

Foundations for statistical inference - Sampling distributions

Jimmy Ng

In this lab, we 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.

The data

We consider real estate data from the city of Ames, Iowa. The details of every real estate transaction in Ames is recorded by the City Assessor's office. Our particular focus for this lab will be all residential home sales in Ames between 2006 and 2010. This collection represents our population of interest. In this lab we would like to learn about these home sales by taking smaller samples from the full population. Let's load the data.

We see that there are quite a few variables in the data set, enough to do a very in-depth analysis. For this lab, we'll restrict our attention to just two of the variables: the above ground living area of the house in square feet (`Gr.Liv.Area`) and the sale price (`SalePrice`). To save some effort throughout the lab, create two variables with short names that represent these two variables.

Let's look at the distribution of area in our population of home sales by calculating a few summary statistics and making a histogram.

##	Min.	1st Qu.	Median	Mean	3rd Qu.	Max.
##	334	1126	1442	1500	1743	5642



1. Describe this population distribution. JN: The distribution appears to look symmetrical; however, it's slightly skewed to the right with mean larger than median.

The unknown sampling distribution

In this lab we 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 we were interested in estimating the mean living area in Ames based on a sample, we can use the following command to survey the population.

```
samp1 <- sample(area, 50)
```

This command collects a simple random sample of size 50 from the vector `area`, which is assigned to `samp1`. This is like going into the City Assessor's database and pulling up the files on 50 random home sales. Working with these 50 files would be considerably simpler than working with all 2930 home sales.

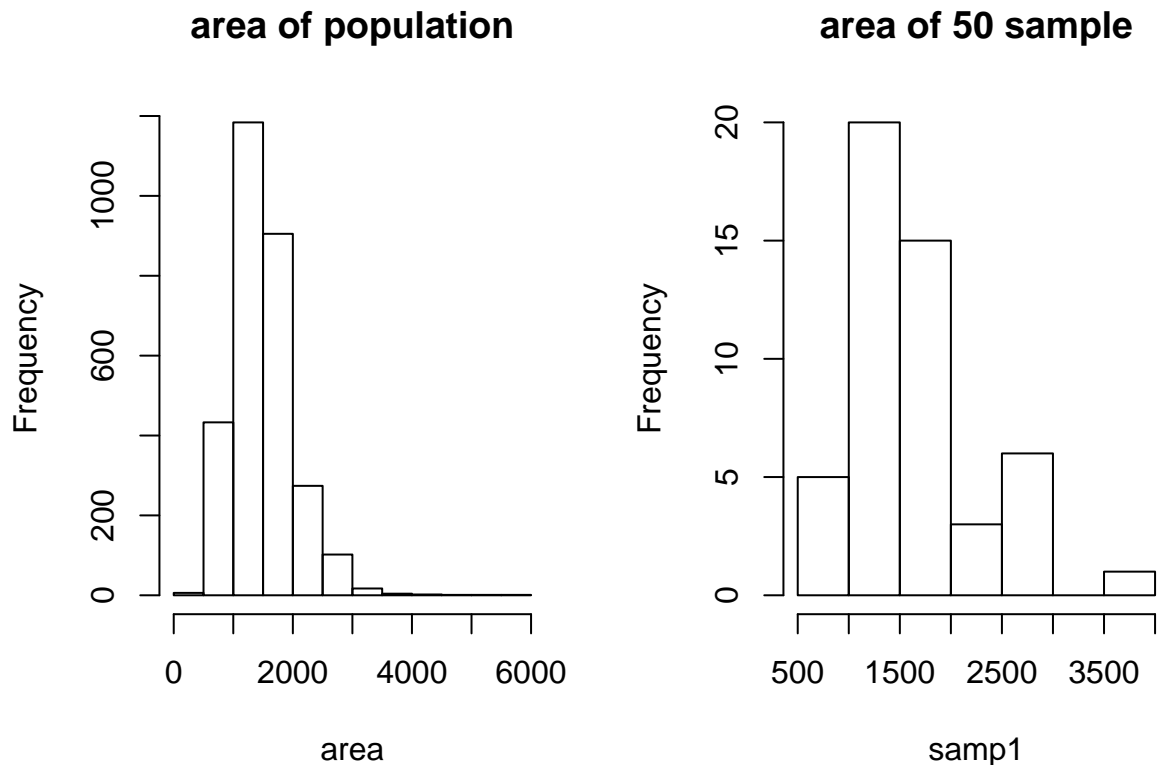
2. Describe the distribution of this sample. How does it compare to the distribution of the population?

```
library(purrr)
map(1:2, function(x) summary(list(area, samp1)[[x]]))
```

```
## [[1]]
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##      334   1126   1442   1500   1743   5642
##
## [[2]]
```

```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##      893    1180    1469    1601    1725    3627
```

```
par(mfrow = c(1, 2))
hist(area, main = 'area of population')
hist(samp1, main = 'area of 50 sample')
```



JN: The distribution looks similar between the overall population and the sample, both are slightly skewed to the right and both have mean larger than median.

