

## Inference for Differences in Proportion

In previous lectures we learned about hypothesis testing and inference for one sample mean or one sample proportion. But there are a whole litany of questions that could involve **two** proportions or **two** means:

- Is there a difference in the proportion of men that die from prostate cancer for those that undergo a surgery versus those that do not?
- How much taller, on average, are adult males than adult females? ETC

		<b>Group</b>	<b>Categorical Variable</b>
	<b>Person</b>	Surgery?	Died from prostate cancer?
<b>Group 1</b>	<b>1</b>	<b>Yes</b>	<b>Yes</b>
	<b>2</b>	<b>Yes</b>	<b>No</b>
	<b><math>\vdots</math></b>	<b><math>\vdots</math></b>	<b><math>\vdots</math></b>
	<b><math>n_1</math></b>	<b>Yes</b>	<b>Yes</b>
<b>Group 2</b>	<b>1</b>	<b>No</b>	<b>Yes</b>
	<b>2</b>	<b>No</b>	<b>No</b>
	<b><math>\vdots</math></b>	<b><math>\vdots</math></b>	<b><math>\vdots</math></b>
	<b><math>n_2</math></b>	<b>No</b>	<b>No</b>

In this case the summary statistics are as follows:

$$\hat{p}_1 =$$

$$\hat{p}_2 =$$

### Structure of the Data

**Grouping Variable** ( a categorical variable)

- 
- 
- 

**Variable of Interest**

- If categorical:
- If quantitative:

## Notation

### Population 1

- $p_1$ :
- $n_1$ :
- $\hat{p}_1$ :

### Population 2

- $p_2$ :
- $n_2$ :
- $\hat{p}_2$ :

### Population 1

- $\mu_1$ :
- $n_1$ :
- $\bar{y}_1$ :
- $s_1$ :

### Population 2

- $\mu_2$ :
- $n_2$ :
- $\bar{y}_2$ :
- $s_2$ :

**Example** Determine whether the following situations involve one group or two groups. If it involves two independent groups, identify the groups.

1. An educator wants to determine the average reading comprehension scores of her students
2. An educator assigns half the class to one reading activity and the other half of the class to another reading activity. She wants to determine if the average reading comprehension scores are different between the activities.
3. We want to compare the proportion of in-state students who get financial aid to the proportion of out-of-state students who get financial aid.
4. We want to determine if the proportion of students at a university that are in-state students is higher than the national average.

We will consider the following two types of inferences for difference in proportions.

**Confidence Intervals** for the difference in population proportions:

CI for  $p_1 - p_2$

**Hypothesis Test** for the difference in population proportions:

HT for  $p_1 - p_2$

In this situation, the parameter and statistic are:

parameter:

statistic:

## Confidence Interval for Differences in Proportions

### Conditions

1. **Randomization condition:**
2. **10% condition:**
3. **Success/Failure condition:**
4. **Independent Groups:**

### Formula

If the above conditions are met, the C% confidence interval for  $p_1 - p_2$  is:

Here the  $z^*$  is chosen based on the desired C% confidence level:

Confidence Level	80%	90%	95%	98%	99%
$z^*$	1.282	1.645	1.96	2.326	2.576

**Example** *There has been debate among doctors over whether surgery can prolong life among men suffering from prostate cancer. In a 2003 study published by the New England Journal of Medicine, men diagnosed with prostate cancer were randomly selected that either underwent surgery or not. Men were then followed to see if they died from prostate cancer.*

*Find a 95% confidence interval for the difference in population proportions of men that die from prostate cancer for those that undergo surgery versus those that do not.*

- **Group 1 – no surgery**

- $n_1 = 348$
- $y_1 = 31$
- $\hat{p}_1 = \frac{31}{348} = 0.0891$

- **Group 2 – surgery**

- $n_2 = 347$
- $y_2 = 16$
- $\hat{p}_2 = \frac{16}{347} = 0.0461$

## Hypothesis Test for Differences in Proportions

### Step 1: Hypotheses

#### Null Hypothesis

- 
- 

**Note:** We could instead test whether the difference is equal to a particular value, but this is rather uncommon.

#### Alternative Hypothesis

- 
- 

### Step 2: Assumptions

Check the following conditions:

1. **Randomization condition:**
2. **10% condition:**
3. **Success/Failure condition:**
4. **Independent Groups:**

**Step 3: Test Statistic**

Because we are dealing with two sample proportions, we need to create a pooled sample proportion:

Then our z-score is calculated as follows:

**Step 4: Find p-value**

We have three different options based on our alternative hypotheses:

$$H_a : p_1 < p_2$$

$$H_a : p_1 > p_2$$

$$H_a : p_1 \neq p_2$$

**Step 5: List your decision**

<u>P-value</u>	<u>Evidence (against Ho)</u>
Greater than .10	Little to no evidence
Between .05 and .10	Weak evidence
Between .01 and .05	Moderate Evidence
Less than .01	Strong evidence

**Step 6: Conclusion**

Make a statement about the relationship between  $p_1$  and  $p_2$  given the information from the hypothesis test.

Be sure to include:

- 
- 
-

**Example** *There has been debate among doctors over whether surgery can prolong life among men suffering from prostate cancer. In a 2003 study published by the New England Journal of Medicine, men diagnosed with prostate cancer were randomly selected that either underwent surgery or not. Men were then followed to see if they died from prostate cancer.*

*Perform a hypothesis test to determine if the proportion who died from cancer that received the surgery was lower than the proportion who died from cancer that did not receive the surgery. Use  $\alpha = 0.05$ .*

- **Group 1 – no surgery**

- $n_1 = 348$

- $y_1 = 31$

- $\hat{p}_1 = \frac{31}{348} = 0.0891$

- **Group 2 – surgery**

- $n_2 = 347$

- $y_2 = 16$

- $\hat{p}_2 = \frac{16}{347} = 0.0461$



## Inference for Difference in Means

Recall the notation for comparing means that we covered in part one:

### Population 1

- $\mu_1$ :
- $n_1$ :
- $\bar{y}_1$ :
- $s_1$ :

### Population 2

- $\mu_2$ :
- $n_2$ :
- $\bar{y}_2$ :
- $s_2$ :

We will consider two type of inference for difference in means:

### Confidence Interval for the Differene in Population Means

CI for  $\mu_1 - \mu_2$

### Hypothesis Test for the Differene in Population Means

HT for  $\mu_1 - \mu_2$

In these situations, the parameter and statistic are:

parameter:

statistic:

## Confidence Interval for Difference in Means

### Conditions

1. **Randomization condition:**
2. **10% condition:**
3. **Nearly normal condition:**
4. **Independent Groups:**

### Formula

If the conditions above are met, the C% confidence interval for  $\mu_1 - \mu_2$  is:

the  $t^*$  value has degrees of freedom computed using the formula below:

**Example** At the beginning of the semester for several years, students in Stat 101 completed a survey. In this survey, the sex and height (in inches) of the students were recorded.

Calculate a 95% CI for the mean difference in heights between males and females of the population of Stat 101 students.

