

Problem Set 4

Daniel Shapiro

9/22/2022

Question 1 Background:

Download the `demo.csv` dataset from the course website. The dataset contains information from a sample of countries in the year 2000, taken from the Democracy Time-Series Dataset. It includes the following variables:

- **Nation:** country name
- **GDP:** GDP per capita in constant US dollars
- **FHouse:** Freedom House rating (a measure of the level of political and civil liberties in a country, on a scale from 1.0 (most free) to 7.0 (least free))
- **OECD:** a dummy variable indicating OECD status
- **regime:** a variable coded from the Freedom House rating that indicates whether a country is free (1.0-2.5), partly free (2.51-5.5), or not free (5.51-7.0).

1a) Using `stargazer()`, show the summary statistics for your dataset. Briefly interpret for the GDP variable.

```
# I like putting this sort of thing in a separate chunk so I don't have to run it a billion times.
```

```
demo <- read.csv("demo.csv")
```

```
# Had to put some extra things into the title brackets to get the table to show up as  
# more than just LaTeX. Also put header = FALSE to suppress the initial lines.
```

```
stargazer(demo, header = FALSE)
```

Table 1:

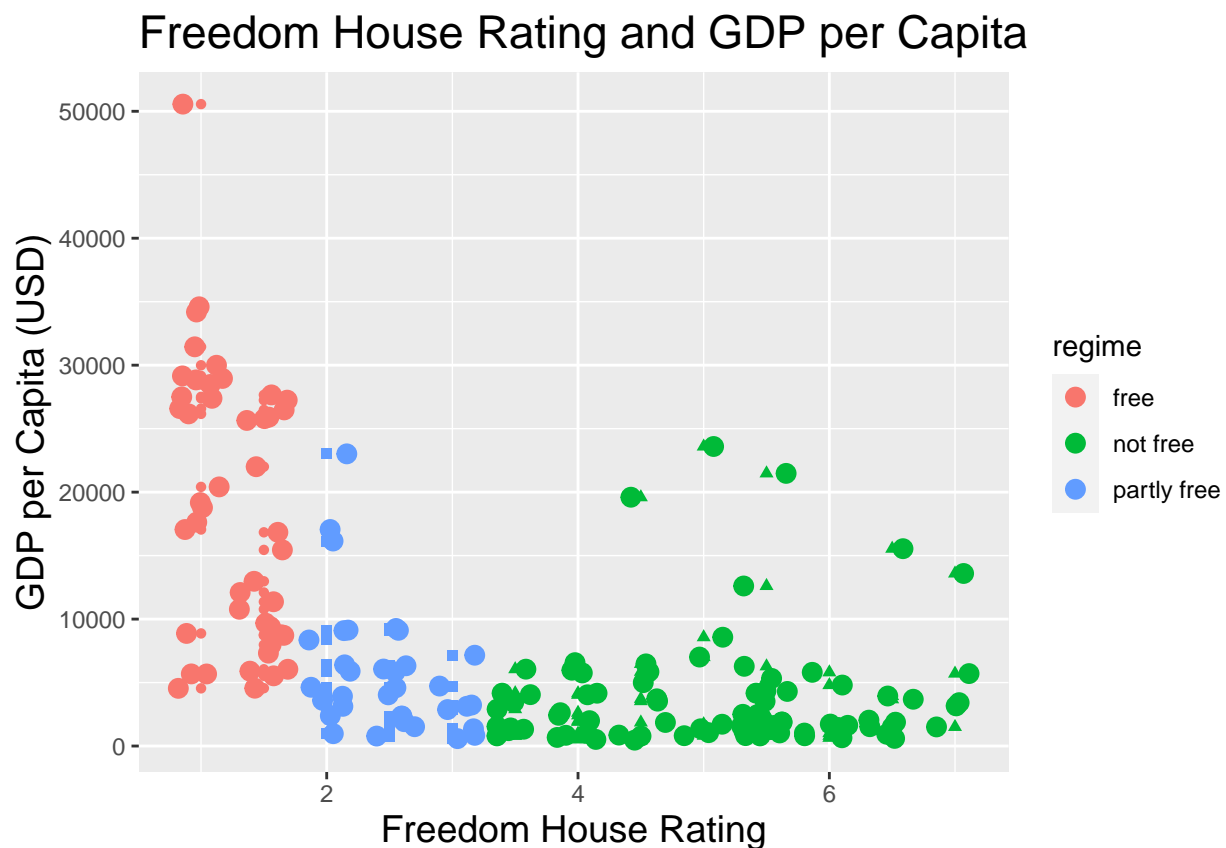
Statistic	N	Mean	St. Dev.	Min	Max
GDP	161	8,468.658	9,516.346	463	50,564
FHouse	161	3.460	1.897	1.000	7.000
OECD	161	0.186	0.391	0	1

