

# Data Science Inference and Modeling

The textbook for the Data Science course series is [freely available online](#).

This course corresponds to the textbook chapters [Statistical Inference](#) and [Statistical Models](#).

## Learning Objectives

- The concepts necessary to define estimates and margins of errors of populations, parameters, estimates, and standard errors in order to make predictions about data
- How to use models to aggregate data from different sources
- The very basics of Bayesian statistics and predictive modeling

## Course Overview

### Section 1: Parameters and Estimates

You will learn how to estimate population parameters.

### Section 2: The Central Limit Theorem in Practice

You will apply the central limit theorem to assess how close a sample estimate is to the population parameter of interest.

### Section 3: Confidence Intervals and p-Values

You will learn how to calculate confidence intervals and learn about the relationship between confidence intervals and p-values.

### Section 4: Statistical Models

You will learn about statistical models in the context of election forecasting.

### Section 5: Bayesian Statistics

You will learn about Bayesian statistics through looking at examples from rare disease diagnosis and baseball.

### Section 6: Election Forecasting

You will learn about election forecasting, building on what you've learned in the previous sections about statistical modeling and Bayesian statistics.

## Section 7: Association Tests

You will learn how to use association and chi-squared tests to perform inference for binary, categorical, and ordinal data through an example looking at research funding rates.

## Introduction to Inference

The textbook for this section is available [here](#)

In this course, we will learn:

- *statistical inference*, the process of deducing characteristics of a population using data from a random sample
- the statistical concepts necessary to define *estimates* and *margins of errors*
- how to *forecast future results* and estimate the precision of our forecast
- how to calculate and interpret *confidence intervals* and *p-values*

### Key points

- Information gathered from a small random sample can be used to infer characteristics of the entire population.
- Opinion polls are useful when asking everyone in the population is impossible.
- A common use for opinion polls is determining voter preferences in political elections for the purposes of forecasting election results.
- The *spread* of a poll is the estimated difference between support two candidates or options.

## Section 1 Overview

Section 1 introduces you to parameters and estimates.

After completing Section 1, you will be able to:

- Understand how to use a sampling model to perform a poll.
- Explain the terms **population**, **parameter**, and **sample** as they relate to statistical inference.
- Use a sample to estimate the population proportion from the sample average.
- Calculate the expected value and standard error of the sample average.

## Sampling Model Parameters and Estimates

The textbook for this section is available [here](#) and [here; first part](#)

### Key points

- The task of statistical inference is to estimate an unknown population parameter using observed data from a sample.
- In a sampling model, the collection of elements in the urn is called the *population*.
- A *parameter* is a number that summarizes data for an entire population.
- A *sample* is observed data from a subset of the population.
- An *estimate* is a summary of the observed data about a parameter that we believe is informative. It is a data-driven guess of the population parameter.
- We want to predict the proportion of the blue beads in the urn, the parameter  $p$ . The proportion of red beads in the urn is  $1 - p$  and the *spread* is  $2p - 1$ .

- The sample proportion is a random variable. Sampling gives random results drawn from the population distribution.

*Code: Function for taking a random draw from a specific urn*

The **dslabs** package includes a function for taking a random draw of size  $n$  from the urn:

```
if(!require(tidyverse)) install.packages("tidyverse")
```

```
## Loading required package: tidyverse
```

```
## -- Attaching packages -----
```

```
## v ggplot2 3.3.2      v purrr   0.3.4
## v tibble  3.0.3      v dplyr   1.0.1
## v tidyr   1.1.1      v stringr 1.4.0
## v readr   1.3.1      v forcats 0.5.0
```

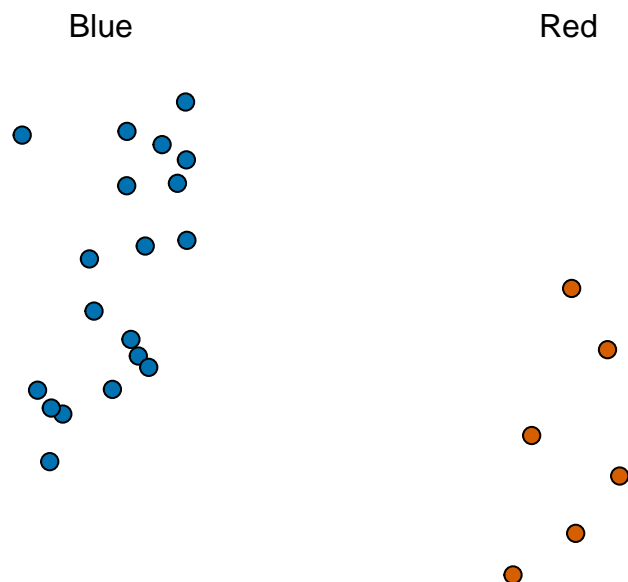
```
## -- Conflicts -----
```

```
## x dplyr::filter() masks stats::filter()
## x dplyr::lag()    masks stats::lag()
```

```
if(!require(dslabs)) install.packages("dslabs")
```

```
## Loading required package: dslabs
```

```
library(tidyverse)
library(dslabs)
take_poll(25)    # draw 25 beads
```



## The Sample Average

The textbook for this section is available [here](#) and [here](#)

### Key points

- Many common data science tasks can be framed as estimating a parameter from a sample.
- We illustrate statistical inference by walking through the process to estimate  $p$ . From the estimate of  $p$ , we can easily calculate an estimate of the spread,  $2p - 1$ .
- Consider the random variable  $X$  that is 1 if a blue bead is chosen and 0 if a red bead is chosen. The proportion of blue beads in  $N$  draws is the average of the draws  $X_1, \dots, X_N$ .
- $\bar{X}$  is the *sample average*. In statistics, a bar on top of a symbol denotes the average.  $\bar{X}$  is a random variable because it is the average of random draws - each time we take a sample,  $\bar{X}$  is different.

$$\bar{X} = \frac{X_1 + X_2 + \dots + X_N}{N}$$

- The number of blue beads drawn in  $N$  draws,  $N\bar{X}$ , is  $N$  times the proportion of values in the urn. However, we do not know the true proportion: we are trying to estimate this parameter  $p$ .

## Polling versus Forecasting

The textbook for this section is available [here](#)

### Key points

- A poll taken in advance of an election estimates  $p$  for that moment, not for election day.
- In order to predict election results, forecasters try to use early estimates of  $p$  to predict  $p$  on election day. We discuss some approaches in later sections.

## Properties of Our Estimate

The textbook for this section is available [here](#)

### Key points

- When interpreting values of  $\bar{X}$ , it is important to remember that  $\bar{X}$  is a random variable with an expected value and standard error that represents the sample proportion of positive events.
- The expected value of  $\bar{X}$  is the parameter of interest  $p$ . This follows from the fact that  $\bar{X}$  is the sum of independent draws of a random variable times a constant  $1/N$ .

$$E(\bar{X}) = p$$

- As the number of draws  $N$  increases, the standard error of our estimate  $\bar{X}$  decreases. The standard error of the average of  $\bar{X}$  over  $N$  draws is:

$$SE(\bar{X}) = \sqrt{p(1-p)/N}$$

- In theory, we can get more accurate estimates of  $p$  by increasing  $N$ . In practice, there are limits on the size of  $N$  due to costs, as well as other factors we discuss later.
- We can also use other random variable equations to determine the expected value of the sum of draws  $E(S)$  and standard error of the sum of draws  $SE(S)$ .

$$E(S) = Np$$

$$SE(S) = \sqrt{Np(1-p)}$$

## Assessment - Parameters and Estimates

1. Suppose you poll a population in which a proportion  $p$  of voters are Democrats and  $1 - p$  are Republicans.

Your sample size is  $N = 25$ . Consider the random variable  $S$ , which is the **total** number of Democrats in your sample.

What is the expected value of this random variable  $S$ ?

- ☐ A.  $E(S) = 25(1 - p)$
- ☒ B.  $E(S) = 25p$
- ☐ C.  $E(S) = \sqrt{25p(1 - p)}$
- ☐ D.  $E(S) = p$

2. Again, consider the random variable  $S$ , which is the **total** number of Democrats in your sample of 25 voters.

The variable  $p$  describes the proportion of Democrats in the sample, whereas  $1 - p$  describes the proportion of Republicans.

What is the standard error of  $S$ ?

- ☐ A.  $SE(S) = 25p(1 - p)$
- ☐ B.  $SE(S) = \sqrt{25p}$
- ☐ C.  $SE(S) = 25(1 - p)$
- ☒ D.  $SE(S) = \sqrt{25p(1 - p)}$

3. Consider the random variable  $S/N$ , which is equivalent to the sample average that we have been denoting as  $\bar{X}$ .

The variable  $N$  represents the sample size and  $p$  is the proportion of Democrats in the population.

What is the expected value of  $\bar{X}$ ?

- ☒ A.  $E(\bar{X}) = p$
- ☐ B.  $E(\bar{X}) = Np$
- ☐ C.  $E(\bar{X}) = N(1 - p)$
- ☐ D.  $E(\bar{X}) = 1 - p$

4. What is the standard error of the sample average,  $\bar{X}$ ?

The variable  $N$  represents the sample size and  $p$  is the proportion of Democrats in the population.

- ☐ A.  $SE(\bar{X}) = \sqrt{Np(1 - p)}$
- ☒ B.  $SE(\bar{X}) = \sqrt{p(1 - p)/N}$
- ☐ C.  $SE(\bar{X}) = \sqrt{p(1 - p)}$
- ☐ D.  $SE(\bar{X}) = \sqrt{N}$

5. Write a line of code that calculates the standard error **se** of a sample average when you poll 25 people in the population.

