

Day 5: Function, Loops, and Monte Carlo Simulations

Qingyin Cai

Department of Applied Economics
University of Minnesota

WebR Status

 Ready!

Introduction

- In the previous lecture, we learned how to do regression analysis using R, which is a fundamental skill for econometric analysis.
- Today, we will learn how to code Monte Carlo simulations in R. Monte Carlo simulations are a very important tool in learning econometrics and statistics. With Monte Carlo simulations, you can test any kind of statistical theory or property, which is very useful!!

Note

Before we start the Monte Carlo simulation, let's review a few key R concepts:

- `for loop` function: how they work and when to use them.
- Writing your own functions: we won't need this for today's simulation, but it's an important skill to practice for future coding tasks.

🎯 Learning Objectives

- to be able to write code for your own R functions.
- to be able to write code for a Monte Carlo simulation using the loop function.

✳ Reference

- [Section 19 Functions](#)
- [Section 20 Iteration](#) in R for Data Science

☰ Today's Outline

- User-Defined Functions
- Loops
 - Basics
 - How to Save the Loop Output
 - Multiple Outputs
- Introduction to Monte Carlo Simulations
 - Demonstration
- Overall Review

User-Defined Functions

Introduction to User-Defined Functions

Intro

You can define your own functions. The beauty of creating your own functions is that you can reuse the same code over and over again without having to rewrite it. A function is more useful when the task is longer and more complicated.

Example Situations

- When you want to automate the process of data cleaning.
- When you do a complicated simulation or resampling methods, such as bootstrapping or Monte Carlo simulations.

Introduction to User-Defined Functions

Basics

You can define your own functions using the `function()` function.

General Syntax

```
1 function_name <- function(arg1, arg2, ...){  
2   code to be executed  
3  
4   return(output)  
5 }
```

Note

- You need to define the function name (`function_name`), what kind of inputs the function takes (`arg1`, `arg2`, etc.), and how the function processes using the given input objects.
- The `return()` function is used to return the output of the function. By default, the output defined in the last line of the function is returned.

Introduction to User-Defined Functions

Examples

```
1 function_name <- function(arg1, arg2, ...){  
2   code to be executed  
3  
4   return(output)  
5 }
```

1. A simple function

▶ Run Code



```
1 # --- Define a function --- #  
2 sum_of_two_numbers <- function(x, y){  
3   x + y  
4 }  
5  
6 # --- Use the function --- #  
7 sum_of_two_numbers(x = 3, y = 4)
```

2. A function with multiple outputs

▶ Run Code



```
1 # --- Define a function --- #  
2 mean_sd <- function(vec_x){  
3   c(mean = mean(vec_x), sd = sd(vec_x))  
4 }  
5  
6 # --- Use the function --- #
```



```
7  vec_test <- c(3, 5, 7, 8)
8  mean_sd(vec_x = vec_test)
```

Introduction to User-Defined Functions

Default Value

You can set default values for function arguments by `argument = value`.

Example:

▶ Run Code



```
1 # --- Function to convert a Acre of Hectare --- #
2 acre_to_ha <- function(x = 1){
3   x * 0.404686
4 }
5
6 # --- Use the function --- #
7 acre_to_ha() # because x = 1 by default
8
9 acre_to_ha(x = 1)
```

Load Functions from a File

If you have multiple functions or a long function, you might want to save the function in a separate file and load it when you need it.

For example:

- Save the function code file in a `.R` file (`.Rmd`, etc.).
 - Click `File` -> `Save As...`
 - Choose a name, e.g., `my_function.R`.
- Load the function using the `source()` function.

Let's practice this on your Rstudio!

In-class Exercise

Exercise 1

Questions

1. Write a function (you can name it whatever you want) to calculate the area of a circle with a given radius. The function should return the area of the circle. (Hint1: $\text{Area} = \pi \times r^2$; Hint2: Use `pi`, which is a built-in constant for the value of π in R.)
2. Write a function to count the number of NA values in a given vector. (Hint: use the `is.na()` function.)
3. Write a function called `mean_no_na` that calculates the mean (average) of a numeric vector, but ignores any NA values.

▶ Run Code



```
1 # You can write your code here.
```

In-class Exercise

Exercise 1

Answers

1. Write a function (you can name it whatever you want) to calculate the area of a circle with a given radius. The function should return the area of the circle. (Hint1: $\text{Area} = \pi \times r^2$; Hint2: Use `pi`, which is a built-in constant for the value of π in R.)
2. Write a function to count the number of NA values in a given vector. (Hint: use the `is.na()` function.)
3. Write a function called `mean_no_na` that calculates the mean (average) of a numeric vector, but ignores any NA values.

Run Code



```
1 # Question 1
2 circle_area <- function(radius) {
3   area <- pi * radius^2
4   return(area)
5 }
6
7 circle_area(3)
8
9 # Question 2
10 count_na <- function(x) {
11   na_count <- sum(is.na(x))
12   return(na_count)
13 }
```

```
14
15 count_na(c(1, 2, NA, 4, NA, 6))
16
17 # Question 3
18 mean_no_na <- function(x) {
19   avg <- mean(x, na.rm = TRUE)
20   return(avg)
21 }
22
23
24 mean_no_na(c(1, 2, 3, NA, 5))
```

In-class Exercise

Exercise 2 (optional)

Questions

You're a data expert at a store chain. The company needs to study its monthly sales growth to plan better. They expect sales to grow by a fixed percentage each month. Your job is to create an R function that shows sales growth over a year.

