# Hardest Unit ☹

Unit 9 – Inferences from Two Samples All Days

Your Bad Planning Professor Colton

# Unit 9 - Outline

# Inference! Our Third Look

# Full Problem

**Setup**: Is there a difference in the proportion of male and female college students who prefer Starbucks as their favorite coffee spot? We collected data from random samples of students from CSCC and found that 31 out of 114 males and 63 out of 176 females prefer Starbucks. Test an appropriate hypothesis with $\alpha = 0.05$.

## Solution

Hypotheses:

Let $p_1$ = true proportion of males who prefer Starbucks
Let $p_2$ = true proportion of females who prefer Starbucks
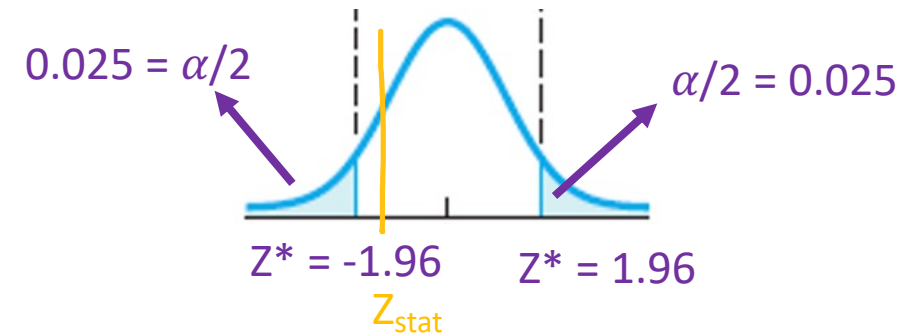
$H_0: p_1 - p_2 = 0$
$H_A: p_1 - p_2 \neq 0$

Set $\alpha = 0.05$

Check Assumptions:
- *Randomization: Random sample of males and females were taken*
- *Independence: Males and females are independent groups*
- *Large enough samples:*
  - *Males → 31 successes and 83 failures, both > 5*
  - *Females →63 successes and 113 failures, both > 5*

- *All conditions are met, appropriate to continue with test!*

Rejection Region:
$Z^* = invNorm(area = 0.05/2, \mu = 0, \sigma = 1) = -1.96$



$0.025 = \alpha/2$     $\alpha/2 = 0.025$

Z* = -1.96     Z* = 1.96

$Z_{stat}$

Test Statistic:
TS = $Z_{stat}$ = 2–PropZTest($x_1$ = 31, $n_1$ = 114, $x_2$ = 63, $n_2$ = 176, $p_1 \neq p_2$) = -1.529

$|Z_{stat}|$ = 1.529 < 1.96 = $|Z^*|$ → Fail to reject $H_0$

Conclusion and Interpretation:
*Because the absolute value our Test Statistic $Z_{stat}$ = 1.529 is less than the absolute value of the Critical Value Z*= 1.96 (5% significance level), we fail to reject the Null hypothesis. We do NOT have sufficient evidence to conclude that the true proportion of male college students who prefer Starbucks is different than that of females.*

# Hypothesis Test Steps – Reminder

1. **State** the Hypotheses
   ○ Define parameter + context.

2. **Check** Assumptions.

3. **Determine** and **Sketch** Rejection Region based of Significance Level

4. **Compute** value of Test Statistic / P-value.

5. **Conclude** and **Interpret**
   ○ State whether you reject $H_0$ or fail to reject $H_0$ AND WHY!
   ○ Interpret your results in the context of the problem

# The Hypothesis Statements – Two Samples

1. State the Hypotheses
   ○ **Define parameter + context.**

## Full Example

Is there a difference in the proportion of male and female college students who prefer Starbucks as their favorite coffee spot? Test an appropriate hypothesis.

Parameters:
*Let $p_1$ = true proportion of males who prefer Starbucks*
*Let $p_2$ = true proportion of females who prefer Starbucks*

## Define Parameters

- Now we have TWO populations and TWO **parameters**!

- These parameters describe the <u>same quantity</u> (ex: 'true proportion who prefer Starbucks') BUT for <u>different groups</u> (ex: males vs females)!

- So we have to CLEARLY <u>define</u> both of them!

  ○ Using subscripts of 1 and 2 will be helpful because that is the notation we will use for the calculations and calculator!
  ○ Order matters and will be important when making our conclusions!

# The Hypothesis Statements – Two Samples

1. **State the Hypotheses**
   o **Define parameter + context.**

Is there a difference in the proportion of male and female college students who prefer Starbucks as their favorite coffee spot?

Hypotheses:
*Let $p_1$ = true proportion of males who prefer Starbucks*
*Let $p_2$ = true proportion of females who prefer Starbucks*

$$H_0 : p_1 - p_2 = 0$$

Null Hypothesis $H_0$

- Now we are comparing some quantity (the same quantity) between TWO populations! So we have TWO parameters!

  o We are not necessarily interested in the specific values of these parameters like we were when testing ONE sample (ex: $H_0$: $p = p_0$)
  o Rather we want to learn about the relationship between the two of them, $p_1$ ?? $p_2$

- We start by assuming both parameters are **equivalent**! So their difference is ZERO!

- This can be written in two ways!

### Option 1

- Directly equating the two parameters:

  $H_0$: $p_1 = p_2$
  $H_0$: $\mu_1 = \mu_2$

### Option 2

- Rewrite as a **difference**!

  $H_0$: $p_1 - p_2 = 0$
  $H_0$: $\mu_1 - \mu_2 = 0$

- Writing our Null hypothesis in this way helps us visualize the Normal curves we use for the CV and TS
- This is an equivalent way to represent it, just a change in perspective to the DIFFERENCE (a single value)

*\*\* Can think of these as a new population of DIFFERENCES!*

-0.15  -0.10  -0.05   0   0.05   0.10   0.15   -3   -2   -1   0   1   2   3

$p_1 - p_2$  → Standardize → $Z$

# The Hypothesis Statements – Two Samples

1. **State the Hypotheses**
   ○ Define parameter + context.

## Alternative Hypothesis $H_A$

- Here is where we state our research interest.

- Again, we are interested in the underlined relationship between our two parameters and how we think they are different

  ○ Our test on the single value of the **difference** may be left-tailed (<), right-tailed (>), or two-tailed (≠).

- Alternative can also be written in two ways:

  Option 1          Option 2

  $H_A$: $p_1 \neq p_2$  →  $p_1 - p_2 \neq 0$
  $H_A$: $p_1 > p_2$  →  $p_1 - p_2 > 0$
  $H_A$: $p_1 < p_2$  →  $p_1 - p_2 < 0$
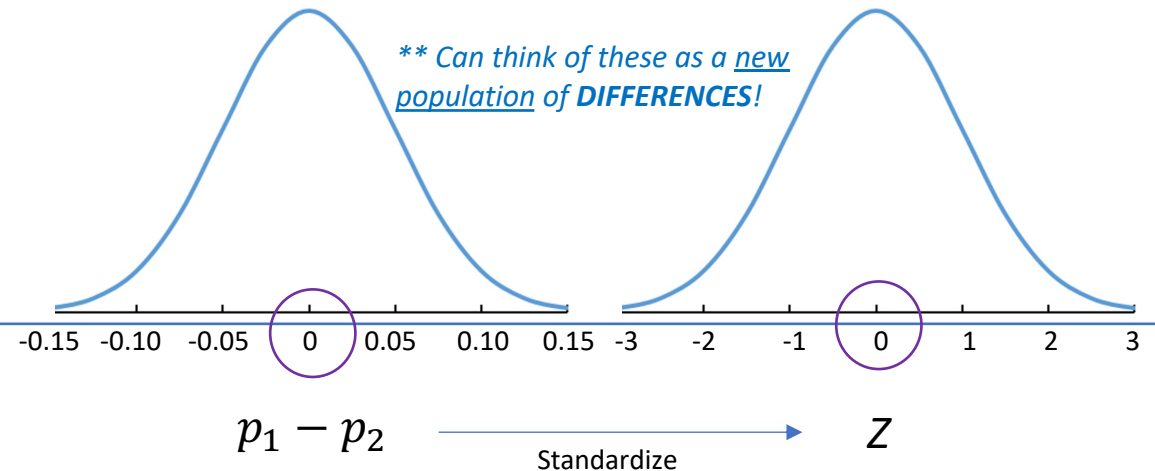
  *(same for $\mu_1$ and $\mu_2$)*

Is there a difference in the proportion of male and female college students who prefer Starbucks as their favorite coffee spot?

Hypotheses:
*Let $p_1$ = true proportion of males who prefer Starbucks*
*Let $p_2$ = true proportion of females who prefer Starbucks*

$H_0: p_1 - p_2 = 0$
$H_A: p_1 - p_2 \neq 0$



| | Reject $H_0$ | Do not reject $H_0$ | Reject $H_0$ |
|---|---|---|---|
| (a) Two tailed | | | |

(b) Left tailed

(c) Right tailed

| | | |
|---|---|---|
| Type of Test: | (a) Two tailed | (b) Left tailed | (c) Right tailed |
| Sign in $H_A$: | ≠ | < | > |
| Rejection Region: | Both sides | Left side | Right side |

# LCQ – Hypotheses

**Problem**: (1) Define the parameters of interest and (2) State the Null and Alternative for the following scenarios:

a) A company has two separate teams (A and B) who complete accounting reports. From a random sample of recent reports from each team, they want to know if these teams perform differently in terms of proportion of reports completed on time.

b) A company randomly selects employees to complete yearly surveys on whether or not they enrolled in at least one wellness class at the company's site. They want to know if a greater percentage took a wellness class **this year compared to last year**.

# LCQ – Hypotheses

**Problem**: (1) Define the parameters of interest and (2) State the Null and Alternative for the following scenarios:

a) A company has two separate teams (A and B) who complete accounting reports. From a random sample of recent reports from each team, they want to know if these teams perform differently in terms of proportion of reports completed on time.

*Let $p_1$ = the true proportion of completed accounting reports from team A*
*Let $p_2$ = the true proportion of completed accounting reports from team B*

$H_0$: $p_1 - p_2 = 0$          OR          $H_0$: $p_1 = p_2$
$H_A$: $p_1 - p_2 \neq 0$          OR          $H_A$: $p_1 \neq p_2$
*Perfect! Both versions are correct, there are just different benefits to each!*
*The difference version (first way) helps us when looking at the result of our TS, but our calculator uses the direct comparison (second way) in the menus*

b) A company randomly selects employees to complete yearly surveys on whether or not they enrolled in at least one wellness class at the company's site. They want to know if a <u>greater</u> percentage took a wellness class **this year compared to last year**.

*Let $p_1$ be the true proportion of employees **enrolled in at least** one wellness class THIS YEAR → Correct! THIS YEAR is our first group, that is the first population*
*Let $p_2$ be the true proportion of employees **that did not enrolled in** at least one wellness class → INCORRECT! Because this is now representing a different quantity than in $p_{1}$ (enrolled vs did NOT enroll); the only thing that should change is the GROUP!*
*Let $p_2$ be the true proportion of employees **enrolled in at least** one wellness class LAST YEAR →Good! LAST YEAR is our second group, now our parameter is for the <u>same quantity</u> but a <u>different POPULATION</u>*

$H_0$: $p_1 - p_2 = 0$
$H_A$: $p_1 - p_2 > 0$
*Very good!*

*What if switch the order of our parameters? Now let:*
*$p_1$ = LAST year*
*$p_2$ = THIS year*

*$H_0$: $p_1 - p_2 = 0$ → same*
*$H_A$: $p_1 - p_2 < 0$ → different! We still want this year to be greater, so now $p_2$ is the larger value making the difference less than zero! Order is important*

# Assumptions

## 2. Check Assumptions.

- For the most part this step is the <u>same as we have seen before</u>!

