

## Study Guide for Stat 200 Exam 1 (Chapters 1-17)

### Part I Study Design Practice Problems

Controlled Experiments—Researchers assigns subjects to treatment and control groups

Main Idea: Treatment and Control should be as much alike as possible

- Randomized, double-blind design is ideal because it eliminates systematic differences (bias). Random differences average out with enough subjects. Blocking reduces random differences that could be a problem for small studies by breaking subjects into similar sub-groups **before** randomization .

Once subjects are randomized into treatment and control, NEVER rearrange them because it will introduce bias. That's why we compare the results of everyone in treatment to everyone in control whether or not they adhered or not.

- Non-randomized controls usually introduce systematic difference between treatment and control groups that could bias the result. These differences are called *confounders*.

Observational Studies—Subjects themselves or simple fate determines treatment and control groups.

Researcher just observes.

Main Idea: Treatment and Control groups are likely to be systematically different, these differences can mix up or confound the results.

- Very difficult to conclude causation from association.
- With observational studies you must always think about what the likely confounders.
- Stratification adjusts for possible confounders by breaking subjects into sub groups where the confounding factor is the same.
- Simpson's Paradox is an example of extreme confounding. It's paradoxical because you get one result before stratification and the opposite afterwards!

#### Question 1

Two experiments were done comparing the effects of listening to classical music versus pop music while studying. All the students in both experimental designs were given an identical 2-hour lesson and then allowed time to study for a short exam.

In **Design A** students themselves chose to study either listening to classical or pop.

In **Design B** the students were randomly assigned to study either listening to classical or pop.

Design A found that the classical study group scored significantly higher on the exam than the pop group did. Design B found no significant difference in exam scores between the 2 groups. **The overall exam average in both designs was the same.**

a) Which design had randomized controls?      A only              B only              Both              Neither

b) Which design is more likely to have confounders?    A                      B                      Both are equally likely

c) Which conclusion is best supported by the evidence?      **Circle one**

- i) Students learn better when they are able to choose their own music while studying.
- ii) Students who choose classical are different in more ways than just their musical tastes than students who choose rap .
- iii) Classical music seems to enhance learning better than pop music.

## Question 2

A study published in the March 4, 2015 issue of the Journal of the American Medical Association evaluated whether peanut consumption might be more effective than peanut avoidance in preventing the development of peanut allergies in infants who are at high risk for the allergy. 640 infants aged 4 to 11 months with severe eczema and egg allergies (high risk indicators for peanut allergy) were **randomly assigned** to either consume (treatment) or avoid peanuts (control) until 5 years of age. The results were striking—17.2% of the children in the peanut-avoidance group tested positive for peanut allergy while only 3.2% of the group in the peanut-consumption group tested positive.

a) Which of the following best describes this study:

- i) A randomized controlled experiment
- ii) An observational study with controls
- iii) A non-randomized controlled experiment

b) Does the study show that eating peanuts helped prevent the children in the study from developing a peanut allergy?

- i) No, it only shows that there is an association between peanut consumption and reduced rate of peanut allergy since many environmental, cultural, social and biological factors contribute to both diet and allergic responses.
- ii) No, simply assigning children to 2 groups without considering the consequences of how peanut consumption or peanut avoidance may confer nutritional advantages limits any causal conclusions.
- iii) Yes, the study is strong evidence that peanut consumption helped prevent peanut allergy in these children although the causal mechanism can only be inferred.

c) Which of the following could confound the results? Circle Yes or No for each.

- i) Cultural/Ethnic differences- Peanuts and peanut oil are popular in West African and Southeast Asian cuisines, groups that have a relatively low incidence of peanut allergies.  
a. Yes                      b. No
- ii) Health Benefits – Peanuts are a relatively healthy snack food. Children who eat peanuts may be healthier in general and less likely to develop allergies.                      a. Yes                      b. No
- iii) Pre-existing Health Problems- The children all had severe health problems to begin with making it difficult to discern whether or not it was the peanuts or pre-existing conditions that led to the development of a peanut allergy.                      a. Yes  
b. No
- iv) Overactive Immune System- Children with overactive immune systems are both more likely to have egg allergies (like the children in the study) and to develop a peanut allergy. a. Yes  
b. No

d) 40 of the 640 infants showed evidence (by a skin-prick test) of already having a peanut allergy before they were even assigned to treatment or control. The researchers want to make sure that the 40 children are exactly evenly divided between the treatment and control groups but they don't want to introduce bias. What should they do?

- i) They should divide the infants into 2 groups (40 with pre-existing peanut allergy, and 600 without). Then randomly assign half of each group to treatment and half to control.
- ii) Randomly assign half of the 640 infants to treatment and half to control. This will ensure the infants will be evenly divided on all characteristics relevant to the response including pre-existing peanut allergy.
- iii) Randomly assign half of the 640 infants to treatment and half to control. In the unlikely event that the 2 groups are not balanced then, the researchers should balance the groups taking into account all variables to be as objective as possible.

**Question 3 pertains to the following study:**

A study was done to test whether Ginkgo biloba (GB) could alleviate symptoms of Alzheimer’s and dementia. The 52-week study randomly assigned half of the patients take GB daily and half to take a placebo. Neither the subjects nor evaluators knew who was in each group. At the end of the study, there was significant evidence that GB improved the cognitive performance and the social functioning of the patients for 6 months to 1 year.

- a) What type of bias could be present in this study **Choose one:**  
i) No systematic bias ii) Subject Bias iii) Evaluator Bias iv) Selection Bias v) ii, iii, and iv
- b) Which of the following could confound the results? **Choose one:**  
i) Forgetfulness- Patients with dementia may forget to take the GB on a regular basis.  
ii) Increased Attention-- Participation in the study increased the attention these patients received. They felt less neglected and therefore more cognitively active.  
iii) More motivated-- Those who volunteered to be in the GB group were probably more conscientious and motivated to begin with since they actively sought a remedy for their condition.  
iv) All of the above  
v) None of the above
- c) Not everyone in the treatment and control group adhered to the program and took their medicine/placebo. Which comparison is best when analyzing the final data?  
i) Compare everyone assigned to take the GB to everyone assigned to take the placebo.  
ii) Compare everyone who actually took the GB to the placebo group.  
iii) Compare only those who took the GB regularly to only those who took the placebo regularly.

**Question 4 pertains to the following study:**

A study was done to test the effectiveness of a new weight loss drug. The subjects were 2000 obese adults. Half were randomly assigned to take the drug every day and half were randomly assigned to take the placebo every day. Neither the subjects nor those who evaluated them knew who was in which group. The subjects were followed for 1 year and the percent of weight they lost or gained was recorded.

- a) Based only on the information above which of the following best describes the study above?  
**Choose one:**  
i) This was a non-randomized controlled experiment with a placebo.  
ii) This was a randomized controlled experiment without a placebo.  
iii) This was an observational study with controls.  
iv) This was a randomized controlled double-blind experiment.

b) The table below gives the average percent weight change of “adherers” and “non-adherers” in both the drug and the placebo group. Adherers regularly took their pills while non-adherers took their pills less than 80% of the time.

	Drug		Placebo	
	Number	%Weight change	Number	%Weight change
Adherers	500	7% loss	502	7.1% lost
Non-Adherers	500	2% gain	498	2.1% gain
Total	1000	2.5 loss	1000	2.52% lost

Based on the results of the table would you conclude there is good evidence for the following statements?

**Circle YES or NO after each statement:**

- i) The drug worked better than the placebo for those who regularly took the medicine. YES NO  
ii) The drug works no better than a placebo. YES NO  
iii) Adherers may be different than non-adherers in ways that help them lose weight. YES NO  
(for example, more responsible about eating balanced meals, exercising regularly, etc.)

### Question 5

A study published in the Feb 18, 2004 issue of the Journal of the American Medical Association compared pharmacy and medical records of 10,219 women and found that women who filled 25 or more prescriptions for antibiotics over a 17 year period received breast cancer diagnoses at twice the rate as those who took no antibiotics. The study concluded that high antibiotic usage increases one's risk of breast cancer.

a) Which of the following statements best describes this study? **Circle one:**

- i) This was a randomized controlled experiment without a placebo.
- ii) This was an observational study with controls.
- iii) This was a randomized controlled double-blind experiment.
- iv) This was a non-randomized controlled experiment with a placebo.

b) Based on the results of this study alone, which of the following statements is best? **Circle one.**

- i) High antibiotic use causes an increased risk of breast cancer.
- ii) High antibiotic use is *associated* with and *may* cause increased breast cancer risk.
- iii) High antibiotic use is *associated* with but does *not* cause increased breast cancer risk.
- iv) Having cancer is likely to cause increased use of antibiotics.

c) Below are either confounders, causal links or neither. Answer based only on given information.

- i) Age of first pregnancy- women who have their first child after the age of 35 are more likely to get breast cancer.  
a) Confounder                      b) Causal Link                      c) Neither
- ii) Destruction of Protective Bacteria- antibiotics kill healthy bacteria that may help prevent breast cancer.  
a) Confounder                      b) Causal Link                      c) Neither
- iii) Underlying Immune Problem- a weak immune system leads both to frequent infections necessitating antibiotics *and* also to a higher cancer risk.  
a) Confounder                      b) Causal Link                      c) Neither
- iv) Regular Check-ups- Women who regularly go to the doctor are both more likely to be prescribed antibiotics and more likely to receive a breast cancer diagnosis (especially for slow growing cancers that are unlikely to lead to serious health problems.)  
a) Confounder                      b) Causal Link                      c) Neither

d) Suppose the researchers thought that income was a possible confounder since high income women tend to take more antibiotics and tend to get more breast cancer. To separate out the effects of income from the effects of antibiotics researchers should ... **Circle one:**

- i) split the data into high, middle and low income groups and compare the antibiotic usage between the 3 groups.
- ii) split the data into high, middle and low income groups and compare the cancer rate of those who took a lot of antibiotics to those who took no antibiotics within each group.
- iii) split the data into high and low antibiotic users and compare the cancer rates between the groups.
- iv) split the data into 2 groups—breast cancer and no breast cancer and compare antibiotic usage between the 2 groups.

