**+**

# Section 10.1
# Comparing Two Proportions

## Learning Objectives

After this section, you should be able to…

- ✓ DETERMINE whether the conditions for performing inference are met.

- ✓ CONSTRUCT and INTERPRET a confidence interval to compare two proportions.

- ✓ PERFORM a significance test to compare two proportions.

- ✓ INTERPRET the results of inference procedures in a randomized experiment.

# ■ <u>Difference Between Pair Differences v. 2 Samples</u>

1) Suppose we want to compare the proportions of individuals with a certain characteristic in Population 1 and Population 2.

> **Let's call these parameters of interest $p_1$ and $p_2$.**

> **The ideal strategy is to take a separate random sample from each population and to compare the sample proportions with that characteristic.**

2) What if we want to compare the effectiveness of Treatment 1 and Treatment 2 in a completely randomized experiment?

> **This time, the parameters $p_1$ and $p_2$ that we want to compare are the true proportions of successful outcomes for each treatment.**

> **We use the proportions of successes in the two treatment groups to make the comparison. Here's a table that summarizes these two situations.**

| Population or treatment | Parameter | Statistic | Sample size |
|---|---|---|---|
| 1 | $p_1$ | $\hat{p}_1$ | $n_1$ |
| 2 | $p_2$ | $\hat{p}_2$ | $n_2$ |

# ■ Describe the Shape, Center, and Spread of the Sampling Distribution of a Difference Between Two Proportions:

## The Sampling Distribution of the Difference Between Sample Proportions

Choose an SRS of size $n_1$ from Population 1 with proportion of successes $p_1$ and an independent SRS of size $n_2$ from Population 2 with proportion of successes $p_2$.

**Shape** When $n_1 p_1$, $n_1(1 - p_1)$, $n_2 p_2$ and $n_2(1 - p_2)$ are all at least 10, the sampling distribution of $\hat{p}_1 - \hat{p}_2$ is approximately Normal.

**Center** The mean of the sampling distribution is $p_1 - p_2$. That is, the difference in sample proportions is an unbiased estimator of the difference in population propotions.

**Spread** The standard deviation of the sampling distribution of $\hat{p}_1 - \hat{p}_2$ is

$$\sqrt{\frac{p_1(1 - p_1)}{n_1} + \frac{p_2(1 - p_2)}{n_2}}$$

as long as each sample is no more than 10% of its population (10% condition).

## Here are the properties for the sampling distribution of a sample proportion:

**Shape** Approximately Normal if $np \geq 10$ and $n(1 - p) \geq 10$

**Center** $\mu_{\hat{p}} = p$

**Spread** $\sigma_{\hat{p}} = \sqrt{\frac{p(1 - p)}{n}}$ if the sample is no more than 10% of the population

# ■ Reviewing of Random Variables

Both $\hat{p}_1$ and $\hat{p}_2$ are random variables. The statistic $\hat{p}_1 - \hat{p}_2$ is the difference of these two random variables.

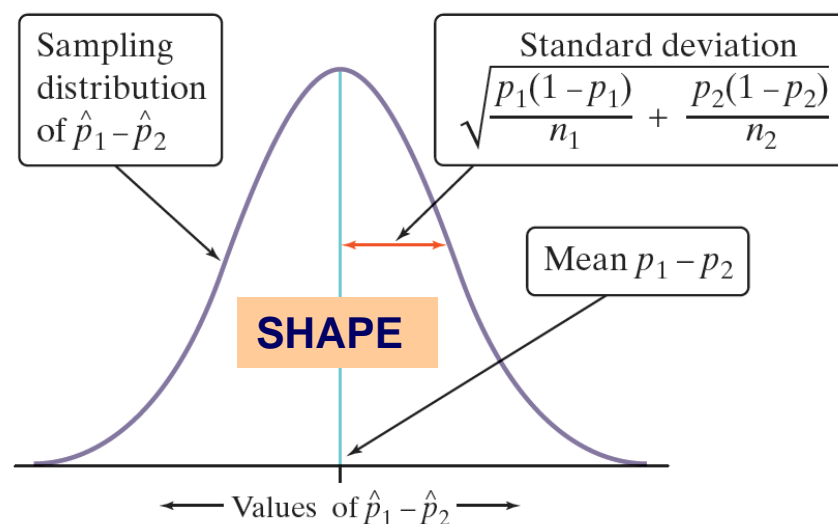We learned that for any 2 independent random variables $X$ and $Y$,

$$\mu_{X-Y} = \mu_X - \mu_Y \quad \text{and} \quad \sigma^2_{X-Y} = \sigma^2_X + \sigma^2_Y$$