Generate a sequence of 100 proportions of Democrats **p** that vary from 0 (no Democrats) to 1 (all Democrats).

Plot **se** versus **p** for the 100 different proportions.

```

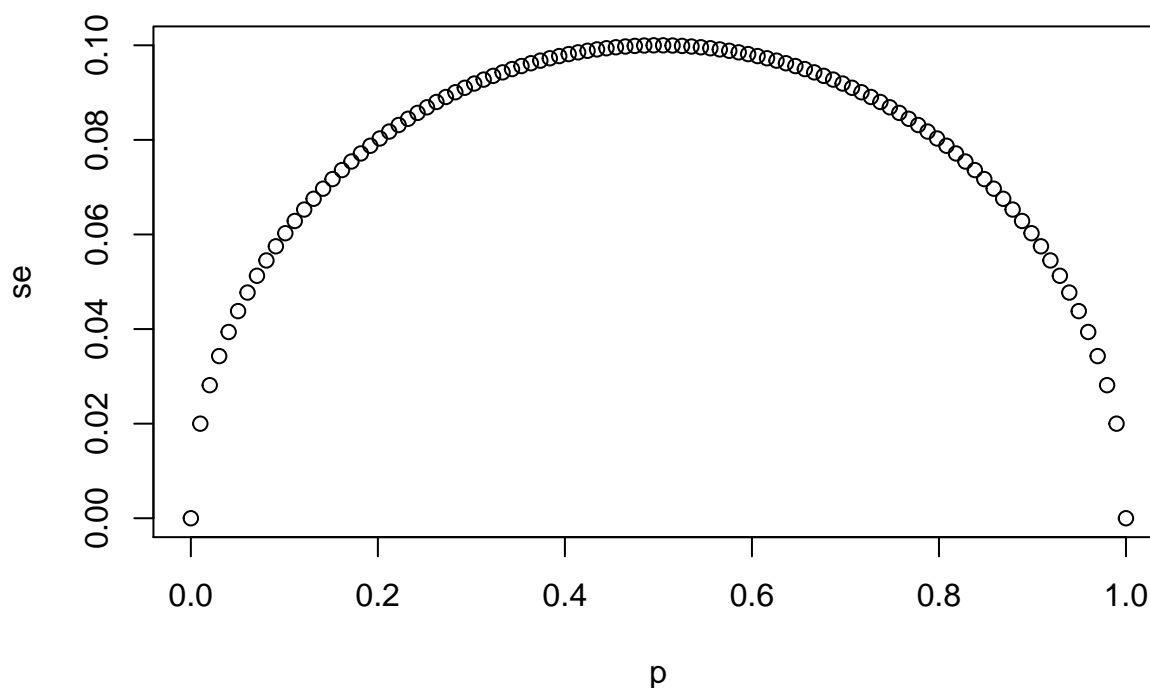
# `N` represents the number of people polled
N <- 25

# Create a variable `p` that contains 100 proportions ranging from 0 to 1 using the `seq` function
p <- seq(0, 1, length.out = 100)

# Create a variable `se` that contains the standard error of each sample average
se <- sqrt(p * (1 - p)/N)

# Plot `p` on the x-axis and `se` on the y-axis
plot(p, se)

```



6. Using the same code as in the previous exercise, create a for-loop that generates three plots of  $p$  versus  $se$  when the sample sizes equal  $N = 25$ ,  $N = 100$ , and  $N = 1000$ .

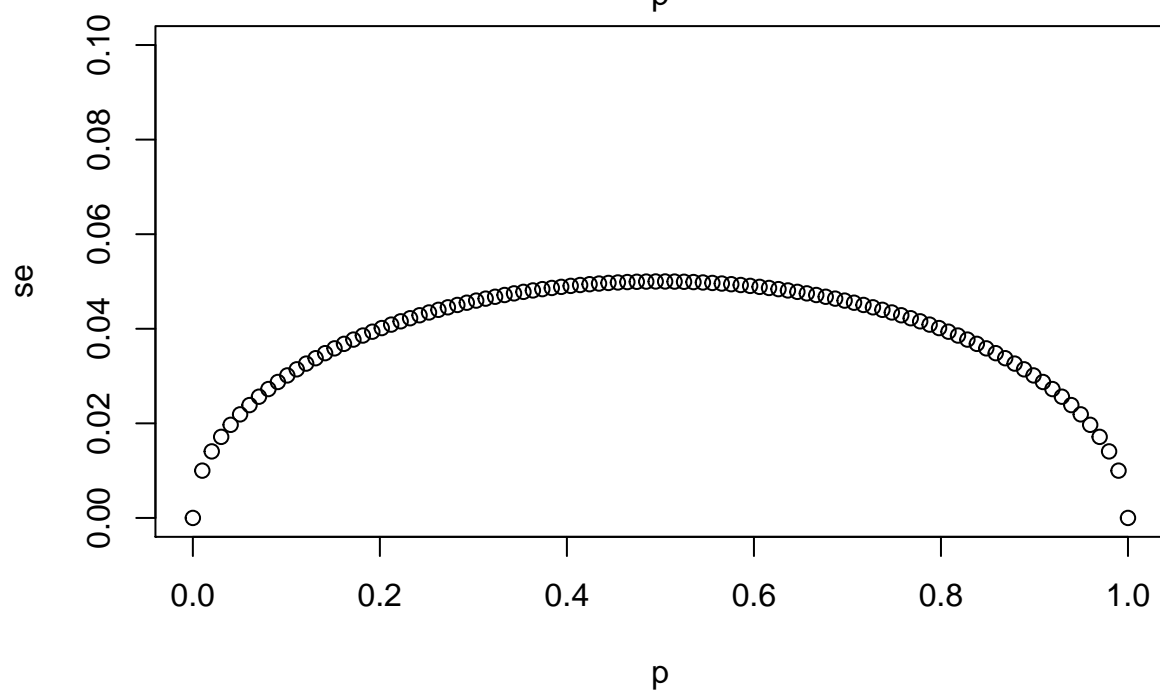
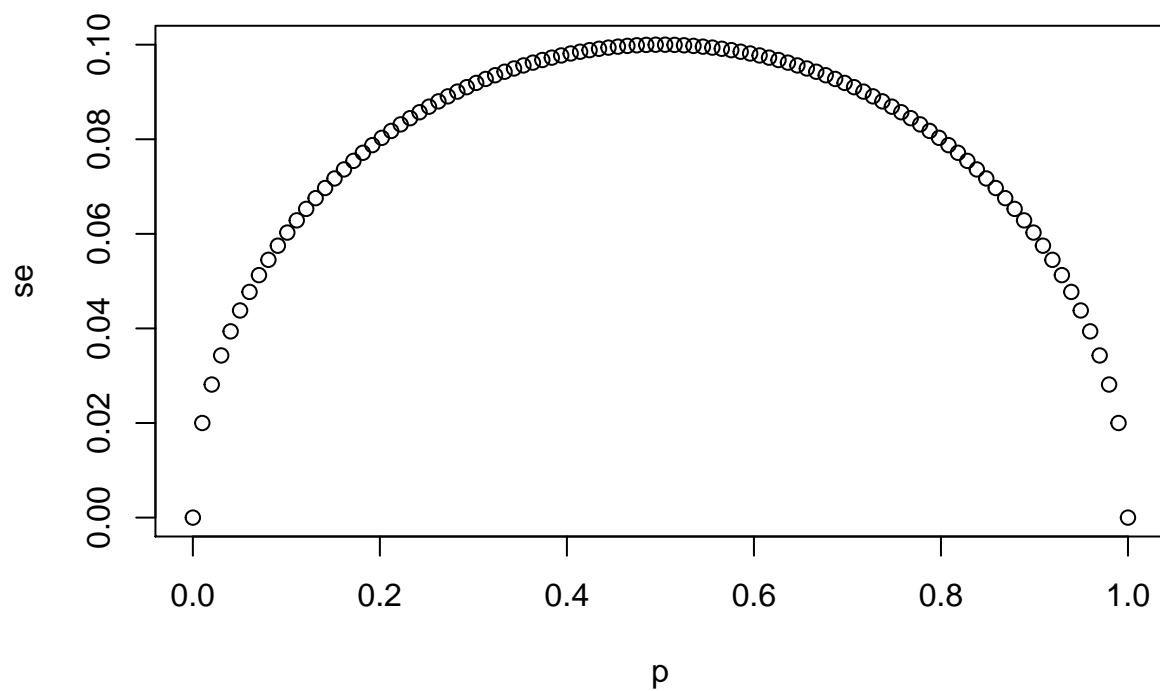
```

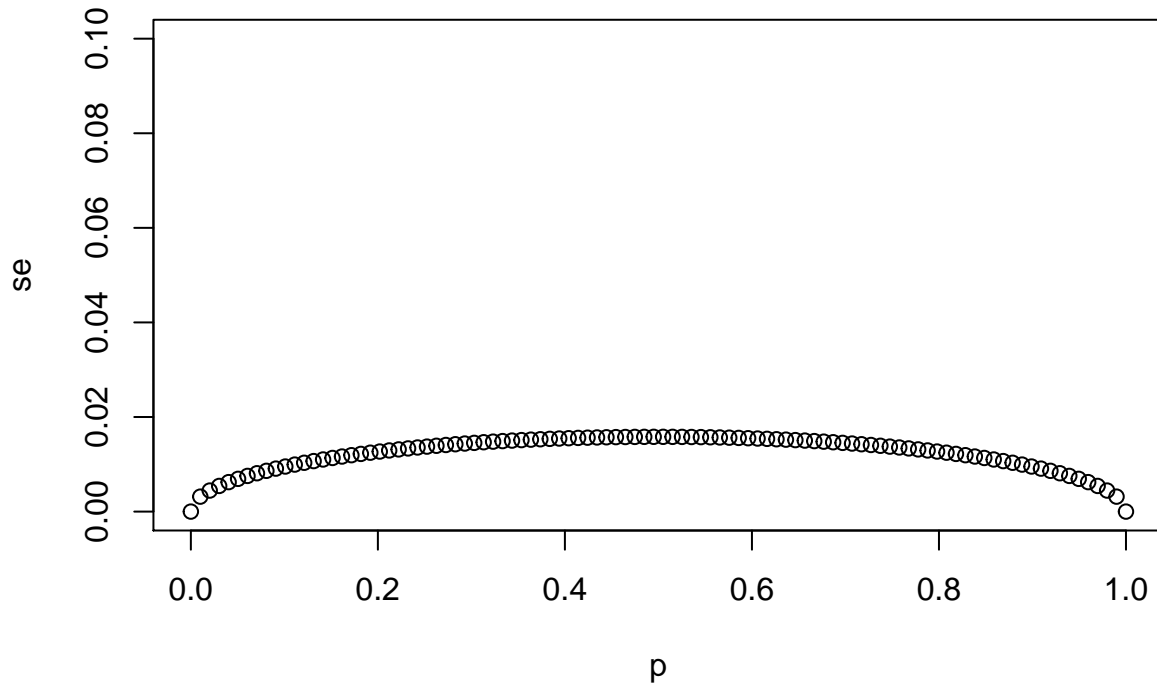
# The vector `p` contains 100 proportions of Democrats ranging from 0 to 1 using the `seq` function
p <- seq(0, 1, length = 100)

# The vector `sample_sizes` contains the three sample sizes
sample_sizes <- c(25, 100, 1000)

# Write a for-loop that calculates the standard error `se` for every value of `p` for each of the three
for (N in sample_sizes)
{
  se <- sqrt(p * (1 - p)/N)
  plot(p, se, ylim = c(0, 0.1))
}

