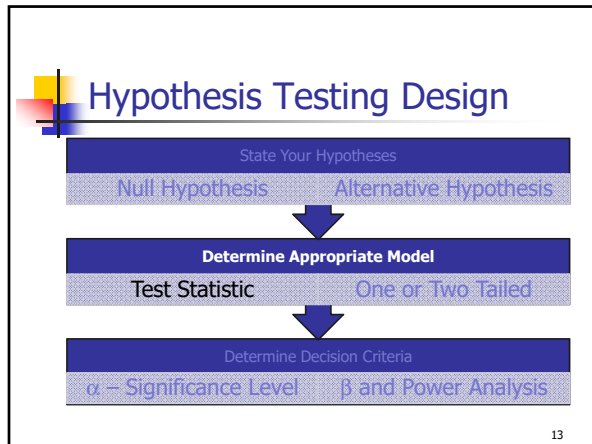
- **Null Hypothesis H_0 :** A statement about the value of a population parameter that is assumed to be true for the purpose of testing.
- **Alternative Hypothesis H_a :** A statement about the value of a population parameter that is assumed to be true if the Null Hypothesis is rejected during testing.

11

EXAMPLE – Stating Hypotheses

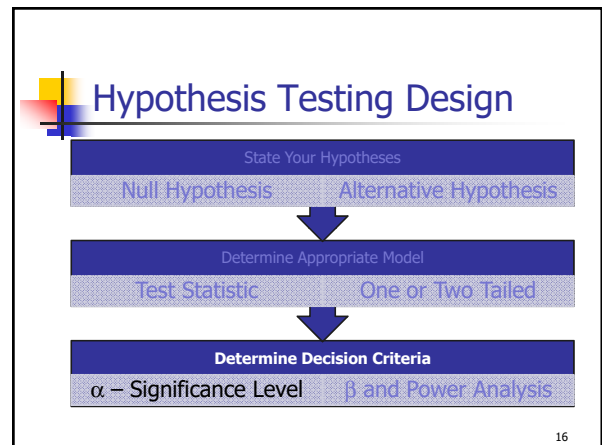
- A food company has a policy that the stated contents of a product match the actual results.
- The quality control statistician decides to test the claim that a 16 ounce bottle of Soy sauce contains on average 16 ounces.
 - H_0 : The mean amount of Soy Sauce is 16 ounces
 - H_a : The mean amount of Soy Sauce is not 16 ounces.
- $H_0: \mu=16$ $H_a: \mu \neq 16$

12



- ### Definitions
- **Statistical Model**: A mathematical model that describes the behavior of the data being tested.
 - **Normal Family** = the Standard Normal Distribution (Z) and functions of independent Standard Normal Distributions (eg: t, χ^2 , F).
 - Most Statistical Models will be from the Normal Family due to the Central Limit Theorem.
 - **Model Assumptions**: Criteria which must be satisfied to appropriately use a chosen Statistical Model.
 - **Test statistic**: A value, determined from sample information, used to determine whether or not to reject the null hypothesis.
- 14

- ### EXAMPLE – Choosing Model
- The quality control statistician decides to test the claim that a 16 ounce bottle of Soy sauce contains on average 16 ounces. We will assume the population standard is known
 - $H_0: \mu=16$ $H_a: \mu \neq 16$
 - Model: One sample test Z test of mean
 - Test Statistic:
$$Z = \frac{\bar{X} - \mu_0}{\frac{\sigma}{\sqrt{n}}}$$
- 15



- ### Definitions
- **Level of Significance**: The probability of rejecting the null hypothesis when it is actually true. (signified by α)
 - **Type I Error**: Rejecting the null hypothesis when it is actually true.
 - **Type II Error**: Failing to reject the null hypothesis when it is actually false.
- 17

Outcomes of Hypothesis Testing

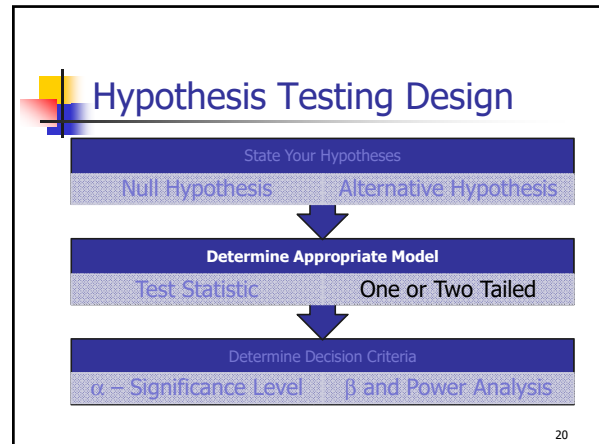
	Fail to Reject H_0	Reject H_0
H_0 is true	Correct Decision	Type I error
H_0 is False	Type II error	Correct Decision

18

EXAMPLE – Type I and Type II Errors

- Ho: The mean amount of Soy Sauce is 16 ounces
- Ha: The mean amount of Soy Sauce is not 16 ounces.
- Type I Error: The researcher **supports** the claim that the mean amount of soy sauce is not 16 ounces when the actual mean is 16 ounces. The company needlessly "fixes" a machine that is operating properly.
- Type II Error: The researcher **fails to support** the claim that the mean amount of soy sauce is not 16 ounces when the actual mean is not 16 ounces. The company fails to fix a machine that is not operating properly.

19



Definitions

- Critical value(s):** The dividing point(s) between the region where the null hypothesis is rejected and the region where it is not rejected. The critical value determines the decision rule.
- Rejection Region:** Region(s) of the Statistical Model which contain the values of the Test Statistic where the Null Hypothesis will be rejected. The area of the Rejection Region = α .

21

One-Tailed Tests of Significance

- A test is one-tailed when the alternate hypothesis, H_a , states a direction, such as:
 - H_0 : The mean income of females is less than or equal to the mean income of males.
 - H_a : The mean income of females is greater than males.
- Equality is part of H_0
- H_a determines which tail to test
 - $H_a: \mu > \mu_0$ means test upper tail.
 - $H_a: \mu < \mu_0$ means test lower tail.

22

One-tailed test

$$H_0 : \mu \leq \mu_0$$

$$H_a : \mu > \mu_0$$

$$\alpha = .05$$

$$Z = \frac{\bar{X} - \mu_0}{\sigma / \sqrt{n}}$$

23

Two-Tailed Tests of Significance

- A test is two-tailed when no direction is specified in the alternate hypothesis H_a , such as:
 - H_0 : The mean income of females is equal to the mean income of males.
 - H_a : The mean income of females is not equal to the mean income of the males.
- Equality is part of H_0
- H_a determines which tail to test
 - $H_a: \mu \neq \mu_0$ means test both tails.

24

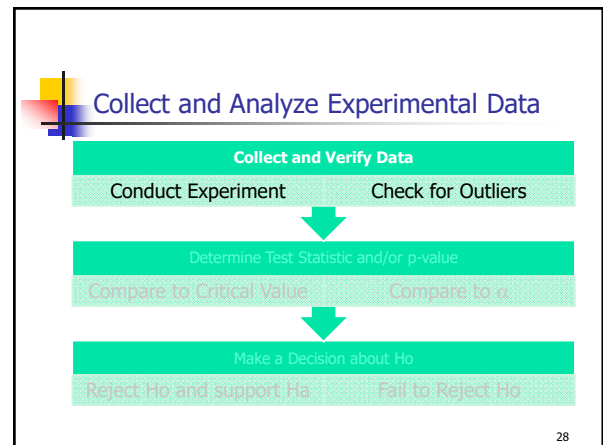
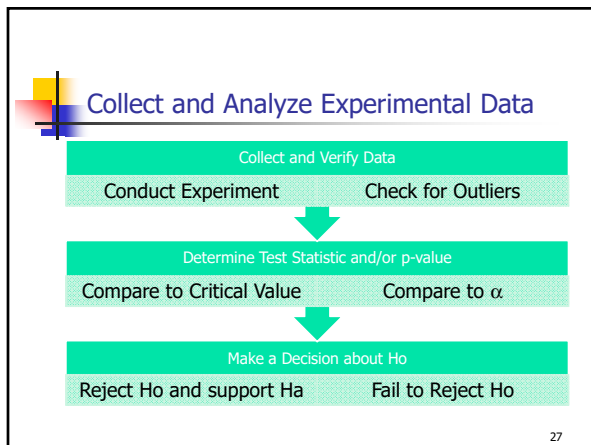
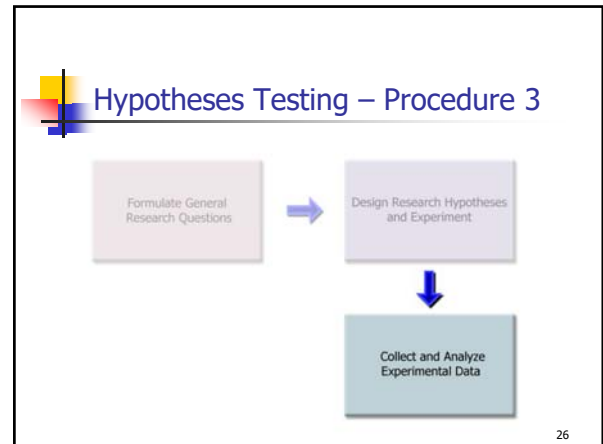
Two-tailed test

$H_0 : \mu = \mu_0$
 $H_a : \mu \neq \mu_0$
 $\alpha = .05 \quad \alpha/2 = .025$
 $Z = \frac{\bar{X} - \mu_0}{\sigma / \sqrt{n}}$

Two-Tailed Test
Regions of Rejection

Critical values: -1.96, 1.96

25



Outliers

- An outlier is data point that is far removed from the other entries in the data set.
- Outliers could be
 - Mistakes made in recording data
 - Data that don't belong in population
 - True rare events

29

Outliers have a dramatic effect on some statistics

- Example quarterly home sales for 10 realtors:

	2	2	3	4	5	5	6	6	7	50
	with outlier					without outlier				
Mean	9.00					4.44				
Median	5.00					5.00				
Std Dev	14.51					1.81				
IQR	3.00					3.50				

30

Using Box Plot to find outliers

- The "box" is the region between the 1st and 3rd quartiles.
- Possible outliers are more than 1.5 IQR's from the box (inner fence)
- Probable outliers are more than 3 IQR's from the box (outer fence)
- In the box plot below, the dotted lines represent the "fences" that are 1.5 and 3 IQR's from the box. See how the data point 50 is well outside the outer fence and therefore an almost certain outlier.

31

Using Z-score to detect outliers

- Calculate the mean and standard deviation without the suspected outlier.
- Calculate the Z-score of the suspected outlier.
- If the Z-score is more than 3 or less than -3, that data point is a probable outlier.

$$Z = \frac{50 - 4.4}{1.81} = 25.2$$

32

Outliers – what to do

- Remove or not remove, there is no clear answer.
- For some populations, outliers don't dramatically change the overall statistical analysis. Example: the tallest person in the world will not dramatically change the mean height of 10000 people.
- However, for some populations, a single outlier will have a dramatic effect on statistical analysis (called "**Black Swan**" by Nicholas Taleb) and inferential statistics may be invalid in analyzing these populations. Example: the richest person in the world will dramatically change the mean wealth of 10000 people.

33

Example – Analyze Data

- In the Soy Sauce Example, a 36 bottles were measured, volume is in fluid ounces

14.51	15.16	15.28	15.33	15.36	15.42
15.43	15.45	15.49	15.59	15.60	15.61
15.62	15.63	15.71	15.81	15.87	16.00
16.01	16.02	16.05	16.06	16.06	16.09
16.09	16.11	16.16	16.16	16.27	16.31
16.35	16.36	16.45	16.72	16.75	16.79

34

Example – Analyze Data

- Although 14.51 might be a possible outlier and the data seems negatively skewed, the Central Limit Theorem assures that the sample mean will have a normal distribution

35

Collect and Analyze Experimental Data

```

    graph TD
      A[Collect and Verify Data] --> B[Conduct Experiment]
      A --> C[Check for Outliers]
      B --> D[Determine Test Statistic and/or p-value]
      C --> D
      D --> E[Compare to Critical Value]
      D --> F[Compare to alpha]
      E --> G[Make a Decision about Ho]
      F --> G
      G --> H[Reject Ho and support Ha]
      G --> I[Fail to Reject Ho]
    
```

36

The logic of Hypothesis Testing

- This is a "Proof" by contradiction.
 - We assume H_0 is true before observing data and design H_a to be the complement of H_0 .
 - Observe the data (evidence). How unusual are these data under H_0 ?
 - If the data are too unusual, we have "proven" H_0 is false: Reject H_0 and go with H_a (Strong Statement)
 - If the data are not too unusual, we fail to reject H_0 . This "proves" nothing and we say data are inconclusive. (Weak Statement)
 - We can never "prove" H_0 , only "disprove" it.
 - "Prove" in statistics means support with $(1-\alpha)100\%$ certainty. (example: if $\alpha=.05$, then we are 95% certain.