  - Random Sample and Large enough sample
  - How we check the Large enough sample assumption depends on the type of test (type of data)

- We just have to do it for <u>both samples</u> now!

- Although now we also have to think about the <u>connection between our two samples</u>

  - Will go over these again when looking at Proportions Tests and Means Tests

<u>Full Example</u>

Is there a difference in the proportion of male and female college students who prefer Starbucks as their favorite coffee spot?
- They random samples of students from CSCC and found that 31 out of 114 males and 63 out of 176 females prefer Starbucks

Check Assumptions:
- *Randomization: Random sample of males and females was taken*
- *Independence: Males and females are independent groups*
- *Large enough samples:*
  - Males → *31 successes and 83 failures, both > 5*
  - Females →*63 successes and 113 failures, both > 5*

- *All conditions are met, appropriate to continue with test!*

\* will go through these with the proportions slide

# Rejection Region and TS – Two Samples

## Rejection Region (RR)

- This is the EXACT same as we have seen in all the other types of tests!

  - Which is because we are framing our tests from the perspective of a difference!

**4. Compute value of Test Statistic / P-value.**

## Test Statistic (TS) and P-Value

- SAME logic as with one sample, just different calculations behind the scenes

- Will go over the specifics on each respective Test's slides

---

Full Example

Is there a difference in the proportion of male and female college students who prefer Starbucks as their favorite coffee spot?

Rejection Region:
*Set $\alpha$ = 0.05 (which was set at beginning)*
*$Z^*$ = invNorm(area = 0.05/2, $\mu$ = 0, $\sigma$ = 1) = -1.96*

$0.025 = \alpha/2$

$\alpha/2 = 0.025$

p-value = 0.126

$Z^*$ = -1.96     $Z^*$ = 1.96

$Z_{stat}$

Test Statistic:
*TS = $Z_{stat}$ = 2–PropZTest($x_1$ = 31, $n_1$ = 114, $x_2$ = 63, $n_2$ = 176, $p_1 \neq p_2$)*
*= -1.529*

*$|Z_{stat}|$ = 1.529 < 1.96 = $|Z^*|$ → Fail to reject $H_0$*

P-Value:
*p-value = 2–PropZTest($x_1$ = 31, $n_1$ = 114, $x_2$ = 63, $n_2$ = 176, $p_1 \neq p_2$) = 0.126*

# Conclude and Interpret – Two Samples

## 5. Conclude and Interpret
- State whether you reject $H_0$ or fail to reject $H_0$ AND WHY!
- Interpret your results in the context of the problem

This has the SAME structure that we had with ONE sample Hypothesis Tests

- We just need to talk about BOTH parameters for the Alternative!

First Part – Decision and Reasoning

- Because (**comparison of TS and CV; OR p-value and $\alpha$**) we (**REJECT** or **FAIL TO REJECT**) the Null Hypothesis.

Second Part – Interpretation

- There (**IS** or **IS NOT**) sufficient evidence to conclude (**THE ALTERNATIVE HYPOTHESIS + CONTEXT**).

- NOTE about wording!

  - When thinking about the alternative as a <u>difference</u> (very literally), our wording could be something similar to "There is / is not sufficient to conclude the true difference in proportion of males who prefer Starbucks and of females is not equal to (or less / greater than) zero."

  - This doesn't read very well and it's not how we would talk about the results in conversation. So try and **reword** it as a <u>direct comparison</u> of the <u>two population parameters</u> like it is in the example (<u>not in terms of the difference and zero</u>)!
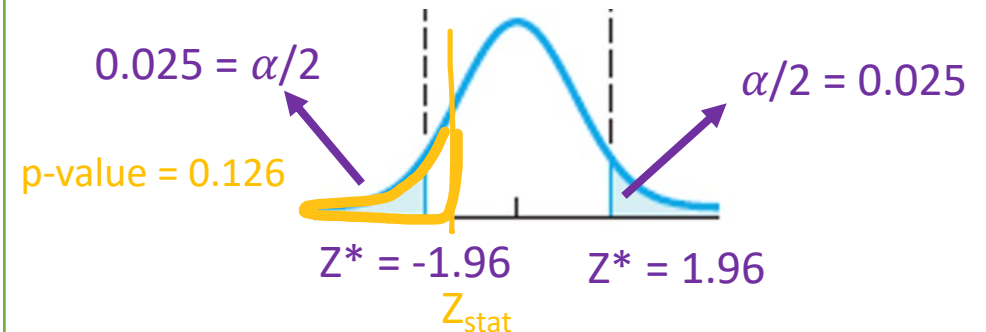
Full Example

Is there a difference in the proportion of male and female college students who prefer Starbucks as their favorite coffee spot?

Conclusion and Interpretation:
*Because*
- *the absolute value of our Test Statistic $|Z_{stat}| = 1.529$ is less than the Critical Value $Z^* = 1.69$ (5% significance level)*
    *OR*
- *our p-value = 0.126 is greater than the significance level 0.01*

*We do have sufficient evidence to conclude that the true proportion of male college students who prefer Starbucks is different than that of females.*

** Give MORE INFORMATION IF POSSIBLE

- If REJECT a two-tailed test, say which of the parameters is l
- Just look at the sign of the TS and the order of our Hypoth
- Ex) If conclude $H_A$: $p_1 - p_2 \neq 0$ and TS = 2.3 $\rightarrow$ $p_1 > p_2$, so sa

# Hypothesis Tests for Proportions – Two Samples!

- Everything above applies, now we are just going to apply it specifically to a Two Sample Proportions Test!

# Proportions Assumptions- Two Samples

## 2. Check Assumptions.

- Some of the same ideas with some new ones as well!

1) Random samples (both of them)

2) Independence

- Because we have <u>two samples</u> now, we need each group to be **independent** (unrelated, no connection).
- So they results from one group should NOT have <u>any effect / impact</u> on the results of the second group!

3) Large Enough Samples (both of them)

- Because we don't have a Null proportion value like before, we can't check $np_0 \geq 5$ AND $n(1 - p_0) = nq_0 \geq 5$
- So all we have to do is make sure *each* sample has at least 5 success and 5 failures

  - Sample 1: $n_1\hat{p}_1 \geq 5$, $n_1(1 - \hat{p}_1) = n\hat{q}_1 \geq 5$ AND Sample 2: $n_2\hat{p}_2 \geq 5$, $n_2(1 - \hat{p}_2) = n\hat{q}_2 \geq 5$

  - So we are actually looking at the <u>sample data</u> here (which was a big no no when we were doing one sample...)

Full Example

Is there a difference in the proportion of male and female college students who prefer Starbucks as their favorite coffee spot?
- They random samples of students from CSCC and found that 31 out of 114 males and 63 out of 176 females prefer Starbucks

Check Assumptions:
- *Randomization: Random sample of males and females was taken*
- *Independence: Males and females are independent groups*
- *Large enough samples:*
  - *Males → 31 successes and 83 failures, both > 5*
  - *Females → 63 successes and 113 failures, both > 5*

- *All conditions are met, appropriate to continue with test!*

---

One Sample Test Conditions

✓Randomization Condition
   Need to have a random sample
✓Large Enough Sample Condition
   $np_0 \geq 5$ AND $n(1 - p_0) = nq_0 \geq 5$ OR
   EXPECT AT LEAST 5 successes and 5 failures

New conditions

**TWO Sample Test Conditions**

✓Randomization Condition
   Need to have two random samples
✓Independence Condition
   Need to independent samples
✓Large Enough Sample Condition

# LCQ – Assumptions

**Problem**: Check the conditions for a Hypothesis Test of the two population proportions for the following scenarios:

a) A company has two separate teams (A and B) who complete accounting reports. From a random sample of 50 recent reports from each team, 40% of Team A's were on time and 36% of Team B's were on time. They want to know if these teams perform differently in terms of proportion of reports completed on time.

b) A company randomly selects employees to complete yearly surveys on whether or not they enrolled in at least one wellness class at the company's site. Last year's survey showed 81 out of 100 employees had taken a wellness class and this year 102 out of 140 had. They want to know if a greater percentage took a wellness class this year compared to last year.

# LCQ – Assumptions

**Problem**: Check the conditions for a Hypothesis Test of the two population proportions for the following scenarios:

a) A company has two separate teams (A and B) who complete accounting reports. From a random sample of 50 recent reports from each team, 40% of Team A's were on time and 36% of Team B's were on time. They want to know if these teams perform differently in terms of proportion of reports completed on time.

*1) Random condition: 'Random sample of 50 reports from each team' Yes! → Have random sample from both groups*
*2) Independence condition: Separate teams so independent, Yes!! → No reason to think these separate teams impact each other*
*3) Large enough samples condition: Yes! → Using sample proportions to check these below*
- *Team A: 50(.4) = 20 > 5 and 50(.6) = 30 > 5*
- *Team B: 50(.36) = 18 > 5 and 50(.64) = 32 > 5*

*ALL conditions are met! Okay to continue with test!*

b) A company randomly selects employees to complete yearly surveys on whether or not they enrolled in at least one wellness class at the company's site. Last year's survey showed 81 out of 100 employees had taken a wellness class and this year 102 out of 140 had. They want to know if a greater percentage took a wellness class this year compared to last year.

*1) Random condition: Company randomly selects individuals to complete the surveys each year, Yes!*
*2) Independence condition: Surveys from different years, Yes! → No reason to think one wellness class enrollment depends on which year it was. And we are not the selecting same people because of randomness*
*3) Large enough samples condition: Yes! → Can directly use the given sample results (no need to get the proportions first) to check these below*
- *This year: 102 successes and 38 failures, both > 5*
- *Last year: 81 successes and 19 failures, both > 5*

*ALL conditions are met! Okay to continue with test!*

# Using Calc - Test Statistic and P-Value for Proportions

## 4. Compute value of Test Statistic / P-value.

**Setup**
A NatGeo Poll interviewed 1200 hiking enthusiasts and 1100 climbers. They asked "Are you more afraid of spiders or snakes???" 768 of the 1200 hikers and 662 of the 1100 climbers, responded "Ewww, snakes…." Is there enough evidence to conclude the proportion of hikers who are more afraid of snakes is different than that of climbers? Use $\alpha = 0.1$

**GOAL**: Conduct a Hypothesis Test!

1.    2–PropZTest
   a)   $x_1$ = number of successes in sample 1
   b)   $n_1$ = sample size 1
   c)   $x_2$ = number of successes in sample 2
   d)   $n_2$ = sample size 2
   e)   $p_1$: Alternative hypothesis with $p_2$ → ** *NOTE this is not in terms of the difference*
   Calculate or Draw

Formula for $Z_{stat}$ by hand:

$$Z = \frac{(\hat{p}_1 - \hat{p}_2) - 0}{\sqrt{\hat{p}(1-\hat{p})\left(\frac{1}{n_1} + \frac{1}{n_2}\right)}}$$

$\hat{p}_1$ = sample proportion from population 1
$\hat{p}_2$ = sample proportion from population 2
$\hat{p}$ = pooled sample proportion
$n_1$ = sample size of group 1
$n_2$ = sample size of group 2

Combined ("pooled")
sample proportion
(no subscript)
$$= \hat{p} = \frac{x_1 + x_2}{n_1 + n_2}$$