If we're interested in estimating the average living area in homes in Ames using the sample, our best single guess is the sample mean.

```
mean(samp1)
```

```
## [1] 1601.02
```

Depending on which 50 homes you selected, your estimate could be a bit above or a bit below the true population mean of 1499.69 square feet. In general, though, the sample mean turns out to be a pretty good estimate of the average living area, and we were able to get it by sampling less than 3% of the population.

3. Take a second sample, also of size 50, and call it **samp2**. How does the mean of **samp2** compare with the mean 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 mean?

```
samp2 <- sample(area, 50)
mean(samp1); mean(samp2)
```

```
## [1] 1601.02
```

```
## [1] 1603.22
```

```
# the mean of samp1 (1593.7) is larger than mean of samp2 (1531.08); both are larger than the mean of p
```

```
set.seed(1234)
samp3 <- sample(area, 100); samp4 <- sample(area, 1000)
library(purrr)
map(1:5, function(x) mean(list(area, samp1, samp2, samp3, samp4)[[x]]))
```

```
## [[1]]
## [1] 1499.69
##
## [[2]]
## [1] 1601.02
##
## [[3]]
## [1] 1603.22
##
## [[4]]
## [1] 1510.92
##
## [[5]]
## [1] 1496.451
```

```
# the largest sample (samp4 with 1000 random samples from the population) has the closest mean to the p
# mean of the population (1499.69) vs. samp4 (1496.451)
```

