

Inference for categorical data

Getting Started

Load packages

In this lab, we will explore and visualize the data using the **tidyverse** suite of packages, and perform statistical inference using **infer**. The data can be found in the companion package for OpenIntro resources, **openintro**.

Let's load the packages.

```
library(tidyverse)
library(openintro)
library(infer)
data('yrbss', package='openintro')
set.seed(1234)
```

The data

You will be analyzing the same dataset as in the previous lab, where you delved into a sample from the Youth Risk Behavior Surveillance System (YRBSS) survey, which uses data from high schoolers to help discover health patterns. The dataset is called **yrbss**.

1. What are the counts within each category for the amount of days these students have texted while driving within the past 30 days? i ### Exercise 1

```
yrbss %>%
  group_by(text_while_driving_30d) %>%
  tally()
```

```
## # A tibble: 9 x 2
##   text_while_driving_30d      n
##   <chr>                <int>
## 1 0                      4792
## 2 1-2                    925
## 3 10-19                  373
## 4 20-29                  298
## 5 3-5                    493
## 6 30                     827
## 7 6-9                    311
## 8 did not drive         4646
## 9 <NA>                   918
```

2. What is the proportion of people who have texted while driving every day in the past 30 days and never wear helmets?

Exercise 2

Proportion of people who texted while driving every day in the past 30 days and never wear helmets is 7.12.

```
yrbss %>%
  filter(! is.na(text_while_driving_30d)) %>% # get rid of NAs
  filter(helmet_12m == "never") %>% # never wear helmet
  count(text_while_driving_30d) %>% # text while driving every day
  mutate(prop = n / sum(n))
```

```
## # A tibble: 8 x 3
##   text_while_driving_30d     n  prop
##   <chr>             <int> <dbl>
## 1 0                 2566 0.395
## 2 1-2                515 0.0792
## 3 10-19              207 0.0318
## 4 20-29              180 0.0277
## 5 3-5                281 0.0432
## 6 30                 463 0.0712
## 7 6-9                175 0.0269
## 8 did not drive     2116 0.325
```

Remember that you can use `filter` to limit the dataset to just non-helmet wearers. Here, we will name the dataset `no_helmet`.

```
no_helmet <- yrbss %>%
  filter(helmet_12m == "never")
```

Also, it may be easier to calculate the proportion if you create a new variable that specifies whether the individual has texted every day while driving over the past 30 days or not. We will call this variable `text_ind`.

```
no_helmet <- no_helmet %>%
  mutate(text_ind = ifelse(text_while_driving_30d == "30", "yes", "no"))
```

Inference on proportions

When summarizing the YRBSS, the Centers for Disease Control and Prevention seeks insight into the population *parameters*. To do this, you can answer the question, “What proportion of people in your sample reported that they have texted while driving each day for the past 30 days?” with a statistic; while the question “What proportion of people on earth have texted while driving each day for the past 30 days?” is answered with an estimate of the parameter.

The inferential tools for estimating population proportion are analogous to those used for means in the last chapter: the confidence interval and the hypothesis test.

```
no_helmet %>%
  filter(! is.na(text_ind)) %>%
  specify(response = text_ind, success = "yes") %>%
  generate(reps = 1000, type = "bootstrap") %>%
  calculate(stat = "prop") %>%
  get_ci(level = 0.95)
```

```
## # A tibble: 1 x 2
##   lower_ci upper_ci
##   <dbl>    <dbl>
## 1    0.0652    0.0777
```

Note that since the goal is to construct an interval estimate for a proportion, it's necessary to both include the `success` argument within `specify`, which accounts for the proportion of non-helmet wearers than have consistently texted while driving the past 30 days, in this example, and that `stat` within `calculate` is here "prop", signaling that you are trying to do some sort of inference on a proportion.

3. What is the margin of error for the estimate of the proportion of non-helmet wearers that have texted while driving each day for the past 30 days based on this survey?

Exercise 3

The 95% confidence interval is (0.0652, 0.0777), so the margin of error is $(0.0777 - 0.0652) / 2$ or 0.00625.

4. Using the `infer` package, calculate confidence intervals for two other categorical variables (you'll need to decide which level to call "success", and report the associated margins of error. Interpret the interval in context of the data. It may be helpful to create new data sets for each of the two countries first, and then use these data sets to construct the confidence intervals.