➢**Therefore, Sampling Distribution of a Difference Between 2 Proportions:**

$$\mu_{\hat{p}_1 - \hat{p}_2} = \mu_{\hat{p}_1} - \mu_{\hat{p}_2} = p_1 - p_2 \quad \textbf{Center}$$

$$\sigma^2_{\hat{p}_1 - \hat{p}_2} = \sigma^2_{\hat{p}_1} + \sigma^2_{\hat{p}_2}$$

$$= \left( \sqrt{\frac{p_1(1-p_1)}{n_1}} \right)^2 + \left( \sqrt{\frac{p_2(1-p_2)}{n_2}} \right)^2$$

$$= \frac{p_1(1-p_1)}{n_1} + \frac{p_2(1-p_2)}{n_2}$$

$$\sigma_{\hat{p}_1 - \hat{p}_2} = \sqrt{\frac{p_1(1-p_1)}{n_1} + \frac{p_2(1-p_2)}{n_2}} \quad \textbf{Spread}$$

Sampling distribution of $\hat{p}_1 - \hat{p}_2$

Standard deviation
$$\sqrt{\frac{p_1(1-p_1)}{n_1} + \frac{p_2(1-p_2)}{n_2}}$$

Mean $p_1 - p_2$

**SHAPE**

← Values of $\hat{p}_1 - \hat{p}_2$ →

# Sampling Distribution of a Difference Between 2 Proportions

**HW EXAMPLE:** To explore the sampling distribution of the difference between two proportions, let's start with two populations having a known proportion of successes.

✓ At School 1, 70% of students did their homework last night

✓ At School 2, 50% of students did their homework last night.

Suppose the counselor at School 1 takes an SRS of 100 students and records the sample proportion that did their homework.

School 2's counselor takes an SRS of 200 students and records the sample proportion that did their homework.

**What can we say about the difference $\hat{p}_1 - \hat{p}_2$ in the sample proportions?**
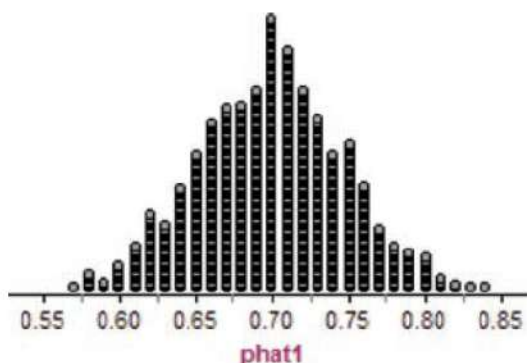
# ■Graphs of Sampling Distributions

**Reviewing GRAPHS of sampling distribution of a sample proportion.**
Each distribution was repeated 1000 times.  The results are below:

- **SRS's of 100 students from School 1.**

- **Separate SRS's of 200 students from School 2.**

- **The difference in sample proportions was then calculated and plotted.**



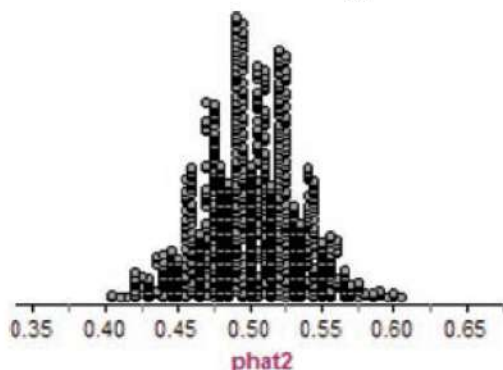(a) Approximate sampling distribution of $\hat{p}_1$

(b) Approximate sampling distribution of $\hat{p}_2$

(c) Approximate sampling distribution of $\hat{p}_1 - \hat{p}_2$

**Shape:** Approximately Normal

**Center:** $\mu_{\hat{p}_1} = p_1 = 0.70$

**Spread:** $\sigma_{\hat{p}_1} = \sqrt{\dfrac{p_1(1-p_1)}{n_1}} = \sqrt{\dfrac{0.7(0.3)}{100}} = 0.0458$
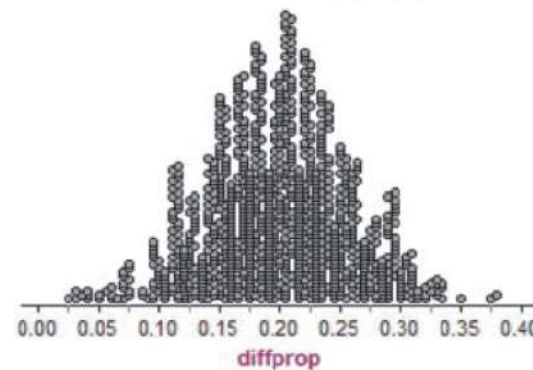
**Shape:** Approximately Normal

**Center:** $\mu_{\hat{p}_2} = p_2 = 0.50$

**Spread:** $\sigma_{\hat{p}_2} = \sqrt{\dfrac{p_2(1-p_2)}{n_2}} = \sqrt{\dfrac{0.5(0.5)}{200}} = 0.0354$

**Shape:** Approximately Normal

**Center:** $\mu_{\hat{p}_1 - \hat{p}_2} = 0.20$

**Spread:** $\sigma_{\hat{p}_1 - \hat{p}_2} = 0.058$

**What do you notice about the shape, center, and spread of the sampling distribution of $\hat{p}_1 - \hat{p}_2$ ?**

# ■ Example: Who Does More Homework?

**EXAMPLE:** Suppose that there are two large high schools, each with more than 2000 students, in a certain town. At School 1, 70% of students did their homework last night. Only 50% of the students at School 2 did their homework last night. The counselor at School 1 takes an SRS of 100 students and records the proportion that did homework. School 2's counselor takes an SRS of 200 students and records the proportion that did homework. School 1's counselor and School 2's counselor meet to discuss the results of their homework surveys. After the meeting, they both report to their principals that $\hat{p}_1 - \hat{p}_2 = 0.10.$

**a) Describe the shape, center, and spread of the sampling distribution of $\hat{p}_1 - \hat{p}_2$.**

- ## <u>Example</u>: Who Does More Homework?

a) **Describe the shape, center, and spread of the sampling distribution of** $\hat{p}_1 - \hat{p}_2$.

