

## Experiment no:1

**Aim:** Installation of R and R studio

### Theory:

#### R:

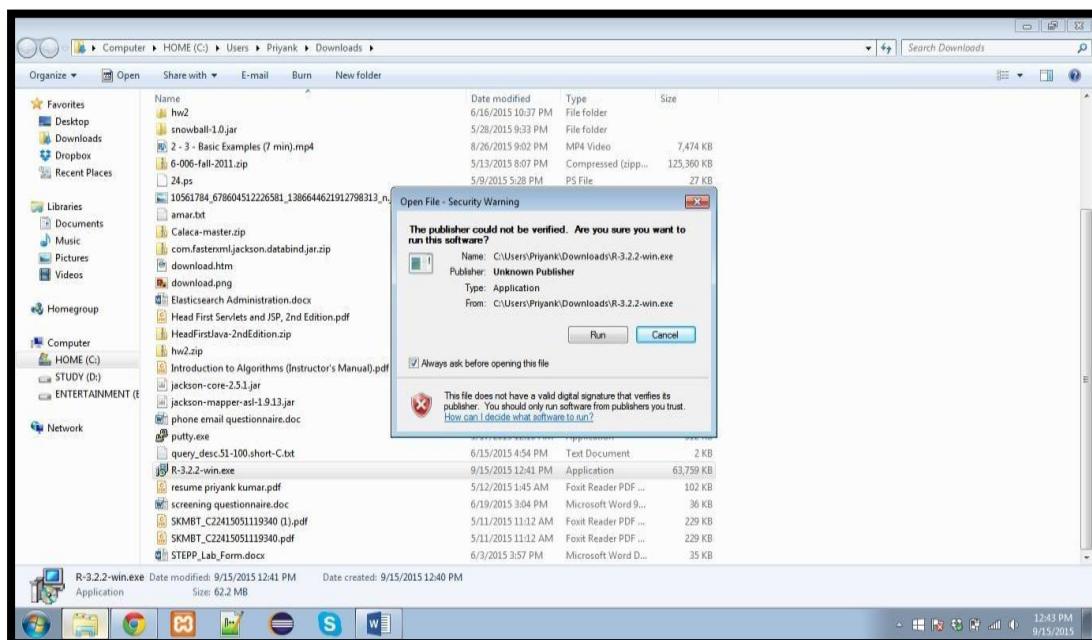
R is a free, open source software program for statistical analysis, based on the S language.

#### R-STUDIO:

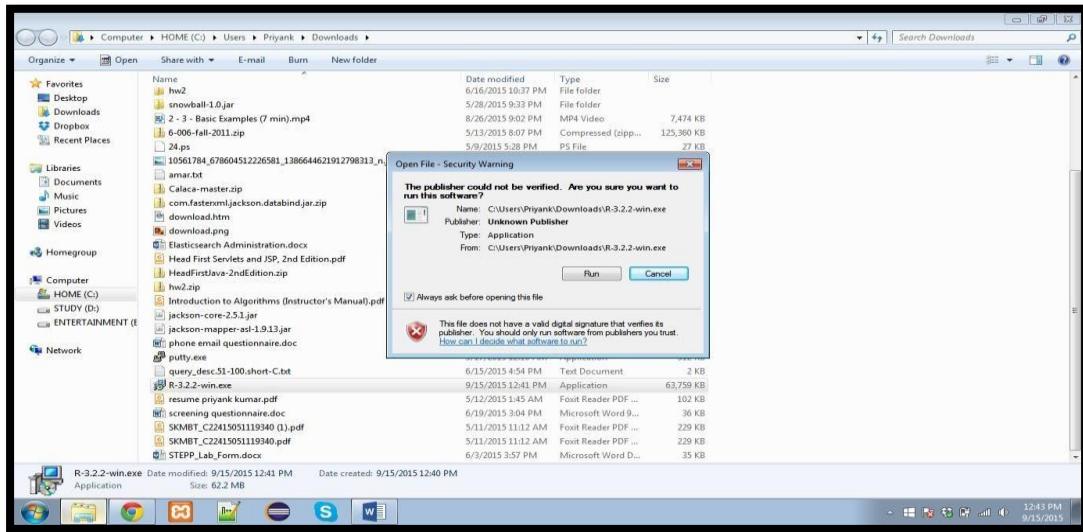
RStudio is a free, open source IDE (integrated development environment) for R. Its interface is organized so that the user can clearly view graphs, data tables, R code, and output all at the same time. It also offers an Import-Wizard-like feature that allows users to import CSV, Excel, SAS (\*.sas7bdat), SPSS (\*.sav), and Stata (\*.dta) files into R without having to write the code to do so.

### INSTALLATION OF R AND R-STUDIO:

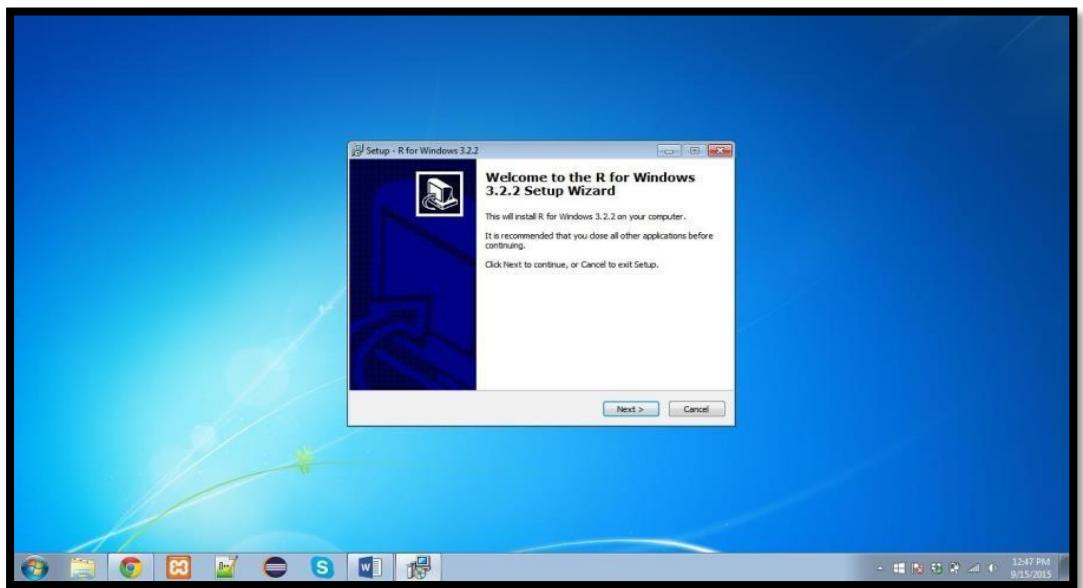
1. Download the installable file from the following link:  
<https://cran.r-project.org/bin/windows/base/>
2. Click on the R 3.2.2.exe file. The 3.2.2 is the version number of the file. The versions can be updated as per the latest releases.



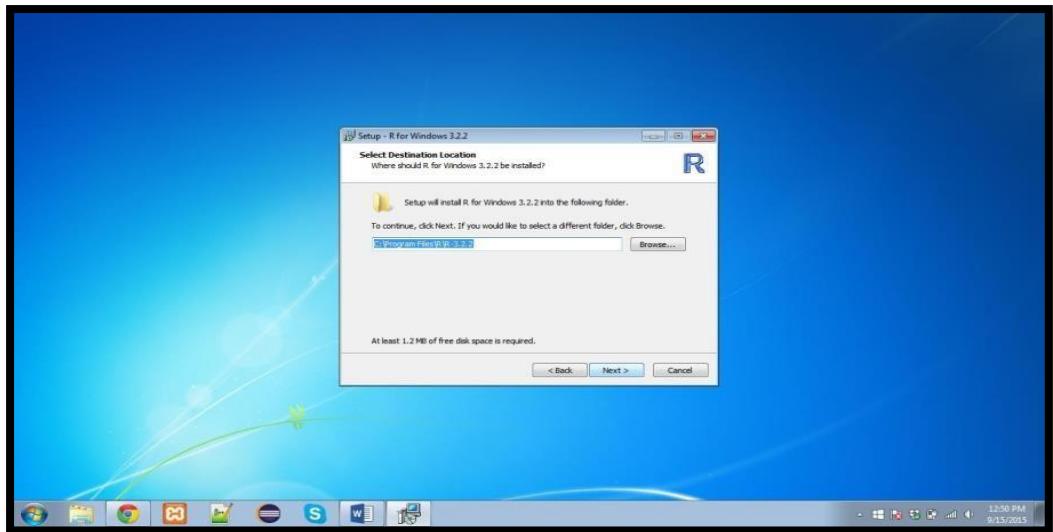
- The Setup will request permission to be installed on the system click yes to proceed.



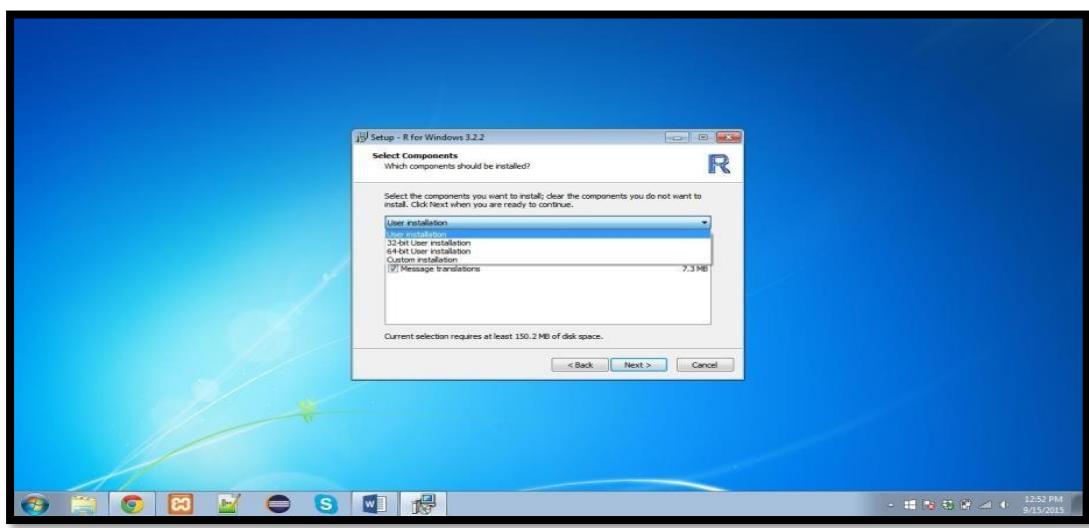
- Select the Preferred language from the dropdown to begin an installation in that preferred language.
- Click next to proceed with the installation.



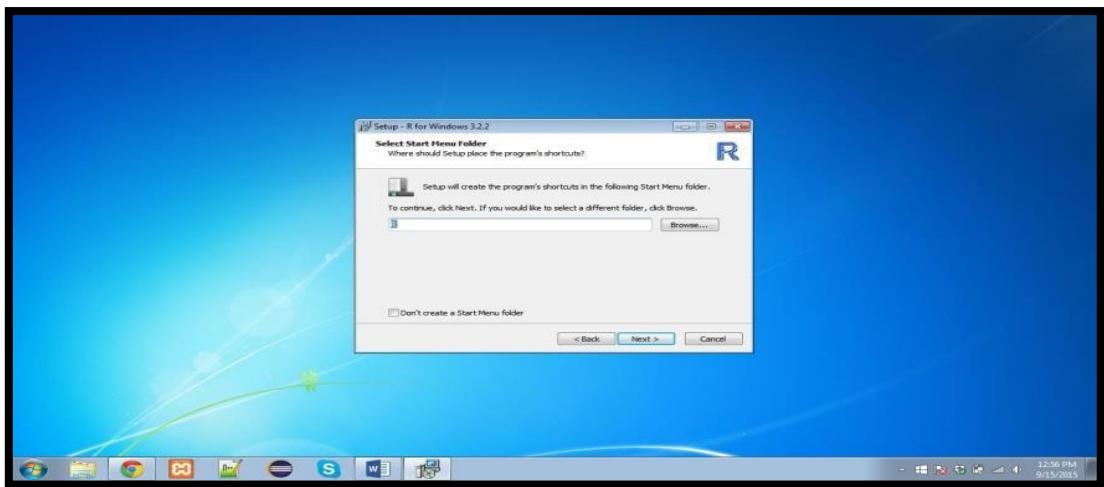
- Choose the path where you wish to install R by clicking on browse and changing the workspace locations. Click next to proceed with the default installation. The minimum space requirements are mentioned at the bottom of the dialog box. Please check you have required amount of free space in your drive.



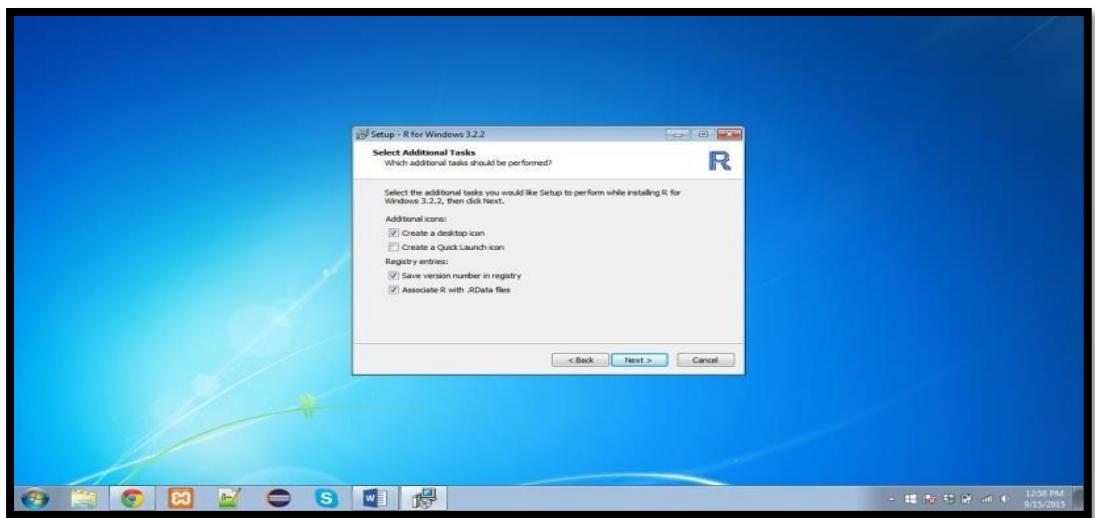
7. Choose the type of installation you require. By default, R installs both the 32- and 64-bit versions on your system. If your system is a 32-bit system you will be requiring a 32-bit installation if the system is a 64-bit system it will be requiring 64-bit installation. Do not uncheck the Core Files and Message Translations. Please make note of the space requirement of the installation.



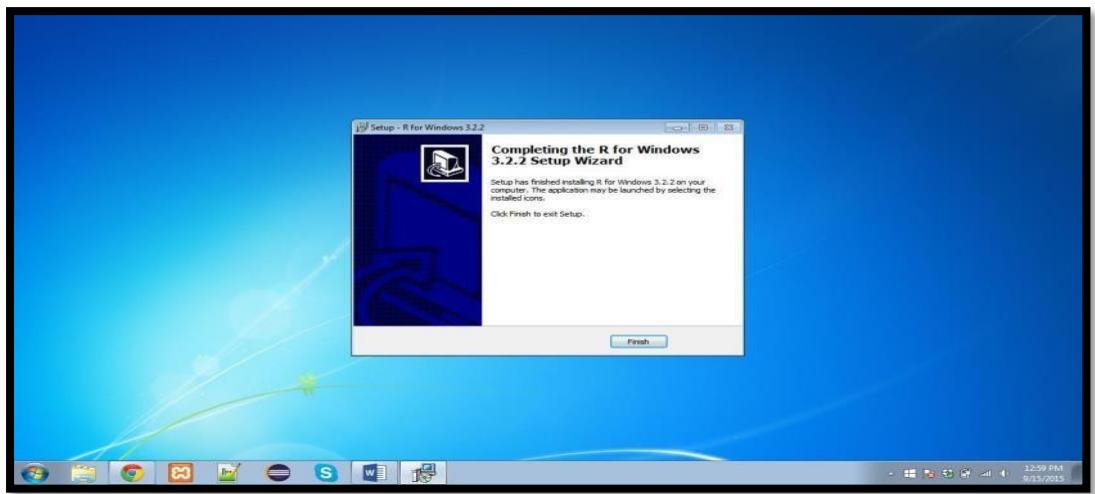
8. To customize the start-up options for R, choose option and customize.  
To proceed with a vanilla installation, use Next.
9. To generate program shortcuts and naming those as per your requirement specify the necessary customizations. To proceed with the default installation hit next.



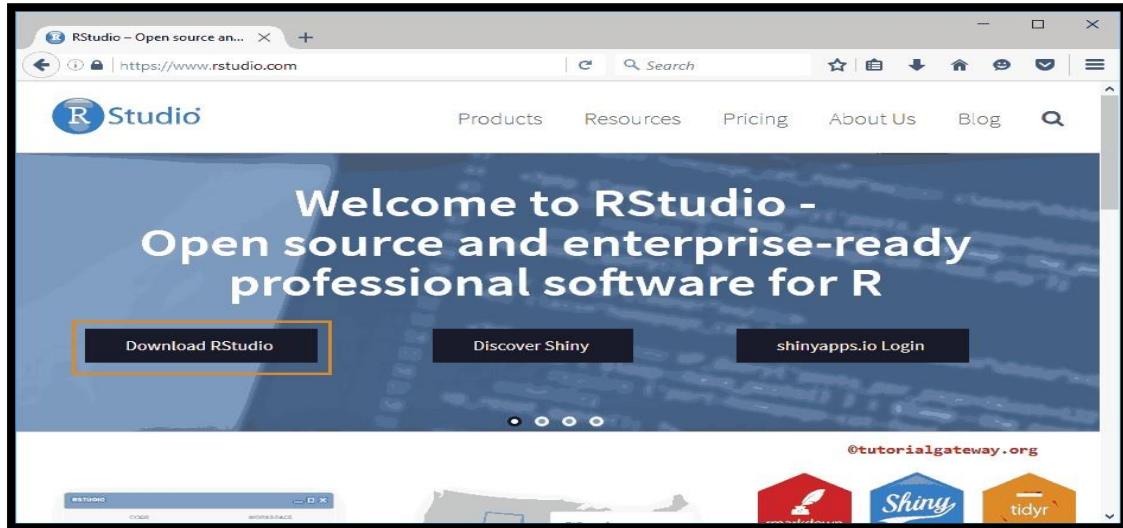
10. Click on the next button to begin your installation.



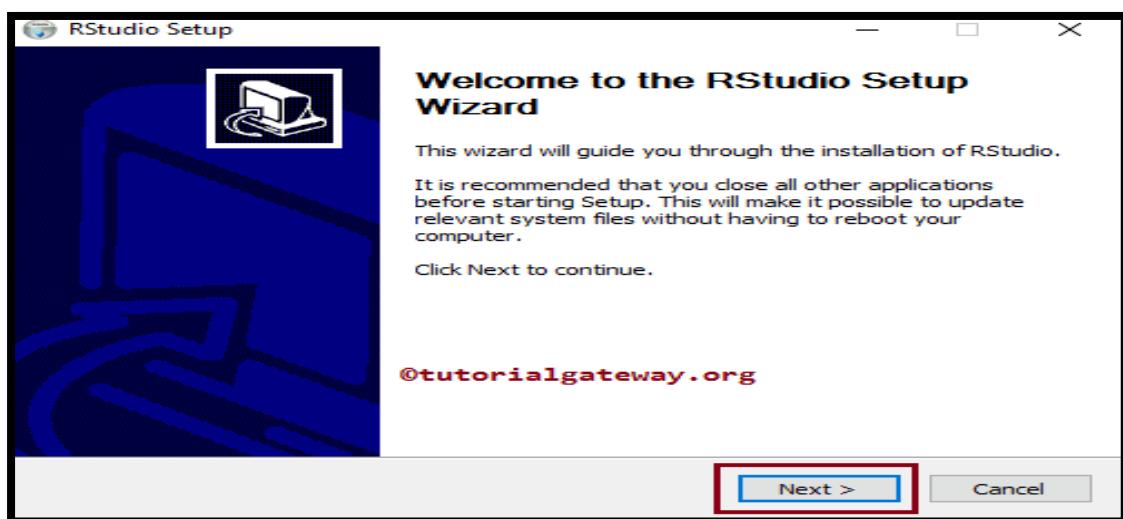
11. After the installation has completed you will see the final screen. Click finish to complete the installation.



