

Lecture 14 - Tests of Two Proportions

Sta102 / BME102

Colin Rundel

October 21, 2015

Difference of two proportions

Results from the GSS & Duke

The GSS asks this question, below is the distribution of responses from the 2010 survey:

| | |
|------------------|-----|
| A great deal | 454 |
| Not a great deal | 226 |
| Total | 680 |

The same question was asked of 88 Duke students, of which 56 said it would bother them a great deal.

We will collapse the data such that we consider only the responses of a great deal or not a great deal.

Difference of two proportions

Example - Melting ice cap survey

Scientists predict that global warming may have big effects on the polar regions within the next 100 years. One of the possible effects is that the northern ice cap may completely melt. Would this bother you a great deal, some, a little, or not at all if it actually happened?

- (a) A great deal
- (b) Some
- (c) A little
- (d) Not at all

Difference of two proportions

Collapsed Results

| | US | Duke | Total |
|------------------|-----|------|-------|
| A great deal | 454 | 56 | 510 |
| Not a great deal | 226 | 32 | 258 |
| Total | 680 | 88 | 768 |

This is an example of a 2×2 contingency table.

We are interested in comparing proportion of Duke students who say it would both them a great deal ($p_{GD|Duke} = 56/88$) to the proportion of all Americans who say it would both them a great deal ($p_{GD|US} = 454/680$).

Condition on what?

Knowing which of the two variables to condition on can be tricky some times.

Ask yourself - which of the two variables is the dependent variable (y) and which is the independent variable (x). In other words, changes in x should cause changes in y (not the other way around).

Once we know this then two proportions of interest are:

$$p_{y_1|x_1} \quad \text{and} \quad p_{y_1|x_2}$$

Inference for comparing proportions

- The details are the same as before...
- CI: $\text{point estimate} \pm \text{critical value} \times \text{std error}$
- HT: Use $Z = \frac{\text{point estimate} - \text{null value}}{\text{std error}}$ to find appropriate p-value.
- We just need the appropriate sampling distribution and standard error of the point estimate.

Parameter and point estimate

- *Parameter of interest*: Difference between the proportions of *all* Duke students and *all* Americans who would be bothered a great deal by the northern ice cap completely melting.

$$PGD|_{Duke} - PGD|_{US}$$

- *Point estimate*: Difference between the proportions of *sampled* Duke students and *sampled* Americans who would be bothered a great deal by the northern ice cap completely melting.

$$\hat{P}GD|_{Duke} - \hat{P}GD|_{US}$$

Sampling Distribution

Last time we saw that the sampling distribution for \hat{p} is a normal with mean p and variance $p(1 - p)/n$.

We combine this with the same approach we used for the test of two means to find the distribution of $\hat{p}_1 - \hat{p}_2$

$$\begin{aligned} E(\hat{p}_1 - \hat{p}_2) &= E(\hat{p}_1) - E(\hat{p}_2) \\ &= p_1 - p_2 \end{aligned}$$

$$\begin{aligned} \text{Var}(\hat{p}_1 - \hat{p}_2) &= \text{Var}(\hat{p}_1) + \text{Var}(\hat{p}_2) \\ &= \frac{p_1(1 - p_1)}{n_1} + \frac{p_2(1 - p_2)}{n} \end{aligned}$$

Note - this variance result requires that p_1 and p_2 are independent.

Conditions for CI for the difference of two proportions

① Independence within groups:

- The US group is sampled randomly and we're assuming that the Duke group represents a random sample as well.
- $n_{Duke} < 10\%$ of all Duke students and $680 < 10\%$ of all Americans.

We can assume that the attitudes of Duke students in the sample are independent of each other, and attitudes of US residents in the sample are independent of each other as well.

② Independence between groups: The sampled Duke students and the US residents are independent of each other.

③ Success-failure:

At least 10 observed successes and 10 observed failures in both groups.