```





7. Our estimate for the difference in proportions of Democrats and Republicans is  $d = \bar{X} - (1 - \bar{X})$ .

Which derivation correctly uses the rules we learned about sums of random variables and scaled random variables to derive the expected value of  $d$

- ☐ A.  $E[\bar{X} - (1 - \bar{X})] = E[2\bar{X} - 1] = 2E[\bar{X}] - 1 = N(2p - 1) = Np - N(1 - p)$
- ☐ B.  $E[\bar{X} - (1 - \bar{X})] = E[\bar{X} - 1] = E[\bar{X}] - 1 = p - 1$
- ☐ C.  $E[\bar{X} - (1 - \bar{X})] = E[2\bar{X} - 1] = 2E[\bar{X}] - 1 = 2\sqrt{p(1-p)} - 1 = p - (1 - p)$
- ☒ D.  $E[\bar{X} - (1 - \bar{X})] = E[2\bar{X} - 1] = 2E[\bar{X}] - 1 = 2p - 1 = p - (1 - p)$

8. Our estimate for the difference in proportions of Democrats and Republicans is  $d = \bar{X} - (1 - \bar{X})$ .

Which derivation correctly uses the rules we learned about sums of random variables and scaled random variables to derive the standard error of  $d$ ?

- ☐ A.  $SE[\bar{X} - (1 - \bar{X})] = SE[2\bar{X} - 1] = 2SE[\bar{X}] = 2\sqrt{p/N}$
- ☐ B.  $SE[\bar{X} - (1 - \bar{X})] = SE[2\bar{X} - 1] = 2SE[\bar{X} - 1] = 2\sqrt{p(1-p)/N} - 1$
- ☒ C.  $SE[\bar{X} - (1 - \bar{X})] = SE[2\bar{X} - 1] = 2SE[\bar{X}] = 2\sqrt{p(1-p)/N}$
- ☐ D.  $SE[\bar{X} - (1 - \bar{X})] = SE[\bar{X} - 1] = SE[\bar{X}] = \sqrt{p(1-p)/N}$

9. Say the actual proportion of Democratic voters is  $p = 0.45$ .

In this case, the Republican party is winning by a relatively large margin of  $d = -0.1$ , or a 10% margin of victory. What is the standard error of the spread  $2\bar{X} - 1$  in this case?

```
# `N` represents the number of people polled
N <- 25
```



```
# `p` represents the proportion of Democratic voters
p <- 0.45

# Calculate the standard error of the spread. Print this value to the console.
2*sqrt((p*(1-p)/N))
```

```
## [1] 0.1989975
```

10. So far we have said that the difference between the proportion of Democratic voters and Republican voters is about 10% and that the standard error of this spread is about 0.2 when  $N = 25$ .

Select the statement that explains why this sample size is sufficient or not.

- ☐ A. This sample size is sufficient because the expected value of our estimate  $2\bar{X} - 1$  is  $d$  so our prediction will be right on.
- ☒ B. This sample size is too small because the standard error is larger than the spread.
- ☐ C. This sample size is sufficient because the standard error of about 0.2 is much smaller than the spread of 10%.
- ☐ D. Without knowing  $p$ , we have no way of knowing that increasing our sample size would actually improve our standard error.

## Section 2 Overview

In Section 2, you will look at the Central Limit Theorem in practice.

After completing Section 2, you will be able to:

- Use the Central Limit Theorem to calculate the probability that a sample estimate  $\bar{X}$  is close to the population proportion  $p$ .
- Run a Monte Carlo simulation to corroborate theoretical results built using probability theory.
- Estimate the spread based on estimates of  $\bar{X}$  and  $\hat{SE}(\bar{X})$ .
- Understand why bias can mean that larger sample sizes aren't necessarily better.

## The Central Limit Theorem in Practice

The textbook for this section is available [here](#)

### Key points

- Because  $\bar{X}$  is the sum of random draws divided by a constant, the distribution of  $\bar{X}$  is approximately normal.
- We can convert  $\bar{X}$  to a standard normal random variable  $Z$ :

$$Z = \frac{\bar{X} - E(\bar{X})}{SE(\bar{X})}$$

- The probability that  $\bar{X}$  is within .01 of the actual value of  $p$  is:

$$Pr(Z \leq .01/\sqrt{p(1-p)/N}) - Pr(Z \leq -.01/\sqrt{p(1-p)/N})$$

- The Central Limit Theorem (CLT) still works if  $\bar{X}$  is used in place of  $p$ . This is called a *plug-in estimate*. Hats over values denote estimates. Therefore:

$$\hat{SE}(\bar{X}) = \sqrt{\bar{X}(1 - \bar{X})/N}$$

Using the CLT, the probability that  $\bar{X}$  is within .01 of the actual value of  $p$  is:

$$Pr(Z \leq .01/\sqrt{\bar{X}(1 - \bar{X})/N}) - Pr(Z \leq -.01/\sqrt{\bar{X}(1 - \bar{X})/N})$$

Code: Computing the probability of  $\bar{X}$  being within .01 of  $p$

```
X_hat <- 0.48
se <- sqrt(X_hat*(1-X_hat)/25)
pnorm(0.01/se) - pnorm(-0.01/se)
```

```
## [1] 0.07971926
```

## Margin of Error

The textbook for this section is available [here](#)

### Key points

- The *margin of error* is defined as 2 times the standard error of the estimate  $\bar{X}$ .
- There is about a 95% chance that  $\bar{X}$  will be within two standard errors of the actual parameter  $p$ .

## A Monte Carlo Simulation for the CLT

The textbook for this section is available [here](#)

### Key points

- We can run Monte Carlo simulations to compare with theoretical results assuming a value of  $p$ .
- In practice,  $p$  is unknown. We can corroborate theoretical results by running Monte Carlo simulations with one or several values of  $p$ .
- One practical choice for  $p$  when modeling is  $\bar{X}$ , the observed value of  $\hat{X}$  in a sample.

Code: Monte Carlo simulation using a set value of  $p$

```
p <- 0.45      # unknown p to estimate
N <- 1000

# simulate one poll of size N and determine x_hat
x <- sample(c(0,1), size = N, replace = TRUE, prob = c(1-p, p))
x_hat <- mean(x)

# simulate B polls of size N and determine average x_hat
B <- 10000     # number of replicates
N <- 1000      # sample size per replicate
x_hat <- replicate(B, {
  x <- sample(c(0,1), size = N, replace = TRUE, prob = c(1-p, p))
  mean(x)
})
```

Code: Histogram and QQ-plot of Monte Carlo results

```
if(!require(gridExtra)) install.packages("gridExtra")
```

```
## Loading required package: gridExtra
```

```
##
```

