

Foundations for statistical inference - Sampling distributions

Souleymane Doumbia

```
knitr::opts_chunk$set(root.dir = getwd())
```

In this lab, you will investigate the ways in which the statistics from a random sample of data can serve as point estimates for population parameters. We're interested in formulating a *sampling distribution* of our estimate in order to learn about the properties of the estimate, such as its distribution.

Setting a seed: We will take some random samples and build sampling distributions in this lab, which means you should set a seed at the start of your lab. If this concept is new to you, review the lab on probability.

Getting Started

Load packages

In this lab, we will explore and visualize the data using the **tidyverse** suite of packages. We will also use the **infer** package for resampling.

Let's load the packages.

```
library(tidyverse)
library(openintro)
library(infer)
library(shiny)
```

The data

A 2019 Gallup report states the following:

The premise that scientific progress benefits people has been embodied in discoveries throughout the ages – from the development of vaccinations to the explosion of technology in the past few decades, resulting in billions of supercomputers now resting in the hands and pockets of people worldwide. Still, not everyone around the world feels science benefits them personally.

Source: World Science Day: Is Knowledge Power?

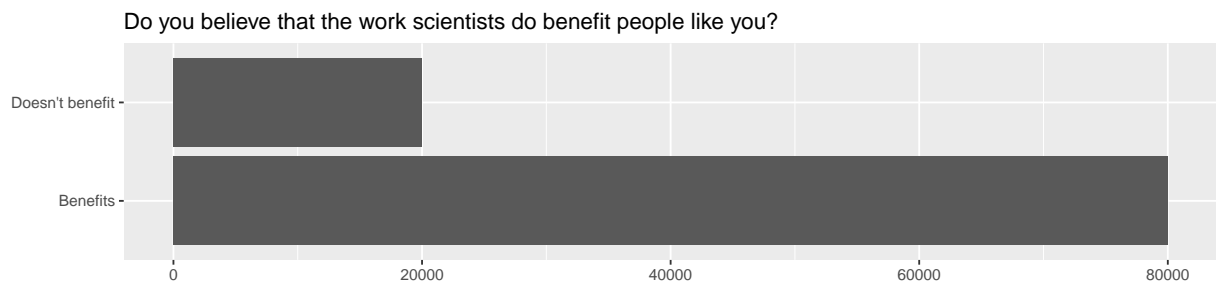
The Wellcome Global Monitor finds that 20% of people globally do not believe that the work scientists do benefits people like them. In this lab, you will assume this 20% is a true population proportion and learn about how sample proportions can vary from sample to sample by taking smaller samples from the population. We will first create our population assuming a population size of 100,000. This means 20,000 (20%) of the population think the work scientists do does not benefit them personally and the remaining 80,000 think it does.

```
global_monitor <- tibble(
  scientist_work = c(rep("Benefits", 80000), rep("Doesn't benefit", 20000))
)
```

The name of the data frame is `global_monitor` and the name of the variable that contains responses to the question “Do you believe that the work scientists do benefit people like you?” is `scientist_work`.

We can quickly visualize the distribution of these responses using a bar plot.

```
ggplot(global_monitor, aes(x = scientist_work)) +
  geom_bar() +
  labs(
    x = "", y = "",
    title = "Do you believe that the work scientists do benefit people like you?"
  ) +
  coord_flip()
```



We can also obtain summary statistics to confirm we constructed the data frame correctly.

```
global_monitor %>%
  count(scientist_work) %>%
  mutate(p = n / sum(n))
```

```
## # A tibble: 2 x 3
##   scientist_work      n      p
##   <chr>          <int> <dbl>
## 1 Benefits       80000  0.8
## 2 Doesn't benefit 20000  0.2
```

The unknown sampling distribution

In this lab, you have access to the entire population, but this is rarely the case in real life. Gathering information on an entire population is often extremely costly or impossible. Because of this, we often take a sample of the population and use that to understand the properties of the population.

