

Inference for a Difference in Proportion

Inference for a Difference in Proportions

Learning Objectives.

- Write out the null and alternative hypothesis for Two Categorical Variables
- Simulate the Null Distribution
- Find the p-value and Assess the Strength of Evidence
- Calculate the z test statistic
- Create and Interpret a Confidence Interval for the difference in proportions

Background

In “Helmet Use and Risk of Head Injuries in Alpine Skiers and Snowboarders” by Sullheim et. al., Journal of the American Medical Association, Vol. 295, No. 8, we can see the results from a random sample 3562 skiers and snowboarders involved in accidents.

	Head Injury	No Head Injury	Total
Wore Helmet	96	656	752
Did Not Wear Helmet	480	2330	2810
Total	576	2986	3562

Is there evidence that safety helmet use reduces the risk of head injury for skiers and snowboarders?

Vocabulary Review

1. What is the explanatory variable?
2. What is the response variable?
3. Is this an experiment or observational study?
4. What is the scope of inference for this study?

5. What is the conditional proportion of skiers/snowboarders with a head injury that wore a helmet?
6. What is the conditional proportion of skiers/snowboarders with a head injury that did not wear a helmet?

```
library(readr)
library(car)
injury <- read.csv("../data/head_injury.csv")
injury$Helmet <- factor(injury$Helmet)
injury$Injury <- factor(injury$Injury)
```

Ask a Research Question

In this study we are looking at the relationship between two groups or two parameters (π_1 and π_2). Remember we define the parameter as the true proportion of observational units that represent the variable of interest.

7. What is the variable of interest in this study?
8. Write the two parameters of interest for this study.

π_1 -

π_2 -

When comparing two groups, we assume the two parameters are equal in the null hypothesis. There is no association between the variables.

9. Write the null hypothesis out in words using your answers to question 8.

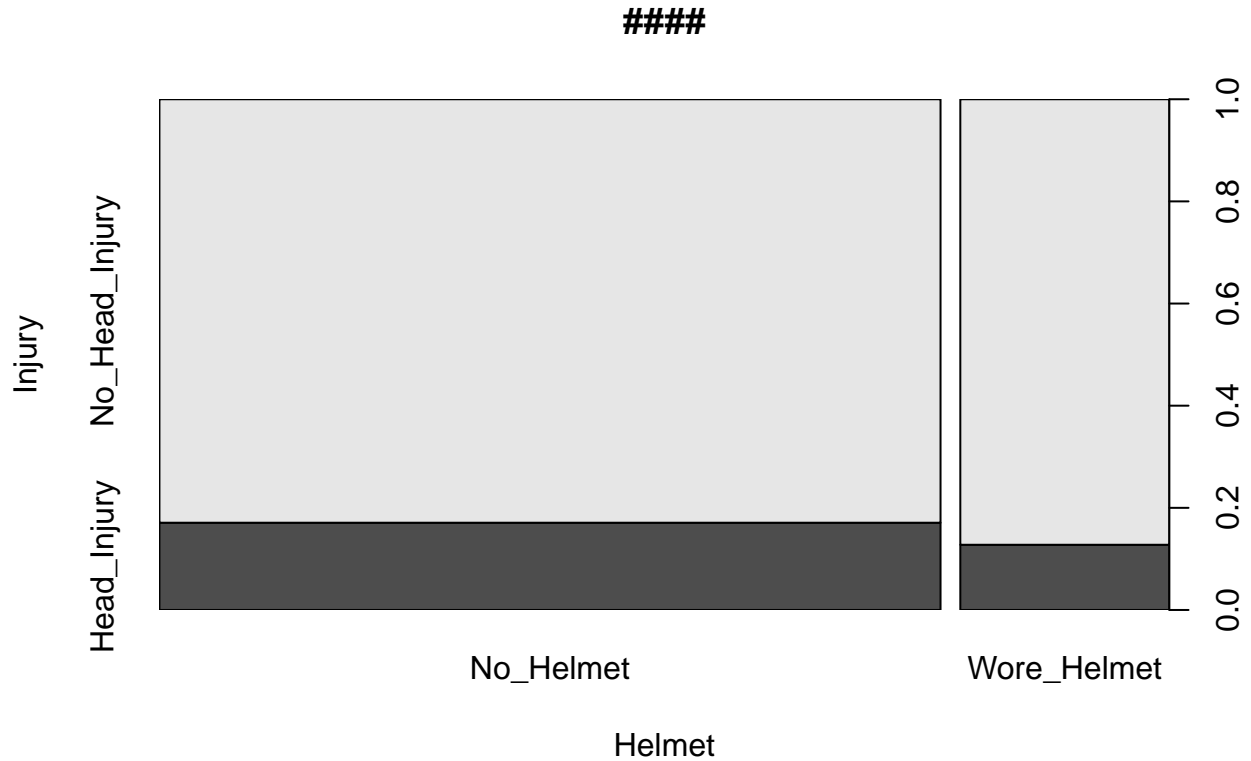
10. What is the research question?

