

Introduction to R

(Selected topics from R manual and online instructions)

1. R - Data Types

Generally, while doing programming in any programming language, you need to use various variables to store various information. Variables are nothing but reserved memory locations to store values. This means that, when you create a variable you reserve some space in memory.

You may like to store information of various data types like character, wide character, integer, floating point, double floating point, Boolean etc. Based on the data type of a variable, the operating system allocates memory and decides what can be stored in the reserved memory.

In contrast to other programming languages like C and java in R, the variables are not declared as some data type. The variables are assigned with R-Objects and the data type of the R-object becomes the data type of the variable. There are many types of R-objects. The frequently used ones are –

- Vectors
- Lists
- Matrices
- Arrays
- Data Frames

The simplest of these objects is the **vector object** and there are six data types of these atomic vectors, also termed as six classes of vectors. The other R-Objects are built upon the atomic vectors.

In R programming, the very basic data types are the R-objects called **vectors** which hold elements of different classes as shown above. Please note in R the number of classes is not confined to only the above six types. For example, we can use many atomic vectors and create an array whose class will become array.

1.1 Vectors

When you want to create vector with more than one element, you should use **c()** function which means to combine the elements into a vector.

```
# Create a vector.
apple <- c('red','green',"yellow")
print(apple)

# Get the class of the vector.
print(class(apple))
```

When we execute the above code, it produces the following result –

```
[1] "red"      "green"    "yellow"
[1] "character"
```

1.2 Lists

A list is an R-object which can contain many different types of elements inside it like vectors, functions and even another list inside it.

```
# Create a list.
list1 <- list(c(2,5,3),21.3,sin)

# Print the list.
print(list1)
```

When we execute the above code, it produces the following result –

```
[[1]]
[1] 2 5 3

[[2]]
[1] 21.3
```

```
[[3]]
function (x) .Primitive("sin")
```

1.3 Matrices

A matrix is a two-dimensional rectangular data set. It can be created using a vector input to the matrix function.

```
# Create a matrix.
M = matrix( c('a','a','b','c','b','a'), nrow = 2, ncol = 3, byrow = TRUE)
print(M)
```

When we execute the above code, it produces the following result –

```
      [,1] [,2] [,3]
[1,] "a"  "a"  "b"
[2,] "c"  "b"  "a"
```

1.4 Arrays

While matrices are confined to two dimensions, arrays can be of any number of dimensions. The array function takes a dim attribute which creates the required number of dimension. In the below example we create an array with two elements which are 3x3 matrices each.

```
# Create an array.
a <- array(c('green','yellow'),dim = c(3,3,2))
print(a)
```

When we execute the above code, it produces the following result –

```
, , 1
      [,1]      [,2]      [,3]
[1,] "green"   "yellow"   "green"
[2,] "yellow"  "green"    "yellow"
[3,] "green"   "yellow"   "green"
, , 2
      [,1]      [,2]      [,3]
[1,] "yellow"  "green"    "yellow"
[2,] "green"   "yellow"   "green"
[3,] "yellow"  "green"    "yellow"
```

1.5 Data Frames

Data frames are tabular data objects. Unlike a matrix in data frame each column can contain different modes of data. The first column can be numeric while the second column can be character and third column can be logical. It is a list of vectors of equal length.

Data Frames are created using the **data.frame()** function.

```
# Create the data frame.
BMI <- data.frame(
  gender = c("Male", "Male","Female"),
  height = c(152, 171.5, 165),
  weight = c(81,93, 78),
  Age = c(42,38,26)
)
print(BMI)
```

When we execute the above code, it produces the following result –

```
  gender height weight Age
1  Male   152.0     81  42
2  Male   171.5     93  38
3 Female   165.0     78  26
```

2. R - Vectors

Vectors are the most basic R data objects and there are six types of atomic vectors. They are logical, integer, double, complex, character and raw.

2.1 Vector Creation

Single Element Vector

Even when you write just one value in R, it becomes a vector of length 1 and belongs to one of the above vector types.

```
# Atomic vector of type character.
print("abc");

# Atomic vector of type double.
print(12.5)

# Atomic vector of type integer.
print(63L)

# Atomic vector of type logical.
print(TRUE)

# Atomic vector of type complex.
print(2+3i)

# Atomic vector of type raw.
print(charToRaw('hello'))
```

When we execute the above code, it produces the following result –

```
[1] "abc"
[1] 12.5
[1] 63
[1] TRUE
[1] 2+3i
[1] 68 65 6c 6c 6f
```

2.2 Multiple Elements Vector

Using colon operator with numeric data

```
# Creating a sequence from 5 to 13.
v <- 5:13
print(v)

# Creating a sequence from 6.6 to 12.6.
v <- 6.6:12.6
print(v)

# If the final element specified does not belong to the sequence then it is discarded.
v <- 3.8:11.4
print(v)
```

When we execute the above code, it produces the following result –

```
[1] 5 6 7 8 9 10 11 12 13
[1] 6.6 7.6 8.6 9.6 10.6 11.6 12.6
[1] 3.8 4.8 5.8 6.8 7.8 8.8 9.8 10.8
```

Using sequence (Seq.) operator

```
# Create vector with elements from 5 to 9 incrementing by 0.4.
print(seq(5, 9, by = 0.4))
```

When we execute the above code, it produces the following result –

```
[1] 5.0 5.4 5.8 6.2 6.6 7.0 7.4 7.8 8.2 8.6 9.0
```

Using the c() function

The non-character values are coerced to character type if one of the elements is a character.

```
# The logical and numeric values are converted to characters.
s <- c('apple', 'red', 5, TRUE)
print(s)
```

When we execute the above code, it produces the following result –

```
[1] "apple" "red"    "5"      "TRUE"
```