If you are interested in estimating the proportion of people who don't think the work scientists do benefits them, you can use the `sample_n` command to survey the population.

```
set.seed(111)
samp1 <- global_monitor %>%
  sample_n(50)
```

This command collects a simple random sample of size 50 from the `global_monitor` dataset, and assigns the result to `samp1`. This is similar to randomly drawing names from a hat that contains the names of all in the population. Working with these 50 names is considerably simpler than working with all 100,000 people in the population.

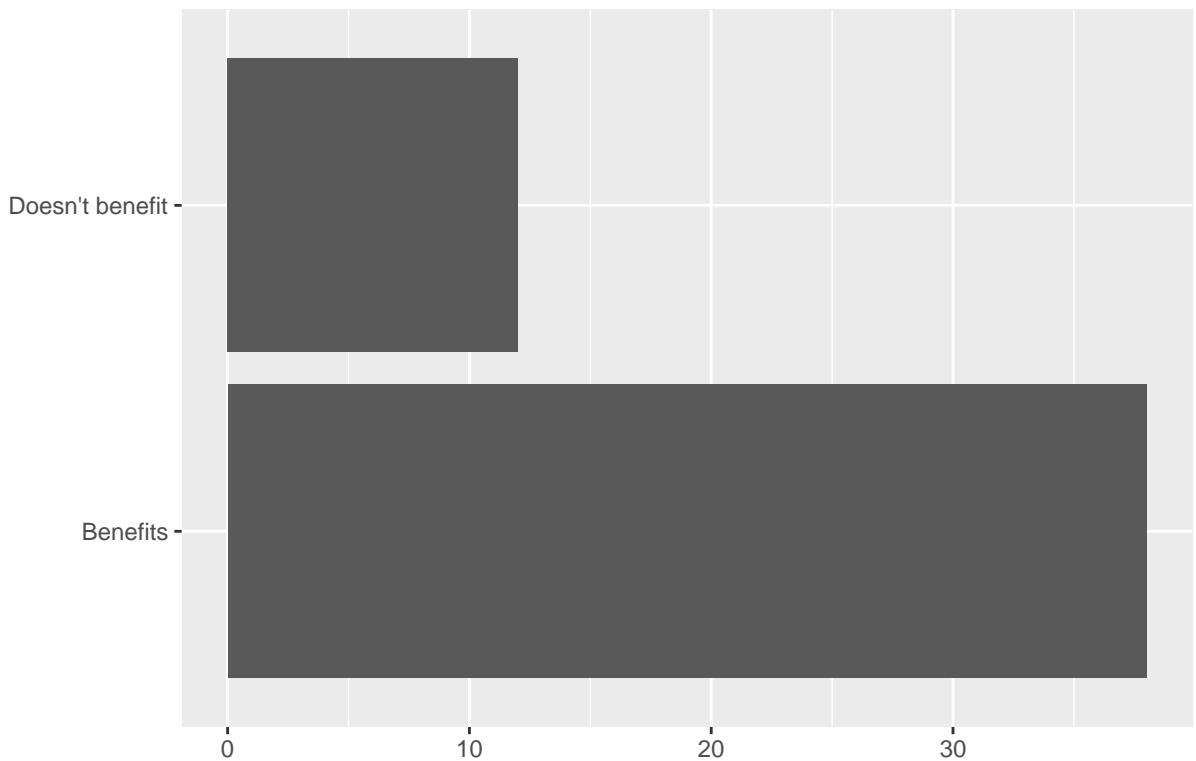
1. Describe the distribution of responses in this sample. How does it compare to the distribution of responses in the population. **Hint:** Although the `sample_n` function takes a random sample of observations (i.e. rows) from the dataset, you can still refer to the variables in the dataset with the same names. Code you presented earlier for visualizing and summarizing the population data will still be useful for the sample, however be careful to not label your proportion `p` since you're now calculating a sample statistic, not a population parameters. You can customize the label of the statistics to indicate that it comes from the sample.

```
samp1 %>%  
  count(scientist_work) %>%  
  mutate(p_samp1 = n / sum(n))
```

```
## # A tibble: 2 x 3  
##   scientist_work      n p_samp1  
##   <chr>          <int>   <dbl>  
## 1 Benefits           38    0.76  
## 2 Doesn't benefit    12    0.24
```

```
ggplot(samp1, aes(x = scientist_work)) +  
  geom_bar() +  
  labs(  
    x = "", y = "",  
    title = "Do you believe that the work scientists do benefit people like you?"  
  ) +  
  coord_flip()
```

Do you believe that the work scientists do benefit people like you?