- **Populations**

—

—

- **Samples**

—

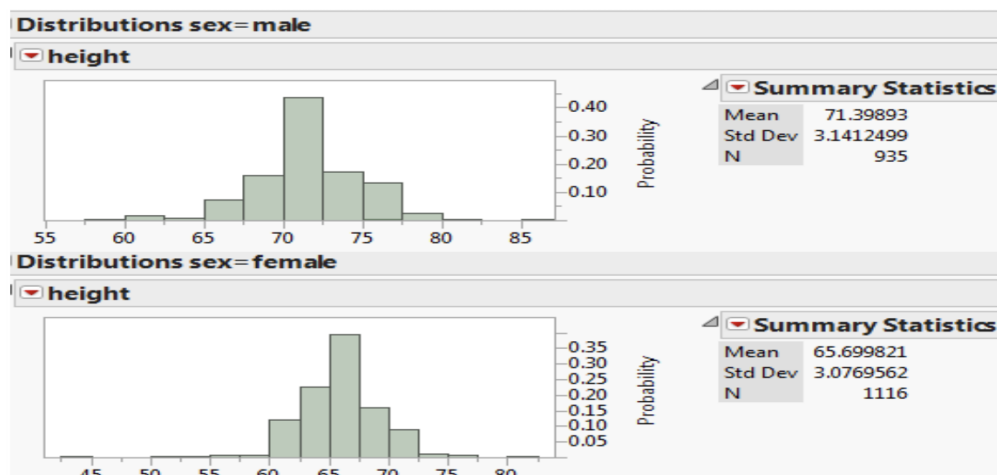
—

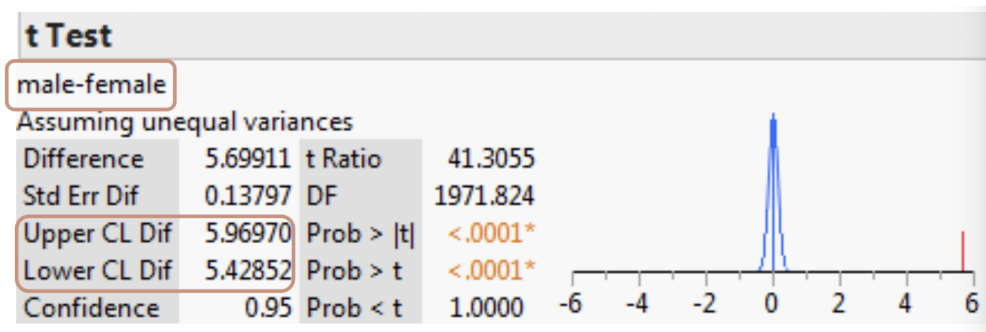
- **Parameter:**

- **Statistic:**

—

—





## Hypothesis Test for the Difference in Means

### Step 1: Hypotheses

#### Null Hypothesis

- 
- 

#### Alternative Hypothesis

- 
- 

### Step 2: Assumptions

Check the following conditions:

1. **Randomization condition:**
2. **10% condition:**
3. **Nearly Normal Condition:**
4. **Independent Groups:**

**Step 3: Test Statistic**

Then our t-statistic is calculated as follows:

**Step 4: Find p-value**

Remember, the p-value is found using a t-distribution with degrees of freedom. To compute the degrees of freedom, we use the following formula:

We have three different options based on our alternative hypotheses:

$$H_a : \mu_1 < \mu_2$$

$$H_a : \mu_1 > \mu_2$$

$$H_a : \mu_1 \neq \mu_2$$

**Step 5: List your decision**

<u>P-value</u>	<u>Evidence (against <math>H_0</math>)</u>
Greater than .10	Little to no evidence
Between .05 and .10	Weak evidence
Between .01 and .05	Moderate Evidence
Less than .01	Strong evidence

**Step 6: Conclusion**

Make a statement about the relationship between  $\mu_1$  and  $\mu_2$  given the information from the hypothesis test.

Be sure to include:

- 
- 
- 

**Example** At ISU, several different intro stats courses are offered. Each course is structured according to a particular audience of majors. At the beginning of the Fall 2006 semester, a "Survey of Attitudes Toward Statistics" was administered to students in Stat 101 and Stat 226. One of the components of this survey is called the "cognitive competence" attitude, which is rated on a scale of 1-7 where:

- 1-3 = negative attitudes
- 4 = neutral attitude
- 5-7 = positive attitudes

We want to determine if there is evidence that stat 226 students have a higher mean attitude towards "cognitive competence" than stat 101 students. There were 396 stat 226 students and 264 stat 101 students sampled. Our parameter of interest is  $\mu_1 - \mu_2$  which means that the population mean attitude score of all stat 226 students minus the population mean attitude score of all stat 101 students.