CI for difference of proportions

Construct a 95% confidence interval for the difference between the proportions of Duke students and Americans who would be bothered a great deal by the melting of the northern ice cap ($P_{GD|Duke} - P_{GD|US}$).

| | Duke | US |
|------------------|------|-----|
| A great deal | 56 | 454 |
| Not a great deal | 32 | 226 |
| Total | 88 | 680 |

Hypotheses for testing the difference of two proportions

Just like the other hypothesis tests we have seen thus far, we formulate our null and alternative hypotheses for testing if the proportion of all Duke students who would be bothered a great deal by the melting of the northern ice cap differs from the proportion of all Americans who do as follows,

$$H_0 : P_{GD|Duke} = P_{GD|US} \Rightarrow P_{GD|Duke} - P_{GD|US} = 0$$

$$H_A : P_{GD|Duke} \neq P_{GD|US} \Rightarrow P_{GD|Duke} - P_{GD|US} \neq 0$$

Flashback to working with one proportion

- When constructing a confidence interval for a population proportion, we check if the *observed* number of successes and failures are at least 10.

$$n\hat{p} \geq 10 \quad n(1 - \hat{p}) \geq 10$$

- When conducting a hypothesis test for a population proportion, we check if the *expected* number of successes and failures are at least 10.

$$np_0 \geq 10 \quad n(1 - p_0) \geq 10$$

A slight wrinkle ...

In the case of comparing two proportions where $H_0 : p_{GD|Duke} = p_{GD|US}$, we haven't fixed either $p_{GD|Duke}$ or $p_{GD|US}$ (only their difference). As such, there isn't a specific null value we can use to calculate the *expected* number of successes and failures in each group or the standard error.

So the null hypothesis doesn't give us either $p_{GD|Duke}$ or $p_{GD|US}$ but what else do we get from it? (What does it mean that $p_{GD|Duke} = p_{GD|US}$?)

Think about them as probabilities, what does it mean that $P(GD|Duke) = P(GD|US)$?

If these two probabilities are equal then we know that global warming concern is *independent* of Duke / US status. Which means

$$P(GD|Duke) = P(GD|US) = P(GD)$$

Pooling

The upshoot of all of this is that our null hypothesis is equivalent to our two variables being independent of each other. So if we are assuming the null hypothesis is true for the test, then we must assume that the two variables are independent.

Under the assumption of independence our best guess for $p_{GD|Duke}$ and $p_{GD|US}$ is \hat{p}_{GD} which is just the sample proportion of all respondents (from Duke or US) who answers "A great deal".

We call this value \hat{p}_{pooled} ,

$$\hat{p}_{pooled} = \frac{\# \text{ of successes in 1} + \# \text{ of successes in 2}}{n_1 + n_2} = \frac{n_1 \hat{p}_1 + n_2 \hat{p}_2}{n_1 + n_2}$$

Pooled estimate of a proportion

Calculate the estimated pooled proportion of Duke students and Americans who would be bothered a great deal by the melting of the northern ice cap.

| | Duke | US | Total |
|------------------|------|-----|-------|
| A great deal | 56 | 454 | 510 |
| Not a great deal | 32 | 226 | 258 |
| Total | 88 | 680 | 788 |

$$\hat{p}_{pooled} = \frac{56 + 454}{88 + 680} = \frac{510}{788} = 0.664$$

Which sample proportion ($\hat{p}_{GD|Duke}$ or $\hat{p}_{GD|US}$) is closer to the pooled estimate? Why?

Implications for the SE

Under the null hypothesis we are stating that $p_1 = p_2$ which does not uniquely identify either p_1 or p_2 . Therefore we are using the pooled proportion (\hat{p}) as our best guess for p_1 and p_2 under the null hypothesis.

For a *confidence interval* we have

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

Therefore, for a *hypothesis test* we use \hat{p}_{pooled} to approximate for p_1 and p_2

$$SE = \sqrt{\frac{p_1(1-p_1)}{n_1} + \frac{p_2(1-p_2)}{n_2}} \approx \sqrt{\frac{\hat{p}_p(1-\hat{p}_p)}{n_1} + \frac{\hat{p}_p(1-\hat{p}_p)}{n_2}} \\ = \sqrt{\hat{p}_p(1-\hat{p}_p) \left(\frac{1}{n_1} + \frac{1}{n_2} \right)}$$

HT for comparing proportions

Do these data suggest that the proportion of all Duke students who would be bothered a great deal by the melting of the northern ice cap differs from the proportion of all Americans who do?

$$\hat{p}_{pooled} = 0.664, \quad n_1 = 88, \quad n_2 = 680$$

$$SE = \sqrt{\hat{p}_p(1 - \hat{p}_p) \left(\frac{1}{n_1} + \frac{1}{n_2} \right)} = \sqrt{0.664(1 - 0.664) \left(\frac{1}{88} + \frac{1}{680} \right)} = 0.0535$$

$$Z = \frac{(56/88 - 454/680) - 0}{0.0535} = -0.59$$

$$\begin{aligned} \text{p-value} &= P(Z < -0.59 \text{ or } Z > 0.59) \\ &= 0.277 + 0.277 = 0.555 \end{aligned}$$

Picking don't care?