In this random sample of 50, 24% of respondents says scientific work does not benefit them. This is different from the distribution of responses in the population. And the probability of people in the sample is less than the proportion of population saying scientific research does benefit them. The distributions are similar in the sense that still more people believes scientific research benefit then.

If you're interested in estimating the proportion of all people who do not believe that the work scientists do benefits them, but you do not have access to the population data, your best single guess is the sample mean.

```
samp1 %>%
  count(scientist_work) %>%
  mutate(p_hat = n / sum(n))
```

```
## # A tibble: 2 x 3
##   scientist_work      n p_hat
##   <chr>          <int> <dbl>
## 1 Benefits         38  0.76
## 2 Doesn't benefit  12  0.24
```

Depending on which 50 people you selected, your estimate could be a bit above or a bit below the true population proportion of 0.24. In general, though, the sample proportion turns out to be a pretty good estimate of the true population proportion, and you were able to get it by sampling less than 1% of the population.

2. Would you expect the sample proportion to match the sample proportion of another student's sample? Why, or why not? If the answer is no, would you expect the proportions to be somewhat different or very different? Ask a student team to confirm your answer.

I do not expect my sample proportion to match that of another student because `sample_n()` produces a random sample each time it is called. By default, `sample_n()` function samples with replacement, meaning that the same individual can be selected multiple times in a single sample. Additionally, the composition of the population consists of two distinct subgroups. Each time you draw a sample without replacement, you are removing individuals from the population, and the composition of the remaining population changes. As a result, different samples may end up including different proportions of the subgroups within the population. Therefore, without setting the same seed in R for both myself and another student, our sampled proportions are likely to be different due to the inherent randomness in the sampling process. However, if we both set the same seed in R before using `sample_n()`, we can ensure that our samples are drawn in the same random order, and therefore, our sampled proportions are likely to be very similar.

3. Take a second sample, also of size 50, and call it `samp2`. How does the sample proportion of `samp2` compare with that of `samp1`? Suppose we took two more samples, one of size 100 and one of size 1000. Which would you think would provide a more accurate estimate of the population proportion?

```
#set.seed(222)
samp2 <- global_monitor %>%
  sample_n(50)

samp2 %>%
  count(scientist_work) %>%
  mutate(p_samp2 = n / sum(n))
```

```
## # A tibble: 2 x 3
##   scientist_work      n p_samp2
##   <chr>          <int>   <dbl>
## 1 Benefits         41    0.82
## 2 Doesn't benefit    9    0.18
```

Comparison of `samp1` and `samp2`: Since `samp1` and `samp2` are both of size 50 and are drawn randomly, their sample proportions can vary. If the random seed was set the same for both samples, then the two samples and their proportions would be identical. However, if different seeds were used (or if no seed was set, relying on true randomness), then the sample proportions could differ due to the inherent variability in random sampling.

Comparison between samples of size 100 and 1000: Typically, a sample of size 1000 will provide a more accurate estimate of the population proportion than one of size 100. This is because a larger sample size reduces the sampling variability, making the estimate more precise. Again, without setting a seed, there's inherent randomness in the results. When a specific seed is set, the results become reproducible.

Not surprisingly, every time you take another random sample, you might get a different sample proportion. It's useful to get a sense of just how much variability you should expect when estimating the population mean this way. The distribution of sample proportions, called the *sampling distribution (of the proportion)*, can help you understand this variability. In this lab, because you have access to the population, you can build up the sampling distribution for the sample proportion by repeating the above steps many times. Here, we use R to take 15,000 different samples of size 50 from the population, calculate the proportion of responses in each sample, filter for only the *Doesn't benefit* responses, and store each result in a vector called `sample_props50`. Note that we specify that `replace = TRUE` since sampling distributions are constructed by sampling with replacement.

```
sample_props50 <- global_monitor %>%
  rep_sample_n(size = 50, reps = 15000, replace = TRUE) %>%
  count(scientist_work) %>%
  mutate(p_hat = n / sum(n)) %>%
  filter(scientist_work == "Doesn't benefit")
```