Because

$n_1 p_1 = 100(0.7) = 70,$

$n_1(1 - p_1) = 100(0.30) = 30,$

$n_2 p_2 = 200(0.5) = 100$

and $n_2(1 - p_2) = 200(0.5) = 100$

are all at least 10, the sampling distribution of $\hat{p}_1 - \hat{p}_2$ is approximately Normal.

Its mean is $p_1 - p_2 = 0.70 - 0.50 = 0.20.$

Its standard deviation is

$$\sqrt{\frac{0.7(0.3)}{100} + \frac{0.5(0.5)}{200}} = 0.058.$$

**b) Find the probability of getting a difference in sample proportions** $\hat{p}_1 - \hat{p}_2$ **of 0.10 or less from the two surveys.**

**c) Does the result in part (b) give us reason to doubt the counselors' reported value?**
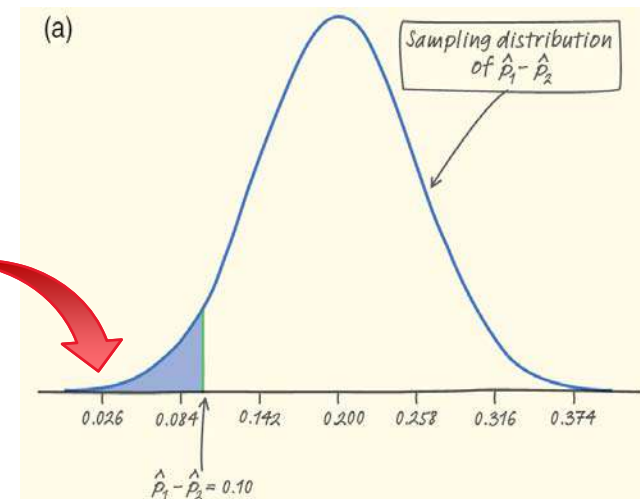
# Example: Who Does More Homework?

**b) Find the probability of getting a difference in sample proportions $\hat{p}_1 - \hat{p}_2$ of 0.10 or less from the two surveys.**

Standardize : When $\hat{p}_1 - \hat{p}_2 = 0.10,$

$$z = \frac{0.10 - 0.20}{0.058} = -1.72$$



(a)

Sampling distribution of $\hat{p}_1 - \hat{p}_2$

0.026    0.084    0.142    0.200    0.258    0.316    0.374

$\hat{p}_1 - \hat{p}_2 = 0.10$

Use Table A : The area to the left of $z = -1.72$ under the standard Normal curve is 0.0427.

**c) Does the result in part (b) give us reason to doubt the counselors' reported value?**

There is only about a 4% chance of getting a difference in sample proportions as small as or smaller than the value of 0.10 reported by the counselors. This does seem suspicious!

# Two-Sample *Z* Confidence Interval for $p_1 - p_2$

## Two-Sample *z* Interval for a Difference Between Proportions

When the Random, Normal, and Independent conditions are met, an approximate level C confidence interval for $(\hat{p}_1 - \hat{p}_2)$ is

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

where $z^*$ is the critical value for the standard Normal curve with area C between $-z^*$ and $z^*$.

**Random** The data are produced by a random sample of size $n_1$ from Population 1 and a random sample of size $n_2$ from Population 2 or by two groups of size $n_1$ and $n_2$ in a randomized experiment.

**Normal** The counts of "successes" and "failures" in each sample or group - - $n_1\hat{p}_1$, $n_1(1 - \hat{p}_1)$, $n_2\hat{p}_2$ and $n_2(1 - \hat{p}_2)$ - - are all at least 10.

**Independent** Both the samples or groups themselves and the individual observations in each sample or group are independent. When sampling without replacement, check that the two populations are at least 10 times as large as the corresponding samples (the 10% condition).