Brief interpretation of the GDP row: First, there's "N" which is just the number of observations – 161. The "Mean" is the mean GDP per capita across all observations. The standard deviation value is relatively high, so that means that the data is spread across a rather wide area. The "Min" and "Max" columns show the lowest GDP per capita measure (463) and the highest (50,564) in the dataset.

1b) Produce an appropriately named and labeled plot of GDP (on the y-axis) against FHouse (on the x-axis) using the `ggplot()` function (including a legend). Do the following:

- 1) Use a different color for data points representing different regimes (i.e. free, partly free, or not free)
- 2) In case some see your plot in black and white, use different point types for each regime category.
- 3) Adjust the size of the axis labels, axis titles, and title to make them more legible.
- 4) Increase the size of your points, and use `geom_jitter()` to make them more legible.

```
ggplot(demo, aes(x = FHouse, y = GDP, color = regime)) +
  geom_point(aes(shape = regime)) +
  labs(title = "Freedom House Rating and GDP per Capita",
       x = "Freedom House Rating",
       y = "GDP per Capita (USD)") +
  theme(plot.title = element_text(size = 17),
        axis.title = element_text(size = 14)) +
  geom_jitter(size = 3)
```



1c) Calculate the conditional expectation and the conditional standard deviation of GDP for the three regime types, using a function that takes as an input the type of regime and returns the conditional mean and standard deviation. What do the conditional summary statistics suggest about the relationship between democracy and wealth? Briefly explain.

To find the expected value of X , we can use the formula: $E[X] = \sum(x)P(x)$.

```
conditional <- function(type){

setup1 <- demo %>%
  filter(regime == "free")
exp1 <- mean(setup1$GDP)
sd1 <- sd(setup1$GDP)

setup2 <- demo %>%
  filter(regime == "partly free")
exp2 <- mean(setup2$GDP)
sd2 <- sd(setup2$GDP)

setup3 <- demo %>%
  filter(regime == "not free")
exp3 <- mean(setup3$GDP)
sd3 <- sd(setup3$GDP)

  if(type == "free"){return <- c(exp1, sd1)}
  else if(type == "partly free"){return <- c(exp2, sd2)}
  else{return <- c(exp3, sd3)}

return
}

conditional("free")
```

```
## [1] 19017.49 10671.70
```

```
conditional("partly free")
```

```
## [1] 5640.853 4975.443
```

```
conditional("not free")
```

```
## [1] 3852.171 4492.659
```

According to these statistics, a higher level of democracy tends to correlate with a higher expected GDP per capita. Regimes labeled “free” have a higher GDP per capita than regimes labeled “partly free,” which in turn have a higher GDP per capita than regimes labeled “not free.”

1d) Using the `geom_histogram()` command in `ggplot()`, produce a density plot of GDP per capita. Overlay two vertical lines, in different colors, for the mean and the median of that variable. Annotate the graph to mark these lines informatively using `geom_text()` (hint: `geom_text` takes a dataframe as an input, so start by making a dataframe of your labels and their desired position). What does the relationship between the mean and the median, as shown on the plot, tell you about the variable GDP per capita?

