

# 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)
```

### 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?

**Insert your answer here** The amount of days are for 0 days a count of 4792 students, for 1-2 days a count of 925 students, for 3-5 days a count of 493 students, for 6-9 days a count of 311 students, for 10-19 days a count of 373 students, for 20-29 days a count of 298 students, for 30 days a count of 827 students, for “did not drive” a count of 4646, for “N” a count of 918.

```
table(yrbss$text_while_driving_30d) |>
  prop.table()
```

```
##
##           0           1-2           10-19           20-29           3-5
## 0.37836557 0.07303593 0.02945124 0.02352941 0.03892617
##           30           6-9 did not drive
## 0.06529807 0.02455586 0.36683774
```

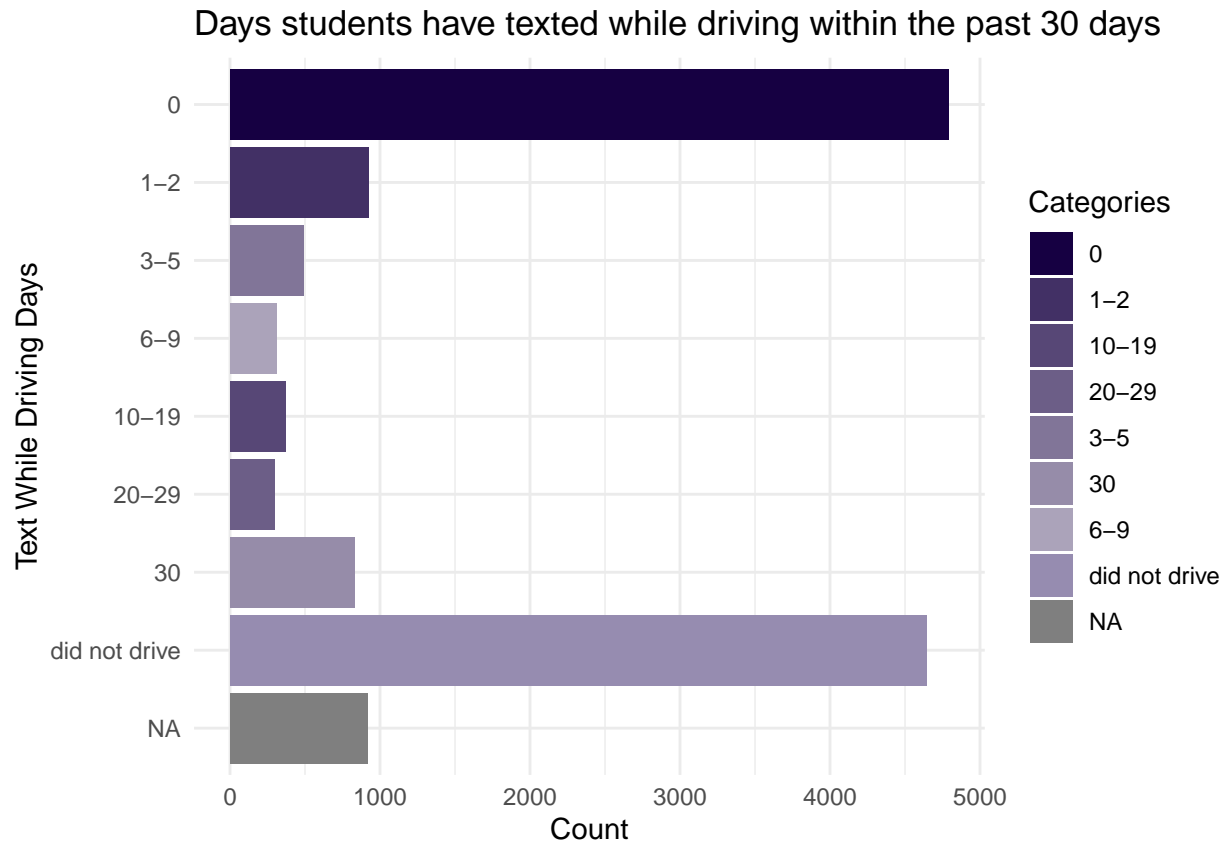
```
summaryt1 <- yrbss %>%
  count(text_while_driving_30d) %>%
  mutate(p = n / sum(n))
#summaryt1

sortindex <- c(1,2,5,6,3,7,4,8,9)
```

```
summarytf <- tibble(summaryt1,sortindex)
summarytf <- summarytf |>
  arrange((sortindex))
summarytf|>
  arrange((sortindex))
```