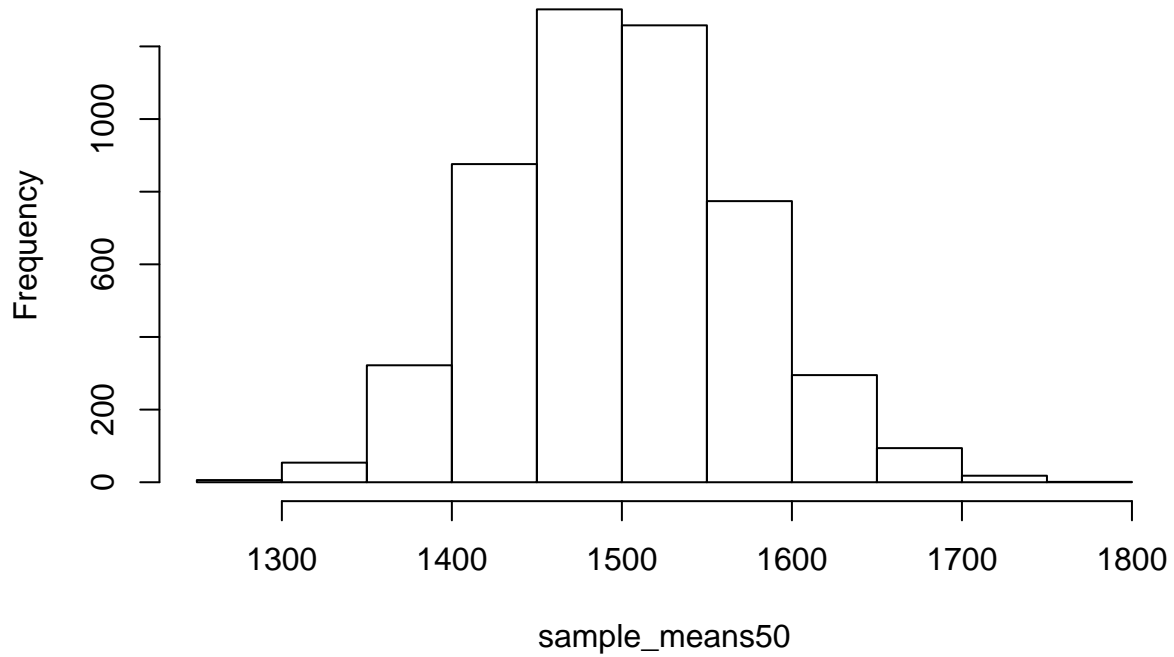
Not surprisingly, every time we take another random sample, we get a different sample mean. It's useful to get a sense of just how much variability we should expect when estimating the population mean this way. The distribution of sample means, called the *sampling distribution*, can help us understand this variability. In this lab, because we have access to the population, we can build up the sampling distribution for the sample mean by repeating the above steps many times. Here we will generate 5000 samples and compute the sample mean of each.

```
sample_means50 <- rep(NA, 5000)

for(i in 1:5000){
  samp <- sample(area, 50)
  sample_means50[i] <- mean(samp)
}

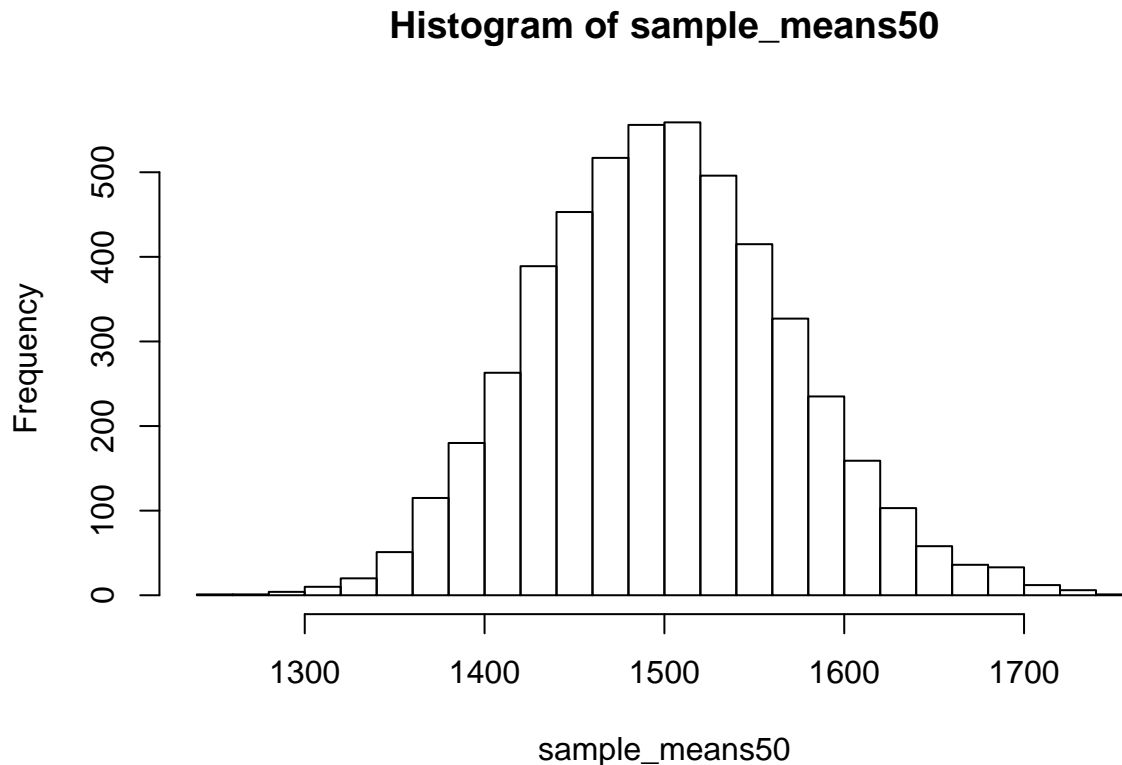
hist(sample_means50)
```

Histogram of sample_means50



If you would like to adjust the bin width of your histogram to show a little more detail, you can do so by changing the `breaks` argument.

```
hist(sample_means50, breaks = 25)
```



Here we use R to take 5000 samples of size 50 from the population, calculate the mean of each sample, and store each result in a vector called `sample_means50`. On the next page, we'll review how this set of code works.

4. How many elements are there in `sample_means50`? Describe the sampling distribution, and be sure to specifically note its center. Would you expect the distribution to change if we instead collected 50,000 sample means?