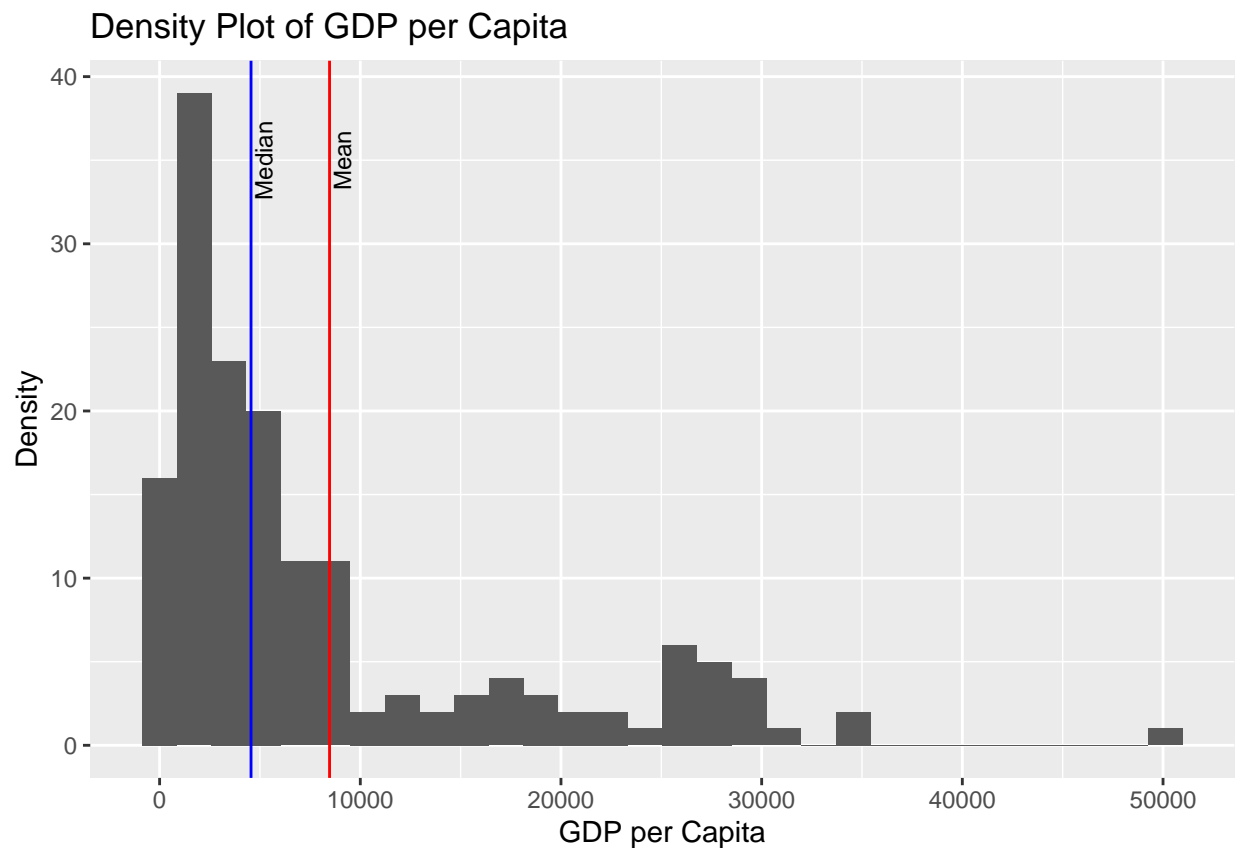
```
xint <- mean(demo$GDP)
medint <- median(demo$GDP)
```

```

textdata <- data.frame(xint, medint)

ggplot(demo, aes(x = GDP)) +
  geom_histogram(bins = 30) +
  geom_vline(xintercept = textdata$xint, color = "red") +
  geom_vline(xintercept = textdata$medint, color = "blue") +
  geom_text(data = textdata, x = xint, y = 35, label = "Mean", size = 3, angle = 90, vjust = 1.25) +
  geom_text(data = textdata, x = medint, y = 35, label = "Median", size = 3, angle = 90, vjust = 1.25) +
  labs(title = "Density Plot of GDP per Capita",
       x = "GDP per Capita",
       y = "Density")

```



Here we see that the median is farther to the left (lower) than the mean. So that means that there are some serious outliers on the right – on the richer side. Most countries are grouped around lower GDP per capita, but there are a few countries that are significantly higher.

