## <u>Unit 11 Slides</u>: Inference for $\mu_1 - \mu_2$ and $p_1 - p_2$





#### Case Studies:

- Is there an association between <u>childhood</u> <u>lead exposure</u> and <u>IQ</u>?
  - What is a <u>plausible range of values</u> for  $\mu_{lo} - \mu_{hi}$ , the **difference** in the <u>average IQ</u> <u>score of children with low lead level</u> <u>exposure</u> and the <u>average IQ score of</u> <u>children with high lead level exposure</u>?
  - Is there sufficient evidence to suggest  $\mu_{lo} \mu_{hi} \neq 0$ , ie. that there is a difference in the average IQ score of children with low lead level exposure and the average IQ score of children with high lead level exposure?
- Is there an association between <u>political</u> <u>party</u> and <u>approval for the direction the</u> country is going in (in 2017)?
  - What is a <u>plausible range of values</u> for *p<sub>dem</sub> - p<sub>rep</sub>*, the *difference* in the *proportion* of democrats that approve vs. the *proportion* of republicans that <u>approve</u>?
  - Is there sufficient evidence to suggest  $p_{dem} - p_{rep} \neq 0$ , ie. that there is a difference in the proportion of democrats that approve vs. the proportion of republicans that approve?

## Purpose of this Lecture:



In this lecture we will cover the following topics in **inference** and **probability**.

- 1. Two main types of inference for unknown population parameters.
- 2. How to better account for the additional uncertainty introduced by having to estimate additional parameters in probability and inference?
  - 2.1. Issues with plugging in s for  $\sigma$
  - 2.2. t-score of a sample mean
  - 2.3. distribution of t-scores (under certain conditions)

#### 2.4. t-distribution

- 3. Properties of the Sampling Distribution of Sample Mean Differences
  - 3.1. Mean
  - 3.2. Standard Deviation
  - 3.3. When is it normal?
- 4. Conducting Inference on a Population Mean Difference  $(\mu_1 \mu_2)$ 
  - 4.1. Creating a confidence interval for  $\mu_1 \mu_2$
  - 4.2. Conducting a hypothesis test to test the claim:  $\mu_1 \mu_2 
    eq 0$
  - 4.3. Conducting a hypothesis test to test the claim:  $\mu_1 \mu_2 \neq 0$  -with a p-value if you know  $\sigma_1$  and  $\sigma_2$
  - 4.4. Conducting a hypothesis test to test the claim:  $\mu_1 \mu_2 \neq 0$  -with a p-value if you DON'T know  $\sigma_1$ and  $\sigma_2$
  - 4.5. Conducting a hypothesis test to test the claim:  $\mu_1 \mu_2 \neq 0$  -with a confidence interval
- 5. Properties of Sampling Distribution of Sample Proportion Differences
  - 5.1. Mean
  - 5.2. Standard Deviation
  - 5.3. When is it normal?
- 6. Conducting Inference on a Population Proportion Difference  $(p_1 p_2)$ 
  - 6.1. Creating a confidence interval for  $p_1 p_2$
  - 6.2. Conducting a hypothesis test to test the claim:  $p_1 p_2 \neq 0$

#### Additional resources:

Sections 5.1, 5.3, 6.2 in Diez, Barr, and Cetinkaya-Rundel, (2015), *OpenIntro Statistics* https://www.openintro.org/download.php?file=os3&redirect=/stat/textbook/os3.php

## 1. Two Main Types of Inference for Unknown Population Parameters

Suppose we were interested in a <u>population parameter</u> of a large, unknown population. But all we can collect is a random sample from this population.

#### What we wish we could know:

Is there an association between <u>lead exposure</u> and <u>childhood IQ</u> for ALL children?



Is  $\mu_{lo} - \mu_{hi} \neq 0$ ?,

- $\mu_{lo}$ : average IQ score of ALL children with <u>low</u> lead level exposure
- $\mu_{hi}$ : average IQ score of ALL children with <u>high</u> lead level exposure

#### What we can answer:

- 1. **Confidence Interval:** What is a <u>plausible range of values</u> for the  $\mu_{lo} \mu_{hi}$ ?
- 2. Hypothesis Test: Is there sufficient evidence to suggest  $\mu_{lo} \mu_{hi} \neq 0$ ?

