Statistical Methods

Nels Grevstad

Metropolitan State University of Denver

ngrevsta@msudenver.edu

September 14, 2019

イロト イポト イヨト イヨト

3

Nels Grevstad

Topics



2 Two-Sample Z Confidence Interval for p_1-p_2

◆□▶ ◆□▶ ◆三▶ ◆三▶ ○○ のへの



Objectives:

- Carry out a two-sample *z* test for two population proportions.
- Compute and interpret a two-sample *z* CI for the difference between two population proportions.

イロト イポト イヨト イヨト

Two-Sample Z Test for Two Population Proportions p_1 and p_2 (9.4)

 Suppose we have random samples of sizes m and n from two populations of successes and failures.

イロト イポト イヨト イヨト

Two-Sample Z Test for Two Population Proportions p_1 and p_2 (9.4)

 Suppose we have random samples of sizes m and n from two populations of successes and failures.

 We'll see how to use the samples to decide if the population proportions of successes p₁ and p₂ are different.

Two-Sample Z Test for Two Population Proportions p_1 and p_2 (9.4)

- Suppose we have random samples of sizes m and n from two populations of successes and failures.
- We'll see how to use the samples to decide if the population proportions of successes p₁ and p₂ are different.

The appropriate test is called the *two-sample* z *test for* $p_1 - p_2$.

ヘロト ヘアト ヘビト ヘビト

• The **null hypothesis** is that no difference between the population proportions p_1 and p_2 :

Null Hypothesis:

$$H_0: p_1 - p_2 = 0$$

イロト イポト イヨト イヨト

• The **alternative hypothesis** will depend on what we're trying to "prove":

Alternative Hypothesis: The alternative hypothesis will be one of

- 1. $H_a: p_1 p_2 > 0$ (one-sided, upper-tailed)
- 2. $H_a: p_1 p_2 < 0$ (one-sided, lower-tailed)
- 3. $H_a: p_1 p_2 \neq 0$ (two-sided, two-tailed)

depending on what we're trying to verify using the data.

イロト イポト イヨト イヨト

The Sampling Distribution of $\hat{P}_1 - \hat{P}_2$

 Suppose we have random samples of sizes m and n from two populations whose proportions of successes are p₁ and p₂.

◆□▶ ◆□▶ ◆三▶ ◆三▶ ● ● ●

The Sampling Distribution of $\hat{P}_1 - \hat{P}_2$

 Suppose we have random samples of sizes m and n from two populations whose proportions of successes are p₁ and p₂.

◆□▶ ◆□▶ ◆三▶ ◆三▶ ● ● ●

The difference P₁ - P₂ between the two sample proportions is an estimator of p₁ - p₂.

Because P₁ and P₂ are (approximately) normal random variables when m and n are both large (Class Notes 4), and linear combinations of normal random variables are themselves normal (Class Notes 1), we have the following fact.

< □ > < 同 > < 三 > <

Proposition

If we have random samples of sizes m and n (drawn *independently* of each other) from **two populations** whose **proportions** of **successes** are p_1 and p_2 , then if m and n are both **large**,

$$\hat{P}_1 - \hat{P}_2 \sim \mathsf{N}\left(p_1 - p_2, \sqrt{\frac{p_1(1-p_1)}{m} + \frac{p_2(1-p_2)}{n}}\right)$$

(approximately). In this case,

$$Z = \frac{\hat{P}_1 - \hat{P}_2 - (p_1 - p_2)}{\sqrt{\frac{p_1(1-p_1)}{m} + \frac{p_2(1-p_2)}{n}}} \sim \mathsf{N}(0,1)$$
(1)

(approximately).

• This follows because

$$\hat{P}_1 \sim N\left(p_1, \sqrt{\frac{p_1(1-p_1)}{m}}\right)$$
 and $\hat{P}_2 \sim N\left(p_2, \sqrt{\frac{p_2(1-p_2)}{n}}\right)$

イロト イポト イヨト イヨト

3

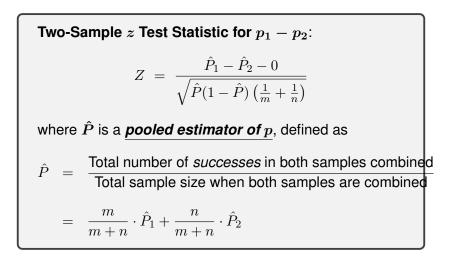
(approximately), and so $\hat{P}_1 - \hat{P}_2$ is a linear combination of two independent normal random variables.

• It follows (from the proposition) that when $H_0: p_1 - p_2 = 0$ is **true**, the random variable

$$Z = \frac{\hat{P}_1 - \hat{P}_2 - 0}{\sqrt{p(1-p)\left(\frac{1}{m} + \frac{1}{n}\right)}} \sim \mathsf{N}(0,1)$$

(approximately), where p is the **common value** of p_1 and p_2 .

◆□▶ ◆□▶ ◆三▶ ◆三▶ ○○ のへの



◆□▶ ◆□▶ ◆三▶ ◆三▶ ● ● ●

• Z measures how many standard errors $\hat{P}_1 - \hat{P}_2$ is away from 0.

- Z measures how many standard errors $\hat{P}_1 \hat{P}_2$ is away from 0.
- $\hat{P}_1 \hat{P}_2$ is an estimator of the unknown difference $p_1 p_2$, so ...

◆□▶ ◆□▶ ◆三▶ ◆三▶ ○○ のへの

- Z measures how many standard errors $\hat{P}_1 \hat{P}_2$ is away from 0.
- $\hat{P}_1 \hat{P}_2$ is an estimator of the unknown difference $p_1 p_2$, so ...
 - 1. Z will be approximately **zero** (most likely) if $p_1 p_2 = 0$.