1e) Write a function that returns the amount of GDP data that falls within 1, 1.96, and 3 standard deviations of the mean. Compare these results with what we would expect if the data were perfectly normally distributed.

```

# Created new column for standard deviation

demo1 <- demo %>%

```

```

mutate("sd" = (GDP - mean(GDP)))

# Set up my function. I use absolute value to show the +/- aspect.

gdpdeviation <- function(data){

initial1 <- data %>%
  dplyr::filter(abs(sd) <= abs(sd(GDP)))
initial2 <- data %>%
  dplyr::filter(abs(sd) <= 1.96*abs(sd(GDP)))
initial3 <- data %>%
  dplyr::filter(abs(sd) <= 3*abs(sd(GDP)))

standard <- nrow(initial1)/nrow(data)
midstandard <- nrow(initial2)/nrow(data)
largestandard <- nrow(initial3)/nrow(data)

c(standard, midstandard, largestandard)

}

gdpdeviation(demo1)

```

```
## [1] 0.8322981 0.9192547 0.9937888
```

In a perfect normal distribution, 68% of the population falls within 1 standard deviation of the mean, while 95% falls within 1.96 standard deviations and 99.7% falls within 3 standard deviations. In this data, it's 83%, 92%, and 99.4%. The 99.4% figure is quite close, 92% is close, and 83% is relatively far away. So the data look rather different here than in a perfect normal distribution.

1f) Now draw 100, 1,000, and 10,000 samples the length of the dataframe, with replacement, from the GDP data (bootstrap). Plot a histogram of the sample means (you can use the above ggplot() code or the hist() function). How well do these sampling distributions approximate the normal distribution? How close are they to the mean value of GDP?

First, I set everything up.

```

demonew <- demo %>%
  select(GDP)

# Not sure if this is a popular choice, but I've used the rep_sample_n() function from
# the infer package before for bootstrapping. It's awesome. Used it here!

first <- demonew %>% rep_sample_n(size = 161, replace = TRUE, reps = 100)

# replace = TRUE ensures that we are bootstrapping.

second <- demonew %>% rep_sample_n(size = 161, replace = TRUE, reps = 1000)
third <- demonew %>% rep_sample_n(size = 161, replace = TRUE, reps = 10000)

# Here, I used nesting and mapping to apply the mean function to all. Love these functions!

```

```

firstmean <- first %>%
  group_by(replicate) %>%
  nest() %>%
  mutate(mean = map(.x = data, .f = ~mean(.x$GDP, na.rm = TRUE))) %>%
  select(-data) %>%
  ungroup()

secondmean <- second %>%
  group_by(replicate) %>%
  nest() %>%
  mutate(mean = map(.x = data, .f = ~mean(.x$GDP, na.rm = TRUE))) %>%
  select(-data) %>%
  ungroup()

thirdmean <- third %>%
  group_by(replicate) %>%
  nest() %>%
  mutate(mean = map(.x = data, .f = ~mean(.x$GDP, na.rm = TRUE))) %>%
  select(-data) %>%
  ungroup()