#### What we wish we could know:

Is there an association between <u>political party</u> and <u>opinion on the direction the country is going in</u> (in 2017) for ALL adults living in the U.S.?



## Is $p_{dem} - p_{rep} \neq 0$ ?,

- *p<sub>dem</sub>*: the **proportion** of ALL <u>democrats</u> living in the U.S. (in 2017) that approve of the direction the party is going in
- *p<sub>rep</sub>*: the **proportion** of ALL <u>republicans</u> living in the U.S. (in 2017) that approve of the direction the party is going in

#### What we can answer:

- 1. **Confidence Interval:** What is a <u>plausible range of values</u> for the  $p_{dem} p_{rep}$ ?
- 2. Hypothesis Test: Is there sufficient evidence to suggest  $p_{dem} p_{rep} \neq 0$ ?





2. How to better account for the additional uncertainty introduced by having to estimate additional parameters in probability and inference?



## BUT..... What if we don't know $\sigma$ ?

- - $\circ$  We can't conduct inference on  $\mu$  when \_\_\_\_\_.
  - ALL approximations where we plug in \_\_\_\_\_ for \_\_\_\_\_ are

#### • How do we make our "work around" for situations like these better?

- Use the **t-scores of sample means.**
- Use the **t-distribution** to calculate probabilities involving these sample means.

## 2.2. t-score of a sample mean

- We define the **t-score of a sample mean**  $\overline{x}_0$  as  $t = \frac{\overline{x}_0 \mu}{\frac{s}{\sqrt{n}}}$
- We define the **t-score of a sample mean random variable**  $\overline{X}$  as  $T = \frac{\overline{X} \mu}{\frac{s}{\sqrt{n}}}$

## 2.3. Distribution of t-scores (under certain conditions)



## 2.4. t-Distribution

#### Random Variable that Follows the t-Distribution:

**Definition**: A continuous random variable is said to follow the **t-distribution** with  $\nu$  degrees of

freedom if it has the following probability density function (pdf).

Short-Hand: \_\_\_\_\_

**Probability Density Function**:

$$f(x) = \frac{\Gamma(\frac{\nu+1}{2})}{\sqrt{\nu\pi}\Gamma(\frac{\nu}{2})} \left(1 + \frac{x^2}{\nu}\right)^{-\frac{\nu+1}{2}}, for - \infty < x < \infty$$

Parameter that Dictates Shape:\_\_\_\_\_

#### Properties:

- Always centered at:\_\_\_\_\_
- As the degrees of freedom increases, the \_\_\_\_\_ of the distribution

\_\_\_\_\_ and the peak of the distribution \_\_\_\_\_\_.

#### Comparison to Standard Normal Distribution:



#### Go to Unit 11 slides (Section 2.4) for these exercises.

**Ex:** Calculate the probability that a t-score (that is an observation from the t-distribution with 20 degrees of freedom) is greater than 1.96.

**Ex:** Calculate the t-score that creates a right tail area of 0.025 under the t-distribution with 20 degrees of freedom.

**Ex:** Suppose we know that the average GPA of ALL UIUC students is 3.3. We then randomly select 20 UIUC students and find that they have a sample mean GPA of 3.5 and a standard deviation of 0.3. Suppose that the distribution of all UIUC student GPAs is approximately normal. Calculate the probability (the most accurate one) of randomly selecting a sample mean that is greater than or equal to the sample mean that we collected.