**Question 6** A study published in the August 15, 2017 issue of *Mayo Clinic Proceedings* tracked 44,000 people aged 20 to 87 for an average of about 16 years and found that those who drank 4 or more cups of coffee a day were 21% more likely to die than those who drank less than 4 cups a day. The risk was 50% higher for heavy coffee drinkers under 55 years of age.

**b)** Which of the following best describes this study?

- i)** An observational study with controls
- ii)** A randomized controlled experiment
- iii)** A non-randomized experiment with historical controls

**c)** Does the study show that drinking 4 or more cups of coffee a day caused the higher death rate?

- i)** No, the study was conducted over such a long time period that it's difficult to determine whether it was the original coffee drinking itself or something *else* about the coffee (for example, the way it was brewed) that caused the higher death rate.
- ii)** Yes, particularly for young people, the study clearly shows that excessive coffee drinking caused an increased risk of death.
- iii)** No, it's possible that coffee drinkers share other traits (besides the coffee) that could put them at a higher risk of dying.
- iv)** No, you cannot conclude causation without a proven causal mechanism. The study does provide strong evidence that it's the coffee that's raising the death rate and not something else, but it fails to explain how or why.

**c)** The study reported that they controlled for cigarette smoking. This means they thought smoking might be a confounder so they eliminated its confounding effect. How did they do that? **Choose one:**

- i)** At the beginning of the study, they divided the patients into smokers and non-smokers and then randomly divided the smokers and non-smokers equally between the coffee and no coffee groups.
- ii)** Throughout the study they eliminated anyone who smoked from the study.
- iii)** At the end of the study, they stratified on smoking, and compared the death rate of coffee drinkers to non-coffee drinkers within each smoking level (non-smokers, light smokers, heavy smokers).

**d)** State whether the following are confounders, causal links, or neither:

- i)** Increased popularity of coffee- The study was conducted over a 16-year time period that coincided with an enormous increase in coffee consumption. **a)** confounder **b)** causal link **c)** neither
  
- ii)** Caffeine—Excessive caffeine intake from 4 cups of coffee per day raises health risks because it increases a person's heart rate and blood pressure, which increase one's risk of death. **a)** confounder **b)** causal link **c)** neither
  
- iii)** Unhealthy Diet – The study stated that people who drank 4 or more cups of coffee were also more likely to have an unhealthy diet that could increase one's risk of death. **a)** confounder **b)** causal link **c)** neither
  
- iv)** Pre-existing-conditions- Some members of the study may have had pre-existing conditions or illness that would cause them to die sooner. **a)** confounder **b)** causal link **c)** neither.

**Question 7** A country club gives a pass-fail golf test every year to professional and amateur golfers. Professionals have a much higher % passing than amateurs. The club members were happy that the overall % passing went up from 68% in 2007 to 70% in 2017 and wanted to know which group contributed to the improved rate.

	2007				2017			
	Number	# Passes	# Failures	% Passing	Number	#Passes	# Failures	% Passing
<b>Professionals</b>	100	92	8	92%	100	90	10	90%
<b>Amateurs</b>	300	180	120	60%	100	50	50	50%
<b>Overall Total</b>	400	272	128	68%	200	140	60	70%

- a) Which group's % passing went up from 2007 to 2017? **Choose one:** a) Prof. b) Amat. c) Neither d) Both
- b) Is it possible for each group's % passing to go down if their overall % passing goes up?
- i) Yes, it's possible because the overall makeup of the club has changed from 25% to 50% professionals which raises the overall % passing even though both groups % passing declined.
- ii) No, it's not possible. If the overall passing rate goes up, then at least one group's passing rates must go up.

**Question 8**

A company has 455 job openings- 70 white collar jobs and 385 blue collar jobs. 600 men and 300 women apply for the new jobs. Here's the data:

	Men			Women		
	# Applied	# Hired	Hiring Rate	# Applied	# Hired	Hiring Rate
<b>White Collar</b>	200	30	15%	200	40	20%
<b>Blue Collar</b>	400	300	75%	100	85	85%
<b>Total</b>	600	330	55%	300	125	41.67%

- a) Overall 55% of the men but only 41.67% of the women who applied were hired, raising the question of sexual discrimination. Assuming that the men and women were equally qualified, which job category was discriminating against women? **Choose one:** i) White Collar Only ii) Blue Collar Only iii) Neither iv) Both
- b) Based only on the data above, if you're applying for a *white collar* job are hiring rates better for males or females? **Choose one:** i) Male ii) Female
- c) Based only on the data above, if you're applying for a *blue collar* job are hiring rates better for males or females? **Choose one:** i) Male ii) Female iii) Not possible to compare rates since 400 men applied, but only 100 women.
- d) Based only on the data above, are hiring rates better for males or females? **Choose one:** i) Better for males. ii) Better for females. iii) Same. iv) Depends on whether it's a white or blue collar job.

**Part II Descriptive Statistics**  
**Chapter 3 – Measures of Center and Spread**

**Question 1** pertains to the following list of 5 numbers: 0, 1, -2, 2, 9

- a) The average is \_\_\_\_\_
- b) The median is \_\_\_\_\_.
- c) The **deviations** from the **average** are \_\_\_\_, \_\_\_\_, \_\_\_\_, \_\_\_\_, \_\_\_\_ (List them in order from smallest to largest).
- d) The sum of the deviations from the average should = \_\_\_\_\_. *Fill in the blank with a number. (1pts.)*
- e) Compute the Standard Deviation. Round your answer to 2 decimal places.  
**Show your work. You may start with the deviations you found in part (c). Circle answer.**

**Question 2**

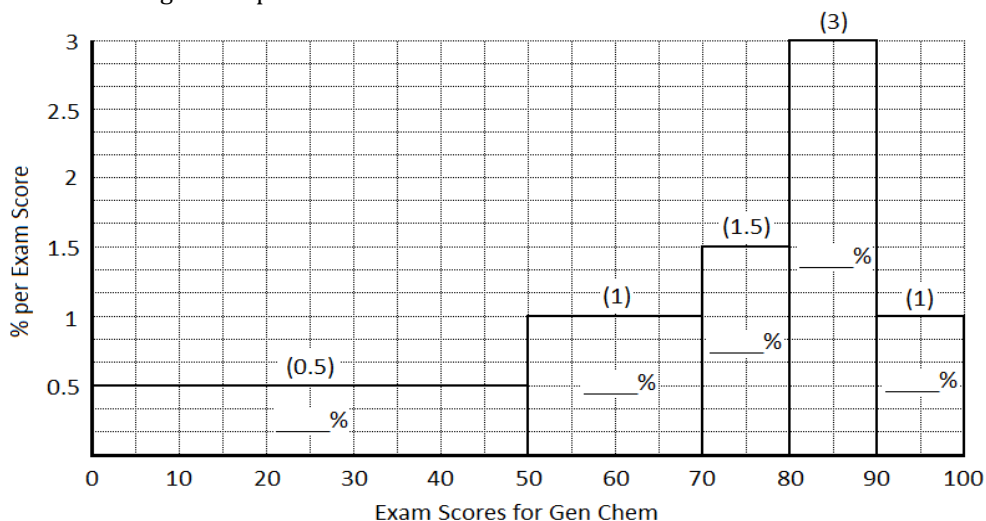
**A list of 10 numbers has an average = 5, median= 4, and SD = 3. Fill in the chart below with numbers.**

For (a-e) below, calculate the new average, median, and SD after the original list has been changed according to the given directions.	<b>New Average</b> (Write a number, not words, like "increase" or "decrease")	<b>New Median</b> (Write a number, not words, like "increase" or "decrease".)	<b>New SD</b> (Write a number, not words, like "increase" or "decrease" except for (e).
<b>a)</b> 5 is added to every number on the original list.			
<b>b)</b> Every number on the original list is multiplied by <b>negative 2</b> .			
<b>c)</b> Every number on the original list is divided by 2.			
<b>d)</b> Subtract 5 from every number on the original list, and then divide every number by 3.			
<b>e)</b> Every number on the original list remains the same, EXCEPT that 10 is added to the largest number.			<b>Choose one:</b> <b>i)</b> Increase <b>ii)</b> Decrease <b>iii)</b> Stays the same

## Chapter 4 Graphical Displays for Numerical Data

### Question 1 pertains to the histogram below.

The figure below is a histogram for the first exam scores of 520 freshmen and sophomores in general chemistry. The height of each block is given in parentheses.



- a) What percent of the students received an exam score between 0 and 50? Write your answer inside the blank provided in the 0-50 interval on the histogram. Do the same for the other 4 intervals. **Fill in ALL 5 blanks in each block of the histogram above with the correct areas.**
- b) The area of the entire histogram is \_\_\_\_\_%
- c) The median exam score is **closest** to:  
**Choose one:**    50    70    73    80    90
- d) Is the median  $>$ ,  $<$ , or  $=$  to the average? \_\_\_\_\_
- e) The percent of students who received exactly 75 on their first exam is closest to (Assume an equal distribution throughout the interval)  
**Choose one:**    0.5%    1%    1.5%    10%    15%
- f) Suppose all the students in the 0-30 range were given extra credit that raised each of their scores 20 points? How would that affect the average, median and Standard Deviation?  
**(Check the appropriate boxes below, check only 1 box per row.)**

	Increase	Decrease	Stay the same	Not enough information
Average would ...				
Median would ...				
Standard Deviation would ...				