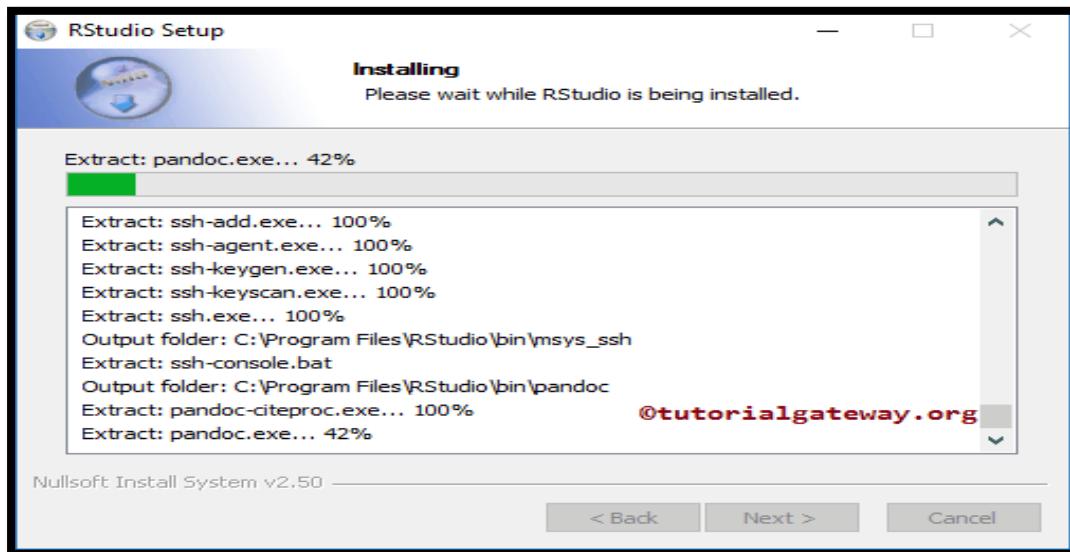
12. Open Start Menu and you will find R in the available set of Programs.
13. Click on the R icon in the menu settings to open R.
14. To download R Studio, first Go to the official R Studio website by clicking this link Download.



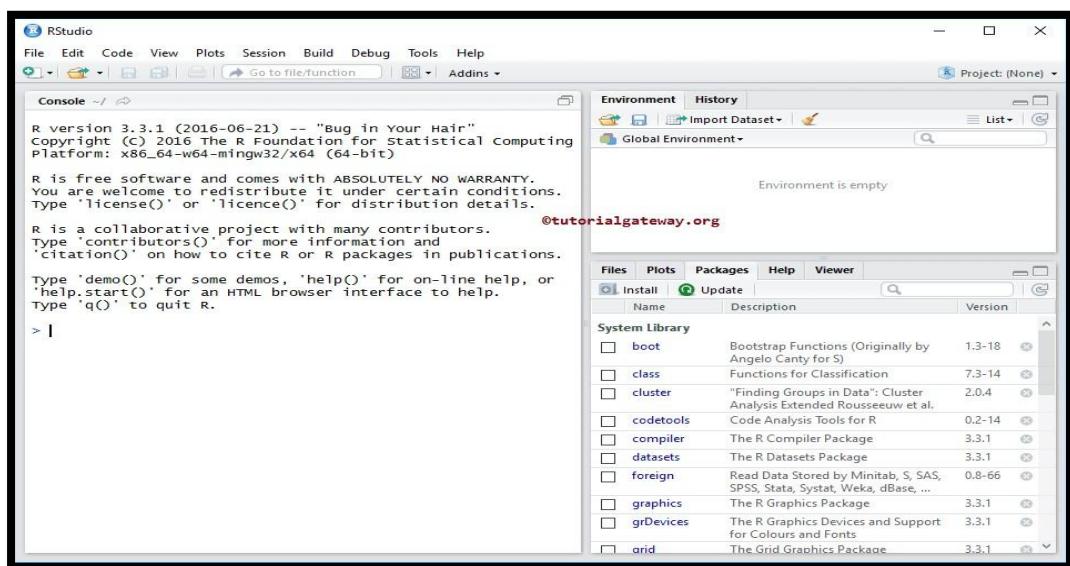
15. Once you click on the specified download RStudio button, the following tab or window opened. Please select the R Studio Desktop (Free License) and click the download button.
16. Or scroll down the page and select the RStudio for your Operating System. From the below screenshot, see that we are selecting the RStudio 0.99.903 – Windows Vista/7/8/10
17. Once you Click on the RStudio hyperlink, a pop-up window opened to save this file. Please select the Save or download R Studio File option.
18. After you click on the install R Studio application, Welcome to the RStudio Setup Wizard window opened. Please select the Next button.



19. Here, you can change the default R Studio installation directory. To do that, we have to click on the Browse... button beside the text box. It opens our file system to navigate the location and select the directory.
20. Here, you can change the R Studio program's shortcuts default directory. For this, we have to select the directory name from the list.
21. Please wait for the extraction and installation of R studio completed. Next, Click the Finish button for completing the Installation of R Studio.



22. Let us open the R Studio to see the R Software environment.



**Conclusion:** Thus, we have successfully installed R and R studio.

## Experiment no:2

**Aim:** Basic functionality in R-variables, basic arithmetic, help command

### Theory:

#### Variables in R:

A variable is a name given to a memory location, which is used to store values in a computer program. Variables in R programming can be used to store numbers (real and complex), words, matrices, and even tables. R is a dynamically programmed language which means that unlike other programming languages, we do not have to declare the data type of a variable before we can use it in our program. For a variable to be valid, it should follow these rules:

1. It should contain letters, numbers, and only dot or underscore characters.
2. It should not start with a number (e.g.: - 2iota)
3. It should not start with a dot followed by a number (e.g.: - .2iota)
4. It should not start with an underscore (e.g.: - \_iota)
5. It should not be a reserved keyword

#### Constants in R:

The entities whose values are fixed are called constants. There are two types of constants in R:

- **Numeric Constants:** All numeric values such as integer, double, or complex fall under this category. Numeric constants followed by 'L' and 'i' are considered as integer and complex respectively. And, numeric constants preceded by 0x/0X are treated as hexadecimal numbers.
- **Character Constants:** These constants are represented by single ('') or double ("") quotes called delimiters.

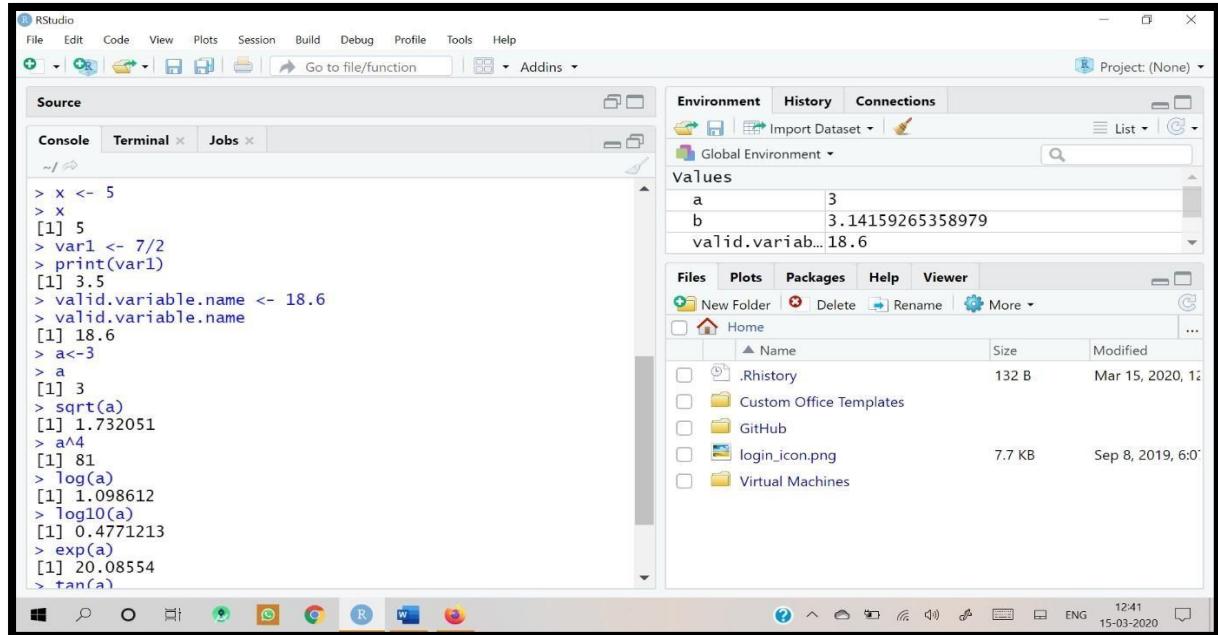
#### Basic Arithmetic in R:

The R Arithmetic operators include operators like Arithmetic Addition, Subtraction, Division, Multiplication, Exponent, Integer Division, and Modulus. All these R arithmetic operators are binary operators, which means they operate on two operands.

## Help command:

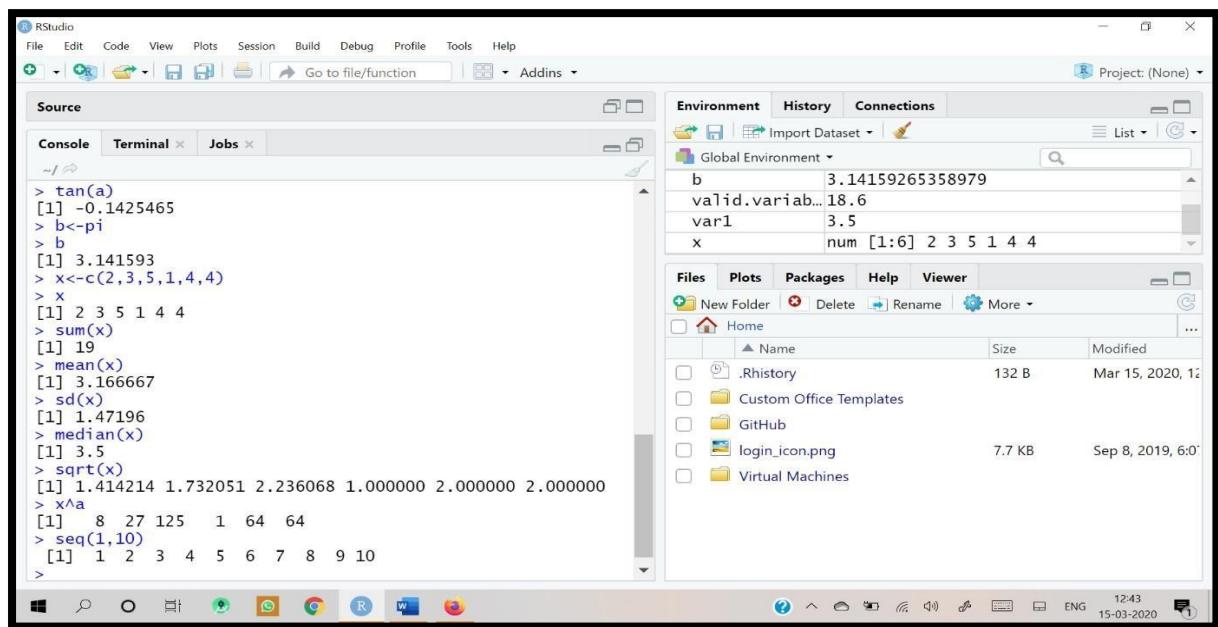
The help() function in R provides access to the documentation pages for R functions, data sets, and other objects, both for packages in the standard R distribution and for contributed packages.

## Output:



A screenshot of the RStudio interface. The Source pane shows R code being run in the Console tab. The code includes assignments like `x <- 5`, `var1 <- 7/2`, and various mathematical operations on `a`. The Global Environment pane shows variables `a` (3), `b` (3.14159265358979), and `valid.variable.name` (18.6). The Files pane shows a folder structure under Home, including files like `.Rhistory` and `login_icon.png`.

```
> x <- 5
> x
[1] 5
> var1 <- 7/2
> print(var1)
[1] 3.5
> valid.variable.name <- 18.6
> valid.variable.name
[1] 18.6
> a<-3
> a
[1] 3
> sqrt(a)
[1] 1.732051
> a^4
[1] 81
> log(a)
[1] 1.098612
> log10(a)
[1] 0.4771213
> exp(a)
[1] 20.08554
> tan(a)
```



A screenshot of the RStudio interface. The Source pane shows R code being run in the Console tab. The code includes mathematical operations like `tan(a)`, `sum(x)`, and `sqrt(x)`, along with sequence generation like `seq(1,10)`. The Global Environment pane shows variables `b` (3.14159265358979), `valid.variable.name` (18.6), `var1` (3.5), and `x` (a numeric vector [1:6] with values 2, 3, 5, 1, 4, 4). The Files pane shows a folder structure under Home, including files like `.Rhistory` and `login_icon.png`.

```
> tan(a)
[1] -0.1425465
> b<-pi
> b
[1] 3.141593
> x<-c(2,3,5,1,4,4)
> x
[1] 2 3 5 1 4 4
> sum(x)
[1] 19
> mean(x)
[1] 3.166667
> sd(x)
[1] 1.47196
> median(x)
[1] 3.5
> sqrt(x)
[1] 1.414214 1.732051 2.236068 1.000000 2.000000 2.000000
> x^a
[1] 8 27 125 1 64 64
> seq(1,10)
[1] 1 2 3 4 5 6 7 8 9 10
>
```

RStudio interface showing a console session and environment viewer.

Console output:

```
~/R
> sqrt(x)
[1] 1.414214 1.732051 2.236068 1.000000 2.000000 2.000000
> x^a
[1] 8 27 125 1 64 64
> seq(1,10)
[1] 1 2 3 4 5 6 7 8 9 10
> seq(1,10,2)
[1] 1 3 5 7 9
> seq(1:10)
[1] 1 2 3 4 5 6 7 8 9 10
> seq(1,10,by=2)
[1] 1 3 5 7 9
> y<-c(1:7)
> y
[1] 1 2 3 4 5 6 7
> z<-1:7
> z
[1] 1 2 3 4 5 6 7
> w<-c(1:12,0,-6)
> w
[1] 1 2 3 4 5 6 7 8 9 10 11 12 0 -6
>
```

Environment viewer:

W	num [1:14]	1 2 3 4 5 6 7 8 9 10 ...
X	num [1:6]	2 3 5 1 4 4
y	int [1:7]	1 2 3 4 5 6 7
z	int [1:7]	1 2 3 4 5 6 7

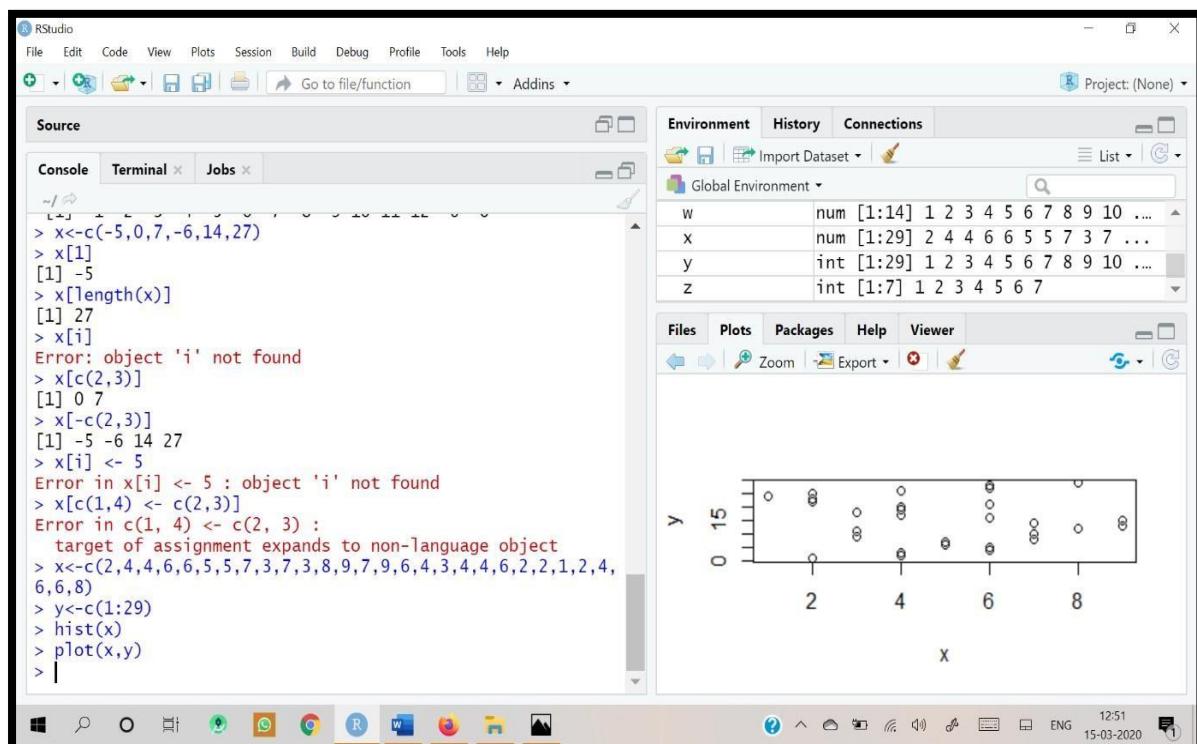
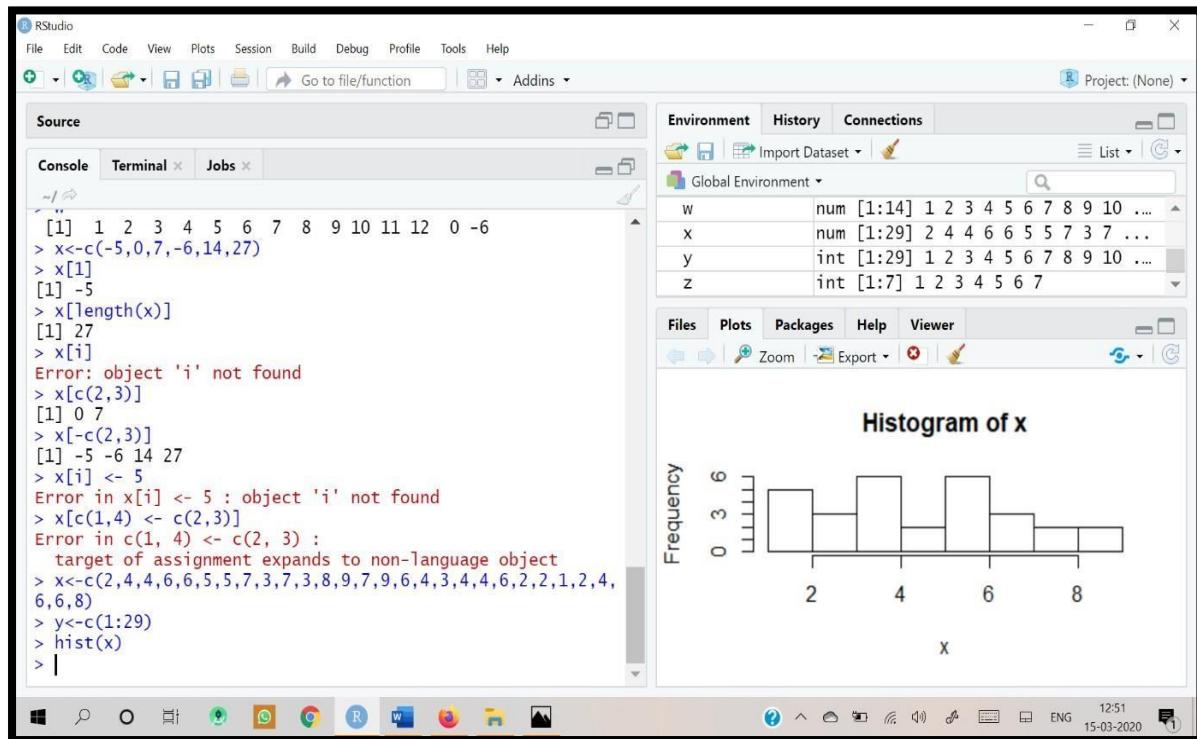
RStudio interface showing a console session and environment viewer.