# Using Calc - Test Statistic and P-Value for Proportions

## 4. Compute value of Test Statistic / P-value.

**Setup**
A NatGeo Poll interviewed 1200 hiking enthusiasts and 1100 climbers. They asked "Are you more afraid of spiders or snakes???" 768 of the 1200 hikers and 662 of the 1100 climbers, responded "Ewww, snakes…." Is there enough evidence to conclude the proportion of hikers who are more afraid of snakes is different than that of climbers? Use $\alpha = 0.1$
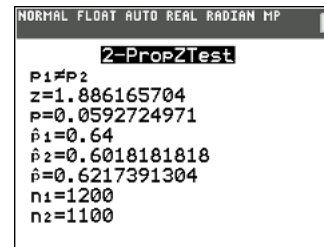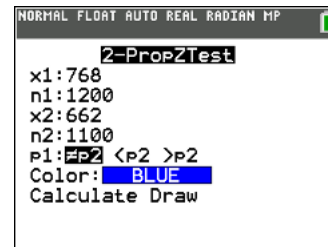
Formula for $Z_{stat}$ by hand:

$$Z = \frac{(\hat{p}_1 - \hat{p}_2) - 0}{\sqrt{\hat{p}(1-\hat{p})\left(\frac{1}{n_1} + \frac{1}{n_2}\right)}}$$

$\hat{p}_1 = $ sample proportion from population 1
$\hat{p}_2 = $ sample proportion from population 2
$\hat{p} = $ pooled sample proportion
$n_1 = $ sample size of group 1
$n_2 = $ sample size of group 2

Combined ("pooled") sample proportion (no subscript) $= \hat{p} = \dfrac{x_1 + x_2}{n_1 + n_2}$

**GOAL**: Conduct a Hypothesis Test!

1. 2–PropZTest

   a) $x_1 = $ number of successes in sample 1
   b) $n_1 = $ sample size 1
   c) $x_2 = $ number of successes in sample 2
   d) $n_2 = $ sample size 2
   e) $p_1$: Alternative hypothesis with $p_2$ → ** *NOTE this is not in terms of the difference*

   Calculate or Draw



NORMAL FLOAT AUTO REAL RADIAN MP
```
2-PropZTest
x1:768
n1:1200
x2:662
n2:1100
p1:≠p2 <p2 >p2
Color:  BLUE
Calculate Draw
```

NORMAL FLOAT AUTO REAL RADIAN MP
```
2-PropZTest
p1≠p2
z=1.886165704
p=0.0592724971
p̂1=0.64
p̂2=0.6018181818
p̂=0.6217391304
n1=1200
n2=1100
```

*(p₁= hikers and p₂ = climbers)*
$H_0: p_1 - p_2 = 0$
$H_A: p_1 - p_2 \neq 0$

**Calculate Output**
$p_1 \neq <> p_2$ Alternative hypothesis
$z = Z_{stat}$
$p = $ p-value
$\hat{p}_1 = $ sample proportion 1
$\hat{p}_2 = $ sample proportion 2
$\hat{p} = $ pooled sample proportion
$n_1 = $ sample size 1
$n_2 = $ sample size 2

NORMAL FLOAT AUTO REAL RADIAN MP
```
2-PropZTest
z=1.8862        p=0.0593
```

**Draw Output**
Plot (and displays values) of p = p-value and z = $Z_{stat}$ on the standard normal curve

# LCQ – Conclusions and Interpretations

5. **Conclude** and **Interpret**
  - ○ State whether you reject $H_0$ or fail to reject $H_0$ AND WHY!
  - ○ Interpret your results in the context of the problem

**Problem**: Write the conclusions and interpretations for the previous scenarios using our results.

**Setup:** A NatGeo Poll interviewed 1200 hiking enthusiasts and 1100 climbers. They asked "Are you more afraid of spiders or snakes???" 768 of the 1200 hikers and 662 of the 1100 climbers, responded "Ewww, snakes…." Is there enough evidence to conclude the proportion of hikers who are more afraid of snakes is different than that of climbers? Use $\alpha = 0.1$

**Solution**:

NORMAL FLOAT AUTO REAL RADIAN MP

2-PropZTest
z=1.8862          p=0.0593

# LCQ – Conclusions and Interpretations
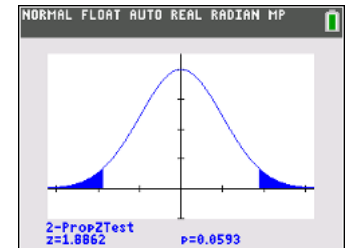
5. **Conclude** and **Interpret**
  ○ State whether you reject $H_0$ or fail to reject $H_0$ AND WHY!
  ○ Interpret your results in the context of the problem

**Problem**: Write the conclusions and interpretations for the previous scenarios using our results.

**Setup:** A NatGeo Poll interviewed 1200 hiking enthusiasts and 1100 climbers. They asked "Are you more afraid of spiders or snakes???" 768 of the 1200 hikers and 662 of the 1100 climbers, responded "Ewww, snakes…." Is there enough evidence to conclude the proportion of hikers who are more afraid of snakes is different than that of climbers? Use $\alpha = 0.1$

**Solution**:
*Need these…*

*Let $p_1$ = the true proportion of hikers who are more afraid of snakes than spiders*
*Let $p_2$ = the true proportion of climbers who are more afraid of snakes than spiders*

$H_0: p_1 - p_2 = 0$
$H_A: p_1 - p_2 \neq 0$
$\alpha = 0.1$

NORMAL FLOAT AUTO REAL RADIAN MP

2-PropZTest
z=1.8862    p=0.0593

*P-Value*
*P-value = 2–PropZTest($x_1$ = 768, $n_1$ = 1200, $x_2$ = 662, $n_2$ = 1100, $p_1 \neq p_2$) = 0.0593*
*p-value = 0.0593 < 0.10 = $\alpha$ → Reject $H_0$!*

*Conclusion and Interpretation*

First part → *Because our p-value = 0.0593 is less than the significance level 0.10, we reject the Null hypothesis.*

Second part → *There IS sufficient evidence to conclude that the true proportion of hikers who are more afraid of snakes than spiders is different than that of climber*

More info → *Further we can say that hikers' proportion is actually greater ($Z_{stat}$ = 1.89).*

# Hypothesis Tests for Means – INDEPENDENT Samples (and Known σ)

- All of the previous Hypothesis tests overview applies, now we are just going to apply it specifically to a Two Sample Means Test!

- And going back to the One Sample Means Test, we still have to determine if the population standard deviation is known or unknown.

  - This tells us if we are doing a Z distribution based Test or a T distribution based Test!

- Now we have to also think about the relationship between our two population data sources → This section is for **independent** samples!

  - And we will start with KNOWN population standard deviations $\sigma_1$ and $\sigma_2$

# LCQ: Independent vs Dependent Samples

How to think about samples
- Independent samples → Groups are unrelated, no connection, no relationship
- Dependent samples → Groups have some relationship between one another, can link the two; PAIRS

**Problem**: Determine if the following scenarios are independent or dependent samples.

1) Comparing the blood pressure of STAT 1450 students before the final exam and after completing the final exam.

2) Seeing if the height of Faculty is shorter than the undergraduate population.

3) Looking to see if there is a difference in the price of the same Video Game Consoles at Target or Walmart.

4) A study is conducted to see what effect a new drug has on dexterity. A random sample of 30 students is chosen. They are given a series of tasks to perform and a score reflecting their performance. A dose of the drug is given to the 30 students and they again perform similar tasks and are scored again.

# LCQ: Independent vs Dependent Samples

**How to think about samples**
- Independent samples → Groups are unrelated, no connection, no relationship
- Dependent samples → Groups have some relationship between one another, can link the two; PAIRS

**Problem**: Determine if the following scenarios are independent or dependent samples.

1) Comparing the blood pressure of STAT 1450 students before the final exam and after completing the final exam.

*Dependent! → There is a relationship between the blood pressure before the final and after the completion of the final. Connection is measuring the SAME student twice*

2) Seeing if the height of Faculty is shorter than the undergraduate population.

*Independent → There is no direct connection (or inherent relationship) between faculty and undergrads*

3) Looking to see if there is a difference in the price of the same Video Game Consoles at Target or Walmart.

*~~Independent??? Two different stores~~*
*Dependent! → No relationship between Target and Walmart, BUT we are looking at the SAME console at the two different stores (groups). So there is a relationship with the consoles (think pairs of X-boxes, one at Walmart and one at Target; same for a PS4)*

4) A study is conducted to see what effect a new drug has on dexterity. A random sample of 30 students is chosen. They are given a series of tasks to perform and a score reflecting their performance. A dose of the drug is given to the 30 students and they again perform similar tasks and are scored again.

*Dependent → SAME students before and after drug. So there is a relationship between the two groups*

# The Hypothesis Statements for Two Samples - Review

1. State the Hypotheses
   ○ **Define parameter + context.**

** Can think of these as a new population of DIFFERENCES!

-30  -20  -10  0  10  20  30    -3  -2  -1  0  1  2  3

$\mu_1 - \mu_2$  →  Z

Standardize
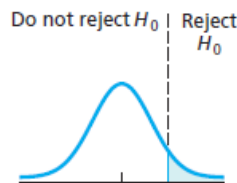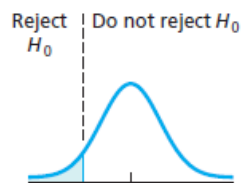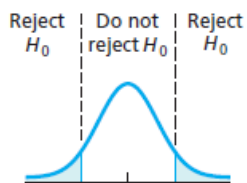
### Define Parameters

- Now we have TWO populations and TWO **parameters**!
- These describe the <u>same quantity</u>, just for <u>different groups</u>!
  ○ Quantitative (numeric) → population means $\mu_1$ and $\mu_2$

### Null Hypothesis $H_0$

- We want to learn about the <u>relationship</u> between the two of them, $\mu_1$ ?? $\mu_2$
- So start by <u>assuming both parameters</u> are **equivalent**! So their difference is ZERO!

### Alternative Hypothesis $H_A$

- This is where we state how we think these means are different
- Our test on the single value of the **difference** may be left-tailed (<), right-tailed (>), or two-tailed (≠).

**How to write hypotheses**
- These can be written in two ways

| Difference | OR | Direct Comparison |
|---|---|---|
| $H_0$: $\mu_1 - \mu_2 = 0$ | | $H_0$: $\mu_1 = \mu_2$ |
| | | |
| $H_A$: $\mu_1 - \mu_2 \neq 0$ | | $H_A$: $\mu_1 \neq \mu_2$ |
| $H_A$: $\mu_1 - \mu_2 < 0$ | | $H_A$: $\mu_1 < \mu_2$ |
| $H_A$: $\mu_1 - \mu_2 > 0$ | | $H_A$: $\mu_1 > \mu_2$ |

Reject $H_0$ | Do not reject $H_0$ | Reject $H_0$

Reject $H_0$ | Do not reject $H_0$

Do not reject $H_0$ | Reject $H_0$

(a) Two tailed      (b) Left tailed      (c) Right tailed

| Type of Test: | (a) Two tailed | (b) Left tailed | (c) Right tailed |
|---|---|---|---|
| Sign in $H_A$: | ≠ | < | > |
| Critical Values: | -Z*    +Z* | -Z* | Z* |

# LCQ – Two Sample Means Hypotheses

**Problem**: (1) Define the parameters of interest and (2) State the Null and Alternative for the following scenarios:

a) A researcher random sampled 15 infants and 20 toddlers to measure their body temperatures. Let's assume that body temperatures for all persons are normally distributed with some known standard deviation. Is there sufficient evidence to conclude that the mean body temperatures for infants and toddlers differ?

b) A prospective home buyer is trying to decide which city to move to. In Washington D.C, a random sample of 32 foreclosed homes sold for an average of $443,705 and in Baltimore and random sample of 45 foreclosed homes sold for $450,000. Real estate experts say the standard deviation for sales across the nation is $190,000. Is there sufficient evidence that the Baltimore housing market is more expensive?

# LCQ – Two Sample Means Hypotheses

**Problem**: (1) Define the parameters of interest and (2) State the Null and Alternative for the following scenarios:

a) A researcher randomly sampled 15 infants and 20 toddlers to measure their body temperatures. Let's assume that body temperatures for all persons are normally distributed with some known standard deviation. Is there sufficient evidence to conclude that the mean body temperatures for infants and toddlers differ?