## 3. Properties of the Sampling Distribution of Sample Mean Differences



**Collect Many Random Samples** (all of size n1=10) drawn with replacement.

Random Sample of n1=10 Course Enrollments (drawn with replacement from population 1)	Random Sample of n1=10 Course Enrollments (drawn with replacement from population 1)		Random Sample of n1=10 Course Enrollments (drawn with replacement from population 1)
197	197		178
178	215		197
178	215		178
197	197		178
215	215		178
178	197	1221	178
215	197		197
197	197		178
178	197		215
178	215		197

#### Sampling Distribution of Sample Mean Differences

Sample Means (from population 1)	Sample Means (from population 2)	Sampling Distribution of Sample Mean Differences
191.1	160.5	30.6
204.2	124	80.2
187.4	142	45.4



P	opul ume	atior	Data	200 1.75 1.10 1.00 1.00 1.00 1.00 1.00 1.00 1.0	n 2 Sundution
	course	section	enrolled	199	
0	cs105	в	345	8 00 300 300 100 Clean	200 250 350 250 Creditment Star
1	cs105	А	201		
2	stat107	А	197	Mean of I	Population
3	stat207	Δ.	53	Distrib	ution= =
7	adu307		37	EIV 1	
	duv507	~	57		
			_	Standard	Deviation of Population
				Distrib	ution= =
San	nple	2 of		$SD[X_2]$	
Bloom	- C		-	Samp	e 2 Distribution
Nur	neri	carD	ata	3.0	
				25	
1	course se	ection enro	lled	Leduc	
0	cs105	в	345	10 CT	
0	cs105	в	345	05	
3	stat207	A	53	03	
0	cs105	B	345	50 100 150 Class	200 250 300 350 Enrollment Size
2	stat107	A	197		
7	adv307	A	37		
2	stat107	A	197	Mean of Sample 2	
2	stat107	A	197	Distribution=	
				Standard Devia	ation of Sample
				2 Distributio	on=
				Size of Sample	2
				Distribution	-
				Distribution	
Paul and a state	Band	. <b>6</b>		Dearline Complex of	
n2=8 Course	n2=5	n Sample of 8 Course		n2=8 Course	Collect Many
Enrollments (drawn	Enrolime	ents (drawr		Enrollments (drawn	Random Samples
with replacement	with re	placement		with replacement	(all of size p2-9)
from population 2)	from po	pulation 2)		from population 2)	(all 01 Size 112=8)
201		3	7	53	drawn with

<u>Mean</u> of <b>Sampling</b> Distribution $\approx 1$	$E[\bar{X}_1 - \bar{X}_2] = \_$
<u>Standard deviation</u> of <b>Sampling</b> Di	stribution $\approx SD[\overline{X}_1 - \overline{X}_2] =$
<u>Shape</u> of <b>Sampling</b> Distribution is _ conditions below are met:	when the 5
1	and
2	and
3	and

and

is independent of

Need to know the following to make an inference about Unknown Population Mean Difference  $\mu_{4} - \mu_{2}$ :

replacement.

4. 5. 201

## **3.1. Mean of the Sampling Distribution of Sample Mean Differences**

Knowns:

- $E[\overline{X}_1] = \mu_1$
- $E[\overline{X}_2] = \mu_2$
- <u>Property</u>: E[aY + bZ] = aE[Y] + bE[Z]

• 
$$E[\overline{X}_1 - \overline{X}_2] =$$
  
=  $\mu_1 - \mu_2$ 

## **3.2. Standard Deviation of the Sampling Distribution of Sample Mean Differences**

Knowns:

• 
$$V[\overline{X}_1] = \frac{\sigma_1^2}{n_1}$$

• 
$$V[\overline{X}_2] = \frac{\sigma_2}{n_2}$$

• Property:  $V[aY + bZ] = a^2 V[Y] + b^2 V[Z]$ 

Proof:

• 
$$V[\overline{X}_1 - \overline{X}_2] =$$
  
=  $\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}$   
•  $SD[\overline{X}_1 - \overline{X}_2] = \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}$ 

## **3.3.** When is the Sampling Distribution of Sample Mean Differences normal?

#### Knowns:

- Central Limit Theorem (for sample means):
  - $\bar{X}_1 \sim N(mean = \mu_1, std = \frac{\sigma_1}{\sqrt{n_1}})$  when:
    - Sample 1 is randomly collected
    - *n*<sub>1</sub> < 10% of population 1
    - $n_1 > 30$  or population 1 distribution is normal
- <u>Central Limit Theorem (for sample means):</u>

• 
$$\bar{X}_2 \sim N(mean = \mu_2, std = \frac{\sigma_2}{\sqrt{n_2}})$$
 when:

- Sample 2 is randomly collected
- *n*<sub>2</sub> < 10% of population 2
- $n_2 > 30$  or population 2 distribution is normal
- <u>Property</u>: If two random variables are \_\_\_\_\_\_ and normal, then the sum (difference) of these random variables is \_\_\_\_\_\_.