For sales growth, use the following formula:

$$S_t = S_0 \times (1 + g)^{t-1}$$

, where S_t is the sales in month t , S_0 is the starting sales, and g is the growth rate.

Create a function called `monthly_sales_growth` with the following three inputs:

- `initial_sales`: Starting sales (in thousands of dollars).
- `growth_rate`: Monthly growth rate (as a decimal, like 0.03 for 3% growth).
- `months`: How many months to predict (usually 12 for a year).

The function should give back a vector of numbers (or it would be even better if you could show in a `data.frame` or `data.table` in which two columns, e.g., month and sales, show the expected sales for each month.)

Run Code



```
1 # You can write your code here.
```

In-class Exercise

Exercise 2 (optional)

Answers

You're a data expert at a store chain. The company needs to study its monthly sales growth to plan better. They expect sales to grow by a fixed percentage each month. Your job is to create an R function that shows sales growth over a year.

For sales growth, use the following formula:

$$S_t = S_0 \times (1 + g)^{t-1}$$

, where S_t is the sales in month t , S_0 is the starting sales, and g is the growth rate.

Create a function called `monthly_sales_growth` with the following three inputs:

- `initial_sales`: Starting sales (in thousands of dollars).
- `growth_rate`: Monthly growth rate (as a decimal, like 0.03 for 3% growth).
- `months`: How many months to predict (usually 12 for a year).

The function should give back a vector of numbers (or it would be even better if you could show in a `data.frame` or `data.table` in which two columns, e.g., month and sales, show the expected sales for each month.)

Run Code



```
1 monthly_sales_growth <- function(initial_sales, growth_rate, months = 12) {  
2   # months: positive integer (e.g., 12 for a year)  
3   t <- seq_len(months) # a vector, or t <- 1: months  
4   sales <- initial_sales * (1 + growth_rate)^(t - 1)
```



```
5   return(sales)
6 }
7
8 # or using data.table/data.frame
9
10 monthly_sales_growth <- function(initial_sales, growth_rate, months = 12) {
11   # months: positive integer (e.g., 12 for a year)
12   t <- seq_len(months) # a vector, or t <- 1: months
13   sales <- initial_sales * (1 + growth_rate)^(t - 1)
14   data.frame(month = t, sales = sales) # or data.table(month = t, sales = sales)
15 }
16
17
18 # Example: $200k starting sales, 3% monthly growth, 12 months
19 monthly_sales_growth(initial_sales = 200, growth_rate = 0.03, months = 12)
```

Loops

Why loop?

Using Loop is useful when you want to repeat the same task (but with a slight change in parameters) over and over again.

Examples

- Downloading the data from the web iteratively.
 - When you want to download the ag-production data from USDA-NASS, you are limited to download 50,000 records per query. You need to repeatedly download the data until you get all the data you need.
 - USDA crop scale data, NOAA weather data, etc.
- Loading multiple data files in a folder.
- Running the same regression analysis for multiple datasets.
- Running simulations or resampling methods, such as bootstrapping or Monte Carlo simulations.

While there are several looping commands in R (e.g., `foreach`, `lapply`, etc.), we will use the `for loop` function, as it is the most basic and widely used looping function in R.

For loops

Basics

The `for loop` function is a built-in R function. The syntax of the `for loop` is very simple.

Syntax:

```
1 for (name in collection) {  
2   # code to run for each element  
3 }
```

Components

- `name` (loop variable): a placeholder that takes on one value from the collection each time the loop runs.
- `collection`: the set of values you want to loop over (usually a vector or list).
- `code block`: the instructions you want R to execute at each step.

Note

- On each iteration, the loop variable (`name`) takes the **next value** from the collection.
- The collection can be:
 - A numeric sequence (e.g., `1:5`)
 - A character vector (e.g., `c("apple", "banana")`)
 - A list of objects (e.g., datasets)

For loops

Examples

1. Print the numbers from 1 to 5.

Loop variable `i` takes each number in the sequence `1:5` in order and print the value of `i`.

▶ Run Code

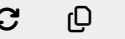


```
1 for (i in 1:5){  
2   print(i)  
3 }
```

2. Print characters in a list.

Loop variable `x` takes each character in the list `list("I", "like", "cats")` in order and print the value of `x`.

▶ Run Code



```
1 for (x in list("I", "like", "cats")){  
2   print(x)  
3 }
```

3. Calculate the mean of each element in a list.

Can you tell me what's going on in the following code?

▶ Run Code



```
1 ls_seq_num <- list(1:3, 4:6, 7:9)  
2  
3 for (seq_num in ls_seq_num){  
4   print(mean(seq_num))  
}
```


In-class Exercise

Exercise 1

Questions

In econometric class, we use the `rnorm()` function a lot! It is a function that generates random numbers from a normal distribution. See `?rnorm` for more details.

The basic syntax is `rnorm(n, mean = 0, sd = 1)`, where `n` is the number of random numbers you want to generate, `mean` is the mean of the normal distribution, and `sd` is the standard deviation of the normal distribution. So `rnorm(n, mean = 0, sd = 1)` generates `n` random numbers from a standard normal distribution.