2.3 Accessing Vector Elements

Elements of a Vector are accessed using indexing. The [] **brackets** are used for indexing. Indexing starts with position 1. Giving a negative value in the index drops that element from result. **TRUE**, **FALSE** or **0** and **1** can also be used for indexing.

```
# Accessing vector elements using position.
t <- c("Sun", "Mon", "Tue", "Wed", "Thurs", "Fri", "Sat")
u <- t[c(2, 3, 6)]
print(u)

# Accessing vector elements using logical indexing.
v <- t[c(TRUE, FALSE, FALSE, FALSE, FALSE, TRUE, FALSE)]
print(v)

# Accessing vector elements using negative indexing.
x <- t[c(-2, -5)]
print(x)

# Accessing vector elements using 0/1 indexing.
y <- t[c(0, 0, 0, 0, 0, 0, 1)]
print(y)
```

When we execute the above code, it produces the following result –

```
[1] "Mon" "Tue" "Fri"
[1] "Sun" "Fri"
[1] "Sun" "Tue" "Wed" "Fri" "Sat"
[1] "Sun"
```

2.4 Vector Manipulation

Vector arithmetic

Two vectors of same length can be added, subtracted, multiplied or divided giving the result as a vector output.

```
# Create two vectors.
v1 <- c(3, 8, 4, 5, 0, 11)
v2 <- c(4, 11, 0, 8, 1, 2)

# Vector addition.
add.result <- v1+v2
print(add.result)
```

```

# Vector subtraction.
sub.result <- v1-v2
print(sub.result)

# Vector multiplication.
multi.result <- v1*v2
print(multi.result)

# Vector division.
divi.result <- v1/v2
print(divi.result)

```

When we execute the above code, it produces the following result –

```

[1] 7 19 4 13 1 13
[1] -1 -3 4 -3 -1 9
[1] 12 88 0 40 0 22
[1] 0.7500000 0.7272727      Inf 0.6250000 0.0000000 5.5000000

```

2.5 Vector element recycling

If we apply arithmetic operations to two vectors of unequal length, then the elements of the shorter vector are recycled to complete the operations.

```

v1 <- c(3,8,4,5,0,11)
v2 <- c(4,11)
# V2 becomes c(4,11,4,11,4,11)

add.result <- v1+v2
print(add.result)

sub.result <- v1-v2
print(sub.result)

```

When we execute the above code, it produces the following result –

```

[1] 7 19 8 16 4 22
[1] -1 -3 0 -6 -4 0

```

2.6 Vector Element Sorting

Elements in a vector can be sorted using the `sort()` function.

```

v <- c(3,8,4,5,0,11, -9, 304)

# Sort the elements of the vector.
sort.result <- sort(v)
print(sort.result)

# Sort the elements in the reverse order.
revsort.result <- sort(v, decreasing = TRUE)
print(revsort.result)

# Sorting character vectors.
v <- c("Red","Blue","yellow","violet")
sort.result <- sort(v)
print(sort.result)

# Sorting character vectors in reverse order.
revsort.result <- sort(v, decreasing = TRUE)
print(revsort.result)

```

When we execute the above code, it produces the following result –

```
[1] -9  0  3  4  5  8 11 304
[1] 304 11  8  5  4  3  0 -9
[1] "Blue" "Red" "violet" "yellow"
[1] "yellow" "violet" "Red" "Blue"
```

3. R - Data Frames

A data frame is a table or a two-dimensional array-like structure in which each column contains values of one variable and each row contains one set of values from each column.

Following are the characteristics of a data frame.

- The column names should be non-empty.
- The row names should be unique.
- The data stored in a data frame can be of numeric, factor or character type.
- Each column should contain same number of data items.

3.1 Create Data Frame

```
# Create the data frame.
emp.data <- data.frame(
  emp_id = c(1:5),
  emp_name = c("Rick", "Dan", "Michelle", "Ryan", "Gary"),
  salary = c(623.3, 515.2, 611.0, 729.0, 843.25),

  start_date = as.Date(c("2012-01-01", "2013-09-23", "2014-11-15", "2014-05-11",
    "2015-03-27")),
  stringsAsFactors = FALSE
)
# Print the data frame.
print(emp.data)
```

When we execute the above code, it produces the following result –

emp_id	emp_name	salary	start_date
1	Rick	623.30	2012-01-01
2	Dan	515.20	2013-09-23
3	Michelle	611.00	2014-11-15
4	Ryan	729.00	2014-05-11
5	Gary	843.25	2015-03-27

3.2 Get the Structure of the Data Frame

The structure of the data frame can be seen by using `str()` function.

```
# Create the data frame.
emp.data <- data.frame(
  emp_id = c(1:5),
  emp_name = c("Rick", "Dan", "Michelle", "Ryan", "Gary"),
  salary = c(623.3, 515.2, 611.0, 729.0, 843.25),

  start_date = as.Date(c("2012-01-01", "2013-09-23", "2014-11-15", "2014-05-11",
    "2015-03-27")),
  stringsAsFactors = FALSE
)
# Get the structure of the data frame.
str(emp.data)
```

When we execute the above code, it produces the following result –

```
'data.frame':  5 obs. of  4 variables:
 $ emp_id   : int  1 2 3 4 5
 $ emp_name : chr  "Rick" "Dan" "Michelle" "Ryan" ...
```

```
$ salary      : num  623 515 611 729 843
$ start_date: Date, format: "2012-01-01" "2013-09-23" "2014-11-15" "2014-05-11" ...
```