#### Central Limit Theorem (for sample mean DIFFERENCES):

If the following conditions hold, then the sampling distribution of sample mean differences will be approximately normal.

$$\bar{X}_1 - \bar{X}_2 \sim N(mean = \mu_1 - \mu_2, std = \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}})$$

- a.  $n_1 > 30$  OR population 1 distribution (equivalently sample 1 distribution) is normal.
- b.  $n_2 > 30$  OR population 2 distribution (equivalently sample 2 distribution) is normal.
- c. Sample 1 is collected randomly and  $n_1 < 10\%$  of population 1 size
- d. Sample 2 is collected randomly and  $n_2 < 10\%$  of population 2 size
- e. Observations in sample 1 and sample 2 are independent (between samples).

## 4.1. Creating a (1- $\alpha$ )100% Confidence Interval for $\mu_1 - \mu_2$

#### 2. Check the Central Limit Theorem conditions for sample mean differences.

- a.  $n_1 > 30$  OR population 1 distribution (equivalently sample 1 distribution) is normal.
- b.  $n_2 > 30$  OR population 2 distribution (equivalently sample 2 distribution) is normal.
- c. Sample 1 is collected randomly and  $n_1 < 10\%$  of population 1 size
- d. Sample 2 is collected randomly and  $n_2 < 10\%$  of population 2 size
- e. Observations in sample 1 and sample 2 are independent (between samples).

If they are met, then your confidence interval interpretations will not be invalid.

#### **3.** If you know $\sigma_1$ and $\sigma_2$ the confidence interval for $\mu_1 - \mu_2$ is calculated by:

 $(point \ estimate) \pm (critical \ value)(standard \ error)$ 

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

4. If you DON'T know both  $\sigma_1$  and  $\sigma_2$  the confidence interval for  $\mu_1 - \mu_2$  is calculated by:

 $(point \ estimate) \pm (critical \ value)(standard \ error)$ 

$$(\bar{x}_1 - \bar{x}_2) \pm t^*_{\min\{n_1 - 1, n_2 - 1\}} \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}$$

Go to Unit 11 slides (Section 4.1) for exercise in confidence interval building.

## 4.2. Conducting Hypothesis Testing for the claim $\mu_1 - \mu_2 = 0$

As with the other population parameters that we discussed in Unit 10 (ie.  $\mu$  and p) we can make a conclusion about a **null** and **alternative hypothesis** 

 $H_0$ : (population parameter) = null value

 $H_A$ : (population parameter)  $\neq$  null value

by using either a:

- 1. p-value,
- 2. test statistic, or
- 3. confidence interval.

In this case,

- our <u>population parameter</u> is  $\mu_1 \mu_2$ , and
- our <u>null value</u> is 0.

Like with  $\mu$  and p, the way that we calculate the confidence interval or p-value is different, but the way that we use them to make a conclusion about our hypotheses is the same.

# **4.3. Conducting Hypothesis Testing for the claim** $\mu_1 - \mu_2 \neq 0$ with a p-value – if you know $\sigma_1$ and $\sigma_2$ .

#### 1. Set up two hypotheses.

 $H_0: \mu_1 - \mu_2 = 0$ 

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

## 2. Check the CLT Conditions (for Sample Means Differences)

- a.  $n_1 > 30$  OR population 1 distribution (equivalently sample 1 distribution) is normal.
- b.  $n_2 > 30$  OR population 2 distribution (equivalently sample 2 distribution) is normal.
- c. Sample 1 is collected randomly and  $n_1 < 10\%$  of population 1 size
- d. Sample 2 is collected randomly and  $n_2 < 10\%$  of population 2 size
- e. Observations in sample 1 and sample 2 are independent (between samples).