37

Test Statistic

- Test Statistic:** A value calculated from the Data under the appropriate Statistical Model from the Data that can be compared to the Critical Value of the Hypothesis test
- If the Test Statistic fall in the Rejection Region, H_0 is rejected.
- The Test Statistic will also be used to calculate the p-value as will be defined next.

38

Example - Testing for the Population Mean Large Sample, Population Standard Deviation Known

- When testing for the population mean from a large sample and the population standard deviation is known, the test statistic is given by:

$$Z = \frac{\bar{X} - \mu}{\sigma / \sqrt{n}}$$

39

p-Value in Hypothesis Testing

- p-Value:** the probability, assuming that the null hypothesis is true, of getting a value of the test statistic at least as extreme as the computed value for the test.
- If the p-value is smaller than the significance level, H_0 is rejected.
- If the p-value is larger than the significance level, H_0 is not rejected.

40

Comparing p-value to α

- Both **p-value** and **α** are probabilities.
- The **p-value** is determined by the **data**, and is the probability of getting results as extreme as the data assuming H_0 is true. Small values make one more likely to reject H_0 .
- α** is determined by **design**, and is the maximum probability the experimenter is willing to accept of rejecting a true H_0 .
- Reject H_0 if $p\text{-value} < \alpha$ for ALL MODELS.

41

Graphic where decision is to Reject H_0

- $H_0: \mu = 10$
 $H_a: \mu > 10$
- Design: Critical Value is determined by significance level α .
- Data Analysis: p-value is determined by Test Statistic
- Test Statistic falls in Rejection Region.
- p-value (blue) $< \alpha$ (purple)
- Reject H_0 .
- Strong statement: Data supports Alternative Hypothesis.

42

Graphic where decision is Fail to Reject Ho

- Ho: $\mu = 10$
Ha: $\mu > 10$
- Design: Critical Value is determined by significance level α .
- Data Analysis: p-value is determined by Test Statistic
- Test Statistic falls in Non-rejection Region.
- p-value (blue) $> \alpha$ (purple)
- Fail to Reject Ho.
- Weak statement: Data is inconclusive and does not support Alternative Hypothesis.

43

EXAMPLE – General Question

- A food company has a policy that the stated contents of a product match the actual results.
- A General Question might be "Does the stated net weight of a food product match the actual weight?"
- The quality control statistician decides to test the 16 ounce bottle of Soy Sauce.

44

EXAMPLE – Design Experiment

- A sample of $n=36$ bottles will be selected hourly and the contents weighed.
- Ho: $\mu=16$ Ha: $\mu \neq 16$
- The Statistical Model will be the one population test of mean using the Z Test Statistic.
- This model will be appropriate since the sample size insures the sample mean will have a Normal Distribution (Central Limit Theorem)
- We will choose a significance level of $\alpha = 5\%$

45

EXAMPLE – Conduct Experiment

- Last hour a sample of 36 bottles had a mean weight of 15.88 ounces.
- From past data, assume the population standard deviation is 0.5 ounces.
- Compute the Test Statistic

$$Z = [15.88 - 16] / [0.5 / \sqrt{36}] = -1.44$$
- For a two tailed test, The Critical Values are at $Z = \pm 1.96$

46

Decision – Critical Value Method

- This two-tailed test has two Critical Value and Two Rejection Regions
- The significance level (α) must be divided by 2 so that the sum of both purple areas is 0.05
- The Test Statistic does not fall in the Rejection Regions.
- Decision is **Fail to Reject Ho.**

47

Computation of the p-Value

- One-Tailed Test: p-Value = $P\{z \geq \text{absolute value of the computed test statistic value}\}$
- Two-Tailed Test: p-Value = $2P\{z \geq \text{absolute value of the computed test statistic value}\}$
- Example: $Z = 1.44$, and since it was a two-tailed test, then p-Value = $2P\{z \geq 1.44\} = 0.0749 = .1498$. Since $.1498 > .05$, do not reject H_0 .

48

Decision – p-value Method

- The p-value for a two-tailed test must include all values (positive and negative) more extreme than the Test Statistic.
- p-value = .1498 which exceeds $\alpha = .05$
- Decision is **Fail to Reject Ho.**

49

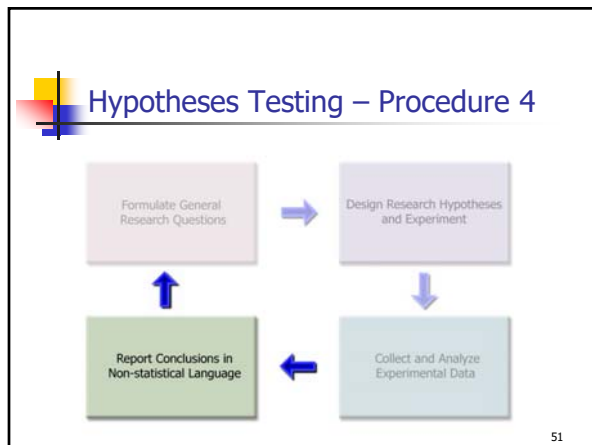
p-value form Minitab (shown as p)

One-Sample Z: weight

Test of $\mu = 16$ vs $\neq 16$
 The assumed standard deviation = 0.5

Variable	N	Mean	StDev	SE Mean	Z	P
weight	36	15.8800	0.4877	0.0833	-1.44	0.150

50



Converting Decision to Conclusion

- Conclusion if Decision is Reject Ho:
 - <Ha in the Context of Problem>
- Conclusion if Decision is Fail to Reject Ho:
 - “There is insufficient evidence to conclude”
 - <Ha in the Context of Problem>

52

Example - Conclusion

- Decision: Fail to Reject Ho
- There is insufficient evidence to conclude that the mean amount of soy sauce being filled into bottles is not 16 ounces.
- There is insufficient evidence to conclude machine that fills 16 ounce soy sauce bottles is operating improperly.

53

Conclusions

- Conclusions need to
 - Be consistent with the results of the Hypothesis Test.
 - Use language that is clearly understood in the context of the problem.
 - Limit the inference to the population that was sampled.
 - Report sampling methods that could question the integrity of the random sample assumption.
 - Conclusions should address the potential or necessity of further research, sending the process back to the first procedure.

54

Conclusions need to be consistent with the results of the Hypothesis Test.

- Rejecting H_0 requires a **strong statement** in support of H_a .
- Failing to Reject H_0 does NOT support H_0 , but requires a **weak statement** of insufficient evidence to support H_a .
- Example:
 - The researcher wants to support the claim that, on average, students send more than 1000 text messages per month
 - $H_0: \mu=1000$ $H_a: \mu>1000$
 - Conclusion if H_0 is rejected: The mean number of text messages sent by students exceeds 1000.
 - Conclusion if H_0 is not rejected: There is insufficient evidence to support the claim that the mean number of text messages sent by students exceeds 1000.

55

Conclusions need to use language that is clearly understood in the context of the problem.

- Avoid technical or statistical language.
- Refer to the language of the original general question.
- Compare these two conclusions from a test of correlation between home prices square footage and price.

Conclusion 1: By rejecting the Null Hypothesis we are inferring that the Alternative Hypothesis is supported and that there exists a significant correlation between the independent and dependent variables in the original problem comparing home prices to square footage.

Conclusion 2: Homes with more square footage generally have higher prices.

56

Conclusions need to limit the inference to the population that was sampled.

- If a survey was taken of a sub-group of population, then the inference applies to the subgroup.
- Example
 - Studies by pharmaceutical companies will only test adult patients, making it difficult to determine effective dosage and side effects for children.
 - "In the absence of data, doctors use their medical judgment to decide on a particular drug and dose for children. 'Some doctors stay away from drugs, which could deny needed treatment,' Blumer says. 'Generally, we take our best guess based on what's been done before.'"
 - "The antibiotic chloramphenicol was widely used in adults to treat infections resistant to penicillin. But many newborn babies died after receiving the drug because their immature livers couldn't break down the antibiotic."

source: FDA Consumer Magazine – Jan/Feb 2003

57

Conclusions need to report sampling methods that could question the integrity of the random sample assumption.

- Be aware of how the sample was obtained. Here are some examples of pitfalls:
 - Telephone polling was found to under-sample young people during the 2008 presidential campaign because of the increase in cell phone only households. Since young people were more likely to favor Obama, this caused bias in the polling numbers.
 - Sampling that didn't occur over the weekend may exclude many full time workers.
 - Self-selected and unverified polls (like ratemyprofessors.com) could contain immeasurable bias.

58

Conclusions should address the potential or necessity of further research, sending the process back to the first procedure.

- Answers often lead to new questions.
- If changes are recommended in a researcher's conclusion, then further research is usually needed to analyze the impact and effectiveness of the implemented changes.
- There may have been limitations in the original research project (such as funding resources, sampling techniques, unavailability of data) that warrants more a comprehensive study.

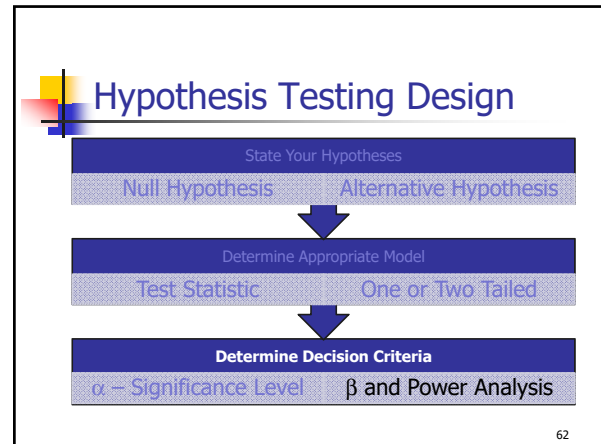
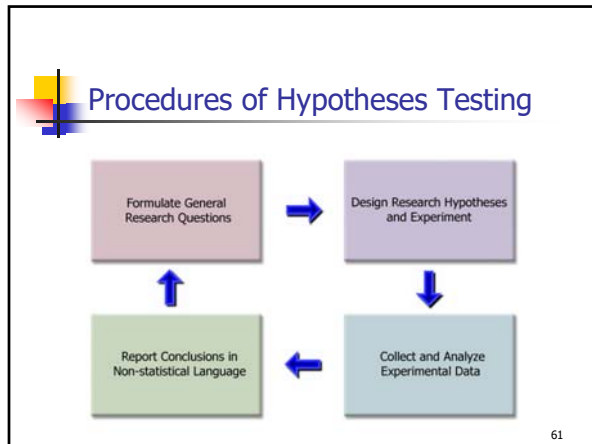
Example: A math department modifies its curriculum based on a performance statistics for an experimental course. The department would want to do further study of student outcomes to assess the effectiveness of the new program.

59

Soy Sauce Example - Conclusion

- There is insufficient evidence to conclude that the machine that fills 16 ounce soy sauce bottles is operating improperly.
- This conclusion is based on 36 measurements taken during a single hour's production run.
- We recommend continued monitoring of the machine during different employee shifts to account for the possibility of potential human error.