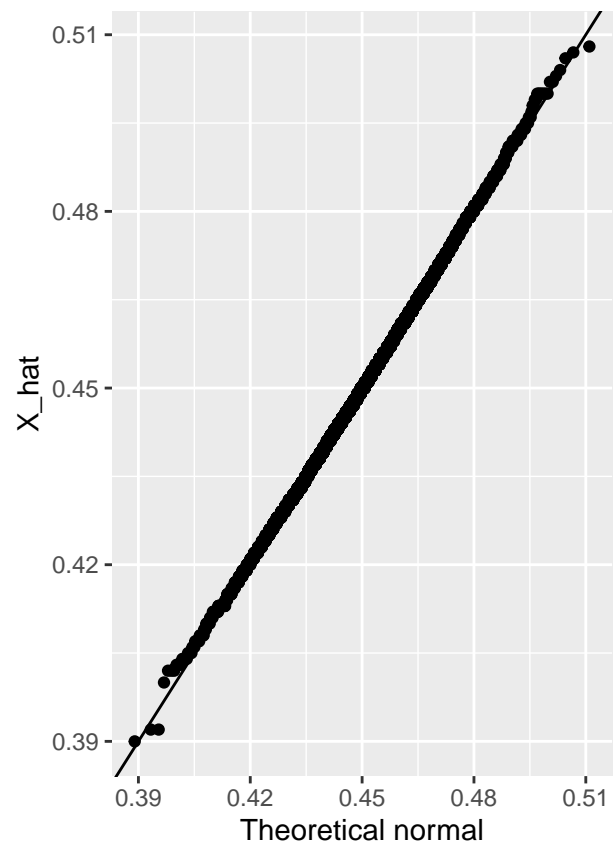
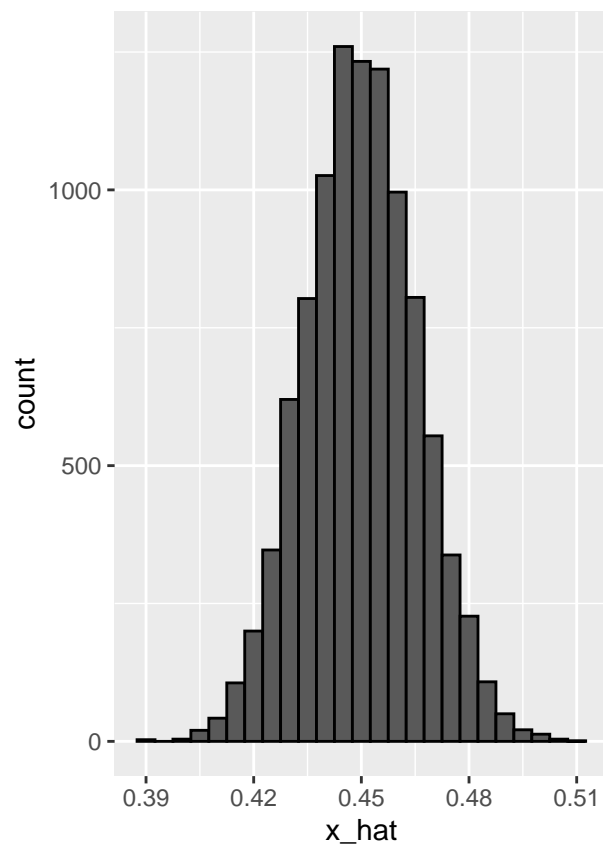
```
## Attaching package: 'gridExtra'
```

```
## The following object is masked from 'package:dplyr':
```

```
##
```

```
##      combine
```

```
library(gridExtra)
p1 <- data.frame(x_hat = x_hat) %>%
  ggplot(aes(x_hat)) +
  geom_histogram(binwidth = 0.005, color = "black")
p2 <- data.frame(x_hat = x_hat) %>%
  ggplot(aes(sample = x_hat)) +
  stat_qq(dparams = list(mean = mean(x_hat), sd = sd(x_hat))) +
  geom_abline() +
  ylab("X_hat") +
  xlab("Theoretical normal")
grid.arrange(p1, p2, nrow=1)
```



## The Spread

The textbook for this section is available [here](#)

### Key points

- The spread between two outcomes with probabilities  $p$  and  $1 - p$  is  $2p - 1$ .
- The expected value of the spread is  $2\bar{X} - 1$ .
- The standard error of the spread is  $2\hat{SE}(\bar{X})$ .
- The margin of error of the spread is 2 times the margin of error of  $\bar{X}$ .

## Bias: Why Not Run a Very Large Poll?

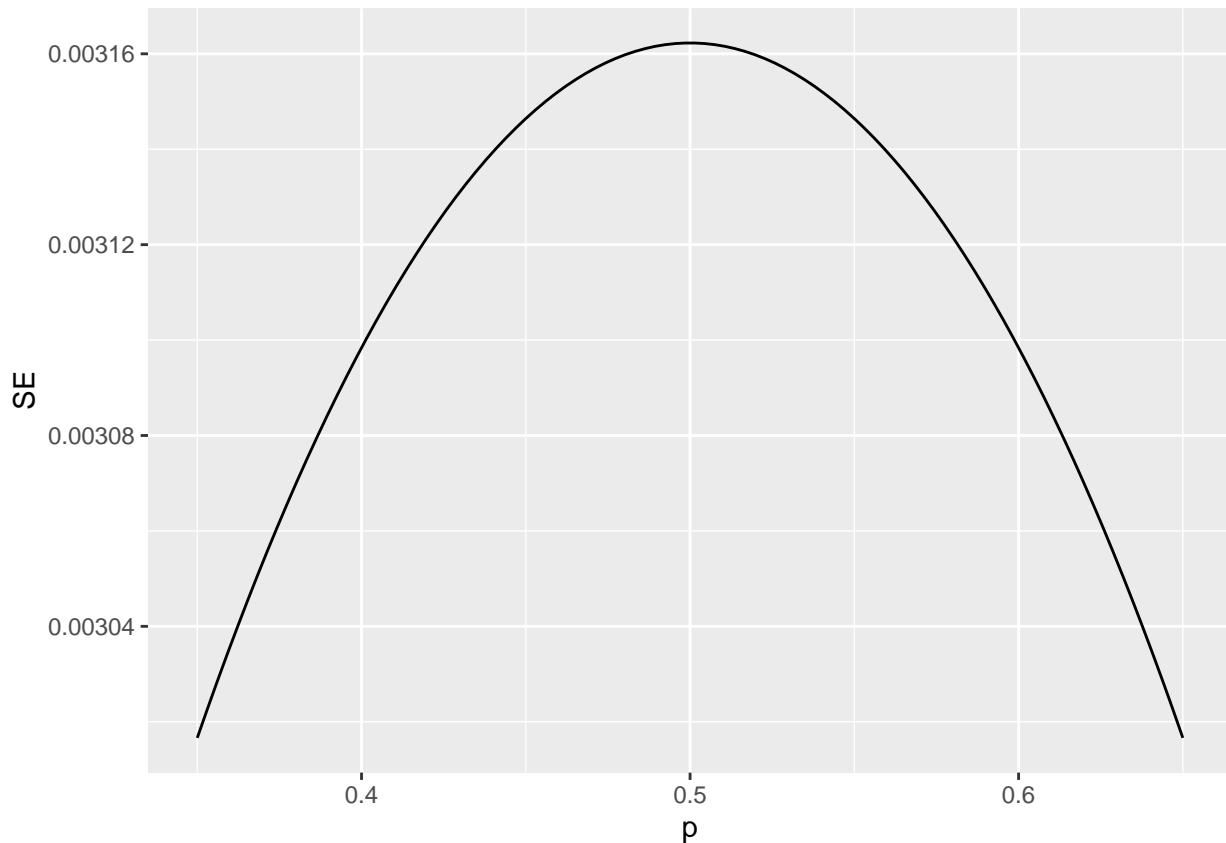
The textbook for this section is available [here](#)

### Key points

- An extremely large poll would theoretically be able to predict election results almost perfectly.
- These sample sizes are not practical. In addition to cost concerns, polling doesn't reach everyone in the population (eventual voters) with equal probability, and it also may include data from outside our population (people who will not end up voting).
- These systematic errors in polling are called *bias*. We will learn more about bias in the future.

*Code: Plotting margin of error in an extremely large poll over a range of values of  $p$*

```
N <- 100000
p <- seq(0.35, 0.65, length = 100)
SE <- sapply(p, function(x) 2*sqrt(x*(1-x)/N))
data.frame(p = p, SE = SE) %>%
  ggplot(aes(p, SE)) +
  geom_line()
```



## Assessment - Introduction to Inference

1. Write function called `take_sample` that takes the proportion of Democrats  $p$  and the sample size  $N$  as arguments and returns the sample average of Democrats (1) and Republicans (0).

Calculate the sample average if the proportion of Democrats equals 0.45 and the sample size is 100.

```
# Write a function called `take_sample` that takes `p` and `N` as arguments and returns the average value
take_sample <- function(p, N) {
  x <- sample(c(0,1), size = N, replace = TRUE, prob = c(1-p, p))
  return(mean(x))
}

# Use the `set.seed` function to make sure your answer matches the expected result after random sampling.
set.seed(1)

# Define `p` as the proportion of Democrats in the population being polled
p <- 0.45

# Define `N` as the number of people polled
N <- 100

# Call the `take_sample` function to determine the sample average of `N` randomly selected people from the population
take_sample(p, N)

## [1] 0.46
```

2. Assume the proportion of Democrats in the population  $p$  equals 0.45 and that your sample size  $N$  is 100 polled voters.

The `take_sample` function you defined previously generates our estimate,  $\bar{X}$ .