# ■ Standard Deviation vs. Standard Error for $p_1 - p_2$

1) When the Independent condition is met, the standard deviation of the statistic $\hat{p}1 - \hat{p}2$ is :

$$\sigma_{\hat{p}_1 - \hat{p}_2} = \sqrt{\frac{p1(1-p1)}{n1} + \frac{p2(1-p2)}{n2}}$$

statistic $\pm$ (critical value) $\cdot$ (standard deviation of statistic)

2) Because we don't know the values of the parameters $p1$ and $p2$, we replace them in the standard deviation formula with the sample proportions. The result is the

**standard error** of the statistic $\hat{p}1 - \hat{p}2$ : $\sqrt{\frac{\hat{p}1(1-\hat{p}1)}{n1} + \frac{\hat{p}2(1-\hat{p}2)}{n2}}$

- **Example:** *Example: Teens and Adults on Social Networks*
- *Construct a  Confidence Interval*

**EXAMPLE:** As part of the Pew Internet and American Life Project, researchers conducted two surveys in late 2009. The first survey asked a random sample of 800 U.S. teens about their use of social media and the Internet. A second survey posed similar questions to a random sample of 2253 U.S. adults. In these two studies, 73% of teens and 47% of adults said that they use social-networking sites. Use these results to construct and interpret a 95% confidence interval for the difference between the proportion of all U.S. teens and adults who use social-networking sites.

**1)  Define Parameters:**

**2) Check Conditions:**

- **Example:** *Example: Teens and Adults on Social Networks*
- *Construct a  Confidence Interval*

**State:** Our parameters of interest are

$p_1$ = the proportion of all U.S. teens who use social networking sites

$p_2$ = the proportion of all U.S. adults who use social-networking sites.

We want to estimate the difference $p_1 - p_2$ at a 95% confidence level.

**Conditions:** We should use a two-sample z interval for $p_1 - p_2$ if the conditions are satisfied.

✓ **Random** The data come from a random sample of 800 U.S. teens and a separate random sample of 2253 U.S. adults.

✓ **Normal** We check the counts of "successes" and "failures" and note the Normal condition is met since they are all at least 10:

$$n_1\hat{p}_1 = 800(0.73) = 584 \qquad n_1(1- \hat{p}_1) = 800(1-0.73) = 216$$
$$n_2\hat{p}_2 = 2253(0.47) = 1058.91 \Rightarrow 1059 \qquad n_2(1- \hat{p}_2) = 2253(1-0.47) = 1194.09 \Rightarrow 1194$$

✓ **Independent** We clearly have two independent samples—one of teens and one of adults. Individual responses in the two samples also have to be independent. The researchers are sampling without replacement, so we check the 10% condition: there are at least 10(800) = 8000 U.S. teens and at least 10(2253) = 22,530 U.S. adults.

**3)** <u>Calculations:</u>

**4)** <u>Conclusion:</u>

- **Example**: *Example: Teens and Adults on Social Networks*
- *Construct a Confidence Interval*

**Calculations:** Since the conditions are satisfied, we can construct a two-sample *z* interval for the difference $p_1 - p_2$.

$$(\hat{p}_1 - \hat{p}_2) \pm z^* \sqrt{\frac{\hat{p}_1(1-\hat{p}_1)}{n_1} + \frac{\hat{p}_2(1-\hat{p}_2)}{n_2}} = (0.73 - 0.47) \pm 1.96 \sqrt{\frac{0.73(0.27)}{800} + \frac{0.47(0.53)}{2253}}$$

$$= 0.26 \pm 0.037$$

$$= (0.223, \ 0.297)$$

**Conclude:** We are 95% confident that the interval from 0.223 to 0.297 captures the true difference in the proportion of all U.S. teens and adults who use social-networking sites. This interval suggests that more teens than adults in the United States engage in social networking by between 22.3 and 29.7 percentage points.

# ■Finding Sample Sizes:

**Example 3**    Suppose that researchers want to estimate the difference in proportions of people who are against the death penalty in Texas & in California. If the two sample sizes are the same, what size sample is needed to be within 2% of the true difference at 90% confidence?

# ■Finding Sample Sizes:

**Solution:** **Both samples must be the same size (n1=n2=n)**

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

**ME**

$$.02 = 1.645 \sqrt{\frac{.5(.5)}{n} + \frac{.5(.5)}{n}}$$

$$.02 = 1.645 \sqrt{\frac{.25 + .25}{n}}$$

n = 3383

# ■Two-Sample *z* Test for The Difference Between Two Proportions
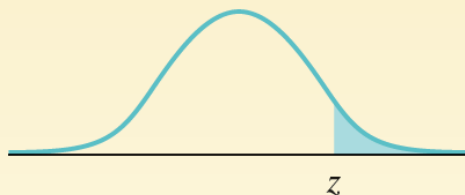
## ■Overview

### Two-Sample *z* Test for the Difference Between Proportions

Suppose the Random, Normal, and Independent conditions are met. To test the hypothesis $H_0 : p_1 - p_2 = 0$, first find the pooled proportion $\hat{p}_C$ of successes in both samples combined. Then compute the $z$ statistic
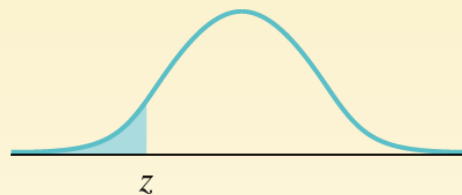
$$z = \frac{(\hat{p}_1 - \hat{p}_2) - 0}{\sqrt{\dfrac{\hat{p}_C(1 - \hat{p}_C)}{n_1} + \dfrac{\hat{p}_C(1 - \hat{p}_C)}{n_2}}}$$