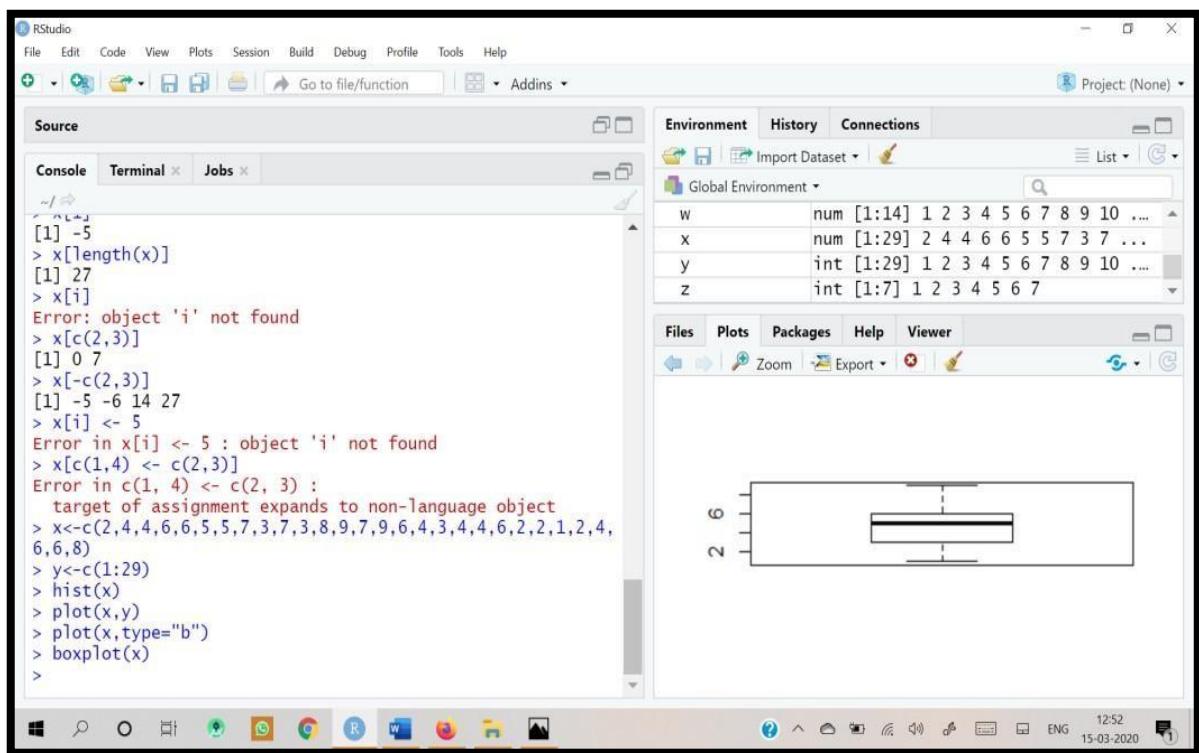
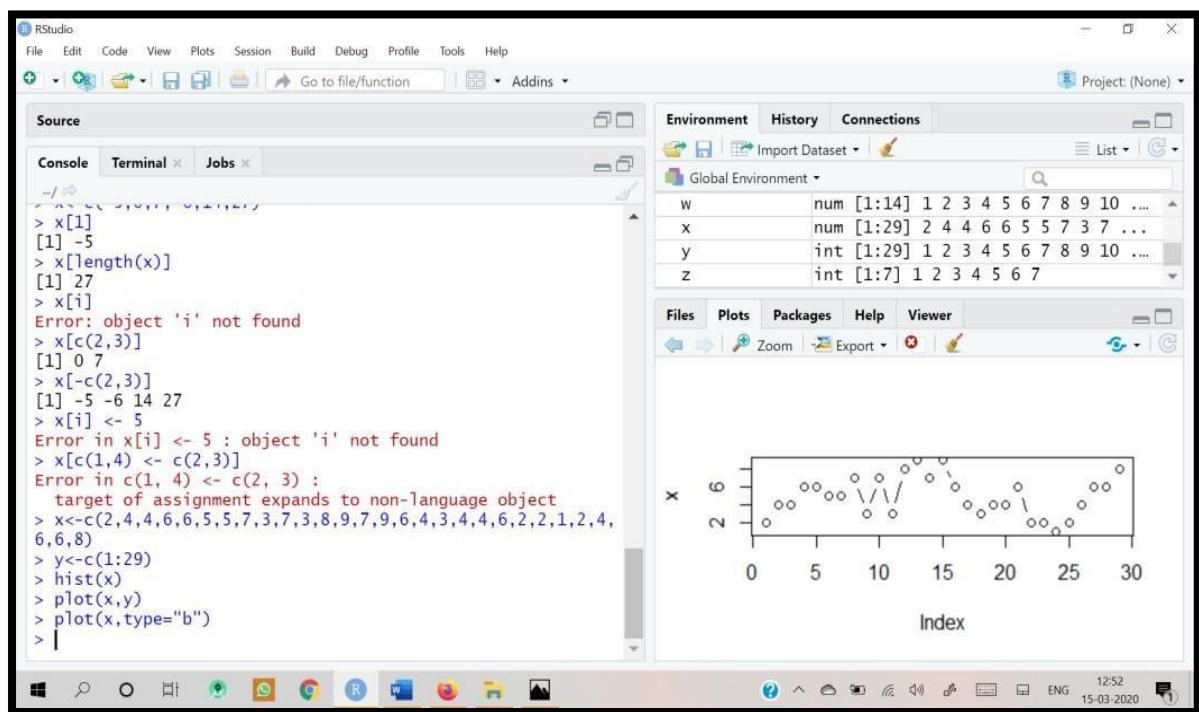
Console output:

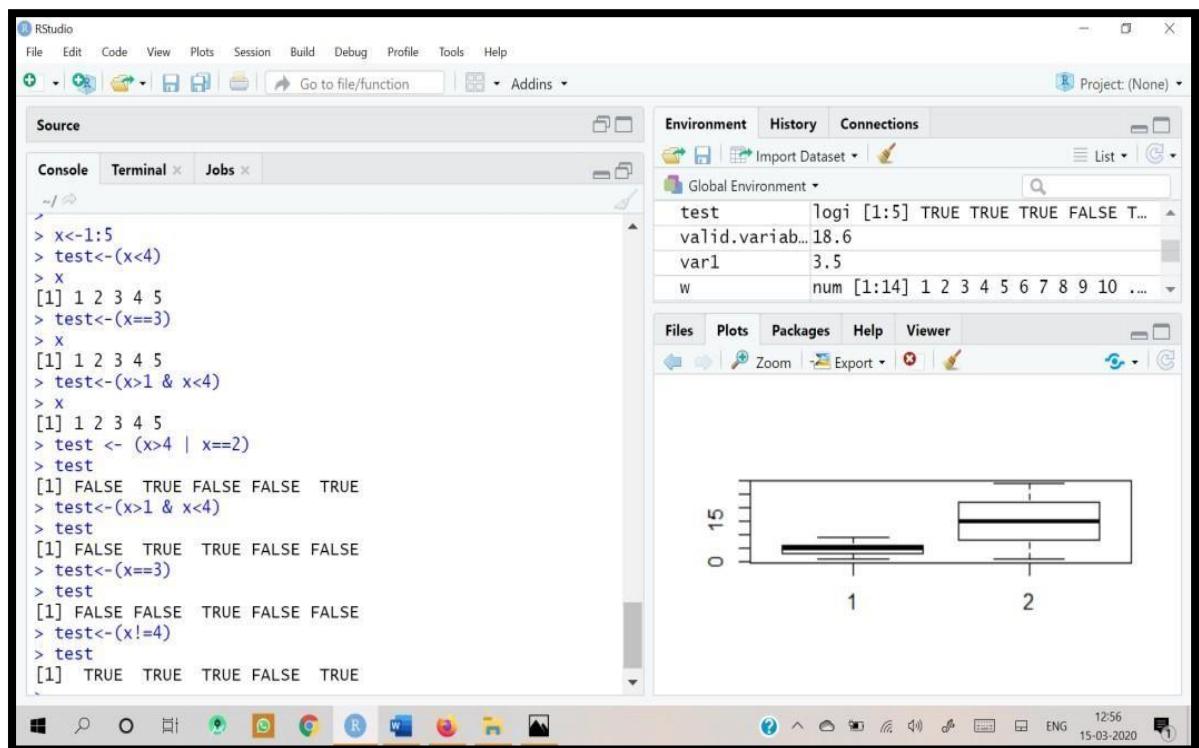
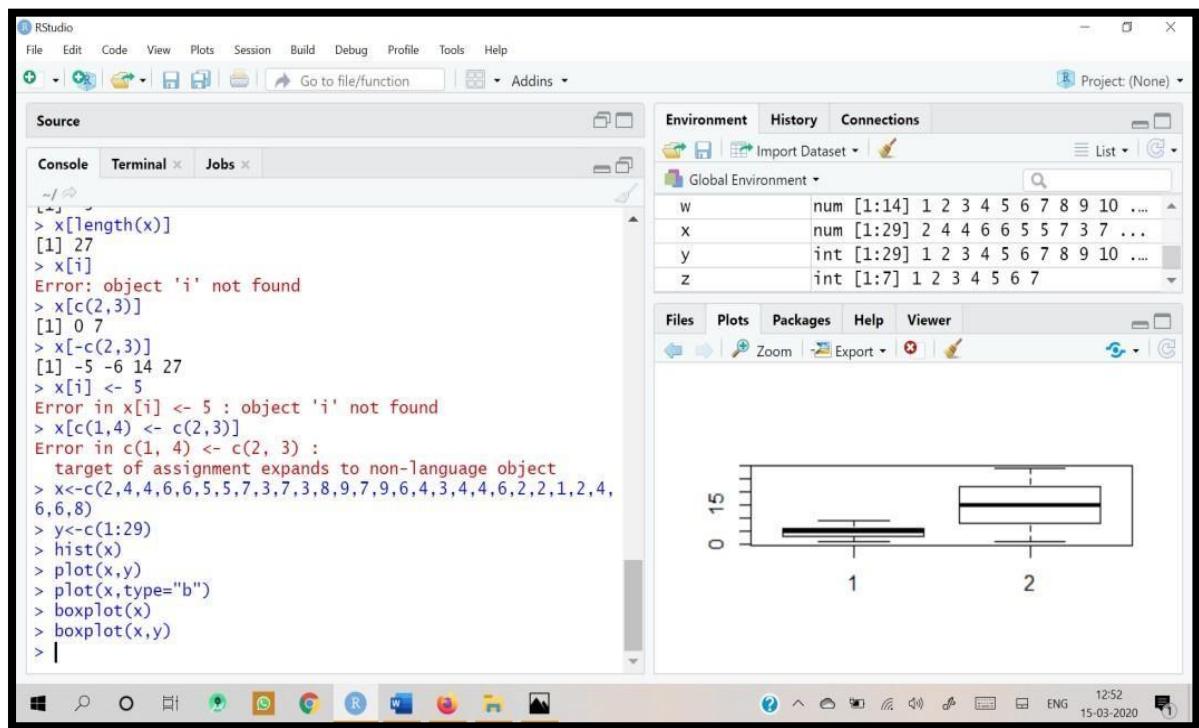
```
~/R
> seq(1,10,2)
[1] 1 3 5 7 9
> seq(1:10)
[1] 1 2 3 4 5 6 7 8 9 10
> seq(1,10,by=2)
[1] 1 3 5 7 9
> y<-c(1:7)
> y
[1] 1 2 3 4 5 6 7
> z<-1:7
> z
[1] 1 2 3 4 5 6 7
> w<-c(1:12,0,-6)
> w
[1] 1 2 3 4 5 6 7 8 9 10 11 12 0 -6
> x<-c(-5,0,7,-6,14,27)
> x[1]
[1] -5
> x[length(x)]
[1] 27
> x[i]
Error: object 'i' not found
> x[2:3]
[1] 7 -6
```

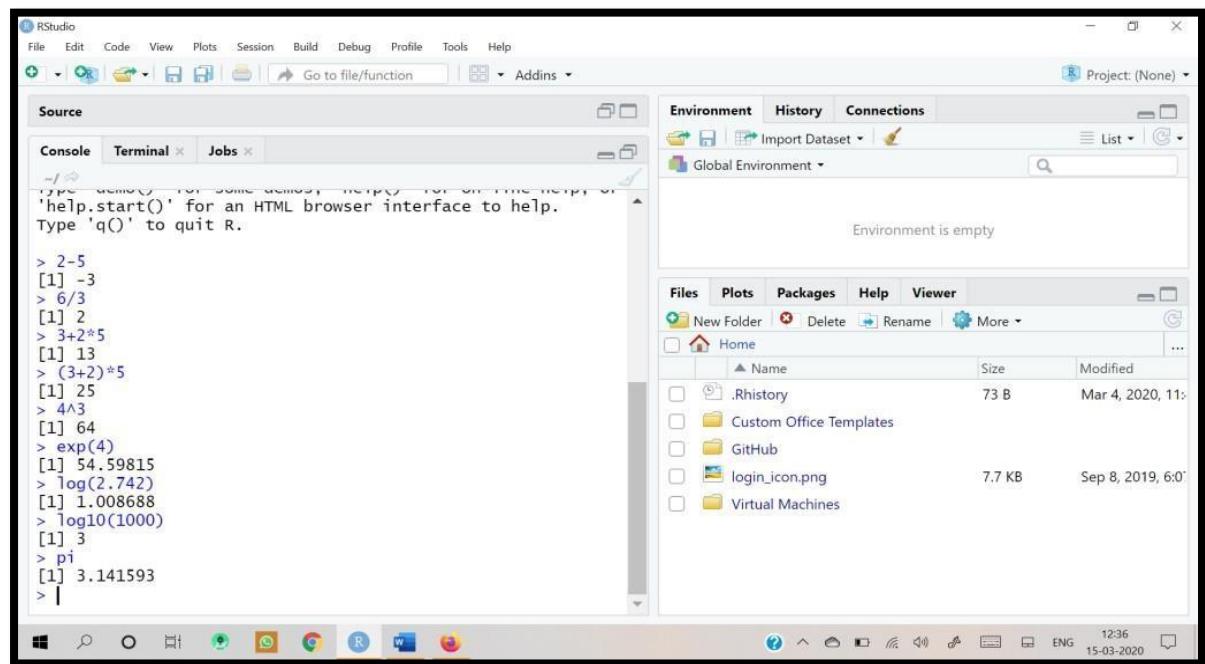
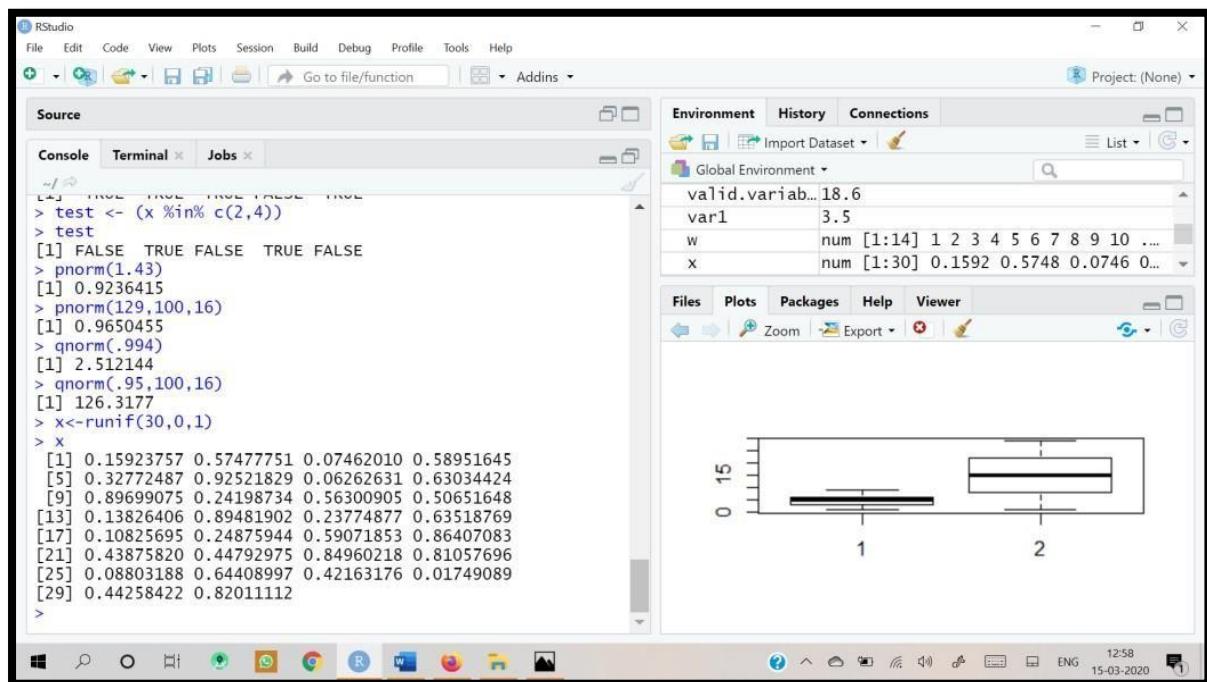
Environment viewer:

W	num [1:14]	1 2 3 4 5 6 7 8 9 10 ...
X	num [1:6]	-5 0 7 -6 14 27
y	int [1:7]	1 2 3 4 5 6 7
z	int [1:7]	1 2 3 4 5 6 7









**Conclusion:** Thus, we have studied the basic functionality of R.

## Experiment no:3

**Aim:** Basic data structures and data types in R-vector, matrices, list, data frames.

### Theory:

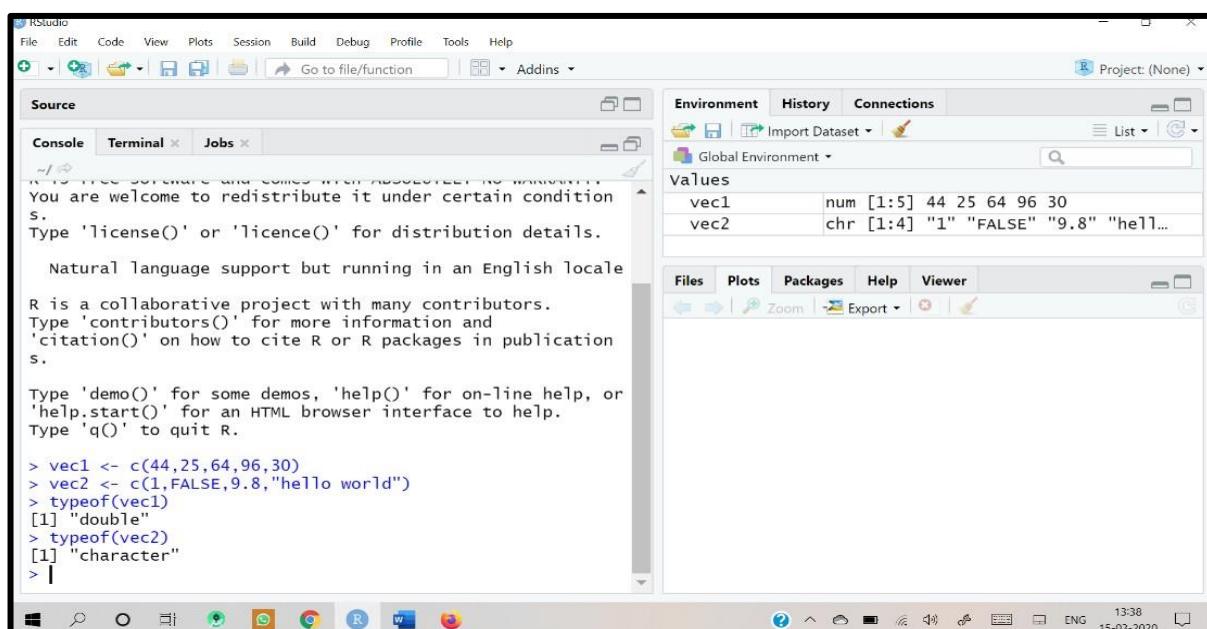
#### Data structures in R:

Data structures are used to store data in an organized fashion in order to make data manipulation and other data operations more efficient. There are five types of Data Structures in R Programming which are mentioned below:

- Vector
- List
- Matrix
- Data Frame

#### A. Vector:

Vector is one of the basic data structures in R programming. It is homogenous in nature, which means that it only contains elements of the same data type. Data types can be numeric, integer, character, complex or logical. The vector in R programming is created using the `c()` function. Coercion takes place in a vector from lower to top, if the elements passed are of different data types from Logical to Integer to Double to Character. The `typeof()` function is used to check the data type of the vector, and `class()` function is used to check the class of a vector.