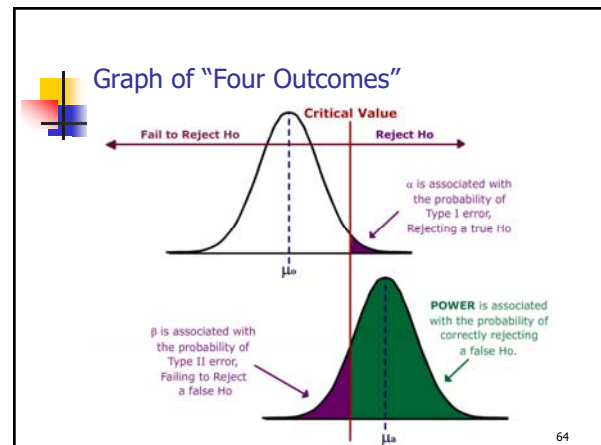
60



Statistical Power and Type II error

	Fail to Reject H_0	Reject H_0
H_0 is true	$1-\alpha$	α Type I error
H_0 is False	β Type II error	$1-\beta$ Power

63



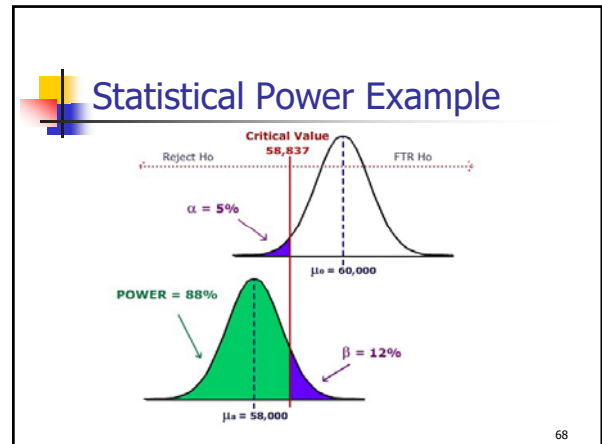
- ### Statistical Power (continued)
- Power is the probability of rejecting a false H_0 , when $\mu = \mu_a$
 - Power depends on:
 - Effect size $|\mu_o - \mu_a|$
 - Choice of α
 - Sample size
 - Standard deviation
 - Choice of statistical test
- 65

- ### Statistical Power Example
- Bus brake pads are claimed to last on average at least 60,000 miles and the company wants to test this claim.
 - The bus company considers a "practical" value for purposes of bus safety to be that the pads at least 58,000 miles.
 - If the standard deviation is 5,000 and the sample size is 50, find the Power of the test when the mean is really 58,000 miles. Assume $\alpha = .05$
- 66

Statistical Power Example

- Set up the test
 - $H_0: \mu \geq 60,000$ miles
 - $H_a: \mu < 60,000$ miles
 - $\alpha = 5\%$
- Determine the Critical Value
 - Reject H_0 if $\bar{X} > 58,837$
- Calculate β and Power
 - $\beta = 12\%$
 - Power = $1 - \beta = 88\%$

67



New Models, Similar Procedures

- The procedures outlined for the test of population mean vs. hypothesized value with known population standard deviation will apply to other models as well.
- Examples of some other one population models:
 - Test of population mean vs. hypothesized value, population standard deviation unknown.
 - Test of population proportion vs. hypothesized value.
 - Test of population standard deviation (or variance) vs. hypothesized value.

69

Testing for the Population Mean: Population Standard Deviation Unknown

- The test statistic for the one sample case is given by:

$$t = \frac{\bar{X} - \mu}{s / \sqrt{n}}$$
- The degrees of freedom for the test is $n-1$.
- The shape of the t distribution is similar to the Z , except the tails are fatter, so the logic of the decision rule is the same.

70

Decision Rules

- Like the normal distribution, the logic for one and two tail testing is the same.
- For a two-tail test using the t -distribution, you will reject the null hypothesis when the value of the test statistic is greater than $t_{df, \alpha/2}$ or if it is less than $-t_{df, \alpha/2}$
- For a left-tail test using the t -distribution, you will reject the null hypothesis when the value of the test statistic is less than $-t_{df, \alpha}$
- For a right-tail test using the t -distribution, you will reject the null hypothesis when the value of the test statistic is greater than $t_{df, \alpha}$

71

Example – one population test of mean, σ unknown

- Humerus bones from the same species have approximately the same length-to-width ratios. When fossils of humerus bones are discovered, archaeologists can determine the species by examining this ratio. It is known that Species A has a mean ratio of 9.6. A similar Species B has a mean ratio of 9.1 and is often confused with Species A.
- 21 humerus bones were unearthed in an area that was originally thought to be inhabited Species A. (Assume all unearthed bones are from the same species.)
- Design a hypotheses where the alternative claim would be the humerus bones were not from Species A.
- Determine the power of this test if the bones actually came from Species B (assume a standard deviation of 0.7)
- Conduct the test using at a 5% significance level and state overall conclusions.

72

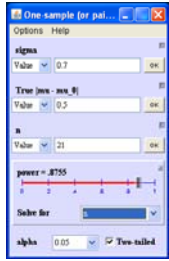
Example – Designing Test

- Research Hypotheses
 - Ho: The humerus bones are from Species A
 - Ha: The humerus bones are not from Species A
- In terms of the population mean
 - Ho: $\mu = 9.6$
 - Ha: $\mu \neq 9.6$
- Significance level
 - $\alpha = .05$
- Test Statistic (Model)
 - t-test of mean vs. hypothesized value.

73

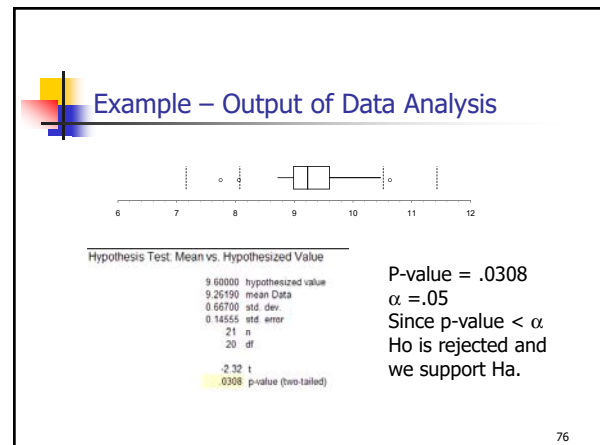
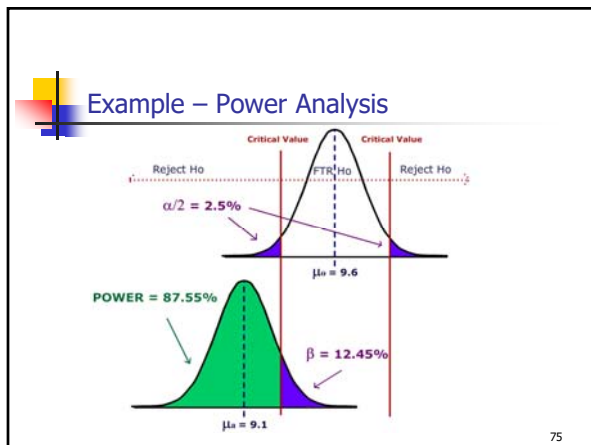
Example - Power Analysis

- Information needed for Power Calculation
 - $\mu_0 = 9.6$ (Species A)
 - $\mu_a = 9.1$ (Species B)
 - Effect Size = $|\mu_0 - \mu_a| = 0.5$
 - $\sigma = 0.7$ (given)
 - $\alpha = .05$
 - $n = 21$ (sample size)
 - Two tailed test
- Results using online Power Calculator*
 - Power = .8755
 - $\beta = 1 - \text{Power} = .1245$
 - If humerus bones are from Species B, test has an 87.55% chance of correctly rejecting Ho and a maximum Type II error of 12.55%



*source: Russ Lenth, University of Iowa - <http://www.stat.uiowa.edu/~rlenth/Power/>

74



Example - Conclusions

- Results:
 - The evidence supports the claim (pvalue<.05) that the humerus bones are not from Species A.
- Sampling Methodology:
 - We are assuming since the bones were unearthed in the same location, they came from the same species.
- Limitations:
 - A small sample size limited the power of the test, which prevented us from making a more definitive conclusion.
- Further Research
 - Test if the bone are from Species B or another unknown species.
 - Test to see if bones are the same age to support the sampling methodology.

77

Tests Concerning Proportion

- Proportion:** A fraction or percentage that indicates the part of the population or sample having a particular trait of interest.
- The population proportion is denoted by p .
- The sample proportion is denoted by \hat{p} where

$$\hat{p} = \frac{\text{number of successes in the sample}}{\text{number sampled}}$$

78

Test Statistic for Testing a Single Population Proportion

- If sample size is sufficiently large, \hat{p} has an approximately normal distribution. This approximation is reasonable if $np(1-p) > 5$

$$z = \frac{\hat{p} - p}{\sqrt{\frac{p(1-p)}{n}}}$$

$p = \text{population proportion}$
 $\hat{p} = \text{sample proportion}$

79

Example

- In the past, 15% of the mail order solicitations for a certain charity resulted in a financial contribution.
- A new solicitation letter has been drafted and will be sent to a random sample of potential donors.
- A hypothesis test will be run to determine if the new letter is more effective.
- Determine the sample size so that:
 - The test can be run at the 5% significance level.
 - If the letter has an 18% success rate, (an effect size of 3%), the power of the test will be 95%
- After determining the sample size, conduct the test.

80


Example – Designing Test

- Research Hypotheses
 - Ho: The new letter is not more effective.
 - Ha: The new letter is more effective.
- In terms of the population proportion
 - Ho: $p = 0.15$
 - Ha: $p > 0.15$
- Significance level
 - $\alpha = .05$
- Test Statistic (Model)
 - Z-test of proportion vs. hypothesized value.

81

Example - Power Analysis

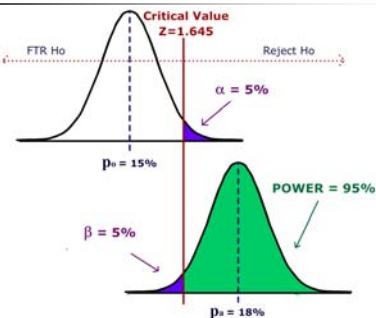
- Information needed for Sample Size Calculation
 - $p_0 = 0.15$ (current letter)
 - $p_a = 0.18$ (potential new letter)
 - Effect Size = $|p_0 - p_a| = 0.03$
 - Desired Power = 0.95
 - $\alpha = .05$
 - One tailed test
- Results using online Power Calculator*
 - Sample size = 1652
 - The charity should send out 1652 new solicitation letters to potential donors and run the test.



*source: Russ Lenth, University of Iowa – <http://www.stat.uiowa.edu/~rlenth/Power/>

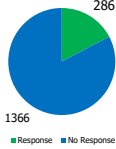
82

Example – Power Analysis



83

Example – Output of Data Analysis



Observed	Hypothesized
0.1731	0.15 p (as decimal)
286/1652	248/1652 p (as fraction)
286	247.8 X
1652	1652 n

0.0088 std. error
2.63 z
.0042 p-value (one-tailed, upper)

- P-value = .0042
- $\alpha = 0.05$
- Since p-value < α , Ho is rejected and we support Ha.

84

EXAMPLE

Critical Value Alternative Method

- Critical Value = 1.645 (95th percentile of the Normal Distribution.)
- H_0 is rejected if $Z > 1.645$
- Test Statistic:
$$Z = \frac{\left(\frac{286}{1652} - .15\right)}{\sqrt{\frac{(.15)(.85)}{1652}}} = 2.63$$
- Since $Z = 2.63 > 1.645$, H_0 is rejected. The new letter is more effective.

85

Example - Conclusions

- Results:
 - The evidence supports the claim ($p\text{-value} < .01$) that the new letter is more effective.
- Sampling Methodology:
 - The 1652 test letters were selected as a random sample from the charity's mailing list. All letters were sent at the same time period.
- Limitations:
 - The letters needed to be sent in a specific time period, so we were not able to control for seasonal or economic factors.
- Further Research
 - Test both solicitation methods over the entire year to eliminate seasonal effects.
 - Send the old letter to another random sample to create a control group.

86

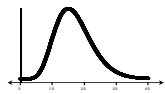
Test for Variance or Standard Deviation vs. Hypothesized Value

- We often want to make a claim about the variability, volatility or consistency of a population random variable.
- Hypothesized values for population variance σ^2 or standard deviation σ are tested with the χ^2 distribution.
- Examples of Hypotheses:
 - $H_0: \sigma = 10$ $H_a: \sigma \neq 10$
 - $H_0: \sigma^2 = 100$ $H_a: \sigma^2 > 100$
- The sample variance s^2 is used in calculating the Test Statistic.

