

Chapter 6

Inference for categorical data¹

Department of Mathematics & Statistics
North Carolina A&T State University

¹These notes use content from OpenIntro Statistics Slides by Mine Cetinkaya-Rundel.

Difference of two proportions

Melting ice cap

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

Results from the GSS

The GSS asks the same question, below are the distributions of responses from the 2020 GSS as well as from a group of introductory statistics students at Duke University:

	GSS	Duke
A great deal	454	69
Some	124	30
A little	52	4
Not at all	50	2
Total	680	105

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.

$$p_{Duke} - p_{US}$$

- ▶ **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.

$$\hat{p}_{Duke} - \hat{p}_{US}$$

Inference for comparing proportions

- ▶ The details are the same as before...
- ▶ CI: *point estimate \pm margin of error*.
- ▶ HT: Use $Z = \frac{\text{point estimate} - \text{null value}}{SE}$ to find appropriate p-value.
- ▶ We just need the appropriate standard error of the point estimate ($SE_{\hat{p}_{Duke} - \hat{p}_{US}}$). which is the only new concept.

Standard error of the difference between two sample proportions

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

Conditions for CI for difference of proportions

- ▶ **Independence within groups:**

- ▶ The US group is sampled randomly and we're assuming that the Duke group represents a random sample as well.
- ▶ $105 < 10\%$ 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.

- ▶ **Independence between groups:** The sampled Duke students and the US residents are independent of each other.
- ▶ **Success-Failure:** At least 10 observed successes and 10 observed failures in the two groups.

Practice

Construct 95% confidence interval for the difference between the proportions of Duke students and Americans who would be bothered a great deal by melting of the northern cap. Which sample proportion (\hat{p}_{Duke} or \hat{p}_{US}) the pooled estimate is closer to? Why?

Data	Duke	US
A great deal	69	454
Not a great deal	36	226
Total	105	680

Practice

Construct 95% confidence interval for the difference between the proportions of Duke students and Americans who would be bothered a great deal by melting of the northern cap. Which sample proportion (\hat{p}_{Duke} or \hat{p}_{US}) the pooled estimate is closer to? Why?

Data	Duke	US
A great deal	69	454
Not a great deal	36	226
Total	105	680
\hat{p}	0.657	0.668

$$(\hat{p}_{Duke} - \hat{p}_{US}) \pm z^* \times \sqrt{\frac{\hat{p}_{Duke}(1 - \hat{p}_{Duke})}{n_{Duke}} + \frac{\hat{p}_{US}(1 - \hat{p}_{US})}{n_{US}}}$$

Practice

Construct 95% confidence interval for the difference between the proportions of Duke students and Americans who would be bothered a great deal by melting of the northern cap. Which sample proportion (\hat{p}_{Duke} or \hat{p}_{US}) the pooled estimate is closer to? Why?

Data	Duke	US
A great deal	69	454
Not a great deal	36	226
Total	105	680
\hat{p}	0.657	0.668

$$\begin{aligned} & (\hat{p}_{Duke} - \hat{p}_{US}) \pm z^* \times \sqrt{\frac{\hat{p}_{Duke}(1 - \hat{p}_{Duke})}{n_{Duke}} + \frac{\hat{p}_{US}(1 - \hat{p}_{US})}{n_{US}}} \\ &= (0.657 - 0.668) \end{aligned}$$

Practice

Construct 95% confidence interval for the difference between the proportions of Duke students and Americans who would be bothered a great deal by melting of the northern cap. Which sample proportion (\hat{p}_{Duke} or \hat{p}_{US}) the pooled estimate is closer to? Why?

Data	Duke	US
A great deal	69	454
Not a great deal	36	226
Total	105	680
\hat{p}	0.657	0.668

$$\begin{aligned} & (\hat{p}_{Duke} - \hat{p}_{US}) \pm z^* \times \sqrt{\frac{\hat{p}_{Duke}(1 - \hat{p}_{Duke})}{n_{Duke}} + \frac{\hat{p}_{US}(1 - \hat{p}_{US})}{n_{US}}} \\ &= (0.657 - 0.668) \pm 1.96 \end{aligned}$$

Practice

Construct 95% confidence interval for the difference between the proportions of Duke students and Americans who would be bothered a great deal by melting of the northern cap. Which sample proportion (\hat{p}_{Duke} or \hat{p}_{US}) the pooled estimate is closer to? Why?