```
## # A tibble: 9 x 4
##   text_while_driving_30d      n      p sortindex
##   <chr>          <int> <dbl>    <dbl>
## 1 0              4792 0.353      1
## 2 1-2            925 0.0681     2
## 3 3-5            493 0.0363     3
## 4 6-9            311 0.0229     4
## 5 10-19          373 0.0275     5
## 6 20-29          298 0.0219     6
## 7 30             827 0.0609     7
## 8 did not drive  4646 0.342      8
## 9 <NA>          918 0.0676     9
```

```
ggplot(yrbss,aes(x=text_while_driving_30d, fill =text_while_driving_30d ))+
  geom_bar()+
  scale_x_discrete(limits = rev(summarytf$text_while_driving_30d))+
  ggtitle('Days students have texted while driving within the past 30 days') +
  ylab('Count')+
  xlab('Text While Driving Days')+
  theme_minimal()+
  scale_fill_manual('Categories', values=c('#160042','#423065','#574776','#6C5E87',
                                           '#817598','#968CA9','#ABA3Ba','#968CB0','#6C5E85'))+
  coord_flip()
```



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

**Insert your answer here** The proportion of people who have texted while driving every day in the past 30 days and never wear helmets for 30 days is 3.41% .

```
#q2_nwh <- yrbss|>
# filter(yrbss$helmet_12m == "never")

#ggplot(q2_nwh,aes(x=text_while_driving_30d, fill =text_while_driving_30d ))+
# geom_bar()+
# scale_x_discrete(limits = rev(summarytf$text_while_driving_30d))+
# ggtitle('Days students have texted while driving within the past 30 days\nand never wear helmets') +
# ylab('Count')+
# xlab('Text While Driving Days')+
# theme_minimal()+
# scale_fill_manual('Categories', values=c('#160042', '#423065', '#574776', '#6C5E87',
#                                           '#817598', '#968CA9', '#ABA3Ba', '#968CB0', '#6C5E85'))+
# coord_flip()

q2 <- yrbss|>
filter(yrbss$helmet_12m == "never",yrbss$text_while_driving_30d=="30" )|>
count(text_while_driving_30d) |>
```

```
mutate(p_hat = n / nrow(yrbss)) #proportion from yrbss

q2
```

```
## # A tibble: 1 x 3
##   text_while_driving_30d    n p_hat
##   <chr>                <int> <dbl>
## 1 30                  463 0.0341
```

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

```
data('yrbss', package='openintro')
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 %>%
  drop_na(text_ind) %>% # Drop missing values
  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 0.0647 0.0780
```

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?

**Insert your answer here** We are currently working with 95% confidence interval so we will be using a Z value of 1.96 in order to calculate the margin of error. 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 is 0.003052.

```
no_helmet_sa1 <- no_helmet |>
  filter(text_ind == "yes") #Last 30 days only (yes)