Find the $P$ - value by calculating the probabilty of getting a $z$ statistic this large or larger in the direction specified by the alternative hypothesis $H_a$ :
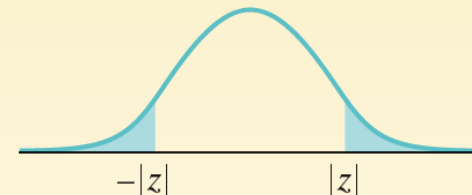
$H_a : p_1 - p_2 > 0$   $H_a : p_1 - p_2 < 0$   $H_a : p_1 - p_2 \neq 0$

# ■ Significance Tests for $p_1 - p_2$

## ■ STEP 1: Check Conditions and Define Hypothesis

- An observed difference between two sample proportions can reflect an actual difference in the parameters, or it may just be due to chance variation in random sampling or random assignment.

- Significance tests help us decide which explanation makes more sense.

- **We'll restrict ourselves to situations in which the hypothesized difference is 0.** Then the null hypothesis says that there is no difference between the two parameters:
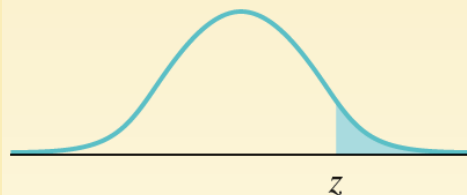
$$H_0: p_1 - p_2 = 0 \text{ or, alternatively, } H_0: p_1 = p_2$$

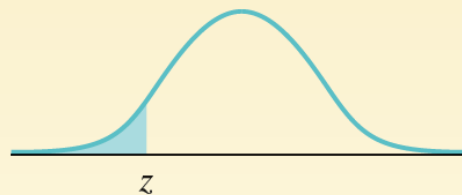. The alternative hypothesis says what kind of difference we expect

$$H_a: p_1 - p_2 > 0, H_a: p_1 - p_2 < 0, \text{ or } H_a: p_1 - p_2 \neq 0$$

Find the $P$ - value by calculating the probabilty of getting a $z$ statistic this large or larger in the direction specified by the alternative hypothesis $H_a$:
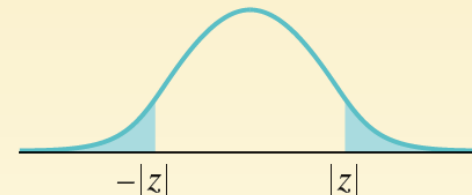
$H_a: p_1 - p_2 > 0$      $H_a: p_1 - p_2 < 0$      $H_a: p_1 - p_2 \neq 0$

$z$      $z$      $-|z|$   $|z|$

# ■ Significance Tests for $p_1 - p_2$
# ■ Define Hypothesis

An observed difference between two sample proportions can reflect an actual difference in the parameters, or it may just be due to chance variation in random sampling or random assignment. Significance tests help us decide which explanation makes more sense. The null hypothesis has the general form

$$H_0: p_1 - p_2 = \text{hypothesized value}$$

We'll restrict ourselves to situations in which the hypothesized difference is 0. Then the null hypothesis says that there is no difference between the two parameters:

$$H_0: p_1 - p_2 = 0 \text{ or, alternatively, } H_0: p_1 = p_2$$

The alternative hypothesis says what kind of difference we expect.

$$H_a: p_1 - p_2 > 0, \ H_a: p_1 - p_2 < 0, \text{ or } H_a: p_1 - p_2 \neq 0$$

If the Random, Normal, and Independent conditions are met, we can proceed with calculations.

# Significance Tests for $p_1 - p_2$
# Check Conditions

If the following conditions are met, we can proceed with a two-sample *z* test for the difference between two proportions:

**Random** The data are produced by a random sample of size $n_1$ from Population 1 and a random sample of size $n_2$ from Population 2 or by two groups of size $n_1$ and $n_2$ in a randomized experiment.

**Normal** The counts of "successes" and "failures" in each sample or group - - $n_1\hat{p}_1$, $n_1(1-\hat{p}_1)$, $n_2\hat{p}_2$ and $n_2(1-\hat{p}_2)$ - - are all at least 10.

**Independent** Both the samples or groups themselves and the individual observations in each sample or group are independent. When sampling without replacement, check that the two populations are at least 10 times as large as the corresponding samples (the 10% condition).
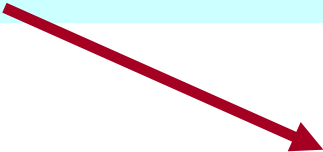
# ■ STEP 2: Calculate Test Statistics

**Two-Sample *z* Test for the Difference Between Proportions**

Suppose the Random, Normal, and Independent conditions are met. To test the hypothesis $H_0 : p_1 - p_2 = 0$, first find the pooled proportion $\hat{p}_C$ of successes in both samples combined. Then compute the $z$ statistic

**a)** **Calculate pooled (or combined) sample proportion.**

**YOU MUST MEMORIZE THIS!**

$$\hat{p}_C = \frac{\text{count of successes in both samples combined}}{\text{count of individuals in both samples combined}} = \frac{X_1 + X_2}{n_1 + n_2}$$