11. Based on the research question fill in the appropriate sign for the alternative hypothesis:

Ha: $\pi_1 \pi_2$

Summarize and Visualize the data

```
plot(Injury~Helmet
     , data = injury, main="####" #Make sure to title your graph
     )
```



12. Fill in the #### with the appropriate variables to plot a segmented bar plot of injury by helmet use.

13. Based on the bar plot, Does there appear to be an association between helmet use and head injury? Explain.

14. Calculate the point estimate for this study. We will use helmet use minus no helmet use as the order of

subtraction.

15. What is the notation used for the value calculated in question 14?

To test the null hypothesis we could use simulation methods as we did with a single categorical variable. In this activity we will focus on theory-based methods. Like with a single proportion, the difference in proportions can be mathematically modeled using the normal distribution if certain conditions are met.

Conditions for the sample distribution of $\hat{p}_1 - \hat{p}_2$

1. *Independence* The data are independent within and between the two groups.
2. *Success-Failure Condition* The success-failure condition holds for each group.
3. Is the independence condition met? Explain your answer.

17. Is the success-failure condition met for each group? Explain your answer.

To calculate the test statistic we use,

$$z = \frac{\text{pointestimate} - \text{nullvalue}}{SE}$$

where the standard error is calculated using the pooled proportion of successes.

$$SE(\hat{p}_1 - \hat{p}_2) = \sqrt{\left(\frac{\hat{p}_{pool}(1-\hat{p}_{pool})}{n_1-1}\right) + \left(\frac{\hat{p}_{pool}(1-\hat{p}_{pool})}{n_2-1}\right)}$$

$$\text{where } \hat{p}_{pool} = \frac{\text{number of "successes"}}{\text{number of cases}} = \frac{\hat{p}_1 n_1 + \hat{p}_2 n_2}{n_1 + n_2}$$

18. Calculate the $SE(\hat{p}_1 - \hat{p}_2)$.

19. Calculate the test statistic.

Not sure how you want to use R to find the p-value

20. How much evidence does the p-value provide against the null hypothesis?

To find a confidence interval for the difference in proportions we will add and subtract the margin of error from the point estimate to find the two endpoints.

$$\hat{p}_1 - \hat{p}_2 \pm SE(\hat{p}_1 - \hat{p}_2)$$

$$\text{where, } SE(\hat{p}_1 - \hat{p}_2) = \sqrt{\left(\frac{\hat{p}_1(1-\hat{p}_1)}{n_1-1} + \frac{\hat{p}_2(1-\hat{p}_2)}{n_2-1}\right)}$$

Note the formula changes When calculating the variability around the statistic in order to calculate a confidence interval. Here use the sample proportions for each group to calculate the standard error for the difference in proportions.

21. Calculate the standard error for a difference in proportions to create a 95% confidence interval.
22. Using the multiplier of 1.96 calculate the 95% confidence interval for the difference in true proportion of head injuries for those that used helmets minus those who did not.
23. Interpret the confidence interval found in question 22 in context of the problem.
24. Write a conclusion to the research question.

Types of Errors

Hypothesis tests are not flawless. In a hypothesis test, there are two competing hypotheses: the null and alternative. We make a decision about which might be true, but we may choose incorrectly.

		Test Conclusion	
Truth	H_0 true	good decision	Type 1 Error
	H_A true	Type 2 Error	good decision

A Type 1 Error is rejecting the null hypothesis when H_0 is actually true. A Type 2 Error is failing to reject the null hypothesis when the alternative is actually true.

25. Using a significance level of 0.05, what decision do you make in regards to the null hypothesis?

26. What type of error could we have made?

27. Write this error in context of the problem.