**1**
Random Sample
of n=50 college students



Sociology Major?

**1**
Sample Percentage = 4.0%

**1**
Random Sample
of n=10 college students



Average GPA

**1**
Sample Mean = 3.0

**1**
Random Sample
of n=50 college students

**1**
Sample
Percentage = 4.0%



Percent Sociology Majors

**1**
Random Sample
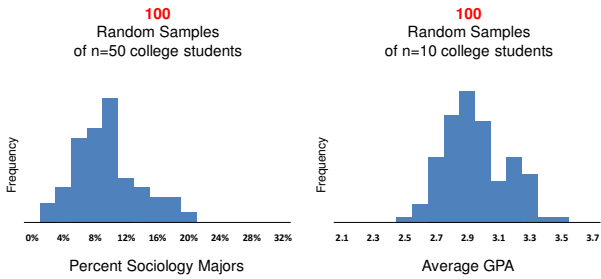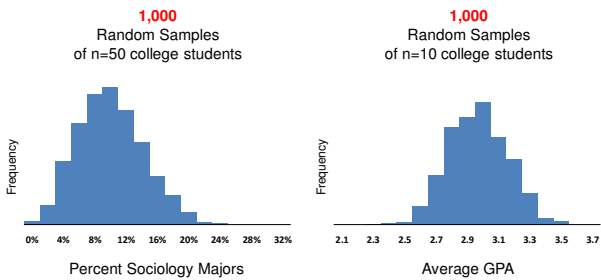of n=10 college students

**1**
Sample
Mean = 3.0



Average GPA

*Note*: These are **sampling distributions** … distributions of sample means and percentages… <u>not</u> distributions of GPA or whether people are sociology majors

**5**
Random Samples
of n=50 college students



Percent Sociology Majors

**5**
Random Samples
of n=10 college students



Average GPA

*Note*: These are **sampling distributions** … distributions of sample means and percentages… <u>not</u> distributions of GPA or whether people are sociology majors
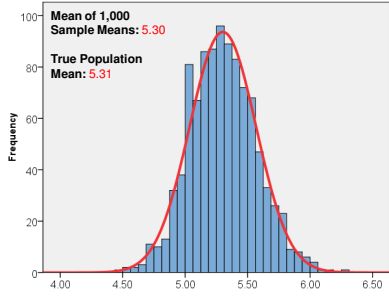
**20**
Random Samples
of n=50 college students



Percent Sociology Majors

**20**
Random Samples
of n=10 college students



Average GPA

*Note*: These are **sampling distributions** … distributions of sample means and percentages… not distributions of GPA or whether people are sociology majors
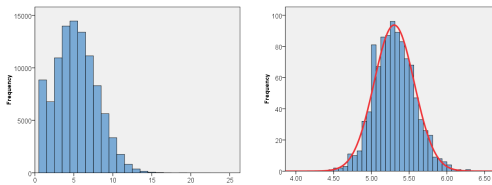
**100**
Random Samples
of n=50 college students



Percent Sociology Majors

**100**
Random Samples
of n=10 college students



Average GPA

*Note*: These are **sampling distributions** … distributions of sample means and percentages… not distributions of GPA or whether people are sociology majors

**1,000**
Random Samples
of n=50 college students



Percent Sociology Majors

**1,000**
Random Samples
of n=10 college students



Average GPA

*Note*: These are **sampling distributions** … distributions of sample means and percentages… not distributions of GPA or whether people are sociology majors

**100,000** Random Samples of n=50 college students

Frequency

0%  4%  8%  12%  16%  20%  24%  28%  32%

Percent Sociology Majors



**100,000** Random Samples of n=10 college students

Frequency

2.1  2.3  2.5  2.7  2.9  3.1  3.3  3.5  3.7

Average GPA

*Note*: These are **sampling distributions** … distributions of sample means and percentages… not distributions of GPA or whether people are sociology majors



**Infinity** Random Samples of n=50 college students

Frequency

0%  4%  8%  12%  16%  20%  24%  28%  32%

Percent Sociology Majors



**Infinity** Random Samples of n=10 college students

Frequency

2.1  2.3  2.5  2.7  2.9  3.1  3.3  3.5  3.7

Average GPA

*Note*: These are **sampling distributions** … distributions of sample means and percentages… not distributions of GPA or whether people are sociology majors