<u>Define Parameters</u>

*Let $\mu_1$ = The TRUE mean body temperature of infants*
*Let $\mu_2$ = The POPULATION mean body temperature of toddlers*

<u>Hypotheses</u>

*$H_0: \mu_1 = \mu_2$   OR    $H_0: \mu_1 - \mu_2 = 0$ → Both CORRECT!*

*$H_A: \mu_1 \neq \mu_2$   OR   $H_A: \mu_1 - \mu_2 \neq 0$ → CORRECT! We just want in to be different, not only interested in less than or greater than*

b) A prospective home buyer is trying to decide which city to move to. In Washington D.C, a random sample of 32 foreclosed homes sold for an average of $443,705 and in Baltimore and random sample of 45 foreclosed homes sold for $450,000. Real estate experts say the standard deviation for sales across the nation is $190,000. Is there sufficient evidence that the Baltimore housing market is more expensive?

<u>Define Parameters</u>

*Let $P_1$ = The average price of homes sold in DC →NOPE! We are talking about MEANS, should NOT be any population proportions (p) in our problems*
*Let $P_2$ = The average price of homes sold in Baltimore. → BE CAREFUL with your notation or shorthand, even if you meant $P_2$ as 'parameter 2', I would interpret this as a proportion (which is WRONG)*

*Let $\mu_1$ = The TRUE average price of homes sold in DC → Yes!*
*Let $\mu_2$ = The POPULATION average price of homes sold in Baltimore → 'true' and 'population' are synonymous here, both indicate the population parameter*

<u>Hypotheses</u>

*$H_0: \mu_1 = \mu_2$   OR    $H_0: \mu_1 - \mu_2 = 0$ → Both CORRECT!*

*$H_A: \mu_1 > \mu_2$   OR   $H_A: \mu_1 - \mu_2 > 0$ → INCORRECT! Because we want Baltimore to be more expensive, so $\mu_2$ should be LARGER*
*$H_A: \mu_1 < \mu_2$    OR   $H_A: \mu_1 - \mu_2 < 0$ → Now CORRECT!*

# Mean Assumptions - Independent Samples

<div style="border: 1px solid orange">

**2. Check Assumptions.**

</div>

- EXACT same Assumptions as for Means Test with one sample, we just have to do it for both!
    - (Remember we need to know / be given the <u>population standard deviations</u> $\sigma_1$ and $\sigma_2$ for now)

- We have to add one more though!

    - Must have an assumption about the <u>connection between our two samples</u>

<u>Independence</u>

- This is a REALLY important assumption now, because we do a <u>different test based on whether this assumption is met or not</u>!

- Right now we are looking at an INDEPENDENT samples test!
    - We need each group to be **independent** (unrelated, no connection).
    - So they results from one group should NOT have <u>any effect / impact</u> on the results of the second group!

<u>One Sample Test Conditions</u>

✓Randomization Condition
   Need to have a random sample
✓Large Enough Sample Condition
   Normal population OR
   n ≥ 30

**New conditions** ➡️

<div style="border: 1px solid blue">

**<u>INDEPENDENT Samples Test Conditions</u>**

✓Randomization Condition
   Need to have two random samples
✓Independence Condition
   Need to independent samples
✓Large Enough Sample Condition
   Normal populations OR
   $n_1 \geq 30$ AND $n_2 \geq 30$

</div>

# LCQ – Assumptions

**Problem**: Check the conditions for a Hypothesis Test of the two population proportions for the following scenarios:

a) A researcher randomly sampled 15 infants and 20 toddlers to measure their body temperatures. Let's assume that body temperatures for all persons are normally distributed with some known standard deviation. Is there sufficient evidence to conclude that the mean body temperatures for infants and toddlers differ?

b) A prospective home buyer is trying to decide which city to move to. In Washington D.C, a random sample of 32 foreclosed homes sold for an average of $443,705 and in Baltimore and random sample of 45 foreclosed homes sold for $450,000. Real estate experts say the standard deviation for sales across the nation is $190,000. Is there sufficient evidence that the Baltimore housing market is more expensive?

# LCQ – Assumptions

**Problem**: Check the conditions for a Hypothesis Test of the two population proportions for the following scenarios:

a) A researcher randomly sampled 15 infants and 20 toddlers to measure their body temperatures. Let's assume that body temperatures for all persons are normally distributed with some known standard deviation. Is there sufficient evidence to conclude that the mean body temperatures for infants and toddlers differ?

*1) Researcher randomly sampled infants and toddlers→ Yes!*
*2) Independent groups, no relationship between infants and toddlers → Good! This is our explanation of why the samples are independent! No mention of measuring the same child once as an infant and years later as a toddler (so reasonable to assume independent!)*
*3) Large enough samples is met because we are assuming body temperatures for everyone are Normal! → Yes! We have Normal populations for BOTH groups, so no need to look at the sample sizes*

b) A prospective home buyer is trying to decide which city to move to. In Washington D.C, a random sample of 32 foreclosed homes sold for an average of $443,705 and in Baltimore and random sample of 45 foreclosed homes sold for $450,000. Real estate experts say the standard deviation for sales across the nation is $190,000. Is there sufficient evidence that the Baltimore housing market is more expensive?

*1) Random samples were taken as stated in problem → Yes!*
*2) Independent groups → NOT ENOUGH! IF this is ALL you wrote, NOT full credit!! Need to EXPLAIN why there are independent groups for this specific problem*
*Independent groups (DC vs Baltimore), there is no relation between DC houses and Baltimore houses → Now this is BETTER!*
*3) Large enough samples → WHY???? How do you know this??? I don't know that you know if this is all that your write on the Test*
*Large enough sample, n ≥ 30 → STILL NOT FULL credit! BE SPECIFIC! (do NOT be general and just write n ≥ 30); and our cutoff is 30 (not 5, which is the check for np ≥ 5 for proportions)*
*        Because $n_1 = 32 \geq 30$ and $n_2 = 45 \geq 30$ → YES!! And we had to look at the sample sizes because we don't have information about these populations already being normally distributed*

## 4. Compute value of Test Statistic / P-value.

**Setup**
A prospective home buyer is trying to decide which city to move to. In Washington D.C, a random sample of 32 foreclosed homes sold for an average of $443,705 and in Baltimore and random sample of 45 foreclosed homes sold for $450,000. Real estate experts say the standard deviation for sales across the nation is $190,000. Is there sufficient evidence that the Baltimore housing market is more expensive? Use $\alpha = 0.1$

**GOAL**: Conduct a Hypothesis Test!

Formula for $Z_{stat}$ by hand:

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

Under the Null hypothesis, the quantity $\mu_1 - \mu_2 = 0$ and everything else is known.

1. 2–SampZTest
   a) Input = Stats
   b) $\sigma_1$ = population 1 SD
   c) $\sigma_2$ = population 2 SD
   d) $\bar{x}_1$ = sample 1 mean
   e) $n_1$ = sample size 1
   f) $\bar{x}_2$ = sample 2 mean
   g) $n_2$ = sample size 2
   h) $\mu_1$: Alternative hypothesis with $\mu_2$ → ** *NOTE this is not in terms of the difference*
   Calculate or Draw

# Using Calc - Test Statistic and P-Value for Ind Means and Known $\sigma$

## 4. Compute value of Test Statistic / P-value.

**Setup**
A prospective home buyer is trying to decide which city to move to. In Washington D.C, a random sample of 32 foreclosed homes sold for an average of $443,705 and in Baltimore and random sample of 45 foreclosed homes sold for $450,000. Real estate experts say the standard deviation for sales across the nation is $190,000. Is there sufficient evidence that the Baltimore housing market is more expensive? Use $\alpha = 0.1$

**GOAL**: Conduct a Hypothesis Test!

Formula for $Z_{stat}$ by hand:

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

Under the Null hypothesis, the quantity $\mu_1 - \mu_2 = 0$ and everything else is known.

1. 2–SampZTest

   a) Input = Stats
   b) $\sigma_1$ = population 1 SD
   c) $\sigma_2$ = population 2 SD
   d) $\bar{x}_1$ = sample 1 mean
   e) $n_1$ = sample size 1
   f) $\bar{x}_2$ = sample 2 mean
   g) $n_2$ = sample size 2
   h) $\mu_1$: Alternative hypothesis with $\mu_2$ → ** NOTE this is *not* in terms of the *difference*

   Calculate or Draw



```
NORMAL FLOAT AUTO REAL RADIAN MP
PRESS [<] OR [>] TO SELECT AN OPTION
         2-SampZTest
Inpt:Data Stats
σ1:190000
σ2:190000
x̄1:443705
n1:32
x̄2:450000
n2:45
μ1:≠μ2 <μ2 >μ2
↓Color:   BLUE   <>
```

```
NORMAL FLOAT AUTO REAL RADIAN MP
         2-SampZTest
μ1<μ2
z=-0.1432775072
P=0.443035484
x̄1=443705
x̄2=450000
n1=32
n2=45
```

($\mu_1$= DC and $\mu_2$ = Baltimore)

$H_0: \mu_1 - \mu_2 = 0$

$H_A: \mu_1 - \mu_2 < 0$

**Calculate Output**
$\mu_1 \neq <> \mu_2$ Alternative hypothesis
$z = Z_{stat}$
$p$ = p-value
$\bar{x}_1$= sample 1 mean
$\bar{x}_2$= sample 2 mean
$n_1$ = sample 1 size
$n_2$ = sample 2 size



```
NORMAL FLOAT AUTO REAL RADIAN MP
2-SampZTest
z=-0.1433        P=0.443
```

**Draw Output**
Plot (and displays values) of p = p-value and z = $Z_{stat}$ on the standard normal curve

# LCQ – Conclusions and Interpretations

5. **Conclude** and **Interpret**
- State whether you reject $H_0$ or fail to reject $H_0$ AND WHY!
- Interpret your results in the context of the problem

**Problem**: Write the conclusions and interpretations for the previous scenarios using our results.

**Setup:** A prospective home buyer is trying to decide which city to move to. In Washington D.C, a random sample of 32 foreclosed homes sold for an average of $443,705 and in Baltimore and random sample of 45 foreclosed homes sold for $450,000. Real estate experts say the standard deviation for sales across the nation is $190,000. Is there sufficient evidence that the Baltimore housing market is more expensive? Use $\alpha = 0.1$

**Solution**:

```
NORMAL FLOAT AUTO REAL RADIAN MP
        2-SampZTest
µ1<µ2
z=-0.1432775072
p=0.443035484
x̄1=443705
x̄2=450000
n1=32
n2=45
```

# LCQ – Conclusions and Interpretations

**5. Conclude and Interpret**
- State whether you reject $H_0$ or fail to reject $H_0$ AND WHY!
- Interpret your results in the context of the problem

**Problem**: Write the conclusions and interpretations for the previous scenarios using our results.
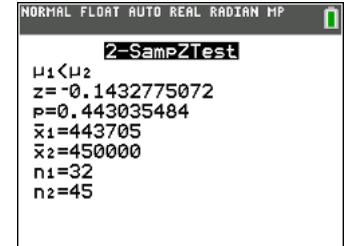
**Setup:** A prospective home buyer is trying to decide which city to move to. In Washington D.C, a random sample of 32 foreclosed homes sold for an average of $443,705 and in Baltimore and random sample of 45 foreclosed homes sold for $450,000. Real estate experts say the standard deviation for sales across the nation is $190,000. Is there sufficient evidence that the Baltimore housing market is more expensive? Use $\alpha = 0.1$

**Solution**:
*Need these…*

Let $\mu_1$ = the true mean sale price of foreclosed homes in DC
Let $\mu_2$ = the true mean sale price of foreclosed homes in Baltimore