And we can visualize the distribution of these proportions with a histogram.

```
set.seed(333)
ggplot(data = sample_props50, aes(x = p_hat)) +
  geom_histogram(binwidth = 0.02) +
  labs(
    x = "p_hat (Doesn't benefit)",
    title = "Sampling distribution of p_hat",
    subtitle = "Sample size = 50, Number of samples = 15000"
  )
```

Next, you will review how this set of code works.

4. How many elements are there in `sample_props50`? Describe the sampling distribution, and be sure to specifically note its center. Make sure to include a plot of the distribution in your answer.

```
sample_props50 %>% nrow()
```

```
## [1] 15000
```

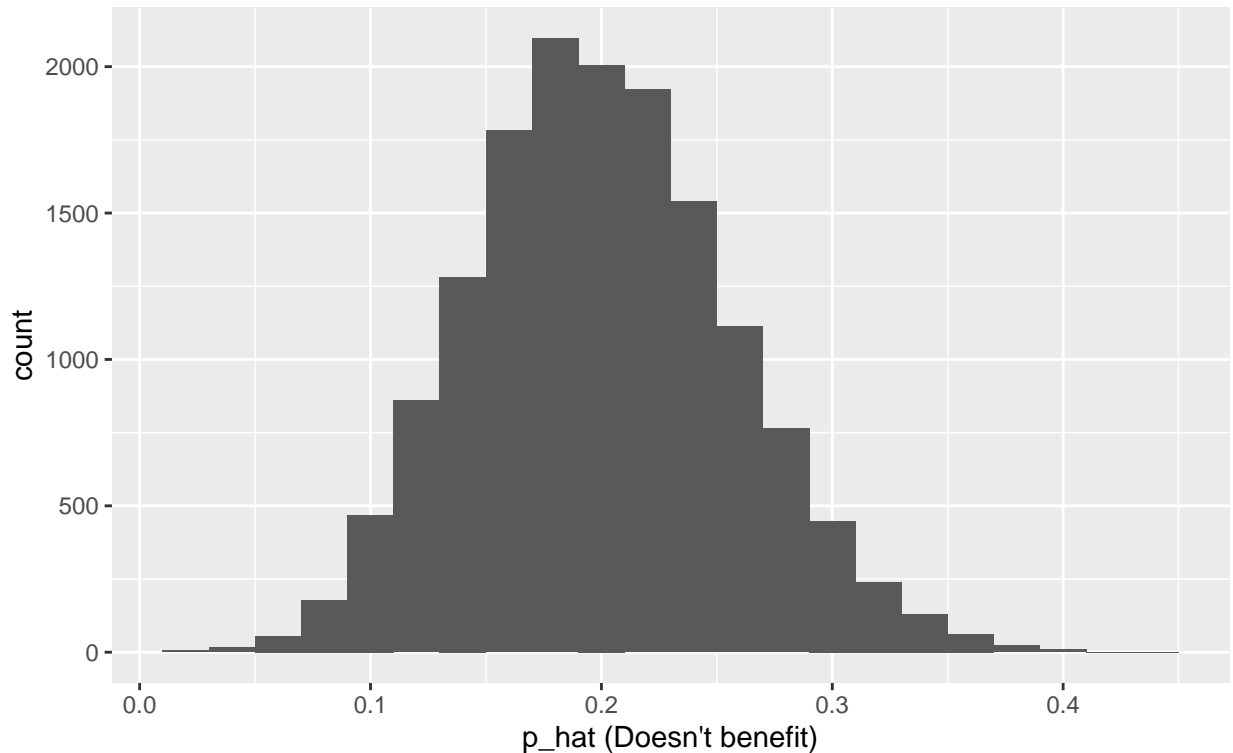
The code above tells us that there 15,000 elements in `sample_props50`.