# Sampling Distributions

**Sampling Distribution**

A theoretical distribution of means or proportions, taken from an infinite number of independent random samples of size n

Sampling distributions of means and proportions are normal regardless of the shape of the distribution of the variable that produces the mean or proportion

**Central Limit Theorem**

If n is sufficiently large, then the sample means from many random samples from a population with mean $\mu$ and variance $\sigma^2$ are approximately normally distributed with mean $\mu$ and variance

$$\sigma^2/\sqrt{n}$$

## Sampling Distributions

**1**
Random Sample
of n=10 college students

**100,000**
Random Samples
of n=10 college students



Average GPA

Average GPA

Distribution of GPA

Distribution of Sample Means

**1**
Sample Mean = 3.0

Average of all **100,000**
Sample Means = 3.0

## Sampling Distributions

From the 1880 Census, the **population** distributions of…

"Born in the United States?"

"Number of People in the Household"



True Population
Mean $\mu$ = 5.31

True Population
Proportion p = 0.8668

**Number of People**
**in the Household**

**Born in the US?**
**(0=No, 1=Yes)**

## Sampling Distributions

1. I sampled 100 people, and computed…
   … the sample proportion born in the U.S. ($\hat{p}$)
   … the sample mean number of people in the household ($\overline{x}$)
2. I repeated that exercise 999 more times

On the next slides are the distributions of those 1,000 sample proportions ($\hat{p}$) and sample means ($\overline{x}$); all 1,000 samples have n=100

## Sampling Distributions



**Number of People in the Household**
1,000 Sample means, each w/ n=100

## Sampling Distributions



**Distribution of # of People**     **Sampling Distribution**

**Number of People in the Household**

## Sampling Distributions



**Proportion of People Born in the US**
1,000 Sample proportions, each w/ n=100

## Sampling Distributions

**Distribution of Birthplace**    **Sampling Distribution**

**Born in the United States?**

## Worksheet

You select **one** sample of n=10 college students

What is the probability that your sample mean GPA ($\bar{x}$) differs from the population mean by more than ±0.4?

**Infinity**
Random Samples of n=10 college students

Sampling Distribution for Mean GPA

$\mu = 3.0$    $\sigma = 0.2$

## Worksheet

You select **one** sample of n=10 college students

What is the probability that your sample mean GPA ($\bar{x}$) is greater than 3.5?

**Infinity**
Random Samples of n=10 college students

Sampling Distribution for Mean GPA

$\mu = 3.0$    $\sigma = 0.2$

## Sampling Distributions

We are 95% certain that a randomly selected sample mean ($\bar{x}$) will fall within +/- 1.96 standard deviations ($\sigma$) of the true population mean ($\mu$)

Frequency

2.1  2.3  2.5  2.7  2.9  3.1  3.3  3.5  3.7

Sampling Distribution for Means

## Sampling Distributions

So…

For **one** observed sample mean ($\bar{x}$), we are 95% certain that the true population mean ($\mu$) falls within +/- 1.96 standard deviations ($\sigma$) of $\bar{x}$

Frequency

2.1  2.3  2.5  2.7  2.9  3.1  3.3  3.5  3.7

Sampling Distribution for Means

## Sampling Distributions

So…

For **one** observed sample proportion ($\hat{p}$), we are 95% certain that the true population proportion (p) falls within +/- 1.96 standard deviations ($\sigma$) of $\hat{p}$

Frequency

0.00  0.04  0.08  0.12  0.16  0.20  0.24  0.28  0.32

Sampling Distribution for Proportions

## Sampling Distributions

Knowing what we know about…

    sampling,
    Z scores,
    random variables, and
    sampling distributions…

…we can make inferences about population **means** ($\mu$) and **proportions** (p) using sample means ($\bar{x}$) and proportions ( $\hat{p}$)

## Proportions

## Proportions

Call p the population proportion

Call p-hat ($\hat{p}$) the sample proportion

Under conditions described below, if we generate many random sample of the same size, then the distribution of the several p-hats will have a mean of p and variance of

$$\frac{\sigma^2}{\sqrt{n}} = \frac{\sum_{i=1}^{k}(Y_i - p)^2 p_i}{\sqrt{n}} = \frac{(0-p)^2 p_0 + (1-p)^2 p_1}{\sqrt{n}} = \ldots = \frac{p(1-p)}{\sqrt{n}}$$

and standard deviation :