```

Next, I create histograms.

```

# First, I need to set columns as numeric -- otherwise I will get errors.

firstmean$mean <- as.numeric(firstmean$mean)
secondmean$mean <- as.numeric(secondmean$mean)
thirdmean$mean <- as.numeric(thirdmean$mean)

# Now, I want to change the column names back to "GDP" so that my function will work.
# I also need to add an "sd" column. EDIT

# Now, for plots.

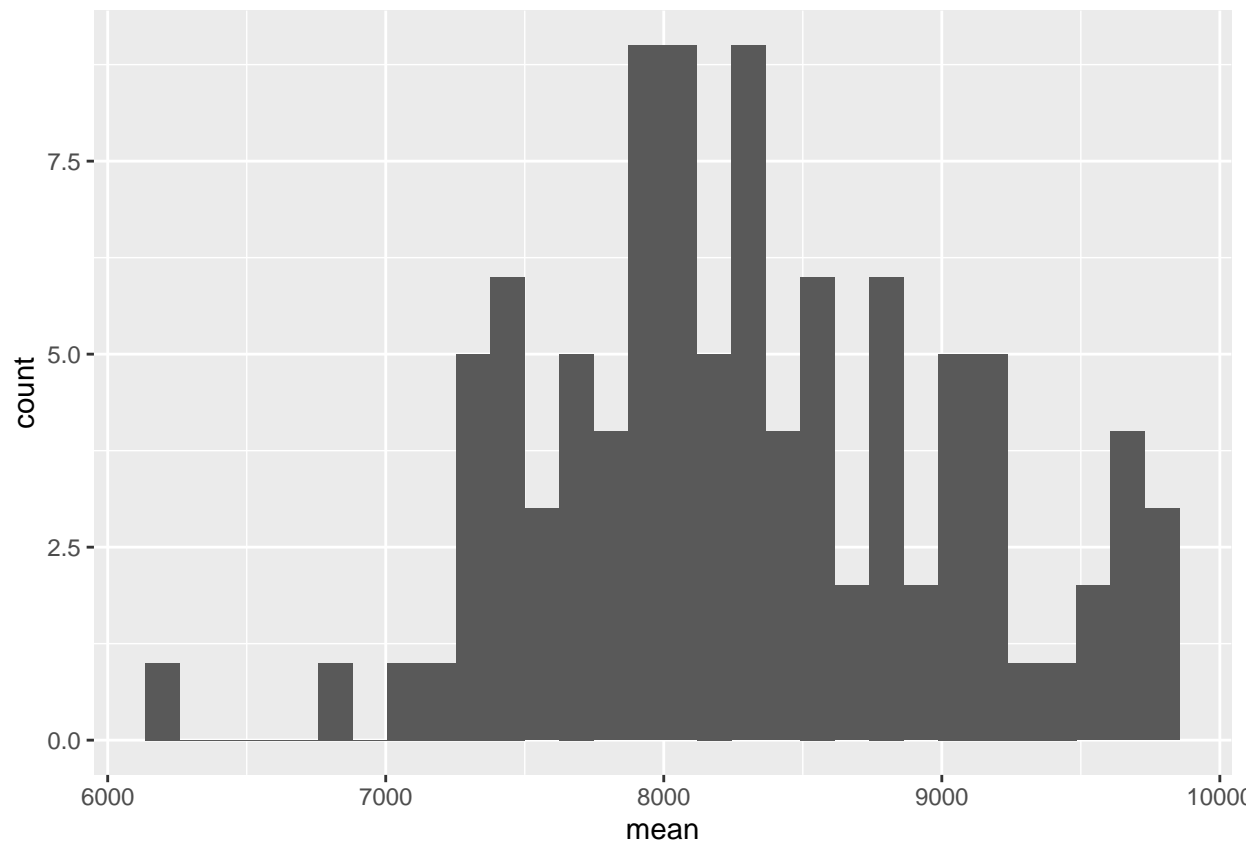
ggplot(firstmean, aes(x = mean)) +
  geom_histogram()

