

Lesson 5 Example Code

Simulating a Null Distribution:

While the applets make it very easy to simulate a null distribution, I want to provide you with some code to demonstrate how to do a similar thing in *R*. I find that more exposure to code make students better coders (radical thought... I know) so don't just skip over this section. I will use the rock-paper-scissors example in your book for this code.

```
library(tidyverse)

replications_dataframe = NULL

num_reps = 1000
sample_size = 12
null_prob = 1/3
sample_stat = 0.167

for (i in 1:num_reps){

  trial = sample(x = c(1,0),
                size = sample_size,
                prob = c(null_prob, 1-null_prob),
                replace = TRUE)

  trial_proportion = sum(trial)/sample_size

  replications_dataframe = rbind(replications_dataframe, data.frame(trial_proportion))
}
```

After turning on the awesome, this code sets up several “parameters” for our simulation. This is “parameter” in the general sense, not to be confused with statistical parameters. We create a dataframe, or storage spot, to put our simulated proportions. We then define the number of repetitions (1000), the size of each simulated repetitions (12... based on the problem in the book), the null probability (probability of getting scissors under the null hypothesis), and the sample statistic (what we observed in our experiment).

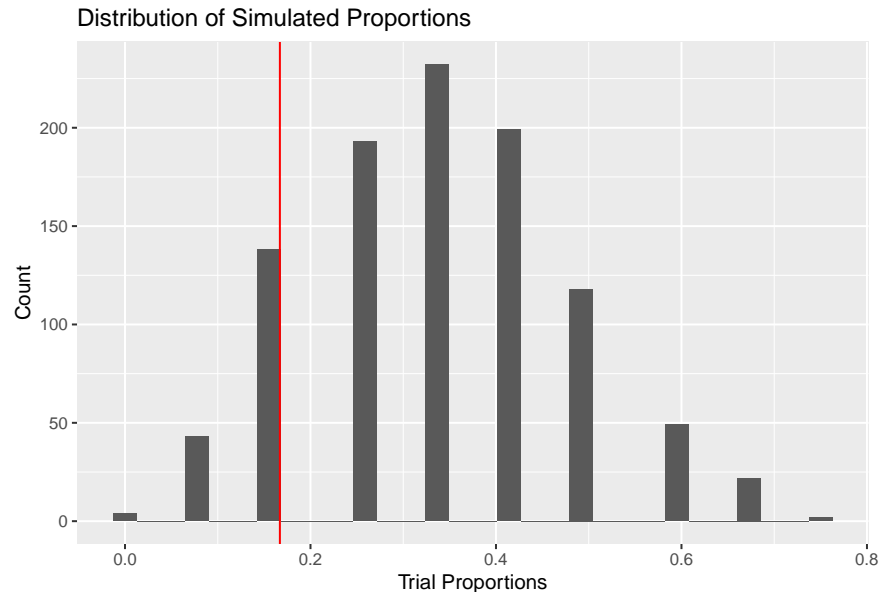
What follows is a **for** loop that serves as the “simulation” in our experiment. The code here is a bit more advanced but nothing insurmountable. The variable *trial* is basically a list of size 12 (sample size) of ones and zeros. This list is populated using the **sample** function which chooses from the vector we give it in *x* with the probability of each elements being chosen as defined by the *prob* parameter in the function. Said another way: the function will choose *1* with a probability of *null_prob* and *0* with the probability of *1-null_prob*. Finally, we need to sample *with replacement* otherwise you can't get twelve numbers out of a list of two.