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

## Proportions

To make inferences about a population proportion based on sample data, the following things have to be true:

1. There is a population with a fixed proportion who have a certain attribute
2. The sample is randomly selected
3. The size of the sample, n, is large … generally, such that both np and n(1-p) equal at least 5 … so how big "relatively large" is depends in part on the population proportion being estimated

## Proportions

1. There is a population with a fixed proportion who have a certain attribute
2. The sample is randomly selected
3. p equals about 0.10, so np=(50)(0.10)=5 and n(1-p)=(50)(0.9)=45



**100,000** Random Samples of n=50 college students

Percent Sociology Majors

## Proportions

If you have many $\hat{p}$, their distribution will be centered over p and will have standard deviation

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

In our example, the center is 0.10 and the standard deviation is

$$\sqrt{\frac{0.1(0.9)}{50}} = 0.042$$



**100,000** Random Samples of n=50 college students

Percent Sociology Majors

## Proportions

But what if I select only **one** random sample, with one $\hat{p}$?

What is my best guess about p?  It is $\hat{p}$

What is my best guess about the standard deviation of the sampling distribution of sample proportions, $\hat{p}$ ?

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

This is called the **standard *error*** of the sampling distribution of $\hat{p}$

## Proportions

If we select one random sample where $\hat{p}$ equals 0.12, our best guess is that p equals 0.12 and

The standard error of the sampling distribution of $\hat{p}$ for our example is:

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

$$\sqrt{\frac{0.12(1-0.12)}{50}} =$$

0.046

## Worksheet

**Work in groups!**

I selected one random sample of 200 people and got a sample proportion $\hat{p}$ of 0.4.

What is the probability that my $\hat{p}$ differs from p by more than 0.035?

## Means

## Means

The same sort of logic can be applied to the sampling distribution of sample means

In the population, the mean is $\mu$ and the standard deviation is $\sigma$

In any sample, the mean is $\bar{x}$ and the standard deviation is s

From the central limit theorem, the sampling distribution of means is centered over $\mu$ and has variance $\sigma^2/\sqrt{n}$

## Means

For our 100,000 sample means of GPA ... each based on a random sample of n=10 college students ... the mean value is 3.0 with a standard deviation of 0.2

**100,000**
Random Samples of n=10 college students

Frequency

2.1  2.3  2.5  2.7  2.9  3.1  3.3  3.5  3.7

Average GPA

## Means

But what if I select only **one** random sample, with one $\bar{x}$ ?

What is my best guess about $\mu$?  It is $\bar{x}$

What is my best guess about the standard deviation of the sampling distribution of sample means, $\bar{x}$ ?

$$\sqrt{s^2/n}$$

This is called the **standard *error*** of the sampling distribution of $\bar{x}$

## Means

If we select one random sample where $\bar{x}$ equals 2.9 and s equals 0.36, our best guess is that $\mu$ equals 2.9 and $\sigma$ equals 0.36

The standard error of the sampling distribution of $\bar{x}$ for our example is:

$$\sqrt{s^2/n} =$$
$$\sqrt{0.36/10} =$$
$$0.190$$

## Worksheet

**Work in groups!**

I selected one random sample of 1,000 people and got a sample mean $\bar{x}$ of 200 with a standard deviation, s, of 50

What is the probability that my $\bar{x}$ differs from $\mu$ by more than 3?

## t distribution

When n is large … say, more than 50 … sampling distributions of means follow the Z distribution

When n is smaller … say, less than 50 … sampling distributions of means follow **t distributions**, not Z distributions

There is a different t distribution for each value of n; each t distribution is defined by its degrees of freedom, df, where df equal n-1

## t distribution

t scores are computed the same way as Z scores

$$Z = \frac{\overline{Y} - \mu_Y}{s_Y / \sqrt{n}} \qquad t = \frac{\overline{Y} - \mu_Y}{s_Y / \sqrt{n}}$$



## t distribution

When n is large, the t and the Z distributions are (approximately) the same and so the area under the curve within any given range of Z or t scores is the same

When n is small, use the t distribution with n-1 degrees of freedom

To be safe, in practice most people always use the t distribution for sampling distributions of means

## Worksheet

**Work in groups!**

I selected one random sample of 20 people and got a sample mean x̄ of 200 with a standard deviation, s, of 50

What is the probability that my x̄ differs from μ by more than 18?