Generate 1000 random numbers from a standard normal distribution and calculate the mean the numbers (use `mean()` function), and print the results. Repeat this process 10 times using the `for` loop.

Run Code



```
1 for (i in 1:10){  
2   # write your code here  
3  
4 }
```

In-class Exercise

Exercise 1

Answers

In econometric class, we use the `rnorm()` function a lot! It is a function that generates random numbers from a normal distribution. See `?rnorm` for more details.

The basic syntax is `rnorm(n, mean = 0, sd = 1)`, where `n` is the number of random numbers you want to generate, `mean` is the mean of the normal distribution, and `sd` is the standard deviation of the normal distribution. So `rnorm(n, mean = 0, sd = 1)` generates `n` random numbers from a standard normal distribution.

Generate 1000 random numbers from a standard normal distribution and calculate the mean the numbers (use `mean()` function), and print the results. Repeat this process 10 times using the `for` loop.

Run Code



```
1 # Set seed for reproducibility
2 set.seed(123)
3
4 # Repeat the process 10 times
5 for (i in 1:10) {
6   # Step 1: Generate 1000 random numbers from standard normal
7   x <- rnorm(1000, mean = 0, sd = 1)
8
9   # Step 2: Calculate the sample mean
10  m <- mean(x)
```

```
11
```



```
11  
12     # Step 3: Print the result  
13     print(m)  
14 }
```

In-class Exercise

[Exercise 2 \(nested loop\)](#)

[Questions](#)

You can nest the `for` loop inside another `for` loop. For example,

```
Run Code
```

```
1 # Outer loop
2 for (i in 1:3) {
3   # Inner loop
4   for (j in 1:2) {
5     print(paste("i =", i, "j =", j))
6   }
7 }
```

Using the above code as a reference, fill in the following empty 3 x 3 matrix with the sum of the row and column indices.

The output should look like this:

	[,1]	[,2]	[,3]
[1,]	2	3	4
[2,]	3	4	5
[3,]	4	5	6

```
Run Code
```

```
1 # Here, is the empty 3 x 3 matrix.
2 empty_matrix <- matrix(NA, nrow = 3, ncol = 3)
```

In-class Exercise

[Exercise 2 \(nested loop\)](#)

[Answers](#)

You can nest the `for` loop inside another `for` loop. For example,

▶ Run Code



```
1 # Outer loop
2 for (i in 1:3) {
3   # Inner loop
4   for (j in 1:2) {
5     print(paste("i =", i, "j =", j))
6   }
7 }
```

Using the above code as a reference, fill in the following empty 3 x 3 matrix with the sum of the row and column indices.

The output should look like this:

	[,1]	[,2]	[,3]
[1,]	2	3	4
[2,]	3	4	5
[3,]	4	5	6

▶ Run Code



```
1 # Here, is the empty 3 x 3 matrix.
2 empty_matrix <- matrix(NA, nrow = 3, ncol = 3)
3 for (i in 1:3){
```

```
4   for (j in 1:3){  
5       empty_matrix[i, j] <- i + j  
6   }  
7 }  
8 print(empty_matrix)
```

For loops: How to Save the loop output?

Introduction

Unlike R functions we have seen so far, `for loop` does not have a return value. It just iterates the process we defined in the loop.

Let's do some experiments:

Experiment 1

▶ Run Code



```
1 ls_seq_num <- list(1:3, 4:6, 7:9)
2
3 for (seq_num in ls_seq_num){
4   # save the results in x
5   x <- mean(seq_num)
6 }
7
8 # What is the value of x?
9 x
```

Every round of the loop, the variables defined inside the loop are updated.

For loops: How to Save the loop output?

Introduction

Unlike R functions we have seen so far, `for` loop does not have a return value. It just iterates the process we defined in the loop.

Let's do some experiments:

Experiment 2

▶ Run Code



```
1 ls_seq_num <- list(1:3, 4:6, 7:9)
2
3 x <- for (seq_num in ls_seq_num){
4   # save the results in x
5   mean(seq_num)
6 }
7
8 # What is the value of x?
9 x
```

- You cannot assign loop to a variable directly (e.g `x <- for (i in 1:3){print(i)}` does not work).

For loops: How to Save the loop output?

Basics

- To save the results of the loop, you need to create an empty object before the loop and save the output in the object in each iteration. (You did this in the exercise 2!)
 - The object can be a vector, a list, a matrix, or a data frame (or data.table), depending on the type of the output you want to save.

Example

Suppose you want to cube each number in the sequence `1:5`.

```
Run Code  ↺  📄  
1  # --- Create an empty vector --- #  
2  output_storage <- rep(0, 5)  
3  
4  # --- for loop --- #  
5  for (i in 1:5){  
6    output_storage[i] <- i^3  
7  }  
8  
9  output_storage
```

Note

- Since the output of each iteration is a number, vector is a good choice for the storage object. (Alternatively you can use a list object.)

Multiple Outputs

What if we want to have multiple outputs from the loop and combine them into a single dataset?

Example

Let's generate 100 random numbers from a standard normal distribution and calculate the mean and the standard deviation of numbers. Repeat this process 10 times using the `for loop` and save the results in a dataset.