Exercise 4.1

Let's compute the 95% confidence interval for students with good sleep habits (8 hours or more per night). First, we create a new variable that indicates good nightly sleep.

```
good_sleep_ind <- yrbss %>%
  mutate(good_sleep = ifelse(school_night_hours_sleep == '8', 'yes',
                             ifelse(school_night_hours_sleep == '9', 'yes',
                                     ifelse(school_night_hours_sleep == '10+', 'yes', 'no'))))
```

Next, compute proportion confidence intervals. The 95% confidence level for the proportion of students with good sleep habits (8 hours or more nightly sleep) based on this survey is (0.298, 0.314).

We can be 95% confident that the population proportion of students with good sleep habits is between 29.8% and 31.4%.

```
good_sleep_ind %>%
  filter(! is.na(good_sleep)) %>%
  specify(response = good_sleep, success = "yes") %>%
  generate(reps = 1000, type = "bootstrap") %>%
  calculate(stat = "prop") %>%
  get_ci(level = 0.95)
```

```
## # A tibble: 1 x 2
##   lower_ci upper_ci
##   <dbl>    <dbl>
## 1    0.298    0.314
```

Exercise 4.2

Let's compute the 95% confidence interval for students with bad TV habits (4 hours or more per night). First, we create a new variable that indicates bad TV habit

```
bad_tv_habit_ind <- yrbss %>%
  mutate(bad_tv_habit = ifelse(hours_tv_per_school_day == '4', 'yes',
                                ifelse(school_night_hours_sleep == '5+', 'yes', 'no')))
```

Next, compute proportion confidence intervals. The 95% confidence level for the proportion of students with bad TV habits (4 hours or more watching TV each night) based on this survey is (0.0799, 0.090).

We can be 95% confident that the population proportion of students with bad sleep habits is between 7.99% and 9.0%.

```
bad_tv_habit_ind %>%
  filter(! is.na(bad_tv_habit)) %>%
  specify(response = bad_tv_habit, success = "yes") %>%
  generate(reps = 1000, type = "bootstrap") %>%
  calculate(stat = "prop") %>%
  get_ci(level = 0.95)
```

```
## # A tibble: 1 x 2
##   lower_ci upper_ci
##   <dbl>    <dbl>
## 1    0.0799    0.0900
```

How does the proportion affect the margin of error?

Imagine you've set out to survey 1000 people on two questions: are you at least 6-feet tall? and are you left-handed? Since both of these sample proportions were calculated from the same sample size, they should have the same margin of error, right? Wrong! While the margin of error does change with sample size, it is also affected by the proportion.

Think back to the formula for the standard error: $SE = \sqrt{p(1-p)/n}$. This is then used in the formula for the margin of error for a 95% confidence interval:

$$ME = 1.96 \times SE = 1.96 \times \sqrt{p(1-p)/n}.$$

Since the population proportion p is in this ME formula, it should make sense that the margin of error is in some way dependent on the population proportion. We can visualize this relationship by creating a plot of ME vs. p .

Since sample size is irrelevant to this discussion, let's just set it to some value ($n = 1000$) and use this value in the following calculations:

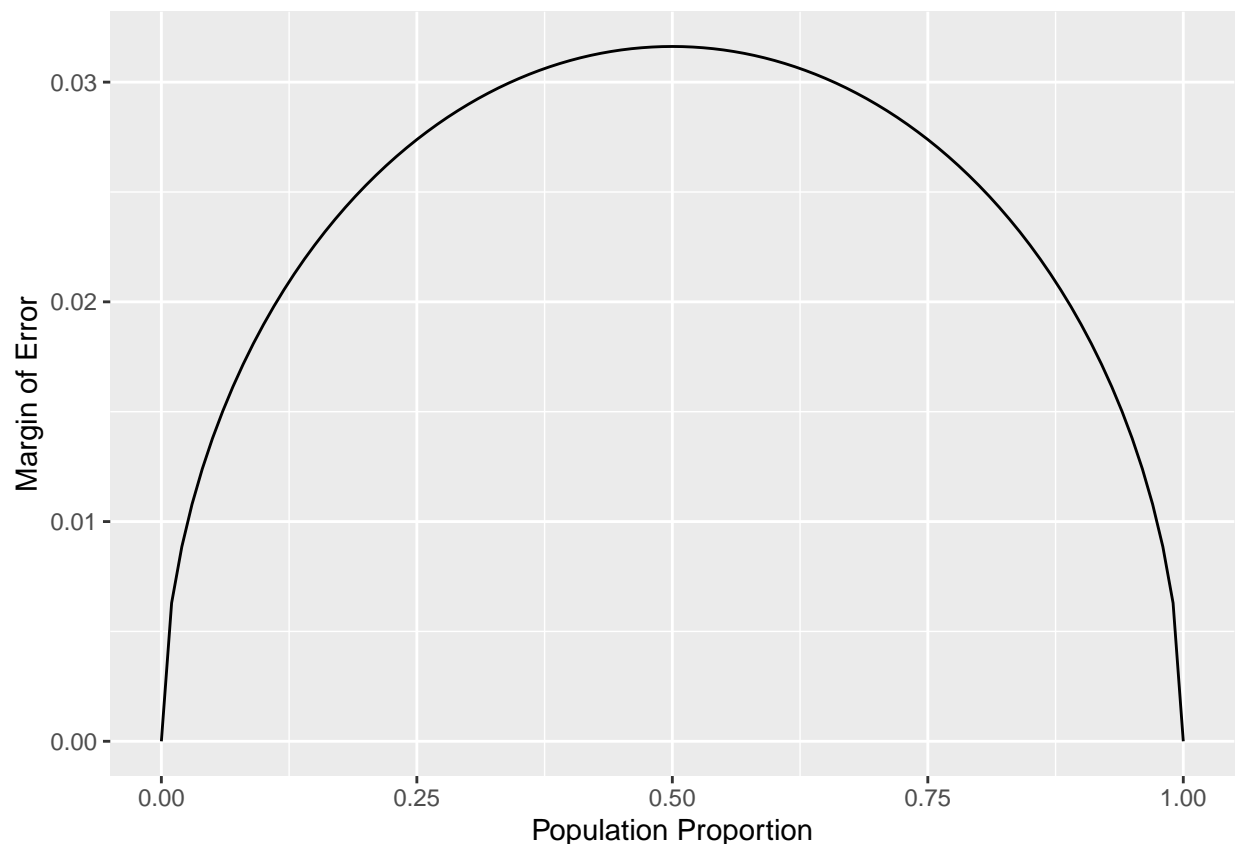
```
n <- 1000
```

The first step is to make a variable p that is a sequence from 0 to 1 with each number incremented by 0.01. You can then create a variable of the margin of error (me) associated with each of these values of p using the familiar approximate formula ($ME = 2 \times SE$).

```
p <- seq(from = 0, to = 1, by = 0.01)
me <- 2 * sqrt(p * (1 - p)/n)
```

Lastly, you can plot the two variables against each other to reveal their relationship. To do so, we need to first put these variables in a data frame that you can call in the `ggplot` function.

```
dd <- data.frame(p = p, me = me)
ggplot(data = dd, aes(x = p, y = me)) +
  geom_line() +
  labs(x = "Population Proportion", y = "Margin of Error")
```



- Describe the relationship between `p` and `me`. Include the margin of error vs. population proportion plot you constructed in your answer. For a given sample size, for which value of `p` is margin of error maximized?

Exercise 5

The relationship is quadratic as shown by the upside down parabola, the ME increases as the population proportion approaches 0.5 from left and from right. The margin of error is maximized at $p = 0.5$.

Success-failure condition

We have emphasized that you must always check conditions before making inference. For inference on proportions, the sample proportion can be assumed to be nearly normal if it is based upon a random sample

of independent observations and if both $np \geq 10$ and $n(1 - p) \geq 10$. This rule of thumb is easy enough to follow, but it makes you wonder: what's so special about the number 10?

The short answer is: nothing. You could argue that you would be fine with 9 or that you really should be using 11. What is the “best” value for such a rule of thumb is, at least to some degree, arbitrary. However, when np and $n(1 - p)$ reaches 10 the sampling distribution is sufficiently normal to use confidence intervals and hypothesis tests that are based on that approximation.

You can investigate the interplay between n and p and the shape of the sampling distribution by using simulations. Play around with the following app to investigate how the shape, center, and spread of the distribution of \hat{p} changes as n and p changes.