_____

_____

_____

_____

_____

_____

_____

## Why is this Useful?

**1. Confidence Intervals**

Based on the distribution of Y in sample data, we are confident that the distribution of Y in the population has particular qualities (e.g., that its mean is within a certain range of values)

*"With 95% certainty, I conclude based on my sample data that between 25% and 35% of everyone in the population has been arrested"*

_____

_____

_____

_____

_____

_____

_____

## Why is this Useful?

**2. Hypothesis Tests**

Based on the distribution of Y in the sample data, we can evaluate the likely truth of theoretically-informed hypotheses about the distribution of Y in the population (e.g., that the mean of X is above some value)

*"With 95% certainty, I reject the claim that fewer than 20% of everyone in the population has ever been arrested"*

_____

_____

_____

_____

_____

_____

# Want More?

Parts A through E of David Lane's book

http://onlinestatbook.com/2/sampling_distributions/sampling_distributions.html

Chapter 6 of Lowry's book

http://vassarstats.net/textbook/

This section of Jerry Dallal's book

http://www.jerrydallal.com/LHSP/meandist.htm

Stat Trek's discussion

http://stattrek.com/sampling/sampling-distribution.aspx

**BREAK**

# Inferences for Population Parameters

**TODAY**

When we compute **confidence intervals** we use sample data to specify a range of values within which we are confident that the population parameter falls

*Example*: "With 95% certainly we conclude that the population proportion of people who own a car is between 0.401 and 0.420"

*Example*: "With 99% certainty we conclude that the mean income in the population is somewhere between $31,200 and $32,000"

## Inferences for Population Parameters

**NEXT WEEK**

When we conduct **hypothesis testing** (or **significance testing**) we use sample data to test particular claims about the value of a population parameter

> *Example*: "Do our sample data support the assertion that more than 41% of people in the population own cars?"

> *Example*: "Do our sample data support the assertion that the mean income in the population is greater than $31,900?"

## Confidence Intervals

**Confidence Interval**
> A range of values that is "likely" to contain the population parameter (for example, a mean or proportion)

Just how "likely" it is that the confidence intervals contains the population proportion is called the *confidence level*
> If our confidence level is C%, then we are saying that if we drew many random samples and computed many confidence intervals, then the true population parameter would be contained within the resulting confidence intervals C% of the time

## Confidence Intervals

We will consider four different sorts of confidence intervals, all of which follow the same logic

**Confidence Intervals for <u>Proportions</u>**
> Use $\hat{p}$ to infer p, the population proportion

**Confidence Intervals for <u>Means</u>**
> Use $\overline{Y}$ to infer $\mu_Y$, the population mean of Y

**Confidence Intervals for <u>Differences in Proportions</u>**
> Use $\hat{p}_1 - \hat{p}_2$ to infer the difference between two population proportions, $p_1$ and $p_2$

**Confidence Intervals for <u>Differences in Means</u>**
> Use $\overline{Y}_1 - \overline{Y}_2$ to infer the difference between two population means, $\mu_{Y1}$ and $\mu_{Y2}$

## Confidence Intervals for Proportions

*EARLIER:* We said that the sample percentage was within plus or minus one "margin of error" of the population proportion 95% of the time

We defined the "conservative margin of error" as:

$$\frac{1}{\sqrt{N}} \times 100\%$$

The "conservative margin of error" is just a quick, informal way of computing a 95% confidence interval

*NOW:* Define the margin of error more formally

## Confidence Intervals for Proportions

For a 95% confidence interval for **proportions**:

$$\text{Margin of Error} = 1.96\sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$$

where

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

is the standard error of the sampling distribution of $\hat{p}$

For a 95% confidence interval, the margin of error equals plus or minus 1.96 standard errors

## Confidence Intervals for Proportions

Based on the logic of sampling distributions:

95% of the sample proportions should fall within plus or minus 1.96 standard errors of p

This is exactly the same as saying that there is a 95% chance that any particular sample proportion falls within plus or minus 1.96 standard errors of p

(Likewise, for a 68% confidence interval we would be 68% confident that the population proportion falls within plus or minus one standard error of any particular sample proportion)

## Confidence Intervals for Proportions

Informally, the formula for the confidence interval for a population proportion is $\hat{p}$ plus or minus the margin of error; we control the size of the margin of error by adjusting the desired confidence level

More formally:

$$\hat{p} \pm Z_{\alpha/2} \times \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$$

The $Z_{\alpha/2}$ is called the "multiplier"

    For a 95% confidence interval, $Z_{\alpha/2}$ equals 1.96

    For a 68% confidence interval, $Z_{\alpha/2}$ equals 1

    Etc.