3.3 Summary of Data in Data Frame

The statistical summary and nature of the data can be obtained by applying **summary()** function.

```
# Create the data frame.
emp.data <- data.frame(
  emp_id = c (1:5),
  emp_name = c("Rick","Dan","Michelle","Ryan","Gary"),
  salary = c(623.3,515.2,611.0,729.0,843.25),

  start_date = as.Date(c("2012-01-01", "2013-09-23", "2014-11-15", "2014-05-11",
    "2015-03-27")),
  stringsAsFactors = FALSE
)
# Print the summary.
print(summary(emp.data))
```

When we execute the above code, it produces the following result –

emp_id	emp_name	salary	start_date
Min. :1	Length:5	Min. :515.2	Min. :2012-01-01
1st Qu.:2	Class :character	1st Qu.:611.0	1st Qu.:2013-09-23
Median :3	Mode :character	Median :623.3	Median :2014-05-11
Mean :3		Mean :664.4	Mean :2014-01-14
3rd Qu.:4		3rd Qu.:729.0	3rd Qu.:2014-11-15
Max. :5		Max. :843.2	Max. :2015-03-27

3.4 Extract Data from Data Frame

Extract specific column from a data frame using column name.

```
# Create the data frame.
emp.data <- data.frame(
  emp_id = c (1:5),
  emp_name = c("Rick","Dan","Michelle","Ryan","Gary"),
  salary = c(623.3,515.2,611.0,729.0,843.25),

  start_date = as.Date(c("2012-01-01","2013-09-23","2014-11-15","2014-05-11",
    "2015-03-27")),
  stringsAsFactors = FALSE
)
# Extract Specific columns.
result <- data.frame(emp.data$emp_name,emp.data$salary)
print(result)
```

When we execute the above code, it produces the following result –

	emp.data.emp_name	emp.data.salary
1	Rick	623.30
2	Dan	515.20
3	Michelle	611.00
4	Ryan	729.00
5	Gary	843.25

Extract the first two rows and then all columns

```
# Create the data frame.
emp.data <- data.frame(
  emp_id = c (1:5),
```

```

emp_name = c("Rick","Dan","Michelle","Ryan","Gary"),
salary = c(623.3,515.2,611.0,729.0,843.25),

start_date = as.Date(c("2012-01-01", "2013-09-23", "2014-11-15", "2014-05-11",
"2015-03-27")),
stringsAsFactors = FALSE
)
# Extract first two rows.
result <- emp.data[1:2,]
print(result)

```

When we execute the above code, it produces the following result –

emp_id	emp_name	salary	start_date
1	Rick	623.3	2012-01-01
2	Dan	515.2	2013-09-23

Extract 3rd and 5th row with 2nd and 4th column

```

# Create the data frame.
emp.data <- data.frame(
  emp_id = c (1:5),
  emp_name = c("Rick","Dan","Michelle","Ryan","Gary"),
  salary = c(623.3,515.2,611.0,729.0,843.25),

  start_date = as.Date(c("2012-01-01", "2013-09-23", "2014-11-15", "2014-05-11",
"2015-03-27")),
stringsAsFactors = FALSE
)
# Extract 3rd and 5th row with 2nd and 4th column.
result <- emp.data[c(3,5),c(2,4)]
print(result)

```

When we execute the above code, it produces the following result –

emp_name	start_date
3 Michelle	2014-11-15
5 Gary	2015-03-27

3.5 Expand Data Frame

A data frame can be expanded by adding columns and rows.

Add Column

Just add the column vector using a new column name.

```

# Create the data frame.
emp.data <- data.frame(
  emp_id = c (1:5),
  emp_name = c("Rick","Dan","Michelle","Ryan","Gary"),
  salary = c(623.3,515.2,611.0,729.0,843.25),

  start_date = as.Date(c("2012-01-01", "2013-09-23", "2014-11-15", "2014-05-11",
"2015-03-27")),
stringsAsFactors = FALSE
)
# Add the "dept" coulumn.
emp.data$dept <- c("IT","Operations","IT","HR","Finance")
v <- emp.data
print(v)

```

When we execute the above code, it produces the following result –

emp_id	emp_name	salary	start_date	dept
1	Rick	623.3	2012-01-01	IT
2	Dan	515.2	2013-09-23	Operations
3	Michelle	611.0	2014-11-15	IT
4	Ryan	729.0	2014-05-11	HR
5	Gary	843.25	2015-03-27	Finance

1	1	Rick	623.30	2012-01-01	IT
2	2	Dan	515.20	2013-09-23	Operations
3	3	Michelle	611.00	2014-11-15	IT
4	4	Ryan	729.00	2014-05-11	HR
5	5	Gary	843.25	2015-03-27	Finance

3.6 Add Row

To add more rows permanently to an existing data frame, we need to bring in the new rows in the same structure as the existing data frame and use the **rbind()** function.

In the example below we create a data frame with new rows and merge it with the existing data frame to create the final data frame.

```
# Create the first data frame.
emp.data <- data.frame(
  emp_id = c(1:5),
  emp_name = c("Rick", "Dan", "Michelle", "Ryan", "Gary"),
  salary = c(623.3, 515.2, 611.0, 729.0, 843.25),

  start_date = as.Date(c("2012-01-01", "2013-09-23", "2014-11-15", "2014-05-11",
    "2015-03-27")),
  dept = c("IT", "Operations", "IT", "HR", "Finance"),
  stringsAsFactors = FALSE
)

# Create the second data frame
emp.newdata <- data.frame(
  emp_id = c(6:8),
  emp_name = c("Rasmi", "Pranab", "Tusar"),
  salary = c(578.0, 722.5, 632.8),
  start_date = as.Date(c("2013-05-21", "2013-07-30", "2014-06-17")),
  dept = c("IT", "Operations", "Fianance"),
  stringsAsFactors = FALSE
)

# Bind the two data frames.
emp.finaldata <- rbind(emp.data, emp.newdata)
print(emp.finaldata)
```