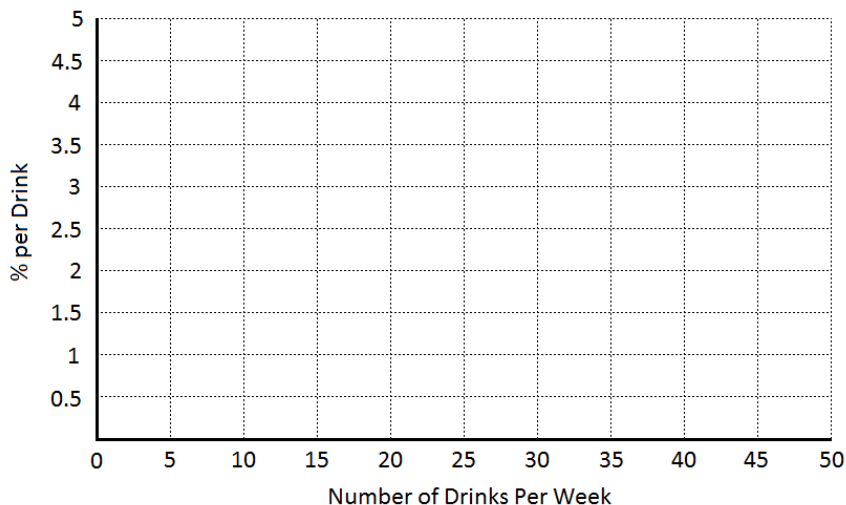


**Question 2**

A distribution table for the number of drinks a past semester of Stat 100 students said they typically consumed per week is shown below. The first row says that 45% of students said they had between 0 and 10 drinks per week. The table has 5 missing blanks. Fill them in with the correct widths, heights, and areas. Then draw the histogram. Write the area of each interval inside the block.

a) Fill in the 5 blanks in the table below and **then draw the histogram on the graph below.**

Interval	Width of Interval	Height (% per Drink)	Area (%)
0 to 10	10		45
10 to 15	5	4	
15 to 20	5	3	15
20 to 30	10		
30 to 50	20		10



b) The area column should sum to \_\_\_\_\_ %. **Fill in blank.**

c) If someone drinks more than 90% of the class, how much does he or she drink per week? \_\_\_\_\_ drinks **Fill in blank.**

d) Would it be appropriate to use a normal approximation for this data?

*Choose one:*

- i) No, the histogram is far from normal, so using a normal approximation would not be appropriate.
- ii) Yes, because converting to z-scores will change the shape and make the histogram normal.
- iii) Yes, because the normal approximation is suitable for all data sets.
- iv) Yes, because we can determine the average and SD from the data.

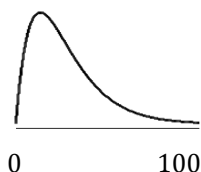
a) The Survey only allowed students to give answers up to 50 drinks. I gave everyone who answered 50 the opportunity to change their answers. A few of them changed their answer from 50 to 60 drinks. How would that affect the average, median and standard deviation? **(Check the appropriate boxes below, check only 1 box per row.)**

	Increase	Decrease	Stay the same	Not enough information
Average would ...				
Median would ...				
Standard Deviation would ...				

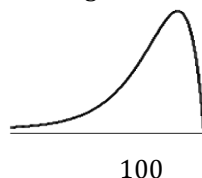
**Question 3**

Below are rough sketches of 2 histograms. One depicts scores on an Easy exam where most students did well. One depicts scores on a hard exam where most students did poorly. The horizontal axis ranges from 0% to 100%.

Histogram A



Histogram B



a) Which histogram depicts the easy exam?

**Choose one:**

- i) Histogram A
- ii) Histogram B

b) In Histogram A, is the average greater than, less than, or equal to the median? **Circle one:** > < =

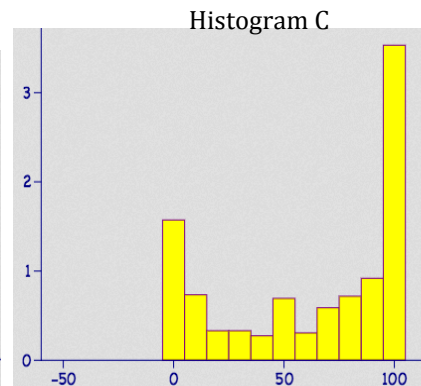
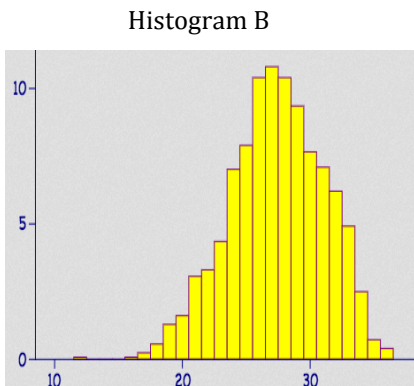
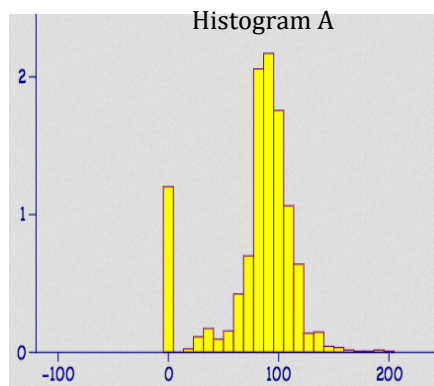
c) In Histogram B, is the average greater than, less than, or equal to the median? **Circle one:** > < =

**Question 4**

If a list of numbers has a SD of 0 then ....

- a) All the numbers on the list must be the same.
- b) The average of the numbers must be 0.
- c) All the numbers on the list must be 0.
- d) There are 0 numbers on the list since the SD can never be 0.

**Question 5 pertains to the 3 histograms below representing your survey responses to 3 questions:** What is your ACT score? What's the fastest speed you've ever driven (in mph)? and What percent of your college costs are your parents paying for?



a) Which graph represents ACT scores? \_\_\_\_\_ Which graph represents fastest speed? \_\_\_\_\_

b) I wrote the average and median of Histogram C down, but I forgot to label them. Here are the 2 numbers: 62.25 and 80. Which is which?

- i) 80 is the median
- ii) 80 is the average
- iii) Cannot be determined

## Chapter 5—Normal Approximation

**Question 1:** According to our survey data, the histogram for the heights of females in our class is close to the normal curve with an **average = 65 inches and a SD = 3 inches.**

a) If a female is below average in height, is her Z score positive or negative?

**Choose One:**

- i) Positive
- ii) Negative
- iii) Not enough information to tell

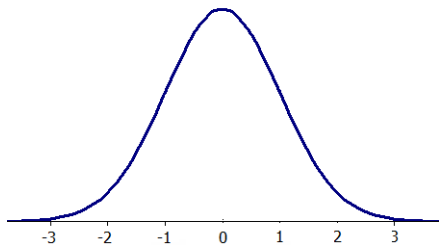
b) If a student is exactly at the 50<sup>th</sup> percentile, her Z score = \_\_\_\_\_ and she is \_\_\_\_\_ inches. (Fill in the 2 blanks above with numbers.)

c) What percent of the females are taller than 66.5 inches? (Use then normal curve provided at the end of this exam, you may round percen on the table to the nearest whole number.)

i) First convert 66.5" to a z score, show work.

Z =

ii) Then mark the Z score on the curve below and shade the area that represents everyone **over** 66.5".



Percent over 66.5" = \_\_\_\_\_%

Write your answer in the blank above.

d) Which of the following is closest to the percentage of females in the class who are between 62" and 68"?

**Choose One:**

- i) 68%
- ii) 82%
- iii) 91%
- iv) 95%

e) Which of the following is closest to the percentage of females in the class who are between 62" and 71"?

**Choose One:**

- i) 68%
- ii) 82%
- iii) 91%
- iv) 95%

f) About 50% of the females are between 63" and 67". Are there more or less females between 65" and 69"?

**Choose One:**

- i) More females are between 65" and 69" than between 63" and 67".
- ii) Less females are between 65" and 69" than between 63" and 67".
- iii) The 2 amounts are the same because the height difference is the same, 4" for both groups.
- iv) There is not enough information to tell.

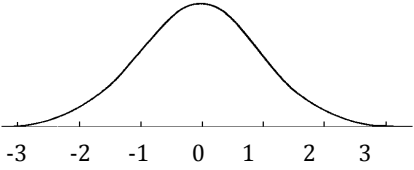
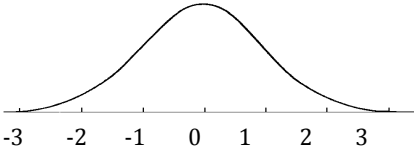
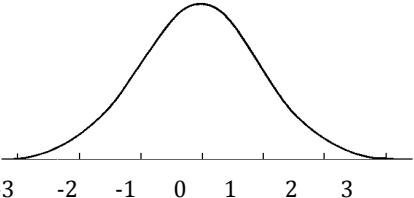
g) Suppose you found out that the heights were far from normally distributed but still had **average = 65"** and **SD = 3"**, would your answers in parts **d, e, f** above change or stay the same?

**Choose One:**

- i) The answers would be the same because the average and SD did not change.
- ii) The answers may change because the distribution is not normal and the table is therefore inaccurate.

**Question 2**

Suppose IQ scores follow the normal curve with an **average=100** and a **SD =16**. In the table below you're given either an IQ score, a Z score or percentile and you have to fill in the missing blanks. **For all these problems, please round the Areas given in the Normal Table to the nearest whole number.**

IQ	Z score	Percentile (% of people with lower IQ scores)
<p>a) Person A has IQ= 108</p>	<p>Z= _____</p> <p>Show work:</p>	<p><i>Person A is in the _____ percentile (Same as asking what % of the area is below z?)</i></p> <p>Mark Z score on curve, shade the area below Z and write the area below z in the blank above</p> 
<p>IQ = _____ (1pt) Do NOT round answer. Show work:</p>	<p>Person B has Z= 1.65</p>	<p><i>Person B is in the _____ percentile.</i></p> <p>Mark Z score on curve, and shade the area below Z and write the area below z in the blank above. (1 pt for correct curve)</p> 
<p>IQ = _____ (1pt) Do NOT round answer. Show work:</p>	<p>Z= _____ (1/2 pt)</p>	<p><i>Person C is in the 8<sup>th</sup> percentile</i></p> <p>What middle area should you look up on the normal table to find the correct Z score? _____% (Fill in blank)</p> <p>Mark the correct Z score on curve, and shade the area below Z.</p> 
<p>IQ= _____ Do NOT round answer. Show work:</p>	<p>Z= _____</p>	<p><b>Person D is in the 92<sup>nd</sup> percentile.</b></p> <p>No work is necessary. Just use the Z score you got for the 8<sup>th</sup> percentile to get the Z score for the 92<sup>nd</sup> percentile. (Hint: Notice how the 2 z scores are symmetrical on the curve?)</p>

**Part III Probability**

The next 6 questions pertain to randomly drawing from the box containing 5 tickets below.