**b)** **Next Calculate the Z statistic**

# ■ Significance Tests for $p_1 - p_2$
# ■ Calculate Test Statistic

To do a test, standardize $\hat{p}_1 - \hat{p}_2$ to get a $z$ statistic :

$$\text{test statistic} = \frac{\text{statistic} - \text{parameter}}{\text{standard deviation of statistic}}$$

$$z = \frac{(\hat{p}_1 - \hat{p}_2) - 0}{\text{standard deviation of statistic}}$$

If $H_0$: $p_1 = p_2$ is true, the two parameters are the same. We call their common value $p$. But now we need a way to estimate $p$, so it makes sense to combine the data from the two samples. This **pooled** (or **combined**) **sample proportion** is:

$$\hat{p}_C = \frac{\text{count of successes in both samples combined}}{\text{count of individuals in both samples combined}} = \frac{X_1 + X_2}{n_1 + n_2}$$

Use $\hat{p}_C$ in place of both $p_1$ and $p_2$ in the expression for the denominator of the test statistic :

$$z = \frac{(\hat{p}_1 - \hat{p}_2) - 0}{\sqrt{\dfrac{\hat{p}_C(1 - \hat{p}_C)}{n_1} + \dfrac{\hat{p}_C(1 - \hat{p}_C)}{n_2}}}$$

# ■ Significance Tests for $p_1 - p_2$

**b) Calculate the Z statistic**

Use $\hat{p}_C$ in place of both $p_1$ and $p_2$ in the expression for the denominator of the test statistic :

$$z = \frac{(\hat{p}_1 - \hat{p}_2) - 0}{\sqrt{\dfrac{\hat{p}_C(1-\hat{p}_C)}{n_1} + \dfrac{\hat{p}_C(1-\hat{p}_C)}{n_2}}}$$

**b) Instead let's use the information from the green sheet to calculate the Z statistic**

## Two-Sample

| Statistic | Standard Deviation of Stastistic |
|---|---|
| Difference of Sample Proportions | $\sqrt{\dfrac{p_1(1-p_1)}{n_1} + \dfrac{p_2(1-p_2)}{n_2}}$ <br><br> Special case when $p_1 = p_2$ <br><br> $\sqrt{p(1-p)}\sqrt{\dfrac{1}{n_1} + \dfrac{1}{n_2}}$ |

**Use for C.I.**

**Use for H.T. and p and q are the pooled Phat**

**Test Tip: write p= $\dfrac{X1+X2}{n1+n2}$**

# ■SUMMARY of Test Statistic for $p_1 - p_2$

## Formula for Hypothesis test:

$$\text{Test statistic} = \frac{\text{statistic} - \text{parameter}}{\text{SD of statistic}}$$

$$Z = \frac{(\hat{p}_1 - \hat{p}_2) - (p_1 - p_2)}{\sqrt{\hat{p}(1-\hat{p})} \cdot \sqrt{\dfrac{1}{n_1} + \dfrac{1}{n_2}}}$$

$p_1 = p_2$

So ...

$p_1 - p_2 = 0$

$$\hat{p} = \frac{X_1 + X_2}{n_1 + n_2}$$

- **Example: *Hungry Children***
- **Significance Tests for $p_1 - p_2$**

> **EXAMPLE:** Researchers designed a survey to compare the proportions of children who come to school without eating breakfast in two low-income elementary schools. An SRS of 80 students from School 1 found that 19 had not eaten breakfast. At School 2, an SRS of 150 students included 26 who had not had breakfast. More than 1500 students attend each school. Do these data give convincing evidence of a difference in the population proportions? Carry out a significance test at the $\alpha = 0.05$ level to support your answer.

**1)  Define Parameters and Hypothesis:**

**2) Check Conditions:**

- **Example:** *Hungry Children*
- **Significance Tests for $p_1 - p_2$**

**State:** Our hypotheses are

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

Where

$p_1$ = the true proportion of students at School 1 who did not eat breakfast, and

$p_2$ = the true proportion of students at School 2 who did not eat breakfast.

**Conditions:** We should perform a two-sample z test for $p_1 - p_2$ if the conditions are satisfied.

✓ **Random** The data were produced using two simple random samples—of 80 students from School 1 and 150 students from School 2.

✓ **Normal** We check the counts of "successes" and "failures" and note the Normal condition is met since they are all at least 10:

$$n_1\hat{p}_1 = 19,\ n_1(1 - \hat{p}_1) = 61,\ n_2\hat{p}_2 = 26,\ n_2(1 - \hat{p}_2) = 124$$

✓ **Independent** We clearly have two independent samples—one from each school. Individual responses in the two samples also have to be independent. The researchers are sampling without replacement, so we check the 10% condition: there are at least 10(80) = 800 students at School 1 and at least 10(150) = 1500 students at School 2.