JN: There are 5000 means collected in the vector of “sample_means50”. There are 5000 sample means, and each mean is generated from 50 random samples from the population. This is a normal distribution (perfectly described by the Central Limit Theorem). The center of this sampling distribution is extremely close to the true mean of the population, i.e. 1499.69. This sampling distribution of 5000 is large enough to approximate the population; running another simulation of collecting 50000 sample means would not change its shape.

Interlude: The for loop

Let's take a break from the statistics for a moment to let that last block of code sink in. You have just run your first `for` loop, a cornerstone of computer programming. The idea behind the `for` loop is *iteration*: it allows you to execute code as many times as you want without having to type out every iteration. In the case above, we wanted to iterate the two lines of code inside the curly braces that take a random sample of size 50 from `area` then save the mean of that sample into the `sample_means50` vector. Without the `for` loop, this would be painful:

```
sample_means50 <- rep(NA, 5000)

samp <- sample(area, 50)
sample_means50[1] <- mean(samp)
```

```
samp <- sample(area, 50)
sample_means50[2] <- mean(samp)

samp <- sample(area, 50)
sample_means50[3] <- mean(samp)

samp <- sample(area, 50)
sample_means50[4] <- mean(samp)
```

and so on...

With the `for` loop, these thousands of lines of code are compressed into a handful of lines. We've added one extra line to the code below, which prints the variable `i` during each iteration of the `for` loop. Run this code.

Let's consider this code line by line to figure out what it does. In the first line we *initialized a vector*. In this case, we created a vector of 5000 zeros called `sample_means50`. This vector will store values generated within the `for` loop.

The second line calls the `for` loop itself. The syntax can be loosely read as, "for every element `i` from 1 to 5000, run the following lines of code". You can think of `i` as the counter that keeps track of which loop you're on. Therefore, more precisely, the loop will run once when `i = 1`, then once when `i = 2`, and so on up to `i = 5000`.

The body of the `for` loop is the part inside the curly braces, and this set of code is run for each value of `i`. Here, on every loop, we take a random sample of size 50 from `area`, take its mean, and store it as the `i`th element of `sample_means50`.

In order to display that this is really happening, we asked R to print `i` at each iteration. This line of code is optional and is only used for displaying what's going on while the `for` loop is running.

The `for` loop allows us to not just run the code 5000 times, but to neatly package the results, element by element, into the empty vector that we initialized at the outset.

5. To make sure you understand what you've done in this loop, try running a smaller version. Initialize a vector of 100 zeros called `sample_means_small`. Run a loop that takes a sample of size 50 from `area` and stores the sample mean in `sample_means_small`, but only iterate from 1 to 100. Print the output to your screen (type `sample_means_small` into the console and press enter). How many elements are there in this object called `sample_means_small`? What does each element represent?

```
sample_means_small <- as.vector(rep(0, 100))

for(i in 1:length(sample_means_small)){
  sample_means_small[i] <- mean(sample(area, 50))
}

sample_means_small
```

```
## [1] 1496.22 1531.72 1497.94 1391.06 1448.08 1440.08 1461.60 1518.32
## [9] 1593.92 1570.14 1480.08 1496.56 1448.12 1465.28 1552.16 1655.90
## [17] 1552.20 1451.66 1448.26 1629.78 1543.44 1524.34 1589.28 1442.90
## [25] 1428.26 1523.96 1572.02 1402.54 1535.02 1621.78 1476.48 1536.56
## [33] 1567.28 1688.98 1516.94 1471.20 1593.12 1523.42 1592.44 1487.62
## [41] 1602.94 1409.30 1475.68 1507.24 1415.44 1432.98 1520.88 1547.48
## [49] 1504.24 1575.34 1513.68 1481.68 1687.90 1578.26 1496.40 1501.30
## [57] 1576.58 1578.90 1534.92 1489.14 1658.96 1565.70 1552.58 1508.34
## [65] 1472.28 1596.28 1614.52 1611.68 1489.70 1511.54 1444.40 1486.84
## [73] 1479.76 1418.54 1490.90 1470.40 1647.56 1572.38 1479.52 1510.32
```