- 1) Two tickets are drawn at random **with** replacement. What is the chance that both tickets shaded?  
 a)  $3/5 \times 2/4$     b)  $3/5 \times 3/5$     c)  $3/5$     d)  $1/5 \times 1/5$     e)  $2/5 \times 1/4$
- 2) Two tickets are drawn at random **without** replacement. What is the chance that both tickets are shaded?  
 a)  $3/5 \times 2/4$     b)  $3/5 \times 3/5$     c)  $3/5$     d)  $1/5 \times 1/5$     e)  $2/5 \times 1/4$
- 3) Five tickets are drawn at random **with** replacement. What is the chance of getting at least one shaded ticket?  
 a)  $1 - (3/5)^5$     b)  $(3/5)^5$     c)  $1 - (4/5)^5$     d)  $(4/5)^5$     e)  $1 - (2/5)^5$
- 4) One ticket is randomly drawn. What is the chance of getting either a shaded ticket or a ticket marked "3"?  
 a)  $2/5$     b)  $4/5$     c)  $3/5$     d)  $1$

The next 4 Questions pertain to rolling fair dice.

- 5) Two dice are rolled. What is the chance that the sum of the spots is 5?  
 i)  $2/36$     ii)  $3/36$     iii)  $4/36$     iv)  $5/36$     v)  $1/6 * 1/6$     vi)  $1/6 + 1/6$
- 6) One die is rolled 3 times. What is the chance of getting all 6's?  
 i)  $(5/6)^3$     ii)  $(1/6)^3$     iii)  $1 - (5/6)^3$     iv)  $1 - (1/6)^3$     v)  $3/6$
- 7) One die is rolled 3 times. What is the chance of not getting all 6's?  
 i)  $(5/6)^3$     ii)  $(1/6)^3$     iii)  $1 - (5/6)^3$     iv)  $1 - (1/6)^3$     v)  $3/6$
- 8) One die is rolled 3 times. What is the chance of getting at least one 6?  
 i)  $(5/6)^3$     ii)  $(1/6)^3$     iii)  $1 - (5/6)^3$     iv)  $1 - (1/6)^3$     v)  $3/6$
- 9) Two dice are rolled. What is the chance of getting a 3 on the first roll or a 4 on the second roll?  
 i)  $1/6$     ii)  $2/6$     iii)  $11/36$     iv)  $13/36$     v)  $1/6 * 1/6$     vi)  $1/6 * 1/6 + 1/36$

The next 2 questions refers to the following medical test:

A screening test for AIDs correctly gives positive results to about 99% of the people who have AIDs and incorrectly gives positive results to about 6% of the people who don't have AIDs. 1% of the population who take the test have AIDs.

Fill in the table below to give the results for 10,000 people.

	Tests Positive	Tests Negative	Total
Has AIDS			
Does Not have AIDS			
Total			10,000

- 10) What's  $P(\text{AIDS} | \text{negative test result})$ ?  
 a)  $99/100$     b)  $99/693$     c)  $9307/10,000$     d)  $1/9307$     e)  $6/100$
- 11) What's  $P(\text{AIDS} | \text{positive test result})$ ?  
 a)  $99/100$     b)  $99/693$     c)  $693/10,000$     d)  $1/9307$     e)  $6/100$

## Part IV: Statistics for Random Variables

### Chapters 6-9 Box Models, EV, SE and Histograms for Random Variables

Translating gambling games into Box models and computing the EV and SE for the sum, average and % of n draws from a box.

- EV of the sum of n draws from a box = n times the average of the box
- Know the 3 SE formulas:

Remember SE = SD either multiplied or divided by  $\sqrt{n}$  (multiply SD by  $\sqrt{n}$  only for SE of sum)

- SE of the sum of n draws from a box =  $SD_{Box} * \sqrt{n}$
- SE of the average of n draws from a box =  $SD_{Box} / \sqrt{n}$
- SE of the % of 1's in n draws from a 0-1 box =  $SD_{Box} / \sqrt{n} (* 100) \%$

(Multiply by 100 to change from a decimal to a percent, for example  $0.1 \times 100 = 10\%$ )

- Know the short-cut formula for the SD of boxes that just have 2 types of tickets on page 50  
If the box has only 1's and 0's this is the same as:

$SD = \sqrt{p*(1-p)}$  where p is the proportion (fraction) of 1's in a 1-0 box.

- Central Limit Theorem—The probability histogram for all possible sums (or averages, or percents) of draws from a box will get closer and closer to the normal curve.
- With enough draws we can use the normal curve to figure the chance that the sum (or average or percent) of the draws will fall within a given range by converting the endpoints of the interval into a Z score  
 $Z = (\text{Value} - \text{Expected Value}) / \text{SE}$

#### Question 1 pertain to the following situation:

A 100 question multiple-choice test awards 4 points for each correct answer and subtracts 1 point for each incorrect answer. Each question has 5 choices.

i) Suppose a student guesses at random on each question, what is the corresponding box model?

- It has two tickets: 1 and 0
- It has 100 tickets: half 1's and half -1's
- It has five tickets: 1, 0, 0, 0, 0
- It has five tickets: 4, 0, 0, 0, 0
- It has five tickets: 4, -1, -1, -1, -1

ii) The expected value for the student's score is

- a) 0                      b) 10                      c) 20                      d) 40                      e) 50

iii) The standard error of the student's score is

- a) 20                      b) .4                      c) 2                      d) .2                      e) not enough info

iv) Now suppose you're just interested in how many correct answers the student would get by guessing, not his score. Then the EV = 20 and the SE = 4. Suppose the student needs to get 27 answers correct in order to pass. What's the probability the student will pass? (Hint: convert to a Z score, and use the normal curve).

- a) 2%                      b) 4%                      c) 8%                      d) 10%                      e) 20%

**Question 2**

A slacker student has 4 Finals. Each Final consists of 100 multiple-choice questions. He knows nothing so he decides to randomly guess on every question so he can complete each Final in less than 5 minutes.

i) To compute the Expected Value (EV) for the student’s score for each Final, you may need additional information. Which of the following do you need to know? **Circle “Yes”** if needed or **“No”** if not.

- a) How many students are taking each final. **Circle one:** Yes No
- b) How many choices there are for each question. **Circle one:** Yes No
- c) How many points are awarded or deducted for each choice. **Circle one:** Yes No
- d) How much time is allotted for the exam. **Circle one:** Yes No

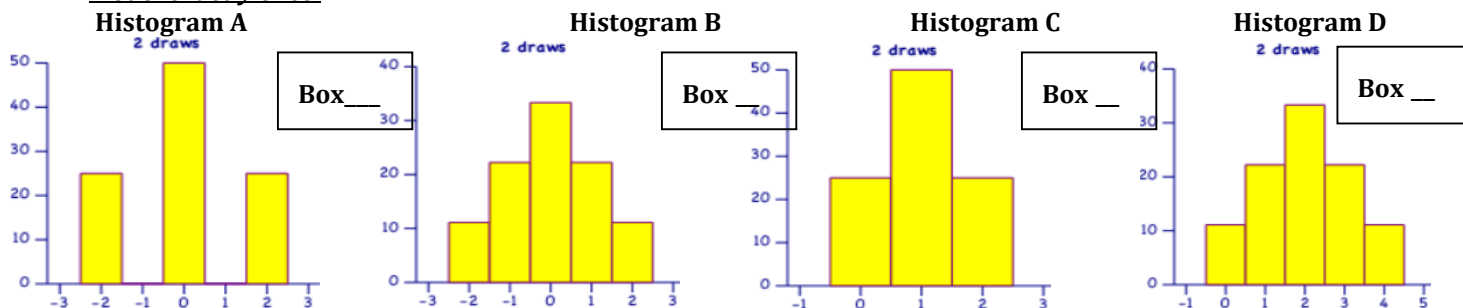
ii) Randomly guessing on all 100 questions corresponds to drawing \_\_\_\_\_times \_\_\_\_\_replacement from the appropriate box model. **(Fill in the first blank with a number and the second with either “with” or “without”.)**

iii) For a-d match the Final exams to their corresponding box models **Use each box model exactly once.**

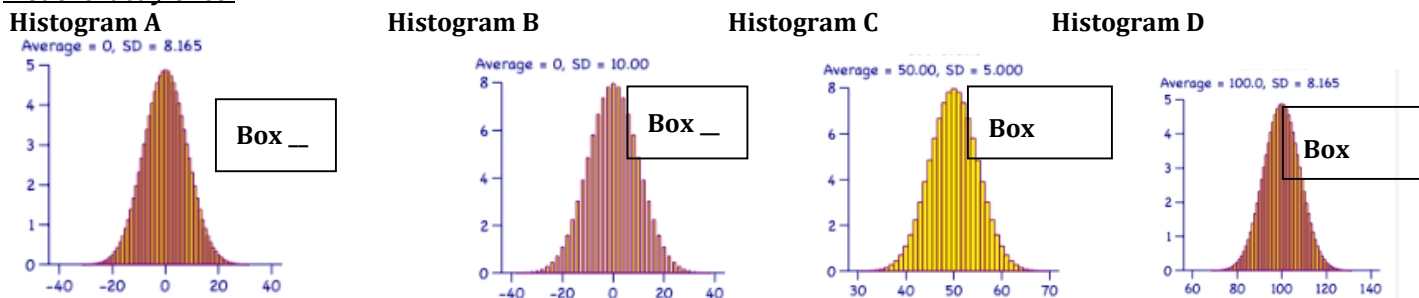
Box I:         Box II:         Box III:        Box IV:

- a) Final A- Each question has 3 choices, one is a right answer, one is a wrong answer and one is an “ I don’t know” answer. Your score is computed as the number of right answers minus the number of wrong answers. The “I don’t know” answers are scored as 0 points. This corresponds to Box... **i) I ii) II iii) III iv) IV**
- b) Final B- Each question has 3 choices, one is the best answer and awarded 2 pts, one is a mediocre answer and awarded 1 pt. and one is a wrong answer and awarded no points. This corresponds to Box... **i) I ii) II iii) III iv) IV**
- c) Final C--Each question is a true/false question. Your score is the number of answers you get right. This corresponds to Box... **i) I ii) II iii) III iv) IV**
- d) Final D-Each question is a true/false question. Your score is the number of answers you get right minus the number of answers you get wrong. This corresponds to Box... **i) I ii) II iii) III iv) IV**

iv) The 4 histograms below represent the probability histogram for the **sum of 2** draws made at random with replacement from each of the boxes in part (iii) above. For each histogram identify the appropriate Box (I, II, III or IV). **Use each box model exactly once.**



v) The 4 histograms below represent the probability histogram for the **sum of 100** draws made at random with replacement from each of the boxes in part (iii) above. For each histogram identify the appropriate Box (I, II, III or IV). **Use each box model exactly once.**

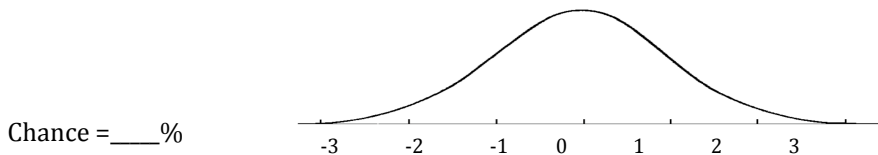


**HINT—The Average and the SD given above each histogram is the EV and the SE of the sum of 100 draws.**

**Question 3**

64 draws are made at random with replacement from the box containing 4 tickets:  $\boxed{2} \boxed{4} \boxed{4} \boxed{10}$

- a) The **smallest** the sum of the 64 draws could possibly be is \_\_\_\_\_ and the **largest** is \_\_\_\_\_.  
*(Fill in the 2 blanks above with the correct numbers.)*
- b) What is the **EV** (expected value) of the **sum** of the **64** draws? (Show work, circle answer.)
- c) What is the **SE** (Standard Error) of the **sum** of the **64** draws? (**SD of box = 3**) (Show work, circle answer.)
- d) Use the normal approximation and to find the **chance** that the sum of **100** draws will be **below 455**? The EVsum= 500 and the SEsum= 30 for 100 draws.
  - i) First calculate the Z score. **Show work. Circle answer.**
  - ii) Now mark the Z score accurately and **shade the area that represents the chance of getting below 455**  
*Round the middle area given in the table to the nearest whole number.*



- e) What is the **EV** of the **average** of the **100** draws? \_\_\_\_\_ (no work is necessary)
- f) What is the **SE** of the **average** of the **100** draws? \_\_\_\_\_ (Show work.)
- g) Now suppose you draw at random with replacement from the same box above, but this time you're only interested in the percent of 4's you get. What is the **EV** and the **SE** of the **percent** of 4's in **100** draws? (**Hint: draw a new box**)
  - i) EV of the **percent** of 4's in **100** draws = \_\_\_\_\_ (no work necessary)
  - ii) SE of the **percent** of 4's in **100** draws = \_\_\_\_\_



**Part V: Sampling and Inference Chapter 10-11**

**Sample Surveys—**

**Random Samples are best for the same 2 reasons that randomized experiments are best:**

1. **They eliminate selection bias**
2. **They can be translated into box models so you can attach SE's to your estimates.**

**Box Model for Sample Surveys:**

- The box has 1 ticket for every person in the population.
- A random sample of n tickets is drawn from the box without replacement (because you don't want to sample the same person twice).
- You know the average or percent of your sample and you use it to estimate the average or percent in the whole population.
- Of course, the average or percent in your sample won't be *exactly* the same as that of the population, because of chance error (samples will vary because of the luck of the draw). As long as the sample size is big enough, the probability histogram for the sample average and percent will follow the normal curve so we can attach SE's to our estimates and build confidence intervals.
- For **small** samples from approximately **normal** populations with **unknown SD**, the probability histogram of the sample average (**not percent**) will follow the **t** distribution, so we can improve our estimates by using the t curves to attach SE's to our estimates to build confidence intervals.

**Note: The size of the population doesn't affect the accuracy of our estimates, only the size of the sample matters. The bigger our sample size, the smaller the SE for averages and percents (smaller by a factor of the square root n).**

**This is apparent in the SE formulas for sample averages and percents because we divide the SD by  $\sqrt{n}$ , where n is the sample size (not the population size)**

**Sample Questions:**

1) City A has **1 million** people and City B has **9 million** people. A simple random sample of **1000** people is taken from City A and a simple random sample of **9000** is taken from City B. Other things being equal the sample from City A is \_\_\_\_\_ the sample from city B.

- a) 9 times *more* accurate   b) 3 times *more* accurate   c) the same accuracy as   d) 9 times *less* accurate   e) 3 times *less* accurate

2) A recent Pew Research Center Poll asked a random sample of 1,211 adults nationwide the following question: "Do you think a woman should be able to get an abortion if she decides she wants one no matter what the reason."

We posted the same question on last semester's Bonus Survey. Here's the results of both surveys:

	Yes	No	Sample Size
Pew Research Center	18%	82%	1211
Bonus Survey	46%	54%	631

a) As you can see, the results of the 2 polls are quite different. Which survey gives a better estimate of the percentage of all US adults who would answer "yes" to this question? **Choose one:**

- i) The Pew Research survey because the sample size was larger.
- ii) The Bonus Survey because we can be sure it was an anonymous survey.
- iii) The Pew Research survey because the people were **randomly** drawn from all adults nation-wide.

b) What is SE of the sample percent for the Pew Poll? **Choose one:**

- i) It's not possible to calculate a SE for this sample because we don't know the SD of the sample.
- ii) It's not possible to calculate a SE for this sample because we don't know the size of the population.
- iii) The SE of the sample percent is approximately 13.4%
- iv) The SE of the sample percent is approximately 1.1%

3) A recent Gallup poll asked a simple random sample of 900 adults nationwide how much they spent on Black Friday. The sample average was \$400 with a SD of \$300.

a) What most closely resembles the relevant box model?

- i) It has 900 tickets marked with "0"s and "1"s.
- ii) It has about millions of tickets marked with "0"s and "1"s.
- iii) It millions of tickets. On each ticket is written a \$ amount. The exact average and SD are unknown but are estimated from the sample.
- iv) It has 900 tickets. The average of the tickets is \$400 and the SD is \$300.

b) 900 draws are made \_\_\_\_\_replacement.

**Choose one:** i) With                      ii) Without

c) What is the SE of the sample average?

- i) \$100              ii) \$10              iii) \$3              iv) \$0.33              v) \$30,000.

d) A 92% CI for the true population average = \$ \_\_\_\_\_ ± \_\_\_\_\_ \* SE.

Fill in the 2 blanks with the correct numbers. (Hint: Use the normal table for the second blank.)

e) To which of the following populations would the above 92% confidence interval apply?

- a) All US females
- b) All US adults
- c) All Illinois adults
- d) All middle class US adults
- e) All of the above

f) How would a 99% CI compare to the 92% CI we calculated in part d?

- a) It would be wider    b) It would be narrower    c) It would be the same.

g) Suppose we wanted to use  $SE^+$ , instead of SE to calculate our CI's, you'd multiply your answer in part c above by

- i)  $\sqrt{\frac{900}{899}}$     ii)  $\sqrt{\frac{899}{900}}$     iii)  $\sqrt{\frac{300}{299}}$     iv)  $\sqrt{\frac{400}{399}}$

v) None of the above because you cannot use  $SE^+$  to calculate a CI if you're using the Normal Curve.

h) Suppose we had a small sample size ( $n < 25$ ) with the same sample average and SD as above .

Should we use the t curves to compute Confidence Intervals?

- i) Yes, because we know the SD of the sample.
- ii) Yes, because we don't know the SD of the population.
- iii) No, because judging from our sample average and SD it's highly unlikely that it comes from a normal population.

4) A CBS News Poll asked a random sample of 1,600 adults nationwide the following question: "Do you think the distribution of money and wealth in this country is fair or do you think wealth should be more evenly distributed among more people?" 26% answered "Fair"

- a) What most closely resembles the relevant box model?  
a) It has 1600 tickets, 26% are marked "1" and 74% are marked "0"  
b) It has 1600 tickets with an average of 0.  
c) It has millions of tickets marked "0" and "1", but the exact percentage of each is unknown and estimated from the sample.
- b) The draws are made \_\_\_\_\_ replacement. a) With b) Without
- c) Which one of the statements below is true?  
a) The expected value for the percent of registered Democrats who would answer "Fair" to the question is 26%.  
b) The expected value for the percent of corporation executives who would answer "Fair" to the question is 26%.  
c) The expected value for the percent of Chicago residents who would answer "Fair" to the question is 26%.  
d) All of the above are true.  
e) None of the above are true.
- d) Is it possible to compute a 95% confidence interval for the percent of all US adults who would answer "Fair" to the question?  
i) Yes, a 95% confidence interval is approximately 26% +/- 1.1%  
ii) Yes, a 95% confidence interval is approximately 26% +/- 2.2%  
iii) No, because we're not given the SD of the sample.  
iv) No, because we cannot infer with 95% confidence the answers of 200 million Americans from data based on a sample of only 1,600 randomly selected Americans.
- e) If 1000 people all took random samples of 1600 and computed 95% CI's, about how many of their intervals would capture the true population percent?  
i) All of them ii) 9999 iii) 995 iii) 950 iv) 10 v) 50 vi) 100 vii) Impossible to predict.
- f) If the researcher decreased his sample size by a factor of 4 (to n=400) then the width of the 95% confidence interval would ...  
i) increase by a factor of 2 ii) increase by a factor of 4 iii) decrease by a factor of 2 iv) decrease by a factor of 4
- g) If our sample size was small (n < 25) would it be appropriate to use the t curves instead of the Normal curve to compute CI's?  
i) Yes ii) No, it's never appropriate to use the t curves with 0-1 data
- 5) To estimate the average IQ of students at a large public high school of 2000 students, a random sample of 17 students is taken. The sample average = 102 with a SD =16. Compute a 95% CI using the t distribution.
- a) SE+ = \_\_\_\_\_
- b) How many degrees of freedom? \_\_\_\_\_
- c) What is the t\* (the critical value of t)? \_\_\_\_\_ (use the t-table in your notes)
- d) 95% CI = (\_\_\_\_\_ to \_\_\_\_\_) Put the lower number first.