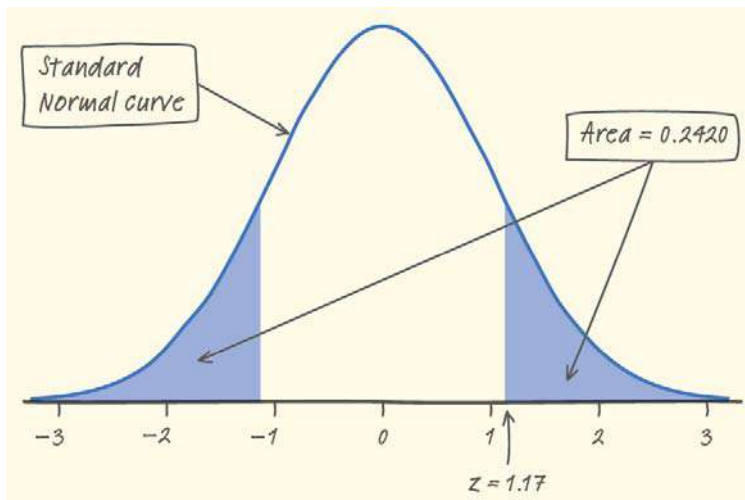**3) Calculations:**

**4) Conclusion:**

- **Example:** *Hungry Children*
- **Significance Tests for $p_1 - p_2$**

**Calculations:** Since the conditions are satisfied, we can perform a two-sample $z$ test for the difference $p_1 - p_2$.

$$\hat{p}_C = \frac{X_1 + X_2}{n_1 + n_2} = \frac{19 + 26}{80 + 150} = \frac{45}{230} = 0.1957$$

**Test statistic :**

$$z = \frac{(\hat{p}_1 - \hat{p}_2) - 0}{\sqrt{\dfrac{\hat{p}_C(1 - \hat{p}_C)}{n_1} + \dfrac{\hat{p}_C(1 - \hat{p}_C)}{n_2}}} = \frac{(0.2375 - 0.1733) - 0}{\sqrt{\dfrac{0.1957(1 - 0.1957)}{80} + \dfrac{0.1957(1 - 0.1957)}{150}}} = 1.17$$



Standard Normal curve

Area = 0.2420

z = 1.17

**P-value** Using Table A or normalcdf, the desired *P*-value is
$2P(z \geq 1.17) = 2(1 - 0.8790) = 0.2420$.

**Conclude:** Since our *P*-value, 0.2420, is greater than the chosen significance level of $\alpha = 0.05$, we fail to reject $H_0$. There is not sufficient evidence to conclude that the proportions of students at the two schools who didn't eat breakfast are different.

- **Example:** *Cholesterol and Heart Attacks*
- ***Significance Test in an Experiment***

**EXAMPLE:** High levels of cholesterol in the blood are associated with higher risk of heart attacks. Will using a drug to lower blood cholesterol reduce heart attacks? The Helsinki Heart Study recruited middle-aged men with high cholesterol but no history of other serious medical problems to investigate this question. The volunteer subjects were assigned at random to one of two treatments: 2051 men took the drug gemfibrozil to reduce their cholesterol levels, and a control group of 2030 men took a placebo. During the next five years, 56 men in the gemfibrozil group and 84 men in the placebo group had heart attacks. Is the apparent benefit of gemfibrozil statistically significant? Perform an appropriate test to find out.

1)  **Define Parameters and Hypothesis:**

2) **Check Conditions:**

- **Example:** *Cholesterol and Heart Attacks*
- *Significance Test in an Experiment*

**State:** Our hypotheses are

$$H_0: p_1 - p_2 = 0 \qquad OR \qquad H_0: p_1 = p_2$$
$$H_a: p_1 - p_2 < 0 \qquad\qquad\qquad H_a: p_1 < p_2$$

Where

- $p_1$ is the actual heart attack rate for middle-aged men like the ones in this study who take gemfibrozil,

- $p_2$ is the actual heart attack rate for middle-aged men like the ones in this study who take only a placebo.

No significance level was specified, so we'll use $\alpha = 0.01$ to reduce the risk of making a Type I error (concluding that gemfibrozil reduces heart attack risk when it actually doesn't).

**3) Calculations: Test statistic and P-value:**

**4) Conclusion:**

- **<u>Example</u>:** *Cholesterol and Heart Attacks*
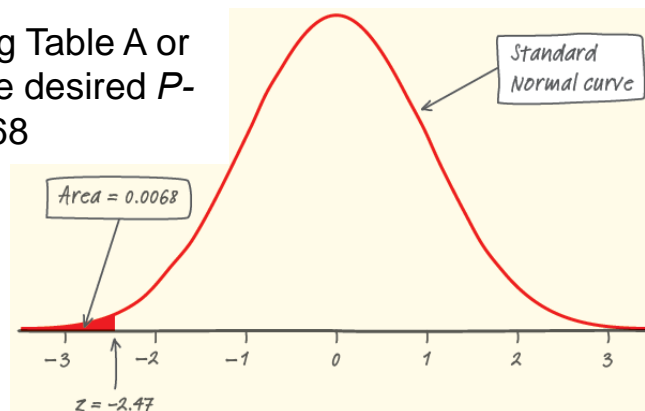- *Significance Test in an Experiment*

**TEST:** We should perform a two-sample z test for $p_1 - p_2$ if the conditions are satisfied.

✓ **Random** The data come from two groups in a randomized experiment

✓ **Normal** The number of successes (heart attacks!) and failures in the two groups are 56, 1995, 84, and 1946. These are all at least 10, so the Normal condition is met.