When we execute the above code, it produces the following result –

emp_id	emp_name	salary	start_date	dept
1	Rick	623.30	2012-01-01	IT
2	Dan	515.20	2013-09-23	Operations
3	Michelle	611.00	2014-11-15	IT
4	Ryan	729.00	2014-05-11	HR
5	Gary	843.25	2015-03-27	Finance
6	Rasmi	578.00	2013-05-21	IT
7	Pranab	722.50	2013-07-30	Operations
8	Tusar	632.80	2014-06-17	Fianance

4. R - CSV Files

In R, we can read data from files stored outside the R environment. We can also write data into files which will be stored and accessed by the operating system. R can read and write into various file formats like csv, excel, xml etc. In this chapter we will learn to read data from a csv file and then write data into a csv file. The file should be present in current working directory so that R can read it. Of course we can also set our own directory and read files from there.

4.1 Getting and Setting the Working Directory

You can check which directory the R workspace is pointing to using the `getwd()` function. You can also set a new working directory using `setwd()` function.

```
# Get and print current working directory.
print(getwd())

# Set current working directory.
setwd("/web/com")

# Get and print current working directory.
print(getwd())
```

When we execute the above code, it produces the following result –

```
[1] "/web/com/1441086124_2016"
[1] "/web/com"
```

This result depends on your OS and your current directory where you are working.

4.2 Input as CSV File

The csv file is a text file in which the values in the columns are separated by a comma. Let's consider the following data present in the file named **input.csv**.

You can create this file using windows notepad by copying and pasting this data. Save the file as **input.csv** using the save As All files(*.*) option in notepad.

```
id,name,salary,start_date,dept
1,Rick,623.3,2012-01-01,IT
2,Dan,515.2,2013-09-23,Operations
3,Michelle,611,2014-11-15,IT
4,Ryan,729,2014-05-11,HR
5,Gary,843.25,2015-03-27,Finance
6,Nina,578,2013-05-21,IT
7,Simon,632.8,2013-07-30,Operations
8,Guru,722.5,2014-06-17,Finance
```

4.3 Reading a CSV File

Following is a simple example of `read.csv()` function to read a CSV file available in your current working directory

```
data <- read.csv("input.csv")
print(data)
```

When we execute the above code, it produces the following result –

	id,	name,	salary,	start_date,	dept
1	1	Rick	623.30	2012-01-01	IT
2	2	Dan	515.20	2013-09-23	Operations
3	3	Michelle	611.00	2014-11-15	IT
4	4	Ryan	729.00	2014-05-11	HR
5	NA	Gary	843.25	2015-03-27	Finance
6	6	Nina	578.00	2013-05-21	IT
7	7	Simon	632.80	2013-07-30	Operations
8	8	Guru	722.50	2014-06-17	Finance

4.4 Analyzing the CSV File

By default the `read.csv()` function gives the output as a data frame. This can be easily checked as follows. Also we can check the number of columns and rows.

```
data <- read.csv("input.csv")

print(is.data.frame(data))
print(ncol(data))
print(nrow(data))
```

When we execute the above code, it produces the following result –

```
[1] TRUE
[1] 5
[1] 8
```

Once we read data in a data frame, we can apply all the functions applicable to data frames as explained in subsequent section.

Get the maximum salary

```
# Create a data frame.
data <- read.csv("input.csv")

# Get the max salary from data frame.
sal <- max(data$salary)
print(sal)
```

When we execute the above code, it produces the following result –

```
[1] 843.25
```

Get the details of the person with max salary

We can fetch rows meeting specific filter criteria similar to a SQL where clause.

```
# Create a data frame.
data <- read.csv("input.csv")

# Get the max salary from data frame.
sal <- max(data$salary)

# Get the person detail having max salary.
retval <- subset(data, salary == max(salary))
print(retval)
```

When we execute the above code, it produces the following result –

	id	name	salary	start_date	dept	
	5	NA	Gary	843.25	2015-03-27	Finance

Get all the people working in IT department

```
# Create a data frame.
data <- read.csv("input.csv")

retval <- subset(data, dept == "IT")
print(retval)
```

When we execute the above code, it produces the following result –

	id	name	salary	start_date	dept
1	1	Rick	623.3	2012-01-01	IT
3	3	Michelle	611.0	2014-11-15	IT
6	6	Nina	578.0	2013-05-21	IT

Get the persons in IT department whose salary is greater than 600

```
# Create a data frame.
data <- read.csv("input.csv")

info <- subset(data, salary > 600 & dept == "IT")
print(info)
```

When we execute the above code, it produces the following result –

	id	name	salary	start_date	dept
1	1	Rick	623.3	2012-01-01	IT
3	3	Michelle	611.0	2014-11-15	IT

Get the people who joined on or after 2014

```
# Create a data frame.
data <- read.csv("input.csv")

retval <- subset(data, as.Date(start_date) > as.Date("2014-01-01"))
print(retval)
```

When we execute the above code, it produces the following result –

	id	name	salary	start_date	dept
3	3	Michelle	611.00	2014-11-15	IT
4	4	Ryan	729.00	2014-05-11	HR
5	NA	Gary	843.25	2015-03-27	Finance
8	8	Guru	722.50	2014-06-17	Finance

4.5 Writing into a CSV File

R can create csv file from existing data frame. The `write.csv()` function is used to create the csv file. This file gets created in the working directory.

```
# Create a data frame.
data <- read.csv("input.csv")
retval <- subset(data, as.Date(start_date) > as.Date("2014-01-01"))

# Write filtered data into a new file.
write.csv(retval, "output.csv")
newdata <- read.csv("output.csv")
print(newdata)
```