The proportions' sampling distribution has a shape that looks normally distributed, like a bell curve. Its middle point is at 0.2 (its center), which is the same as the true proportion of the whole population. This suggests that the proportion we see in 'sample_props50' usually mirrors the actual proportion in the broader population. Take a look at the graph below:

```
ggplot(data = sample_props50, aes(x = p_hat)) +
  geom_histogram(binwidth = 0.02) +
  labs(
    x = "p_hat (Doesn't benefit)",
    title = "Sampling distribution of p_hat",
    subtitle = "Sample size = 50, Number of samples = 15000")
```

Sampling distribution of \hat{p}

Sample size = 50, Number of samples = 15000



Interlude: Sampling distributions

The idea behind the `rep_sample_n` function is *repetition*. Earlier, you took a single sample of size n (50) from the population of all people in the population. With this new function, you can repeat this sampling procedure `rep` times in order to build a distribution of a series of sample statistics, which is called the **sampling distribution**.

Note that in practice one rarely gets to build true sampling distributions, because one rarely has access to data from the entire population.

Without the `rep_sample_n` function, this would be painful. We would have to manually run the following code 15,000 times

```
global_monitor %>%
  sample_n(size = 50, replace = TRUE) %>%
  count(scientist_work) %>%
  mutate(p_hat = n / sum(n)) %>%
  filter(scientist_work == "Doesn't benefit")
```

```
## # A tibble: 1 x 3
##   scientist_work      n p_hat
##   <chr>          <int> <dbl>
## 1 Doesn't benefit    10  0.2
```

as well as store the resulting sample proportions each time in a separate vector.

Note that for each of the 15,000 times we computed a proportion, we did so from a **different** sample!

5. To make sure you understand how sampling distributions are built, and exactly what the `rep_sample_n` function does, try modifying the code to create a sampling distribution of **25 sample proportions** from **samples of size 10**, and put them in a data frame named `sample_props_small`. Print the output. How many observations are there in this object called `sample_props_small`? What does each observation represent?

```
set.seed(444)
sample_props_small <- global_monitor %>%
  rep_sample_n(size = 10, reps = 25, replace = TRUE) %>%
  count(scientist_work) %>%
  mutate(p_hat = n / sum(n)) %>%
  filter(scientist_work == "Doesn't benefit")

sample_props_small %>% nrow()
```

```
## [1] 23
```

```
sample_props_small
```

```
## # A tibble: 23 x 4
## # Groups:   replicate [23]
##   replicate scientist_work      n p_hat
##   <int> <chr>          <int> <dbl>
## 1      1 1 Doesn't benefit      1  0.1
## 2      2 2 Doesn't benefit      3  0.3
## 3      3 3 Doesn't benefit      2  0.2
## 4      4 4 Doesn't benefit      1  0.1
## 5      5 5 Doesn't benefit      5  0.5
## 6      6 6 Doesn't benefit      3  0.3
## 7      7 7 Doesn't benefit      2  0.2
## 8      8 8 Doesn't benefit      1  0.1
## 9      9 9 Doesn't benefit      3  0.3
## 10    10 10 Doesn't benefit      1  0.1
## # i 13 more rows
```

There are 23 observations in ‘`sample_props_small`’. Each observation represent the proportion of people among a population of 10 people who says that the work scientists doesn’t benefit them.

Sample size and the sampling distribution

Mechanics aside, let’s return to the reason we used the `rep_sample_n` function: to compute a sampling distribution, specifically, the sampling distribution of the proportions from samples of 50 people.

```
ggplot(data = sample_props50, aes(x = p_hat)) +
  geom_histogram(binwidth = 0.02)
```

The sampling distribution that you computed tells you much about estimating the true proportion of people who think that the work scientists do doesn’t benefit them. Because the sample proportion is an unbiased estimator, the sampling distribution is centered at the true population proportion, and the spread of the

distribution indicates how much variability is incurred by sampling only 50 people at a time from the population.

In the remainder of this section, you will work on getting a sense of the effect that sample size has on your sampling distribution.