◆□▶ ◆□▶ ◆三▶ ◆三▶ ● ● ●

- 2. It will be **positive** (most likely) if $p_1 p_2 > 0$.
- 3. It will be **negative** (most likely) if $p_1 p_2 < 0$.

- 1. Large positive values of Z provide evidence against H_0 in favor of $H_a: p_1 - p_2 > 0.$
- 2. Large negative values of Z provide evidence against H_0 in favor of

 $H_a: p_1 - p_2 < 0.$

 Large positive and large negative values of Z provide evidence against H₀ in favor of H_a: p₁ − p₂ ≠ 0.

▲□▶ ▲□▶ ▲目▶ ▲目▶ 目 のへで

Sampling Distribution of the Test Statistic Under H_0 : If Z is the two-sample z test statistic, then when m and n are both large and

$$H_0: p_1 - p_2 = 0$$

is true,

$$Z \sim N(0, 1)$$

◆□▶ ◆□▶ ◆三▶ ◆三▶ ● ● ●

(approximately).

• The N(0, 1) curve gives us:

- The N(0, 1) curve gives us:
 - The *rejection region* as the extreme 100α% of z values (in the direction(s) specified by H_a).

イロト イ理ト イヨト イヨト

= 990

- The N(0, 1) curve gives us:
 - The *rejection region* as the extreme 100α% of z values (in the direction(s) specified by H_a).
 - The *p-value* as the tail area(s) beyond the observed z value (in the direction(s) specified by H_a).

A study appearing in the journal *Science* investigated whether there's a link between television violence and aggressive behavior by those who watch a lot of TV.

イロト イポト イヨト イヨト

A study appearing in the journal *Science* investigated whether there's a link between television violence and aggressive behavior by those who watch a lot of TV.

The researchers randomly sampled **707** families in New York state and made follow-up observations over 17 years.

A study appearing in the journal *Science* investigated whether there's a link between television violence and aggressive behavior by those who watch a lot of TV.

The researchers randomly sampled **707** families in New York state and made follow-up observations over 17 years.

The table on the next slide shows results about whether a sampled teenager later conducted any aggressive act against another person.

イロト イポト イヨト イヨト

	Sample	Aggressive Act	
Time Watching TV	Size	Yes	No
Less than 1 hr per day	m = 88	5	83
More than 1 hr per day	n=619	154	465

Nels Grevstad

	Sample	Aggressive Act	
Time Watching TV	Size	Yes	No
Less than 1 hr per day	m = 88	5	83
More than 1 hr per day	n = 619	154	465

The sample proportions that conducted aggressive acts are

$$\hat{P}_1 = \frac{5}{88} = 0.057$$
 and $\hat{P}_2 = \frac{154}{619} = 0.249$

イロト 不得 とくほと くほとう

= 990

Carry out the **two-sample** *z* **test** of the hypotheses:

$$H_0: p_1 - p_2 = 0$$

$$H_a: p_1 - p_2 < 0$$

where p_1 is the **proportion** that conduct aggressive acts in the **population** that watches **less than 1 hour** of TV per day, and p_2 is the **proportion** in the **population** that watches **more than 1 hour** per day.

ヘロン 人間 とくほ とくほ とう

Carry out the **two-sample** *z* **test** of the hypotheses:

$$H_0: p_1 - p_2 = 0$$

$$H_a: p_1 - p_2 < 0$$

where p_1 is the **proportion** that conduct aggressive acts in the **population** that watches **less than 1 hour** of TV per day, and p_2 is the **proportion** in the **population** that watches **more than 1 hour** per day.

Hints: You should get the pooled estimate $\hat{P} = 0.225$, a test statistic Z = -4.04, and a p-value **0.0000**.

ヘロン 人間 とくほ とくほ とう

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

Two-Sample Z CI: For independent samples of sizes m and n from two populations whose proportions of *successes* are p_1 and p_2 , a $100(1 - \alpha)\%$ *two-sample* z *confidence interval for* $p_1 - p_2$ is

$$\hat{P}_1 - \hat{P}_2 \pm z_{\alpha/2} \cdot \sqrt{\frac{\hat{P}_1(1-\hat{P}_1)}{m} + \frac{\hat{P}_2(1-\hat{P}_2)}{n}}.$$

イロト イポト イヨト イヨト 一日

• The CI is valid as long as the sample sizes *m* and *n* are both large

イロト イポト イヨト イヨト

= 990

Consider again the study of TV violence and aggressive behavior in people who watch TV.

イロト イポト イヨト イヨト

= 990

Consider again the study of TV violence and aggressive behavior in people who watch TV.

a) Give a (point) **estimate** of the (unknown) difference $p_1 - p_2$.

イロト イポト イヨト イヨト

Consider again the study of TV violence and aggressive behavior in people who watch TV.

a) Give a (point) **estimate** of the (unknown) difference $p_1 - p_2$.

イロト イポト イヨト イヨト

э.

b) Compute and interpret a **95% CI** for $p_1 - p_2$.

Consider again the study of TV violence and aggressive behavior in people who watch TV.

- a) Give a (point) **estimate** of the (unknown) difference $p_1 p_2$.
- b) Compute and interpret a **95% CI** for $p_1 p_2$.

Hints: The *z* critical value is $z_{0.025} = 1.96$ and you should get $-0.192 \pm 0.059 = (-0.251, -0.133)$.

◆□▶ ◆□▶ ◆三▶ ◆三▶ ● ● ●