What would happen to our analysis if we had picked "Not a great deal"?

| | Duke | US | Total |
|------------------|------|-----|-------|
| A great deal | 56 | 454 | 510 |
| Not a great deal | 32 | 226 | 258 |
| Total | 88 | 680 | 788 |

$$\begin{aligned} H_0 : P_{NGD|Duke} &= P_{NGD|US} & \hat{p}_{pooled} &= \frac{32 + 226}{88 + 680} = \frac{258}{788} = 0.336 \\ H_0 : P_{NGD|Duke} &\neq P_{NGD|US} \end{aligned}$$

$$SE = \sqrt{0.336(1 - 0.336) \left(\frac{1}{88} + \frac{1}{680} \right)} = 0.0535$$

$$\begin{aligned} Z &= \frac{(32/88 - 226/680) - 0}{0.0535} = 0.585 & \text{p-value} &= P(Z < -0.59 \text{ or } Z > 0.59) \\ & & &= 0.277 + 0.277 = 0.555 \end{aligned}$$

Swapping dependent and independent variables?

What would happen to our analysis if we had swapped our independent and dependent variable?

| | Duke | US | Total |
|------------------|------|-----|-------|
| A great deal | 56 | 454 | 510 |
| Not a great deal | 32 | 226 | 258 |
| Total | 88 | 680 | 788 |

$$H_0 : P_{Duke|GD} = P_{Duke|NGD}$$

$$H_0 : P_{Duke|GD} \neq P_{Duke|NGD}$$

$$\hat{p}_{pooled} = \frac{56 + 32}{510 + 258} = \frac{88}{788} = 0.115$$

$$SE = \sqrt{0.115(1 - 0.115) \left(\frac{1}{510} + \frac{1}{258} \right)} = 0.0241$$

$$\begin{aligned} Z &= \frac{(56/510 - 32/258) - 0}{0.0241} = 0.59 & \text{p-value} &= P(Z < -0.59 \text{ or } Z > 0.59) \\ & & &= 0.2775 + 0.2775 = 0.555 \end{aligned}$$

Power

What is the power of our hypothesis test to detect a difference of 0.1?

$$\text{Step 0: } H_0 : p_1 = p_2 \quad H_A : p_1 \neq p_2 \quad \alpha = 0.05 \quad \delta = 0.1 \quad SE = 0.0535 \quad \text{power} = ?$$

Step 1: Find $\hat{p}_1 - \hat{p}_2$ such that we reject H_0 .

$$P(Z < -z \text{ or } Z > z) < 0.05 \Rightarrow z > 1.96$$

$$\frac{(\hat{p}_1 - \hat{p}_2) - 0}{0.0535} < -1.96 \quad \text{or} \quad \frac{(\hat{p}_1 - \hat{p}_2) - 0}{0.0535} > 1.96$$

$$\hat{p}_1 - \hat{p}_2 < -0.105 \quad \text{or} \quad \hat{p}_1 - \hat{p}_2 > 0.105$$

Step 2: Assume $p_1 - p_2 = 0 + \delta = 0.1$ - we no longer assume independence, must use $SE = 0.0537$ from the CI instead.

$$\begin{aligned} &P(\hat{p}_1 - \hat{p}_2 < -0.105 \text{ or } \hat{p}_1 - \hat{p}_2 > 0.105) \\ &= P\left(Z < \frac{-0.105 - 0.1}{0.0527}\right) + P\left(Z > \frac{0.105 - 0.1}{0.0527}\right) \\ &= P(Z < -3.88) + P(Z > 0.094) \\ &= 0 + 0.462 = 0.462 \end{aligned}$$

Planned Parenthood

A Pew Research poll conducted between September 22-27, 2015 asked 805 randomly sampled Americans (who self identify as a Democrat or Republican) and ask about their party affiliation and whether they think any budget agreement must eliminate or maintain funding for Planned Parenthood. The distribution of their responses is shown below.

| | Eliminate | Maintain | Total |
|------------|-----------|----------|-------|
| Democrat | 45 | 378 | 423 |
| Republican | 277 | 105 | 382 |
| Total | 322 | 483 | 805 |

Pew Research Center. *Majority Says Any Budget Deal Must Include Planned Parenthood Funding*. Sep 28, 2015. <http://www.people-press.org/2015/09/28/majority-says-any-budget-deal-must-include-planned-parenthood-funding>.