87

Test Statistic uses χ^2 distribution

s^2 is the test statistic for the population variance. Its sampling distribution is a χ^2 distribution with $n-1$ d.f.



$$\chi^2 = \frac{(n-1)s^2}{\sigma_o^2}$$

88

Example

- A state school administrator claims that the standard deviation of test scores for 8th grade students who took a life-science assessment test is less than 30, meaning the results for the class show consistency.
- An auditor wants to support that claim by analyzing 41 students recent test scores, shown here:

57	75	86	92	101	108	110	120	155
63	77	88	96	102	108	111	122	
66	78	88	96	107	109	115	135	
68	81	92	98	107	109	115	137	
72	82	92	99	107	110	118	139	

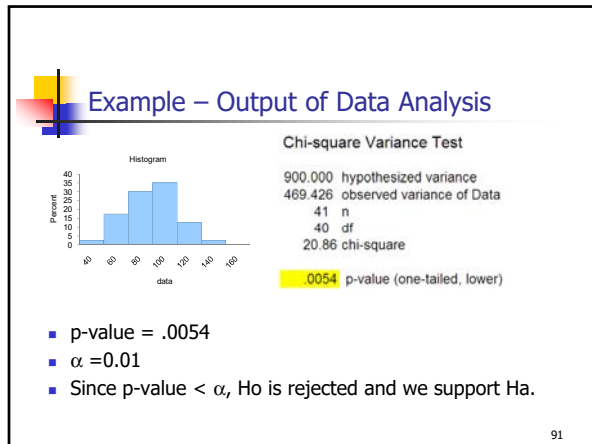
- The test will be run at 1% significance level.

89

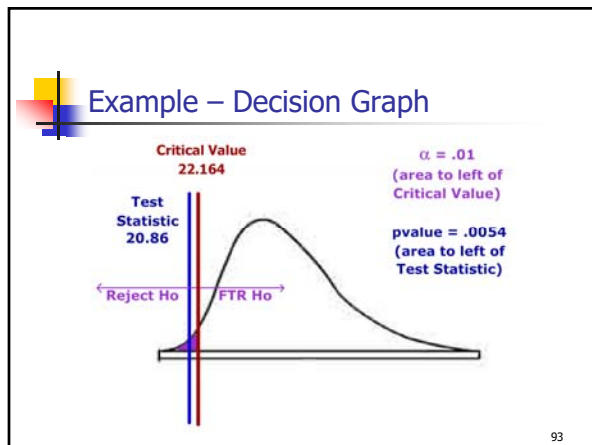
Example – Designing Test

- Research Hypotheses
 - H_0 : Standard deviation for test scores equals 30.
 - H_a : Standard deviation for test scores is less than 30.
- In terms of the population variance
 - $H_0: \sigma^2 = 900$
 - $H_a: \sigma^2 < 900$
- Significance level
 - $\alpha = .01$
- Test Statistic (Model)
 - χ^2 -test of variance vs. hypothesized value.

90




- ### EXAMPLE
- #### Critical Value Alternative Method
- Critical Value = 22.164 (1st percentile of the Chi-square Distribution.)
 - H_0 is rejected if $\chi^2 < 22.164$
 - Test Statistic: $\chi^2 = \frac{(40)(469.426)}{900} = 20.86$
 - Since $Z = 20.86 < 22.164$, H_0 is rejected. The claim that the standard deviation is under 30 is supported.
- 92



- ### Example - Conclusions
- Results:
 - The evidence supports the claim (pvalue<.01) that the standard deviation for 8th grade test scores is less than 30.
 - Sampling Methodology:
 - The 41 test scores were the results of the recently administered exam to the 8th grade students.
 - Limitations:
 - Since the exams were for the current class only, there is no assurance that future classes will achieve similar results.
 - Further Research
 - Compare results to other schools that administered the same exam.
 - Continue to analyze future class exams to see if the claim is holding true.
- 94

Inferential Statistics and Probability a Holistic Approach

Chapter 10 Two Population Inference


This Course Material by Maurice Geraghty is licensed under a Creative Commons Attribution-ShareAlike 4.0 International License. Conditions for use are shown here: <https://creativecommons.org/licenses/by-sa/4.0/>

1

Comparing two population means

- Four models
 - Independent Sampling
 - Known population variances
 - **Two sample Z - test**
 - The 2 population variances are equal
 - **Pooled variance t-test**
 - The 2 population variances are unequal
 - **t-test for unequal variances**
 - Dependent Sampling
 - **Matched Pairs t-test**

2

Independent Sampling

Population 1

μ_1, σ_1

↓

n_1
 \bar{X}_1, s_1

Population 2

μ_2, σ_2

↓

n_2
 \bar{X}_2, s_2

3

Dependent sampling

Population

↓

n

↙ ↘

Measurement 1

minus

Measurement 2

↙ ↘

\bar{X}_d, s_d

4

Difference of Two Population means

- $\bar{X}_1 - \bar{X}_2$ is Random Variable
- $\bar{X}_1 - \bar{X}_2$ is a point estimator for $\mu_1 - \mu_2$
- The standard deviation is given by the formula $\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}$
- If n_1 and n_2 are sufficiently large, follows a normal distribution.

5

Difference between two means – known population variances

- If both σ_1 and σ_2 are known and the two populations are independently selected, this test can be run.
- Test Statistic:

$$Z = \frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}}$$

6

Example 1

- Are larger houses more likely to have pools?
- The housing data square footage (size) was split into two groups by pool (Y/N).
- Test the hypothesis that the homes with pools have more square feet than the homes without pools. Let $\alpha = .01$

EXAMPLE 1 - Design

$H_o : \mu_1 \leq \mu_2$ $H_a : \mu_1 > \mu_2$
 $H_o : \mu_1 - \mu_2 \leq 0$ $H_a : \mu_1 - \mu_2 > 0$

$\alpha = .01$

$$Z = (\bar{X}_1 - \bar{X}_2) / (\sqrt{\sigma_1^2 / n_1 + \sigma_2^2 / n_2})$$

H_0 is rejected if $Z > 2.326$

EXAMPLE 1 Data

Population 1 Size with pool	Population 2 Size without pool
Sample size = 130	Sample size = 95
Sample mean = 26.25	Sample mean = 23.04
Pop Std Dev = 6.93	Pop Std Dev = 4.55

EXAMPLE 1 DATA

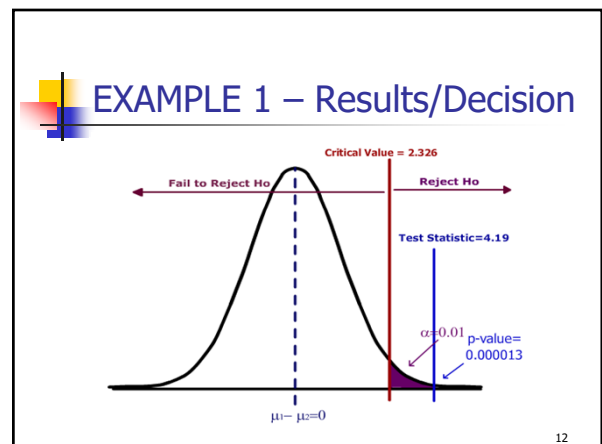
$$Z = \frac{(26.25 - 23.04) - 0}{\sqrt{\frac{6.93^2}{130} + \frac{4.55^2}{95}}} = 4.19$$

- Decision:** Reject H_0
- Conclusion:** Homes with pools have more mean square footage.

EXAMPLE 1 p-value method

Using Technology
Reject H_0 if the p-value $< \alpha$

	Sq ft with pool	Sq ft no pool
Mean	26.25	23.04
Std Dev	6.93	4.55
Observations	130	95
Hypothesized Mean Difference	0	
Z	4.19	
p-value	0.0000137	



Pooled variance t-test

- To conduct this test, three assumptions are required:
 - The populations must be normally or approximately normally distributed (or central limit theorem must apply).
 - The sampling of populations must be **independent**.
 - The **population variances** must be **equal**.

Pooled Sample Variance and Test Statistic

- Pooled Sample Variance:
$$s_p^2 = \frac{(n_1 - 1)s_1^2 + (n_2 - 1)s_2^2}{n_1 + n_2 - 2}$$
- Test Statistic:
$$t = \frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{s_p \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}}$$

$$df = n_1 + n_2 - 2$$

EXAMPLE 2

A recent EPA study compared the highway fuel economy of domestic and imported passenger cars.

- A sample of 12 imported cars revealed a mean of 35.76 mpg with a standard deviation of 3.86.
- A sample of 15 domestic cars revealed a mean of 33.59 mpg with a standard deviation of 2.16 mpg.
- At the .05 significance level can the EPA conclude that the mpg is higher on the imported cars? (Let subscript 2 be associated with domestic cars.)

EXAMPLE 2 – critical value method

- $H_o : \mu_1 \leq \mu_2$ $H_a : \mu_1 > \mu_2$
- $\alpha = .05$
- $t = (\bar{X}_1 - \bar{X}_2) / (s_p \sqrt{1/n_1 + 1/n_2})$
- H_o is rejected if $t > 1.708$, $df = 25$
- $t = 1.85$ H_o is rejected. Imports have a higher mean mpg than domestic cars.

t-test when variances are not equal.

- Test statistic:
$$t' = \frac{(\bar{X}_1 - \bar{X}_2) - (\mu_1 - \mu_2)}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$$
- Degrees of freedom:
$$df = \frac{\left(\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}\right)^2}{\left[\frac{(s_1^2/n_1)^2}{(n_1 - 1)} + \frac{(s_2^2/n_2)^2}{(n_2 - 1)}\right]}$$
- This test (also known as the Welch-Aspin Test) has **less power** than the prior test and should only be used when it is clear the population variances are different.

EXAMPLE 2

- $H_o : \mu_1 \leq \mu_2$ $H_a : \mu_1 > \mu_2$
- $\alpha = .05$
- t' test
- H_o is rejected if $t > 1.746$, $df = 16$
- $t' = 1.74$ H_o is not rejected. There is insufficient sample evidence to claim a higher mpg on the imported cars.

Using Technology

- Decision Rule: Reject H_0 if $p\text{-value} < \alpha$
- Megastat: Compare Two Independent Groups
- Use Equal Variance or Unequal Variance Test
- Use Original Data or Summarized Data

domestic 29.8 33.3 34.7 37.4 34.4 32.7 30.2 36.2 35.5 34.6 33.2 35.1 33.6 31.3 31.9

import 39.0 35.1 39.1 32.2 35.6 35.5 40.8 34.7 33.2 29.4 42.3 32.2

19

Pooled Variance t-test

- Minitab output
- p-value = 0.038
- p-value < $\alpha = .05$
- Reject H_0

```
Two-sample T for MPG
TYPE      N   Mean  StDev  SE Mean
import   12  35.76   3.86    1.1
US car   15  33.59   2.16    0.56

Difference =  $\mu$  (import) -  $\mu$  (US car)
Estimate for difference: 2.16
95% lower bound for difference: 0.16
T-Test of difference = 0 (vs >): T-Value = 1.85
P-Value = 0.038  DF = 25
Both use Pooled StDev = 3.0264
```

20

Unequal Variances t-test

- Minitab output
- p-value = 0.051
- p-value < $\alpha = .05$
- Fail to Reject H_0

```
Two-sample T for MPG
TYPE      N   Mean  StDev  SE Mean
import   12  35.76   3.86    1.1
US car   15  33.59   2.16    0.56

Difference =  $\mu$  (import) -  $\mu$  (US car)
Estimate for difference: 2.16
95% lower bound for difference: -0.01
T-Test of difference = 0 (vs >): T-Value = 1.74
P-Value = 0.051  DF = 16
```

21

Hypothesis Testing - Paired Observations

- Independent samples are samples that are not related in any way.
- Dependent samples are samples that are paired or related in some fashion.
 - For example, if you wished to buy a car you would look at the *same* car at two (or more) *different* dealerships and compare the prices.
- Use the following test when the samples are dependent:

22

Hypothesis Testing Involving Paired Observations

$$t = \frac{\bar{X}_d - \mu_d}{s_d / \sqrt{n}}$$

- where \bar{X}_d is the average of the differences
- s_d is the standard deviation of the differences
- n is the number of pairs (differences)

23

EXAMPLE 3

- An independent testing agency is comparing the daily rental cost for renting a compact car from Hertz and Avis.
- A random sample of 15 cities is obtained and the following rental information obtained.
- At the .05 significance level can the testing agency conclude that there is a difference in the rental charged?