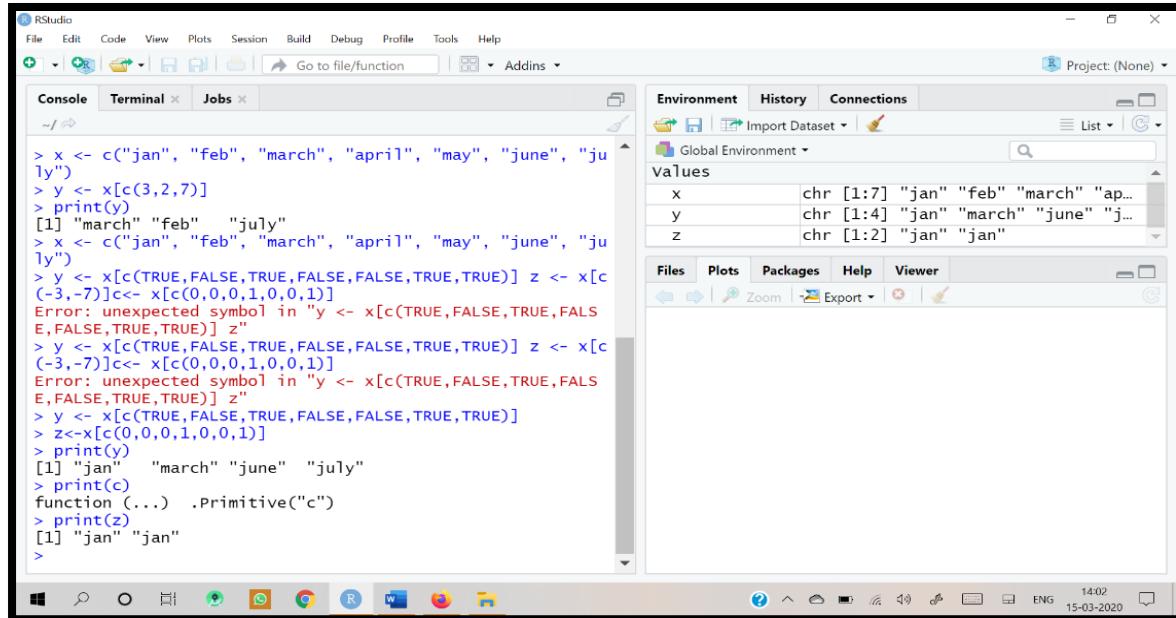


To delete a vector, we simply do the following:

```
Vec1<- NULL
```

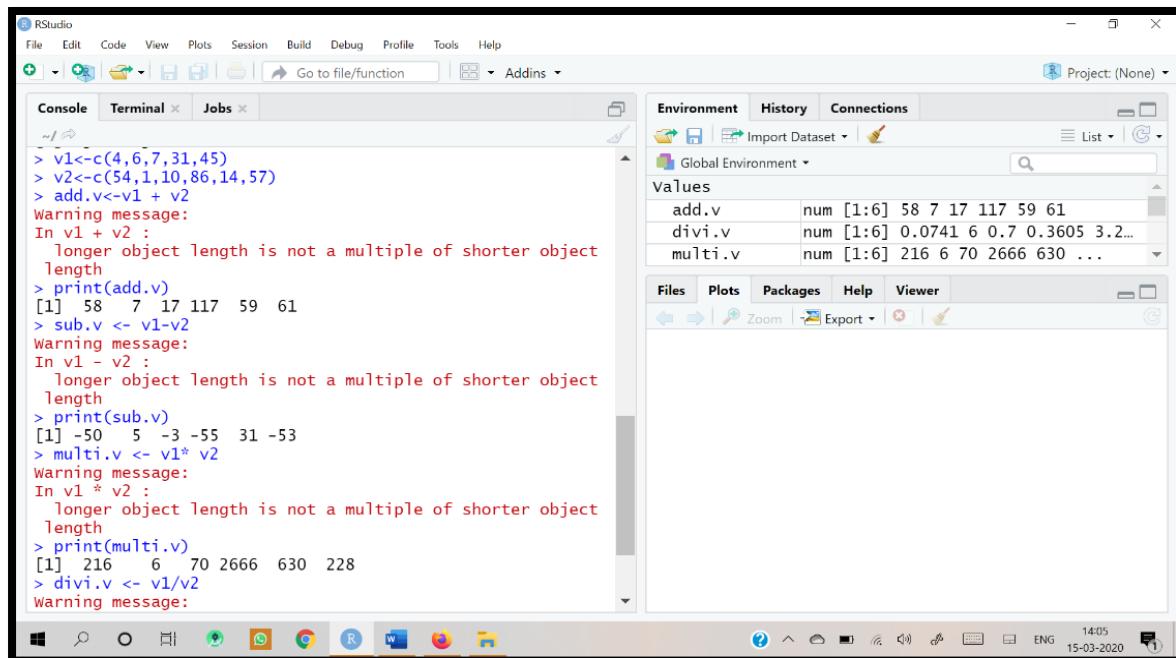
```
Vec2<- NULL
```

Elements of a vector can be accessed by using their respective indexes. [ ] brackets are used to specify indexes of the elements to be accessed.



The screenshot shows the RStudio interface with the following console output:

```
> x <- c("jan", "feb", "march", "april", "may", "june", "july")
> y <- x[c(3,2,7)]
> print(y)
[1] "march" "feb"   "july"
> x <- c("jan", "feb", "march", "april", "may", "june", "july")
> y <- x[c(TRUE,FALSE,TRUE,FALSE,FALSE,TRUE,TRUE)] z <- x[c(-3,-7)]<- x[c(0,0,0,1,0,0,1)]
Error: unexpected symbol in "y <- x[c(TRUE,FALSE,TRUE,FALSE,FALSE,TRUE,TRUE)] z"
> y <- x[c(TRUE,FALSE,TRUE,FALSE,FALSE,TRUE,TRUE)] z <- x[c(-3,-7)]<- x[c(0,0,0,1,0,0,1)]
Error: unexpected symbol in "y <- x[c(TRUE,FALSE,TRUE,FALSE,FALSE,TRUE,TRUE)] z"
> y <- x[c(TRUE,FALSE,TRUE,FALSE,FALSE,TRUE,TRUE)]
> z<-x[c(0,0,0,1,0,0,1)]
> print(y)
[1] "jan"   "march" "june"  "july"
> print(c)
function (...) .Primitive("c")
> print(z)
[1] "jan" "jan"
```



The screenshot shows the RStudio interface with the following console output:

```
> v1<-c(4,6,7,31,45)
> v2<-c(54,1,10,86,14,57)
> add.v<-v1 + v2
Warning message:
In v1 + v2 :
  longer object length is not a multiple of shorter object
  length
> print(add.v)
[1] 58 7 17 117 59 61
> sub.v <- v1-v2
Warning message:
In v1 - v2 :
  longer object length is not a multiple of shorter object
  length
> print(sub.v)
[1] -50 5 -3 -55 31 -53
> multi.v <- v1* v2
Warning message:
In v1 * v2 :
  longer object length is not a multiple of shorter object
  length
> print(multi.v)
[1] 216 6 70 2666 630 228
> divi.v <- v1/v2
Warning message:
```

The screenshot shows the RStudio interface with the following details:

- Console Tab:** Displays R code and its output. The code involves creating vectors, performing arithmetic operations, and sorting them. It also demonstrates how to sort a list of character strings.
- Environment Tab:** Shows the Global Environment pane with objects like `resort.v`, `sort.v`, `sub.v`, and `v`.
- File Bar:** Includes icons for File, Edit, Code, View, Plots, Session, Build, Debug, Profile, Tools, Help, and Addins.
- System Taskbar:** Shows icons for various applications like File Explorer, Task View, and a clock indicating 14:09 on 15-03-2020.

## B. List:

A list in R programming is a non-homogenous data structure, which implies that it can contain elements of different data types. It accepts numbers, characters, lists, and even matrices and functions inside it. It is created using the `list()` function.

The screenshot shows the RStudio interface with the following details:

- Console Tab:** Displays R code for creating a list named `list1` with various elements including character strings, a vector, logical values, and numbers.
- Environment Tab:** Shows the Global Environment pane with two lists: `list` and `list1`.
- File Bar:** Includes icons for File, Edit, Code, View, Plots, Session, Build, Debug, Profile, Tools, Help, and Addins.
- System Taskbar:** Shows icons for various applications like File Explorer, Task View, and a clock indicating 14:30 on 15-03-2020.

The screenshot shows the RStudio interface. In the Console tab, the following R code is run:

```
> list2<-list(matrix(c(3,9,5,1,-2,8),nrow=2) , c("jan","feb","mar"),list(3,4,5))
> print(list2[1])
[[1]]
 [,1] [,2] [,3]
[1,] 3 5 -2
[2,] 9 1 8

> print(list2[2])
[[1]]
[1] "green"

> print(list2[3])
[[1]]
[[1]][[1]]
[1] 3

[[1]][[2]]
[1] 4

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

The Environment pane shows three objects: list (List of 7), list1 (List of 7), and list2 (List of 3). The Global Environment pane shows the same three objects.

The screenshot shows the RStudio interface. In the Console tab, the following R code is run:

```
[[1]][[2]]
[1] 4

[[1]][[3]]
[1] 5

> list2<-list(matrix(c(3,9,5,1,-2,8),nrow=2) , c("jan","feb","mar"),list(3,4,5))
> list2[4]<-NULL
> print(list2[4])
[[1]]
[1] "hello"

> list2[4]<-NULL
> print(list2[4])
[[1]]
NULL

> list2[3]<-"element updated"
> print(list2[3])
[[1]]
[1] "element updated"
```

The Environment pane shows three objects: list (List of 7), list1 (List of 7), and list2 (List of 3). The Global Environment pane shows the same three objects.

## C. Matrix:

The matrix in R programming is a 2-dimensional data structure that is homogenous in nature, which means that it only accepts elements of the same data type. Coercion takes place if elements of different data types are passed. It is created using the `matrix()` function. The basic syntax to create a matrix is given below: `matrix(data, nrow, ncol, byrow, dimnames)`

where,

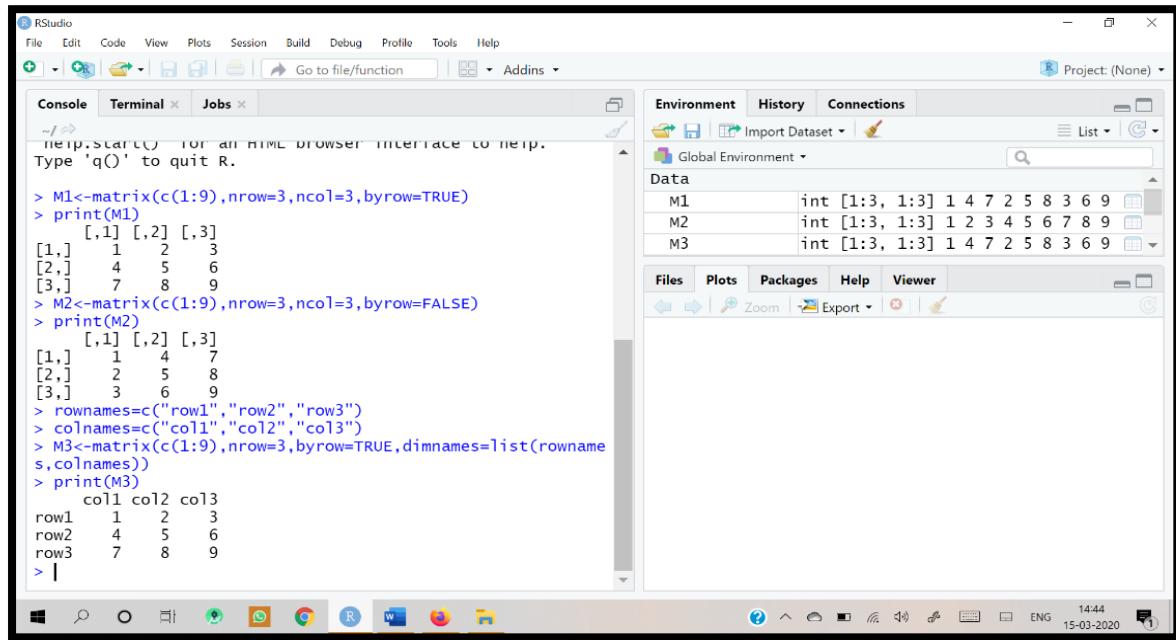
data = the input element of a matrix given as a vector.

nrow = the number of rows to be created.

ncol = the number of columns to be created.

byrow = the row-wise arrangement of the elements instead of column-wise

dimnames = the names of columns/rows to be created.

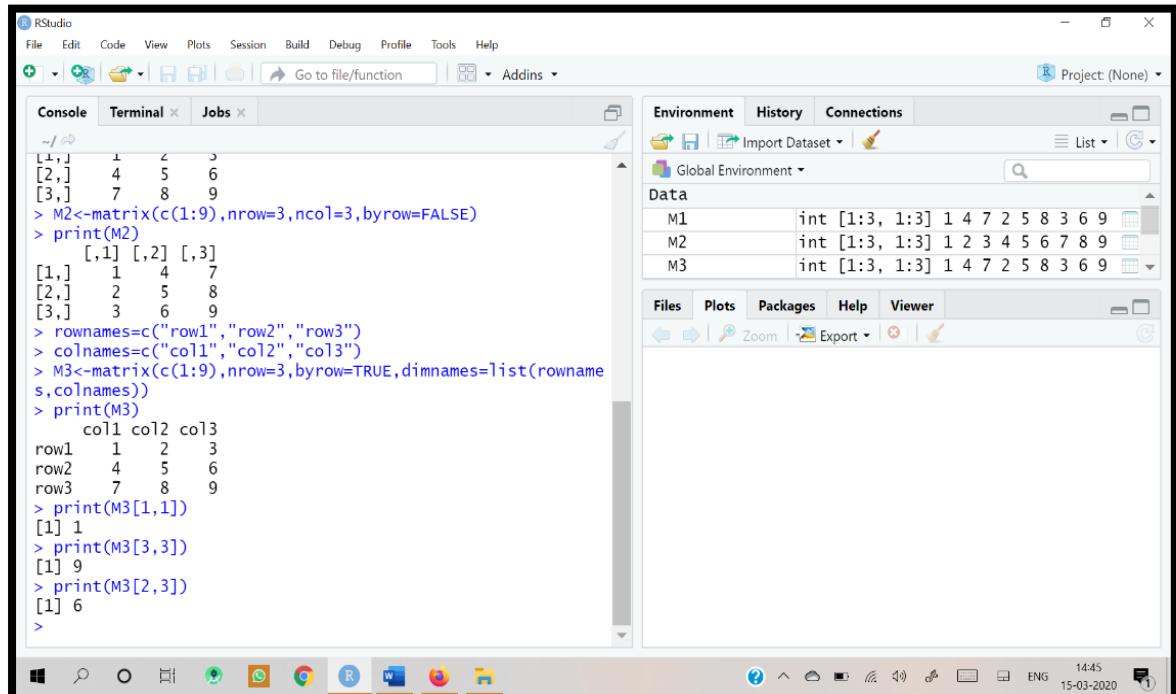


The screenshot shows the RStudio interface with the following code in the Console:

```
> M1<-matrix(c(1:9),nrow=3,ncol=3,byrow=TRUE)
> print(M1)
 [,1] [,2] [,3]
[1,] 1 2 3
[2,] 4 5 6
[3,] 7 8 9
> M2<-matrix(c(1:9),nrow=3,ncol=3,byrow=FALSE)
> print(M2)
 [,1] [,2] [,3]
[1,] 1 4 7
[2,] 2 5 8
[3,] 3 6 9
> rownames=c("row1","row2","row3")
> colnames=c("col1","col2","col3")
> M3<-matrix(c(1:9),nrow=3,byrow=TRUE,dimnames=list(rownames,colnames))
> print(M3)
     col1 col2 col3
row1 1 2 3
row2 4 5 6
row3 7 8 9
> |
```

The Global Environment pane shows three matrices:

Matrix	Type	Dimensions	Elements
M1	int	[1:3, 1:3]	1 4 7 2 5 8 3 6 9
M2	int	[1:3, 1:3]	1 2 3 4 5 6 7 8 9
M3	int	[1:3, 1:3]	1 4 7 2 5 8 3 6 9



The screenshot shows the RStudio interface with the following code in the Console:

```
> L1<
[1,] 1 2 3
[2,] 4 5 6
[3,] 7 8 9
> M2<-matrix(c(1:9),nrow=3,ncol=3,byrow=FALSE)
> print(M2)
 [,1] [,2] [,3]
[1,] 1 4 7
[2,] 2 5 8
[3,] 3 6 9
> rownames=c("row1","row2","row3")
> colnames=c("col1","col2","col3")
> M3<-matrix(c(1:9),nrow=3,byrow=TRUE,dimnames=list(rownames,colnames))
> print(M3)
     col1 col2 col3
row1 1 2 3
row2 4 5 6
row3 7 8 9
> print(M3[1,1])
[1] 1
> print(M3[3,3])
[1] 9
> print(M3[2,3])
[1] 6
>
```

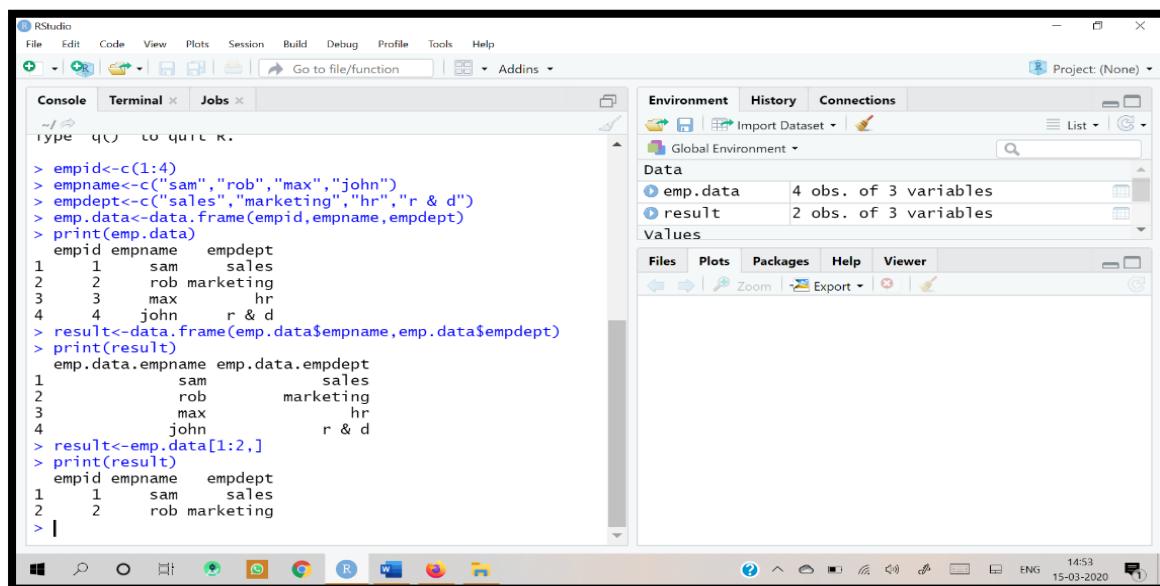
The Global Environment pane shows three matrices:

Matrix	Type	Dimensions	Elements
M1	int	[1:3, 1:3]	1 4 7 2 5 8 3 6 9
M2	int	[1:3, 1:3]	1 2 3 4 5 6 7 8 9
M3	int	[1:3, 1:3]	1 4 7 2 5 8 3 6 9

## D. Data Frame:

A data frame in R programming is a 2-dimensional array-like structure that also resembles a table, in which each column contains values of one variable and each row contains one set of values from each column. A data frame has the following characteristics:

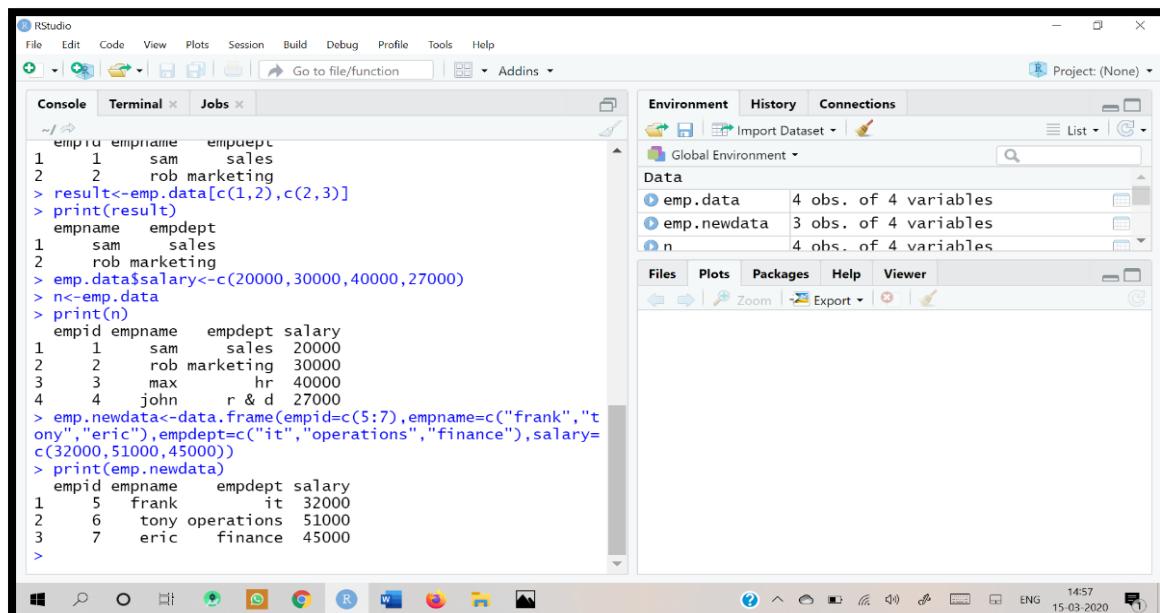
- The column names of a data frame should not be empty.
- Row names should be unique.
- Data stored in a data frame can be numeric, factor or character type.
- Each column should contain the same number of data items.



The screenshot shows an RStudio interface. In the Console tab, the following R code is run:

```
type q() to quit R.  
> empid<-c(1:4)  
> empname<-c("sam","rob","max","john")  
> empdept<-c("sales","marketing","hr","r & d")  
> emp.data<-data.frame(empid,empname,empdept)  
> print(emp.data)  
  empid empname   empdept  
1     1    sam    sales  
2     2    rob  marketing  
3     3    max       hr  
4     4   john  r & d  
> result<-data.frame(emp.data$empname,emp.data$empdept)  
> print(result)  
  emp.data.empname emp.data.empdept  
1             sam           sales  
2             rob         marketing  
3             max            hr  
4             john          r & d  
> result<-emp.data[1:2,]  
> print(result)  
  empid empname   empdept  
1     1    sam    sales  
2     2    rob  marketing
```

In the Environment pane, two objects are listed: `emp.data` (4 obs. of 3 variables) and `result` (2 obs. of 3 variables). The system tray at the bottom right shows the date as 15-03-2020 and the time as 14:53.



The screenshot shows an RStudio interface. In the Console tab, the following R code is run:

```
empid empname   empdept  
1     1    sam    sales  
2     2    rob  marketing  
> result<-emp.data[c(1,2),c(2,3)]  
> print(result)  
  empname   empdept  
1    sam    sales  
2   rob  marketing  
> emp.data$salary<-c(20000,30000,40000,27000)  
> n<-emp.data  
> print(n)  
  empid empname   empdept salary  
1     1    sam    sales  20000  
2     2    rob  marketing 30000  
3     3    max       hr  40000  
4     4   john  r & d  27000  
> emp.newdata<-data.frame(empid=c(5:7),empname=c("frank","tony","eric"),empdept=c("it","operations","finance"),salary=c(32000,51000,45000))  
> print(emp.newdata)  
  empid empname   empdept salary  
1     5   frank      it  32000  
2     6   tony operations 51000  
3     7    eric   finance  45000
```

In the Environment pane, three objects are listed: `emp.data` (4 obs. of 4 variables), `emp.newdata` (3 obs. of 4 variables), and `n` (4 obs. of 4 variables). The system tray at the bottom right shows the date as 15-03-2020 and the time as 14:57.

The screenshot shows the RStudio interface. The console window displays R code and its output. The code creates three data frames: emp.data, emp.newdata, and emp.finaldata. The emp.data frame has 4 observations and 4 variables. The emp.newdata frame has 3 observations and 4 variables. The emp.finaldata frame has 7 observations and 4 variables. The Data View window on the right shows the contents of these three data frames.

```

> print(emp)
   empid empname    empdept salary
1      1     sam    sales 20000
2      2     rob  marketing 30000
3      3     max        hr 40000
4      4    john  r & d 27000
> emp.newdata<-data.frame(empid=c(5:7),empname=c("frank","tony","eric"),empdept=c("it","operations","finance"),salary=c(32000,51000,45000))
> print(emp.newdata)
   empid empname    empdept salary
1      5   frank        it 32000
2      6   tony  operations 51000
3      7    eric   finance 45000
> emp.finaldata<-rbind(emp.data,emp.newdata)
> print(emp.finaldata)
   empid empname    empdept salary
1      1     sam    sales 20000
2      2     rob  marketing 30000
3      3     max        hr 40000
4      4    john  r & d 27000
5      5   frank        it 32000
6      6   tony  operations 51000
7      7    eric   finance 45000
>

```

## Data-Types in R:

There are many basic data types in R, which are of frequent occurrence in coding R calculations and programs. Here is the list of all the data types provided by R:

- Numeric
- Integer
- Complex
- Logical
- Character

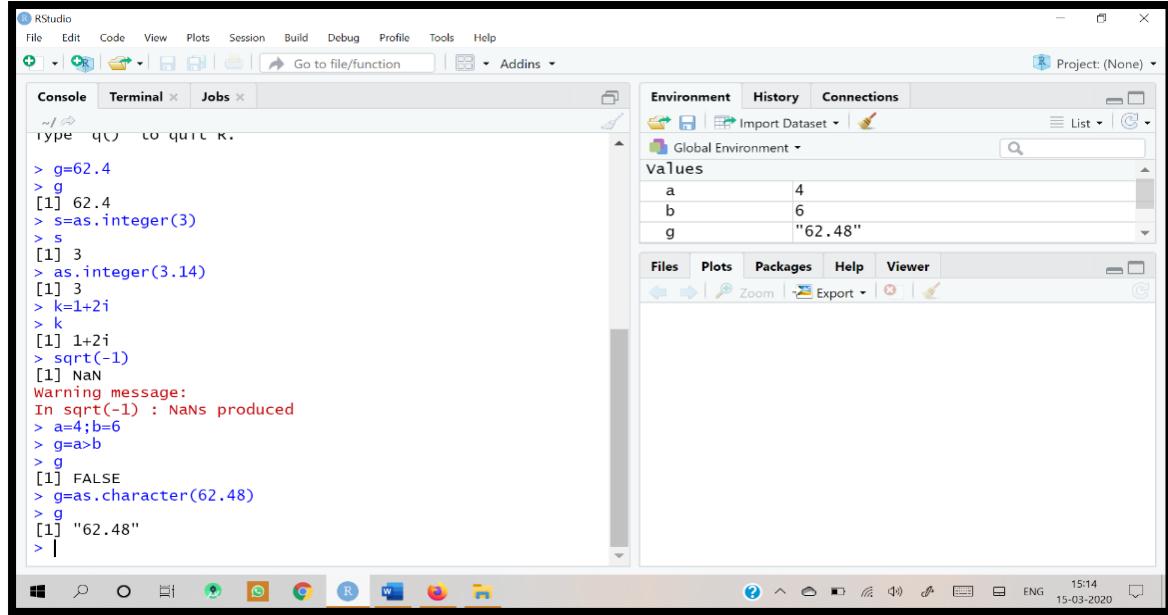
**a) Number data type:** Decimal values are referred to as numeric data types in R. This is the default working out data type. If you assign a decimal value for any variable x, x will become a numeric type.

**b) Integer data type:** If you want to create an integer variable in R, you have to invoke the `as.integer()` function to define any integer type data. You can be certain that y is definitely an integer by applying the `is.integer()` function.

**c) Complex data type:** A complex value for coding in R can be defined using the pure imaginary values 'i'.

**d) Logical data type:** A logical value is mostly created when a comparison between variables are done.

**e) Character data type:** A character object can be used for representing string values in R. You have to convert objects into character values using the `as.character()` function within your code.

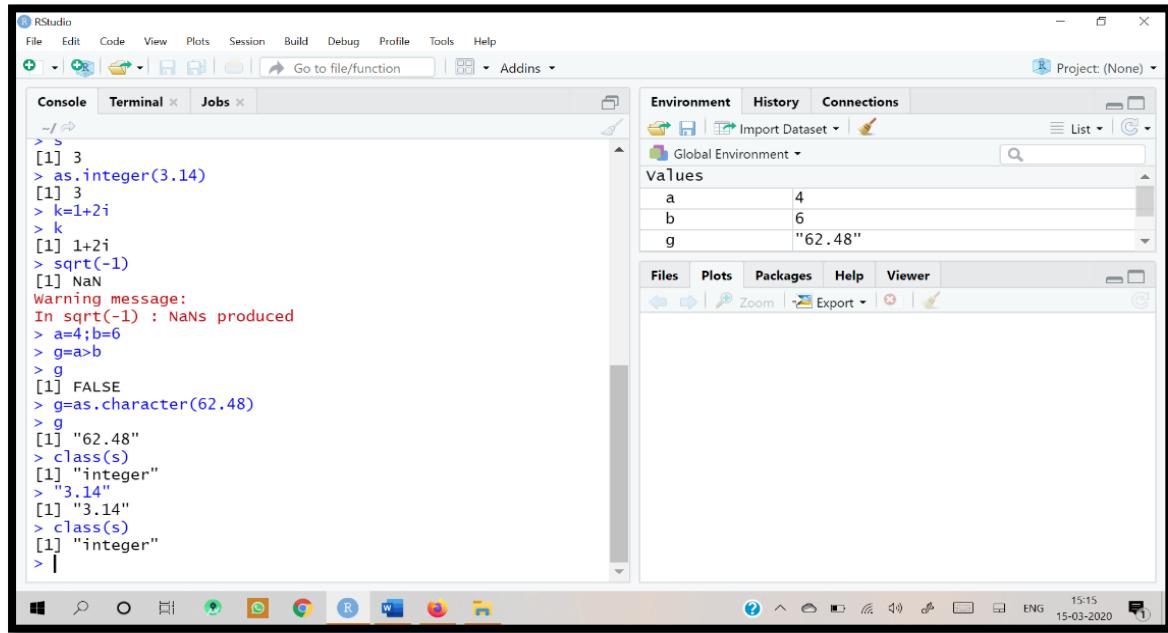


The screenshot shows the RStudio interface with the following session history in the Console tab:

```
> g=62.4
> g
[1] 62.4
> s=as.integer(3)
> s
[1] 3
> as.integer(3.14)
[1] 3
> k=1+2i
> k
[1] 1+2i
> sqrt(-1)
[1] NaN
Warning message:
In sqrt(-1) : NaNs produced
> a=4;b=6
> g=a>b
> g
[1] FALSE
> g=as.character(62.48)
> g
[1] "62.48"
> |
```

The Global Environment pane shows variables `a`, `b`, and `g`:

Values	a	4
b	6	
g	"62.48"	



The screenshot shows the RStudio interface with the following session history in the Console tab:

```
> s
[1] 3
> as.integer(3.14)
[1] 3
> k=1+2i
> k
[1] 1+2i
> sqrt(-1)
[1] NaN
Warning message:
In sqrt(-1) : NaNs produced
> a=4;b=6
> g=a>b
> g
[1] FALSE
> g=as.character(62.48)
> g
[1] "62.48"
> class(s)
[1] "integer"
> "3.14"
[1] "3.14"
> class(s)
[1] "integer"
> |
```

The Global Environment pane shows variables `a`, `b`, and `g`:

Values	a	4
b	6	
g	"62.48"	

## **Conclusion:**

Thus, we have studied data types and data structures provided by R.

## Experiment no:4

**Aim:** Programming constructs in R-loops, conditional executions.

### Theory:

#### A. For Loop:

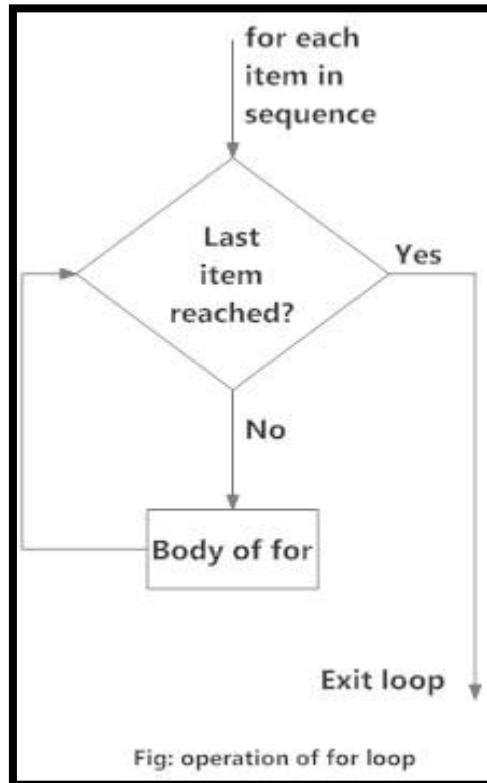
A for loop is used to iterate over a vector in programming.

#### Syntax:

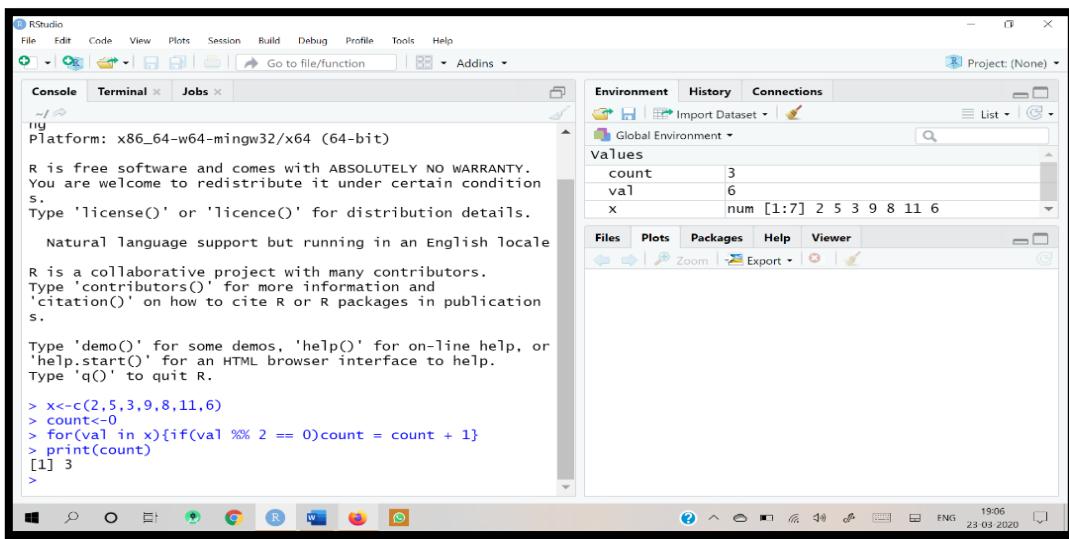
```
for (val in sequence)
{
  statement
}
```

Where, sequence is a vector and val takes on each of its value during the loop. In each iteration, statement is evaluated.

#### Flowchart for For Loop:



## Output:



The screenshot shows the RStudio interface. The left pane displays the R console with the following output:

```
R> Platform: x86_64-w64-mingw32/x64 (64-bit)
R is free software and comes with ABSOLUTELY NO WARRANTY.
You are welcome to redistribute it under certain conditions.
Type 'license()' or 'licence()' for distribution details.

Natural language support is running in an English locale.

R is a collaborative project with many contributors.
Type 'contributors()' for more information and
'citation()' on how to cite R or R packages in publications.

Type 'demo()' for some demos, 'help()' for on-line help, or
'help.start()' for an HTML browser interface to help.
Type 'q()' to quit R.

> x<-c(2,5,3,9,8,11,6)
> count<-0
> for(val in x){if(val %% 2 == 0)count = count + 1}
> print(count)
[1] 3
>
```

The right pane shows the Environment tab with the following variables:

Values	count	3
val	6	
x	num [1:7]	2 5 3 9 8 11 6

## B. R while Loop:

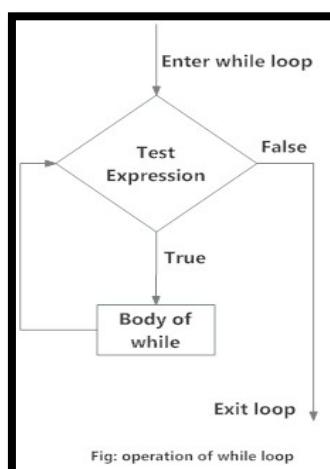
In R programming, while loops are used to loop until a specific condition is met.

### Syntax:

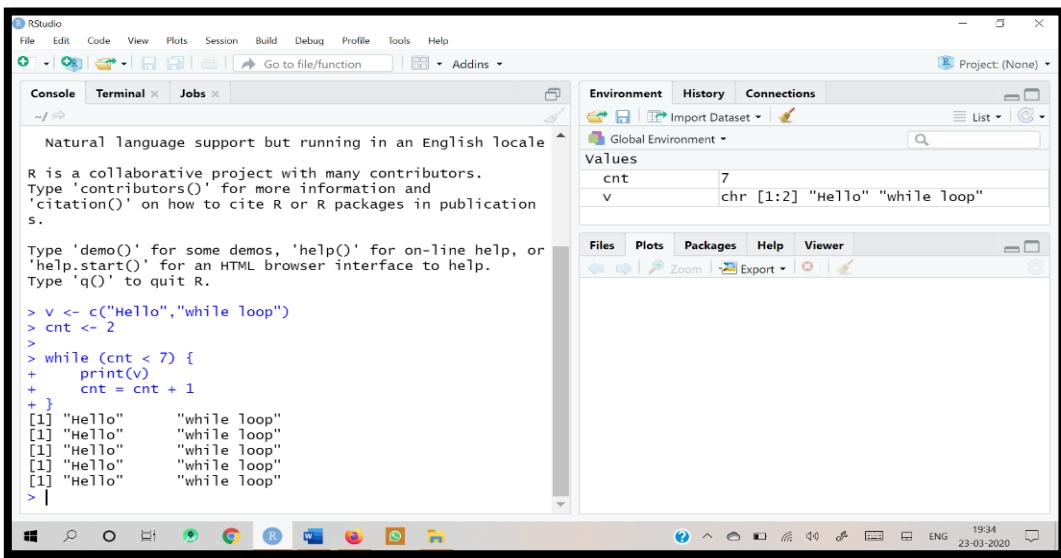
```
while (test_expression)
{
  statement
}
```

Where test\_expression is evaluated and the body of the loop is entered if the result is TRUE. The statements inside the loop are executed and the flow returns to evaluate the test\_expression again. This is repeated each time until test\_expression evaluates to FALSE, in which case, the loop exits.

### Flowchart for while loop:



## Output:



The screenshot shows the RStudio interface. The console tab is active, displaying the following R code and its output:

```
Natural language support but running in an English locale
R is a collaborative project with many contributors.
Type 'contributors()' for more information and
'citation()' on how to cite R or R packages in publications.