6. Use the app below to create sampling distributions of proportions of *Doesn't benefit* from samples of size 10, 50, and 100. Use 5,000 simulations. What does each observation in the sampling distribution represent? How does the mean, standard error, and shape of the sampling distribution change as the sample size increases? How (if at all) do these values change if you increase the number of simulations? (You do not need to include plots in your answer.)

Each observation in the sampling distribution gives a fraction that represents how many individuals in that sample think the work of scientists doesn't benefit them. As we collected data from larger groups, the average of these fractions started approaching 0.2, which matches the known proportion of the whole population believing that scientists' work doesn't benefit them. It's noteworthy that the standard error, which tells us how much these fractions vary across samples, decreases when we consider more substantial sample sizes. As our sample sizes grew, the distribution of these fractions started to take on a bell-shaped curve, suggesting a normal distribution. Yet, when we increased the number of times we took these samples (simulations), neither the average value nor its standard error changed significantly.

More Practice

So far, you have only focused on estimating the proportion of those you think the work scientists doesn't benefit them. Now, you'll try to estimate the proportion of those who think it does.

Note that while you might be able to answer some of these questions using the app, you are expected to write the required code and produce the necessary plots and summary statistics. You are welcome to use the app for exploration.

7. Take a sample of size 15 from the population and calculate the proportion of people in this sample who think the work scientists do enhances their lives. Using this sample, what is your best point estimate of the population proportion of people who think the work scientists do enhances their lives?

```
set.seed(555)
samp3 <- global_monitor %>%
  sample_n(15) %>%
  count(scientist_work) %>%
  mutate(p_samp3 = n/sum(n)) %>%
  filter(scientist_work == 'Benefits')
samp3
```

```
## # A tibble: 1 x 3
##   scientist_work      n p_samp3
##   <chr>          <int>   <dbl>
## 1 Benefits         11    0.733
```

My best point estimate of people who think the work scientists do enhance their lives using a sample of size 15 is 73.3%.

8. Since you have access to the population, simulate the sampling distribution of proportion of those who think the work scientists do enhances their lives for samples of size 15 by taking 2000 samples from the population of size 15 and computing 2000 sample proportions. Store these proportions in as `sample_props15`. Plot the data, then describe the shape of this sampling distribution. Based on this sampling distribution, what would you guess the true proportion of those who think the work scientists do enhances their lives to be? Finally, calculate and report the population proportion.

```
set.seed(666)
set.seed(8000)
sample_props15 <- global_monitor %>%
  rep_sample_n(size = 15, reps = 2000, replace = TRUE) %>%
  count(scientist_work) %>%
  mutate(p_hat = n / sum(n)) %>%
  filter(scientist_work == "Benefits")
glimpse(sample_props15)

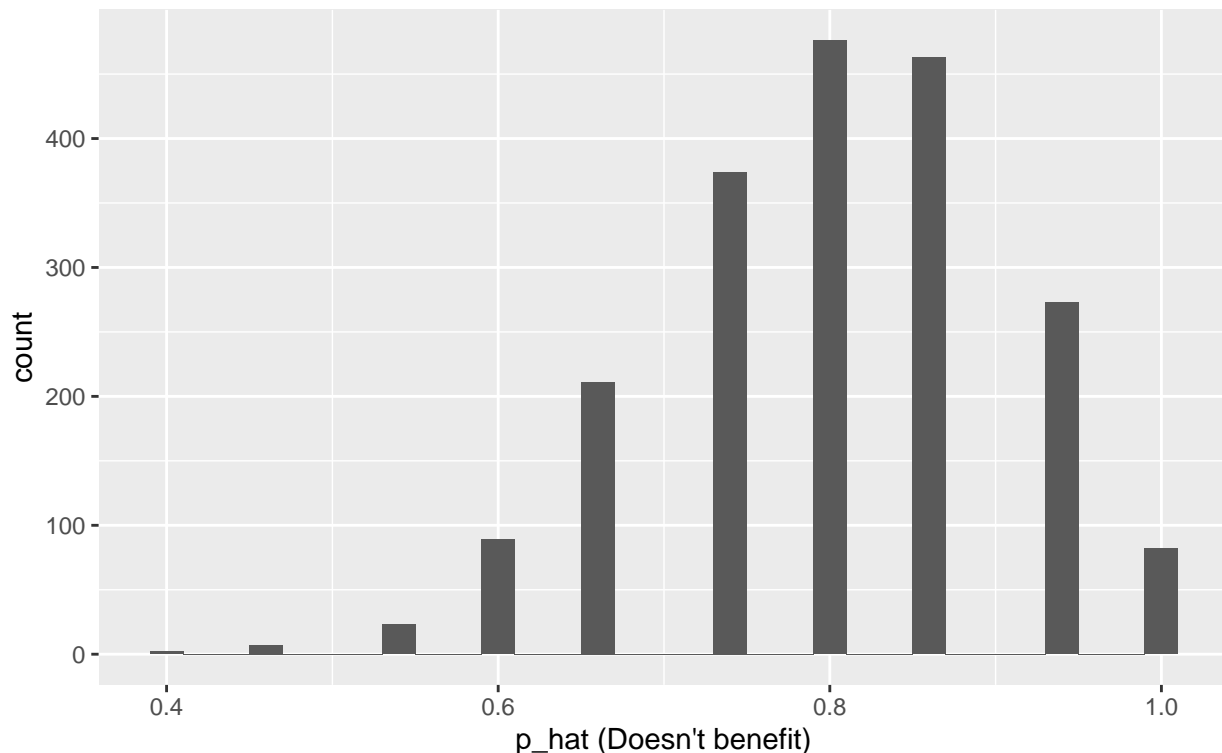
## Rows: 2,000
## Columns: 4
## Groups: replicate [2,000]
## $ replicate      <int> 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, ~
## $ scientist_work <chr> "Benefits", "Benefits", "Benefits", "Benefits", "Benefi~
## $ n              <int> 10, 11, 14, 14, 15, 11, 9, 13, 13, 12, 12, 13, 14, 12, ~
## $ p_hat          <dbl> 0.6666667, 0.7333333, 0.9333333, 0.9333333, 1.0000000, ~
```