## Confidence Intervals for Proportions

**What proportion of American adults (age 25+) has a physical disability?**

I drew one random sample with n=1,000; in this sample, the proportion with a disability was 0.121

What is the 95% confidence interval for the population proportion of people with disabilities?

$$\hat{p} \pm Z_{\alpha/2} \times \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$$

$$0.121 \pm 1.96 \times \sqrt{\frac{0.121(0.879)}{1,000}}, \text{ or } 0.121 \pm 0.020$$

## Confidence Intervals for Proportions

Thus with 95% certainty we conclude that…

    …the population proportion of people with a disability equals 0.121 plus or minus 0.020

    …the population proportion likely falls between 0.101 and 0.141

What would a 68% confidence interval look like?

$$\hat{p} \pm Z_{\alpha/2} \times \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$$

$$0.121 \pm 1 \times \sqrt{\frac{0.121(0.879)}{1,000}}, \text{ or } 0.121 \pm 0.010$$

Note that the margin of error got smaller

## Confidence Intervals for Proportions

For any sample of size n, the width of the confidence interval increases (or becomes less precise) as the confidence level increases



## Confidence Intervals for Proportions

We would like to have a very precise confidence interval
> It wouldn't be very useful, for example, to say that we are confident that the population proportion of people with disabilities is between 0.01 and 0.23

We can have a more precise confidence interval if we accept a lower confidence level
> But it also wouldn't very useful to say, for example, that we are 10% confident that the population proportion of people with disabilities falls between 0.120 and 0.122

How can we have a high confidence level **and** a narrow (that is to say, precise) confidence interval?

## Confidence Intervals for Proportions

The formula for the confidence interval for a population proportion includes n

$$\hat{p} \pm Z_{\alpha/2} \times \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$$

The larger the value of n, the smaller the margin of error
> We set the confidence level
>
> If n is large, we can have a high confidence level <u>and</u> a narrow confidence interval
>
> Of course, a larger n usually costs more money…

## Confidence Intervals for Proportions

Returning to our example, what would 95% confidence intervals for the population proportion of people with disabilities look like with different sample sizes?

| n | p-hat | margin of error |
|---|---|---|
| 10 | 0.121 | 0.206 |
| 100 | 0.121 | 0.065 |
| 1,000 | 0.121 | 0.020 |
| 10,000 | 0.121 | 0.007 |

## Confidence Intervals for Proportions

In fact many researchers decide how many people to include in their sample based on how precise of a confidence interval they desire

Imagine that we want to be able to construct a 95% confidence interval that has a margin of error of plus or minus 0.01

What sample size should we choose?

We would first need to make a guess about the value of *p*-hat (let's say 0.12, using the disability example)

$$0.01 = 1.96 \times \sqrt{\frac{0.12(0.88)}{n}}, \text{so } n = 4,057$$

## Worksheet

How many Americans cannot name the governor of the state in which they live?

In 1987, the General Social Survey asked 1,819 people for the name of the governor of their state

447 people gave incorrect answers

Construct and interpret a 99% confidence interval for p

## Confidence Intervals More Generally

So far we have examined confidence intervals for proportions, but the fundamental logic is the same for other types of confidence intervals

All confidence intervals can be written generally as:

$$\text{Sample Estimate} \pm \text{Multiplier} \times \text{Standard Error}$$

Or:

$$\bar{Y} \pm (Z_{\alpha/2})(\sigma_{\bar{Y}})$$

## Confidence Intervals More Generally

$$\text{Sample Estimate} \pm \text{Multiplier} \times \text{Standard Error}$$

For **proportions**: $\quad se_{\hat{p}} = \sqrt{\dfrac{\hat{p}(1-\hat{p})}{n}}$

For **means**: $\quad se_{\bar{Y}} = \dfrac{s_Y}{\sqrt{n}}$

For **differences in proportions**: $\quad se_{\hat{p}_1 - \hat{p}_2} = \sqrt{\dfrac{\hat{p}_1(1-\hat{p}_1)}{n_1} + \dfrac{\hat{p}_2(1-\hat{p}_2)}{n_2}}$