Run Code



```
1  # --- Number of iterations --- #
2  B <- 10
3
4  # --- Create an empty data.table (or data.frame for matrix) --- #
5  output_storage <- data.table(mean = rep(0, B), sd = rep(0, B))
6
7  # --- for loop --- #
8  for(i in 1:B){
9    # --- Generate random numbers --- #
10   random_numbers <- rnorm(100, mean = 0, sd = 1)
11
12   # --- Calculate the mean and the standard deviation --- #
13   output_storage[i, mean := mean(random_numbers)]
14   output_storage[i, sd := sd(random_numbers)]
15 }
16
17 #or output_storage[i, `:=`(
18 #  mean = mean(random_numbers),
19 #  sd   = sd(random_numbers)
```


20 #)]

21

22 output_storage

In-class Exercise

Exercise 1

Questions

- 1. Using the `for loop`, calculate the sum of the first `n` numbers for `n = 1, 2, ..., 10`. For example, the sum of the first 3 numbers (`n=3`) is `1 + 2 + 3 = 6`. Save the results in a vector object.
- 2. Fibonacci sequence is a series of numbers in which each number is the sum of the two preceding ones (e.g. 0, 1, 1, 2, 3, 5, 8, 13, 21, 34, ...). Write a function that generates the first `n` numbers in the Fibonacci sequence. (You use the `for loop` function inside the function.) For example, when `n = 5`, the function should return `c(0, 1, 1, 2, 3)`. For simplicity, let's assume that (You don't need to consider the case where `n = 1` or `2`).

$$n \geq 3$$

▶ Run Code

1 # You can write your code here

In-class Exercise

Exercise 1

Answers

1. Using the `for loop`, calculate the sum of the first `n` numbers for `n = 1, 2, ..., 10`. For example, the sum of the first 3 numbers (`n=3`) is $1 + 2 + 3 = 6$. Save the results in a vector object.
2. Fibonacci sequence is a series of numbers in which each number is the sum of the two preceding ones (e.g. 0, 1, 1, 2, 3, 5, 8, 13, 21, 34, ...). Write a function that generates the first `n` numbers in the Fibonacci sequence. (You use the `for loop` function inside the function.) For example, when `n = 5`, the function should return `c(0, 1, 1, 2, 3)`. For simplicity, let's assume that (You don't need to consider the case where `n = 1` or 2).
 $n \geq 3$

Run Code



```
1 #-----Part1-----#
2 sums <- rep(0, 10) # or sums <- numeric(10)
3
4 # Loop over n = 1 to 10
5 for (n in 1:10) {
6   sums[n] <- sum(1:n)
7 }
8
9 # Print results
10 sums
11
```

```
12 #-----Part2-----#
13 fibonacci <- function(n) {
14   # Create a vector to store the sequence
15   fib <- rep(0, n)
16
17   # First two numbers of Fibonacci sequence
18   fib[1] <- 0
19   fib[2] <- 1
20
21   # Generate the rest using a for loop
22   for (i in 3:n) {
23     fib[i] <- fib[i - 1] + fib[i - 2]
24   }
25
26   return(fib)
27 }
28
29 # Example: first 5 Fibonacci numbers
30 fibonacci(5)
```

In-class Exercise

Exercise 2 (optional)

Questions

NOTE: It's okay if you cannot solve this problem! I will show two approaches to solve this problem, and do speed comparison between the two approaches.

Pythagorean triples are sets of three **positive integers** that satisfy the equation $a^2 + b^2 = c^2$, who are named after the ancient Greek mathematician Pythagoras.

Let's take this concept further. Suppose Pythagoras challenges you to find all possible Pythagorean triples where a and b are less than or equal to a given number n . To address this problem, let's create an R function that will produce all such triples.

1. Write a function that takes one argument n , an integer, representing the maximum value for a and b . The function should return a data frame with columns a , b , and c , containing all Pythagorean triples where $a \leq n$ and $b \leq n$ and $a^2 + b^2 = c^2$.

Hints:

- Consider using nested loops to iterate through all possible values of a and b up to n .
- Use the `sqrt()` function to calculate the potential value of c , and check if it's an integer.
- Use the `floor()` function to round down the value of c .

▶ Run Code

1 # You can write your code here.

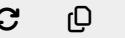
Reference: [Pythagorean Triples](#)

In-class Exercise

Exercise 2 (optional)

Answers

Run Code



```
1  /*-----*/
2  #' ## Approach 1 (loops)
3  /*-----*/
4  find_pythagorean_triples <- function(n){
5    # --- prepare an empty storages --- #
6    triples_storage_dt <-
7      data.frame(a = integer(), b = integer(), c = integer())
8
9    for(a in 1:n){
10     # for example a = 4
11     for(b in 1:a){
12       # for example b = 3
13       c <- sqrt(a^2 + b^2)
14       if(c == floor(c)){
15         res_triples_dt <- data.frame(a = a, b = b, c = c)
16         triples_storage_dt <- rbind(triples_storage_dt, res_triples_dt)
17       }
18     }
19   }
20   return(triples_storage_dt)
```

```

21 }
22
23 # --- test --- #
24 find_pythagorean_triples(n = 5)
25
26
27 /*-----*/
28 #' ## Approach 2 (Vectorized)
29 /*-----*/
30 library(data.table)
31 library(dplyr)
32
33 find_pythagorean_triples_vec <- function(n){
34   # for example, n = 5
35   # --- create all combinations of (a, b) --- #
36   abc_all_dt <-
37     CJ(a = 1:n, b = 1:n) %>%
38     .[, c := sqrt(a^2 + b^2)] %>%
39     # find a row where c is an integer
40     .[c == floor(c), ]
41
42   return(abc_all_dt)
43 }
44
45 # NOTE: #CJ() is a function from data.table package that creates all combinations of two vectors. (It is
46   similar to expand.grid() function in base R.)
47
48 # --- test --- #