### Choosing how many people to poll

- 6) In a pre-election poll in a close race, how many people would you have to poll to get... (Assume SD= 0.5)
- a) a 95% CI with a 3% margin of error?
  - b) an 80% CI with a 3% margin of error?
  - c) Let's say the SD = 0.4, would we need more or less people than we did assuming the SD=0.5 ?  
i) More ii) Less iii) Same

**Part VI Significance Tests – are statistical checks to decide whether some difference we observe is “real” (due to some particular cause) is just due to chance variation.**

### Chapters 12-The one sample Z Test

$$Z \text{ test-statistic} = \frac{\text{Observed} - \text{Expected}}{\text{SE}}$$

Look at the sampling distribution of Z under the null and see how likely it would be to get our data or something even more extreme if the null were true. That's called the p-value.

The convention is to reject the null when  $p < 5\%$  and call the result “statistically significant” and when  $p < 1\%$  call the result “highly significant”. There's no particular justification for those values. In other words, a p-value of 4.9% isn't really much different than a p-value of 5.1%, people just like to draw the line somewhere.

1) Ellen thinks she has no musical ability but Karle thinks she does. To find out Ellen took a musical memory test online that had 36 questions. For each question she had to choose whether a sequence of notes were the same or different. She answered 24 of the 36 questions correctly. The null hypothesis is that she was just guessing.

- a) Which of the following most accurately describes the null box?
- i) It has 36 tickets, 24 marked "1" and 12 marked "0"
  - ii) It has 36 tickets marked either "1" or "0" but the exact percentage of each is unknown.
  - iii) It has 2 tickets, 1 marked "1" and 1 marked "0"

- b) The draws are made \_\_\_\_\_ replacement.                      i) with                      ii) without

Assuming the null hypothesis to be true, you would expect Ellen to answer \_\_\_\_ questions correct, give or take \_\_\_\_ questions.

- c) Fill in the first blank in the above sentence with the correct expected value.

- i) 12                      ii) 18                      iii) 21                      iv) 24                      v) 18

- d) Fill in the second blank in the above sentence with the correct SE.

- i) 1                      ii) 2                      iii) 3                      iv) 4                      v) 5

- e) The Z -statistic for testing the null hypothesis is

- i) 6/SE for the average                      ii) 6/ SE for the sum                      iii) 7/SE for sum                      iv) 6/SD of the box

- f) The p-value for the one-sided alternative is ...

- i) 2.5%                      ii) 5%                      iii) 16%                      iv) 21%                      v) 11.5%

- g) Suppose our sample size was  $< 25$  would it be appropriate to use a t-test here?

- a) Yes    b) No

2) An internet access company that serves millions of customers claims that it takes an average of only 1.8 attempts to connect with their service, but customers think it takes more. To test the company's claim, a consumer advocate looked at a random sample of 400 connections and recorded the number of attempts required to establish each connection. The average of the 400 observations is 2.1 and the SD is 5.0.

a) What is the null hypothesis?

- i)  $\mu = 1.8$  ii)  $\mu > 1.8$  iii)  $\mu \neq 1.8$  iv)  $\bar{x} = 1.8$  v)  $\bar{x} > 1.8$  vi)  $\bar{x} \neq 1.8$

b) What is the alternative hypothesis?

- i)  $\mu = 1.8$  ii)  $\mu > 1.8$  iii)  $\mu \neq 1.8$  iv)  $\bar{x} = 1.8$  v)  $\bar{x} > 1.8$  vi)  $\bar{x} \neq 1.8$

c) The null hypothesis box is best described as:

- i) containing millions of tickets, each marked 1 or 0, where 1 denotes that a connection was made.  
 ii) containing 400 tickets, each marked 1 or 0, where 1 denotes that a connection was made.  
 iii) containing millions of tickets with whole number values such as 1, 3, 5, 2, ...  
 iv) containing 400 tickets with whole number values such as 1, 3, 5, 2 ...

d) The average of the null hypothesis box is: a) 1.8 b) 2.1

e) The SE of the sample average is closest to:  
 a) 0.05 b) 0.25 c) 0.50 d) 5.0 e) 20.0

f) The Z-statistic is closest to:  
 a) 0.15 b) 0.12 c) 0.6 d) 1.2 e) 6.0

g) The p-value is closest to: a) 77% b) 23% c) 11.5%

h) The critical value ( $Z^*$ ) to reject the null (for a one-sided test) at significance level  $\alpha = 0.05$  is closest to ...  
 a) 1 b) 1.3 c) 1.65 d) 2 e) 2.5

g) Conclusion

- a) Reject the null, there is very strong evidence that the company's claim is false and the average number of attempts is greater than 1.8  
 b) Cannot reject the null, it's reasonable to think that observed difference could be simply due to chance.

i) Would it be wrong to do a t-test here?

- a) Yes, because the sampling distribution of the mean never follows a t-distribution when  $n=400$ .  
 b) Not wrong, just unnecessary, because the t-distribution is extremely close to the normal distribution when  $n=400$ .

### Chapter 13: The t test

We use the SE+ and the t distribution when we have:

1.  $\sigma$ , the SD of the population (null box), is unknown, all you know is the SD of the observed sample.
2. The population (contents of the box) roughly follows the normal curve
3. A small sample  $n \leq 25$  (You can use t with a larger sample but the difference between t and z becomes negligible as n gets larger.)

This means you NEVER use the t test when the population (null box) is 1's and 0's since the population isn't normal and  $\sigma$  is tied to the sample percent so it's not completely unknown.

When the sample size is small, using the sample SD to estimate the SD of the box is not very accurate. It's likely to be too low so we use  $SD^+ = \frac{\sqrt{n}}{\sqrt{n-1}} \times SD$  instead.  $SD^+ > SD$  but the difference becomes negligible as n gets large.

$$t - \text{statistic} = \frac{\text{Observed avg} - \text{Expected avg}}{SE_{\text{avg}}^+} \quad \text{where } SE_{\text{avg}}^+ = SE_{\text{avg}} = \frac{SD^+}{\sqrt{n}} = \frac{SD}{\sqrt{n-1}}$$

When the null is true the sample distribution of the t statistic follows the t curve with **n-1 degrees of freedom**.

1) A factory that packages corn flakes is supposed to put the flakes in the boxes so that the boxes weigh an average of 16 ounces and a standard deviation of 1 ounce. An inspector randomly chose 12 boxes from one day's output of 2500 boxes. These 12 had an average weight of 15 ounces. The inspector wishes to test the null hypothesis that the factory is doing what it is supposed to on this day.

- a. Which of the following best describes the null box?
- i) The box has 12 tickets, with an average of  $180/12 = 15$  ounces.
  - ii) The box has 12 tickets, with an average of 16 ounces.
  - iii) The box has 2500 tickets, but we do not know exactly the average.
  - iv) The box has 2500 tickets, with 16% 1's and 84% 0's.
  - v) The box has 2500 tickets, with an average of 16 ounces.

- b. The SE for the average of the draws is closest to
- i) 0.367
  - ii) 0.288
  - iii) 3.46
  - iv) 4
  - v) .02

- c. What test statistic would you use?
- i) z-statistic
  - ii) t-statistic

- d. The test statistic is -3.47. What conclusion do you draw?
- i) Accept the null hypothesis.
  - ii) There is not enough evidence to suspect there is anything wrong.
  - iii) Reject the null hypothesis, there is strong evidence that the factory is not doing what it is supposed to.
  - iv) The p-value is larger than 5%.

2) Now suppose the factory makes the same claim as above, that the boxes weigh 16 ounces on the average, but the factory doesn't make any claim about the SD. Instead, the inspector computes the SD of the 12 boxes and finds the SD = 1 ounce

- a) What is the best estimate of the SD of the 2500 boxes?
- i) 1 ounce
  - ii) 1.049 ounces
  - iii) 1.4 ounces

- b) What test statistic should the inspector now use?
- i) z-statistic
  - ii) t-statistic

- c) If he decides to use the t-statistic, how many degrees of freedom are there?
- i) 2499
  - ii) 12
  - iii) 11
  - iv) 6

- d) What is the value of the t-statistic?
- i) -3.3
  - ii) -3.47
  - iii) -3.9

- e) Which test yields a larger p-value for the same data, the t-test or the z-test?
- i) t test
  - ii) Z test
  - iii) they always yield exactly the same p-value

**Chapter 14 and 15**— 2 sample Z and t-tests used to compare averages of 2 populations. Only the Z test can be used to compare to population percents.