$H_0: \mu_1 - \mu_2 = 0$
$H_A: \mu_1 - \mu_2 < 0$     $\alpha = 0.1$

```
NORMAL FLOAT AUTO REAL RADIAN MP
          2-SampZTest
µ1<µ2
z=-0.1432775072
p=0.443035484
x̄1=443705
x̄2=450000
n1=32
n2=45
```

<u>P-Value</u>

P-value = 2–SampZTest($\sigma_1$ = 190000, $\sigma_2$ = 190000, $\bar{x}_1$ = 443705, $n_1$ = 32, $\bar{x}_2$ = 450000, $n_2$ = 45, $\mu_1 < \mu_2$) = 0.443
p-value = 0.443 < 0.10 = $\alpha$ → Fail to reject $H_0$

*There is two parts that we need 1) Conclusion and 2) Interpretation*

<u>Conclusion and Interpretation</u>
*We fail to reject the null hypothesis because the p-value is greater than the significance level →This is the correct decision, BUT if this is all you write, you are MISSING the entire INTERPRETATION part; and should be MORE SPECIFIC!! What are the the p-value and significance level???*

*This would be better for the CONCLUSION (only) → We fail to reject the null hypothesis because the p-value = 0.443 is greater than the significance level 0.1 → SHOW ME YOU KNOW WHAT YOU'RE DOING!*

*Can also word the CONCLUSION part like this, which is correct as long as we have all the needed info! → Because the p-value 0.443 is greater than $\alpha$ = 0.1, we fail to reject the null Hypothesis*

*Now here is the Interpretation part, which needs to be right after our CORRECT conclusion from above*
*We do not have sufficient evidence to conclude that the prices in Baltimore are greater than DC → Almost there! Correctly said that there is NOT sufficient evidence and talked about the Alternative Hypothesis very good and good context. But MISSING the PARAMETER TRUE MEAN*
*There is NOT sufficient evidence to claim that the homes in Baltimore are more expensive than homes sold in DC → Same thing, MISSING TRUE AVERAGE*
*There is NOT sufficient evidence to conclude that the true mean sale price of foreclosed homes in Baltimore is more than that of DC → NOW this is CORRECT!*
*We do not have sufficient evidence to conclude that the true average price of Baltimore homes are greater than DC homes → This would also be CORRECT!*

# Hypothesis Tests for Means – INDEPENDENT Samples (and Unknown $\sigma$)

- Now we will go over <u>UNKNOWN population standard deviations!</u>

    - This is the scenario when ONLY SAMPLE standard deviations are given

- All of the previous Two Sample overview applies and the Means Hypotheses, Conditions and Interpretations are the same

    - Our test is just based on the T distribution now!

- This is still for **independent** samples!

# Using Calc - Test Statistic and P-Value for Ind Means and Unknown $\sigma$

### 4. Compute value of Test Statistic / P-value.

**Setup**
A prospective home buyer is trying to decide which city to move to. In Washington D.C, a random sample of 32 foreclosed homes sold for an average of $443,705 with standard deviation $150,000 and in Baltimore and random sample of 45 foreclosed homes sold for $450,000 with standard deviation $130,000. Is there sufficient evidence that the Baltimore housing market is more expensive? Use $\alpha = 0.1$

<u>GOAL</u>: Conduct a Hypothesis Test!

1. 2–SampTTest
   a) Input = Stats
   b) $\bar{x}_1$ = sample 1 mean
   c) $Sx_1$ = sample 1 SD
   d) $n_1$ = sample size 1
   e) $\bar{x}_2$ = sample 2 mean
   f) $Sx_2$ = sample 2 SD
   g) $n_2$ = sample size 2
   h) $\mu_1$: Alternative hypothesis with $\mu_2$ → ** *NOTE this is <u>not</u> in terms of the <u>difference</u>*
   i) Pooled: No → ** *ALWAYS keep as NO for this class*
   Calculate or Draw

Formula for $t_{stat}$ by hand:

$$t = \frac{(\bar{x}_1 - \bar{x}_2) - (\mu_1 - \mu_2)}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}$$

Under the Null hypothesis, the quantity $\mu_1 - \mu_2 = 0$ and everything else is known.

*** <u>Pooled Standard Deviation</u>
- We have the option to use a weighted average of the two sample standard deviations
- This is appropriate if the two values are similar and results in a slightly better test
- **But we are going to keep it simple and ALWAYS NOT pool the SDs**

*** <u>Degrees of Freedom</u>
- There is a difficult fancy way to determine the DF for two sample T-Tests (which our calc does) that we aren't going to do
- **So only going to be making conclusions using the p-value method**

# Using Calc - Test Statistic and P-Value for Ind Means and Unknown $\sigma$

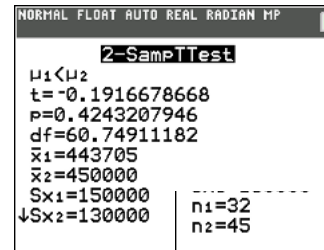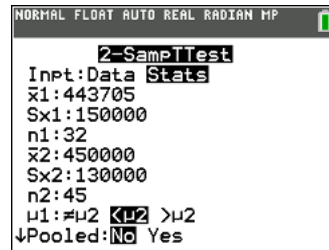## 4. Compute value of Test Statistic / P-value.

**Setup**
A prospective home buyer is trying to decide which city to move to. In Washington D.C, a random sample of 32 foreclosed homes sold for an average of $443,705 with standard deviation $150,000 and in Baltimore and random sample of 45 foreclosed homes sold for $450,000 with standard deviation $130,000. Is there sufficient evidence that the Baltimore housing market is more expensive? Use $\alpha = 0.1$

Formula for $t_{stat}$ by hand:

$$t = \frac{(\bar{x}_1 - \bar{x}_2) - (\mu_1 - \mu_2)}{\sqrt{\dfrac{s_1^2}{n_1} + \dfrac{s_2^2}{n_2}}}$$

Under the Null hypothesis, the quantity $\mu_1 - \mu_2 = 0$ and everything else is known.

**GOAL**: Conduct a Hypothesis Test!
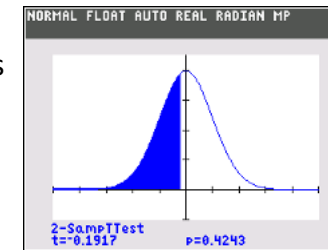
1. 2–SampTTest
   a) Input = Stats
   b) $\bar{x}_1$ = sample 1 mean
   c) $Sx_1$ = sample 1 SD
   d) $n_1$ = sample size 1
   e) $\bar{x}_2$ = sample 2 mean
   f) $Sx_2$ = sample 2 SD
   g) $n_2$ = sample size 2
   h) $\mu_1$: Alternative hypothesis with $\mu_2$ → ** NOTE this is not in terms of the difference
   i) Pooled: No → ** ALWAYS keep as NO for this class
   Calculate or Draw

NORMAL FLOAT AUTO REAL RADIAN MP
2-SampTTest
Inpt:Data **Stats**
x̄1:443705
Sx1:150000
n1:32
x̄2:450000
Sx2:130000
n2:45
μ1:≠μ2 **<μ2** >μ2
↓Pooled:**No** Yes

NORMAL FLOAT AUTO REAL RADIAN MP
2-SampTTest
μ1<μ2
t=-0.1916678668
p=0.4243207946
df=60.74911182
x̄1=443705
x̄2=450000
Sx1=150000
↓Sx2=130000
n1=32
n2=45

$(\mu_1 = DC \text{ and } \mu_2 = Baltimore)$
$H_0: \mu_1 - \mu_2 = 0$
$H_A: \mu_1 - \mu_2 < 0$

**Calculate Output**
$\mu_1 \neq <> \mu_2$ Alternative hypothesis
$t = t_{stat}$
$p$ = p-value
df = pooled degrees of freedom
$\bar{x}_1$ = sample 1 mean
$\bar{x}_2$ = sample 2 mean
$Sx_1$ = sample  2 SD
$Sx_2$ = sample  2 SD
$n_1$ = sample size 1
$n_2$ = sample size 2

NORMAL FLOAT AUTO REAL RADIAN MP
2-SampTTest
t=-0.1917    p=0.4243

**Draw Output**
Plot (and displays values) of p = p-value and t = $t_{stat}$ on the standard normal curve

*** Degrees of Freedom
- There is a difficult fancy way to determine the DF for two sample T-Tests (which our calc does) that we aren't going to do
- **So only going to be making conclusions using the p-value method**

*** Pooled Standard Deviation
- We have the option to use a weighted average of the two sample standard deviations
- This is appropriate if the two values are similar and results in a slightly better test
- **But we are going to keep it simple and ALWAYS NOT pool the SDs**

# Confidence Intervals – Two Samples!

- Everything we learned about Confidence Intervals (the different pieces, interpretation, etc.) still applies!

- Now we are just trying to estimate the DIFFERENCE between the two parameters!

# Structure of a Two Sample Confidence Interval

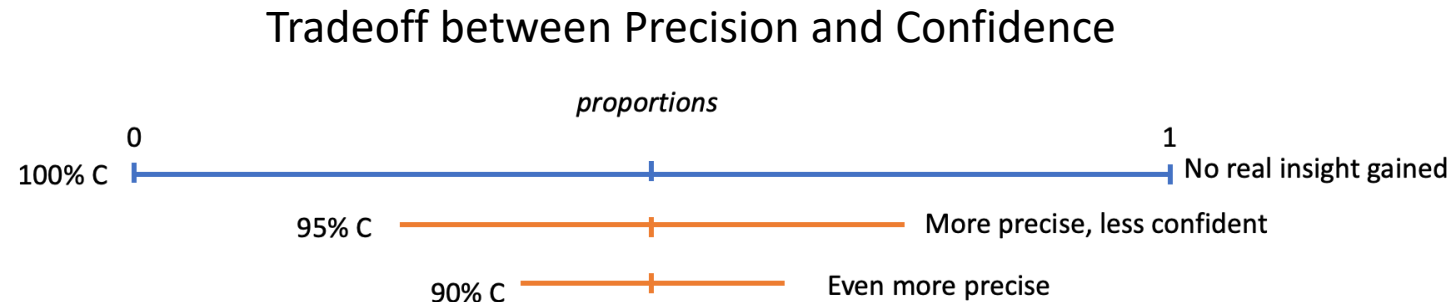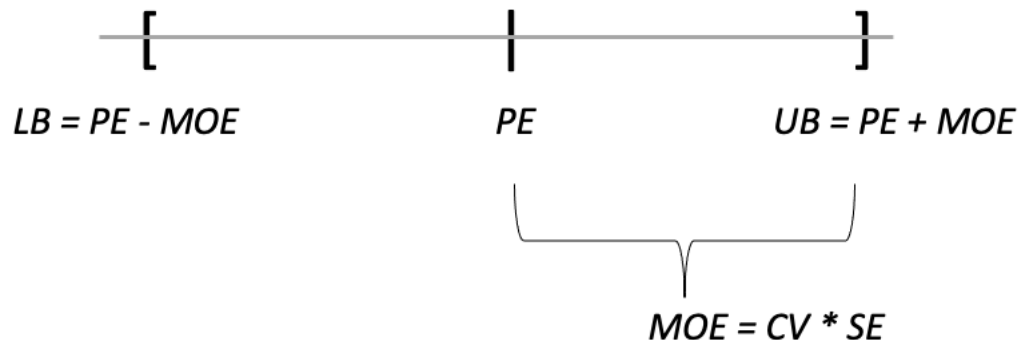SAME structure as the One Sample Confidence Intervals we learned previously

## C.I. = Point Estimate ± Margin of Error

<u>But Now with Two Samples</u>

- Point Estimate is your best guess of the DIFFERENCE; at the center of the interval.