We calculate the proportion of scissors (ones) in our vector by summing the elements of the vector (again, a bunch of ones and zeros) and dividing by the sample size. Finally, we take this *trial_proportion* and “bind” it to the end of our big list of simulated proportions. What we are left with is 1000 proportions that are simulated under the assumption that the probability of picking scissors is 1/3... the null hypothesis. We can visualize the distribution of these numbers using a histogram.

```

replications_dataframe %>%
  ggplot(aes(x = trial_proportion)) +
  geom_histogram() +
  labs(x = "Trial Proportions", y = "Count", title = "Distribution of Simulated Proportions") +
  geom_vline(xintercept = sample_stat, color = "red")

```



```

#Find the p-value
replications_dataframe %>%
  summarise(pvalue = sum(trial_proportion <= sample_stat) / n())

```

```

## pvalue
## 1 0.185

```

One-proportion z-test

Here is some example code for conducting a one-proportion z-test in *R*. I will use the “Halloween Treats” example from your textbook. Before we can use a theory-based test, we need to check the validity conditions. With sample we clearly have at least 10 successes and 10 failures.

```

sample_proportion = 148/(135+148)

std_null = sqrt(0.5 * (1-0.5) / 283)

z = (sample_proportion - 0.5) / std_null

p = 2 * (1 - pnorm(abs(z)))

p

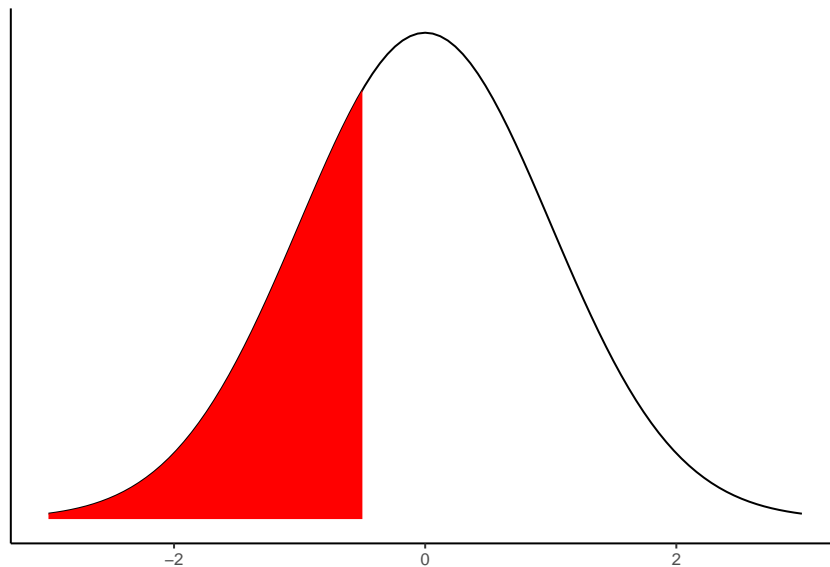
```

```
## [1] 0.4396586
```

Hopefully, the first three steps here are pretty clear from your reading in your book. The last equation (to get the p-value) may be a little new because of the **pnorm** function. Please spend some time experimenting with this function because it is going to become very useful to you in this class. You can see this function as saying “give me an x and I’ll give you an area under the normal curve from negative infinity to that x .” Now maybe you think I’m crazy because functions talk to me but that’s ok as long as you remember that.

Here are a couple examples of using the **pnorm** function. Note: You won't be responsible for knowing all functions below but I'm using a few new ones to illustrate my point.

```
ggplot(NULL, aes(x = c(-3, 3))) +  
  stat_function(fun = dnorm,  
               args = list(mean = 0, sd = 1)) +  
  stat_function(fun = dnorm,  
               args = list(mean = 0, sd = 1),  
               xlim = c(-3, -0.5), #This is (-3, 0) for the purposes of plotting.  
               geom = "area", fill = "red") +  
  labs(x = "", y = "") + theme_classic() +  
  scale_y_continuous(breaks = NULL)
```