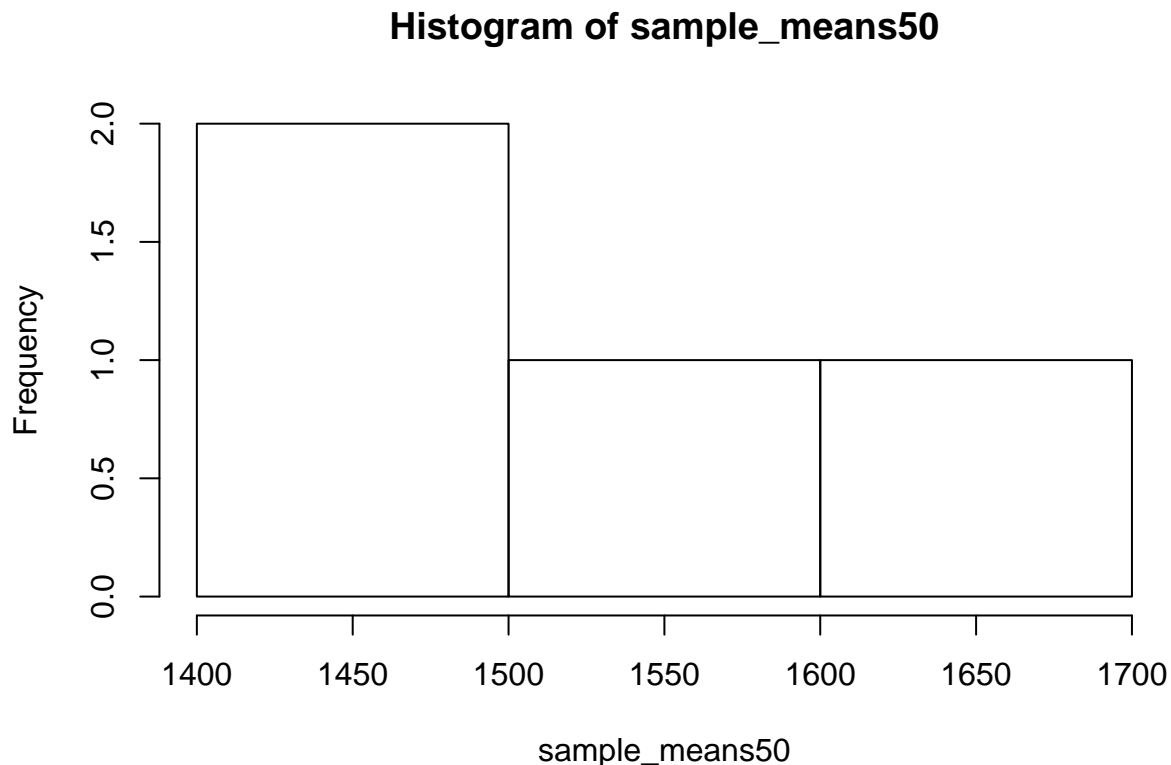
```
## [81] 1537.42 1642.64 1443.86 1539.26 1534.26 1521.34 1515.78 1366.04
## [89] 1418.84 1449.26 1431.66 1411.42 1485.18 1431.16 1527.62 1590.80
## [97] 1418.12 1418.80 1496.76 1609.56
```

JN: There are 100 elements and each represents a mean from a randomly selected sample (size = 50).

Sample size and the sampling distribution

Mechanics aside, let's return to the reason we used a `for` loop: to compute a sampling distribution, specifically, this one.

```
hist(sample_means50)
```



The sampling distribution that we computed tells us much about estimating the average living area in homes in Ames. Because the sample mean is an unbiased estimator, the sampling distribution is centered at the true average living area of the the population, and the spread of the distribution indicates how much variability is induced by sampling only 50 home sales.

To get a sense of the effect that sample size has on our distribution, let's build up two more sampling distributions: one based on a sample size of 10 and another based on a sample size of 100.

```
sample_means10 <- rep(NA, 5000)
sample_means100 <- rep(NA, 5000)

for(i in 1:5000){
  samp <- sample(area, 10)
  sample_means10[i] <- mean(samp)
  samp <- sample(area, 100)
```



```

sample_means100[i] <- mean(samp)
}

```

Here we're able to use a single `for` loop to build two distributions by adding additional lines inside the curly braces. Don't worry about the fact that `samp` is used for the name of two different objects. In the second command of the `for` loop, the mean of `samp` is saved to the relevant place in the vector `sample_means10`. With the mean saved, we're now free to overwrite the object `samp` with a new sample, this time of size 100. In general, anytime you create an object using a name that is already in use, the old object will get replaced with the new one.

To see the effect that different sample sizes have on the sampling distribution, plot the three distributions on top of one another.