24

Example 3 – continued

- Data for Hertz
 - $\bar{X}_1 = 46.67$
 - $s_1 = 5.23$
- Data for Avis
 - $\bar{X}_2 = 44.87$
 - $s_2 = 5.62$

City	Hertz	Avis
Atlanta	42	40
Baltimore	51	47
Boston	46	42
Chicago	56	52
Cleveland	45	43
Denver	48	48
Dallas	56	54
Honolulu	37	32
Los Angeles	51	48
Kansas City	45	48
Miami	41	39
New York	44	42
San Francisco	48	45
Seattle	46	50
Washington DC	44	43

25

Example 3 - continued

By taking the difference of each pair, variability (measured by standard deviation) is reduced.

$$\bar{X}_d = 1.80$$

$$s_d = 2.513$$

$$n = 15$$

City	Hertz	Avis	Difference
Atlanta	42	40	2
Baltimore	51	47	4
Boston	46	42	4
Chicago	56	52	4
Cleveland	45	43	2
Denver	48	48	0
Dallas	56	54	2
Honolulu	37	32	5
Los Angeles	51	48	3
Kansas City	45	48	-3
Miami	41	39	2
New York	44	42	2
San Francisco	48	45	3
Seattle	46	50	-4
Washington DC	44	43	1

26

EXAMPLE 3 continued

- $H_0: \mu_d = 0$ $H_1: \mu_d \neq 0$
- $\alpha = .05$
- Matched pairs t test, $df = 14$
- H_0 is rejected if $t < -2.145$ or $t > 2.145$
- $t = (1.80) / [2.513 / \sqrt{15}] = 2.77$
- Reject H_0 .
- There is a difference in mean price for compact cars between Hertz and Avis. Avis has lower mean prices.

27

Megastat Output – Example 3

Hypothesis Test: Paired Observations

```

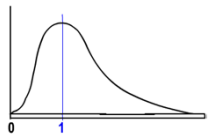
0.000 hypothesized value
46.667 mean Hertz
44.867 mean Avis
1.800 mean difference (Hertz - Avis)
2.513 std. dev.
0.649 std. error
15 n
14 df

2.77 t
.0149 p-value (two-tailed)
    
```

28

Characteristics of F-Distribution

- There is a "family" of F Distributions.
- Each member of the family is determined by two parameters: the numerator degrees of freedom and the denominator degrees of freedom.
- F cannot be negative, and it is a continuous distribution.
- The F distribution is positively skewed.
- Its values range from 0 to ∞ . As $F \rightarrow \infty$ the curve approaches the X-axis.



29

Test for Equal Variances

- For the two tail test, the test statistic is given by:

$$F = \frac{S_i^2}{S_j^2}$$
- s_i^2 and s_j^2 are the sample variances for the two populations.
- There are 2 sets of degrees of freedom: $n_i - 1$ for the numerator, $n_j - 1$ for the denominator

30

EXAMPLE 4

- A stockbroker at brokerage firm, reported that the mean rate of return on a sample of 10 software stocks was 12.6 percent with a standard deviation of 4.9 percent.
- The mean rate of return on a sample of 8 utility stocks was 10.9 percent with a standard deviation of 3.5 percent.
- At the .05 significance level, can the broker conclude that there is more variation in the software stocks?

31

Test Statistic depends on Hypotheses

Hypotheses	Test Statistic
$H_o : \sigma_1 \geq \sigma_2$	$F = \frac{s_2^2}{s_1^2}$ use α table
$H_a : \sigma_1 < \sigma_2$	
$H_o : \sigma_1 \leq \sigma_2$	$F = \frac{s_1^2}{s_2^2}$ use α table
$H_a : \sigma_1 > \sigma_2$	
$H_o : \sigma_1 = \sigma_2$	$F = \frac{\max(s_1^2, s_2^2)}{\min(s_1^2, s_2^2)}$ use $\alpha/2$ table
$H_a : \sigma_1 \neq \sigma_2$	

32

EXAMPLE 4 continued

- $H_o : \sigma_1 \leq \sigma_2$ $H_a : \sigma_1 > \sigma_2$
- $\alpha = .05$
- F-test
- H_o is rejected if $F > 3.68$, $df=(9,7)$
- $F = 4.9^2/3.5^2 = 1.96 \rightarrow$ Fail to Reject H_o .
- There is insufficient evidence to claim more variation in the software stock.

33

Excel Example

- Using Megastat – Test for equal variances under two population independent samples test and click the box to test for equality of variances
- The default p-value is a two-tailed test, so take one-half reported p-value for one-tailed tests
- Example – Domestic vs Import Data
- $H_o : \sigma_1 = \sigma_2$ $H_a : \sigma_1 \neq \sigma_2$
- $\alpha = .10$
- Reject H_o means use unequal variance t-test
- FTR H_o means use pooled variance t-test

34

Excel Output

F-test for equality of variance
 14.894 variance: import
 4.654 variance: domestic
 3.20 F
 .0438 p-value

pvalue < .10, Reject H_o

Use unequal variance t-test to compare means.

35

Comparing two proportions

- Suppose we take a sample of n_1 from population 1 and n_2 from population 2.
- Let X_1 be the number of success in sample 1 and X_2 be the number of success in sample 2.
- The sample proportions are then calculated for each group.

$$\hat{p}_1 = \frac{X_1}{n_1} \qquad \hat{p}_2 = \frac{X_2}{n_2}$$

36

Hypothesis testing for 2 Proportions

- In conducting a Hypothesis test where the Null hypothesis assumes equal proportions, it is best practice to pool or combine the sample proportions into a single estimated proportion, and use an estimated standard error.

$$\bar{p} = \frac{X_1 + X_2}{n_1 + n_2} \quad s_{\hat{p}_1 - \hat{p}_2} = \sqrt{\frac{\bar{p}(1-\bar{p})}{n_1} + \frac{\bar{p}(1-\bar{p})}{n_2}}$$

37

Hypothesis testing for 2 Proportions

- The test statistic will have a Normal Distribution as long as there are at least 10 successes and 10 failures in both samples.

$$Z = \frac{(\hat{p}_1 - \hat{p}_2) - (p_1 - p_2)}{\sqrt{\frac{\bar{p}(1-\bar{p})}{n_1} + \frac{\bar{p}(1-\bar{p})}{n_2}}}$$

38

Example

- In an August 2016 Study, Pew Research asked the sampled Americans if background checks required at gun stores should be made universal extended to all sales of guns between private owners or at gun shows.
- 772 out 990 men said yes, while 857 out of 1020 women said yes.
- Is there a difference in the proportion of men and women who support universal background checks for purchasing guns? Design and conduct the test with a significance level of 1%.

39

Example (Design)

- Ho:** $p_m = p_w$ (There is no difference in the proportion of support for background checks by gender)
- Ha:** $p_m \neq p_w$ (There is a difference in the proportion of support for background checks by gender)
- Model:** Two proportion Z test. This is a two-tailed test with $\alpha = 0.01$.
- Model Assumptions:** for men there are 772 yes and 218 no. For women there are 857 yes and 163 no. Since all these numbers exceed 10, the model is appropriate.
- Decision Rules:** Critical Value Method - Reject Ho if $Z > 2.58$ or $Z < -2.58$. P-value method - Reject Ho if p-value < 0.01

40

Example (Results)

$$\hat{p}_m = \frac{772}{990} = 0.780 \quad \hat{p}_w = \frac{857}{1020} = 0.840$$

$$\bar{p} = \frac{772+857}{990+1020} = 0.810 \quad Z = \frac{(0.780 - 0.840) - 0}{\sqrt{\frac{0.810(1-0.810)}{990} + \frac{0.810(1-0.810)}{1020}}} = -3.45$$

p-value = 0.0005 $< \alpha$
Reject Ho Under both methods.

Conclusion: There is a difference in the proportion of support for background checks by gender. Women are more likely to support background checks.

41

Inferential Statistics and Probability
a Holistic Approach

Chapter 11
Chi-square Tests for
Categorical Data

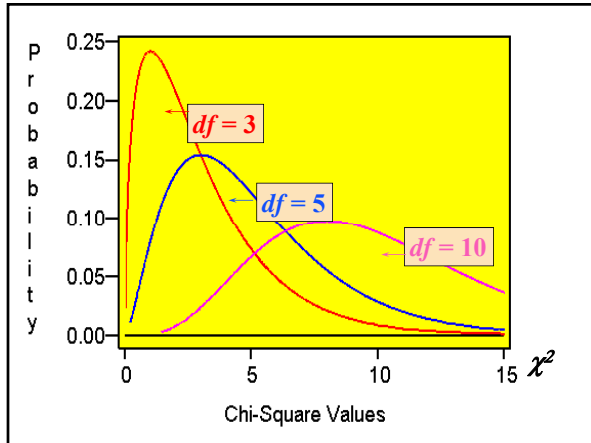
This Course Material by Maurice Geraghty is licensed under a Creative Commons Attribution-ShareAlike 4.0 International License. Conditions for use are shown here: <https://creativecommons.org/licenses/by-sa/4.0/>

1

Characteristics of the Chi-Square Distribution

- The major characteristics of the chi-square distribution are:
 - It is positively skewed
 - It is non-negative
 - It is based on degrees of freedom
 - When the degrees of freedom change a new distribution is created

2



Goodness-of-Fit Test: Equal Expected Frequencies

- Let O_i and E_i be the observed and expected frequencies respectively for each category.
- H_0 : there is no difference between Observed and Expected Frequencies
- H_a : there is a difference between Observed and Expected Frequencies
- The test statistic is: $\chi^2 = \sum \frac{(O_i - E_i)^2}{E_i}$
- The critical value is a chi-square value with $(k-1)$ degrees of freedom, where k is the number of categories

4

EXAMPLE 1

The following data on absenteeism was collected from a manufacturing plant. At the .01 level of significance, test to determine whether there is a difference in the absence rate by day of the week.

Day	Frequency
Monday	95
Tuesday	65
Wednesday	60
Thursday	80
Friday	100

5

EXAMPLE 1 *continued*

- Assume equal expected frequency: $(95+65+60+80+100)/5=80$

Day	O	E	$(O-E)^2/E$
Mon	95	80	2.8125
Tues	65	80	2.8125
Wed	60	80	5.0000
Thur	80	80	0.0000
Fri	100	80	5.0000
Total	400	400	15.625

6

EXAMPLE 1 *continued*

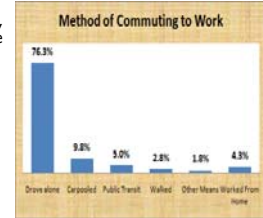
- H_0 : there is no difference between the observed and the expected frequencies of absences.
- H_a : there is a difference between the observed and the expected frequencies of absences.
- Test statistic: $\chi^2 = \sum(O-E)^2/E = 15.625$
- Decision Rule: reject H_0 if test statistic is greater than the critical value of 13.277. (4 df, $\alpha = .01$)
- Conclusion: reject H_0 and conclude that there is a difference between the observed and expected frequencies of absences.

7

Goodness-of-Fit Test: Unequal Expected Frequencies

EXAMPLE 2

- In the 2010 United States census, data was collected on how people get to work -- their method of commuting.
- Suppose you wanted to know if people who live in the San Jose metropolitan area (Santa Clara County) commute with similar proportions as the United States.
- Design and conduct a hypothesis test at the 5% significance level.



8

EXAMPLE 2 *continued*

Method Of Commuting	Observed Frequency O_i	Expected Proportion p_i	Expected Frequency E_i	$\sum \frac{(O - E)^2}{E}$
Drive Alone	764	0.763	763	0.0013
Carpooled	105	0.098	98	0.5000
Public Transit	34	0.050	50	5.1200
Walked	20	0.028	28	2.2857
Other Means	30	0.018	18	8.0000
Worked from Home	47	0.043	43	0.3721
TOTAL	1000	1.000	1000	16.2791

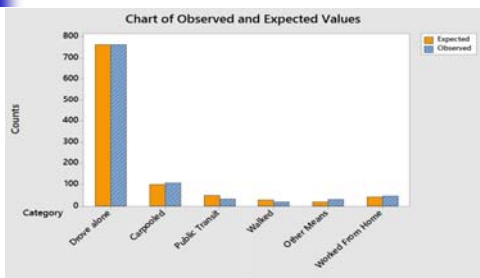
9

EXAMPLE 2 *continued*

- Design:**
 - Ho:** $p_1 = .763$ $p_2 = .098$ $p_3 = .050$ $p_4 = .028$ $p_5 = .018$ $p_6 = .043$
 - Ha:** At least one p_i is different than what was stated in Ho
 - $\alpha = .05$
 - Model: Chi-Square Goodness of Fit, $df = 5$
 - H_0 is rejected if $\chi^2 > 11.071$
- Data:**
 - $\chi^2 = 16.2791$, Reject H_0
- Conclusion:**
 - Workers in Santa Clara County do not have the same frequencies of method of commuting as workers in the entire United States.