```

```
48 find_pythagorean_triples_vec(n = 5)
49
50 #/*-----*/
51 #' ## Speed Competition
52 #/*-----*/
53 system.time(find_pythagorean_triples(n = 10^4))
54 system.time(find_pythagorean_triples_vec(n = 10^4))
55
56 # Lesson: Vectorized operations are much faster than loops in R!
```


Check Point

Up to this point, as long as you understand the following points, you are good to go!

- You know how to use `function()` to define a simple function yourself.
- You know how to use `for loop` (i.e., syntax, which argument you need to define).
- You know that you need to prepare an empty object to save the output of the loop.

Introduction to Monte Carlo Simulations

Introduction to Monte Carlo Simulations

Analogy

- Imagine you want to know the chance of winning a dice game. Instead of trying to calculate the exact probability with formulas, you could just roll the dice thousands of times, record the outcomes, and then see how often you win.
- That's basically what Monte Carlo Simulation does: it uses repeated random experiments to approximate answers to problems that are hard to solve analytically.
- Key ideas:
 - We model uncertainty with randomness.
 - We simulate many scenarios (sometimes thousands or millions).
 - We average the results to estimate probabilities, risks, or outcomes.

Introduction to Monte Carlo Simulations

What is it?

Monte Carlo simulation is a technique to approximate a likelihood of possible outcomes (e.g., predictions, estimates) from a model by iteratively running the model on artificially created datasets. In every iteration, the datasets are randomly sampled from the assumed data generating process, so it varies every iteration.

- The incorporation of randomness in the simulation is the key feature of the Monte Carlo simulation. *It is mimicking the randomness of real-world phenomena.*

Introduction to Monte Carlo Simulations

Monte Carlo simulation in Econometrics

In econometrics, the Monte Carlo simulation is used to evaluate the performance of a statistical procedure or the validity of theories in a realistic setting.

For example

Suppose that a researcher came up with a new estimator to estimate the coefficients of a regression model.

- An estimator (e.g, sample mean, standard error, OLS) is a function of a random variable, therefore it is also a random variable.
- A random variable has its own probability distribution.
- So, to understand the performance of the estimator (e.g., unbiasedness and efficiency), we need to examine the properties of the probability distribution of the estimator.
- We use Monte Carlo simulation to approximate the probability distribution of the estimator!

Example: Binomial Distribution

Think about the following example.

Example

- Suppose that we flip a coin $n = 10$ times and count the number of heads. Let's denote the number of heads X .
- The coin is not fair, however. The probability of getting a head is $p = \text{Pr}[\text{heads}] = 1/3$.
- Suppose that you repeat this experiment 1000 times. What is the mean and the variance of X ?

This kind of experiment is modeled by the binomial distribution. According to the theory, it is predicted that

- Mean of X is $E[X] = np = 10 \times 1/3 = 3.33$
- Variance of X is $\text{Var}[X] = np(1 - p) = 10 \times 1/3 \times 2/3 = 2.22$

Is it true? Let's check this using a Monte Carlo simulation!

Monte Carlo Simulation: Steps

step 1: Specify the data generating process.

- You need to pick a specific probability distribution to generate a random number.

step 2: Repeat:

- step 2.1: generate a (pseudo) random sample data based on the data generating process.
- step 2.2: get an outcome you are interested in based on the generated data.

step 3: compare your estimates with the true parameter

Demonstration: Binomial Distribution

A Single Iteration

Let's start writing code for a single iteration to get an idea of the Monte Carlo simulation process in R.

- We want to repeat this 1000 times.

Run Code



```
1  # For the reproducibility purpose, set the seed.
2  set.seed(1843)
3
4  # === Step 1: Specify the Data Generating Process === #
5  # parameters
6  p <- 1/3 # the probability of getting a head
7  n <- 10 # the number of trials
8
9  # === Step 2.1: Generate Data === #
10 # Here I randomly picked 1 or 0 for n = 10 times, with the probability of p = 1/3 for 1.
11 seq_x <- sample(c(1,0), size=n, prob=c(p, 1-p), replace=TRUE) #replace = TRUE means sampling with
    replacement
12
13 # === Step 2.2: Get an Outcome You are Interested in === #
14 # count the number of heads
15 sum(seq_x)
```