Type 'demo()' for some demos, 'help()' for on-line help, or
'help.start()' for an HTML browser interface to help.
Type 'q()' to quit R.

> v <- c("Hello","while loop")
> cnt <- 2
>
> while (cnt < 7) {
+   print(v)
+   cnt = cnt + 1
+ }
[1] "Hello"      "while loop"
>
```

The environment pane shows two variables: 'cnt' with a value of 7 and 'v' which is a character vector containing the strings "Hello" and "while loop".

## C. R ifelse () Function:

Vectors form the basic building block of R programming. Most of the functions in R take vector as input and output a resultant vector. This vectorization of code, will be much faster than applying the same function to each element of the vector individually. Similar to this concept, there is a vector equivalent form of the if...else statement in R, the ifelse() function.

### Syntax:

ifelse (test\_expression, x, y)

Here, test\_expression must be a logical vector (or an object that can be coerced to logical). The return value is a vector with the same length as test\_expression.

This returned vector has element from x if the corresponding value of test\_expression is TRUE or from y if the corresponding value of test\_expression is FALSE.

This is to say, the i-th element of result will be x[i] if test\_expression[i] is TRUE else it will take the value of y[i].

The vectors x and y are recycled whenever necessary.

## Output:

R version 3.6.2 (2019-12-12) -- "Dark and Stormy Night"  
Copyright (C) 2019 The R Foundation for Statistical Computing  
Platform: x86\_64-w64-mingw32/x64 (64-bit)  
R is free software and comes with ABSOLUTELY NO WARRANTY.  
You are welcome to redistribute it under certain conditions.  
Type 'license()' or 'licence()' for distribution details.  
Natural language support but running in an English locale  
R is a collaborative project with many contributors.  
Type 'contributors()' for more information and  
'citation()' on how to cite R or R packages in publications.  
Type 'demo()' for some demos, 'help()' for on-line help, or  
'help.start()' for an HTML browser interface to help.  
Type 'q()' to quit R.

```
> a = c(5,7,2,9)
> ifelse(a %% 2 == 0,"even","odd")
[1] "odd" "odd" "even" "odd"
>
```

## D. R break and next statement:

In R programming, a normal looping sequence can be altered using the break or the next statement.

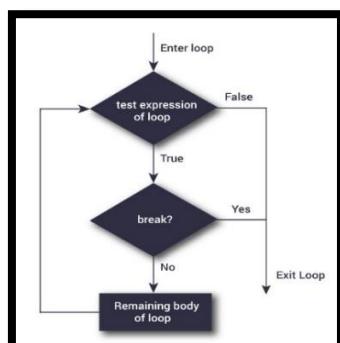
### i. Break statement:

A break statement is used inside a loop (repeat, for, while) to stop the iterations and flow the control outside of the loop. In a nested looping situation, where there is a loop inside another loop, this statement exits from the innermost loop that is being evaluated.

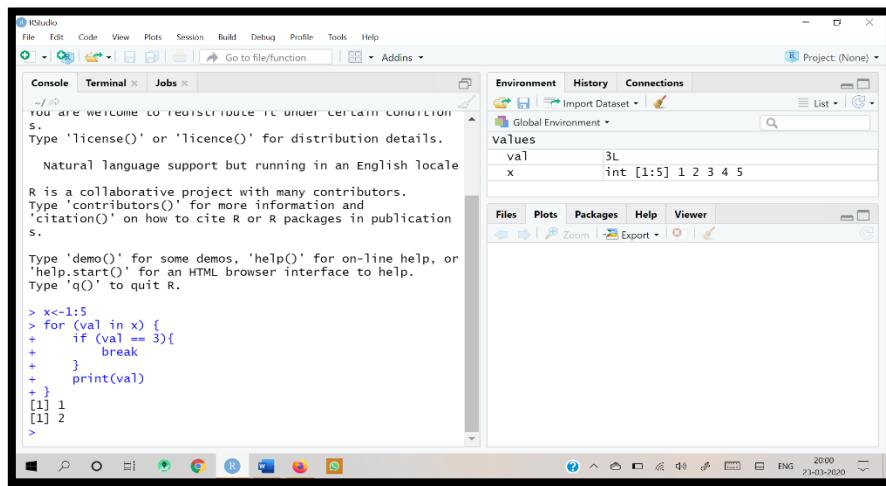
### Syntax:

```
if (test_expression) {
  break
}
```

### Flowchart for break statement:



## Output:



The screenshot shows the RStudio interface with the following session output in the Console tab:

```
R> x<-1:5
R> for (val in x) {
+   if (val == 3){
+     break
+   }
+   print(val)
+ }
[1] 1
[1] 2
```

The Environment tab shows the variable `x` as an integer vector from 1 to 5, and `val` as 3L.

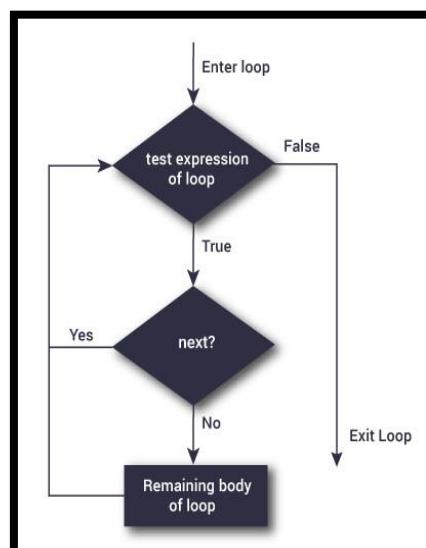
### **ii. Next Statement:**

A next statement is useful when we want to skip the current iteration of a loop without terminating it. On encountering next, the R parser skips further evaluation and starts next iteration of the loop.

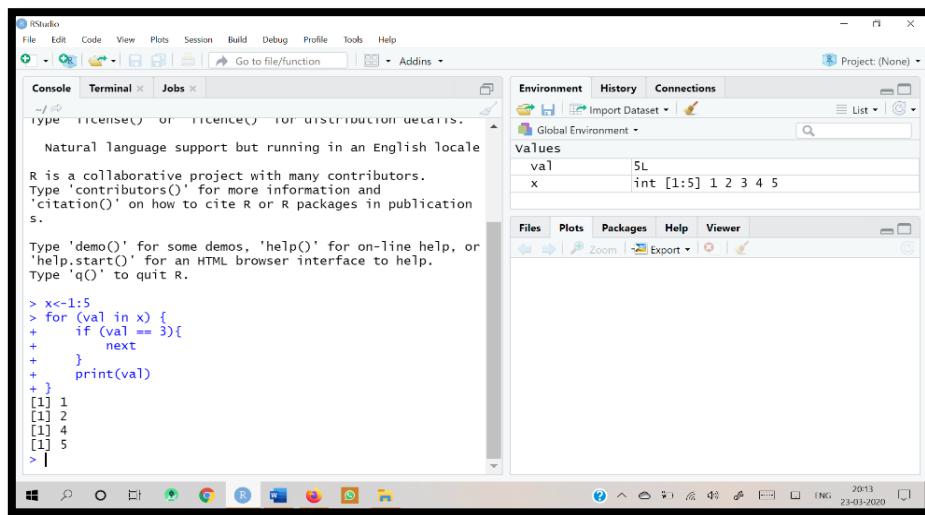
#### **Syntax:**

```
if (test_condition) {
  next
}
```

#### **Flowchart for next statement:**



## Output:



A screenshot of the RStudio interface. The console tab shows the following R code and its output:

```
> x<-1:5
> for (val in x) {
+   if (val == 3){
+     next
+   }
+   print(val)
+ }
[1] 1
[1] 2
[1] 4
[1] 5
```

The environment pane shows:

val	5L
x	int [1:5] 1 2 3 4 5

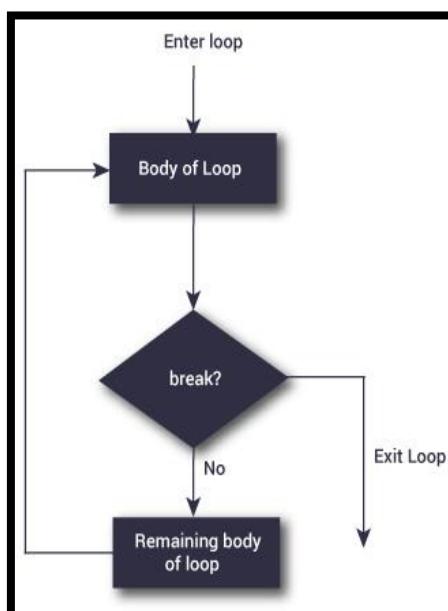
## E. R repeat loop:

A repeat loop is used to iterate over a block of code multiple number of times. There is no condition check in repeat loop to exit the loop. We must ourselves put a condition explicitly inside the body of the loop and use the break statement to exit the loop. Failing to do so will result into an infinite loop.

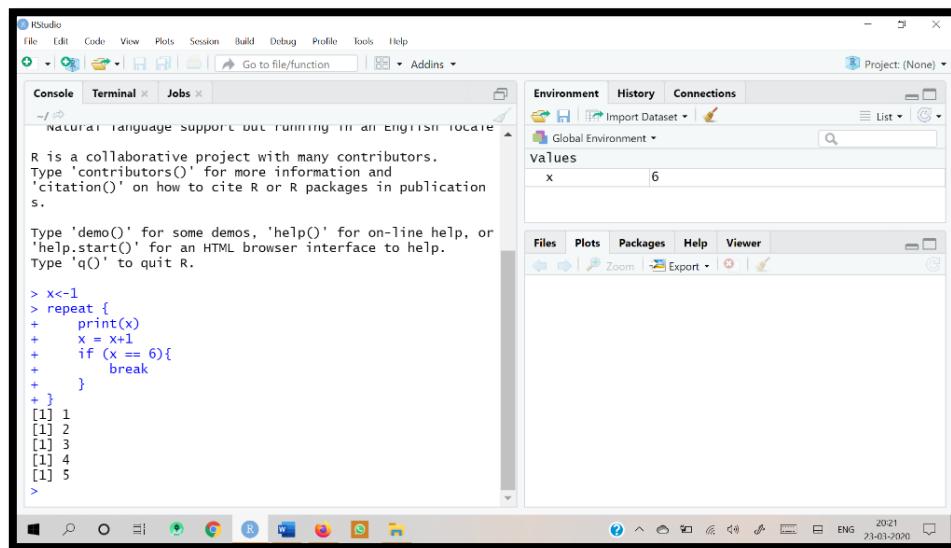
### Syntax:

```
repeat {
  statement
}
```

### Flowchart:



## Output:



## R Functions:

Functions are used to logically break our code into simpler parts which become easy to maintain and understand.

## Syntax:

```
func_name <- function (argument) {
  statement
}
```

- Here, we can see that the reserved word `function` is used to declare a function in R.
- The statements within the curly braces form the body of the function. These braces are optional if the body contains only a single expression.
- Finally, this function object is given a name by assigning it to a variable, `func_name`.

## Conclusion:

Thus, we have studied loops and functions provided by R.

## Experiment no:5

**Aim:** Tables with labels in R - multiple response variables excel function translations.

### Theory:

expss computes and displays tables with support for ‘SPSS’-style labels, multiple / nested banners, weights, multiple-response variables and significance testing. There are facilities for nice output of tables in ‘knitr’, R notebooks, ‘Shiny’ and ‘Jupyter’ notebooks. Proper methods for labelled variables add value labels support to base R functions and to some functions from other packages. Additionally, the package offers useful functions for data processing in marketing research / social surveys - popular data transformation functions from ‘SPSS’ Statistics (‘RECODE’, ‘COUNT’, ‘COMPUTE’, ‘DO IF’, etc.) and ‘Excel’ (‘COUNTIF’, ‘VLOOKUP’, etc.). Package is intended to help people to move data processing from ‘Excel’/ ‘SPSS’ to R. See examples below. You can get help about any function by typing `function_name` in the R console.

### Installation:

expss is on CRAN, so for installation you can print in the console `install.packages("expss")`.

```
library(expss)
data(mtcars)
mtcars = apply_labels(mtcars,
  mpg = "Miles/(US) gallon",
  cyl = "Number of cylinders",
  disp = "Displacement (cu.in.)",
  hp = "Gross horsepower",
  drat = "Rear axle ratio",
  wt = "Weight (1000 lbs)",
  qsec = "1/4 mile time",
  vs = "Engine",
  vs = c("V-engine" = 0,
  "Straight engine" = 1),
  am = "Transmission",
  am = c("Automatic" = 0,
  "Manual"=1),
  gear = "Number of forward gears",
  carb = "Number of carburetors"
)
```

For quick cross-tabulation there are fre and cro family of function. For simplicity we demonstrate here only cro\_cpct which calculates column percent. Documentation for other functions, such as cro\_cases for counts, cro\_rpct for row percent, cro\_tpct for table percent and cro\_fun for custom summary functions can be seen by typing ?cro and ?cro\_fun in the console.

```
# 'cro' examples
# just simple crosstabulation, similar to base R 'table' function
cro(mtcars$am, mtcars$vs)

# Table column % with multiple banners
cro_cpct(mtcars$cyl, list(total(), mtcars$am, mtcars$vs))

# or, the same result with another notation
mtcars %>% calc_cro_cpct(cyl, list(total(), am, vs))
```

We have more sophisticated interface for table construction with magrittr piping. Table construction consists of at least of three functions chained with pipe operator: %>%.

At first, we need to specify variables for which statistics will be computed with tab\_cells. Secondary, we calculate statistics with one of the tabs\_stat\_\* functions. And last, we finalize table creation with tab\_pivot, e. g.: dataset %>% tab\_cells(variable) %>% tab\_stat\_cases() %>% tab\_pivot(). After that we can optionally sort table with tab\_sort\_asc, drop empty rows/columns with drop\_rc and transpose with tab\_transpose. Resulting table is just a data. Frame so we can use usual R operations on it. Detailed documentation for table creation can be seen via ?tables. For significance testing see ?significance. Generally, tables automatically translated to HTML for output in knitr or Jupyter notebooks. However, if we want HTML output in the R notebooks or in the RStudio viewer we need to set options for that: expss\_output\_rnotebook() or expss\_output\_viewer().

## Output:

```
Console Terminal Jobs
~ / ~
> library(expss)
> data(mtcars)
> mtcars = apply_labels(mtcars,
+   mpg = "Miles/(US) gallon",
+   cyl = "Number of cylinders",
+   disp = "Displacement (cu.in.)",
+   hp = "Gross horsepower",
+   drat = "Rear axle ratio",
+   wt = "Weight (1000 lbs)",
+   qsec = "1/4 mile time",
+   vs = "Engine",
+   vs = c("V-engine" = 0,
+         "Straight engine" = 1),
+   am = "Transmission",
+   am = c("Automatic" = 0,
+         "Manual"=1),
+   gear = "Number of forward gears",
+   carb = "Number of carburetors"
+ )
> cro(mtcars$am, mtcars$vs)
| | | Engine | |
| | | V-engine | Straight engine |
|-----|-----|
| Transmission | Automatic | 12 | 7 |
| | Manual | 6 | 7 |
| | #Total cases | 18 | 14 |
> cro_cpct(mtcars$cyl, list(total(), mtcars$am, mtcars$vs))
| | | #Total | Transmission | | Engine | |
| | | | Automatic | | Manual | V-engine | Straight engine |
|-----|-----|-----|-----|-----|-----|
| Number of cylinders | | | | | |
| | 4 | 34.4 | 15.8 | 61.5 | 5.6 | 71.4 |
| | 6 | 21.9 | 21.1 | 23.1 | 16.7 | 28.6 |
| | 8 | 43.8 | 63.2 | 15.4 | 77.8 | |
| | #Total cases | 32.0 | 19.0 | 13.0 | 18.0 | 14.0 |
>
```

```
Console Terminal Jobs
~ / ~
> mtcars %>% calc_cro_cpct(cyl, list(total(), am, vs))
| | | #Total | Transmission | | Engine | |
| | | | Automatic | | Manual | V-engine | Straight engine |
|-----|-----|-----|-----|-----|-----|
| Number of cylinders | | | | | |
| | 4 | 34.4 | 15.8 | 61.5 | 5.6 | 71.4 |
| | 6 | 21.9 | 21.1 | 23.1 | 16.7 | 28.6 |
| | 8 | 43.8 | 63.2 | 15.4 | 77.8 | |
| | #Total cases | 32.0 | 19.0 | 13.0 | 18.0 | 14.0 |
> mtcars %>% calc_cro_cpct(cyl, list(total(), am %nest% vs))
| | | #Total | Transmission | | | |
| | | | Automatic | | | |
| | | | Engine | | | |
| | | | V-engine | | | |
|-----|-----|-----|-----|-----|-----|
| Number of cylinders | | | | | |
| | 4 | 34.4 | 42.9 | 16.7 | |
| | 6 | 21.9 | 57.1 | 50.0 | |
| | 8 | 43.8 | 100 | 33.3 | |
| | #Total cases | 32.0 | 12 | 7.0 | 6.0 |
|
|-----|-----|
| Straight engine | | | |
| | 100 | | |
| | 7 | | |
>
> mtcars %>%
+   tab_cells(cyl) %>%
+   tab_cols(total(), am) %>%
+   tab_stat_cpct() %>%
+   tab_pivot()
```

Console Terminal Jobs

```

Number of cylinders | #Total | Transmission | Manual |
|-----|-----|-----|
| 4 | 34.4 | 15.8 | 61.5 |
| 6 | 21.9 | 21.1 | 23.1 |
| 8 | 43.8 | 63.2 | 15.4 |
| #Total cases | 32.0 | 19.0 | 13.0 |

> mtcars %>%
+ tab_cells(mpg, disp, hp, wt, qsec) %>%
+ tab_cols(total(), am) %>%
+ tab_stat_mean_sd_n() %>%
+ tab_last_sig_means(subtable_marks = "both") %>%
+ tab_pivot() %>%
+ set_caption("Table with summary statistics and significance marks.")

Table with summary statistics and significance marks.