Data	Duke	US
A great deal	69	454
Not a great deal	36	226
Total	105	680
\hat{p}	0.657	0.668

$$\begin{aligned} & (\hat{p}_{Duke} - \hat{p}_{US}) \pm z^* \times \sqrt{\frac{\hat{p}_{Duke}(1 - \hat{p}_{Duke})}{n_{Duke}} + \frac{\hat{p}_{US}(1 - \hat{p}_{US})}{n_{US}}} \\ = & (0.657 - 0.668) \pm 1.96 \times \sqrt{\frac{0.657 \times 0.343}{105} + \frac{0.668 \times 0.332}{680}} \end{aligned}$$

Practice

Construct 95% confidence interval for the difference between the proportions of Duke students and Americans who would be bothered a great deal by melting of the northern cap. Which sample proportion (\hat{p}_{Duke} or \hat{p}_{US}) the pooled estimate is closer to? Why?

Data	Duke	US
A great deal	69	454
Not a great deal	36	226
Total	105	680
\hat{p}	0.657	0.668

$$\begin{aligned}& (\hat{p}_{Duke} - \hat{p}_{US}) \pm z^* \times \sqrt{\frac{\hat{p}_{Duke}(1 - \hat{p}_{Duke})}{n_{Duke}} + \frac{\hat{p}_{US}(1 - \hat{p}_{US})}{n_{US}}} \\&= (0.657 - 0.668) \pm 1.96 \times \sqrt{\frac{0.657 \times 0.343}{105} + \frac{0.668 \times 0.332}{680}} \\&= -0.011 \pm\end{aligned}$$

Practice

Construct 95% confidence interval for the difference between the proportions of Duke students and Americans who would be bothered a great deal by melting of the northern cap. Which sample proportion (\hat{p}_{Duke} or \hat{p}_{US}) the pooled estimate is closer to? Why?

Data	Duke	US
A great deal	69	454
Not a great deal	36	226
Total	105	680
\hat{p}	0.657	0.668

$$\begin{aligned}& (\hat{p}_{Duke} - \hat{p}_{US}) \pm z^* \times \sqrt{\frac{\hat{p}_{Duke}(1 - \hat{p}_{Duke})}{n_{Duke}} + \frac{\hat{p}_{US}(1 - \hat{p}_{US})}{n_{US}}} \\&= (0.657 - 0.668) \pm 1.96 \times \sqrt{\frac{0.657 \times 0.343}{105} + \frac{0.668 \times 0.332}{680}} \\&= -0.011 \pm 1.96 \times 0.0497\end{aligned}$$

Practice

Construct 95% confidence interval for the difference between the proportions of Duke students and Americans who would be bothered a great deal by melting of the northern cap. Which sample proportion (\hat{p}_{Duke} or \hat{p}_{US}) the pooled estimate is closer to? Why?

Data	Duke	US
A great deal	69	454
Not a great deal	36	226
Total	105	680
\hat{p}	0.657	0.668

$$\begin{aligned}& (\hat{p}_{Duke} - \hat{p}_{US}) \pm z^* \times \sqrt{\frac{\hat{p}_{Duke}(1 - \hat{p}_{Duke})}{n_{Duke}} + \frac{\hat{p}_{US}(1 - \hat{p}_{US})}{n_{US}}} \\&= (0.657 - 0.668) \pm 1.96 \times \sqrt{\frac{0.657 \times 0.343}{105} + \frac{0.668 \times 0.332}{680}} \\&= -0.011 \pm 1.96 \times 0.0497 = -0.011 \pm 0.097\end{aligned}$$

Practice

Construct 95% confidence interval for the difference between the proportions of Duke students and Americans who would be bothered a great deal by melting of the northern cap. Which sample proportion (\hat{p}_{Duke} or \hat{p}_{US}) the pooled estimate is closer to? Why?

Data	Duke	US
A great deal	69	454
Not a great deal	36	226
Total	105	680
\hat{p}	0.657	0.668

$$\begin{aligned}& (\hat{p}_{Duke} - \hat{p}_{US}) \pm z^* \times \sqrt{\frac{\hat{p}_{Duke}(1 - \hat{p}_{Duke})}{n_{Duke}} + \frac{\hat{p}_{US}(1 - \hat{p}_{US})}{n_{US}}} \\&= (0.657 - 0.668) \pm 1.96 \times \sqrt{\frac{0.657 \times 0.343}{105} + \frac{0.668 \times 0.332}{680}} \\&= -0.011 \pm 1.96 \times 0.0497 = -0.011 \pm 0.097 = (-0.108, 0.086)\end{aligned}$$

Practice

Which of the following is the correct set of 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?