For **differences in means**: $\quad se_{Y\text{-}X} = \sqrt{\dfrac{s_Y^2}{n_Y} + \dfrac{s_X^2}{n_X}}$

## Confidence Intervals for Means

Given this general formula for any confidence interval:

$$\text{Sample Estimate} \pm \text{Multiplier} \times \text{Standard Error}$$

we can then specify a confidence interval for a population mean as:

$$\bar{Y} \pm t_{\alpha/2} \dfrac{s_Y}{\sqrt{n}}$$

Recall that to use the sample standard deviation ($s_Y$) in place of the unknown population standard deviation ($\sigma_Y$) we must use the t-distribution with n-1 degrees of freedom instead of the Z distribution

## Confidence Intervals for Means

The sample must be representative of the population from which it was drawn

One of these two things must be true:

The random sample is small (<30) and symmetrically distributed with no outliers, and the population of measurements is bell-shaped

--- or ---

The size of the random sample is large (≥30), regardless of the shape of the distribution of the measurements in the population

## Confidence Intervals for Means

1,000 adults are randomly selected from the population of the United States. Their mean personal income was $33,529, with a standard deviation of $40,609

In general, a confidence interval for μ equals

$$\$33,529 \pm t_{\alpha/2} \frac{\$40,609}{\sqrt{1000}}$$

According to a **t Table**,

$t_{\alpha/2}$ = 1.645 yields a confidence level of 0.90
$t_{\alpha/2}$ = 1.960 yields a confidence level of 0.95
$t_{\alpha/2}$ = 2.576 yields a confidence level of 0.99

## Confidence Intervals for Means

A 90% confidence interval for the mean equals:

$$\$33,529 \pm 1.645 \times \frac{\$40,609}{\sqrt{1000}}, \text{or } \$33,529 \pm \$2,119$$

A 95% confidence interval for the mean equals:

$$\$33,529 \pm 1.960 \times \frac{\$40,609}{\sqrt{1000}}, \text{or } \$33,529 \pm \$2,517$$

A 99% confidence interval for the mean equals:

$$\$33,529 \pm 2.576 \times \frac{\$40,609}{\sqrt{1000}}, \text{or } \$33,529 \pm \$3,313$$

## Worksheet

We sampled 1,600 people and observed their IQ scores. We got a sample mean of 103 and a standard deviation of 15.

Construct a 90% confidence interval for the population mean $\mu$

_____

_____

_____

_____

_____

_____

_____

_____

## Confidence Intervals for Differences in Proportions

Given this general formula for any confidence interval:

$$\text{Sample Estimate} \pm \text{Multiplier} \times \text{Standard Error}$$

we can then specify a confidence interval for the difference between two population proportions as:

$$\hat{p}_1 - \hat{p}_2 \pm Z_{\alpha/2}\sqrt{\frac{\hat{p}_1(1-\hat{p}_1)}{n_1} + \frac{\hat{p}_2(1-\hat{p}_2)}{n_2}}$$

Note that here we use the $z_{\alpha/2}$ multiplier and the standard normal distribution instead of the $t_{\alpha/2}$ multiplier and the $t$-distribution

_____

_____

_____

_____

_____

_____

_____

_____

## Confidence Intervals for Differences in Proportions

When computing a confidence interval for differences in proportions, it must be the case that the two samples are independent such that measures in one sample are not related to measures in the other sample

   If our samples were from populations of (1) men and (2) women, then measures of blood pressure might be independent

   If our samples were from populations of (1) wives and (2) their husbands, then measure of blood pressure might be dependent

Also, $n_1\hat{p}_1$, $n_1(1-\hat{p}_1)$, $n_2\hat{p}_2$, and $n_2(1-\hat{p}_2)$ must all be at least 5 (and preferably 10)

_____

_____

_____

_____

_____

_____

_____

## Confidence Intervals for Differences in Proportions

Do high school dropouts and high school graduates experience different rates of disability?

I randomly selected 1,022 people from the 2000 U.S. Census

Of the 1,022 people, $n_1$=141 had completed less than high school and $n_2$=881 had at least completed high school

Among the 141 high school non-completers, 20.2% had a disability

Among the 881 high school completers, 9.1% had a disability

In the population, how do rates of disability differ between people who have completed high school and those who have not?