When we execute the above code, it produces the following result –

X	id	name	salary	start_date	dept
1	3	Michelle	611.00	2014-11-15	IT
2	4	Ryan	729.00	2014-05-11	HR
3	5	Gary	843.25	2015-03-27	Finance
4	8	Guru	722.50	2014-06-17	Finance

Here the column X comes from the data set newper. This can be dropped using additional parameters while writing the file.

```
# Create a data frame.
data <- read.csv("input.csv")
retval <- subset(data, as.Date(start_date) > as.Date("2014-01-01"))

# Write filtered data into a new file.
write.csv(retval, "output.csv", row.names = FALSE)
newdata <- read.csv("output.csv")
print(newdata)
```

When we execute the above code, it produces the following result –

	id	name	salary	start_date	dept
1	3	Michelle	611.00	2014-11-15	IT
2	4	Ryan	729.00	2014-05-11	HR
3	NA	Gary	843.25	2015-03-27	Finance
4	8	Guru	722.50	2014-06-17	Finance

5. R - Packages

R packages are a collection of R functions, compiled code and sample data. They are stored under a directory called **"library"** in the R environment. By default, R installs a set of packages during installation. More packages are added later, when they are needed for some specific purpose. When we start the R console, only the default packages are available by default. Other packages which are already installed have to be loaded explicitly to be used by the R program that is going to use them.

All the packages available in R language are listed at [R Packages](#).

Below is a list of commands to be used to check, verify and use the R packages.

5.1 Check Available R Packages

Get library locations containing R packages

```
.libPaths()
```

When we execute the above code, it produces the following result. It may vary depending on the local settings of your pc.

```
[2] "C:/Program Files/R/R-3.2.2/library"
```

5.2 Get the list of all the packages installed

```
library()
```

When we execute the above code, it produces the following result. It may vary depending on the local settings of your pc.

```
Packages in library 'C:/Program Files/R/R-3.2.2/library':

base           The R Base Package
boot           Bootstrap Functions (Originally by Angelo Canty
              for S)
class          Functions for Classification
cluster        "Finding Groups in Data": Cluster Analysis
              Extended Rousseeuw et al.
codetools      Code Analysis Tools for R
compiler       The RCompiler Package
```

Get all packages currently loaded in the R environment

```
search()
```

When we execute the above code, it produces the following result. It may vary depending on the local settings of your pc.

```
[1] ".GlobalEnv"      "package:stats"    "package:graphics"
[4] "package:grDevices" "package:utils"    "package:datasets"
[7] "package:methods" "Autoloads"        "package:base"
```

5.3 Install a New Package

There are two ways to add new R packages. One is installing directly from the CRAN directory and another is downloading the package to your local system and installing it manually.

Install directly from CRAN

The following command gets the packages directly from CRAN webpage and installs the package in the R environment. You may be prompted to choose a nearest mirror. Choose the one appropriate to your location.

```
install.packages("Package Name")

# Install the package named "XML".
install.packages("XML")
```

Install package manually

Go to the link [R Packages](#) to download the package needed. Save the package as a **.zip** file in a suitable location in the local system.

Now you can run the following command to install this package in the R environment.

```
install.packages(file_name_with_path, repos = NULL, type = "source")

# Install the package named "XML"
install.packages("E:/XML_3.98-1.3.zip", repos = NULL, type = "source")
```

5.4 Load Package to Library

Before a package can be used in the code, it must be loaded to the current R environment. You also need to load a package that is already installed previously but not available in the current environment.

A package is loaded using the following command –

```
library("package Name", lib.loc = "path to library")

# Load the package named "XML"
install.packages("E:/XML_3.98-1.3.zip", repos = NULL, type = "source")
```

6. R - Linear Regression

6.1 Steps to Establish a Regression

A simple example of regression is predicting weight of a person when his height is known. To do this we need to have the relationship between height and weight of a person.

The steps to create the relationship is –

- Carry out the experiment of gathering a sample of observed values of height and corresponding weight.
- Create a relationship model using the **lm()** functions in R.
- Find the coefficients from the model created and create the mathematical equation using these
- Get a summary of the relationship model to know the average error in prediction. Also called **residuals**.
- To predict the weight of new persons, use the **predict()** function in R.

Input Data

Below is the sample data representing the observations –

```
# Values of height
151, 174, 138, 186, 128, 136, 179, 163, 152, 131

# Values of weight.
63, 81, 56, 91, 47, 57, 76, 72, 62, 48
```

lm() Function

This function creates the relationship model between the predictor and the response variable.

Syntax

The basic syntax for **lm()** function in linear regression is –

```
lm(formula,data)
```

Following is the description of the parameters used –

- **formula** is a symbol presenting the relation between x and y.
- **data** is the vector on which the formula will be applied.

Create Relationship Model & get the Coefficients

```
x <- c(151, 174, 138, 186, 128, 136, 179, 163, 152, 131)
y <- c(63, 81, 56, 91, 47, 57, 76, 72, 62, 48)

# Apply the lm() function.
relation <- lm(y~x)

print(relation)
```

When we execute the above code, it produces the following result –

```
Call:
lm(formula = y ~ x)

Coefficients:
(Intercept)          x
-38.4551          0.6746
```

Get the Summary of the Relationship

```
x <- c(151, 174, 138, 186, 128, 136, 179, 163, 152, 131)
y <- c(63, 81, 56, 91, 47, 57, 76, 72, 62, 48)

# Apply the lm() function.
relation <- lm(y~x)
```

```
print(summary(relation))
```

When we execute the above code, it produces the following result -

```
Call:
lm(formula = y ~ x)

Residuals:
    Min       1Q   Median       3Q      Max
-6.3002  -1.6629   0.0412   1.8944   3.9775

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept) -38.45509     8.04901  -4.778  0.00139 **
x             0.67461     0.05191  12.997 1.16e-06 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 3.253 on 8 degrees of freedom
Multiple R-squared:  0.9548,    Adjusted R-squared:  0.9491
F-statistic: 168.9 on 1 and 8 DF,  p-value: 1.164e-06
```

6.2 predict() Function

Syntax

The basic syntax for predict() in linear regression is -

```
predict(object, newdata)
```

Following is the description of the parameters used -

- **object** is the formula which is already created using the lm() function.
- **newdata** is the vector containing the new value for predictor variable.