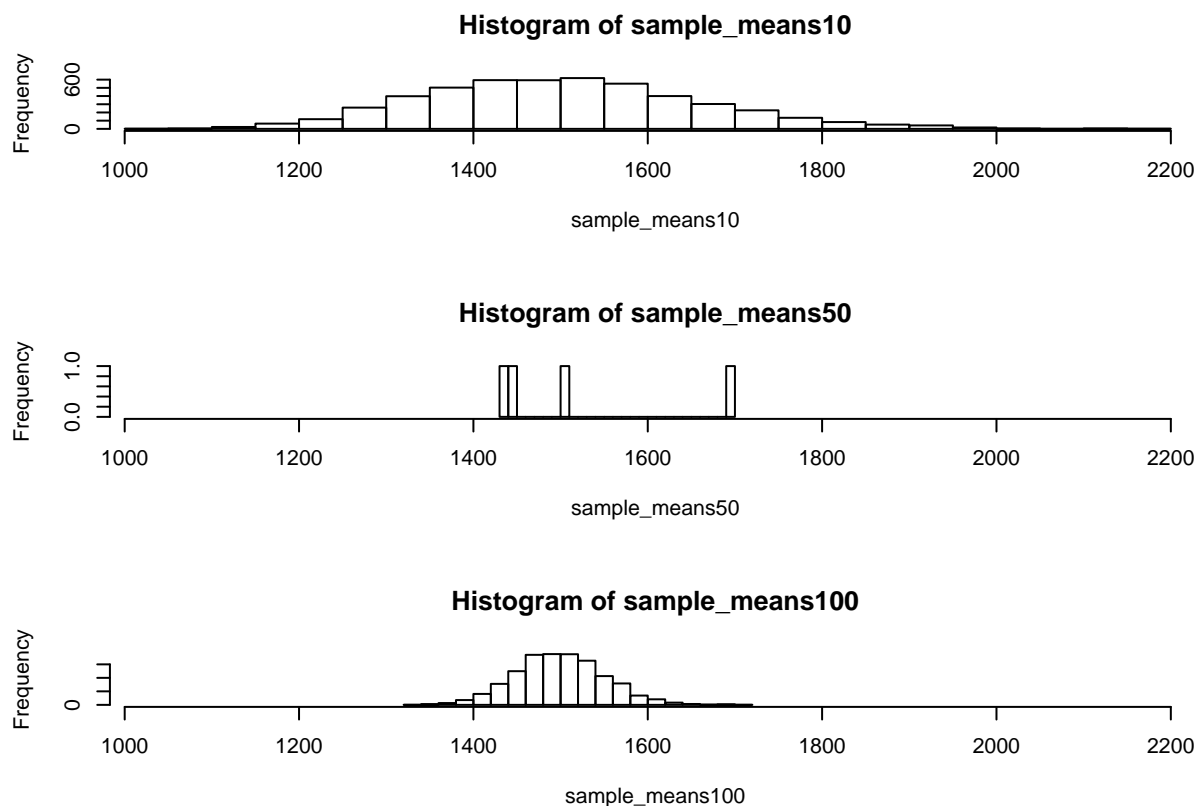
```

par(mfrow = c(3, 1))

xlimits <- range(sample_means10)

hist(sample_means10, breaks = 20, xlim = xlimits)
hist(sample_means50, breaks = 20, xlim = xlimits)
hist(sample_means100, breaks = 20, xlim = xlimits)

```



The first command specifies that you'd like to divide the plotting area into 3 rows and 1 column of plots (to return to the default setting of plotting one at a time, use `par(mfrow = c(1, 1))`). The `breaks` argument specifies the number of bins used in constructing the histogram. The `xlim` argument specifies the range of the x-axis of the histogram, and by setting it equal to `xlimits` for each histogram, we ensure that all three histograms will be plotted with the same limits on the x-axis.

6. When the sample size is larger, what happens to the center? What about the spread? JN: When the

sample size gets larger, the distribution gets tighter around the center, the spread gets smaller and the distribution looks closer to a normal distribution. In other words, sample size does matter. A large enough sample size can help us closely approximate a normal distribution for doing estimation for a population.

On your own

So far, we have only focused on estimating the mean living area in homes in Ames. Now you'll try to estimate the mean home price.