Replicate the random sampling 10,000 times and calculate  $p - \bar{X}$  for each random sample. Save these differences as a vector called `errors`. Find the average of `errors` and plot a histogram of the distribution.

```
# Define `p` as the proportion of Democrats in the population being polled
p <- 0.45

# Define `N` as the number of people polled
N <- 100

# The variable `B` specifies the number of times we want the sample to be replicated
B <- 10000

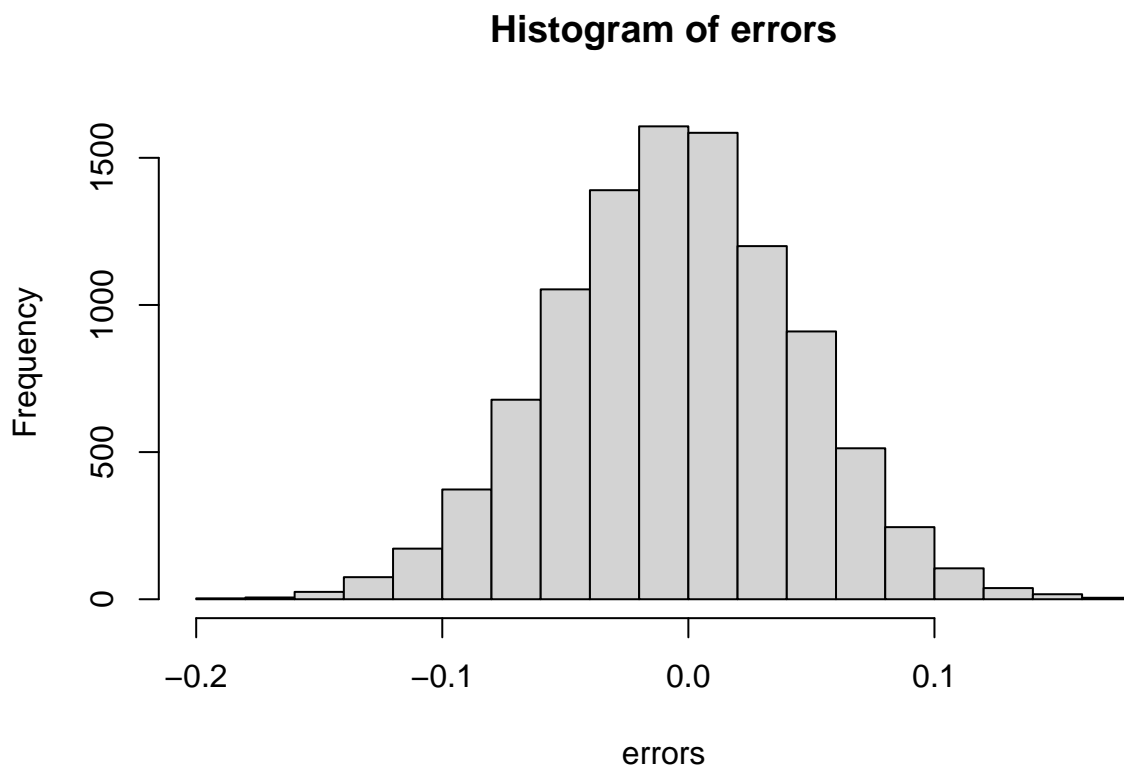
# Use the `set.seed` function to make sure your answer matches the expected result after random sampling
set.seed(1)

# Create an objected called `errors` that replicates subtracting the result of the `take_sample` function
errors <- replicate(B, p - take_sample(p, N))

# Calculate the mean of the errors. Print this value to the console.
mean(errors)
```

```
## [1] -4.9e-05
```

```
hist(errors)
```



3. In the last exercise, you made a vector of differences between the actual value for  $p$  and an estimate,  $\bar{X}$ .

We called these differences between the actual and estimated values **errors**.

The **errors** object has already been loaded for you. Use the **hist** function to plot a histogram of the values contained in the vector **errors**. Which statement best describes the distribution of the errors?

- ☐ A. The errors are all about 0.05.
- ☐ B. The errors are all about -0.05.
- ☒ C. The errors are symmetrically distributed around 0.
- ☐ D. The errors range from -1 to 1.

4. The error  $p - \bar{X}$  is a random variable.

In practice, the error is not observed because we do not know the actual proportion of Democratic voters,  $p$ . However, we can describe the size of the error by constructing a simulation.

What is the average size of the error if we define the size by taking the absolute value  $|p - \bar{X}|$ ?

```
# Define `p` as the proportion of Democrats in the population being polled
p <- 0.45

# Define `N` as the number of people polled
N <- 100

# The variable `B` specifies the number of times we want the sample to be replicated
B <- 10000

# Use the `set.seed` function to make sure your answer matches the expected result after random sampling
set.seed(1)

# We generated `errors` by subtracting the estimate from the actual proportion of Democratic voters
errors <- replicate(B, p - take_sample(p, N))

# Calculate the mean of the absolute value of each simulated error. Print this value to the console.
mean(abs(errors))
```

```
## [1] 0.039267
```

5. The standard error is related to the typical **size** of the error we make when predicting.

We say **size** because, as we just saw, the errors are centered around 0. In that sense, the typical error is 0. For mathematical reasons related to the central limit theorem, we actually use the standard deviation of **errors** rather than the average of the absolute values.

As we have discussed, the standard error is the square root of the average squared distance  $(\bar{X} - p)^2$ . The standard deviation is defined as the square root of the distance squared.

Calculate the standard deviation of the spread.

```

# Define `p` as the proportion of Democrats in the population being polled
p <- 0.45

# Define `N` as the number of people polled
N <- 100

# The variable `B` specifies the number of times we want the sample to be replicated
B <- 10000

# Use the `set.seed` function to make sure your answer matches the expected result after random sampling
set.seed(1)

# We generated `errors` by subtracting the estimate from the actual proportion of Democratic voters
errors <- replicate(B, p - take_sample(p, N))

# Calculate the standard deviation of `errors`
sqrt(mean(errors^2))

## [1] 0.04949939

```

6. The theory we just learned tells us what this standard deviation is going to be because it is the standard error of  $\bar{X}$ .

Estimate the standard error given an expected value of 0.45 and a sample size of 100.

```

# Define `p` as the expected value equal to 0.45
p <- 0.45

# Define `N` as the sample size
N <- 100

# Calculate the standard error
sqrt(p*(1-p)/N)

## [1] 0.04974937

```

7. In practice, we don't know  $p$ , so we construct an estimate of the theoretical prediction based by plugging in  $\bar{X}$  for  $p$ . Calculate the standard error of the estimate:  $\hat{SE}(\bar{X})$

```

# Define `p` as a proportion of Democratic voters to simulate
p <- 0.45

# Define `N` as the sample size
N <- 100

# Use the `set.seed` function to make sure your answer matches the expected result after random sampling
set.seed(1)

# Define `X` as a random sample of `N` voters with a probability of picking a Democrat ('1') equal to `p`
X <- sample(c(0,1), size = N, replace = TRUE, prob = c(1-p, p))

```



```
# Define `X_bar` as the average sampled proportion
X_bar <- mean(X)

# Calculate the standard error of the estimate. Print the result to the console.
se <- sqrt((X_bar*(1-X_bar)/N))
se

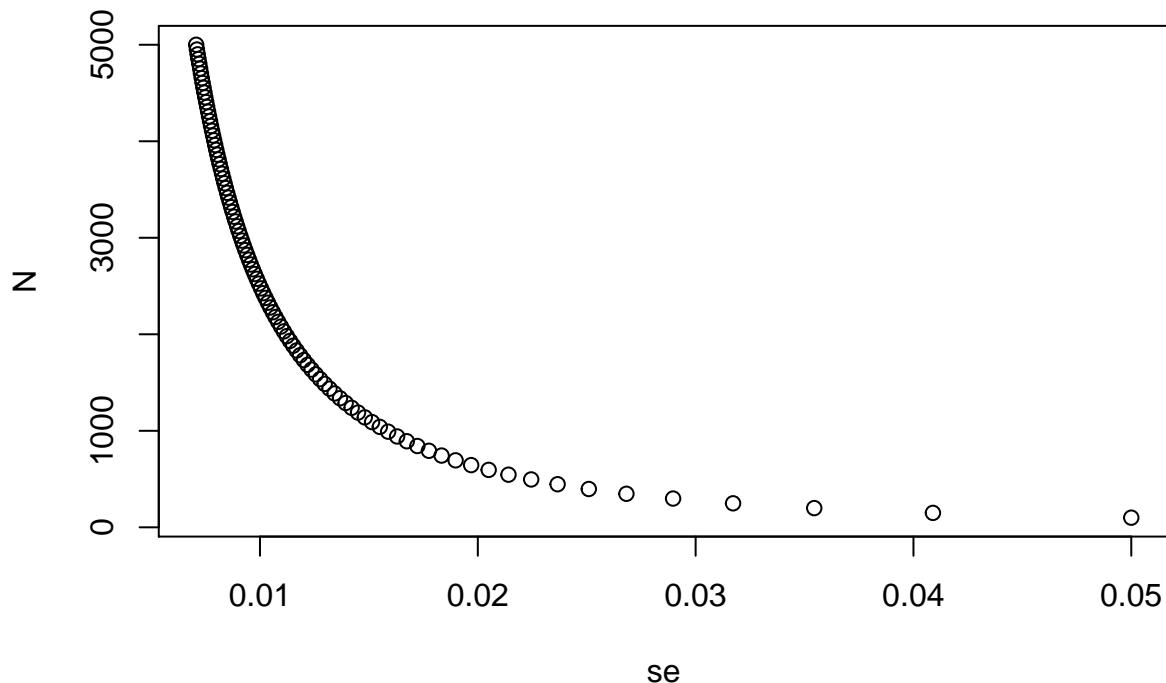
## [1] 0.04983974
```

8. The standard error estimates obtained from the Monte Carlo simulation, the theoretical prediction, and the estimate of the theoretical prediction are all very close, which tells us that the theory is working.

This gives us a practical approach to knowing the typical error we will make if we predict  $p$  with  $\hat{X}$ . The theoretical result gives us an idea of how large a sample size is required to obtain the precision we need. Earlier we learned that the largest standard errors occur for  $p = 0.5$ .