Demonstration: Binomial Distribution

Multiple Iterations

Run Code



```
1  # For the reproducibility purpose, set the seed.
2  set.seed(1843)
3
4  # === Step 1: Specify the Data Generating process === #
5  # parameters
6  p <- 1/3 # the probability of getting a head
7  n <- 10 # the number of trials
8  B <- 1000 # the number of iterations
9
10 # Don't forget to create an empty object to save the output!
11 X <- rep(0, B)
12
13 # === Step 2: Iteration === #
14 for (i in 1:B){
15   # --- Step 2.1: Generate Data --- #
16   # Here I randomly picked 1 or 0 for n = 10 times, with the probability of p = 1/3 for 1.
17   seq_x <- sample(c(1,0), size=n, prob=c(p, 1-p), replace=TRUE)
18
19   # --- Step 2.2: Get an Outcome You are Interested in --- #
20   # count the number of heads and save the result
21   X[i] <- sum(seq_x)
22 }
```

```
23
24 # === Step 3 === #
25 # Compute the mean of X
26 mean(X)
27
28 # Compute the variance of X
29 var(X)
```

In-class Exercise Problem

Exercise Problem: (Weak) Law of Large Number

Questions

Weak law of large number states that the sample mean converges (in probability) to the population mean as the sample size increases. In other words, the sample mean more accurately estimates the population mean as the sample size increases.

$$\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i \xrightarrow{p} E[X]$$

Lets check this using Monte Carlo simulation! We compare the distribution of sample mean with different sample size. Let's compare two sample sizes: and .

$n = 100$ $n = 1000$

Process

Repeat 1 and 2 for times.

$B = 1000$

1. Using a normal distribution with mean and , generate random numbers for and . e.g. `rnorm(n = 10, mean = 5, sd = 10)`.

2. Compute sample mean for each sample data, and save them.

Finally,

3. Plot histograms of the sample means obtained from the two samples.

Exercise Problem: (Weak) Law of Large Number

Answers

Run Code



```
1  set.seed(1843)
2
3  # --- Parameters --- #
4  mu <- 5
5  sd <- 10
6
7  B <- 1000
8
9  # --- Create an empty object to save the output --- #
10 sample_mean_100 <- rep(0, B)
11 sample_mean_1000 <- rep(0, B)
12
13 # --- for loop --- #
14 for (i in 1:B){
15   # --- Generate random numbers --- #
16   sample_100 <- rnorm(n = 100, mean = mu, sd = sd)
17   sample_1000 <- rnorm(n = 1000, mean = mu, sd = sd)
18
19   # --- Compute sample mean --- #
20   sample_mean_100[i] <- mean(sample_100)
21   sample_mean_1000[i] <- mean(sample_1000)
22 }
```

```
23
24
25 # --- Plot --- #
26 ggplot() +
27   geom_histogram(
28     aes(x = sample_mean_100, fill = "Sample size: 100"),
29     alpha = 0.7
30   ) +
31   geom_histogram(
32     aes(x = sample_mean_1000, fill = "Sample size: 1000"),
33     alpha = 0.7
34   ) +
35   geom_vline(
36     xintercept = mu,
37     linetype = "dashed",
38   ) +
39   theme_bw()
```

Exercise Problem: Two estimators to estimate the population mean? (optional)

Questions

Suppose you're interested in estimating the unknown population mean of men's heights (i.e., μ) in the US. We have randomly sampled data with the size of $n = 1000$. Let X_i denote the individual i 's height in the sample. **How should we use the sample data to estimate the population mean?**

Your friends suggested two different estimators:

Estimator 1. Use the sample mean:

$$\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$$

Estimator 2. Use only the first observation:

$$X_1$$

Theoretically, both estimators are unbiased estimators (i.e., if we repeat the estimation process many times on a different sample data every time, the average of the estimates will be equal to the population mean):

$$E[\bar{X}_n] = E\left[\frac{1}{n} \sum_{i=1}^n X_i\right] = \frac{1}{n} E\left[\sum_{i=1}^n X_i\right] = \frac{1}{n} \sum_{i=1}^n E[X_i] = \frac{1}{n} \cdot n \cdot \mu = \mu$$

$$E[X_1] = \mu$$

Questions:

- Is it true that both estimators are correctly estimating the population mean, on average?
- Which one is more accurate in estimating the population mean?

Using Monte Carlo simulation, let's examine these questions!

Exercise Problem: Two estimators to estimate the population mean? (optional)

Hint

1. Repeat the following processes 1000 times:

step 1. Draw random numbers with known mean μ and standard deviation σ . This will be the sample data.

step 2. Get (1) $n = 1000$ the mean of the sample and (2) the value of the first observation, and save the results.

1. The previous iterations produce 1000 estimates of the population mean for estimator 1 and estimator 2, respectively. Compute the means for each estimator. Are they both close to the true population mean? Compute the variance of the estimates. Which one has a smaller variance?

If you could also visually compare the distribution of estimates from the two estimators, that would be great!