If these conditions hold, then the claims that you make with this analysis will be valid.

### 3. Collect an Observed Sample Statistic:

Collect a random sample from population 1 and population 2 and calculate the sample mean difference

 $\bar{x}_1 - \bar{x}_2$ 

### 4. Calculate the p-value

$$\underline{\text{One way:}} \ p - value = \begin{cases} 2P(\bar{X} \ge (\bar{x}_1 - \bar{x}_2)), & \text{ if } (\bar{x}_1 - \bar{x}_2) \ge 0 \\ 2P(\bar{X} \le (\bar{x}_1 - \bar{x}_2)), & \text{ if } (\bar{x}_1 - \bar{x}_2) \le 0 \end{cases},$$

assuming  $\overline{X} \sim N(mean = \_\_, std \ deviation = \_\_]$ 



Another way: 
$$p - value = 2P(Z \ge |\frac{(\bar{x}_1 - \bar{x}_2) - 0}{\sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}}|),$$
  
assuming  $Z \sim N(mean = \_\_\_, std deviation = \_\_]$ 

#### 5. Make a Decision

- a. If  $p value < \alpha$ , then we "reject the null hypothesis." And we say that "there IS sufficient evidence to suggest the alternative hypothesis."
- b. If  $p value \ge \alpha$ , then we "fail to reject the null hypothesis." And we say that "there IS NOT sufficient evidence to suggest the alternative hypothesis."

# **4.4. Conducting Hypothesis Testing for the claim** $\mu_1 - \mu_2 \neq 0$ with a p-value – if you DON'T know both $\sigma_1$ and $\sigma_2$ .

#### 1. Set up two hypotheses.

 $H_0: \mu_1 - \mu_2 = 0$ 

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

## 2. Check the CLT Conditions (for Sample Means Differences)

- a.  $n_1 > 30$  OR population 1 distribution (equivalently sample 1 distribution) is normal.
- b.  $n_2 > 30$  OR population 2 distribution (equivalently sample 2 distribution) is normal.
- c. Sample 1 is collected randomly and  $n_1 < 10\%$  of population 1 size
- d. Sample 2 is collected randomly and  $n_2 < 10\%$  of population 2 size
- e. Observations in sample 1 and sample 2 are independent (between samples).

#### 3. Collect an Observed Sample Statistic:

Collect a random sample from population 1 and population 2 and calculate the sample mean difference

 $\bar{x}_1 - \bar{x}_2$ 

#### 4. Calculate the p-value

$$p - value = 2P(T_{\min\{n_1-1,n_2-1\}} \ge |\frac{(\bar{x}_1 - \bar{x}_2) - 0}{\sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}}|),$$

assuming  $T \sim t - distribution with \min\{n_1 - 1, n_2 - 1\}$  degrees of freedom



#### 5. Make a Decision

- a. If  $p value < \alpha$ , then we "reject the null hypothesis." And we say that "there IS sufficient evidence to suggest the alternative hypothesis."
- b. If  $p value \ge \alpha$ , then we "fail to reject the null hypothesis." And we say that "there IS NOT sufficient evidence to suggest the alternative hypothesis."

# **4.5. Conducting Hypothesis Testing for the claim** $\mu_1 - \mu_2 \neq 0$ with a confidence interval.

#### 1. Set up two hypotheses.

 $H_0: \mu_1 - \mu_2 = 0$ 

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

## 2. Check the CLT Conditions (for Sample Means Differences)

- a.  $n_1 > 30$  OR population 1 distribution (equivalently sample 1 distribution) is normal.
- b.  $n_2 > 30$  OR population 2 distribution (equivalently sample 2 distribution) is normal.
- c. Sample 1 is collected randomly and  $n_1 < 10\%$  of population 1 size
- d. Sample 2 is collected randomly and  $n_2 < 10\%$  of population 2 size
- e. Observations in sample 1 and sample 2 are independent (between samples).