- Take a random sample of size 50 from `price`. Using this sample, what is your best point estimate of the population mean?

```
set.seed(1234)
print(paste("The best point estimate of the population mean using a random sample of size 50 is ",
            mean(sample(price, 50)),
            ".",
            sep = ""))
```

```
## [1] "The best point estimate of the population mean using a random sample of size 50 is 173386.16."
```

“The best point estimate of the population mean using a random sample of size 50 is 173386.16.”

- Since you have access to the population, simulate the sampling distribution for \bar{x}_{price} by taking 5000 samples from the population of size 50 and computing 5000 sample means. Store these means in a vector called `sample_means50`. Plot the data, then describe the shape of this sampling distribution. Based on this sampling distribution, what would you guess the mean home price of the population to be? Finally, calculate and report the population mean.

```
simulation <- function(var, iteration = 5000, size = 50, seed = 1234){
  if(!require(purrr)){install.packages("purrr"); require(purrr)}
  if(!require(plyr)){install.packages("plyr"); require(plyr)}

  set.seed(seed)
  sim <- as.vector(rep(NA, iteration))
  sim <- map(1:iteration, function(x) sim[x] <- mean(sample(var, size)))
  sim <- sim %>%
    plyr::ldply(., data.frame)
  sim <- as.vector(sim[,1])
}
```

```
sample_means50 <- simulation(price, iter = 5000, size = 50, seed = 1234)
```

```
## Loading required package: plyr
```

```
##
```

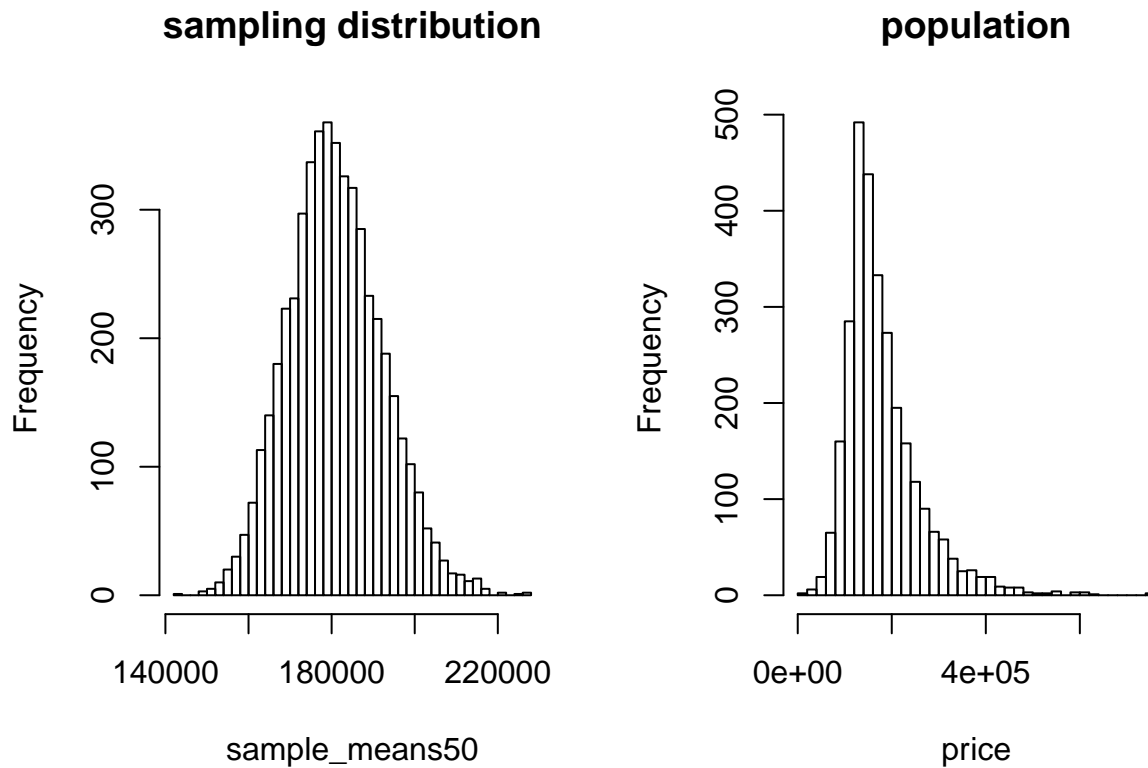
```
## Attaching package: 'plyr'
```

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

```
##
```

```
## compact
```

```
par(mfrow = c(1, 2))
hist(sample_means50, breaks = 50, main = "sampling distribution")
hist(price, breaks = 50, main = "population")
```



```
mean(sample_means50) # [1] 180955.9
```

```
## [1] 180955.9
```

```
mean(price) # [1] 180796.1
```