6. Describe the sampling distribution of sample proportions at $n = 300$ and $p = 0.1$. Be sure to note the center, spread, and shape.

Exercise 6

The sampling distribution of sample proportions at $n = 300$ and $p = 0.1$ is Normal but not smooth. The center is around 0.1 and min/max is approximately (0.01, 0.17)

7. Keep n constant and change p . How does the shape, center, and spread of the sampling distribution vary as p changes. You might want to adjust min and max for the x -axis for a better view of the distribution.

Exercise 7

Increasing the population proportion makes the distribution closer to Normal up to 0.5. Increasing beyond 0.5 makes the distribution no longer be Normal.

8. Now also change n . How does n appear to affect the distribution of \hat{p} ?

Exercise 8

As the sample size increased \hat{p} distribution gets closer to Normal with smaller spread. This is to be expected as the standard error decreases.

More Practice

For some of the exercises below, you will conduct inference comparing two proportions. In such cases, you have a response variable that is categorical, and an explanatory variable that is also categorical, and you are comparing the proportions of success of the response variable across the levels of the explanatory variable. This means that when using `infer`, you need to include both variables within `specify`.

9. Is there convincing evidence that those who sleep 10+ hours per day are more likely to strength train every day of the week? As always, write out the hypotheses for any tests you conduct and outline the status of the conditions for inference. If you find a significant difference, also quantify this difference with a confidence interval.

Exercise 9

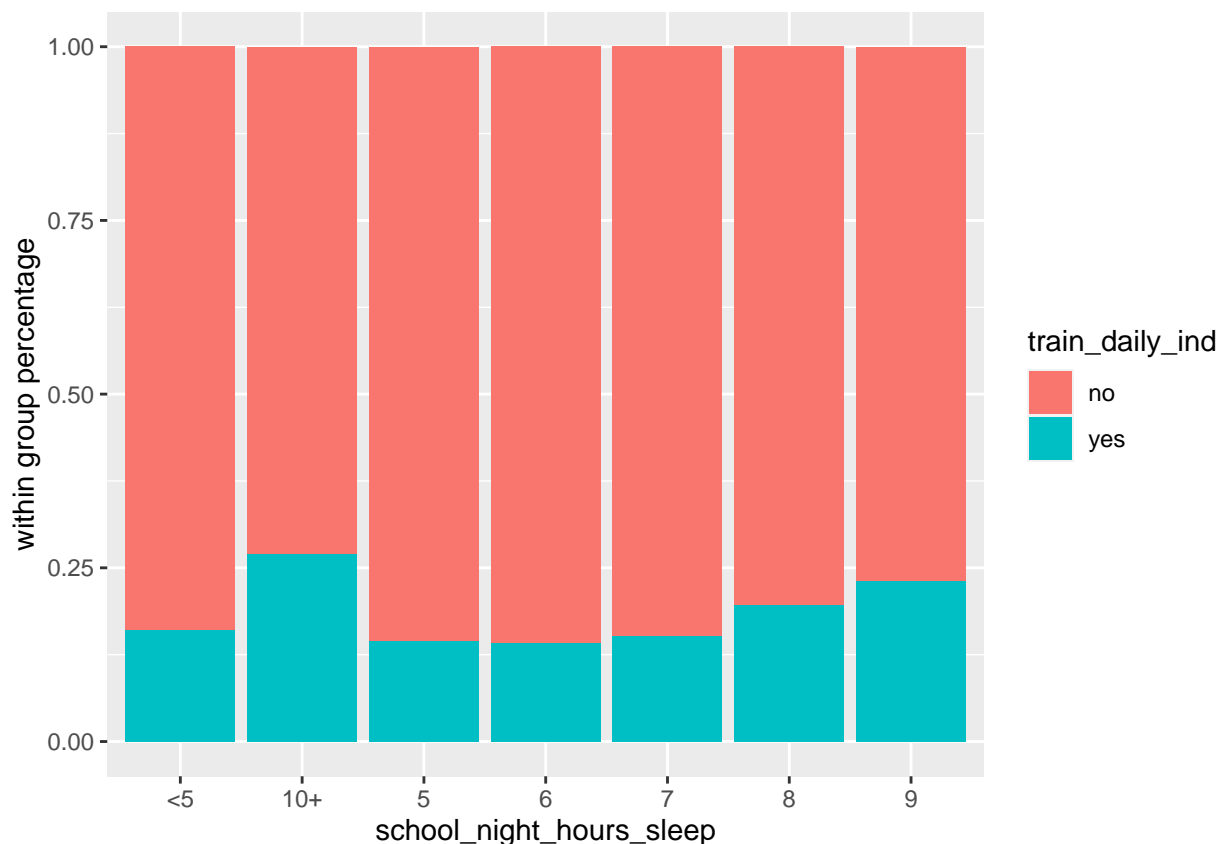
Explanatory variable is 'those who sleep 10+ hours per day', response variable is 'strength train every day of the week'.