Create a plot of the largest standard error for  $N$  ranging from 100 to 5,000.

```
N <- seq(100, 5000, len = 100)
p <- 0.5
se <- sqrt(p*(1-p)/N)
plot(se, N)
```



Based on this plot, how large does the sample size have to be to have a standard error of about 1%?

- ☐ A. 100
- ☐ B. 500
- ☒ C. 2,500
- ☐ D. 4,000

9. For  $N = 100$ , the central limit theorem tells us that the distribution of  $\hat{X}$  is...

- ☐ A. practically equal to  $p$ .
- ☒ B. approximately normal with expected value  $p$  and standard error  $\sqrt{p(1-p)/N}$ .
- ☐ C. approximately normal with expected value  $\bar{X}$  and standard error  $\sqrt{\bar{X}(1-\bar{X})/N}$ .
- ☐ D. not a random variable.

10. We calculated a vector `errors` that contained, for each simulated sample, the difference between the actual value  $p$  and our estimate  $\hat{X}$ .

The errors  $\bar{X} - p$  are:

- ☐ A. practically equal to 0.
- ☒ B. approximately normal with expected value 0 and standard error  $\sqrt{p(1-p)/N}$ .
- ☐ C. approximately normal with expected value  $p$  and standard error  $\sqrt{p(1-p)/N}$ .
- ☐ D. not a random variable.

11. Make a qq-plot of the `errors` you generated previously to see if they follow a normal distribution.

```
# Define `p` as the proportion of Democrats in the population being polled
p <- 0.45

# Define `N` as the number of people polled
N <- 100

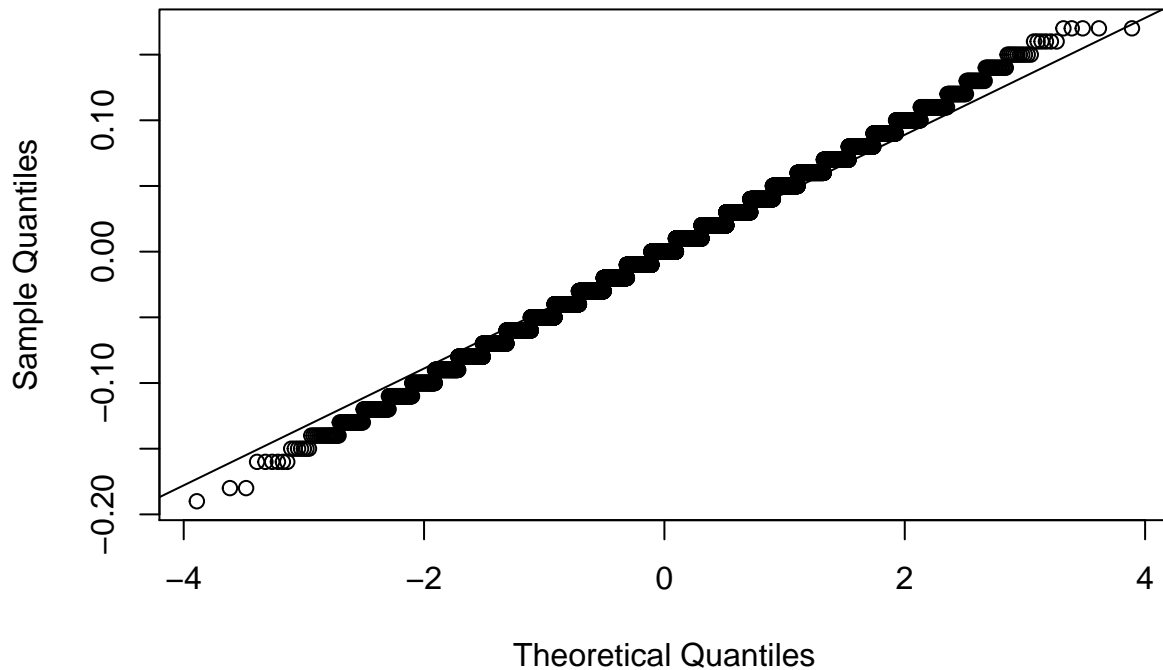
# The variable `B` specifies the number of times we want the sample to be replicated
B <- 10000

# Use the `set.seed` function to make sure your answer matches the expected result after random sampling
set.seed(1)

# Generate `errors` by subtracting the estimate from the actual proportion of Democratic voters
errors <- replicate(B, p - take_sample(p, N))

# Generate a qq-plot of `errors` with a qq-line showing a normal distribution
qqnorm(errors)
qqline(errors)
```

## Normal Q-Q Plot



12. If  $p = 0.45$  and  $N = 100$ , use the central limit theorem to estimate the probability that  $\bar{X} > 0.5$ .

```
# Define `p` as the proportion of Democrats in the population being polled
p <- 0.45

# Define `N` as the number of people polled
N <- 100

# Calculate the probability that the estimated proportion of Democrats in the population is greater than 0.5
1-pnorm(0.5, p, sqrt(p*(1-p)/N))

## [1] 0.1574393
```

13. Assume you are in a practical situation and you don't know  $p$ .

Take a sample of size  $N = 100$  and obtain a sample average of  $\bar{X} = 0.51$ .

What is the CLT approximation for the probability that your error size is equal or larger than 0.01?

```
# Define `N` as the number of people polled
N <- 100

# Define `X_hat` as the sample average
X_hat <- 0.51

# Define `se_hat` as the standard error of the sample average
se_hat <- sqrt(X_hat*(1-X_hat)/N)
```

```
# Calculate the probability that the error is 0.01 or larger
1-pnorm(0.01,0,se_hat) + pnorm(-0.01,0,se_hat)
```

```
## [1] 0.8414493
```

## Section 3 Overview

In Section 3, you will look at confidence intervals and p-values.

After completing Section 3, you will be able to:

- Calculate confidence intervals of difference sizes around an estimate.
- Understand that a confidence interval is a random interval with the given probability of falling on top of the parameter.
- Explain the concept of “power” as it relates to inference.
- Understand the relationship between p-values and confidence intervals and explain why reporting confidence intervals is often preferable.

## Confidence Intervals

The textbook for this section is available [here](#)

### Key points

- We can use statistical theory to compute the probability that a given interval contains the true parameter  $p$ .
- 95% confidence intervals are intervals constructed to have a 95% chance of including  $p$ . The margin of error is approximately a 95% confidence interval.
- The start and end of these confidence intervals are random variables.
- To calculate any size confidence interval, we need to calculate the value  $z$  for which  $Pr(-z \leq Z \leq z)$  equals the desired confidence. For example, a 99% confidence interval requires calculating  $z$  for  $Pr(-z \leq Z \leq z) = 0.99$ .
- For a confidence interval of size  $q$ , we solve for  $z = 1 - \frac{1-q}{2}$ .
- To determine a 95% confidence interval, use `z <- qnorm(0.975)`. This value is slightly smaller than 2 times the standard error.

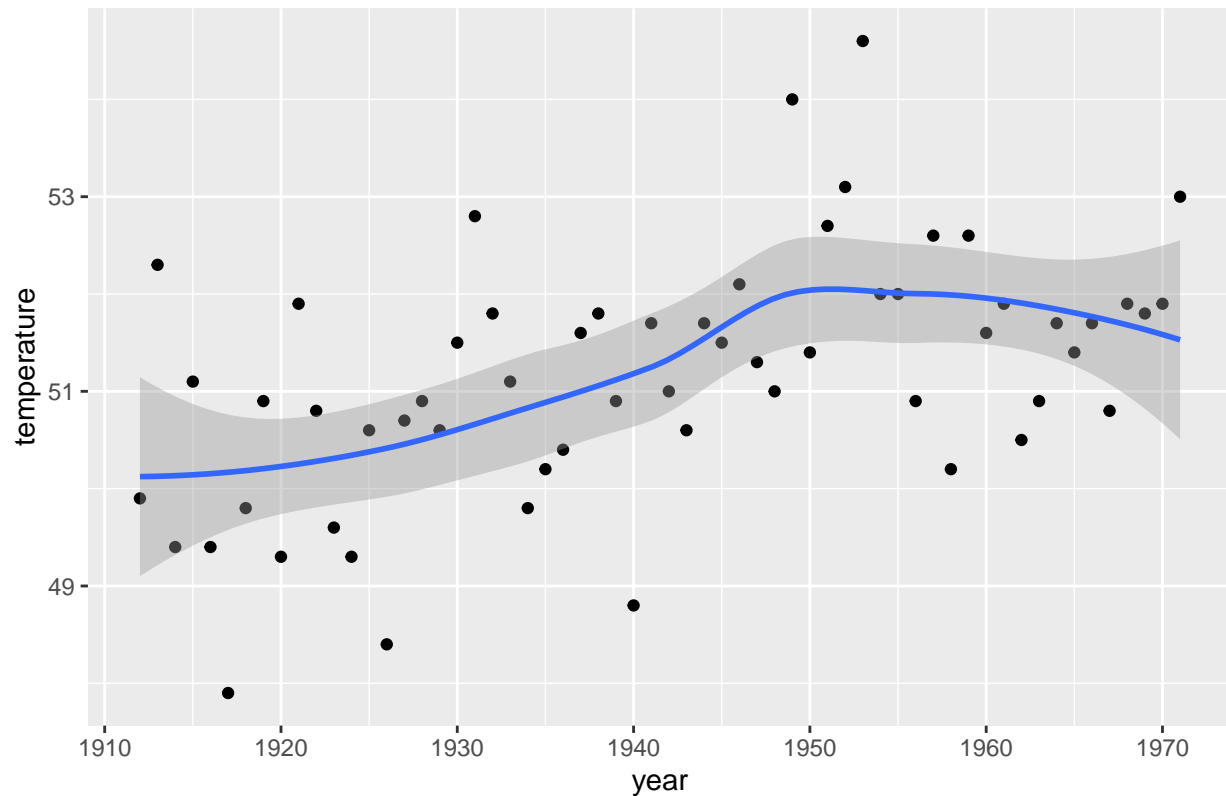
*Code: geom\_smooth confidence interval example*