```

```

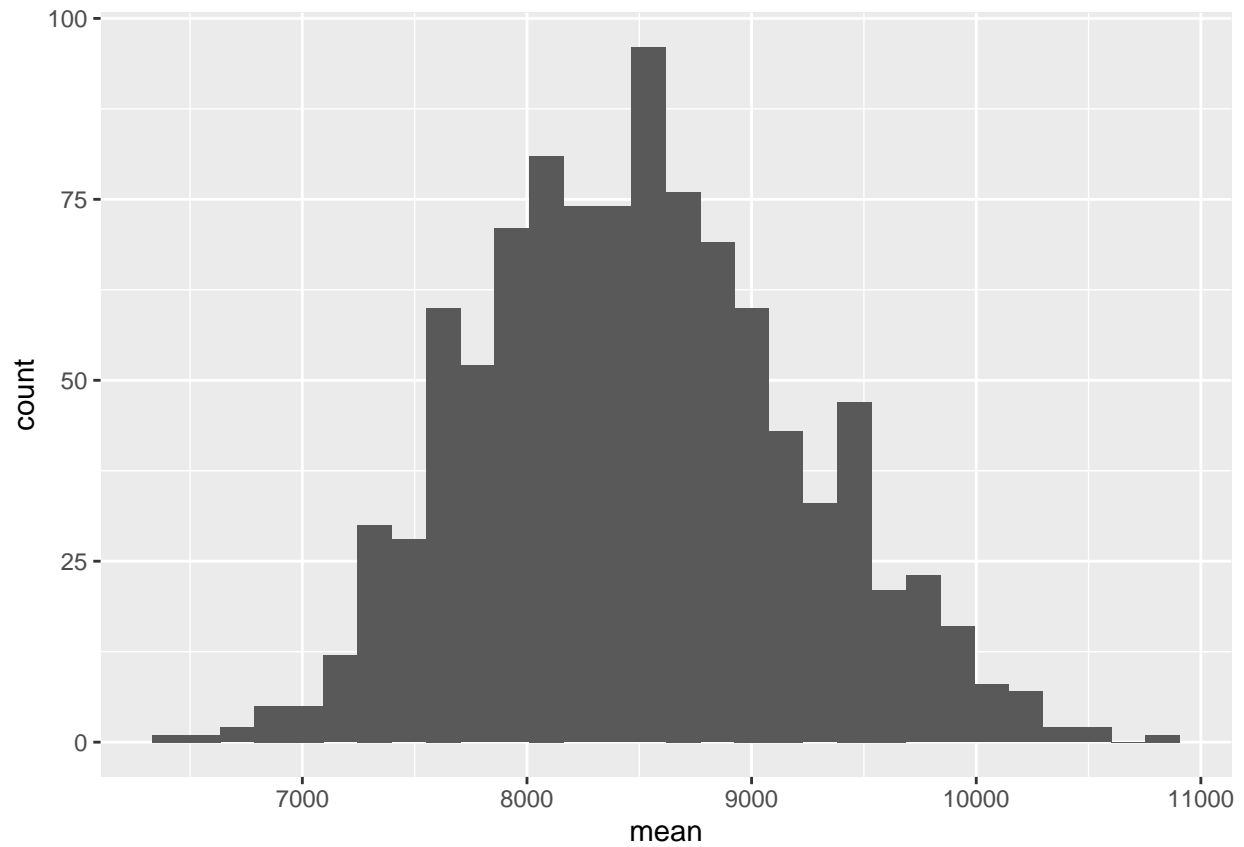
## 'stat_bin()' using 'bins = 30'. Pick better value with 'binwidth'.