10

EXAMPLE 2 *continued*



11

Contingency Table Analysis

- Contingency table analysis is used to test whether two traits or variables are related.
- Each observation is classified according to two variables.
- The usual hypothesis testing procedure is used.
- The *degrees of freedom* is equal to: (number of rows-1)(number of columns-1).
- The expected frequency is computed as: $\text{Expected Frequency} = (\text{row total})(\text{column total})/\text{grand total}$

12

EXAMPLE 3

- In May 2014, Colorado became the first state to legalize the recreational use of marijuana.
- A poll of 1000 adults were classified by gender and their opinion about legalizing marijuana
- At the .05 level of significance, can we conclude that gender and the opinion about legalizing marijuana for recreational use are dependent events?

Marijuana should be	Men	Women	Total
Legal	270	230	500
Not Legal	205	245	450
Unsure	25	25	50
Total	500	500	1000

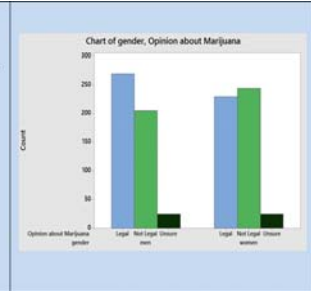
13

EXAMPLE 3 *continued*

Rows: opinion about Marijuana
Columns: gender

1st Value = Observed
2nd Value = Expected
3rd Value = Contribution to Chi-square

	men	women	All
Legal	270 250 1.600	230 250 1.600	500
Not Legal	205 225 1.778	245 225 1.778	450
Unsure	25 0.000	25 0.000	50
All	500	500	1000



14

EXAMPLE 3 *continued*


- Design:** H_0 : Gender and Opinion are independent.
 H_a : Gender and Opinion are dependent.
- $\alpha = .05$
- Model: Chi-Square Test for Independence, $df=2$
- H_0 is rejected if $\chi^2 > 5.99$
- Data:** $\chi^2 = 6.756$, Reject H_0
- Conclusion:** Gender and opinion are dependent variables. Men are more likely to support legalizing marijuana for recreational use.

15

Inferential Statistics and Probability a Holistic Approach

Chapter 12

One Factor Analysis of Variance (ANOVA)


This Course Material by Maurice Geraghty is licensed under a Creative Commons Attribution-ShareAlike 4.0 International License. Conditions for use are shown here: <https://creativecommons.org/licenses/by-sa/4.0/>

1

Underlying Assumptions for ANOVA

- The F distribution is also used for testing the equality of more than two means using a technique called analysis of variance (ANOVA). ANOVA requires the following conditions:
 - The populations being sampled are normally distributed.
 - The populations have equal standard deviations.
 - The samples are randomly selected and are independent.

2

ANOVA Definitions


- **Factor** – categorical variable that defines the populations.
- **Response** – variable that is being measured.
- **Levels** – the number of choices for the factor, represented by k
- **Replicates** – the sample size for each level, n_1, n_2, \dots, n_k .
- If $n_1 = n_2 = \dots = n_k$, then the design is **balanced**.

- **H₀**: There is no difference in the mean <response in context> due to the <factor in context>.
- **H_a**: There is a difference in the mean <response in context> due to the <factor in context>.

3

Characteristics of F-Distribution

- There is a "family" of F Distributions.
- Each member of the family is determined by two parameters: the numerator degrees of freedom and the denominator degrees of freedom.
- F cannot be negative, and it is a continuous distribution.
- The F distribution is positively skewed.
- Its values range from 0 to ∞ . As $F \rightarrow \infty$ the curve approaches the X -axis.



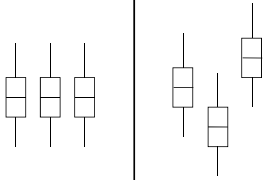
4

Analysis of Variance Procedure

- **The Null Hypothesis**: the population means are the same.
- **The Alternative Hypothesis**: at least one of the means is different.
- **The Test Statistic**: $F = (\text{between sample variance}) / (\text{within sample variance})$.
- **Decision rule**: For a given significance level α , reject the null hypothesis if $F(\text{computed})$ is greater than $F(\text{table})$ with numerator and denominator degrees of freedom.

5

ANOVA – Null Hypothesis



Ho is true -all means the same

Ho is false -not all means the same

6

ANOVA NOTES

- If there are k populations being sampled, then the df (numerator)=k-1
- If there are a total of n sample points, then df (denominator) = n-k
- The test statistic is computed by: $F = [(SS_F)/(k-1)] / [(SS_E)/(N-k)]$.
- SS_F represents the factor (between) sum of squares.
- SS_E represents the error (within) sum of squares.
- Let T_c represent the column totals, n_c represent the number of observations in each column, and ΣX represent the sum of all the observations.
- These calculations are tedious, so technology is used to generate the **ANOVA table**.

7

Formulas for ANOVA

$$SS_{Total} = \Sigma(X^2) - \frac{(\Sigma X)^2}{n}$$

$$SS_{Factor} = \Sigma \left(\frac{T_c^2}{n_c} \right) - \frac{(\Sigma X)^2}{n}$$

$$SS_{Error} = SS_{Total} - SS_{Factor}$$

8

ANOVA Table

Source	SS	df	MS	F
Factor	SS_{Factor}	k-1	SS_F/df_F	MS_F/MS_E
Error	SS_{Error}	n-k	SS_E/df_E	
Total	SS_{Total}	n-1		

9

EXAMPLE

Party Pizza specializes in meals for students. Hsieh Li, President, recently developed a new tofu pizza.

- Before making it a part of the regular menu she decides to test it in several of her restaurants. She would like to know if there is a difference in the mean number of tofu pizzas sold per day at the Cupertino, San Jose, and Santa Clara pizzerias for sample of five days.
- At the .05 significance level can Hsieh Li conclude that there is a difference in the mean number of tofu pizzas sold per day at the three pizzerias?

10

Example

	Cupertino	San Jose	Santa Clara	Total
	13	10	18	
	12	12	16	
	14	13	17	
	12	11	17	
			17	
T	51	46	85	182
n	4	4	5	13
Means	12.75	11.5	17	14
Σ^2	653	534	1447	2634

11

Example continued

$$SS_{Total} = 2634 - \frac{182^2}{13} = 86$$

$$SS_{Factor} = 2624.25 - \frac{182^2}{13} = 76.25$$

$$SS_{Error} = 86 - 76.25 = 9.75$$

12

Example 4 *continued*

ANOVA TABLE

Source	SS	df	MS	F
Factor	76.25	2	38.125	39.10
Error	9.75	10	0.975	
Total	86.00	12		

13

EXAMPLE 4 *continued*

- **Design:** $H_0: \mu_1 = \mu_2 = \mu_3$
 H_a : Not all the means are the same
- $\alpha = .05$
- Model: One Factor ANOVA
- H_0 is rejected if $F > 4.10$
- **Data:** Test statistic: $F = [76.25/2]/[9.75/10] = 39.1026$
- H_0 is rejected.
- **Conclusion:** There is a difference in the mean number of pizzas sold at each pizzeria.

14

One-way ANOVA: Cupertino, San Jose, Santa Clara

Source	DF	SS	MS	F	P
Factor	2	76.250	38.125	39.10	0.000
Error	10	9.750	0.975		
Total	12	86.000			

S = 0.9874 R-Sq = 88.66% R-Sq(adj) = 86.40%

Individual 95% CIs For Mean Based on Pooled StDev

Level	N	Mean	StDev
Cupertino	4	12.750	0.957
San Jose	4	11.500	1.291
Santa Clara	5	17.000	0.707

15

Post Hoc Comparison Test

- Used for pairwise comparison
- Designed so the **overall** significance level is 5%.
- Use technology.
- Refer to **Tukey Test** Material in Supplemental Material.

16

Post Hoc Comparison Test

Grouping Information Using Tukey Method

	N	Mean	Grouping
Santa Clara	5	17.0000	A
Cupertino	4	12.7500	B
San Jose	4	11.5000	B

Means that do not share a letter are significantly different.

17


Post Hoc Comparison Test

Individual Value Plot of Cupertino, San Jose, Santa Clara

18

Inferential Statistics and Probability a Holistic Approach

Chapter 13 Correlation and Linear Regression


This Course Material by Maurice Geraghty is licensed under a Creative Commons Attribution-ShareAlike 4.0 International License. Conditions for use are shown here: <https://creativecommons.org/licenses/by-sa/4.0/>

1

Mathematical Model

- You have a small business producing custom t-shirts.
- Without marketing, your business has revenue (sales) of \$1000 per week.
- Every dollar you spend marketing will increase revenue by 2 dollars.
- Let variable X represent amount spent on marketing and let variable Y represent revenue per week.
- Write a mathematical model that relates X to Y

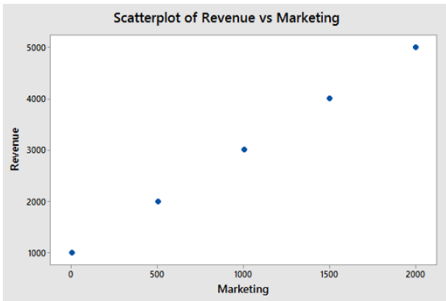
2

Mathematical Model - Table

X=marketing	Y=revenue
\$0	\$1000
\$500	\$2000
\$1000	\$3000
\$1500	\$4000
\$2000	\$5000

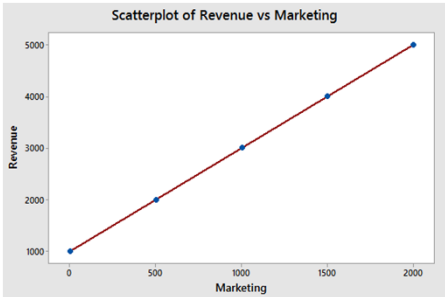
3

Mathematical Model - Scatterplot



4

Mathematical Model - Linear



5

Mathematical Linear Model

Linear Model	Example
$Y = \beta_0 + \beta_1 X$	$Y = 1000 + 2X$
Y : Dependent Variable	Y : Revenue
X : Independent Variable	X : Marketing
β_0 : Y-intercept	β_0 : \$1000
β_1 : Slope	β_1 : \$2 per \$1 marketing

6

Statistical Model

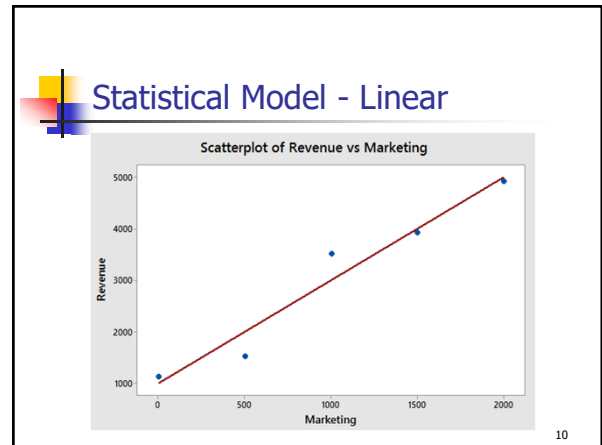
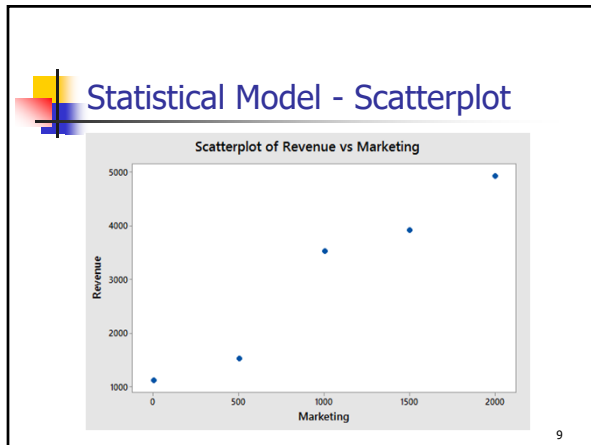
- You have a small business producing custom t-shirts.
- Without marketing, your business has revenue (sales) of \$1000 per week.
- Every dollar you spend marketing will increase revenue by an expected value of 2 dollars.
- Let variable X represent amount spent on marketing and let variable Y represent revenue per week.
- Let ϵ represent the difference between Expected Revenue and Actual Revenue (Residual Error)
- Write a statistical model that relates X to Y

7

Statistical Model - Table

X=Marketing	Expected Revenue	Y=Actual Revenue	ϵ =Residual Error
\$0	\$1000	\$1100	+\$100
\$500	\$2000	\$1500	-\$500
\$1000	\$3000	\$3500	+\$500
\$1500	\$4000	\$3900	-\$100
\$2000	\$5000	\$4900	-\$100

8



Statistical Linear Model

Regression Model	Example
$Y = \beta_0 + \beta_1 X + \epsilon$	$Y = 1000 + 2X + \epsilon$
Y : Dependent Variable	Y : Revenue
X : Independent Variable	X : Marketing
β_0 : Y-intercept	β_0 : \$1000
β_1 : Slope	β_1 : \$2 per \$1 marketing
ϵ : Normal(0, σ)	

11

Regression Analysis

- Purpose:** to determine the regression equation; it is used to predict the value of the dependent response variable (Y) based on the independent explanatory variable (X).
- Procedure:**
 - select a sample from the population
 - list the paired data for each observation
 - draw a scatter diagram to give a visual portrayal of the relationship
 - determine the regression equation.