#Total | Transmission | Manual |
|-----|-----|-----| | |
|---|---|---|---|---|
| Miles/(US) gallon | Mean | 20.1 | 17.1 < B | 24.4 > A |
| | Std. dev. | 6.0 | 3.8 | 6.2 |
| | Unw. valid N | 32.0 | 19.0 | 13.0 |
| Displacement (cu.in.) | Mean | 230.7 | 290.4 > B | 143.5 < A |
| | Std. dev. | 123.9 | 110.2 | 87.2 |
| | Unw. valid N | 32.0 | 19.0 | 13.0 |
| Gross horsepower | Mean | 146.7 | 160.3 | 126.8 |
| | Std. dev. | 68.6 | 53.9 | 84.1 |
| | Unw. valid N | 32.0 | 19.0 | 13.0 |
| Weight (1000 lbs) | Mean | 3.2 | 3.8 > B | 2.4 < A |
| | Std. dev. | 1.0 | 0.8 | 0.6 |
| | Unw. valid N | 32.0 | 19.0 | 13.0 |
| 1/4 mile time | Mean | 17.8 | 18.2 | 17.4 |
| | Std. dev. | 1.8 | 1.8 | 1.8 |
| | Unw. valid N | 32.0 | 19.0 | 13.0 |

> # Table with the same summary statistics. Statistics labels in columns.
> mtcars %>%
+ tab_cells(mpg, disp, hp, wt, qsec) %>%
+ tab_cols(total(label = "#Total"), am) %>%
+ tab_stat_fun(Mean = w_mean, "Std. dev." = w_sd, "Valid N" = w_n, method = list) %>%
+ tab_pivot()

```

Console Terminal Jobs

```

#Total | Mean | Std. dev. | Valid N | Transmission |
|-----|-----|-----|-----|-----|
|-----|-----|-----|-----|-----|
| Miles/(US) gallon | 20.1 | 6.0 | 32 | Automatic |
| | Mean | 17.1 | | |
| Displacement (cu.in.) | 230.7 | 123.9 | 32 | 290.4 |
| | Mean | 143.5 | | |
| Gross horsepower | 146.7 | 68.6 | 32 | 160.3 |
| | Mean | 126.8 | | |
| Weight (1000 lbs) | 3.2 | 1.0 | 32 | 3.8 |
| | Mean | 2.4 | | |
| 1/4 mile time | 17.8 | 1.8 | 32 | 18.2 |
| | Mean | 17.4 | | |

Std. dev. | Valid N | Manual | Mean | Std. dev. | Valid N |
|-----|-----|-----|-----|-----|-----|
| 3.8 | 19 | 24.4 | 6.2 | 13 | |
| 110.2 | 19 | 143.5 | 87.2 | 13 | |
| 53.9 | 19 | 126.8 | 84.1 | 13 | |
| 0.8 | 19 | 2.4 | 0.6 | 13 | |
| 1.8 | 19 | 17.4 | 1.8 | 13 | |

> # Different statistics for different variables.
> mtcars %>%
+ tab_cols(total(), vs) %>%
+ tab_cells(mpg) %>%
+ tab_stat_mean() %>%
+ tab_stat_valid_n() %>%
+ tab_cells(am) %>%
+ tab_stat_cpct(total_row_position = "none", label = "col %") %>%
+ tab_stat_rpct(total_row_position = "none", label = "row %") %>%
+ tab_stat_tpct(total_row_position = "none", label = "table %") %>%
+ tab_pivot(stat_position = "inside_rows")

#Total | Engine |
|-----|-----| | | | |
|---|---|---|---|---|---|
| Miles/(US) gallon | Mean | 20.1 | 16.6 |
| | Valid N | 32.0 | 18.0 |
| | Automatic | col % | 59.4 | 66.7 |
| | | row % | 100.0 | 63.2 |
| | | table % | 59.4 | 37.5 |
| | | | col % | 40.6 | 33.3 |
| | | | row % | 100.0 | 46.2 |
| | | | table % | 40.6 | 18.8 |

```

```
Console Terminal Jobs ~/  
Straight engine  
----  
24.6  
14.0  
50.0  
36.8  
21.9  
50.0  
53.8  
21.9  
> # Table with split by rows and with custom totals.  
> mtcars %>%  
+ tab_cells(cyl) %>%  
+ tab_cols(total(), vs) %>%  
+ tab_rows(am) %>%  
+ tab_stat_cpct(total_row_position = "above",  
+                 total_label = c("number of cases", "row %"),  
+                 total_statistic = c("u_cases", "u_rpct")) %>%  
+ tab_pivot()  
  
#Total  
-----  
Transmission | Automatic | Number of cylinders | #number of cases | #row % |  
-----  
              |           | 4                   | 19.0          | 100.0   |  
              |           | 6                   | 15.8          | 41.1    |  
              |           | 8                   | 21.1          | 23.1    |  
              |           | 13.0               | 63.2          | 15.4    |  
              |           | 100.0              | 13.0          | 100.0   |  
              |           | 4                   | 61.5          | 41.1    |  
              |           | 6                   | 23.1          | 23.1    |  
              |           | 8                   | 15.4          | 15.4    |  
-----  
Manual | Number of cylinders | #number of cases | #row % |  
-----  
              | 4                   | 19.0          | 100.0   |  
              | 6                   | 15.8          | 41.1    |  
              | 8                   | 21.1          | 23.1    |  
              | 13.0               | 63.2          | 15.4    |  
              | 100.0              | 13.0          | 100.0   |  
              | 4                   | 61.5          | 41.1    |  
              | 6                   | 23.1          | 23.1    |  
              | 8                   | 15.4          | 15.4    |  
-----  
Engine | Straight engine  
-----  
V-engine |  
-----  
12.0      | 7.0  
63.2      | 36.8  
          | 42.9  
          | 57.1  
100.0     | 7.0  
6.0       |  
46.2     |
```

```
Console Terminal Jobs
~/
> data(product_test)
>
> w = product_test # shorter name to save some keystrokes
>
> # here we recode variables from first/second tested product to separate variables for each product according to their cells
> # 'h' variables - VSX123 sample, 'p' variables - 'SDF456' sample
> # also we recode preferences from first/second product to true names
> # for first cell there are no changes, for second cell we should exchange 1 and 2.
> w = w %>%
+   do_if(cell == 1, {
+     recode(a1_1 %to% a1_6, other ~ copy) %into% (h1_1 %to% h1_6)
+     recode(b1_1 %to% b1_6, other ~ copy) %into% (p1_1 %to% p1_6)
+     recode(a2_1, other ~ copy) %into% h2_2
+     recode(b2_2, other ~ copy) %into% p2_2
+     clr = cl
+   }) %%
+   do_if(cell == 2, {
+     recode(a1_1 %to% a1_6, other ~ copy) %into% (p1_1 %to% p1_6)
+     recode(b1_1 %to% b1_6, other ~ copy) %into% (h1_1 %to% h1_6)
+     recode(a2_2, other ~ copy) %into% p2_2
+     recode(b2_2, other ~ copy) %into% h2_2
+     recode(cl, 1 ~ 2, 2 ~ 1, other ~ copy) %into% clr
+   }) %%
+   compute({
+     # recode age by groups
+     age_cat = recode(s2a, 10 %thru% 25 ~ 1, 10 %thru% hi ~ 2)
+     # count number of likes
+     # codes 2 and 99 are ignored.
+     h_likes = count_row_if(1 | 3 %thru% 98, h1_1 %to% h1_6)
+     p_likes = count_row_if(1 | 3 %thru% 98, p1_1 %to% p1_6)
+   })
>
> # here we prepare labels for future usage
> codeframe_likes = num_lab(""
+   1 Liked everything
+   2 Disliked everything
+   3 Chocolate
+   4 Appearance
+   5 Taste
+   6 Stuffing
+   7 Nuts
+   8 Consistency
+   98 Other
+   99 Hard to answer
```

**Conclusion:** Thus, we have studied tables with labels in R.

## Experiment no:6

**Aim:** Working with large datasets with dplyr and data-table.

### Theory:

dplyr is the next iteration of plyr, focussed on tools for working with data frames (hence the d in the name). It has three main goals:

- Identify the most important data manipulation tools needed for data analysis and make them easy to use from R.
- Provide blazing fast performance for in-memory data by writing key pieces in C++.
- Use the same interface to work with data no matter where it's stored, whether in a data frame, a data table or database.

### One can install:

- The latest released version from CRAN with

```
install.packages("dplyr")
```

- The latest development version from github with

```
if (packageVersion("devtools") < 1.6)
{
  install.packages("devtools")
}
devtools::install_github("hadley/lazyeval")
devtools::install_github("hadley/dplyr")
```

### tbls:

The key object in dplyr is a `tbl`, a representation of a tabular data structure. Currently dplyr supports:

- data frames
- [data tables]
- SQLite
- PostgreSQL
- MySQL/MariaDB

- Bigquery
- MonetDB
- data cubes with arrays (partial implementation)

```
library(dplyr)
# for functions library(hflights)
# for data head(hflights)

# Coerce to data table
hflights_dt <- tbl_dt(hflights)

# Caches data in local SQLite db
hflights_db1 <- tbl(hflights_sqlite(), "hflights")

# Caches data in local postgres db
hflights_db2 <- tbl(hflights_postgres(), "hflights")
Each tbl also comes in a grouped variant which allows one to easily perform
operations "by group":
```

```
carriers_df <- group_by(hflights, UniqueCarrier)
carriers_dt <- group_by(hflights_dt, UniqueCarrier)
carriers_db1 <- group_by(hflights_db1, UniqueCarrier)
carriers_db2 <- group_by(hflights_db2, UniqueCarrier)
```

dplyr can be effectively used to

- Pick observations by their values (filter()).
- Reorder the rows (arrange()).
- Pick variables by their names (select()).
- Create new variables with functions of existing variables (mutate())
- Collapse many values down to a single summary (summarise()).

These can all be used in conjunction with group\_by() which changes the scope of each function from operating on the entire dataset to operating on its group-by-group. These six functions provide the verbs for a language of data manipulation.

### **All verbs work similarly:**

1. The first argument is a data frame.

2. The subsequent arguments describe what to do with the data frame, using the variable names (without quotes).

3. The result is a new data frame.

## Output:

```
> library(mtcars)
> library(tidyverse)
-- Attaching packages -- tidyverse 1.3.0 --
#> # ggplot2 3.2.1   # purrr  0.3.3
#> # tibble  2.1.3   # dplyr   0.8.3
#> # tidyr   1.0.2   # stringr 1.4.0
#> # readr   1.3.1   # forcats 0.4.0
-- Conflicts -- tidyverse_conflicts() --
#> # dplyr::between()    masks expss::between()
#> # dplyr::compute()    masks expss::compute()
#> # dplyr::contains()    masks tidyrr::contains(), expss::contains()
#> # dplyr::filter()     masks tidyrr::filter()
#> # dplyr::first()      masks expss::first()
#> # stringr::fixed()    masks expss::fixed()
#> # purrr::keep()       masks expss::keep()
#> # dplyr::lag()        masks stats::lag()
#> # dplyr::last()       masks expss::last()
#> # stringr::modify()   masks tidyrr::modify()
#> # purrr::modify_if()  masks expss::modify_if()
#> # dplyr::nest()       masks expss::nest()
#> # dplyr::recode()    masks expss::recode()
#> # stringr::regex()    masks expss::regex()
#> # purrr::transpose()  masks expss::transpose()
#> # dplyr::vars()       masks ggplot2::vars(), expss::vars()
#> # purrr::when()       masks expss::when()
> flights
#> # A tibble: 336,776 x 19
#> #   year month day dep_time sched_dep_time dep_delay arr_time sched_arr_time
#> #   <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
#> 1 2013 1 1 517 2 830 819
#> 2 2013 1 1 533 4 850 830
#> 3 2013 1 1 542 549 2 923 850
#> 4 2013 1 1 544 545 -.1 1004 1022
#> 5 2013 1 1 554 600 -.6 812 837
#> 6 2013 1 1 554 558 -.4 740 728
#> 7 2013 1 1 555 600 -.5 913 854
#> 8 2013 1 1 557 600 -.3 769 723
#> 9 2013 1 1 557 600 -.3 838 846
#> 10 2013 1 1 558 600 -.3 753 745
#> # ... with 336,766 more rows, and 11 more variables: arr_delay <dbl>,
#> #   carrier <chr>, flight <int>, tailnum <chr>, origin <chr>, dest <chr>,
#> #   air_time <dbl>, distance <dbl>, hour <dbl>, minute <dbl>, time_hour <dttm>
#> |
#> # filter(flights, month == 1, day == 1)
#> # A tibble: 842 x 19
#> #   year month day dep_time sched_dep_time dep_delay arr_time sched_arr_time
#> #   <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
#> 1 2013 1 1 517 515 2 830 819
#> 2 2013 1 1 533 529 4 850 830
#> 3 2013 1 1 542 549 2 923 850
#> 4 2013 1 1 544 545 -.1 1004 1022
#> 5 2013 1 1 554 600 -.6 812 837
#> 6 2013 1 1 554 558 -.4 740 728
#> 7 2013 1 1 555 600 -.5 913 854
#> 8 2013 1 1 557 600 -.3 769 723
#> 9 2013 1 1 557 600 -.3 838 846
#> 10 2013 1 1 558 600 -.3 753 745
#> # ... with 832 more rows, and 11 more variables: arr_delay <dbl>,
#> #   carrier <chr>, flight <int>, tailnum <chr>, origin <chr>, dest <chr>,
#> #   air_time <dbl>, distance <dbl>, hour <dbl>, minute <dbl>, time_hour <dttm>
#> # filter(flights, month == 1, day == 1)
#> # A tibble: 842 x 19
#> #   year month day dep_time sched_dep_time dep_delay arr_time sched_arr_time
#> #   <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl> <dbl>
#> 1 2013 11 1 5 2359 6 352 345
#> 2 2013 11 1 35 2256 105 123 236
#> 3 2013 11 1 559 545 506 -.5 641 651
#> 4 2013 11 1 559 545 506 -.5 536 537
#> 5 2013 11 1 542 545 506 -.5 831 855
#> 6 2013 11 1 549 600 -.11 912 923
#> 7 2013 11 1 558 600 -.10 765 659
#> 8 2013 11 1 554 600 -.6 659 701
#> 9 2013 11 1 554 600 600 -.6 826 827
#> 10 2013 11 1 554 600 600 -.6 749 751
#> # ... with 55,393 more rows, and 11 more variables: arr_delay <dbl>,
#> #   carrier <chr>, flight <int>, tailnum <chr>, origin <chr>, dest <chr>,
#> #   air_time <dbl>, distance <dbl>, hour <dbl>, minute <dbl>, time_hour <dttm>
#> |
#> |
```

70% 01:06 Sunday, 09.21

```
> select(flights, year, month, day)
#> # A tibble: 336,776 x 3
#> #   year month day
#> #   <dbl> <dbl> <dbl>
#> 1 2013 1 1
#> 2 2013 1 1
#> 3 2013 1 1
#> 4 2013 1 1
#> 5 2013 1 1
#> 6 2013 1 1
#> 7 2013 1 1
#> 8 2013 1 1
#> 9 2013 1 1
#> 10 2013 1 1
#> # ... with 336,766 more rows
#> # select(flights, year:day)
#> # A tibble: 336,776 x 3
#> #   year month day
#> #   <dbl> <dbl> <dbl>
#> 1 2013 5 1
#> 2 2013 1 1
#> 3 2013 1 1
#> 4 2013 1 1
#> 5 2013 1 1
#> 6 2013 1 1
#> 7 2013 1 1
#> 8 2013 1 1
#> 9 2013 1 1
#> 10 2013 1 1
#> # ... with 336,766 more rows
#> # filter(flights, year:day)
#> # A tibble: 336,776 x 16
#> #   year month day dep_time sched_dep_time dep_delay arr_time sched_arr_time
#> #   <dbl> <dbl>
#> 1 2013 1 1 517 515 2 830 819 11 UA
#> 2 2013 1 1 533 529 4 850 830
#> 3 2013 1 1 542 549 2 923 850
#> 4 2013 1 1 544 545 -.1 1004 1022
#> 5 2013 1 1 554 600 -.6 812 837
#> 6 2013 1 1 554 558 -.4 740 728
#> 7 2013 1 1 555 600 -.5 913 854
#> 8 2013 1 1 557 600 -.3 769 723
#> 9 2013 1 1 557 600 3 838 846
#> 10 2013 1 1 558 600 -.2 753 745 8 AA
#> # ... with 336,766 more rows, and 9 more variables: flight <chr>, tailnum <chr>,
#> #   origin <chr>, dest <chr>, air_time <dbl>, distance <dbl>, hour <dbl>,
#> #   minute <dbl>, time_hour <dttm>
#> |
```

74% 01:06 Sunday, 09.29

```

> library(nycflights13)
> library(tidyverse)
-- Attaching packages tidyverse 1.3.0 --
  # ggplot2 3.2.1   purrr  0.3.3
  # tibble   2.1.3   dplyr  0.8.3
  # tidyr   1.0.2   stringr 1.4.0
  # readr   1.3.1   forcats 0.4.0
-- Conflicts tidyverse_conflicts() --
  # dplyr::between()    masks expss::between()
  # dplyr::compute()    masks expss::compute()
  # dplyr::contains()    masks tidyv::contains(), expss::contains()
  # dplyr::filter()     masks stats::filter()
  # dplyr::first()      masks expss::first()
  # stringr::fixed()   masks expss::fixed()
  # purrr::keep()       masks expss::keep()
  # dplyr::lag()        masks stats::lag()
  # dplyr::last()       masks expss::last()
  # purrr::modify()    masks expss::modify()
  # purrr::modify_if() masks expss::modify_if()
  # dplyr::na_if()      masks expss::na_if()
  # tidyv::nest()       masks expss::nest()
  # dplyr::recode()    masks expss::recode()
  # stringr::regex()   masks expss::regex()
  # purrr::transpose() masks expss::transpose()
  # dplyr::vars()       masks ggplot2::vars(), expss::vars()
  # purrr::when()       masks expss::when()

# A tibble: 336,776 x 19
  year month   day dep_time sched_dep_time dep_delay arr_time sched_arr_time
  <int> <int> <int> <int> <int> <int> <dbl> <int> <int>
1 2013     1     1   517          830        -2       830         819
2 2013     1     1   533          850        -4       850         830
3 2013     1     1   542          850        -2       923         850
4 2013     1     1   544          845        -1      1094        1024
5 2013     1     1   554          600        -6       612         637
6 2013     1     1   554          558        -8       600         728
7 2013     1     1   555          600        -5       913         854
8 2013     1     1   557          600        -3       799         723
9 2013     1     1   557          600        -3       600         846
10 2013    1     1   558          600        -2       753         745
# ... with 336,766 more rows, and 11 more variables: arr_delay <dbl>,
# carrier <chr>, flight <int>, tailnum <chr>, origin <chr>, dest <chr>,
# air_time <dbl>, distance <dbl>, hour <dbl>, minute <dbl>, time_hour <dtm>,
# >

# A tibble: 336,776 x 19
  year month   day dep_time sched_dep_time dep_delay arr_time sched_arr_time
  <int> <int> <int> <int> <int> <int> <dbl> <int> <int>
1 2013     1     1   517          515        -2       830         819
2 2013     1     1   533          529        -4       850         830
3 2013     1     1   542          542        -1       923         850
4 2013     1     1   544          544        -1      1094        1024
5 2013     1     1   554          600        -8       612         637
6 2013     1     1   555          558        -8       600         728
7 2013     1     1   557          600        -5       913         854
8 2013     1     1   558          600        -3       799         723
9 2013     1     1   559          600        -3       600         846
10 2013    1     1   560          600        -2       753         745
# ... with 336,766 more rows, and 11 more variables: arr_delay <dbl>,
# carrier <chr>, flight <int>, tailnum <chr>, origin <chr>, dest <chr>,
# air_time <dbl>, distance <dbl>, hour <dbl>, minute <dbl>, time_hour <dtm>,
# >

# A tibble: 336,776 x 19
  year month   day dep_time sched_dep_time dep_delay arr_time sched_arr_time
  <int> <int> <int> <int> <int> <int> <dbl> <int> <int>
1 2013     1     1   941          900        -4       1381        1242
2 2013     6    15   1432         1935        -1       1137        1667
3 2013     1    10   1121         1635        -6       1126        1239
4 2013     9    20   1139         1845        -1       1014        1457
5 2013     7    22   845          600        -2       1085        1044
6 2013     4    18   1188         1086        -2       968        1342
7 2013     3    17   2411         2110        -3       2111        1715
8 2013     3    27   2590         1909        -6       2209        2266
9 2013     7    22   2557         759        -4       988        121
10 2013    12    5    756         1700        -4       896        1958
# ... with 336,766 more rows, and 11 more variables: arr_delay <dbl>,
# carrier <chr>, flight <int>, tailnum <chr>, origin <chr>, dest <chr>,
# air_time <dbl>, distance <dbl>, hour <dbl>, minute <dbl>, time_hour <dtm>,
# >