▶ Run Code

1

Exercise Problem: Two estimators to estimate the population mean? (optional)

Answers

Run Code



```
1  # --- Parameter setting --- #
2  B = 1000
3  n = 1000
4  mu = 170 #(cm)
5
6  # --- Prepare a storage object --- #
7  res_data <-
8    data.table(
9      x_bar = rep(0, B),
10     x_single = rep(0, B)
11   )
12
13 # --- Monte Carlo simulation --- #
14 for(i in 1:B){
15
16   # --- Create sample data --- #
17   data <- rchisq(n, df=mu)  #rchisq() is a random number generator for the chi-squared distribution.
18
19   # --- Compute the sample mean --- #
20   x_bar <- mean(data)
21   # --- Get the 1st value --- #
22   x_single <- data[1]
```

```
23
24 # --- store the results --- #
25 res_data[i, "x_bar"] <- x_bar
26 res_data[i, "x_single"] <- x_single
27 }
28
29 # === Visualization === #
30 # Just for convenience, transform the data to long-format (you don't need to do this though)
31 vis_sampling <-
32   ggplot(res_data) +
33     geom_histogram(aes(x=x_bar, fill="x_bar"), alpha = 0.5)+
34     geom_histogram(aes(x=x_single, fill="x_single",), alpha = 0.5)+
35     labs(title="Sample size 1000")+
36     theme_bw() +
37     # modify x-label and legend name
38     labs(x = expression ( $\hat{\beta}$ ), fill = "Estimator")
39
40 vis_sampling
```

Appendix

foreach Function

Basics

The `foreach` function is a function of the `foreach` package. It is used to iterate the same process over and over again, similar to the `for loop` function.

Basic Syntax

While there are some differences, the basic syntax of the `foreach` function is pretty much similar to the `for loop` function.

```
1 foreach(variable = collection_of_objects) %do% {  
2   the code to be executed in each iteration  
3  
4   return(output)  
5 }
```

Note

- Differences between `for loop` and `foreach` function:
 - use `=` instead of `in`.
 - You need to use `%do%` operator.
 - `foreach` function has a return value, while `for loop` does not. By default, the output is returned as a list.
- `foreach` function also supports parallel processing. (we will not cover this in this class.)

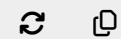
foreach Function

Example

Using the `for` loop and `foreach` function, let's calculate the square of the numbers from 1 to 10, respectively.

foreach

▶ Run Code



```
1 # --- Load the package --- #
2 library(foreach)
3
4 results <- foreach(i = 1:10) %do% {
5   i_square <- i^2
6   return(i_square)
7 }
8
9 # see the output
10 results
```

for loop

▶ Run Code



```
1 # --- Create an empty object --- #
2 results <- rep(0, 10)
3
4 # --- for loop --- #
5 for (i in 1:10){
6   results[i] <- i^2
7 }
8
9 # see the output
10 results
```

foreach Function

Change the Output Format

By default, `foreach` function returns each iteration's result as a list. But you can choose the format of the output by using the `.combine` argument.

- `.combine = c` combines each iteration's result as a vector (like `c()` to create a vector).
- `.combine = rbind` combines each iteration's result by row.
- `.combine = cbind` combines each iteration's result by column.

The last two options are used when the output is a `matrix` or a `data.frame` (or `data.table`).

Example

- Try different `.combine` options in the following code.

Run Code



```
1 # --- Load the package --- #
2 library(foreach)
3
4 results <- foreach(i = 1:10, .combine = c) %do% {
5   i_square <- i^2
6   return(i_square)
7 }
8
9 # see the output
```


Overall Review

Review

Basics of R

- Coding rules & projects: file paths, working directories, using `.Rproj`.
- Data types: numeric, character, logical; handling `NA`.
- Core structures: vectors, matrices, data frames — creation, indexing, subsetting.
- Base math & objects: assignment, arithmetic, functions; saving/loading data.

Review

Data Manipulation with `data.table`

- `DT[i, j, by]` pattern and chaining.
- Row filtering and column selection; create/modify columns with `:=`.
- Grouped summaries with `by`; ordering and efficient aggregation.
- Reshaping with `melt()` and `dcast()` (wide ↔ long).

Review

Visualization with ggplot2

- Grammar of Graphics: `data` + `aes()` + `geom_*`.
- Common plots: scatter, line, bar, histogram, box, density; faceting.
- Aesthetics: color/size/shape, grouping, collective geoms.
- Labels/themes: axes, titles, legends; layering multiple datasets.

Review

Regression & Reporting

- Linear models with `lm()`.
- `modelsummary` for publish-ready regression tables; `datasummary` for descriptives.
- Customization: model labels, stats formatting, notes.
- Quarto reporting: render to HTML/PDF for reproducible outputs.

Review

Functions, Loops, & Monte Carlo

- User-defined functions `function()`: arguments, returns, defaults; save in .R and `source()`.
- Loops function `for()`: basics, saving outputs, multiple outputs.
- Randomness & reproducibility: `set.seed()`, `rnorm()`, `sample(..., replace=TRUE)`.

Review Exercise

Instructions

This exercise integrates all major concepts from our R course. We will use the built-in R datasets `mtcars` provided.

Setup (Run this first!)

▶ Run Code



```
1 # Load required packages
2 library(data.table)
3 library(ggplot2)
4 library(modelsummary)
5
6 # Load and convert datasets
7 data(mtcars)
8 data(airquality)
9 mtcars_dt <- as.data.table(mtcars)
10 air_dt <- as.data.table(airquality)
11
12 # Add car names to mtcars
13 mtcars_dt[, car_name := rownames(mtcars)]
```

Review Exercise

Questions

Data Structures and Basic Operations

- 1. Show the structure of mtcars_dt
- 2. Calculate the mean mpg (miles per gallon) for all cars
- 3. Create a logical vector showing which cars have both mpg > 20 AND hp < 150
- 4. Extract the car names of cars with the highest mpg

▶ Run Code

↺

📄

1

#Write your code here

Review Exercise

Questions

data.table Operations

1. Create a new column 'power_ratio' = hp/wt (horsepower per weight) and add it to mtcars_dt
2. Replace all mpg values less than 15 with NA
3. Filter mtcars_dt to show only cars with automatic transmission (am == 0)
4. Sort mtcars_dt by mpg in descending order and show top 5 cars
5. Calculate average mpg by number of cylinders (cyl)
6. For each transmission type (am), calculate: count of cars, mean mpg, mean hp

▶ Run Code



```
1 #Write your code here
```

Review Exercise

Questions

Data Visualization

- Create a scatter plot of wt vs mpg, with:
 - Different colors for transmission type (am).
 - Different shapes for number of cylinders (cyl).
 - A smooth trend line for each transmission type.