12

Simple Linear Regression Model

$$Y = \beta_0 + \beta_1 X + \varepsilon$$

Y : *Dependent Variable*
 X : *Independent Variable*
 β_0 : *Y-intercept*
 β_1 : *Slope*
 ε : *Normal (0, σ)*

13

Estimation of Population Parameters

- From sample data, find statistics that will estimate the 3 population parameters
- Slope parameter
 - b_1 will be an estimator for β_1
- Y-intercept parameter
 - b_0 will be an estimator for β_0
- Standard deviation
 - s_e will be an estimator for σ

14

Regression Analysis

- the regression equation: $\hat{Y} = b_0 + b_1 X$, where:
 - \hat{Y} is the average predicted value of Y for any X .
 - b_0 is the Y-intercept, or the estimated Y value when $X=0$
 - b_1 is the slope of the line, or the average change in \hat{Y} for each change of one unit in X
 - the least squares principle is used to obtain b_1 and b_0

$$SSX = \sum X^2 - \frac{1}{n}(\sum X)^2 \qquad b_1 = \frac{SSXY}{SSX}$$

$$SSY = \sum Y^2 - \frac{1}{n}(\sum Y)^2$$

$$SSXY = \sum XY - \frac{1}{n}(\sum X \cdot \sum Y) \qquad b_0 = \bar{Y} - b_1 \bar{X}$$

15

Assumptions Underlying Linear Regression

- For each value of X , there is a group of Y values, and these Y values are *normally distributed*.
- The *means* of these normal distributions of Y values all lie on the straight line of regression.
- The *standard deviations* of these normal distributions are equal.
- The Y values are statistically independent. This means that in the selection of a sample, the Y values chosen for a particular X value do not depend on the Y values for any other X values.

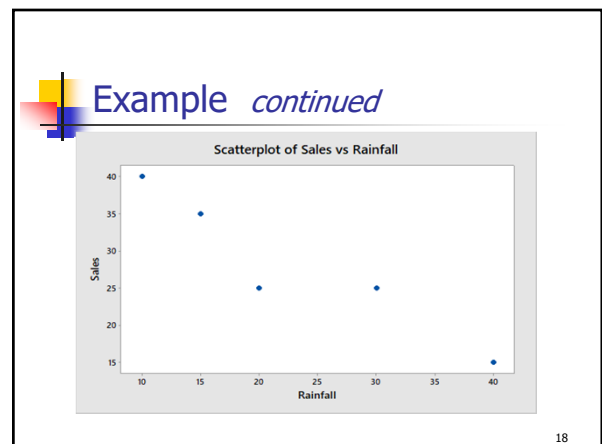
16

Example

- X = Average Annual Rainfall (Inches)
- Y = Average Sale of Sunglasses/1000
 - Make a Scatterplot
 - Find the least square line

X	10	15	20	30	40
Y	40	35	25	25	15

17



Example *continued*

	X	Y	X ²	Y ²	XY
	10	40	100	1600	400
	15	35	225	1225	525
	20	25	400	625	500
	30	25	900	625	750
	40	15	1600	225	600
Σ	115	140	3225	4300	2775

19

- ### Example *continued*
- Find the least square line
 - SSX = 580
 - SSY = 380
 - SSXY = -445
 - $b_1 = -.767$
 - $b_0 = 45.647$
 - $\hat{Y} = 45.647 - .767X$
- 20

The Standard Error of Estimate

- The **standard error of estimate** measures the scatter, or dispersion, of the observed values around the line of regression
- The formulas that are used to compute the standard error:

$$SSR = b_1 \cdot SSXY$$

$$SSE = \sum (y - \hat{y})^2 = SSY - SSR$$

$$MSE = \frac{SSE}{(n - 2)}$$

$$s_e = \sqrt{MSE}$$

21

Example *continued*

- Find SSE and the standard error:

x	y	\hat{y}	$y - \hat{y}$	$(y - \hat{y})^2$
10	40	37.97	2.03	4.104
15	35	34.14	0.86	0.743
20	25	30.30	-5.30	28.108
30	25	22.63	2.37	5.620
40	15	14.96	0.04	0.002
			Total	38.578

 - SSR = 341.422
 - SSE = 38.578
 - MSE = 12.859
 - $s_e = 3.586$

22

Correlation Analysis

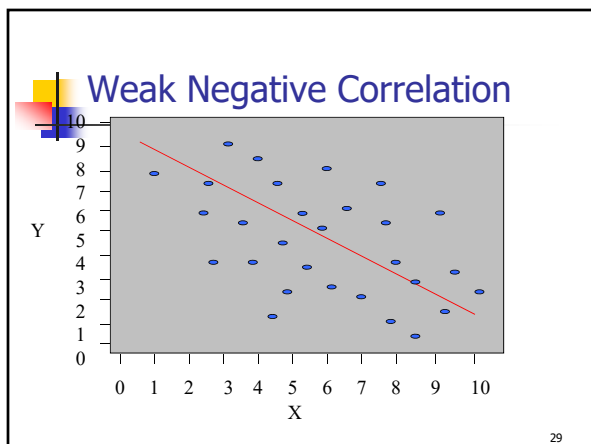
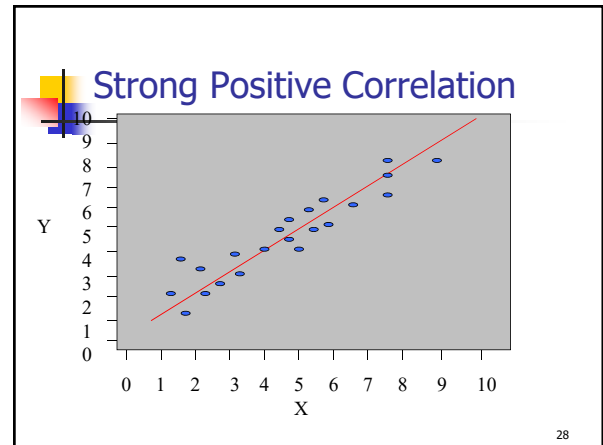
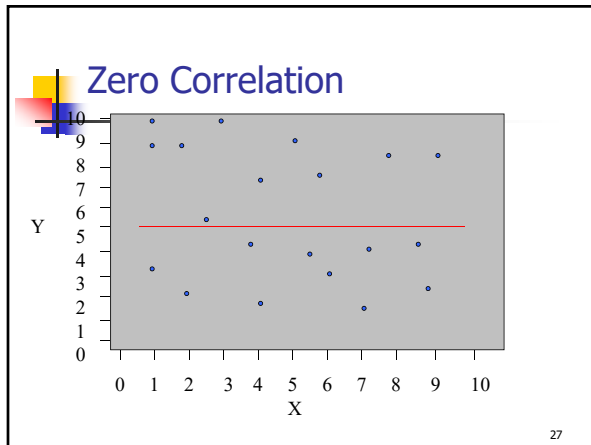
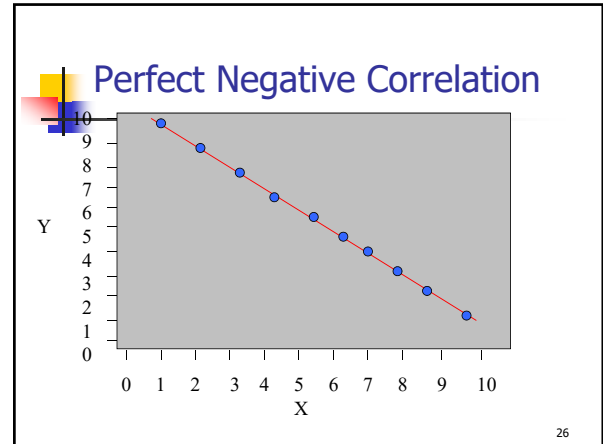
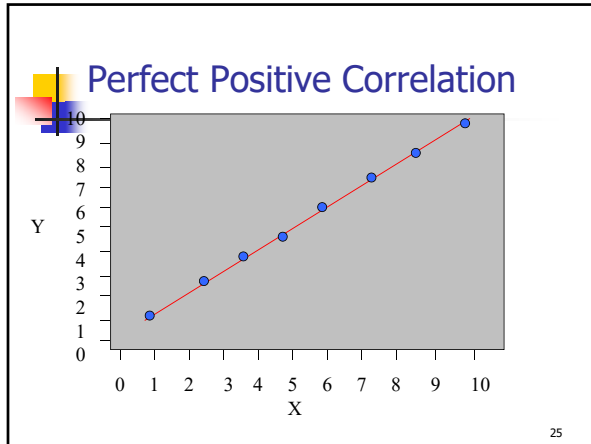
- **Correlation Analysis:** A group of statistical techniques used to measure the strength of the relationship (correlation) between two variables.
- **Scatter Diagram:** A chart that portrays the relationship between the two variables of interest.
- **Dependent Variable:** The variable that is being predicted or estimated. "Effect"
- **Independent Variable:** The variable that provides the basis for estimation. It is the predictor variable. "Cause?" (Maybe!)

23

The Coefficient of Correlation, r

- The **Coefficient of Correlation** (r) is a measure of the **strength** of the relationship between two variables.
 - It requires interval or ratio-scaled data (variables).
 - It can range from -1.00 to 1.00.
 - Values of -1.00 or 1.00 indicate perfect and strong correlation.
 - Values close to 0.0 indicate weak correlation.
 - Negative values indicate an inverse relationship and positive values indicate a direct relationship.

24



- ### Causation
- Correlation does not necessarily imply causation.
 - There are 4 possibilities if X and Y are correlated:
 1. X causes Y
 2. Y causes X
 3. X and Y are caused by something else.
 4. Confounding - The effect of X and Y are hopelessly mixed up with other variables.
- 30

Causation - Examples

- City with more police per capita have more crime per capita.
- As Ice cream sales go up, shark attacks go up.
- People with a cold who take a cough medicine feel better after some rest.

31

r^2 : Coefficient of Determination

- r^2 is the proportion of the total variation in the dependent variable Y that is explained or accounted for by the variation in the independent variable X.
- The coefficient of determination is the square of the coefficient of correlation, and ranges from 0 to 1.

32

Formulas for r and r^2

$$r = \frac{SSXY}{\sqrt{SSX \cdot SSY}} \quad r^2 = \frac{SSR}{SSY}$$

$$SSX = \sum X^2 - \frac{1}{n}(\sum X)^2$$

$$SSY = \sum Y^2 - \frac{1}{n}(\sum Y)^2$$

$$SSXY = \sum XY - \frac{1}{n}(\sum X \cdot \sum Y)$$

$$SSR = SSY - \left(\frac{SSXY^2}{SSX} \right)$$

33

Example

- X = Average Annual Rainfall (Inches)
- Y = Average Sale of Sunglasses/1000

X	10	15	20	30	40
Y	40	35	25	25	15

34

Example *continued*

- Make a Scatter Diagram
- Find r and r^2

35

Example *continued*

X	Y	X^2	Y^2	XY
10	40	100	1600	400
15	35	225	1225	525
20	25	400	625	500
30	25	900	625	750
40	15	1600	225	600
115	140	3225	4300	2775

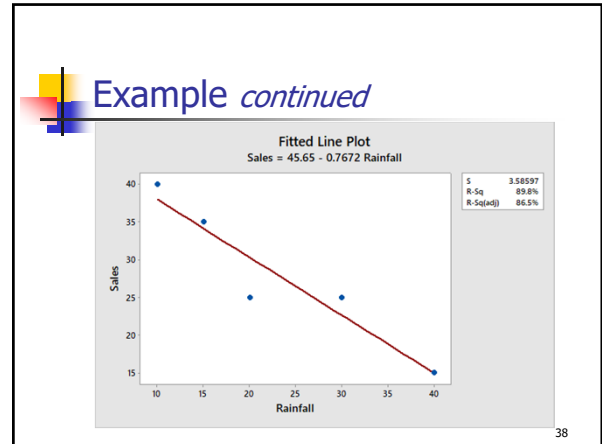
- $SSX = 3225 - 115^2/5 = 580$
- $SSY = 4300 - 140^2/5 = 380$
- $SSXY = 2775 - (115)(140)/5 = -445$

36

Example *continued*

- $r = -445/\sqrt{580 \times 330} = -.9479$
 - Strong negative correlation
- $r^2 = .8985$
 - About 89.85% of the variability of sales is explained by rainfall.