If they are met, then your confidence interval interpretations will not be invalid.

## 3. Collect an Observed Sample Proportion Difference $\bar{x}_1 - \bar{x}_2$ and Create a Confidence Interval Around it:

### 4. Make a decision

- a. If the null value (0 in this case) is \_\_\_\_\_\_, then we "reject the null hypothesis." And we say that "there IS sufficient evidence to suggest the alternative hypothesis."
- b. If the null value (0 in this case) is \_\_\_\_\_\_, then we "fail to reject the null hypothesis." And we say that "there IS NOT sufficient evidence to suggest the alternative hypothesis."

Go to Unit 11 slides (Section 4.5) for exercise in confidence interval building.

## 5. Properties of the Sampling Distribution of Sample Proportion Differences



Collect Many Random Samples (all of size n1=10) drawn with replacement.

Random Sample of n1=10 Values (drawn with replacement from population 1)	Random Sample of n1=10 Values (drawn with replacement from population 1)		Random Sample of n1=10 Values (drawn with replacement from population 1)
1	1		1
1	0		1
0	0		1
0	0		(
1	1		C
0	1		1
0	1	1.000	C
1	0	1.000	C
1	0	1.000	C
1	0		0

Random Sample of n2=8 Values (drawn with replacement from population 2)	Random Sample of n2=8 Values (drawn with replacement from population 3)		Random Sample of n2=8 Values (drawn with replacement from population 2)	
jioni population 2j	from population 2)		jion population 2)	
1	1		1	
1	0		1	
0	0	200	1	
0	0		0	
0	1		1	
1	1		1	
1	1	110	1	
1	0		0	

Collect Many Random Samples (all of size **n2=8)** drawn with replacement.

Sampling Distribution of
Sample Proportion
Differences

Sample Proportions (from population 1)	Sample Proportions (from population 2)	Sampling Distribution of Sample Proportion Differences
0.6	0.625	-0.025
0.4	0.5	-0.1
0.4	0.75	-0.35



Need to know the following to make an inference about Unknown				
Population Proportion Difference $p_1 - p_2$ :				
<u>Mean</u> of Sampling Distribution $\approx E[\hat{P}_1 - \hat{P}_2] =$				

<u>Standard deviation</u> of **Sampling** Distribution  $\approx SD[\hat{P}_1 - \hat{P}_2] = \_$ 

Shape of Sampling D conditions below a	when the 5	
1	OR	
2	OR	
3.	and	
4	and	
5.	is independent of	

## **5.1. Mean of the Sampling Distribution of Sample Proportion Differences**

Knowns:

- $E[\widehat{P}_1] = p_1$
- $E[\widehat{P}_2] = p_2$
- Property: E[aY + bZ] = aE[Y] + bE[Z]

Proof:

• 
$$E[\widehat{P}_1 - \widehat{P}_2] =$$
  
=  $p_1 - p_2$ 

## **5.2. Standard Deviation of the Sampling Distribution of Sample Proportion Differences**

Knowns:

• 
$$V[\hat{P}_1] = \frac{p_1(1-p_1)}{n_1}$$
  
•  $V[\hat{P}_2] = \frac{p_2(1-p_2)}{n_2}$ 

• <u>Property</u>:  $V[aY + bZ] = a^2V[Y] + b^2V[Z]$ 

Proof:

• 
$$V[\hat{P}_1 - \hat{P}_2] =$$
  
=  $\frac{p_1(1 - p_1)}{n_1} + \frac{p_2(1 - p_2)}{n_2}$   
•  $SD[\hat{P}_1 - \hat{P}_2] = \sqrt{\frac{p_1(1 - p_1)}{n_1} + \frac{p_2(1 - p_2)}{n_2}}$ 

## **3.3.** When is the Sampling Distribution of Sample Proportion Differences normal?

#### Knowns:

0

• Central Limit Theorem (for sample proportions):

$$\hat{P}_1 \sim N(mean = p_1, std = \sqrt{\frac{p_1(1-p_1)}{n_1}})$$
 when:

- Sample 1 is randomly collected
- *n*<sub>1</sub> < 10% of population 1
- $n_1p_1 \ge 10$  and  $n_1(1-p_1) \ge 10$
- <u>Central Limit Theorem (for sample proportions):</u>

• 
$$\hat{P}_2 \sim N(mean = p_2, std = \sqrt{\frac{p_2(1-p_2)}{n_2}})$$
 when:

- Sample 2 is randomly collected
- *n*<sub>2</sub> < 10% of population 2
- $n_2 p_2 \ge 10$  and  $n_2(1-p_2) \ge 10$
- <u>Property</u>: If two random variables are \_\_\_\_\_\_ and normal, then the sum/difference of these random variables is \_\_\_\_\_\_.

#### Central Limit Theorem (for sample proportion DIFFERENCES):

If the following conditions hold, then the sampling distribution of sample proportion differences will be approximately normal

(ie. 
$$\hat{P}_1 - \hat{P}_2 \sim N(mean = p_1 - p_2, std = \sqrt{\frac{p_1(1-p_1)}{n_1} + \frac{p_2(1-p_2)}{n_2}})$$

- a.  $n_1 p_1 \ge 10$  and  $n_1 (1 p_1) \ge 10$
- b.  $n_2 p_2 \ge 10$  and  $n_2 (1 p_2) \ge 10$
- c. Sample 1 is randomly selected and  $n_1 < 10\%$  of population 1 size
- d. Sample 2 is randomly selected and  $n_2 < 10\%$  of population 2 size
- e. Sample 1 is independent of sample 2.

## 6.1. Creating a (1- $\alpha$ )100% Confidence Interval for $p_1 - p_2$

- 2. Check the Central Limit Theorem conditions for sample mean differences.
  - a.  $n_1 p_1 \ge 10$  and  $n_1 (1 p_1) \ge 10$
  - b.  $n_2 p_2 \ge 10$  and  $n_2 (1 p_2) \ge 10$
  - c. Sample 1 is randomly selected and  $n_1 < 10\%$  of population 1 size
  - d. Sample 2 is randomly selected and  $n_2 < 10\%$  of population 2 size
  - e. Sample 1 is independent of sample 2.

If they are met, then your confidence interval interpretations will not be invalid.

**3.** The confidence interval for  $p_1 - p_2$  is calculated by:

 $(point \ estimate) \pm (critical \ value)(standard \ error)$ 

$$(\hat{p}_1 - \hat{p}_2) \pm z^* \sqrt{\frac{p_1(1-p_1)}{n_1} + \frac{p_2(1-p_2)}{n_2}}$$

**Note**: Usually we don't know  $p_1$  and  $p_2$  to plug into the conditions and standard error. So what we can do is plug in  $\hat{p}_1$  for  $p_1$  and  $\hat{p}_2$  for  $p_2$ .

## 6.2. Conducting Hypothesis Testing for the claim $p_1 - p_2 \neq 0$

As with the other population parameters that we discussed in Unit 10 (ie.  $\mu$  and p) we can make a conclusion about a **null** and **alternative hypothesis** 

 $H_0$ : (population parameter) = null value

 $H_A$ : (population parameter)  $\neq$  null value

by using either a:

- 1. p-value,
- 2. test statistic, or
- 3. confidence interval.

In this case,

- our <u>population parameter</u> is  $p_1 p_2$ , and
- our <u>null value</u> is 0.

Like with  $\mu$  and p, the way that we calculate the confidence interval or p-value is different, but the way that we use them to make a conclusion about our hypotheses is the same.

# **6.3. Conducting Hypothesis Testing for the claim** $p_1 - p_2 \neq 0$ with a p-value

#### 1. Set up two hypotheses.

 $H_0: p_1 - p_2 = 0$ 

 $H_A: p_1 - p_2 \neq 0$ 

## 2. Check the CLT Conditions (for Sample Prportion Differences)

- a.  $n_1 p_1 \ge 10$  and  $n_1 (1 p_1) \ge 10$
- b.  $n_2 p_2 \ge 10$  and  $n_2 (1 p_2) \ge 10$
- c. Sample 1 is randomly selected and  $n_1 < 10\%$  of population 1 size
- d. Sample 2 is randomly selected and  $n_2 < 10\%$  of population 2 size
- e. Sample 1 is independent of sample 2.