- Margin of Error (MOE) = Critical Value (CV) * Standard Error (SE).
  - CV are the exact same!
  - SE formulas are slightly different because we have an additional sample

- Same relationships with MOE: Smaller MOE, more precise your estimate of the difference is.
  - The more confident, the wider your interval is (if everything else stays the same)

LB = PE - MOE    PE    UB = PE + MOE

MOE = CV * SE

Tradeoff between Precision and Confidence

*proportions*

0                                               1

100% C |————————————————|————————————————| No real insight gained

95% C    |————————————|————————————| More precise, less confident

90% C       |————————|————————| Even more precise

# Final Confidence Interval for $p_1 - p_2$

## 2 Proportion Z Interval

**Recall**: Our point estimate is the sample proportion $\hat{p}_1 = \frac{x_1}{n_1}$, which represents the number of success divided by the sample size, same for the second sample

C.I. = Point Estimate ± Margin of Error

$= (\hat{p}_1 - \hat{p}_2) \pm Z^* \sigma_{\hat{p}_1 - \hat{p}_2}$

$= (\hat{p}_1 - \hat{p}_2) \pm Z^* \sqrt{\dfrac{\hat{p}_1(1-\hat{p}_1)}{n_1} + \dfrac{\hat{p}_2(1-\hat{p}_2)}{n_2}}$

**TWO Sample Interval Conditions**

✓Randomization Condition
　　　Need to have two random samples
✓Independence Condition
　　　Need to independent samples
✓Large Enough Sample Condition
　　　AT LEAST 5 successes and 5 failures in EACH collected sample

## Using Calc

**GOAL**: Find the Two Sample Confidence Interval for Difference in Proportions!

**Setup**
A NatGeo Poll interviewed 1200 hiking enthusiasts and 1100 climbers. They asked "Are you more afraid of spiders or snakes???" 768 of the 1200 hikers and 662 of the 1100 climbers, responded "Ewww, snakes…." **Calculate** a 95% Confidence Interval for the difference in proportions.

*** Have to state which parameter is 1 and which is 2*

2-PropZInt

　a)　　$x_1$ = # of successes (people that said yes)  in sample 1

　b)　　$n_1$ = sample size

　c)　　$x_2$ = # of successes in sample 2

　d)　　$n_2$ = sample size 2

　e)　　C-Level = Confidence level (as a decimal or whole number, both work)

# Final Confidence Interval for $p_1 - p_2$

## 2 Proportion Z Interval

Recall: Our point estimate is the sample proportion $\hat{p}_1 = \frac{x_1}{n_1}$, which represents the number of success divided by the sample size, same for the second sample

C.I. = Point Estimate ± Margin of Error

$$= (\hat{p}_1 - \hat{p}_2) \pm Z^* \sigma_{\hat{p}_1 - \hat{p}_2}$$

$$= (\hat{p}_1 - \hat{p}_2) \pm Z^* \sqrt{\frac{\hat{p}_1(1-\hat{p}_1)}{n_1} + \frac{\hat{p}_2(1-\hat{p}_2)}{n_2}}$$

**TWO Sample Interval Conditions**

✓Randomization Condition
- Need to have two random samples

✓Independence Condition
- Need to independent samples

✓Large Enough Sample Condition
- AT LEAST 5 successes and 5 failures in EACH collected sample

## Using Calc

**GOAL**: Find the Two Sample Confidence Interval for Difference in Proportions!

**Setup**
A NatGeo Poll interviewed 1200 hiking enthusiasts and 1100 climbers. They asked "Are you more afraid of spiders or snakes???" 768 of the 1200 hikers and 662 of the 1100 climbers, responded "Ewww, snakes...." **Calculate** a 95% Confidence Interval for the difference in proportions.

*** Have to state which parameter is 1 and which is 2*

$p_1 \rightarrow$ hikers
$p_2 \rightarrow$ climbers

2-PropZInt

a) $x_1$ = # of successes (people that said yes) in sample 1

b) $n_1$ = sample size

c) $x_2$ = # of successes in sample 2

d) $n_2$ = sample size 2

e) C-Level = Confidence level (as a decimal or whole number, both work)



NORMAL FLOAT AUTO REAL RADIAN MP
2-PropZInt
x1:768
n1:1200
x2:662
n2:1100
C-Level:95
Calculate

NORMAL FLOAT AUTO REAL RADIAN MP
2-PropZInt
(-0.0015,0.07786)
p̂1=0.64
p̂2=0.6018181818
n1=1200
n2=1100

# Final Confidence Interval for $\mu_1 - \mu_2$ INDEPENDENT Samples

<u>2 Sample Z Interval</u> – KNOWN **σ**s

C.I. = Point Estimate ± Margin of Error

$$= (\bar{x}_1 - \bar{x}_2) \pm Z^* \sigma_{\bar{x}_1 - \bar{x}_2}$$

$$= (\bar{x}_1 - \bar{x}_2) \pm Z^* \sqrt{\frac{\sigma_1}{n_1} + \frac{\sigma_2}{n_2}}$$

<u>Using Calc</u>

<u>**GOAL**</u>: Find the Two Sample Confidence Interval for Difference in Means!!

2-SampZInt

a)   Input = Stats
b)   **σ**$_1$ = population 1 standard deviation
c)   **σ**$_2$ = population 1 standard deviation
d)   $\bar{x}_1$ = sample mean 1
e)   $n_1$ = sample size
f)   $\bar{x}_2$ = sample mean 2
g)   $n_2$ = sample size
h)   C-Level = Confidence level (as a decimal or whole number, both work)

---

* Same Critical Value as with a
2 Proportion Z Interval

---

**<u>INDEPENDENT Samples Interval Conditions</u>**

✓Randomization Condition
   Need to have two random samples
✓Independence Condition
   Need to independent samples
✓Large Enough Sample Condition
   Normal populations OR
   $n_1 \geq 30$ AND $n_2 \geq 30$

---

<u>Setup</u>
A prospective home buyer is trying to decide which city to move to. In Washington D.C, a random sample of 32 foreclosed homes sold for an average of $443,705 and in Baltimore and random sample of 45 foreclosed homes sold for $450,000. Real estate experts say the standard deviation for sales across the nation is $190,000. **Calculate** a 85% Confidence Interval for the difference in means.

*** Have to state which parameter is 1 and which is 2*

# Final Confidence Interval for $\mu_1 - \mu_2$ INDEPENDENT Samples

## 2 Sample Z Interval – KNOWN $\sigma$s

C.I. = Point Estimate ± Margin of Error

$$= (\bar{x}_1 - \bar{x}_2) \pm Z^* \sigma_{\bar{x}_1 - \bar{x}_2}$$

$$= (\bar{x}_1 - \bar{x}_2) \pm Z^* \sqrt{\frac{\sigma_1}{n_1} + \frac{\sigma_2}{n_2}}$$

## Using Calc

**GOAL**: Find the Two Sample Confidence Interval for Difference in Means!!

2-SampZInt

a)   Input = Stats
b)   $\sigma_1$ = population 1 standard deviation
c)   $\sigma_2$ = population 1 standard deviation
d)   $\bar{x}_1$ = sample mean 1
e)   $n_1$ = sample size
f)   $\bar{x}_2$ = sample mean 2
g)   $n_2$ = sample size
h)   C-Level = Confidence level (as a decimal or whole number, both work)

---

* Same Critical Value as with a
2 Proportion Z Interval

---

**INDEPENDENT Samples Interval Conditions**

✓ Randomization Condition
    Need to have two random samples
✓ Independence Condition
    Need to independent samples
✓ Large Enough Sample Condition
    Normal populations OR
    $n_1 \geq 30$ AND $n_2 \geq 30$

---

**Setup**
A prospective home buyer is trying to decide which city to move to. In Washington D.C, a random sample of 32 foreclosed homes sold for an average of $443,705 and in Baltimore and random sample of 45 foreclosed homes sold for $450,000. Real estate experts say the standard deviation for sales across the nation is $190,000. **Calculate** a 85% Confidence Interval for the difference in means.

*\*\* Have to state which parameter is 1 and which is 2*

$\mu_1 \rightarrow$ DC
$\mu_2 \rightarrow$ Baltimore



NORMAL FLOAT AUTO REAL RADIAN MP
2-SampZInt
Inpt:Data **Stats**
σ1:190000
σ2:190000
x̄1:443705
n1:32
x̄2:450000
n2:45
C-Level:85
Calculate

NORMAL FLOAT AUTO REAL RADIAN MP
2-SampZInt
(-69542,56952)
x̄1=443705
x̄2=450000
n1=32
n2=45

# Final Confidence Interval for $\mu_1 - \mu_2$ INDEPENDENT Samples

## 2 Sample t Interval – UNKNOWN σs

C.I. = Point Estimate ± Margin of Error

$$= (\bar{x}_1 - \bar{x}_2) \pm t^* \, \sigma_{\bar{x}_1 - \bar{x}_2}$$

$$= (\bar{x}_1 - \bar{x}_2) \pm t^* \sqrt{\frac{s_1}{n_1} + \frac{s_2}{n_2}}$$

Using Calc

**GOAL**: Find the Two Sample Confidence Interval for Difference in Means!!

2-SameTInt

a)   Input = Stats
b)   $\bar{x}_1$ = sample mean 1
c)   $Sx_1$ = population 1 standard deviation
d)   $n_1$ = sample size 1
e)   $\bar{x}_2$ = sample mean 2
f)   $Sx_2$ = population 2 standard deviation
g)   $n_2$ = sample size 2
h)   Pooled = No
i)   C-Level = Confidence level (as a decimal or whole number, both work)

* Same Critical Value is based on t distribution now

*** Degrees of Freedom
• There is a difficult fancy way to determine the DF for two sample, (which our calc does) that we wouldn't do by hand

## INDEPENDENT Samples Interval Conditions

✓Randomization Condition
    Need to have two random samples
✓Independence Condition
    Need to independent samples
✓Large Enough Sample Condition
    Normal populations OR
    $n_1 \geq 30$ AND $n_2 \geq 30$

**Setup**
A prospective home buyer is trying to decide which city to move to. In Washington D.C, a random sample of 32 foreclosed homes sold for an average of $443,705 with standard deviation $150,000 and in Baltimore and random sample of 45 foreclosed homes sold for $450,000 with standard deviation $130,000. **Calculate** a 85% Confidence Interval for the difference in means.

*** Have to state which parameter is 1 and which is 2*

*** Pooled Standard Deviation
• We have the option to use a weighted average of the two sample standard deviations
• This is appropriate if the two values are similar and results in a slightly more precise interval
• **But we are going to keep it simple and ALWAYS NOT pool the SDs**

# Final Confidence Interval for $\mu_1 - \mu_2$ INDEPENDENT Samples

## 2 Sample t Interval – UNKNOWN σs

C.I. = Point Estimate ± Margin of Error

$$= (\bar{x}_1 - \bar{x}_2) \pm t^* \, \sigma_{\bar{x}_1 - \bar{x}_2}$$

$$= (\bar{x}_1 - \bar{x}_2) \pm t^* \sqrt{\frac{s_1}{n_1} + \frac{s_2}{n_2}}$$

Using Calc

**GOAL**: Find the Two Sample Confidence Interval for Difference in Means!!

2-SameTInt

- a) Input = Stats
- b) $\bar{x}_1$ = sample mean 1
- c) $Sx_1$ = population 1 standard deviation
- d) $n_1$ = sample size 1
- e) $\bar{x}_2$ = sample mean 2
- f) $Sx_2$ = population 2 standard deviation
- g) $n_2$ = sample size 2
- h) Pooled = No
- i) C-Level = Confidence level (as a decimal or whole number, both work)

*** Pooled Standard Deviation
- We have the option to use a weighted average of the two sample standard deviations
- This is appropriate if the two values are similar and results in a slightly more precise interval
- **But we are going to keep it simple and ALWAYS NOT pool the SDs**

---

* Same Critical Value is based on t distribution now

*** Degrees of Freedom
- There is a difficult fancy way to determine the DF for two sample, (which our calc does) that we wouldn't do by hand

---

**INDEPENDENT Samples Interval Conditions**

✓ Randomization Condition
   Need to have two random samples
✓ Independence Condition
   Need to independent samples
✓ Large Enough Sample Condition
   Normal populations OR
   $n_1 \geq 30$ AND $n_2 \geq 30$

---

**Setup**
A prospective home buyer is trying to decide which city to move to. In Washington D.C, a random sample of 32 foreclosed homes sold for an average of $443,705 with standard deviation $150,000 and in Baltimore and random sample of 45 foreclosed homes sold for $450,000 with standard deviation $130,000. **Calculate** a 85% Confidence Interval for the difference in means.