## Confidence Intervals for Differences in Proportions

Given this formula for a confidence interval for differences in proportions:

$$\hat{p}_1 - \hat{p}_2 \pm Z_{\alpha/2}\sqrt{\frac{\hat{p}_1(1-\hat{p}_1)}{n_1} + \frac{\hat{p}_2(1-\hat{p}_2)}{n_2}}$$

we can construct a 95% confidence interval using a $Z_{\alpha/2}$ value of 1.96. With $n_1$=141, $n_2$=881, $\hat{p}_1$=0.202, and $\hat{p}_2$= 0.091, we have

$$0.202 - 0.091 \pm 1.96 \times \sqrt{\frac{0.202(1-0.202)}{141} + \frac{0.091(1-0.091)}{881}}$$

$$0.111 \pm 0.069 \text{ (between } 0.042 \text{ and } 0.180)$$

## Worksheet

Every year the General Social Surveys asks people whether they agree or disagree that "a working mother can establish just as warm and secure a relationship with her children as a mother who does not work."

In 1977, 735 of 1,503 respondents agreed. In 2012, 939 of 1,301 respondents agreed.

Construct a 99% confidence interval for the difference in population proportions between 1977 and 2012

## Confidence Intervals for Differences in Means

Given this general formula for any confidence interval:

$$\text{Sample Estimate} \pm \text{Multiplier} \times \text{Standard Error}$$

we can then specify a confidence interval for the difference between two population means as:

$$\overline{Y} - \overline{X} \pm t_{\alpha/2}\sqrt{\frac{s_Y^2}{n_Y} + \frac{s_X^2}{n_X}}$$

What value do we use for the df of the t-distribution?

    The mathematically correct number of degrees of freedom for the t-distribution in this case is complex

    Good Approximation: Use the smaller of $n_Y-1$ or $n_X-1$

---

## Confidence Intervals for Differences in Means

When computing a confidence interval for differences in means, it must again be the case that the two samples are <u>independent</u> such that measures in one sample are not related to measures in the other sample

One of these two things must be true for each sample:

    The random sample is small (<30) and symmetrically distributed with no outliers, and the population of measurements is bell-shaped

        --- or ---

    The size of the random sample is large ($\geq$30), regardless of the shape of the distribution of the measurements in the population

---

## Confidence Intervals for Differences in Means

Do high school dropouts and high school graduates have different mean earnings?

    I randomly selected 1,022 people from the 2000 U.S. Census

    Among the $n_Y$=141 high school non-completers, income averaged $16,259 with a standard deviation of $14,672.

    Among the $n_X$=881 high school completers, income averaged $36,284 with a standard deviation of $42,697

In the population, how do earnings differ between people who have completed high school and those who have not?

## Confidence Intervals for Differences in Means

Given this formula for a confidence interval for differences in means:

$$\overline{Y} - \overline{X} \pm t_{\alpha/2} \sqrt{\frac{s_Y^2}{n_Y} + \frac{s_X^2}{n_X}}$$

we can construct a 95% confidence interval using a $t_{\alpha/2}$ value of 1.98. With $n_Y$=141, $n_X$=881, $\overline{Y}$ =$16,259, $s_Y$=$14,672, $\overline{X}$ =$36,284, and $s_X$=$42,697, we have

$$\$16,259 - \$36,284 \pm 1.98 \times \sqrt{\frac{\$14,672^2}{141} + \frac{\$42,697^2}{881}}$$

$$-\$20,025 \pm \$3,755 \, (\text{between} - \$16,270 \, \text{and} - \$23,780)$$

## Worksheet

Every year the General Social Surveys administers a 10-item vocabulary knowledge test

In 1974, the 1,447 respondents had a mean score of 6.0 (our of 10) with a standard deviation of 2.2

In 2012, the 1,280 respondents had a mean score of 5.9 with a standard deviation of 2.0

Construct a 99% confidence interval for the difference in population means between 1974 and 2012

## Want More?

Parts 1 through 5 of David Lane's book
http://onlinestatbook.com/2/estimation/confidence_ov.html
This section of Jerry Dallal's book
http://www.jerrydallal.com/LHSP/ci.htm
Stat Trek's discussion
http://stattrek.com/estimation/confidence-interval.aspx