- A) $H_0 : p_{Duke} = p_{US}$
 $H_A : p_{Duke} \neq p_{US}$
- B) $H_0 : \hat{p}_{Duke} = \hat{p}_{US}$
 $H_A : \hat{p}_{Duke} \neq \hat{p}_{US}$
- C) $H_0 : p_{Duke} - p_{US} = 0$
 $H_A : p_{Duke} - p_{US} \neq 0$
- D) $H_0 : p_{Duke} = p_{US}$
 $H_A : p_{Duke} < p_{US}$

Practice

Which of the following is the correct set of 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?

- A) $H_0 : p_{Duke} = p_{US}$
 $H_A : p_{Duke} \neq p_{US}$
- B) $H_0 : \hat{p}_{Duke} = \hat{p}_{US}$
 $H_A : \hat{p}_{Duke} \neq \hat{p}_{US}$
- C) $H_0 : p_{Duke} - p_{US} = 0$
 $H_A : p_{Duke} - p_{US} \neq 0$
- D) $H_0 : p_{Duke} = p_{US}$
 $H_A : p_{Duke} < p_{US}$

Both A) and C) are correct.

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 failure are at least 10.

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

Pooled estimate of a proportion

- ▶ In the case of comparing two proportions where $h_0 : p_1 = p_2$, there isn't a given null value we can use to calculate the **expected** number of successes and failures in each sample.
- ▶ Therefore, we need to first find a common (**pooled**) proportion for the two groups, and use that in our analysis.
- ▶ This simply means finding the proportion of total successes among the total number of observations.

Pooled estimate of a proportion

$$\hat{p} = \frac{\# \text{ of successes}_1 + \# \text{ of successes}_2}{n_1 + n_2}$$

Practice

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

Data	Duke	US
A great deal	69	454
Not a great deal	36	226
Total	105	680
\hat{p}	0.657	0.668

Practice

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

Data	Duke	US
A great deal	69	454
Not a great deal	36	226
Total	105	680
\hat{p}	0.657	0.668

$$\hat{p} = \frac{\# \text{ of successes}_1 + \# \text{ of successes}_2}{n_1 + n_2}$$

Practice

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

Data	Duke	US
A great deal	69	454
Not a great deal	36	226
Total	105	680
\hat{p}	0.657	0.668

$$\begin{aligned}\hat{p} &= \frac{\# \text{ of successes}_1 + \# \text{ of successes}_2}{n_1 + n_2} \\ &= \frac{69 + 454}{105 + 680}\end{aligned}$$

Practice

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

Data	Duke	US
A great deal	69	454
Not a great deal	36	226
Total	105	680
\hat{p}	0.657	0.668

$$\begin{aligned}\hat{p} &= \frac{\# \text{ of successes}_1 + \# \text{ of successes}_2}{n_1 + n_2} \\ &= \frac{69 + 454}{105 + 680} = \frac{523}{785}\end{aligned}$$

Practice

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

Data	Duke	US
A great deal	69	454
Not a great deal	36	226
Total	105	680
\hat{p}	0.657	0.668

$$\begin{aligned}\hat{p} &= \frac{\# \text{ of successes}_1 + \# \text{ of successes}_2}{n_1 + n_2} \\ &= \frac{69 + 454}{105 + 680} = \frac{523}{785} = 0.666\end{aligned}$$

Practice

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? Calculate the test statistic, the p-value, and interpret your conclusion in context of the data.

Data	Duke	US
A great deal	69	454
Not a great deal	36	226
Total	105	680
\hat{p}	0.657	0.668

Practice

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? Calculate the test statistic, the p-value, and interpret your conclusion in context of the data.

Data	Duke	US
A great deal	69	454
Not a great deal	36	226
Total	105	680
\hat{p}	0.657	0.668

$$Z = \frac{(\hat{p}_{Duke} - \hat{p}_{US})}{\sqrt{\frac{\hat{p}(1-\hat{p})}{n_{Duke}} + \frac{\hat{p}(1-\hat{p})}{n_{US}}}}$$

Practice

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? Calculate the test statistic, the p-value, and interpret your conclusion in context of the data.

Data	Duke	US
A great deal	69	454
Not a great deal	36	226
Total	105	680
\hat{p}	0.657	0.668

$$\begin{aligned} Z &= \frac{(\hat{p}_{Duke} - \hat{p}_{US})}{\sqrt{\frac{\hat{p}(1-\hat{p})}{n_{Duke}} + \frac{\hat{p}(1-\hat{p})}{n_{US}}}} \\ &= \frac{(0.657 - 0.668)}{\sqrt{\frac{0.666 \times 0.334}{105} + \frac{0.666 \times 0.334}{680}}} = \end{aligned}$$

Practice

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? Calculate the test statistic, the p-value, and interpret your conclusion in context of the data.