If they are met, then your confidence interval interpretations will not be invalid.

#### 3. Collect an Observed Sample Statistic:

Collect a random sample from population 1 and population 2 and calculate the sample proportion difference

$$\hat{p}_1 - \hat{p}_2$$

#### 4. Calculate the p-value

$$\underline{\text{One way}}: p - value = \begin{cases} 2P(\hat{P}_1 - \hat{P}_2 \ge (\hat{p}_1 - \hat{p}_2)), & \text{ if } (\hat{p}_1 - \hat{p}_2) \ge 0 \\ 2P(\hat{P}_1 - \hat{P}_2 \le (\hat{p}_1 - \hat{p}_2)), & \text{ if } (\hat{p}_1 - \hat{p}_2) \le 0 \end{cases},$$

assuming  $\hat{P}_1 - \hat{P}_2 \sim N(mean = \_\_, std \ deviation = \_\_])$ 

We call 
$$z = \frac{(\hat{p}_1 - \hat{p}_2) - 0}{\sqrt{\frac{p_1(1-p_1)}{n_1} + \frac{p_2(1-p_2)}{n_2}}}$$
  
the **test statistic** for this  
hypothesis test.

Another way: 
$$p - value = 2P(Z \ge |\frac{(\hat{p}_1 - \hat{p}_2) - 0}{\sqrt{\frac{p_1(1-p_1)}{n_1} + \frac{p_2(1-p_2)}{n_2}}}|),$$
  
assuming  $Z \sim N(mean = \_\_\_, std deviation = \_\_\_]$ 

**Note**: Usually we don't know  $p_1$  and  $p_2$  to plug into the conditions and standard error. So what we can do is plug in  $\hat{p}_1$  for  $p_1$  and  $\hat{p}_2$  for  $p_2$ .

#### 5. Make a Decision

- a. If  $p value < \alpha$ , then we "reject the null hypothesis." And we say that "there IS sufficient evidence to suggest the alternative hypothesis."
- b. If  $\mathbf{p} \mathbf{value} \ge \alpha$ , then we "fail to reject the null hypothesis." And we say that "there IS NOT sufficient evidence to suggest the alternative hypothesis."

# **6.4. Conducting Hypothesis Testing for the claim** $p_1 - p_2 \neq 0$ with a confidence interval.

#### 1. <u>Set up two hypotheses.</u>

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

## 2. Check the CLT Conditions (for Sample Proportion Differences)

- a.  $n_1 p_1 \ge 10$  and  $n_1 (1 p_1) \ge 10$
- b.  $n_2 p_2 \ge 10$  and  $n_2 (1 p_2) \ge 10$
- c. Sample 1 is randomly selected and  $n_1 < 10\%$  of population 1 size
- d. Sample 2 is randomly selected and  $n_2 < 10\%$  of population 2 size
- e. Sample 1 is independent of sample 2.

If they are met, then your confidence interval interpretations will not be invalid.

## 3. <u>Collect an Observed Sample Proportion Difference</u> $\hat{p}_1 - \hat{p}_2$ and Create a <u>Confidence Interval Around it:</u>

**Note**: Usually we don't know  $p_1$  and  $p_2$  to plug into the conditions and standard error. So what we can do is plug in  $\hat{p}_1$  for  $p_1$  and  $\hat{p}_2$  for  $p_2$ .

#### 4. Make a decision

- a. If the null value (0 in this case) is \_\_\_\_\_\_, then we "reject the null hypothesis." And we say that "there IS sufficient evidence to suggest the alternative hypothesis."
- b. If the null value (0 in this case) is \_\_\_\_\_\_, then we "fail to reject the null hypothesis." And we say that "there IS NOT sufficient evidence to suggest the alternative hypothesis."

<u>Go to Unit 11 slides (Section 6.4) for exercise in hypothesis testing and confidence interval building.</u>