37



Characteristics of F-Distribution

- There is a "family" of F Distributions.
- Each member of the family is determined by two parameters: the numerator degrees of freedom and the denominator degrees of freedom.
- F cannot be negative, and it is a continuous distribution.
- The F distribution is positively skewed.
- Its values range from 0 to ∞ . As $F \rightarrow \infty$ the curve approaches the X-axis.

39

Hypothesis Testing in Simple Linear Regression

- The following Tests are equivalent:
 - H_0 : X and Y are uncorrelated
 - H_a : X and Y are correlated
 - H_0 : $\beta_1 = 0$
 - H_a : $\beta_1 \neq 0$
- Both can be tested using ANOVA

40

ANOVA Table for Simple Linear Regression

Source	SS	df	MS	F
Regression	SSR	1	SSR/dfR	MSR/MSE
Error/Residual	SSE	n-2	SSE/dfE	
TOTAL	SSY	n-1		

41

Example *continued*

- Test the Hypothesis $H_0: \beta_1 = 0, \alpha = 5\%$

Source	SS	df	MS	F	p-value
Regression	341.422	1	341.422	26.551	0.0142
Error	38.578	3	12.859		
TOTAL	380.000	4			

- Reject H_0 p-value $< \alpha$

42

Confidence Interval

- The confidence interval for the mean value of Y for a given value of X is given by:

$$\hat{Y} \pm t \cdot s_e \cdot \sqrt{\frac{1}{n} + \frac{(X - \bar{X})^2}{SSX}}$$

- Degrees of freedom for t = n-2

43

Prediction Interval

- The prediction interval for an individual value of Y for a given value of X is given by:

$$\hat{Y} \pm t \cdot s_e \cdot \sqrt{1 + \frac{1}{n} + \frac{(X - \bar{X})^2}{SSX}}$$

- Degrees of freedom for t = n-2

44

Example *continued*

- Find a 95% Confidence Interval for Sales of Sunglasses when rainfall = 25 inches.
- Find a 95% Prediction Interval for Sales of Sunglasses when rainfall = 25 inches.

45

Example – Minitab output

- Sales = 45.65 - 0.767 Rainfall
- Variable Setting
- Rainfall 25

Fit	SE Fit	95% CI	95% PI
26.4655	1.63111	(21.2746, 31.6564)	(13.9282, 39.0028)

46

Example *continued*

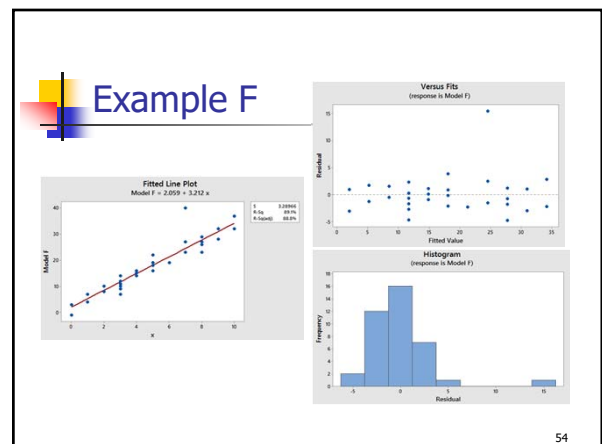
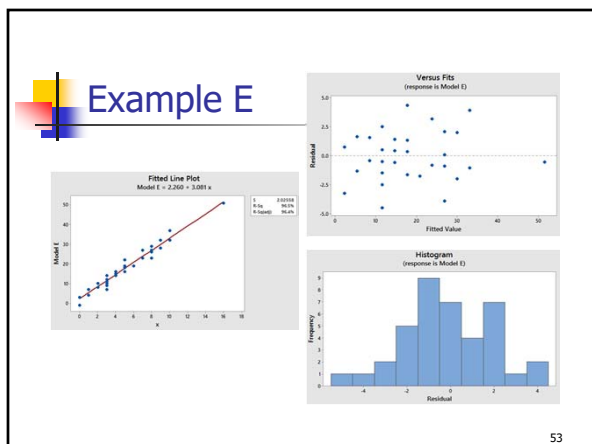
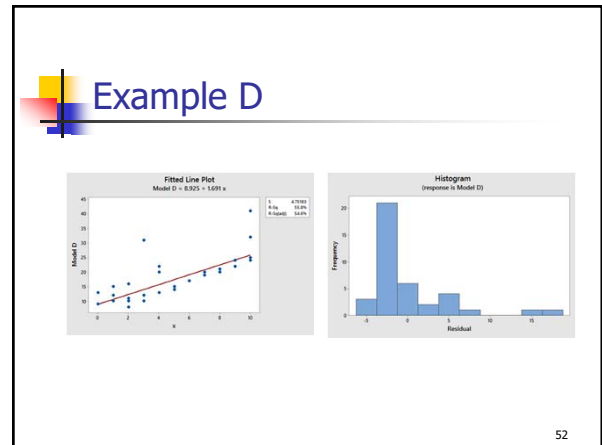
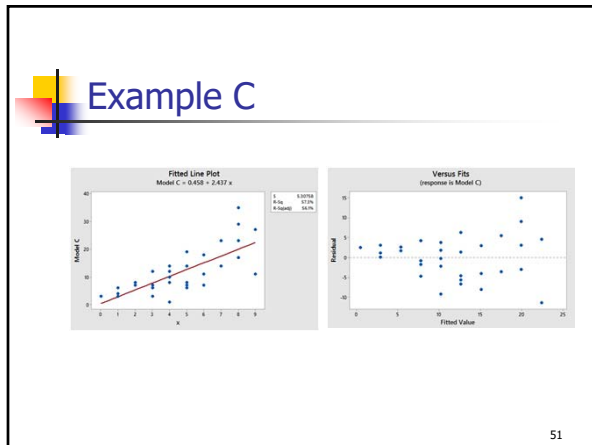
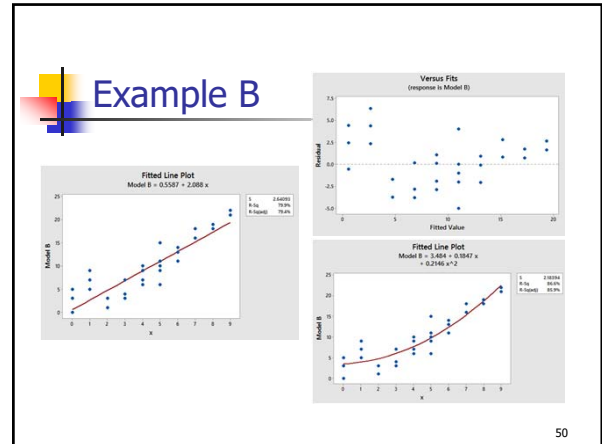
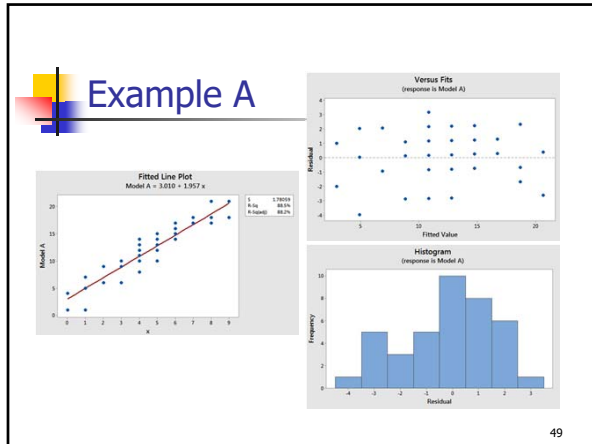
- 95% Confidence Interval
22.63 ± 6.60
- 95% Prediction Interval
22.63 ± 13.18

47

Residual Analysis

- Residuals should
 - have a normal distribution with constant σ
 - be mutually independent
 - not follow a pattern
 - be checked for outliers
 - with respect the line
 - with respect to X

48



Using Minitab to Run Regression

- Data shown is engine size in cubic inches (X) and MPG (Y) for 20 cars.

x	y	x	y
400	15	104	25
455	14	121	26
113	24	199	21
198	22	360	10
199	18	307	10
200	21	318	11
97	27	400	9
97	26	97	27
110	25	140	28
107	24	400	15

55

Using Minitab to Run Regression

Select Graphs>Scatterplot with regression line

56

Using Minitab to Run Regression

Select Statistics>Regression>Regression, then choose the Response (Y-variable) and model (X-variable)

57

Using Minitab to Run Regression

Click the results box, and choose the fits and residuals to get all predictions.

58

Using Minitab to Run Regression

The results at the beginning are the regression equation, the intercept and slope, the standard error of the residuals, and the r^2

The regression equation is
 $mpg = 30.2 - 0.0466 \text{ EngineSize}$

Predictor	Coef	SE Coef	T	P
Constant	30.203	1.361	22.20	0.000
EngineSize	-0.046598	0.005378	-8.66	0.000

S = 2.95688 R-Sq = 80.7% R-Sq(adj) = 79.6%

59

Using Minitab to Run Regression

Next is the ANOVA table, which tests the significance of the regression model.

Source	DF	SS	MS	F	P
Regression	1	656.42	656.42	75.08	0.000
Residual Error	18	157.38	8.74		
Total	19	813.80			

60

Using Minitab to Run Regression

Finally, the residuals show the potential outliers.

Obs	EngineSize	mpg	Fit	SE Fit	Residual	St Resid
1	400	15.000	11.564	1.167	3.436	1.26
2	455	14.000	9.001	1.421	4.999	1.93
3	113	24.000	24.937	0.880	-0.937	-0.33
4	198	22.000	20.976	0.673	1.024	0.36
5	199	18.000	20.930	0.672	-2.930	-1.02
6	200	21.000	20.883	0.671	0.117	0.04
7	97	27.000	25.683	0.939	1.317	0.47
8	97	26.000	25.683	0.939	0.317	0.11
9	110	25.000	25.077	0.891	-0.077	-0.03
10	107	24.000	25.217	0.902	-1.217	-0.43
11	104	25.000	25.357	0.913	-0.357	-0.13
12	121	26.000	24.565	0.853	1.435	0.51
13	199	21.000	20.930	0.672	0.070	0.02
14	360	10.000	13.427	0.998	-3.427	-1.23
15	307	10.000	15.897	0.807	-5.897	-2.07R
16	318	11.000	15.385	0.842	-4.385	-1.55
17	400	9.000	11.564	1.167	-2.564	-0.94
18	97	27.000	25.683	0.939	1.317	0.47
19	140	28.000	23.679	0.792	4.321	1.52
20	400	15.000	11.564	1.167	3.436	1.26

61

Using Minitab to Run Regression

- Find a 95% confidence interval for the **expected** MPG of a car with an engine size of 250 ci.
- Find a 95% prediction interval for the **actual** MPG of a car with an engine size of 250 ci.

mpg = 30.20 - 0.04660 EngineSize

Variable	Setting	Fit	SE Fit	95% CI	95% PI
EngineSize	250	18.5533	0.679201	(17.1264, 19.9803)	(12.1793, 24.9273)

62

Residual Analysis

- Residuals for Simple Linear Regression
 - The residuals should represent a linear model.
 - The standard error (standard deviation of the residuals) should not change when the value of X changes.
 - The residuals should follow a normal distribution.
 - Look for any potential extreme values of X.
 - Look for any extreme residual errors

63