n_yrbss <- nrow(yrbss)
p_no_helmet_sa1 <- nrow(no_helmet_sa1) / n_yrbss

p_no_helmet_sa1 # p
```

```
## [1] 0.03408673
```

```
n_yrbss # n
```

```
## [1] 13583
```

```
1.96*sqrt(p_no_helmet_sa1*(1-p_no_helmet_sa1)/n_yrbss)
```

```
## [1] 0.003051546
```

```
#90% - 1.64, 95% - 1.96, 99% - 2.58
```

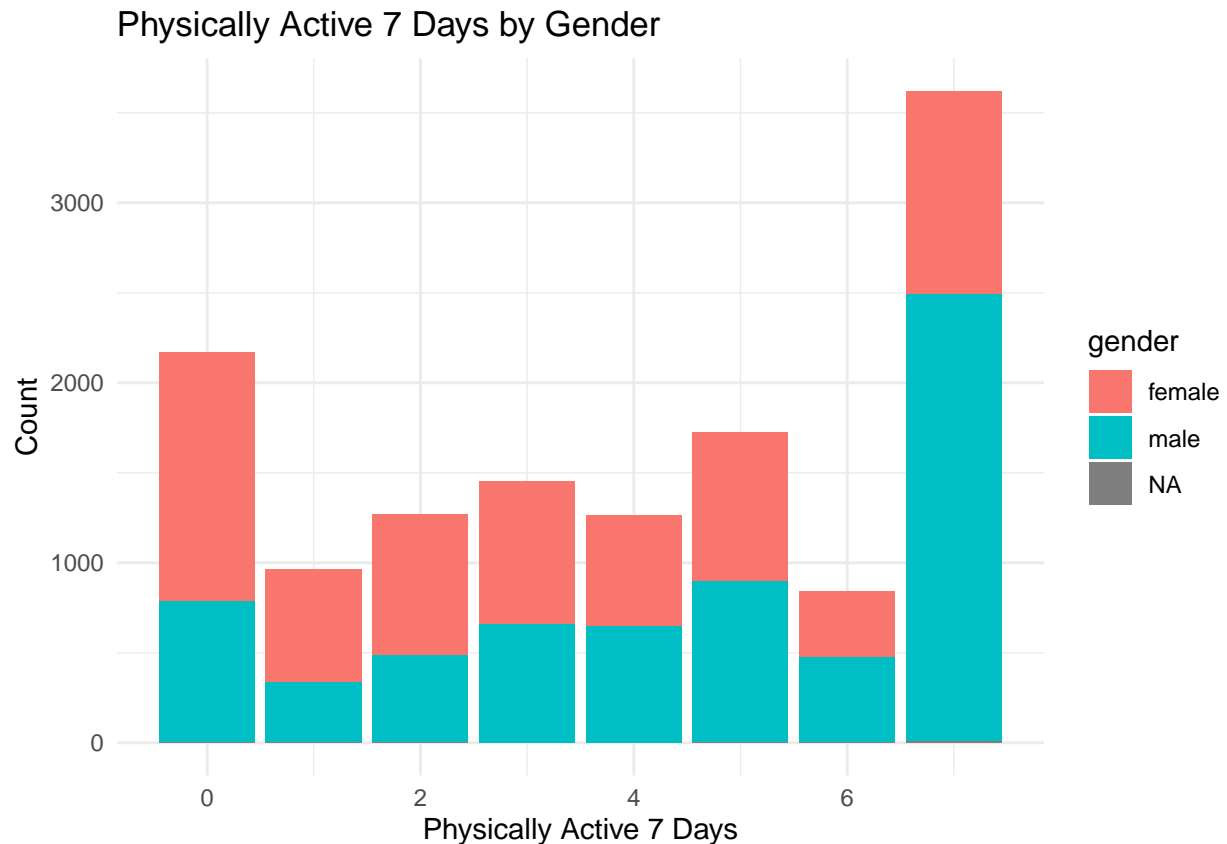
- Using the `infer` package, calculate confidence intervals for two other categorical variables (you'll need to decide which level to call "success", and report the associated margins of error. Interpret the interval in context of the data. It may be helpful to create new data sets for each of the two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

**Insert your answer here** We are testing how many people don't exercise for a period of 7 days of a female gender. There is 95% confident that the true population proportion of a person that don't exercise is between 61.5% and 65.6% with 0.0058085733 margin of error.