Predict the weight of new persons

```
# The predictor vector.
x <- c(151, 174, 138, 186, 128, 136, 179, 163, 152, 131)

# The response vector.
y <- c(63, 81, 56, 91, 47, 57, 76, 72, 62, 48)

# Apply the lm() function.
relation <- lm(y~x)

# Find weight of a person with height 170.
a <- data.frame(x = 170)
result <- predict(relation,a)
print(result)
```

When we execute the above code, it produces the following result -

```
1
76.22869
```

6.3 Visualize the Regression Graphically

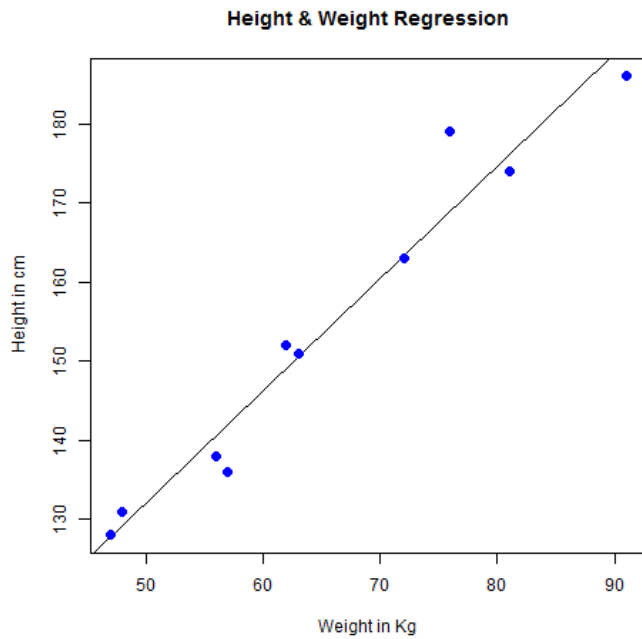
```
# Create the predictor and response variable.
x <- c(151, 174, 138, 186, 128, 136, 179, 163, 152, 131)
y <- c(63, 81, 56, 91, 47, 57, 76, 72, 62, 48)
relation <- lm(y~x)

# Give the chart file a name.
```



```
png(file = "linearregression.png")  
  
# Plot the chart.  
plot(y,x,col = "blue",main = "Height & Weight Regression",  
abline(lm(x~y)),cex = 1.3,pch = 16,xlab = "Weight in Kg",ylab = "Height in  
cm")  
  
# Save the file.  
dev.off()
```

When we execute the above code, it produces the following result –



7. R - Time Series Analysis

Time series is a series of data points in which each data point is associated with a timestamp. A simple example is the price of a stock in the stock market at different points of time on a given day. Another example is the amount of rainfall in a region at different months of the year. R language uses many functions to create, manipulate and plot the time series data. The data for the time series is stored in an R object called **time-series object**. It is also a R data object like a vector or data frame.

The time series object is created by using the **ts()** function.

Syntax

The basic syntax for **ts()** function in time series analysis is –

```
timeseries.object.name <- ts(data, start, end, frequency)
```

Following is the description of the parameters used –

- **data** is a vector or matrix containing the values used in the time series.
- **start** specifies the start time for the first observation in time series.
- **end** specifies the end time for the last observation in time series.
- **frequency** specifies the number of observations per unit time.

Except the parameter "data" all other parameters are optional.

Example

Consider the annual rainfall details at a place starting from January 2012. We create an R time series object for a period of 12 months and plot it.

```
# Get the data points in form of a R vector.
rainfall <- c(799,1174.8,865.1,1334.6,635.4,918.5,685.5,998.6,784.2,985,882.8,1071)

# Convert it to a time series object.
rainfall.timeseries <- ts(rainfall,start = c(2012,1),frequency = 12)

# Print the timeseries data.
print(rainfall.timeseries)

# Give the chart file a name.
png(file = "rainfall.png")

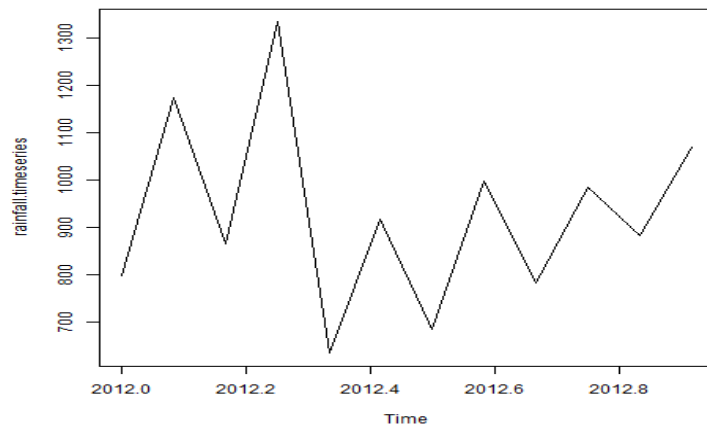
# Plot a graph of the time series.
plot(rainfall.timeseries)

# Save the file.
dev.off()
```

When we execute the above code, it produces the following result and chart –

Jan	Feb	Mar	Apr	May	Jun	Jul	Aug	Sep	
2012	799.0	1174.8	865.1	1334.6	635.4	918.5	685.5	998.6	784.2
	Oct	Nov	Dec						
2012	985.0	882.8	1071.0						

The Time series chart –



Different Time Intervals

The value of the **frequency** parameter in the `ts()` function decides the time intervals at which the data points are measured. A value of 12 indicates that the time series is for 12 months. Other values and its meaning is as below –

- **frequency = 12** pegs the data points for every month of a year.
- **frequency = 4** pegs the data points for every quarter of a year.
- **frequency = 6** pegs the data points for every 10 minutes of an hour.
- **frequency = 24*6** pegs the data points for every 10 minutes of a day.