▶ Run Code

↺

📄

1

#Write your code here

Review Exercise

Questions

Regression Analysis

- 1. Run a regression: `mpg ~ hp + wt + factor(cyl)`
- 2. Extract the coefficient for `hp` from your model
- 3. Run a second model: `mpg ~ hp + wt + factor(am)`, and create a comparison table using `modelsummary()` for both models

▶ Run Code

1 #Write your code here

Review Exercise

Questions

Programming

1. Write a function that calculates mpg per cylinder (mpg/cyl)

▶ Run Code



```
1 calc_mpg_per_cyl <- function(mpg_val, cyl_val) {  
2   # Your code here  
3  
4 }  
5  
6 # Test: calc_mpg_per_cyl(21, 6) should return 3.5
```

2. Use a for loop to calculate the mean mpg for cars with 4, 6, and 8 cylinders. Store results in a vector called 'cyl_means'

▶ Run Code



```
1 cyl_values <- c(4, 6, 8)  
2 cyl_means <- rep(0, 3)  
3  
4 # write your code here
```

Review Exercise

Answers

Data Structures and Basic Operations

1. Show the structure of mtcars_dt
2. Calculate the mean mpg (miles per gallon) for all cars
3. Find how many cars have mpg > 20

▶ Run Code



```
1 #-----part 1-----#
2 str(mtcars_dt) # or ?mtcars_dt
3
4 #-----part 2-----#
5 mean(mtcars_dt$mpg)
6
7 #-----part 3-----#
8 sum(mtcars_dt$mpg > 20)
```

Review Exercise

Answers

data.table Operations

1. Create a new column 'power_ratio' = hp/wt (horsepower per weight) and add it to mtcars_dt
2. Filter mtcars_dt to show only cars with automatic transmission (am == 0)
3. Sort mtcars_dt by mpg in descending order and show top 5 cars
4. For each transmission type (am), calculate: count of cars, mean mpg, mean hp

Run Code

```
1  #-----part 1-----#
2  mtcars_dt[, power_ratio := hp/wt]
3
4  #-----part 2-----#
5  mtcars_dt[am == 0]
6
7  #-----part 3-----#
8  mtcars_dt[order(-mpg)][1:5]
9
10 #-----part 4-----#
11 mtcars_dt[, .(count = .N, avg_mpg = mean(mpg, na.rm = TRUE), avg_hp = mean(hp)), by = am]
```

Review Exercise

Answers

Data Visualization

- Create a scatter plot of wt vs mpg, with:
 - Different colors for transmission type (am).
 - A smooth trend line for each transmission type.

▶ Run Code



```
1  ggplot(mtcars_dt, aes(x = wt, y = mpg, color = factor(am))) +  
2    geom_point() +  
3    geom_smooth(method = "lm") +  
4    labs(title = "Weight vs MPG by Transmission and Cylinders",  
5          x = "Weight (1000 lbs)",  
6          y = "Miles per Gallon",  
7          color = "Transmission") +  
8    theme_bw() +  
9    theme(legend.position = "right")  
10
```

Review Exercise

Answers

Regression Analysis

1. Run a regression: `mpg ~ hp + wt + factor(cyl)`
2. Extract the coefficient for `hp` from your model
3. Run a second model: `mpg ~ hp + wt + factor(am)`, and create a comparison table using `modelsummary()` for both models

Run Code



```
1 #-----part 1-----#
2 model1 <- lm(mpg ~ hp + wt + factor(cyl), data = mtcars_dt)
3
4 #-----part 2-----#
5 coef(model1)["hp"]
6
7 #-----part 3-----#
8 model2 <- lm(mpg ~ hp + wt + factor(am), data = mtcars_dt)
9 modelsummary(list("Model 1" = model1, "Model 2" = model2))
```


Review Exercise

Answers

Programming

1. Write a function that calculates mpg per cylinder (mpg/cyl)

▶ Run Code



```
1 calc_mpg_per_cyl <- function(mpg_val, cyl_val) {  
2   return(mpg_val / cyl_val)  
3 }  
4  
5 # Test: calc_mpg_per_cyl(21, 6) should return 3.5
```

2. Use a for loop to calculate the mean mpg for cars with 4, 6, and 8 cylinders. Store results in a vector called 'cyl_means'

▶ Run Code



```
1 cyl_values <- c(4, 6, 8)  
2 cyl_means <- rep(0, 3)  
3  
4 for (i in 1:3) {  
5   cyl_means[i] <- mean(mtcars_dt[cyl == cyl_values[i], mpg], na.rm = TRUE)  
6 }
```