Same conditions hold for 2 sample t-tests as for 1 sample t-tests:

- Unknown SDs in the population,
- Close to normal populations
- Small sample sizes. (It's not wrong to use t for larger samples, it's just not needed usually. )

**H<sub>0</sub> : 2 populations have the SAME average or percent (only use Z for percent);**  
**2-sided H<sub>a</sub> is that they're not the same, 1-sided H<sub>a</sub> specifies which is larger.**

$$Z \text{ stat} = \frac{\text{Observed difference in the 2 samples} - \text{Expected difference}}{SE_{\text{difference}}} = \frac{\text{Observed Difference}}{SE_{\text{difference}}} \text{ since Exp diff}=0$$

The t stat is the same except it uses SE<sup>+</sup> instead of SE and df= n-1 where n is the size of the smaller sample.

SE difference is the square root of the sum of the squares of each sample's SE (or SE+)

For 2 independent samples one randomly drawn from population 1 and the other from population 2 the SE of their difference is:

$$SE_{\text{diff}} = \sqrt{SE_1^2 + SE_2^2} = \sqrt{\left(\frac{SD_1}{\sqrt{n_1}}\right)^2 + \left(\frac{SD_2}{\sqrt{n_2}}\right)^2} = \sqrt{\frac{SD_1^2}{n_1} + \frac{SD_2^2}{n_2}}$$

$$SE_{\text{diff}}^+ = \sqrt{(SE_1^+)^2 + (SE_2^+)^2} = \sqrt{\left(\frac{SD_1}{\sqrt{n_1 - 1}}\right)^2 + \left(\frac{SD_2}{\sqrt{n_2 - 1}}\right)^2} = \sqrt{\frac{SD_1^2}{n_1 - 1} + \frac{SD_2^2}{n_2 - 1}}$$

1) A study to see if teenage girls spend more time on their phones than teenage boys do took a nation-wide random sample of 25 girls and 20 boys and found the following:

	Girls	Boys
Average hours per day	5.5 hrs.	4.5 hrs.
SD	2 hrs.	2 hrs.

a) The null hypothesis is that the average phone time per day for girls \_\_\_\_\_ average phone time for boys in the \_\_\_\_\_. Fill in the first blank with i) = ii) > iii) < and the second blank with i) sample ii) population.

b) Which of the following most accurately describes the null box(es)?

- i) There is one null box with 45 tickets, 25 marked "1" and 24 marked "0"
- ii) There is one null box with millions of tickets each marked with the amount of phone hours per day
- iii) There are 2 null boxes, each with millions of tickets. One box has an average =5 and the other has average =4
- iv) There are 2 null boxes, each with millions of tickets. On each ticket is written the amount of phones hours per day. The 2 boxes have the same average.
- v) There are 2 null boxes, each with millions of tickets marked "0" and "1".

c) First we'll do a Z-test even though the samples are relatively small. The SE of the difference of the 2 sample averages is  
 i) 0.42 ii) 0.18 iii) 0.3 iv) 0.36 v) 0.6

d) The Z statistic for testing the null hypothesis is closest to  
 i) 0 ii) 1 iii) 1.63 iv) 1.67 v) 1.71

- e) The p-value is 4.75%. If the significance level is set at 5%, we would
- i) Reject the null and conclude girls spend more time on the phone than boys do 95.25% of the time.
  - ii) Reject the null and conclude that if the average phone time were the same for girls and boys in the population, the probability that we'd see a one hour higher average or more for girls in our sample is less than 5%.
  - iii) There's good evidence that there is no difference between the amount of time boys and girls spend on their phones in the population.

f) Suppose we had chosen a 2-sided alternative hypothesis at the start of the problem. What would be our p-value?

2) Suppose you wanted to use SE+ and the t-test instead of the Z test to test the same null and alternative as in (1) above.

$$H_0 : \mu_{\text{Girls}} = \mu_{\text{Boys}} \quad H_a : \mu_{\text{Girls}} > \mu_{\text{Boys}}$$

a)  $SE_{\text{diff}}^+ =$

i)  $\sqrt{\frac{2}{20} + \frac{2}{25}}$     ii)  $\sqrt{\frac{4}{20} + \frac{4}{25}}$     iii)  $\sqrt{\frac{2}{19} + \frac{2}{24}}$     iv)  $\sqrt{\frac{4}{19} + \frac{4}{24}}$     v)  $\sqrt{\frac{4}{19^2} + \frac{4}{24^2}}$

b) The t statistic for testing the null hypothesis is closest to    i) 0    ii) 1    iii) 1.63    iv) 1.67    v) 1.71

c) To find the p-value you'd look at the t curve with \_\_\_\_\_ df.

d) The critical value t\* for rejecting the null at  $\alpha = 0.05$  for a one sided alternative = \_\_\_\_\_.

e) Your t stat \_\_\_\_\_ t\* , so you \_\_\_\_\_ the null.

Fill in the first blank with > or < and the second blank with reject or cannot reject.

f) The p-value using the t test is \_\_\_\_\_ the p-value using the Z test.

i) >    ii) <    iii) =    iv) Not enough info



3) Gallup asked a random sample of 400 men and 400 women nationwide the following question: "If you were taking a new job and had your choice of a boss, would you prefer to work for a man or a woman?"

$H_0$ : % of all US women who would prefer a male boss = % of all US men who would prefer a male boss.

$H_a$ : % of all US women who would prefer a male boss  $\neq$  % of all US men who would prefer a male boss.

In our sample we found 50% of the women and 45% of the men said they would prefer a male boss.

a) Which of the following most accurately describes the null box(es)?

- i) There is one null box with 800 tickets, marked with "0"s and "1"s
- ii) There is one null box with millions of tickets, marked with "0"s and "1"s
- iii) There are 2 null boxes, each with millions of tickets. One box has 45% "1"s and 55% "0"s and the other has 50% "1"s and 50% "0"s
- iv) There are 2 null boxes, each with millions of tickets. The 2 boxes have the same percentage of "1"s and "0"s.

b) The SE for the 2 sample percentages are both about 2.5%.

The SE for the difference of the 2 sample percentages is closest to

- a) 2.5%
- b) 0%
- c) 5%
- d) 3.5%

c) The p-value for testing the null hypothesis is closest to

- a) 0%
- b) 2%
- c) 8%
- d) 16%
- e) 84%

## Chapter 16

### Part I-The Chi-Square Goodness-of-Fit Test

Used to decide whether the observed data fits a specified model when the model has more than 2 categories.

With 2 categories (0-1 box) we use the one sample z test.

Null Hypothesis: The observed data fits the model "good". (The difference between the observed and expected is just due to chance.)

Alternative Hypothesis: The observed data does NOT fit the model "good". (The difference between the observed and expected are too big to be due to chance.)

Chi-Square Statistic =  $\text{sum of (observed frequency - expected frequency)}^2 / \text{expected frequency}$

Degrees of freedom = # of categories - 1

### Part II- The Chi-Square Independence Test

Use to compare the percent composition of 2 or more variable when each variable has 2 or more categories. With 2 variables and 2 categories you can use either a 2 sample z-test or a chi-sq ind test.

(You can think of the Chi-Square Goodness-of-fit Test as a 1 sample test, comparing the sample percents to a null box that has multiple categories and you can think of the Chi-Square Independence Test as a 2 sample test, comparing the percent composition of 2 populations when each population has multiple categories.)

Null Hypothesis: The 2 variables are independent. (The 2 populations have the SAME percent composition; the difference between observed and expected frequencies are just due to chance.)

Alternative Hypothesis: The 2 variables are dependent. (The 2 populations have different percent compositions; the difference between observed and expected are too big to be due to chance.)

Chi-Square Statistic =  $\text{sum of (observed frequency - expected frequency)}^2 / \text{expected frequency}$

Degrees of freedom = (# of rows - 1) x (# of columns - 1)

\*To figure the expected frequency for each cell: multiply the row total x column total/overall total

1) A certain University has 10% freshman, 20% sophomores, 30% juniors and 40% seniors. A group of 200 students are chosen for a survey. The group has 30 freshman, 40 sophomores, 60 juniors and 70 seniors. The null hypothesis is the students were chosen at random.

	Expected Percents	Observed #	Expected #
Freshman	10%	30	
Sophomores	20%	40	
Juniors	30%	60	60
Seniors	40%	70	80
Total	100%	200	200

- a) To test the null hypothesis that the students were chosen at random we'd do  
 the chi-square test for "goodness -of-fit"  
 the chi-square test for independence  
 the one-sample z test  
 the two-sample z test

The table above is missing 3 values. Fill in the missing values by answering the following 3 questions:

- b) What is the **expected** number of *freshman*?  
 i) 10 ii) 20 iii) 30 iv) 40 v) 50

- c) What is the **expected** number of *sophomores*?  
 i) 10 ii) 20 iii) 30 iv) 40 v) 50

- d) To compute the proper test statistic you'd sum 4 terms:  $5 + 0 + 0 + \underline{\quad}$ . The term for **seniors** is missing, what should it be?  
 i) 0 ii) 1 iii) 1.25 iv) 1.43 v) 2.5

- e) The number of degrees of freedom is  
 i) 2 ii) 3 iii) 4 iv) 5 v) 6

- f) What do you conclude?  
 i) Reject the null because  $p < 5\%$  ii) Reject the null because  $p > 5\%$  iii) Cannot reject the null because  $p < 5\%$

2) A simple random sample of 148 Stat 100 students were asked whether or not they thought they would ever use statistics again in their lives. Assume the students were chosen from a population of 2000. The following table gives the results:

	Would use	Would not use
Men	47	21
Women	64	16

The chi-square statistic to test the null hypothesis that sex and anticipated use are independent is 2.32.

- a. To compute this statistic, expected frequencies were calculated. What is the expected frequency for the men who answer "would use"?

a) 51                      b) 47                      c) 44

- b. How many degree of freedom does the chi-square statistic have?

a) 1                      b) 2                      c) 3                      d) 4

- c. Can you reject the null hypothesis?

a) Yes                      b) No

- d. If we had done a 2 sample z test with a 2 sided  $H_a$ , would we have gotten the same exact p-value?

i) Yes ii) No, we would have gotten half the p-value iii) No, we would have gotten twice the p-value

**3)** The table below shows the results of a recent nationwide poll of Hispanic adults who were asked; "All in all, do you think the situation for the younger generation of Hispanic or Latino Americans is better, worse, or about the same as their parents' situation was when they were the same age?" You may assume that the data are from a simple random sample of 200 people, of whom 100 were over 35 years old and 100 were 18-34 years old.