Multiple Time Series

We can plot multiple time series in one chart by combining both the series into a matrix.

```
# Get the data points in form of a R vector.
rainfall1 <- c(799,1174.8,865.1,1334.6,635.4,918.5,685.5,998.6,784.2,985,882.8,1071)
rainfall2 <-
c(655,1306.9,1323.4,1172.2,562.2,824,822.4,1265.5,799.6,1105.6,1106.7,1337.8)

# Convert them to a matrix.
combined.rainfall <- matrix(c(rainfall1,rainfall2),nrow = 12)

# Convert it to a time series object.
rainfall.timeseries <- ts(combined.rainfall,start = c(2012,1),frequency = 12)

# Print the timeseries data.
print(rainfall.timeseries)

# Give the chart file a name.
png(file = "rainfall_combined.png")

# Plot a graph of the time series.
plot(rainfall.timeseries, main = "Multiple Time Series")

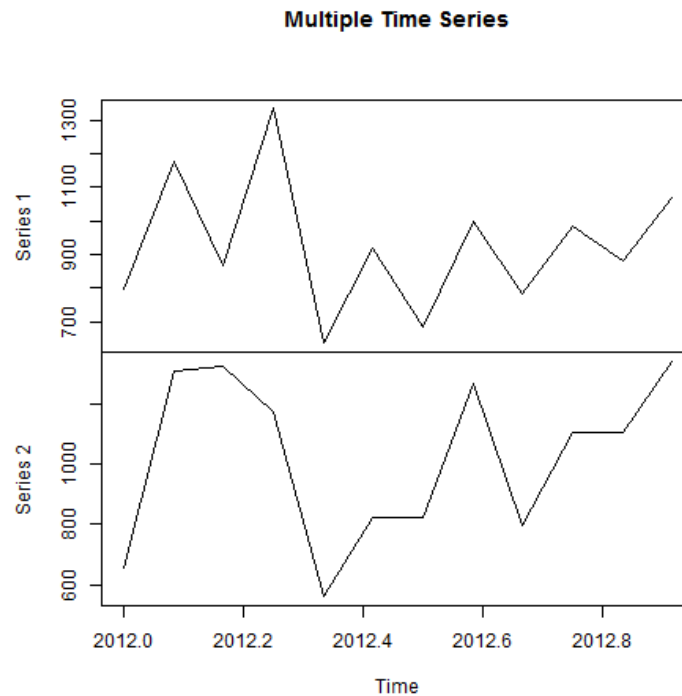
# Save the file.
dev.off()
```

When we execute the above code, it produces the following result and chart –

	Series 1	Series 2
Jan 2012	799.0	655.0
Feb 2012	1174.8	1306.9
Mar 2012	865.1	1323.4
Apr 2012	1334.6	1172.2
May 2012	635.4	562.2
Jun 2012	918.5	824.0

Jul 2012	685.5	822.4
Aug 2012	998.6	1265.5
Sep 2012	784.2	799.6
Oct 2012	985.0	1105.6
Nov 2012	882.8	1106.7
Dec 2012	1071.0	1337.8

The Multiple Time series chart –



Overview of Time Series Objects in R

The core data object for holding data in R is the data.frame object. A data.frame is a rectangular data object whose columns can be of different types (e.g., numeric, character, logical, Date, etc.). The data.frame object, however, is not designed to work efficiently with time series data. In particular, sub-setting and merging data based on a time index is cumbersome and transforming and aggregating data based on a time index is not at all straightforward.

Furthermore, the default plotting methods in R are not designed for handling time series data. Hence, there is a need for a flexible time series class in R with a rich set of methods for manipulating and plotting time series data.

Base R has limited functionality for handling general time series data. For example, univariate and multivariate regularly spaced calendar time series data can be represented using the ts and mts classes, respectively. These classes have a limited set of method functions for manipulating and plotting time series data. However, these classes cannot adequately represent more general irregularly spaced non-calendar time series such as intra-day transactions

level financial price and quote data. Fortunately, there are several R packages that can be used to handle general time series data.

The table below lists the main time series objects that are available in R and their respective packages.