```



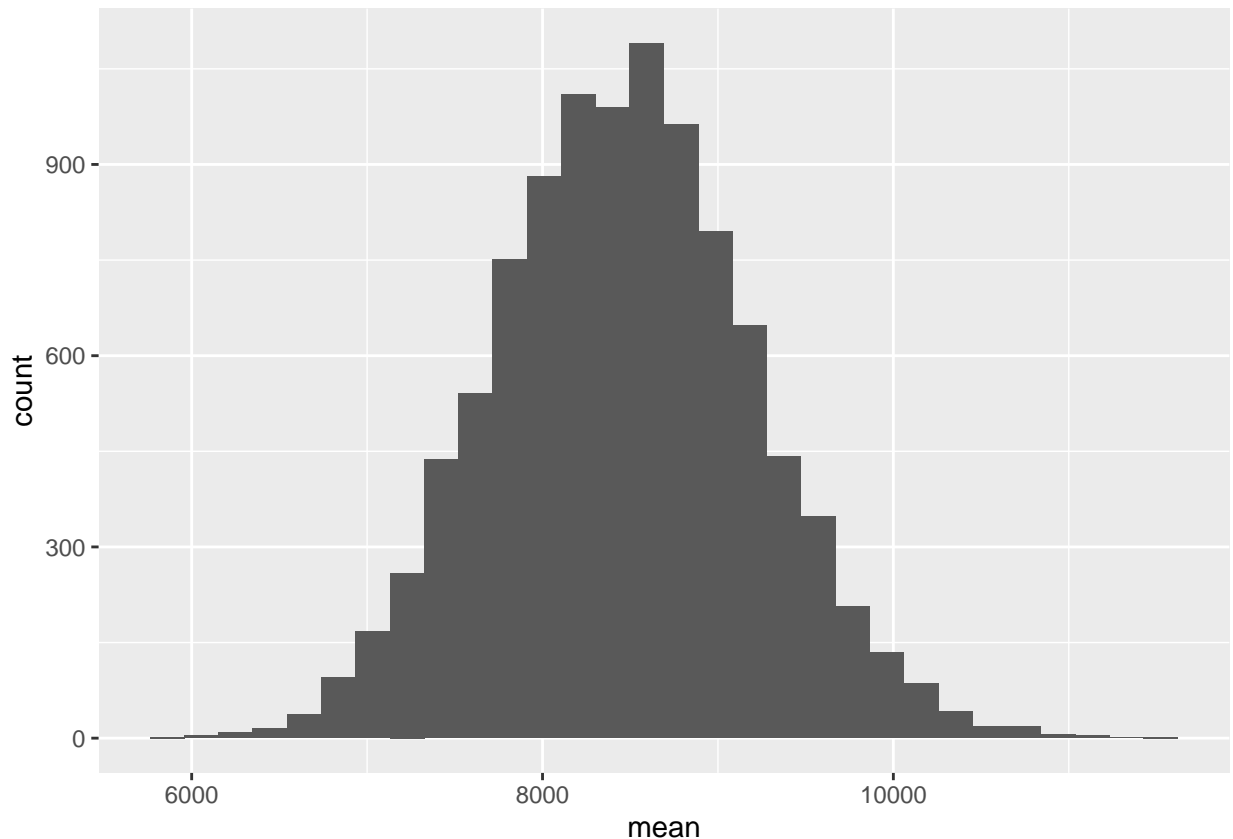
```
ggplot(secondmean, aes(x = mean)) +  
  geom_histogram()
```

```
## 'stat_bin()' using 'bins = 30'. Pick better value with 'binwidth'.
```



```
ggplot(thirdmean, aes(x = mean)) +  
  geom_histogram()
```

```
## 'stat_bin()' using 'bins = 30'. Pick better value with 'binwidth'.
```

So after the histograms, we see that the higher the number of samples, the more even the distribution of means is.

Now, I want to look at the standard deviation and see how well the data approximates the normal distribution. Just looking at the graphs, they certainly look much closer to the normal distribution, because they have more data. This is due to the Central Limit Theorem. To see if this bears true in the data, I run my `gdpdeviation()` function that I created earlier. As we can see below, the distributions of the means do tend to get closer and closer to the normal distribution. Additionally, the means do tend to be rather close to the initial mean, which was 8468.6583851. We can see this graphically and in the tables.

Below is my code.

```
# I have to change the names of the "mean" column for it to work in my function.
```

```
setnames(firstmean, "mean", "GDP")
setnames(secondmean, "mean", "GDP")
setnames(thirdmean, "mean", "GDP")
```

```
# Here, I will use my gdpdeviation() function. I have to create a SD column for each.
```

```
firstmean1 <- firstmean %>%
  mutate("sd" = (GDP - mean(GDP)))

secondmean1 <- secondmean %>%
  mutate("sd" = (GDP - mean(GDP)))

thirdmean1 <- thirdmean %>%
```

```
mutate("sd" = (GDP - mean(GDP)))  
gdpdeviation(firstmean1)
```

```
## [1] 0.63 0.95 1.00
```

```
gdpdeviation(secondmean1)
```

```
## [1] 0.665 0.957 0.999
```

```
gdpdeviation(thirdmean1)
```

```
## [1] 0.6856 0.9498 0.9963
```

1g) Calculate the standard error of the mean, using the mathematical formula from lecture and using the three sampling distributions you created above. Explain what the standard error means here.