```
#the proportion for physically_active_7d
yrbss %>%
  count(physically_active_7d) %>%
  mutate(p = n / sum(n))
```

```
## # A tibble: 9 x 3
##   physically_active_7d     n     p
##           <int> <int> <dbl>
## 1             0  2172 0.160
## 2             1   962 0.0708
## 3             2  1270 0.0935
## 4             3  1451 0.107
## 5             4  1265 0.0931
## 6             5  1728 0.127
## 7             6   840 0.0618
## 8             7  3622 0.267
## 9            NA   273 0.0201
```

```
ggplot(yrbss,aes(x=physically_active_7d, fill =gender ))+
  geom_bar()+
  ggtitle('Physically Active 7 Days by Gender') +
  ylab('Count')+
  xlab('Physically Active 7 Days')+
  theme_minimal()
```



```
#people with no activity
data('yrbss', package='openintro')
nphysically_active <- yrbss %>%
  filter(physically_active_7d == "0") #no activity

#gender of people with no activity
nphysically_active <- nphysically_active %>%
  mutate(text_ind = ifelse(gender == "female", "yes", "no"))

#Interval confidence
nphysically_active %>%
  drop_na(text_ind) %>% # Drop missing values
  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    0.617    0.656

#Margin of error
nphysically_active_sa1 <- nphysically_active |>
  filter(text_ind == "yes") #Only women (yes)

n_yrbss <- nrow(yrbss) #n
p_nphysically_active_sa1 <- nrow(nphysically_active_sa1) / n_yrbss #p

1.96*sqrt(p_nphysically_active_sa1*(1-p_nphysically_active_sa1)/n_yrbss)

## [1] 0.005085733
```

## 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:

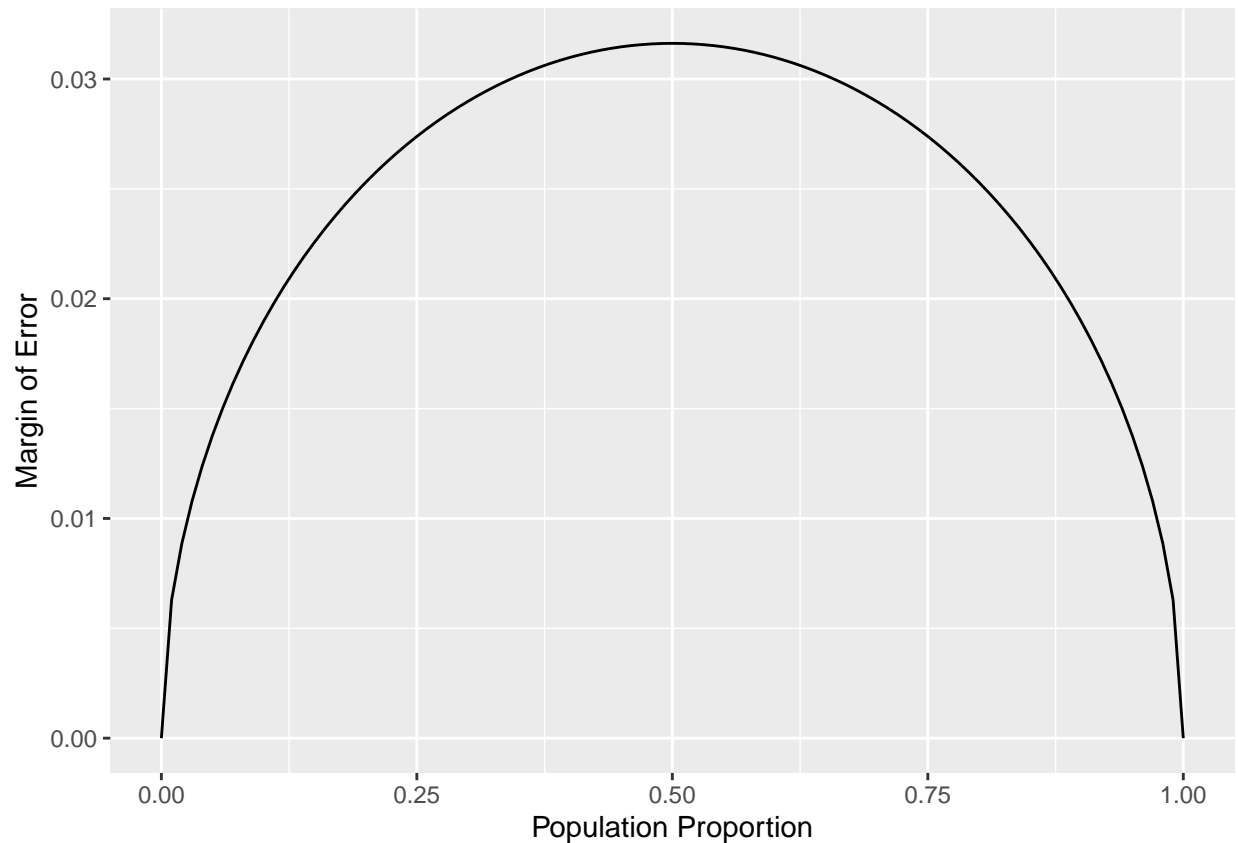
```
n <- 1000
```

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")
```