Time Series Object	Package	Description
<code>fts</code>	fts	An R interface to <code>tslib</code> (a time series library in C++)
<code>its</code>	its	An S4 class for handling irregular time series
<code>irts</code>	tseries	<code>irts</code> objects are irregular time-series objects. These are scalar or vector valued time series indexed by a time-stamp of class "POSIXct".
<code>timeSeries</code>	timeSeries	Rmetrics package of time series tools and utilities. Similar to the Tibco S-PLUS <code>timeSeries</code> class
<code>ti</code>	tis	Functions and S3 classes for time indexes and time indexed series, which are compatible with FAME frequencies
<code>ts, mts</code>	stats	Regularly spaced time series objects
<code>zoo</code>	zoo	S3 class of indexed totally ordered observations which includes irregular time series.
<code>xts</code>	xts	Extension of the <code>zoo</code> class

The `ts` and `mts` classes in base R are suitable for representing regularly spaced calendar time series such as monthly sales or quarterly real GDP. In addition, several of the time series modeling functions in base R and in several R packages take `ts` and `mts` objects as data inputs. For handling more general irregularly spaced financial time series, by far the most used packages are **timeSeries**, **zoo** and **xts**. The **timeSeries** package is part of the suite of **Rmetrics** packages for financial data analysis and computational finance created by Diethelm Weurtz and his colleagues at ETZ Zurich (see www.Rmetrics.org). In these packages, `timeSeries` objects are the core data objects. However, outside of **Rmetrics**, `timeSeries` objects are not as frequently used as `zoo` and `xts` objects for representing time series data. Hence, in this tutorial I will focus mostly on using `zoo` and `xts` objects for handling general time series.

Time series data represented by `timeSeries`, `zoo` and `xts` objects have a similar structure: the time index is stored as a vector in some (typically ordered) date-time object, and the data is stored in some rectangular data object. The resulting `timeSeries`, `zoo` or `xts` objects combine the time index and data into a single object. These objects can then be manipulated and visualized using various method functions.

Here is the example to load data from FRED (Federal Reserve Economic Data)

```
# install.packages("fImport")
# install.packages("quantmod")
# install.packages("fBasics")
# install.packages("car")
# install.packages("xts")
# install.packages("lmtest")
#install.packages("sandwich")

# Clear data and Memory
rm(list=ls())

library(fImport)
library(quantmod)
library(fBasics)
library(car)
library(xts)
library(lmtest)
library(sandwich)
#environment in which to store data
data <- new.env()

# FRED data in R
# set tickers
tickers <- c("FEDFUNDS", "UNRATE", "CPIAUCSL")
getSymbols(tickers, src = "FRED")
INFLATION <- 100*((1+diff(log(CPIAUCSL)))^12) - 100
colnames(INFLATION) <- "INFLATION"

# Time Series Plot
chartSeries(FEDFUNDS,theme="white")
chartSeries(UNRATE,theme="white")
chartSeries(INFLATION,theme="white")
chartSeries(INFLATION,theme="white", subset='2007-01::2017-07')

# Handling Missing observations
dat1 <- merge(FEDFUNDS, UNRATE, INFLATION)
dat2 <- dat1[complete.cases(dat1),]
head(dat1)
head(dat2)
class(dat2)

# Multiple Time Series Graph
par(mfcol=c(3,1))
plot(dat2[,1],main=dimnames(dat2)[[2]][1],col="blue")
plot(dat2[,2],main=dimnames(dat2)[[2]][2],col="blue")
plot(dat2[,3],main=dimnames(dat2)[[2]][3],col="blue")

# Subsample for XTS data
dat2a <- window(dat2,begin="1954-01-01",end="1970-12-31")
dat2b <- dat2["1971-02-01::"]
head(dat2b)

# Simple Regression
# Before 1970
OLS1 <- lm(UNRATE ~ INFLATION, data=dat2a)
summary(OLS1)
dwtest(OLS1)
coefstest(OLS1, df = Inf, vcov = NeweyWest)
plot(OLS1)
```

```

par(mfcol=c(1,1))
scatterplot(as.vector(dat2a$INFLATION),as.vector(dat2a$UNRATE))

# Whole Periods
OLS2 <- lm(UNRATE ~ INFLATION, data=dat2)
summary(OLS2)
dwtest(OLS2)
coeftest(OLS2, df = Inf, vcov = NeweyWest)

scatterplot(as.vector(dat2$INFLATION),as.vector(dat2$UNRATE))

```

Here is the example to load data from Yahoo Finance

```

# Yahoo Data in R

getSymbols("AAPL",from="2005-01-02", to="2016-12-31")
getSymbols("IBM",from="2005-01-02", to="2016-12-31")
getSymbols("^GSPC",from="2005-01-02", to="2016-12-31")
getSymbols("^VIX",from="2005-01-02", to="2016-12-31")

chartSeries(GSPC, theme="white")
chartSeries(AAPL, theme="white")
chartSeries(IBM, theme="white",subset='2007-01::2010-01')
chartSeries(VIX, theme="white")

AAPL.rtn = 100*diff(log(AAPL$AAPL.Adjusted))
IBM.rtn = 100*diff(log(IBM$IBM.Adjusted))
GSPC.rtn = 100*diff(log(GSPC$GSPC.Adjusted))

chartSeries(AAPL.rtn, theme="white",subset='2007::2010')
chartSeries(IBM.rtn, theme="white",subset='2007::2010')
chartSeries(GSPC.rtn, theme="white",subset='2007::2010')

par(mfcol=c(2,1))
plot(AAPL$AAPL.Adjusted,col="blue",main="Apple Stock Price")
plot(AAPL.rtn,col="blue",main="Apple Stock Daily Return")

# Descriptive Statistics and Histogram
basicStats(AAPL.rtn)
hist(AAPL.rtn,nclass=15)
basicStats(IBM.rtn)
hist(IBM.rtn,nclass=15)

apple.capm <- lm(AAPL.rtn ~ GSPC.rtn)
summary(apple.capm)

ibm.capm <- lm(IBM.rtn ~ GSPC.rtn)
summary(ibm.capm)

# Merged Data and use data
findata <- cbind(IBM.rtn, AAPL.rtn)
class(findata)
colnames(findata) <- c("ibm","apple")
class(findata)
head(findata)

# Plot subset of Data
par(mfcol=c(2,1))

```

```
plot(findata$ibm["2007::2012"],col="blue")  
plot(findata$apple["2007::2012"],col="red")
```