	18-34	Over 35
Better	49%	39%
Worse	37%	45%
About the same	14%	16%

To answer the question of whether the answers are really different for young and old adults, you use

- i) the one-sample z test
- ii) the two-sample z test
- iii) the chi-square test for "goodness-of-fit" which specifies the contents of the box
- iv) the chi-square test for independence

**Chapter 17--**Significance tests can only tell you whether or not a difference is likely to be due to chance, not whether a difference was important or what caused the difference, or whether the experiment was properly designed

By definition, significant Results will appear by chance with enough tests. A p-value of 5% means that even when the null is true, you'll reject it 5% of the time.

**1)** Which of the following does a test of significance deal with?

- a. Is the difference due to chance?
- b. Is the difference important?
- c. Was the experiment properly designed?
- d. What are the probable causes of the difference?

**2)** 100 investigators each set out to test a different null hypothesis. Unknown to them, all the null hypotheses happen to be true.

**a.** About how many of them would you expect to get statistically significant results?

- i. None, if they did the test correctly they would all confirm that the null hypothesis is true.
- ii. 1
- iii. 5
- iv. 95
- v. Impossible to predict.

**b.** About how many of them would you expect to get highly statistically significant results?

- i. None, if they did the test correctly they would all confirm that the null hypothesis is true.
- ii. 1
- iii. 5
- iv. 95
- v. Impossible to predict.

3) A significance test is performed to analyze the results of a randomized experiment to determine whether students learn more or less from watching a lecture online compared to attending the same lecture in person. Subjects are randomly assigned to treatment (online lecture) and control (in person lecture) and then given the same exam afterwards.

a) What are the null and alternative hypotheses?

$H_0$ : Choose one: i)  $\mu_t - \mu_c = 0$  ii)  $\mu_t - \mu_c > 0$  iii)  $\mu_t - \mu_c < 0$  iv)  $\mu_t - \mu_c \neq 0$

$H_a$ : Choose one: i)  $\mu_t - \mu_c = 0$  ii)  $\mu_t - \mu_c > 0$  iii)  $\mu_t - \mu_c < 0$  iv)  $\mu_t - \mu_c \neq 0$

b) A significance level of  $\alpha = 0.02$  means when the null is true the probability of making a Type I error = \_\_\_\_\_

**Circle one:** i) 0% ii) 1% iii) 2% iv) 4% v) 96% vi) 98% vii) not enough info

(A Type I error is rejecting the null when the null is true.)

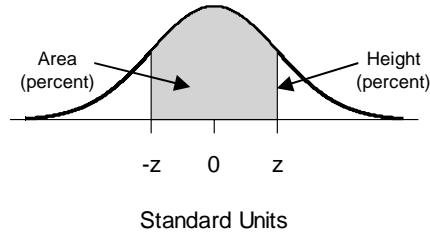
c) If we set  $\alpha = 0.05$  (null cut-off at 5%) for a 2-sided  $H_a$  then the critical value of our test-statistic,  $Z^* =$  \_\_\_\_\_

Choose closest answer. i) 0.85 ii) 1.3 iii) 1.65 iv) 2 v) 2.35 vi) 2.6

d) Repeat (c) above with a 1-sided  $H_a$  keeping all else the same. Choose closest answer.

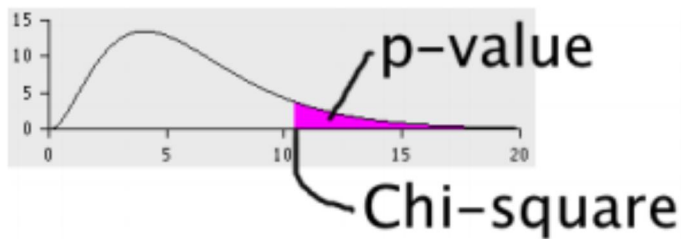
i) 0.85 ii) 1.3 iii) 1.65 iv) 2 v) 2.35 vi) 2.6

## STANDARD NORMAL TABLE



<i>z</i>	<i>Area</i>		<i>z</i>	<i>Area</i>		<i>z</i>	<i>Area</i>
0.00	0.00		1.50	86.64		3.00	99.730
0.05	3.99		1.55	87.89		3.05	99.771
0.10	7.97		1.60	89.04		3.10	99.806
0.15	11.92		1.65	90.11		3.15	99.837
0.20	15.85		1.70	91.09		3.20	99.863
0.25	19.74		1.75	91.99		3.25	99.885
0.30	23.58		1.80	92.81		3.30	99.903
0.35	27.37		1.85	93.57		3.35	99.919
0.40	31.08		1.90	94.26		3.40	99.933
0.45	34.73		1.95	94.88		3.45	99.944
0.50	38.29		2.00	95.45		3.50	99.953
0.55	41.77		2.05	95.96		3.55	99.961
0.60	45.15		2.10	96.43		3.60	99.968
0.65	48.43		2.15	96.84		3.65	99.974
0.70	51.61		2.20	97.22		3.70	99.978
0.75	54.67		2.25	97.56		3.75	99.982
0.80	57.63		2.30	97.86		3.80	99.986
0.85	60.47		2.35	98.12		3.85	99.988
0.90	63.19		2.40	98.36		3.90	99.990
0.95	65.79		2.45	98.57		3.95	99.992
1.00	68.27		2.50	98.76		4.00	99.9937
1.05	70.63		2.55	98.92		4.05	99.9949
1.10	72.87		2.60	99.07		4.10	99.9959
1.15	74.99		2.65	99.20		4.15	99.9967
1.20	76.99		2.70	99.31		4.20	99.9973
1.25	78.87		2.75	99.40		4.25	99.9979
1.30	80.64		2.80	99.49		4.30	99.9983
1.35	82.30		2.85	99.56		4.35	99.9986
1.40	83.85		2.90	99.63		4.40	99.9989
1.45	85.29		2.95	99.68		4.45	99.9991

# Chi-Square Table

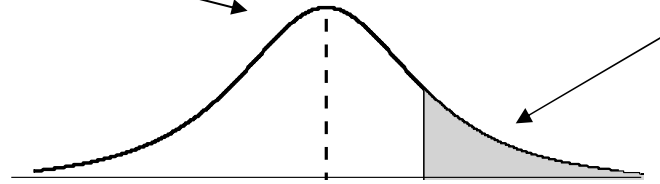


Degrees of freedom ↓	30%	10%	5%	1%	0.1%	← p-value
1	1.07	2.71	3.84	6.63	10.83	← Chi-square
2	2.41	4.61	5.99	9.21	13.82	
3	3.66	6.25	7.81	11.34	16.27	
4	4.88	7.78	9.49	13.28	18.47	
5	6.06	9.24	11.07	15.09	20.52	
6	7.23	10.64	12.59	16.81	22.46	
7	8.38	12.02	14.07	18.48	24.32	
8	9.52	13.36	15.51	20.09	26.12	
9	10.66	14.68	16.92	21.67	27.88	
10	11.78	15.99	18.31	23.21	29.59	
11	12.90	17.28	19.68	24.72	31.26	
12	14.01	18.55	21.03	26.22	32.91	
13	15.12	19.81	22.36	27.69	34.53	
14	16.22	21.06	23.68	29.14	36.12	
15	17.32	22.31	25.00	30.58	37.70	
16	18.42	23.54	26.30	32.00	39.25	
17	19.51	24.77	27.59	33.41	40.79	
18	20.60	25.99	28.87	34.81	42.31	
19	21.69	27.20	30.14	36.19	43.82	
20	22.77	28.41	31.41	37.57	45.31	
21	23.86	29.62	32.67	38.93	46.80	
22	24.94	30.81	33.92	40.29	48.27	
23	26.02	32.01	35.17	41.64	49.73	
24	27.10	33.20	36.42	42.98	51.18	

Student's curve, with degrees of freedom shown at the left of the table

### Student's *t*-TABLE

The shaded area is shown along the top of the table



*t* is shown in the body of the table

<i>Degrees of freedom</i>	<b>25%</b>	<b>10%</b>	<b>5%</b>	<b>2.5%</b>	<b>1%</b>	<b>0.5%</b>
1	1.00	3.08	6.31	12.71	31.82	63.66
2	0.82	1.89	2.92	4.30	6.96	9.92
3	0.76	1.64	2.35	3.18	4.54	5.84
4	0.74	1.53	2.13	2.78	3.75	4.60
5	0.73	1.48	2.02	2.57	3.36	4.03
6	0.72	1.44	1.94	2.45	3.14	3.71
7	0.71	1.41	1.89	2.36	3.00	3.50
8	0.71	1.40	1.86	2.31	2.90	3.36
9	0.70	1.38	1.83	2.26	2.82	3.25
10	0.70	1.37	1.81	2.23	2.76	3.17
11	0.70	1.36	1.80	2.20	2.72	3.11
12	0.70	1.36	1.78	2.18	2.68	3.05
13	0.69	1.35	1.77	2.16	2.65	3.01
14	0.69	1.35	1.76	2.14	2.62	2.98
15	0.69	1.34	1.75	2.13	2.60	2.95
16	0.69	1.34	1.75	2.12	2.58	2.92
17	0.69	1.33	1.74	2.11	2.57	2.90
18	0.69	1.33	1.73	2.10	2.55	2.88
19	0.69	1.33	1.73	2.09	2.54	2.86
20	0.69	1.33	1.72	2.09	2.53	2.85
21	0.69	1.32	1.72	2.08	2.52	2.83
22	0.69	1.32	1.72	2.07	2.51	2.82
23	0.69	1.32	1.71	2.07	2.50	2.81
24	0.68	1.32	1.71	2.06	2.49	2.80
25	0.68	1.32	1.71	2.06	2.49	2.79