✓ **Independent** Due to the random assignment, these two groups of men can be viewed as independent. Individual observations in each group should also be independent: knowing whether one subject has a heart attack gives no information about whether another subject does.

**Calculations:** Since the conditions are satisfied, we can perform a two-sample *z* test for the difference $p_1 - p_2$. **Test statistic:**

$$\hat{p}_C = \frac{X_1 + X_2}{n_1 + n_2} = \frac{56 + 84}{2051 + 2030}$$
$$= \frac{140}{4081} = 0.0343$$

$$z = \frac{(\hat{p}_1 - \hat{p}_2) - 0}{\sqrt{\dfrac{\hat{p}_C(1 - \hat{p}_C)}{n_1} + \dfrac{\hat{p}_C(1 - \hat{p}_C)}{n_2}}} = \frac{(0.0273 - 0.0414) - 0}{\sqrt{\dfrac{0.0343(1 - 0.0343)}{2051} + \dfrac{0.0343(1 - 0.0343)}{2030}}} = -2.47$$

*P*-value Using Table A or normalcdf, the desired *P*-value is 0.0068



Standard Normal curve

Area = 0.0068

z = -2.47

**Conclude:** Since the *P*-value, 0.0068, is less than 0.01, the results are statistically significant at the $\alpha = 0.01$ level. We can reject $H_0$ and conclude that there is convincing evidence of a lower heart attack rate for middle-aged men like these who take gemfibrozil than for those who take only a placebo.

# + Comparing Two Proportions

## Summary

In this section, we learned that…

✓ Choose an SRS of size $n_1$ from Population 1 with proportion of successes $p_1$ and an independent SRS of size $n_2$ from Population 2 with proportion of successes $p_2$.

**Shape** When $n_1 p_1$, $n_1(1-p_1)$, $n_2 p_2$ and $n_2(1-p_2)$ are all at least 10, the sampling distribution of $\hat{p}_1 - \hat{p}_2$ is approximately Normal.

**Center** The mean of the sampling distribution is $p_1 - p_2$. That is, the difference in sample proportions is an unbiased estimator of the difference in population proportions.

**Spread** The standard deviation of the sampling distribution of $\hat{p}_1 - \hat{p}_2$ is

$$\sqrt{\frac{p_1(1-p_1)}{n_1} + \frac{p_2(1-p_2)}{n_2}}$$

as long as each sample is no more than 10% of its population (10% condition).

✓ Confidence intervals and tests to compare the proportions $p_1$ and $p_2$ of successes for two populations or treatments are based on the difference between the sample proportions.

✓ When the Random, Normal, and Independent conditions are met, we can use two-sample z procedures to estimate and test claims about $p_1 - p_2$.

# **+** Comparing Two Proportions

## Summary

In this section, we learned that…

✓ The conditions for two-sample **z** procedures are:

**Random** The data are produced by a random sample of size $n_1$ from Population 1 and a random sample of size $n_2$ from Population 2 or by two groups of size $n_1$ and $n_2$ in a randomized experiment.

**Normal** The counts of "successes" and "failures" in each sample or group - - $n_1\hat{p}_1$, $n_1(1-\hat{p}_1)$, $n_2\hat{p}_2$ and $n_2(1-\hat{p}_2)$ - - are all at least 10.

**Independent** Both the samples or groups themselves and the individual observations in each sample or group are independent. When sampling without replacement, check that the two populations are at least 10 times as large as the corresponding samples (the 10% condition).

✓ An approximate level **C** confidence interval for **$p_1$ - $p_2$** is

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

where **z\*** is the standard Normal critical value. This is called a **two-sample z interval** for **$p_1$ - $p_2$**.

# + Comparing Two Proportions

## Summary

In this section, we learned that…

- ✓ **Significance tests** of $H_0$: $p_1 - p_2 = 0$ use the **pooled (combined) sample proportion**

$$\hat{p}_C = \frac{\text{count of successes in both samples combined}}{\text{count of individuals in both samples combined}} = \frac{X_1 + X_2}{n_1 + n_2}$$

- ✓ The **two-sample z test for $p_1$ - $p_2$** uses the test statistic

$$z = \frac{(\hat{p}_1 - \hat{p}_2) - 0}{\sqrt{\dfrac{\hat{p}_C(1-\hat{p}_C)}{n_1} + \dfrac{\hat{p}_C(1-\hat{p}_C)}{n_2}}}$$

  with *P*-values calculated from the standard Normal distribution.

- ✓ Inference about the difference $p_1$ - $p_2$ in the effectiveness of two treatments in a completely randomized experiment is based on the **randomization distribution** of the difference of sample proportions. When the Random, Normal, and Independent conditions are met, our usual inference procedures based on the sampling distribution will be approximately correct.