Visualization:

```
ggplot(data = sample_props15, aes(x = p_hat)) +
  geom_histogram(binwidth = 0.02) +
  labs(
    x = "p_hat (Doesn't benefit)",
    title = "Sampling distribution of p_hat",
    subtitle = "Sample size = 15, Number of samples = 2000"
  )
```

Sampling distribution of \hat{p}

Sample size = 15, Number of samples = 2000

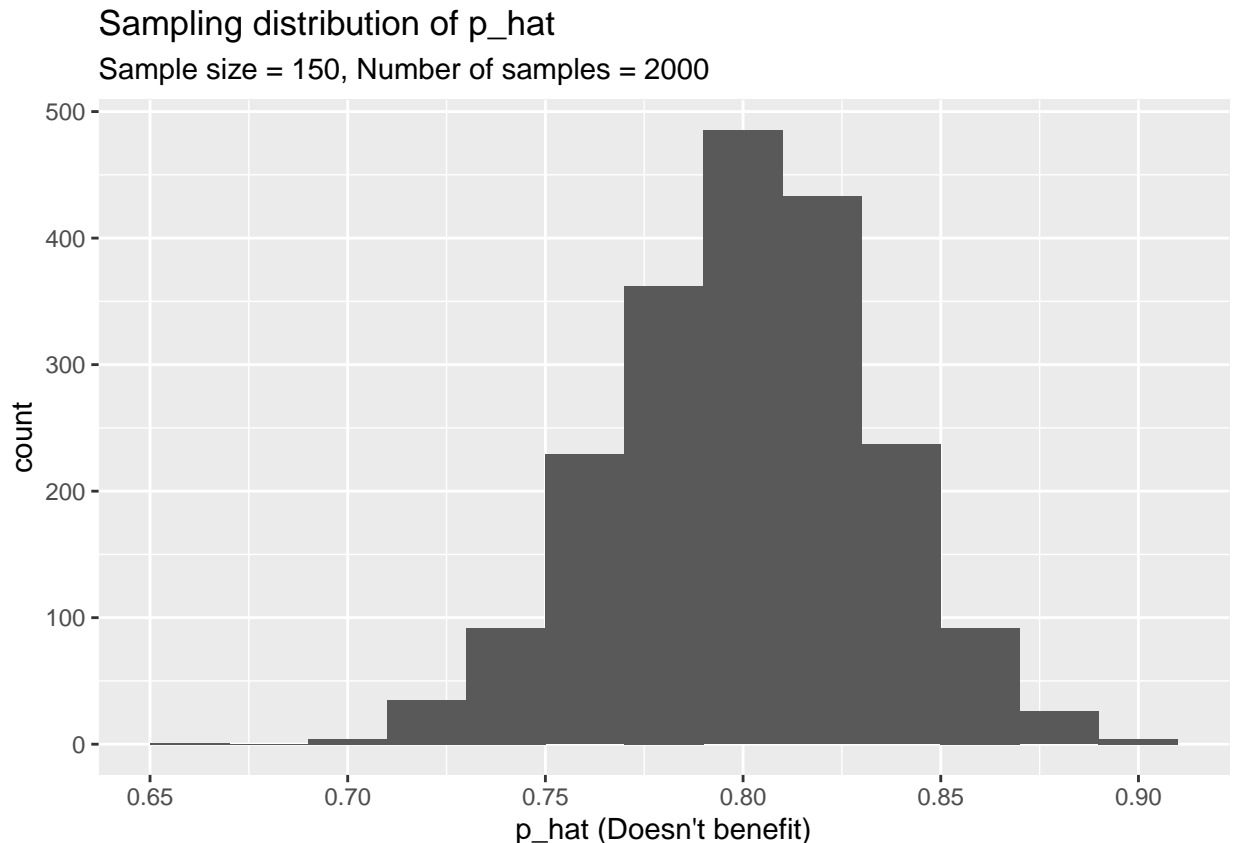


The distributions seems to be bell shaped, so tends to be normally distributed. From the graph, the true proportion of those who think the work scientists do enhances their lives is 0.8.

9. Change your sample size from 15 to 150, then compute the sampling distribution using the same method as above, and store these proportions in a new object called `sample_props150`. Describe the shape of this sampling distribution and compare it to the sampling distribution for a sample size of 15. Based on this sampling distribution, what would you guess to be the true proportion of those who think the work scientists do enhances their lives?

```
set.seed(777)
set.seed(8000)
sample_props150 <- global_monitor %>%
  rep_sample_n(size = 150, reps = 2000, replace = TRUE) %>%
  count(scientist_work) %>%
  mutate(p_hat = n / sum(n)) %>%
  filter(scientist_work == "Benefits")

ggplot(data = sample_props150, aes(x = p_hat)) +
  geom_histogram(binwidth = 0.02) +
  labs(
    x = "p_hat (Doesn't benefit)",
    title = "Sampling distribution of p_hat",
    subtitle = "Sample size = 150, Number of samples = 2000"
  )
```



Again, the distributions seems to be bell shaped, so tends to be normally distributed. From the plot, The center of the distribution of the sample proportions is centered on 0.8 which is equal to the true proportion of the population. The true proportion of those who think the work scientists do enhances their lives is 80% same as the sample size of 15. .

10. Of the sampling distributions from 2 and 3, which has a smaller spread? If you're concerned with making estimates that are more often close to the true value, would you prefer a sampling distribution with a large or small spread?

```
Parmeters_15_150 <- data.frame(
  Sample_Error = c(sample_props15 = sd(sample_props15$p_hat), sample_props150 = sd(sample_props150$p_hat)),
  Sample_Mean = c(mean(sample_props15$p_hat), mean(sample_props150$p_hat)))
Parmeters_15_150
```

```
##           Sample_Error Sample_Mean
## sample_props15    0.1046482  0.8017667
## sample_props150    0.0328532  0.8005833
```

From the result above, we can recall that sample_props150 has the smaller spread and smaller sample mean. I would prefer a sampling distribution with a small spread if I am concerned with making estimates that are more often close to the true value. Also, we can easily see that as sample size increase the parameter values decrease, approaching the true parameters of the population.