Ho: there is no difference in strength training every day of the week between people who sleep 10+ hours per day and those who don't
Ha: those who sleep 10+ hours per day are more likely to strength train every day of the week

```
sleep_exercise <- yrbss %>%  
  filter(! is.na(school_night_hours_sleep)) %>%  
  filter(! is.na(strength_training_7d)) %>%  
  mutate(sleep_long_ind = ifelse(school_night_hours_sleep == "10+", "yes", "no")) %>%  
  mutate(train_daily_ind = ifelse(strength_training_7d == "7", "yes", "no"))
```

Let's see a quick barplot with normalized within-group percentages. It appears there is a difference in strength training every day of the week between people who sleep 10+ hours per day and those who don't.

```
sleep_exercise %>%  
  select(train_daily_ind, school_night_hours_sleep) %>%  
  ggplot(aes(x=school_night_hours_sleep, fill = train_daily_ind)) +  
  geom_bar(position = "fill") +  
  ylab("within group percentage")
```



Since both variables, *train_daily_ind* and *sleep_long_ind*, are categorical we can use the Chi-square test framework. The assumption for Chi-square testing are:

1. Both variables are categorical. Check.
2. All observations are independent. Check.
3. It's assumed that individuals can only belong to one cell in the contingency table. That is, cells in the table are mutually exclusive – an individual cannot belong to more than one cell. Check.
4. There must be at least 5 frequencies in each cell. Check.

Calculate observed Chisq statistic. It is 23.2

```
observed_indep_statistic <- sleep_exercise %>%
  specify(train_daily_ind ~ sleep_long_ind, success = "yes") %>%
  #specify(train_daily_ind ~ school_night_hours_sleep) %>%
  hypothesize(null = "independence") %>%
  calculate(stat = "Chisq")
```

```
observed_indep_statistic
```

```
## Response: train_daily_ind (factor)
## Explanatory: sleep_long_ind (factor)
## Null Hypothesis: independence
## # A tibble: 1 x 1
##   stat
##   <dbl>
## 1  23.2
```