*** Have to state which parameter is 1 and which is 2*

$\mu_1 \rightarrow$ DC
$\mu_2 \rightarrow$ Baltimore

NORMAL FLOAT AUTO REAL RADIAN MP
```
2-SampTInt
Inpt:Data Stats
x̄1:443705
Sx1:150000
n1:32
x̄2:450000
Sx2:13000
n2:45
C-Level:85
↓Pooled:No Yes
```

NORMAL FLOAT AUTO REAL RADIAN MP
```
2-SampTInt
(-45530.32940)
df=31.33141118
x̄1=443705
x̄2=450000
Sx1=150000
Sx2=13000
n1=32
n2=45
```

# Hypothesis Tests and Confidence Intervals

What each type of inference tells us

- For One Sample
    - Hypothesis tests tell us how a parameter compares to a specific value (greater than, less than, or not equal to)
    - Confidence intervals give us a range of plausible values
- For Two Sample
    - Hypothesis Tests tell us if there is a difference between the two parameters
    - Confidence Intervals give us a range of plausible values for this difference

Both of these inference methods can be used together!

Example

- A hypothesis test on a random sample of 200 American adults found that greater than 50% of them have tried marijuana
    - *This conclusion just tells us that the true proportion is **somewhere above 50%***

- A confidence interval could be constructed to find how much more than 50% of American adults have tried marijuana
    - *Lets say CI = (0.52, 0.60), then we have a specific estimate as to where the true proportion actually is, from 52% to 60%. This tells us **exactly where** we think the parameter is and **how much greater than 50%** it is! MORE INFORMATION*

CI for Testing

- Confidence Intervals actually give us enough information to say whether or not we would reject a corresponding Hypothesis Test!

# LCQ – Confidence Intervals for Testing

**Problem**: Based on the parameters below, make a conclusion whether we would reject or fail to reject the Hypothesis Test below based on each of the following Confidence Intervals.

Let $\mu_1$ = population mean height of football teams in meters

Let $\mu_2$ = population mean height of soccer teams in meters

$$H_0: \mu_1 - \mu_2 = 0$$
$$H_A: \mu_1 - \mu_2 \neq 0$$

a) 85% CI for $\mu_1 - \mu_2 \rightarrow$ (-0.3, -0.05)

b) 90% CI for $\mu_1 - \mu_2 \rightarrow$ (-0.1, 0.2)

c) 90% CI for $\mu_1 - \mu_2 \rightarrow$ (0.03, 0.36)

# LCQ – Confidence Intervals for Testing

**Problem**: Based on the parameters below, **make** a conclusion whether we would reject or fail to reject the Hypothesis Test below based on each of the following Confidence Intervals and **explain** why.

Let $\mu_1$ = population mean height of football teams in meters
Let $\mu_2$ = population mean height of soccer teams in meters

$H_0$: $\mu_1 - \mu_2 = 0$
$H_A$: $\mu_1 - \mu_2 \neq 0$

$(- , -)$

a) 85% CI for $\mu_1$ - $\mu_2$ → (-0.3, -0.05)
*REJECT→ because the entire interval is below zero!* Here's the long description of why we would reject:

- *A confidence interval gives a range of plausible values for the parameter we are estimating, in this case it is the difference between $\mu_1$ and $\mu_2$*
- *Under the Null, we are assuming this difference is zero (so they are equivalent)*
- *Well our entire CI is below this null difference of zero, which means it is NOT a plausible value based on our results!*
- *So this would be enough evidence to show that these two parameters are different*
- *And also we know that the entire interval is negative. So based on the order of subtraction in the hypotheses, this indicates that $\mu_2$ is LARGER! Average heights of soccer teams are larger based on this interval.*
- *This conclusion would MATCH the conclusion if we actually did the Hypothesis Test for this. And we would get a negative Test Statistic, also indicating that $\mu_2$ is bigger*

$\mu_1 < \mu_2$    $H_0: \mu_1 = \mu_2$

-0.3    -0.175    -0.05    0

Lower bound    $\bar{x}_1 - \bar{x}_2$    Upper bound    $\mu_1 - \mu_2$

$(-, +)$

b) 90% CI for $\mu_1$ - $\mu_2$ → (-0.1, 0.2)
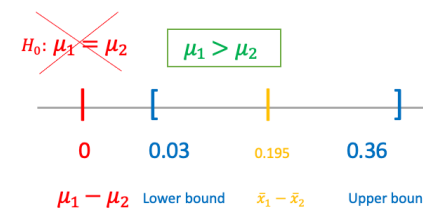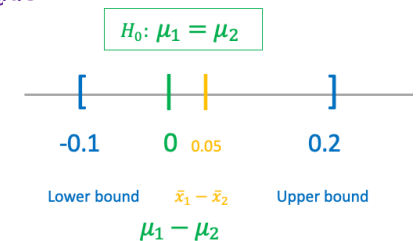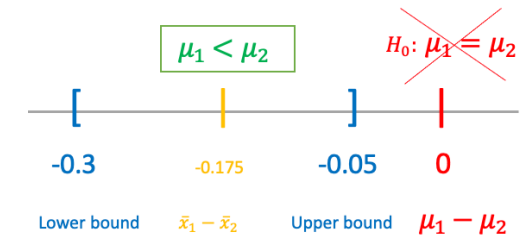*FAIL TO REJECT → because the entire interval is above zero!*
- *Now our interval starts negative, crosses zero and end positive (so zero is contained in the interval)*
- *This says that zero is a plausible value for the difference of these two parameters, which means it's possible that they are equivalent!*
- *We can NOT say that one is ALWAYS larger than the other, so we would NOT be able to conclude they are significantly different → so would fail to reject the Null*

$H_0: \mu_1 = \mu_2$

-0.1    0  0.05    0.2

Lower bound    $\bar{x}_1 - \bar{x}_2$    Upper bound
$\mu_1 - \mu_2$

$(+ , +)$

c) 95% CI for $\mu_1$ - $\mu_2$ → (0.03, 0.36)
*REJECT→ because the entire interval is above zero!*
- *Our interval suggest that zero is not a possible value for the difference of these two parameters, so it's not possible that they are equivalent based on our results*
- *That is exactly what rejecting the null hypothesis tells us!*
- *And because the interval is positive, this time $\mu_2$ is LARGER (subtracting a smaller number from a larger number results in a positive difference) so the average heights of football teams are greater than soccer teams*

$H_0: \mu_1 = \mu_2$    $\mu_1 > \mu_2$

0    0.03    0.195    0.36

$\mu_1 - \mu_2$  Lower bound    $\bar{x}_1 - \bar{x}_2$    Upper boun

# Two Sample Confidence Interval Interpretations

General Structure

- Here was the structure for ONE sample CI:
    - I am % confident that the true/population parameter + context is between (lower bound) and (upper bound).
- It is not as structured when discussing **TWO parameters**, but we still have the same key parts below

3 Pieces

1. 95% Confident: This is a Confidence Statement
    - Tells us what percent off ALL possible samples result in a CI that captures the true proportion.
2. Parameters + Context: We are talking about TWO parameters.
    - But what parameters??? We ALWAYS need context.
    - Now because we have TWO parameters, we can be specific about which one is greater / less than!
3. Interval: The range of plausible values for the DIFFERENCE!
    - Uses the difference of our sample statistics and the MOE based on Two Sample Standard Errors.

** Wording gets a little tricky when the CI of the difference contains zero (negative lower bound and positive upper bound)… but nonetheless same logic

Example

Let $p_1$ = true proportion of Columbus males who enjoy running
Let $p_2$ = true proportion of Columbus females who enjoy running
95% CI for $p_1 - p_2$ = (0.05, 0.25)

- We are 95% confident that the true proportion of all Columbus males who enjoy running is between 0.05 and 0.25 greater than the true proportion of females who enjoy running.
- *Or equivalently but slightly shorter* → We are 95% confident that the true proportion of all Columbus males who enjoy running is between 0.05 and 0.25 greater than that of females.

# LCQ – Two Sample CI Interpretations

**Problem**: Based on the parameters below, interpret each of the following Confidence Intervals.

Let $\mu_1$ = population mean height of football teams (meters)

Let $\mu_2$ = population mean height of soccer teams (meters)

a) 85% CI for $\mu_1 - \mu_2 \rightarrow$ (-0.3, -0.05)

b) 90% CI for $\mu_1 - \mu_2 \rightarrow$ (-0.1, 0.2)

c) 95% CI for $\mu_1 - \mu_2 > 0 \rightarrow$ (0.03, 0.36)

# LCQ – Two Sample CI Interpretations

**Problem**: Based on the parameters below, interpret each of the following Confidence Intervals.

Let $\mu_1$ = population mean height of football teams (meters)
Let $\mu_2$ = population mean height of soccer teams (meters)

a) 85% CI for $\mu_1 - \mu_2 \rightarrow$ (-0.3, -0.05) $\rightarrow$ *We already talked how this interval indicates $\mu_1 < \mu_2$. So we can <u>phrase our CI interpretation using this knowledge</u>:*

We are 85% confident that the population mean height of <u>football teams</u> is between 0.05 and 0.3 meters <u>less than</u> the population mean height of <u>soccer teams</u>

*Or equivalently (rearranging with $\mu_2 > \mu_1$):* We are 85% confident that the population mean height of <u>soccer teams</u> is between 0.05 and 0.3 meters <u>greater than</u> the population mean height of <u>football teams</u>

b) 90% CI for $\mu_1 - \mu_2 \rightarrow$ (-0.1, 0.2) $\rightarrow$ *Here we can't say definitively that one mean is greater, so our wording needs to reflect that $\mu_1$ can be less than OR greater than $\mu_2$ (this wording is a little less straightforward than before)*

We are 90% the that true mean height of <u>football teams</u> is between <u>0.1 meters less than</u> OR <u>0.2 meters greater than</u> the true mean height of <u>soccer teams</u>

*Or equivalently (rearranging with $\mu_2 - \mu_1$):* We are 90% confident that the population mean height of <u>soccer teams</u> is between <u>0.2 meters less than</u> OR <u>0.1 meters greater than</u> the population mean height of <u>football teams</u>

c) 95% CI for $\mu_1 - \mu_2 > 0 \rightarrow$ (0.03, 0.36) $\rightarrow$ *We know this interval indicates $\mu_1 > \mu_2$. So we can again <u>phrase our CI interpretation in that way</u>:*

We are 95% confident that the true mean height of <u>football teams</u> is between 0.03 and 0.36 meters <u>taller than</u> the true mean height of <u>soccer teams</u>

*Or equivalently (rearranging with $\mu_2 < \mu_1$):* We are 95% confident that the population mean height of <u>soccer players</u> is between 0.03 and 0.36 meters <u>shorter than</u> the population mean height of <u>football players</u>

# LCQ 2 – Two Sample CI Interpretations

**Problem**: Based on the parameters below, interpret each of the following Confidence Intervals. Then determine if you would reject or fail to the corresponding Hypothesis Test.