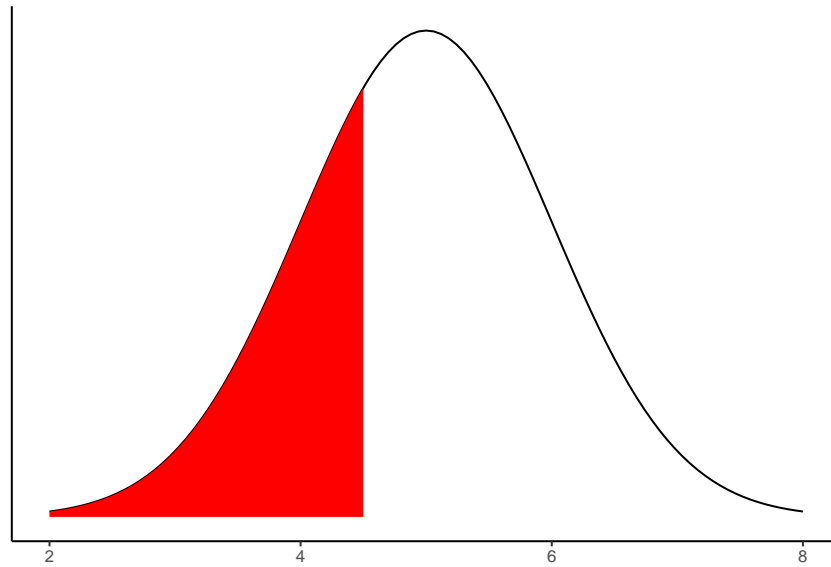


```
pnorm(-0.5)
```

```
## [1] 0.3085375
```

In this case, **pnorm** gives us the area of the region shaded in red above. When we use **pnorm** without any other parameters, it defaults to the *standard normal distribution* which has a mean of 0 and a standard deviation of 1. One of the great things about *R*, however, is that we don't have to "standardize" our sample statistics to calculate our areas. We can also do something like the following two examples if one or both of our parameters are different than the standard normal.

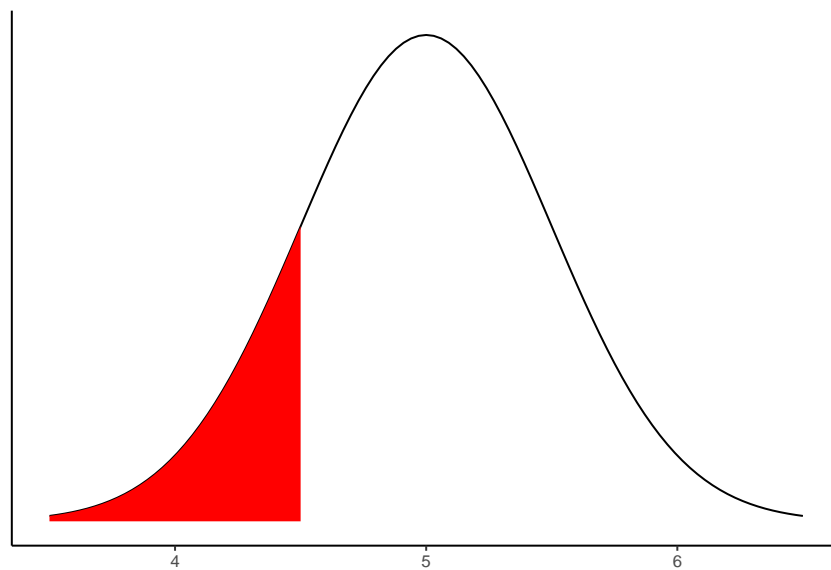
```
ggplot(NULL, aes(x = c(2, 8))) +  
  stat_function(fun = dnorm,  
               args = list(mean = 5, sd = 1)) +  
  stat_function(fun = dnorm,  
               args = list(mean = 5, sd = 1),  
               xlim = c(2, 4.5),  
               geom = "area", fill = "red") +  
  labs(x = "", y = "") + theme_classic() +  
  scale_y_continuous(breaks = NULL)
```



```
pnorm(4.5, mean = 5, sd = 1)
```

```
## [1] 0.3085375
```

```
ggplot(NULL, aes(x = c(3.5, 6.5))) +  
  stat_function(fun = dnorm,  
               args = list(mean = 5, sd = 0.5)) +  
  stat_function(fun = dnorm,  
               args = list(mean = 5, sd = 0.5),  
               xlim = c(3.5, 4.5),  
               geom = "area", fill = "red") +  
  labs(x = "", y = "") + theme_classic() +  
  scale_y_continuous(breaks = NULL)
```



```
pnorm(4.5, mean = 5, sd = 0.5)
```

```
## [1] 0.1586553
```