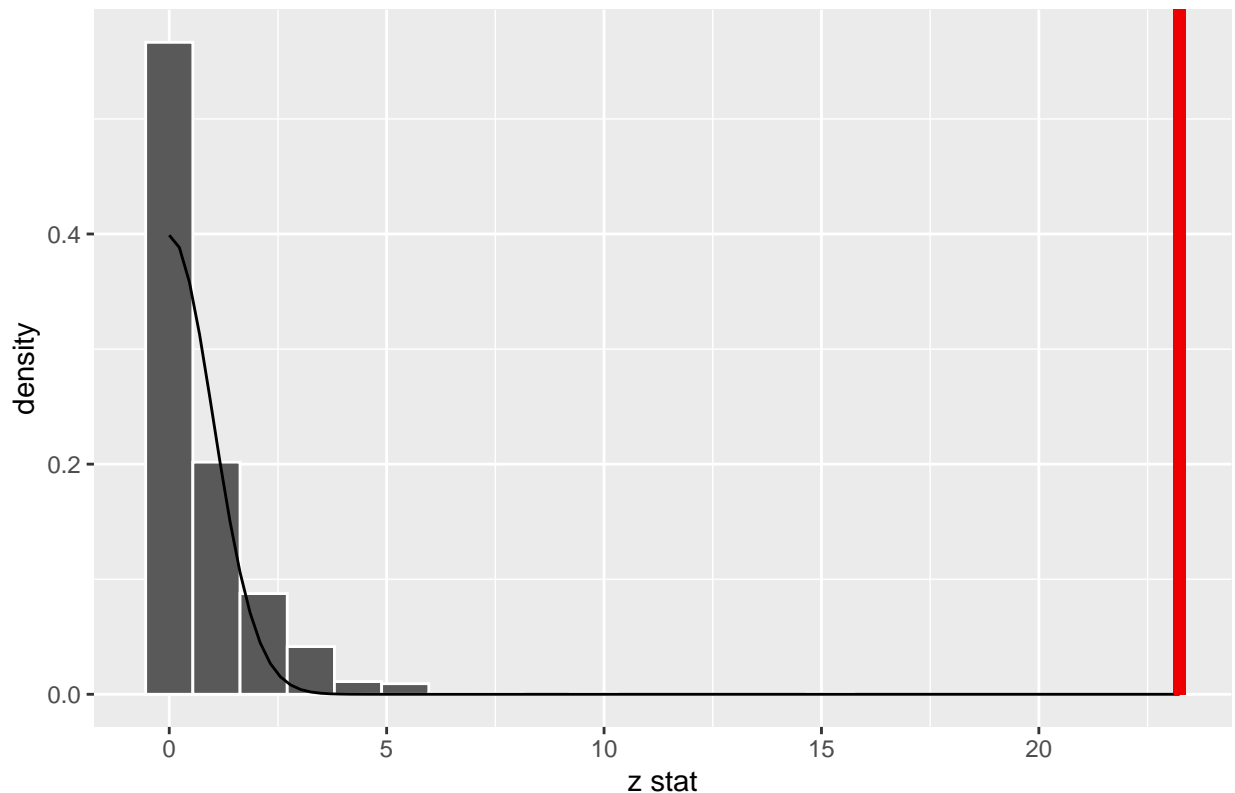
Generate the null distribution by permuting the response and explanatory variables. The null distribution is what we would expect to see if the two variables were truly independent.

```
null_dist_sim <- sleep_exercise %>%
  specify(train_daily_ind ~ sleep_long_ind, success = "yes") %>%
  hypothesize(null = "independence") %>%
  generate(reps = 1000, type = "permute") %>%
  calculate(stat = "Chisq")
```

Let's plot the simulated null distribution.

```
null_dist_sim %>%
  visualize(method = "both") +
  shade_p_value(observed_indep_statistic,
    direction = "greater")
```


Simulation-Based and Theoretical z Null Distributions



Plot of simulated null distribution shows that the observed test statistic would be unlikely if there was no association between 10+ hours of sleep each day and exercising 7 days a week.

```
p_value_independence <- null_dist_sim %>%  
  get_p_value(obs_stat = observed_indep_statistic,  
             direction = "greater")  
p_value_independence
```

```
## # A tibble: 1 x 1  
##   p_value  
##   <dbl>  
## 1      0
```

10. Let's say there has been no difference in likeliness to strength train every day of the week for those who sleep 10+ hours. What is the probability that you could detect a change (at a significance level of 0.05) simply by chance? *Hint:* Review the definition of the Type 1 error.

Exercise 10

The probability that we could detect a change (at a significance level of 0.05) simply by chance. i.e. a Type 1 error (rejecting the null Hypothesis when it is true), could occur with a probability of 5%.

11. Suppose you're hired by the local government to estimate the proportion of residents that attend a religious service on a weekly basis. According to the guidelines, the estimate must have a margin of

error no greater than 1% with 95% confidence. You have no idea what to expect for p . How many people would you have to sample to ensure that you are within the guidelines?

Hint: Refer to your plot of the relationship between p and margin of error. This question does not require using a dataset.

Exercise 11

From earlier plot we know that worst case ME is when p is 0.5. Let's use $ME = 0.01$ and $p = 0.5$ to solve for n . Based on these assumptions we would want a sample size that is the *minimum* of 9604 and 10% of total population.

$$ME = 1.96 \times SE = 1.96 \times \sqrt{p(1-p)/n}$$

$$0.01 = 1.96 \times \sqrt{\frac{0.5(1-0.5)}{n}}$$

$$0.01^2 = 1.96^2 \times \frac{0.5(1-0.5)}{n}$$

$$n = \frac{1.96^2 \times 0.5(1-0.5)}{0.01^2}$$

$$n = \frac{3.8416 \times 0.25}{0.0001}$$

$$n = 9604$$