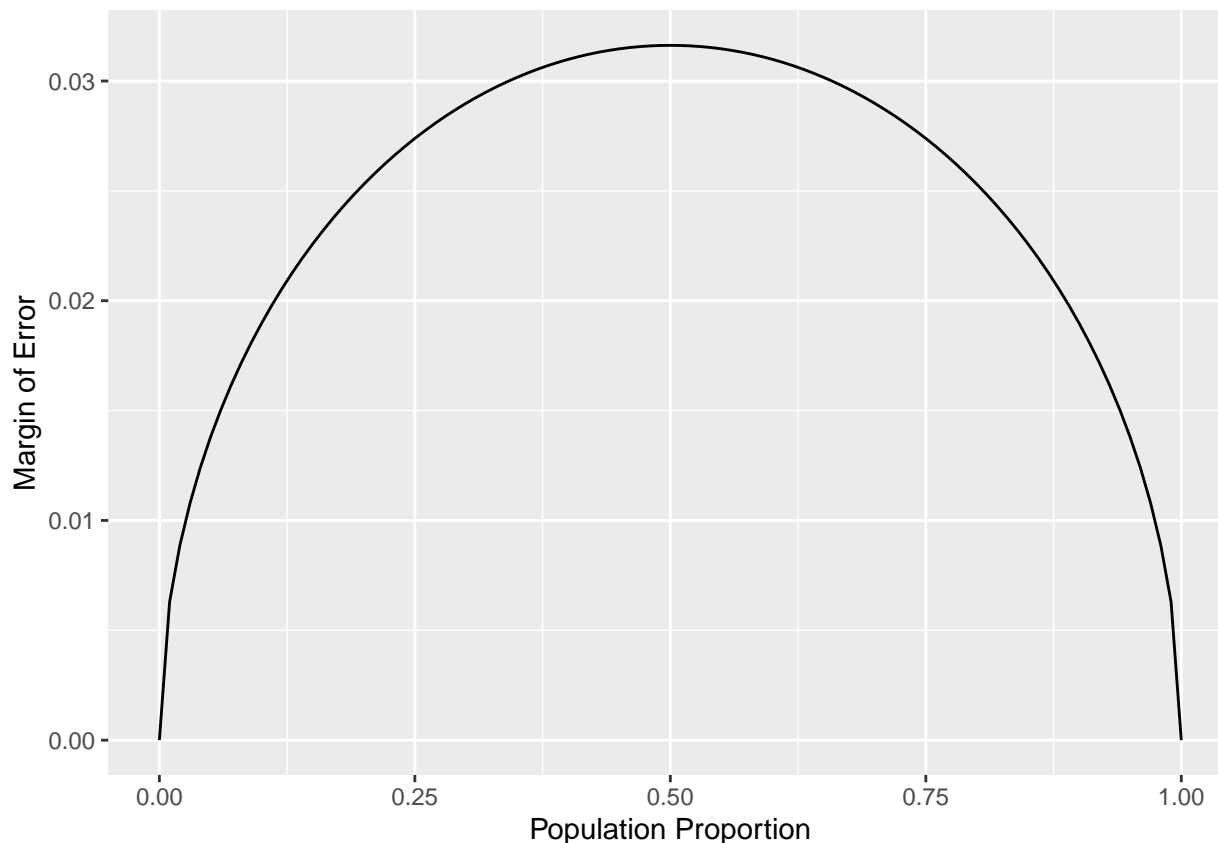


5. Describe the relationship between **p** and **me**. Include the margin of error vs. population proportion plot you constructed in your answer. For a given sample size, for which value of **p** is margin of error maximized?

**Insert your answer here** The vertical axis represents the Margin of Error while the horizontal axis indicates the percentages in the population. Based on plot the highest margin of error occurs where we have 50% of the population proportion, and it disappears in both directions when the population proportion is 0 or 1.

```
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")
```





### 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.

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

**Insert your answer here** the sampling distribution of sample proportions at  $n=300$  and  $p = 0.1$  appears to be symmetric with the center and uni modal.

- 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.

**Insert your answer here** The shape, center, and spread of the sampling distribution doesn't vary with a  $n=300$  and  $p=0.68$ , the sample distribution still appears to be symmetric with the center shifting somewhat to the left and uni modal.

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

**Insert your answer here** when I increased I dont see much change however when I reduced the value for  $n$  it seems less symmetric.

#much of a spread

## 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.

**Insert your answer here** We have the following Hypotheses  $H_0$ : people who sleep 10+ hours per day are more likely to strength train every day of the week  $H_A$ : people who sleep 10+ hours per day are not more likely to strength train every day of the week

We will reject  $H_0$  because there is no enough evidence that those who sleep 10+ hours per day are more likely to strength train every day of the week.