The shaded area around the curve is related to the concept of confidence intervals.

```
data("nhtemp")
data.frame(year = as.numeric(time(nhtemp)), temperature = as.numeric(nhtemp)) %>%
  ggplot(aes(year, temperature)) +
  geom_point() +
  geom_smooth() +
  ggtitle("Average Yearly Temperatures in New Haven")
```

```
## `geom_smooth()` using method = 'loess' and formula 'y ~ x'
```

## Average Yearly Temperatures in New Haven



Code: Monte Carlo simulation of confidence intervals

Note that to compute the exact 95% confidence interval, we would use `qnorm(.975)*SE_hat` instead of `2*SE_hat`.

```
p <- 0.45
N <- 1000
X <- sample(c(0,1), size = N, replace = TRUE, prob = c(1-p, p)) # generate N observations
X_hat <- mean(X) # calculate X_hat
SE_hat <- sqrt(X_hat*(1-X_hat)/N) # calculate SE_hat, SE of the mean of N observations
c(X_hat - 2*SE_hat, X_hat + 2*SE_hat) # build interval of 2*SE above and below mean
```

```
## [1] 0.4135691 0.4764309
```

Code: Solving for  $z$  with `qnorm`

```
z <- qnorm(0.995) # calculate z to solve for 99% confidence interval
pnorm(qnorm(0.995)) # demonstrating that qnorm gives the z value for a given probability
```

```
## [1] 0.995
```

```
pnorm(qnorm(1-0.995)) # demonstrating symmetry of 1-qnorm
```

```
## [1] 0.005
```

```
pnorm(z) - pnorm(-z)    # demonstrating that this z value gives correct probability for interval
```

```
## [1] 0.99
```

## A Monte Carlo Simulation for Confidence Intervals

The textbook for this section is available [here](#)

### Key points

- We can run a Monte Carlo simulation to confirm that a 95% confidence interval contains the true value of  $p$  95% of the time.
- A plot of confidence intervals from this simulation demonstrates that most intervals include  $p$ , but roughly 5% of intervals miss the true value of  $p$ .

*Code: Monte Carlo simulation*

Note that to compute the exact 95% confidence interval, we would use `qnorm(.975)*SE_hat` instead of `2*SE_hat`.

```
B <- 10000
inside <- replicate(B, {
  X <- sample(c(0,1), size = N, replace = TRUE, prob = c(1-p, p))
  X_hat <- mean(X)
  SE_hat <- sqrt(X_hat*(1-X_hat)/N)
  between(p, X_hat - 2*SE_hat, X_hat + 2*SE_hat)    # TRUE if p in confidence interval
})
mean(inside)
```

```
## [1] 0.9566
```

## The Correct Language

The textbook for this section is available [here](#)

### Key points

- The 95% confidence intervals are random, but  $p$  is not random.
- 95% refers to the probability that the random interval falls on top of  $p$ .
- It is technically incorrect to state that  $p$  has a 95% chance of being in between two values because that implies  $p$  is random.

## Power

The textbook for this section is available [here](#)

### Key points

- If we are trying to predict the result of an election, then a confidence interval that includes a spread of 0 (a tie) is not helpful.

- A confidence interval that includes a spread of 0 does not imply a close election, it means the sample size is too small.
- Power is the probability of detecting an effect when there is a true effect to find. Power increases as sample size increases, because larger sample size means smaller standard error.

*Code: Confidence interval for the spread with sample size of 25*

Note that to compute the exact 95% confidence interval, we would use `c(-qnorm(.975), qnorm(.975))` instead of 1.96.

```
N <- 25
X_hat <- 0.48
(2*X_hat - 1) + c(-2, 2)*2*sqrt(X_hat*(1-X_hat)/N)
```

```
## [1] -0.4396799  0.3596799
```

## p-Values

The textbook for this section is available [here](#)

### Key points

- The null hypothesis is the hypothesis that there is no effect. In this case, the null hypothesis is that the spread is 0, or  $p = 0.5$ .
- The p-value is the probability of detecting an effect of a certain size or larger when the null hypothesis is true.
- We can convert the probability of seeing an observed value under the null hypothesis into a standard normal random variable. We compute the value of  $z$  that corresponds to the observed result, and then use that  $z$  to compute the p-value.
- If a 95% confidence interval does not include our observed value, then the p-value must be smaller than 0.05.
- It is preferable to report confidence intervals instead of p-values, as confidence intervals give information about the size of the estimate and p-values do not.

*Code: Computing a p-value for observed spread of 0.02*

```
N <- 100      # sample size
z <- sqrt(N) * 0.02/0.5      # spread of 0.02
1 - (pnorm(z) - pnorm(-z))
```

```
## [1] 0.6891565
```

## Another Explanation of p-Values

The p-value is the probability of observing a value as extreme or more extreme than the result given that the null hypothesis is true.

In the context of the normal distribution, this refers to the probability of observing a Z-score whose absolute value is as high or higher than the Z-score of interest.

Suppose we want to find the p-value of an observation 2 standard deviations larger than the mean. This means we are looking for anything with  $|z| \geq 2$ .

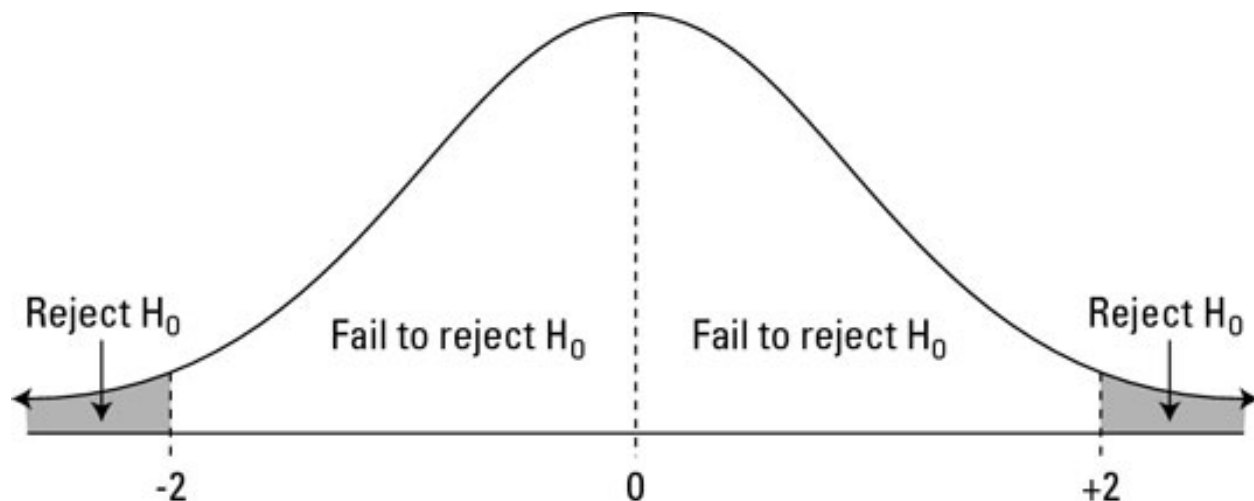


Figure 1: Standard normal distribution (centered at  $z=0$  with a standard deviation of 1)

Graphically, the p-value gives the probability of an observation that's at least as far away from the mean or further. This plot shows a standard normal distribution (centered at  $z = 0$ ) with a standard deviation of 1). The shaded tails are the region of the graph that are 2 standard deviations or more away from the mean.

The p-value is the proportion of area under a normal curve that has z-scores as extreme or more extreme than the given value - the tails on this plot of a normal distribution are shaded to show the region corresponding to the p-value.

The right tail can be found with `1-pnorm(2)`. We want to have both tails, though, because we want to find the probability of any observation as far away from the mean or farther, in either direction. (This is what's meant by a two-tailed p-value.) Because the distribution is symmetrical, the right and left tails are the same size and we know that our desired value is just `2*(1-pnorm(2))`.

Recall that, by default, `pnorm()` gives the CDF for a normal distribution with a mean of  $\mu = 0$  and standard deviation of  $\sigma = 1$ . To find p-values for a given z-score  $z$  in a normal distribution with mean `mu` and standard deviation `sigma`, use `2*(1-pnorm(z, mu, sigma))` instead.

### Assessment 3.1: Confidence Intervals and p-Values

#### 1. Confidence interval for p

For the following exercises, we will use actual poll data from the 2016 election. The exercises will contain pre-loaded data from the `dslabs` package.

```
library(dslabs)
library(dplyr)
library(ggplot2)
data(polls_us_election_2016)
```

We will use all the national polls that ended within a few weeks before the election.

Assume there are only two candidates and construct a 95% confidence interval for the election night proportion  $p$ .