Data	Duke	US
A great deal	69	454
Not a great deal	36	226
Total	105	680
\hat{p}	0.657	0.668

$$\begin{aligned} Z &= \frac{(\hat{p}_{Duke} - \hat{p}_{US})}{\sqrt{\frac{\hat{p}(1-\hat{p})}{n_{Duke}} + \frac{\hat{p}(1-\hat{p})}{n_{US}}}} \\ &= \frac{(0.657 - 0.668)}{\sqrt{\frac{0.666 \times 0.334}{105} + \frac{0.666 \times 0.334}{680}}} = \frac{-0.011}{0.0495} \end{aligned}$$

Practice

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? Calculate the test statistic, the p-value, and interpret your conclusion in context of the data.

Data	Duke	US
A great deal	69	454
Not a great deal	36	226
Total	105	680
\hat{p}	0.657	0.668

$$\begin{aligned} Z &= \frac{(\hat{p}_{Duke} - \hat{p}_{US})}{\sqrt{\frac{\hat{p}(1-\hat{p})}{n_{Duke}} + \frac{\hat{p}(1-\hat{p})}{n_{US}}}} \\ &= \frac{(0.657 - 0.668)}{\sqrt{\frac{0.666 \times 0.334}{105} + \frac{0.666 \times 0.334}{680}}} = \frac{-0.011}{0.0495} = -0.22 \end{aligned}$$

Practice

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? Calculate the test statistic, the p-value, and interpret your conclusion in context of the data.

Data	Duke	US
A great deal	69	454
Not a great deal	36	226
Total	105	680
\hat{p}	0.657	0.668

$$\begin{aligned} Z &= \frac{(\hat{p}_{Duke} - \hat{p}_{US})}{\sqrt{\frac{\hat{p}(1-\hat{p})}{n_{Duke}} + \frac{\hat{p}(1-\hat{p})}{n_{US}}}} \\ &= \frac{(0.657 - 0.668)}{\sqrt{\frac{0.666 \times 0.334}{105} + \frac{0.666 \times 0.334}{680}}} = \frac{-0.011}{0.0495} = -0.22 \end{aligned}$$

$$p\text{-value} = 2 \times P(Z < -0.22)$$

Practice

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? Calculate the test statistic, the p-value, and interpret your conclusion in context of the data.

Data	Duke	US
A great deal	69	454
Not a great deal	36	226
Total	105	680
\hat{p}	0.657	0.668

$$\begin{aligned} Z &= \frac{(\hat{p}_{Duke} - \hat{p}_{US})}{\sqrt{\frac{\hat{p}(1-\hat{p})}{n_{Duke}} + \frac{\hat{p}(1-\hat{p})}{n_{US}}}} \\ &= \frac{(0.657 - 0.668)}{\sqrt{\frac{0.666 \times 0.334}{105} + \frac{0.666 \times 0.334}{680}}} = \frac{-0.011}{0.0495} = -0.22 \end{aligned}$$

$$p\text{-value} = 2 \times P(Z < -0.22) = 2 \times 0.41 = 0.82$$

Recap - comparing two proportions

- ▶ Population parameter: $(p_1 - p_2)$, point estimate: $(\hat{p}_1 - \hat{p}_2)$

Recap - comparing two proportions

- ▶ Population parameter: $(p_1 - p_2)$, point estimate: $(\hat{p}_1 - \hat{p}_2)$
- ▶ Conditions:

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
 - if not \rightarrow randomization (Section 6.4)

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
 - if not \rightarrow randomization (Section 6.4)
- ▶ $SE_{(\hat{p}_1 - \hat{p}_2)} = \sqrt{\frac{p_1(1-p_1)}{n_1} + \frac{p_2(1-p_2)}{n_2}}$
 - ▶ 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

Reference - Standard error calculations

	one sample	two samples
mean	$SE = \frac{s}{\sqrt{n}}$	$SE = \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}$
proportion	$SE = \sqrt{\frac{p(1-p)}{n}}$	$SE = \sqrt{\frac{p_1(1-p_1)}{n_1} + \frac{p_2(1-p_2)}{n_2}}$

Reference - Standard error calculations

	one sample	two samples
mean	$SE = \frac{s}{\sqrt{n}}$	$SE = \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}$
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 .

Reference - Standard error calculations

	one sample	two samples
mean	$SE = \frac{s}{\sqrt{n}}$	$SE = \sqrt{\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}}$
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 .
- ▶ When working with proportions,
 - ▶ if doing a hypothesis test, p comes from the null hypothesis
 - ▶ if constructing a confidence interval, use \hat{p} instead