```
#sleep more than 10 hours
data('yrbss', package='openintro')
school_night_hours_sleep <- yrbss %>%
  filter(school_night_hours_sleep == "10+") #sleep more than 10

#strength train every day of the week
school_night_hours_sleep_active <- school_night_hours_sleep %>%
  mutate(text_ind = ifelse(strength_training_7d == "7", "yes", "no"))

school_night_hours_sleep_active %>%
  count(text_ind) %>%
  mutate(p_train = n / sum(n))%>%
  drop_na(text_ind) # the proportion is low 0.266
```

```
## # A tibble: 2 x 3
##   text_ind      n p_train
##   <chr>    <int>   <dbl>
## 1 no       228   0.722
## 2 yes       84   0.266
```

```
#difference in p_train_no(related with Ha) and p_train_y(related H0)
0.722-0.266
```

```
## [1] 0.456
```

```
#Interval confidence
school_night_hours_sleep_active %>%
  drop_na(text_ind) %>% # Drop missing values
  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     0.218     0.321
```

```
table(school_night_hours_sleep$school_night_hours_sleep,school_night_hours_sleep$strength_training_7d) |>
  prop.table()
```

```
##
##           0           1           2           3           4           5
## 10+ 0.32051282 0.05448718 0.09935897 0.09935897 0.05769231 0.07371795
##
##           6           7
## 10+ 0.02564103 0.26923077
```

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.

**Insert your answer here** A type 1 error is when you reject the null hypothesis when  $H_0$  is actually true (5% for  $\alpha=0.05$ ).

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.

**Insert your answer here** In order to to ensure that you are within the guidelines, we need 9604 people.

```
p_ni_w <-0.5
margin_error <-0.01
Z <- 1.96 #95%
Z^2*p_ni_w*(1-p_ni_w)/margin_error^2
```

```
## [1] 9604
```

```
#Using Margin of Error Formula  $ME=z*sqrt(p*(1-p)/n) \Rightarrow n=(p(1-p)/e^2)*z^2$ 
```