Instructions - Use `filter` to subset the data set for the poll data you want. Include polls that ended on or after October 31, 2016 (`enddate`). Only include polls that took place in the United States. Call this filtered



object polls. - Use nrow to make sure you created a filtered object polls that contains the correct number of rows. - Extract the sample size N from the first poll in your subset object polls. - Convert the percentage of Clinton voters (rawpoll\_clinton) from the first poll in polls to a proportion, X\_hat. Print this value to the console. - Find the standard error of X\_hat given N. Print this result to the console. - Calculate the 95% confidence interval of this estimate using the qnorm function. - Save the lower and upper confidence intervals as an object called ci. Save the lower confidence interval first.

```
# Load the data
```

```
data(polls_us_election_2016)
```

```
# Generate an object `polls` that contains data filtered for polls that ended on or after October 31, 2016
polls <- filter(polls_us_election_2016, enddate >= "2016-10-31" & state == "U.S.")
```

```
# How many rows does `polls` contain? Print this value to the console.
nrow(polls)
```

```
## [1] 70
```

```
# Assign the sample size of the first poll in `polls` to a variable called `N`. Print this value to the console.
N <- head(polls$samplesize,1)
N
```

```
## [1] 2220
```

```
# For the first poll in `polls`, assign the estimated percentage of Clinton voters to a variable called X_hat
X_hat <- (head(polls$rawpoll_clinton,1)/100)
X_hat
```

```
## [1] 0.47
```

```
# Calculate the standard error of `X_hat` and save it to a variable called `se_hat`. Print this value to the console.
se_hat <- sqrt(X_hat*(1-X_hat)/N)
se_hat
```

```
## [1] 0.01059279
```

```
# Use `qnorm` to calculate the 95% confidence interval for the proportion of Clinton voters. Save the lower and upper bounds of the confidence interval as an object called ci.
qnorm(0.975)
```

```
## [1] 1.959964
```

```
ci <- c(X_hat - qnorm(0.975)*se_hat, X_hat + qnorm(0.975)*se_hat)
```

## 2. Pollster results for p

Create a new object called pollster\_results that contains the pollster's name, the end date of the poll, the proportion of voters who declared a vote for Clinton, the standard error of this estimate, and the lower and upper bounds of the confidence interval for the estimate.

Instructions - Use the mutate function to define four new columns: X\_hat, se\_hat, lower, and upper. Temporarily add these columns to the polls object that has already been loaded for you. - In the X\_hat

column, convert the raw poll results for Clinton to a proportion. - In the se\_hat column, calculate the standard error of X\_hat for each poll using the sqrt function. - In the lower column, calculate the lower bound of the 95% confidence interval using the qnorm function. - In the upper column, calculate the upper bound of the 95% confidence interval using the qnorm function. - Use the select function to select the columns from polls to save to the new object pollster\_results.

```
# The `polls` object that filtered all the data by date and nation has already been loaded. Examine it
head(polls)
```

```
      state  startdate  enddate
      <fctr> <date>    <date>
1  U.S.    2016-11-03 2016-11-06
2  U.S.    2016-11-01 2016-11-07
3  U.S.    2016-11-02 2016-11-06
4  U.S.    2016-11-04 2016-11-07
5  U.S.    2016-11-03 2016-11-06
6  U.S.    2016-11-03 2016-11-06
6 rows | 1-4 of 16 columns
```

```
# Create a new object called `pollster_results` that contains columns for pollster name, end date, X_hat
polls <- mutate(polls, X_hat = polls$rawpoll_clinton/100, se_hat = sqrt(X_hat*(1-X_hat)/polls$samplesize))
pollster_results <- select(polls, pollster, enddate, X_hat, se_hat, lower, upper)
```

### 3. Comparing to actual results - p

The final tally for the popular vote was Clinton 48.2% and Trump 46.1%. Add a column called hit to pollster\_results that states if the confidence interval included the true proportion p=0.482 or not. What proportion of confidence intervals included p?

Instructions - Use the mutate function to define a new variable called 'hit'. - Use logical expressions to determine if each values in lower and upper span the actual proportion. - Use the mean function to determine the average value in hit and summarize the results using summarize. - Save the result as an object called avg\_hit.

```
# The `pollster_results` object has already been loaded. Examine it using the `head` function.
head(pollster_results)
```

```
      pollster  enddate  X_hat
      <fctr>    <date>    <dbl>
1 ABC News/Washington Post 2016-11-06 0.4700
2 Google Consumer Surveys 2016-11-07 0.3803
3 Ipsos 2016-11-06 0.4200
4 YouGov 2016-11-07 0.4500
5 Gravis Marketing 2016-11-06 0.4700
6 Fox News/Anderson Robbins Research/Shaw & Company Research 2016-11-06 0.4800
6 rows | 1-4 of 7 columns
```

```
# Add a logical variable called `hit` that indicates whether the actual value exists within the confidence interval
avg_hit <- pollster_results %>% mutate(hit=(lower<0.482 & upper>0.482)) %>% summarize(mean(hit))
avg_hit
```

```

                                mean(hit)
                                <dbl>
                                0.3142857
1 row

```

#### 4. Theory of confidence intervals

If these confidence intervals are constructed correctly, and the theory holds up, what proportion of confidence intervals should include  $p$ ?

Possible Answers - ☐ A. 0.05 - ☐ B. 0.31 - ☐ C. 0.50 - ☒ D. 0.95

#### 5. Confidence interval for $d$

A much smaller proportion of the polls than expected produce confidence intervals containing  $p$ . Notice that most polls that fail to include  $p$  are underestimating. The rationale for this is that undecided voters historically divide evenly between the two main candidates on election day.

In this case, it is more informative to estimate the spread or the difference between the proportion of two candidates  $d$ , or  $0.482 - 0.461 = 0.021$  for this election.

Assume that there are only two parties and that  $d = 2p - 1$ . Construct a 95% confidence interval for difference in proportions on election night.

Instructions - Use the mutate function to define a new variable called 'd\_hat' in polls. The new variable subtract the proportion of Trump voters from the proportion of Clinton voters. - Extract the sample size  $N$  from the first poll in your subset object polls. - Extract the difference in proportions of voters  $d\_hat$  from the first poll in your subset object polls. - Use the formula above to calculate  $p$  from  $d\_hat$ . Assign  $p$  to the variable  $X\_hat$ . - Find the standard error of the spread given  $N$ . - Calculate the 95% confidence interval of this estimate of the difference in proportions,  $d\_hat$ , using the qnorm function. - Save the lower and upper confidence intervals as an object called ci. Save the lower confidence interval first.

```

# Add a statement to this line of code that will add a new column named `d_hat` to `polls`. The new column
polls <- polls_us_election_2016 %>% filter(enddate >= "2016-10-31" & state == "U.S.") %>%
  mutate(d_hat = rawpoll_clinton/100 - rawpoll_trump/100)

# Assign the sample size of the first poll in `polls` to a variable called `N`. Print this value to the console
N <- polls$samplesize[1]

# For the difference `d_hat` of the first poll in `polls` to a variable called `d_hat`. Print this value to the console
d_hat <- polls$d_hat[1]
d_hat

## [1] 0.04

# Assign proportion of votes for Clinton to the variable `X_hat`.
X_hat <- (d_hat+1)/2

# Calculate the standard error of the spread and save it to a variable called `se_hat`. Print this value to the console
se_hat <- 2*sqrt(X_hat*(1-X_hat)/N)
se_hat

## [1] 0.02120683

```

```
# Use `qnorm` to calculate the 95% confidence interval for the difference in the proportions of voters.
ci <- c(d_hat - qnorm(0.975)*se_hat, d_hat + qnorm(0.975)*se_hat)
```

## 6. Pollster results for d

Create a new object called `pollster_results` that contains the pollster's name, the end date of the poll, the difference in the proportion of voters who declared a vote either, the standard error of this estimate, and the lower and upper bounds of the confidence interval for the estimate.

Instructions - Use the `mutate` function to define four new columns: 'X\_hat', 'se\_hat', 'lower', and 'upper'. Temporarily add these columns to the `polls` object that has already been loaded for you. - In the `X_hat` column, calculate the proportion of voters for Clinton using `d_hat`. - In the `se_hat` column, calculate the standard error of the spread for each poll using the `sqrt` function. - In the lower column, calculate the lower bound of the 95% confidence interval using the `qnorm` function. - In the upper column, calculate the upper bound of the 95% confidence interval using the `qnorm` function. - Use the `select` function to select the columns from `polls` to save to the new object `pollster_results`.

```
# The subset `polls` data with 'd_hat' already calculated has been loaded. Examine it using the `head` :
head(polls)
```

	state	startdate	enddate
	<fctr>	<date>	<date>
1	U.S.	2016-11-03	2016-11-06
2	U.S.	2016-11-01	2016-11-07
3	U.S.	2016-11-02	2016-11-06
4	U.S.	2016-11-04	2016-11-07
5	U.S.	2016-11-03	2016-11-06
6	U.S.	2016-11-03	2016-11-06

6 rows | 1-4 of 17 columns

```
# Create a new object called `pollster_results` that contains columns for pollster name, end date, d_hat
pollster_results <- polls %>% mutate(X_hat = (d_hat + 1) / 2) %>% mutate(se_hat = 2 * sqrt(X_hat * (1 -
pollster_results
```

pollster	enddate	d_hat
<fctr>	<date>	<dbl>
ABC News/Washington Post	2016-11-06	0.0400
Google Consumer Surveys	2016-11-07	0.0234
Ipsos	2016-11-06	0.0300
YouGov	2016-11-07	0.0400
Gravis Marketing	2016-11-06	0.0400
Fox News/Anderson Robbins Research/Shaw & Company Research	2016-11-06	0.0400
CBS News/New York Times	2016-11-06	0.0400
NBC News/Wall Street Journal	2016-11-05	0.0400
IBD/TIPP	2016-11-07	-0.0150
Selzer & Company	2016-11-06	0.0300

1-10 of 70 rows | 1-3 of 5 columns

## 7. Comparing to actual results - d

What proportion of confidence intervals for the difference between the proportion of voters included d, the actual difference in election day?