Analysis - HT

Is there evidence that a greater percentage of Democrats support maintaining funding for Planned Parenthood than Republicans?

$$H_0 : p_{m|D} = p_{m|R}$$

$$H_A : p_{m|D} > p_{m|R}$$

$$\hat{p}_{pooled} = \frac{378 + 105}{423 + 382} = 0.6$$

$$SE = \sqrt{0.6(1 - 0.6) \left(\frac{1}{423} + \frac{1}{382} \right)} = 0.0346$$

$$Z = \frac{(378/423 - 105/382) - 0}{0.0346} = 17.88$$

$$p\text{-value} = P(Z > 17.88) \approx 0$$

Analysis - CI

Is there evidence that a greater percentage of Democrats support maintaining funding for Planned Parenthood than Republicans?

$$\hat{p}_{m|D} = 378/423 = 0.894 \quad \hat{p}_{m|R} = 105/382 = 0.275$$

$$SE = \sqrt{\frac{378/423(1 - 378/423)}{423} + \frac{105/382(1 - 105/382)}{382}} = 0.0273$$

$$\begin{aligned} CI &= (\hat{p}_{m|D} - \hat{p}_{m|R}) \pm Z^* SE \\ &= (378/423 - 105/382) \pm 1.96 \times 0.0273 \\ &= (0.565, 0.672) \end{aligned}$$

Power

What is the power of our hypothesis test to detect a difference of 0.62?

$$\text{Step 0: } H_0 : p_1 = p_2 \quad H_A : p_1 > p_2 \quad \alpha = 0.05 \quad \delta = 0.62 \quad SE = 0.0346 \quad \text{power} = ?$$

Step 1: Find $\hat{p}_1 - \hat{p}_2$ such that we reject H_0 .

$$P(Z > z) < 0.05 \Rightarrow z > 1.644$$

$$\frac{(\hat{p}_1 - \hat{p}_2) - 0}{0.0346} > 1.644$$

$$\hat{p}_1 - \hat{p}_2 > 0.0569$$

Step 2: Assume $p_1 - p_2 = 0 + \delta = 0.62$ - we no longer assume independence, must use $SE = 0.0273$ from the CI instead.

$$\begin{aligned} &P(\hat{p}_1 - \hat{p}_2 > 0.0569) \\ &= P\left(Z > \frac{0.0569 - 0.62}{0.0273}\right) \\ &= P(Z > -20) \\ &\approx 1 \end{aligned}$$

Recap - inference for one proportion

- Population parameter: p , point estimate: \hat{p}
- Conditions:
 - independence
 - random sample and 10% condition
 - at least 10 successes and failures
 - *observed* for CI
 - *expected* for HT
- Standard error: $SE = \sqrt{\frac{p(1-p)}{n}}$
 - for CI: use \hat{p}
 - for HT: use p_0
 - for Power:
 - Step 1 - use p_0
 - Step 2 - use $p_A = p_0 + \delta$

Recap - comparing two proportions

- Population parameter: $(p_1 - p_2)$, point estimate: $(\hat{p}_1 - \hat{p}_2)$
- Conditions:
 - independence within groups
 - random sample and 10% condition met for both groups
 - independence between groups
 - at least 10 successes and failures in each group
 - *observed* for CI
 - *expected* for HT
- $SE = \sqrt{\frac{p_1(1-p_1)}{n_1} + \frac{p_2(1-p_2)}{n_2}}$
 - for CI: use \hat{p}_1 and \hat{p}_2
 - for HT:
 - when $H_0 : p_1 = p_2$: use $\hat{p}_{pool} = \frac{\#suc_1 + \#suc_2}{n_1 + n_2}$
 - when $H_0 : p_1 - p_2 = (\text{some value other than } 0)$: use \hat{p}_1 and \hat{p}_2
 - this is pretty rare
 - for Power:
 - Step 1 - use \hat{p}_{pool}
 - Step 2 - use \hat{p}_1 and \hat{p}_2

Reference - standard error calculations

| | one sample | two samples |
|------------|--------------------------------|---|
| mean | $SE = \frac{\sigma}{\sqrt{n}}$ | $SE = \sqrt{\frac{\sigma_1^2}{n_1} + \frac{\sigma_2^2}{n_2}}$ |
| proportion | $SE = \sqrt{\frac{p(1-p)}{n}}$ | $SE = \sqrt{\frac{p_1(1-p_1)}{n_1} + \frac{p_2(1-p_2)}{n_2}}$ |

- When working with means, it's very rare that σ is known, so we usually use s as an approximation.
- When working with proportions, we will not know p therefore
 - if doing a hypothesis test, p comes from the null hypothesis
 - if constructing a confidence interval, use \hat{p} instead