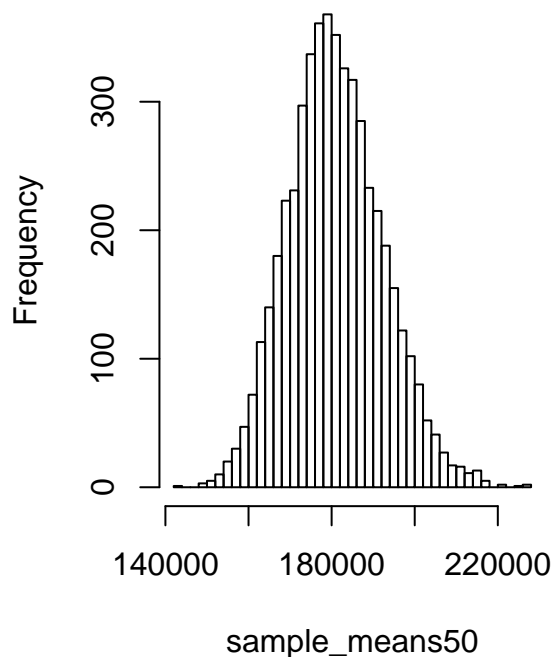
```
## [1] 180796.1
```

The shape of the sampling distribution resembles a normal distribution. The shape and spread look normal with center around 181000 with few extreme outliers on the right side. The mean home price (population) should be very close to the mean of sampling distribution. In fact, the mean of sampling distribution (with seed 1234) is 180955.9, whereas the population mean is 180796.1.

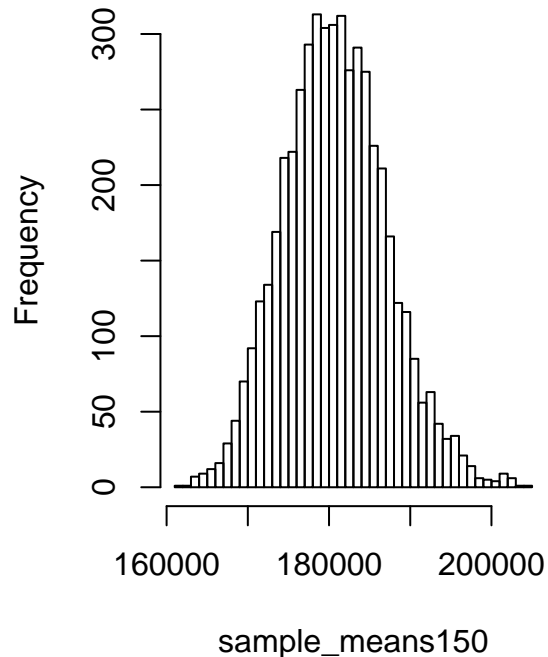
- Change your sample size from 50 to 150, then compute the sampling distribution using the same method as above, and store these means in a new vector called `sample_means150`. Describe the shape of this sampling distribution, and compare it to the sampling distribution for a sample size of 50. Based on this sampling distribution, what would you guess to be the mean sale price of homes in Ames?

```
sample_means150 <- simulation(price, iter = 5000, size = 150, seed = 1234)
par(mfrow = c(1, 2))
hist(sample_means50, breaks = 50, main = "sampling distribution - size 50")
hist(sample_means150, breaks = 50, main = "sampling distribution - size 150")
```

sampling distribution – size 50



sampling distribution – size 150



Both sampling distributions look highly similar to each other and both resemble a normal distribution by looking at their center, spread. Based on either of the sampling distributions, the population mean would most likely fall around 181000.

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

```
sd(sample_means50)
```

```
## [1] 11332.55
```

```
sd(sample_means150)
```

```
## [1] 6363.008
```

The sample_means150 (with larger sample size) has a smaller spread as shown by sd from above calculation. Choosing larger sample size for each sampling would get us closer to the true population mean. We prefer a sampling distribution with a smaller spread because that would give us more confidence in finding a closer range for estimating a population parameter.

This is a product of OpenIntro that is released under a Creative Commons Attribution-ShareAlike 3.0 Unported. This lab was written for OpenIntro by Andrew Bray and Mine Çetinkaya-Rundel.