Let $\mu_1$ = population mean price ($) at an Italian restaurant

Let $\mu_2$ = population mean price ($) at a Mexican restaurant

$H_0$: $\mu_1 - \mu_2 = 0$

$H_A$: $\mu_1 - \mu_2 \neq 0$

a) 85% CI for $\mu_1 - \mu_2 \rightarrow$ (-30, -3)

b) 90% CI for $\mu_1 - \mu_2 \rightarrow$ (-5, 20)

c) 95% CI for $\mu_1 - \mu_2 > 0 \rightarrow$ (3, 10)

# LCQ 2 – Two Sample CI Interpretations

**Problem**: Based on the parameters below, interpret each of the following Confidence Intervals. Then determine if you would reject or fail to the corresponding Hypothesis Test.

Let $\mu_1$ = population mean price ($) at an Italian restaurant

Let $\mu_2$ = population mean price ($) at a Mexican restaurant

$H_0$: $\mu_1 - \mu_2 = 0$

$H_A$: $\mu_1 - \mu_2 \neq 0$

a) 85% CI for $\mu_1 - \mu_2 \rightarrow$ (-30, -3)   **Reject test**! Difference of zero $\leftrightarrow$ equal is NOT plausible based on interval (not captured)

-30   -3   0

*Options*

1.   *I am 85% confident that the true mean price at the Italian restaurant is less than that of the true mean price at the Mexican restaurant $\rightarrow$ Correct but MISSING the values of the CI! The whole goal of a CI is to find the range of plausible values, so use them in the interpretation!!*

2.   *I am 85% confident that the true mean price at the Italian restaurant is less (-30, -3) than that of the true mean price at the Mexican restaurant $\rightarrow$ Correct, but the phrasing for the values should be IMPROVED!! This is NOT how we would try to talk to someone when comparing the prices of two restaurants! Make it flow, how we would naturally speak it*

3.   *We are 85% confident that the true mean price of dinner at an Italian restaurant is between $3 and $30 less expensive than the true mean price of dinner at a Mexican restaurant $\rightarrow$ PERFECT!!!, saying 'less than' takes care of the negatives and now it reads much better!*

4.   *We are 85% confident that the true mean Mexican food price is $3 to $30 less expensive than the true mean of Italian food price $\rightarrow$ WRONG!! We have a negative interval for the subtraction, so the second parameter must be larger LARGER (order matters in our interpretation). We could flip it to say it like option 5*

5.   *We are 85% confident that the true mean Mexican food price is $3 to $30 MORE expensive than the true mean of Italian food price $\rightarrow$ Now this is CORRECT!! Italian less expensive is equivalent to Mexican more expensive!*

# LCQ 2 – Two Sample CI Interpretations

**Problem**: Based on the parameters below, interpret each of the following Confidence Intervals. Then determine if you would reject or fail to the corresponding Hypothesis Test.
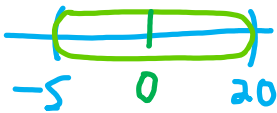
Let $\mu_1$ = population mean price ($) at an Italian restaurant
Let $\mu_2$ = population mean price ($) at a Mexican restaurant

$H_0: \mu_1 - \mu_2 = 0$
$H_A: \mu_1 - \mu_2 \neq 0$

b) 90% CI for $\mu_1 - \mu_2 \rightarrow$ (-5, 20)          Fail to Reject test! Difference of zero $\leftrightarrow$ equal IS plausible based on interval (contained within)

*Process*

*This wording is more tricky than before because we have a negative lower bound and positive upper bound, not strictly less expensive or strictly more expensive.*
*But we can word our interpretation one bound at a time like so and put the pieces together:*
1. *Start with $\mu_1$ = Italian → We are 90% confident that the true mean price of an Italian restaurant is between*
2. *Now talk about the negative lower bound -5 → (Italian) $5 less expensive (than Mexican)*
3. *Now positive upper bound 20 → $20 more expensive (than Mexican)*
4. *End with $\mu_2$ = Mexican → than the true mean price of a Mexican restaurant*
*All  together 1+ 2 + 3 + 4 → We are 90% confident that the true mean price of an Italian restaurant is between $5 less expensive and $20 more expensive than the true mean price at a Mexican restaurant*

c) 95% CI for $\mu_1 - \mu_2 > 0 \rightarrow$ (3, 10)          Reject test! Difference of zero $\leftrightarrow$ equal is NOT plausible based on interval (NOT contained inside)

*Process*

*We can even do an entirely positive interval like this as well, one bound at a time like so and put the pieces together:*
1. *Start with $\mu_1$ = Italian → We are 95% confident that the true mean price of an Italian restaurant is between*
2. *Now talk about the negative lower bound 3 → (Italian) $3 more expensive (than Mexican)*
3. *Now positive upper bound 10 → $10 more expensive (than Mexican)*
4. *End with $\mu_2$ = Mexican → than the true mean price of a Mexican restaurant*
*All  together 1+ 2 + 3 + 4 → We are 95% confident that the true mean price of an Italian restaurant is between $3 more expensive and $10 more expensive than the true mean price at a Mexican restaurant*
*Can simplify wording a bit cause both are more expensive → We are 95% confident that the true mean price of an Italian restaurant is between $3 and $10 more expensive than the true mean price at a Mexican restaurant*

# LCQ 2 – Two Sample CI Interpretations

**Problem**: Based on the parameters below, interpret each of the following Confidence Intervals.

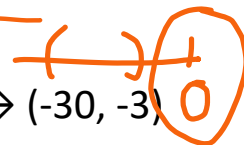Let $\mu_1$ = population mean price ($) at an Italian restaurant

Let $\mu_2$ = population mean price ($) at a Mexican restaurant

H0: $\mu_1 - \mu_2 = 0$
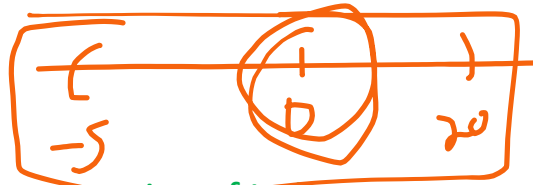
HA: $\mu_1 - \mu_2 \neq 0$   $\rightarrow$  $N \, b_T$

a) 85% CI for $\mu_1 - \mu_2 \rightarrow$ (-30, -3)  0

1. I am 85% confident that the true mean price at the Italian restaurant is less than that of the true mean price at the Mexican restaurant

2. We are 85% confident that the true mean price of dinner at an Italian restaurant is between $3 and $30 less expensive than the true mean price of dinner at a Mexican restaurant

3. We are 85% confident that the true mean Mexican food price is $3 to $30 MORE expensive than the true mean of Italian food price

b) 90% CI for $\mu_1 - \mu_2 \rightarrow$ (-5, 20)

- We are 90% confident that the true mean price of Italian restaurant is between $5 less expensive and $20 more expensive than the true mean price of Mexican

- We are 90% confident that the true mean price of a dinner at an Italian restaurant is between $5 less and $20 more than the true mean

# Problem Session!!!

# Problem 1

Are people waiting longer to marry? In 2007, a random sample of young adults (ages 18-31) showed that 468 of 1872 of those surveyed were married. In 2012, 581 of the 1940 young adults (ages 18-31) randomly surveyed were married. Is there any evidence to suggest the true proportion of young adults who are married has decreased? Use $\alpha$ = 0.10.

a) Define the parameters and state the hypotheses.

b) Check the conditions to run this test

c) Carry out the test

d) Calculate and interpret a 80% confidence interval for the difference in proportions

# Problem 1 - Solution

a) Let p1 = true proportion of married young adults in 2007 and
p2 = true proportion of married young adults in 2012

H0: p1 = p2                    vs.            Ha: p1 > p2
Also acceptable would be H0: p1 - p2 = 0 vs. Ha: p1 - p2 > 0

b) Check the assumptions:
- Whether or not someone is married is a categorical variable.
- It is stated that we have two random samples of young adults (ages 18-31)
- Whether or not one person is married does not affect whether or not others are married, so the groups are independent.
- The number of successes (those married) and failures (those not married) are at least 5 for both samples:
  - Sample 1 has 468 successes and 1872 – 468 = 1404 failures.
  - Sample 2 has 581 successes and 1940 – 581 = 1359 failures.

Significance Level: α = 0.10

# Problem 1 - Solution

Hypothesis test results:
$p_1$ : proportion of successes for population 1
$p_2$ : proportion of successes for population 2
$p_1 - p_2$ : Difference in proportions
$H_0 : p_1 - p_2 = 0$
$H_A : p_1 - p_2 > 0$

| Difference | Count1 | Total1 | Count2 | Total2 | Sample Diff. | Std. Err. | Z-Stat | P-value |
|---|---|---|---|---|---|---|---|---|
| $p_1 - p_2$ | 468 | 1872 | 581 | 1940 | -0.049484536 | 0.014469314 | -3.4199641 | 0.9997 |

c) The z test statistic = -3.420

The P-value = 0.9997

Since the P-value is greater than our significance level of 0.10, we fail to reject the null hypothesis.  There is not sufficient evidence to conclude that the true proportion of young adults who are married has decreased.

d) 80% Confidence Interval

I am 80% confident that  the true proportion of young adults who were married in 2007 is between 3.1% and 6.8% lower than the true proportion of young adults who were married in 2012.

Note: The result of the test would have been significant at $\alpha$ = 0.10 if the statement had called for a left-tailed alternative hypothesis as opposed to right-tailed! That is, if the question posed has been "Is there evidence to suggest the true proportion of young adults who are married has **increased**?"

# Problem 3

Is there a difference in the proportion of males and females who participate in Greek life at Miami University? A researcher collected data from a representative sample of students from Miami University and found that 31 out of 114 males and 63 out of 176 females participated in Greek life.  If appropriate, Test an appropriate hypothesis with α = 0.05.

# Problem 3 - Solution

Test to be ran: Two Proportion z-test

Parameters:
- $P_1$ = Population Proportion of Males that participate in Greek life at Miami
- $P_2$ = Population Proportion of Females that participate in Greek life at Miami

Hypotheses:
- $H_0$: $P_1 = P_2$      vs.      $H_a$: $P_1 \neq P_2$
- Or... $H_0$: $P_1 - P_2 = 0$      vs.      $H_a$: $P_1 - P_2 \neq 0$

Assumptions:
- Random sample of 114 males and 176 females
- Males and females are do not affect each other's probability of joining a fraternity or sorority. (They are independent)
- Each group has 5 successes and 5 failures in their sample
  - Males have 31 successes and 5 failures
  - Females have 63 successes and 113 failures

Significance Level: 5% or 0.05

# Problem 3 - Solution

Test Stat: -1.529

P-value: 0.1263

**Two sample proportion summary hypothesis test:**
$p_1$ : proportion of successes for population 1
$p_2$ : proportion of successes for population 2
$p_1 - p_2$ : Difference in proportions
$H_0 : p_1 - p_2 = 0$
$H_A : p_1 - p_2 \neq 0$

**Hypothesis test results:**

| Difference | Count1 | Total1 | Count2 | Total2 | Sample Diff. | Std. Err. | Z-Stat | P-value |
|---|---|---|---|---|---|---|---|---|
| $p_1 - p_2$ | 31 | 114 | 63 | 176 | -0.086024721 | 0.056270944 | -1.5287591 | 0.1263 |

Because our p-value of 0.1263 is greater than our significance level of 0.05, we fail to reject the null hypothesis. There is not sufficient evidence that the true population proportion of male Miami students participating in Greek life is different than the true population proportion of female Miami students participating in Greek life.