# A tibble: 336,776 x 19
  year month   day dep_time sched_dep_time dep_delay arr_time sched_arr_time
  <int> <int> <int> <int> <int> <int> <dbl> <int> <int>
1 2013     1     1   941          900        -4       1381        1242
13 2013     1     1   1533         1529        -4       1137        1667
3 2013     1    10   1121         1635        -6       1126        1239
4 2013     9    20   1139         1845        -1       1014        1457
5 2013     7    22   845          600        -2       1085        1044
6 2013     4    18   1188         1086        -2       968        1342
7 2013     3    17   2411         2110        -3       2111        1715
8 2013     3    27   2590         1909        -6       2209        2266
9 2013     7    22   2557         759        -4       988        121
10 2013    12    5    756         1700        -4       896        1958
# ... with 336,766 more rows, and 11 more variables: arr_delay <dbl>,
# carrier <chr>, flight <int>, tailnum <chr>, origin <chr>, dest <chr>,
# air_time <dbl>, distance <dbl>, hour <dbl>, minute <dbl>, time_hour <dtm>,
# >

# A tibble: 336,776 x 19
  year month   day dep_time sched_dep_time dep_delay arr_time sched_arr_time
  <int> <int> <int> <int> <int> <int> <dbl> <int> <int>
1 2013     1     1   941          900        -4       1381        1242
2 2013     1     1   1533         1529        -4       1137        1667
3 2013     1    10   1121         1635        -6       1126        1239
4 2013     1    11   1139         1845        -1       1014        1457
5 2013     1    12   845          600        -2       1085        1044
6 2013     1    13   1188         1086        -2       968        1342
7 2013     1    14   2411         2110        -3       2111        1715
8 2013     1    15   2590         1909        -6       2209        2266
9 2013     1    16   2557         759        -4       988        121
10 2013    12    5    756         1700        -4       896        1958
# ... with 336,766 more rows, and 11 more variables: arr_delay <dbl>,
# carrier <chr>, flight <int>, tailnum <chr>, origin <chr>, dest <chr>,
# air_time <dbl>, distance <dbl>, hour <dbl>, minute <dbl>, time_hour <dtm>,
# >
```

69% | Sunday, 09:25

```

> library(mvflights13)
> library(tidyverse)
- Attaching packages
  ✓ ggplot2 3.2.1   ✓ purrr  0.3.3
  ✓ tibble  2.1.3   ✓ dplyr  0.8.3
  ✓ tidyver 1.0.2   ✓ stringr 1.4.0
  ✓ readr  1.3.1    ✓ forcats 0.4.0
- Conflicts
  #::dplyr::between() masks expss::between()
  #::dplyr::compute() masks expss::compute()
  #::dplyr::contains() masks tidyverse::contains(), expss::contains()
  #::dplyr::filter()  masks stats::filter()
  #::dplyr::first()   masks expss::first()
  #::stringr::fixed() masks expss::fixed()
  #::purrr::keep()   masks expss::keep()
  #::purrr::last()   masks expss::last()
  #::dplyr::last()   masks expss::last()
  #::purrr::modify() masks expss::modify()
  #::purrr::modify_if() masks expss::modify_if()
  #::dplyr::na_if()   masks expss::na_if()
  #::tidyver::nest()  masks expss::nest()
  #::dplyr::recode()  masks expss::recode()
  #::stringr::regex() masks expss::regex()
  #::purrr::transpose() masks expss::transpose()
  #::dplyr::vars()   masks ggplot2::vars(), expss::vars()
  #::purrr::when()   masks expss::when()
> flights
# A tibble: 336,776 x 19
  year month day dep_time sched_dep_time arr_time delay arr_delay distance air_time gain
  <int> <int> <int> <int> <int> <int> <dbl> <int> <int> <dbl> <dbl> <dbl> <dbl>
1 2013 1 1 517 515 2 830 819
2 2013 1 1 533 529 4 850 830
3 2013 1 1 542 540 2 923 850
4 2013 1 1 544 545 1 - 1004 1922
5 2013 1 1 554 600 .6 812 837
6 2013 1 1 554 558 -.4 740 728
7 2013 1 1 555 600 -.5 913 854
8 2013 1 1 557 600 -.3 793 723
9 2013 1 1 557 600 -.3 838 846
10 2013 1 1 558 600 -.2 753 745
# ... with 336,766 more rows, and 11 more variables: arr_delay <dbl>, carrier <chr>, flight <int>, tailnum <chr>, origin <chr>, dest <chr>, air_time <dbl>, distance <dbl>, hour <dbl>, minute <dbl>, time_hour <dttm>
> |
+ air_time
+ mutate(flights_gnl,
+       gain = dep_delay - arr_delay,
+       speed = distance / air_time * 60
+     )
# A tibble: 336,776 x 9
  year month day dep_time sched_dep_time arr_time delay arr_delay air_time gain
  <int> <int> <int> <int> <int> <int> <dbl> <int> <dbl>
1 2013 1 1 517 515 2 830 819 227 3.78
2 2013 1 1 533 529 4 850 830 227 3.78
3 2013 1 1 542 540 2 923 850 160 -31 408
4 2013 1 1 544 545 1 - 1004 1922 309 17 517
5 2013 1 1 554 600 .6 812 837 257 19 304
6 2013 1 1 555 600 -.5 913 854 12 719 150 -16 288
7 2013 1 1 557 600 -.3 793 723 19 1065 158 -24 304
8 2013 1 1 557 600 -.3 838 846 14 229 53 11 259
9 2013 1 1 558 600 -.2 753 745 8 944 140 5 405
10 2013 1 1 558 600 -.2 753 745 8 733 138 -10 319
# ... with 336,766 more rows, and 1 more variable: gain_per_hour <dbl>
> |

```

79% | Sunday, 09:33

81% | Sunday, 09:34

**Conclusion:** Thus, we have worked with large datasets.

## Experiment no:7

**Aim:** EDA with R- Range, Mean, Variance, Median, Standard Deviation, Histogram, Box-Plot, Scatter-Plot.

### Theory:

#### Exploratory Data Analysis:

**Exploratory Data Analysis (EDA)** is the process of analysing and visualizing the data to get a better understanding of the data and glean insight from it. There are various steps involved when doing EDA but the following are the common steps that a data analyst can take when performing EDA:

1. Import the data
2. Clean the data
3. Process the data
4. Visualize the data

#### A. Range in R:

The range of an observation variable is the difference of its largest and smallest data values. It is a measure of how far apart the entire data spreads in value.

$$\text{Range} = \text{Largest Value} - \text{Smallest Value}$$

### Output:

The screenshot shows the RStudio interface with the R console tab active. The console output displays the R version information and the calculation of the range for the 'duration' variable from the 'faithful' dataset.

```
R version 3.6.2 (2019-12-12) -- "Dark and Stormy Night"
Copyright (C) 2019 The R Foundation for Statistical Computing
Platform: x86_64-w64-mingw32/x64 (64-bit)

R is free software and comes with ABSOLUTELY NO WARRANTY.
You are welcome to redistribute it under certain conditions.
Type 'license()' or 'licence()' for distribution details.

Natural language support but running in an English locale

R is a collaborative project with many contributors.
Type 'contributors()' for more information and
'citation()' on how to cite R or R packages in publications.

Type 'demo()' for some demos, 'help()' for on-line help, or
'help.start()' for an HTML browser interface to help.
Type 'q()' to quit R.

> duration=faithful$eruptions
> max(duration)-min(duration)
[1] 3.5
>
```

## **B. Mean in R:**

It is calculated by taking the sum of the values and dividing with the number of values in a data series. The function **mean ()** is used to calculate this in R.

### **Syntax:**

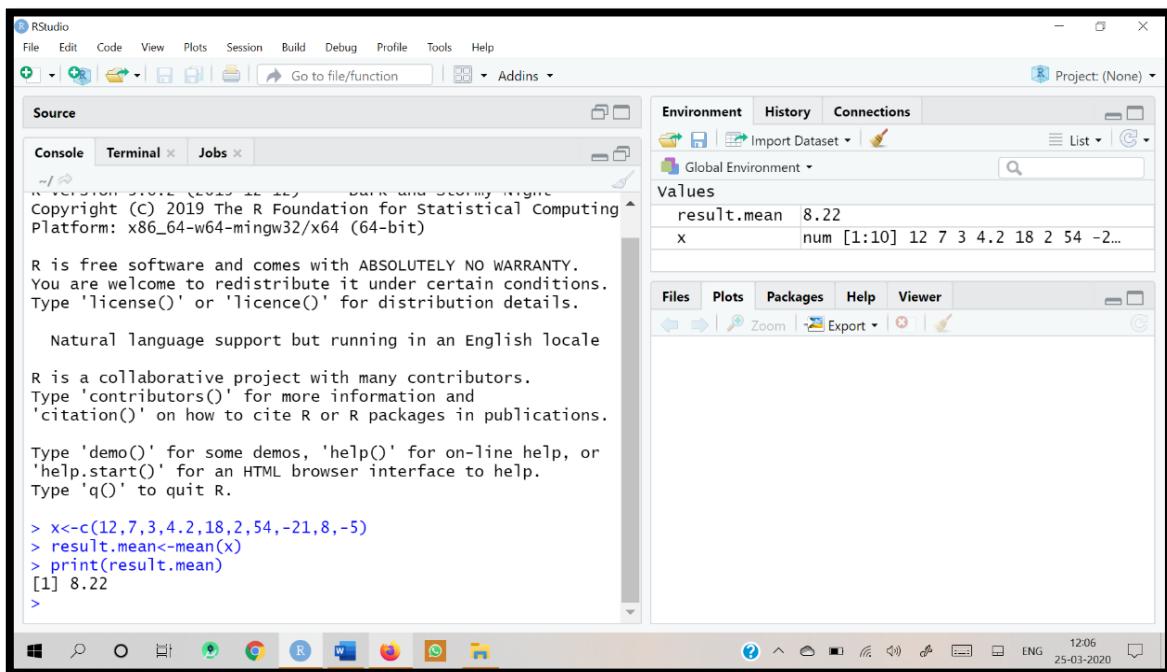
The basic syntax for calculating mean in R is –

```
mean (x, trim = 0, na.rm = FALSE, ...)
```

Following is the description of the parameters used –

- **x** is the input vector.
- **trim** is used to drop some observations from both end of the sorted vector.
- **na.rm** is used to remove the missing values from the input vector.

### **Output:**



A screenshot of the RStudio interface. The top menu bar includes File, Edit, Code, View, Plots, Session, Build, Debug, Profile, Tools, and Help. The main window has tabs for Source, Console, Terminal, and Jobs. The Console tab shows R startup messages, including the R version (4.0.2), copyright notice, and platform (x86\_64-w64-mingw32/x64). It also shows the command `> x<-c(12,7,3,4.2,18,2,54,-21,8,-5)`, the assignment of `result.mean<-mean(x)`, and the output `[1] 8.22`. The Environment tab in the right sidebar shows the variable `result.mean` with the value `8.22` and the variable `x` with the value `num [1:10] 12 7 3 4.2 18 2 54 -2...`.

## **C. Variance in R:**

The variance is a numerical measure of how the data values is dispersed around the mean. In particular, the sample variance is defined as:

$$s^2 = \frac{1}{n-1} \sum_{i=1}^{n-1} (x_i - \bar{x})^2$$

## Output:

```
R version 3.6.2 (2019-12-12) -- "Dark and Stormy Night"
Copyright (C) 2019 The R Foundation for Statistical Computing
Platform: x86_64-w64-mingw32/x64 (64-bit)

R is free software and comes with ABSOLUTELY NO WARRANTY.
You are welcome to redistribute it under certain conditions.
Type 'license()' or 'licence()' for distribution details.

Natural language support but running in an English locale

R is a collaborative project with many contributors.
Type 'contributors()' for more information and
'citation()' on how to cite R or R packages in publications.

Type 'demo()' for some demos, 'help()' for on-line help, or
'help.start()' for an HTML browser interface to help.
Type 'q()' to quit R.

> duration=faithful$eruptions
> var(duration)
[1] 1.302728
>
```

## D. Median in R:

The median of an observation variable is the value at the middle when the data is sorted in ascending order. It is an ordinal measure of the central location of the data values.

## Output:

```
R version 3.6.2 (2019-12-12) -- "Dark and Stormy Night"
Copyright (C) 2019 The R Foundation for Statistical Computing
Platform: x86_64-w64-mingw32/x64 (64-bit)

R is free software and comes with ABSOLUTELY NO WARRANTY.
You are welcome to redistribute it under certain conditions.
Type 'license()' or 'licence()' for distribution details.

Natural language support but running in an English locale

R is a collaborative project with many contributors.
Type 'contributors()' for more information and
'citation()' on how to cite R or R packages in publications.

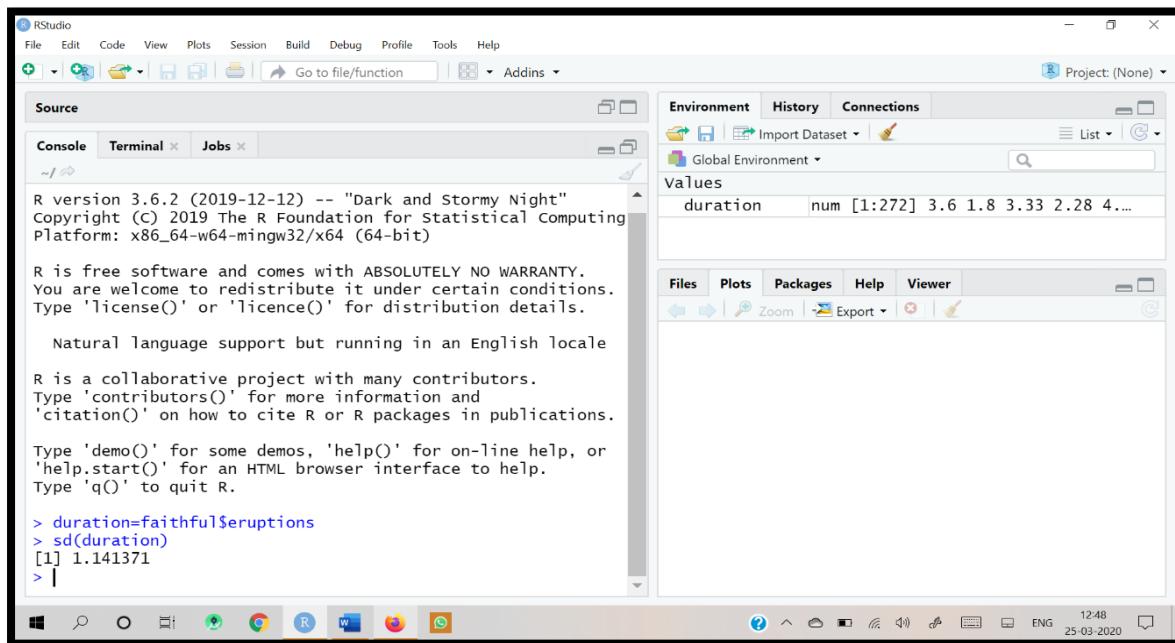
Type 'demo()' for some demos, 'help()' for on-line help, or
'help.start()' for an HTML browser interface to help.
Type 'q()' to quit R.

> duration=faithful$eruptions
> median(duration)
[1] 4
>
```

## **E. Standard Deviation in R:**

The standard deviation of an observation variable is the square root of its variance.

### **Output:**



A screenshot of the RStudio interface. The top menu bar includes File, Edit, Code, View, Plots, Session, Build, Debug, Profile, Tools, and Help. The bottom taskbar shows various application icons. The main area has tabs for Console, Terminal, and Jobs. The Console tab displays the R startup message and a command-line session:  
R version 3.6.2 (2019-12-12) -- "Dark and Stormy Night"  
Copyright (C) 2019 The R Foundation for Statistical Computing  
Platform: x86\_64-w64-mingw32/x64 (64-bit)  
  
R is free software and comes with ABSOLUTELY NO WARRANTY.  
You are welcome to redistribute it under certain conditions.  
Type 'license()' or 'licence()' for distribution details.  
  
Natural language support but running in an English locale  
  
R is a collaborative project with many contributors.  
Type 'contributors()' for more information and  
'citation()' on how to cite R or R packages in publications.  
  
Type 'demo()' for some demos, 'help()' for on-line help, or  
'help.start()' for an HTML browser interface to help.  
Type 'q()' to quit R.  
  
> duration=faithful\$eruptions  
> sd(duration)  
[1] 1.141371  
>

The Global Environment pane on the right shows a variable named 'duration' with a value of num [1:272] 3.6 1.8 3.33 2.28 4....

## **F. Histogram in R:**

A histogram represents the frequencies of values of a variable bucketed into ranges. Histogram is similar to bar chart but the difference is it groups the values into continuous ranges. Each bar in histogram represents the height of the number of values present in that range. R creates histogram using **hist()** function. This function takes a vector as an input and uses some more parameters to plot histograms.

### **Syntax:**

The basic syntax for creating a histogram using R is –

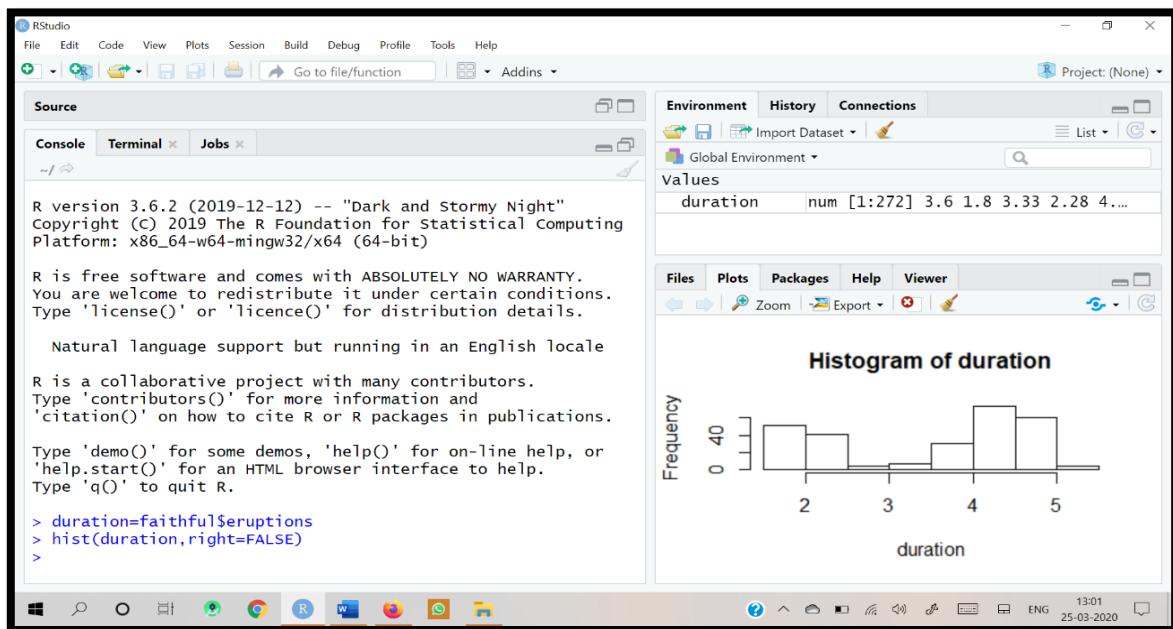
```
hist (v, main, xlab, xlim, ylim, breaks, col, border)
```

Following is the description of the parameters used –

- **v** is a vector containing numeric values used in histogram.
- **main** indicates title of the chart.

- **col** is used to set color of the bars.
- **border** is used to set border color of each bar.
- **xlab** is used to give description of x-axis.
- **xlim** is used to specify the range of values on the x-axis.
- **ylim** is used to specify the range of values on the y-axis.
- **breaks** are used to mention the width of each bar.

## Output:



## G. Box-plot in R:

Boxplots are a measure of how well distributed is the data in a data set. It divides the data set into three quartiles. This graph represents the minimum, maximum, median, first quartile and third quartile in the data set. It is also useful in comparing the distribution of data across data sets by drawing boxplots for each of them. Boxplots are created in R by using the **boxplot()** function.

## Syntax:

The basic syntax to create a boxplot in R is –

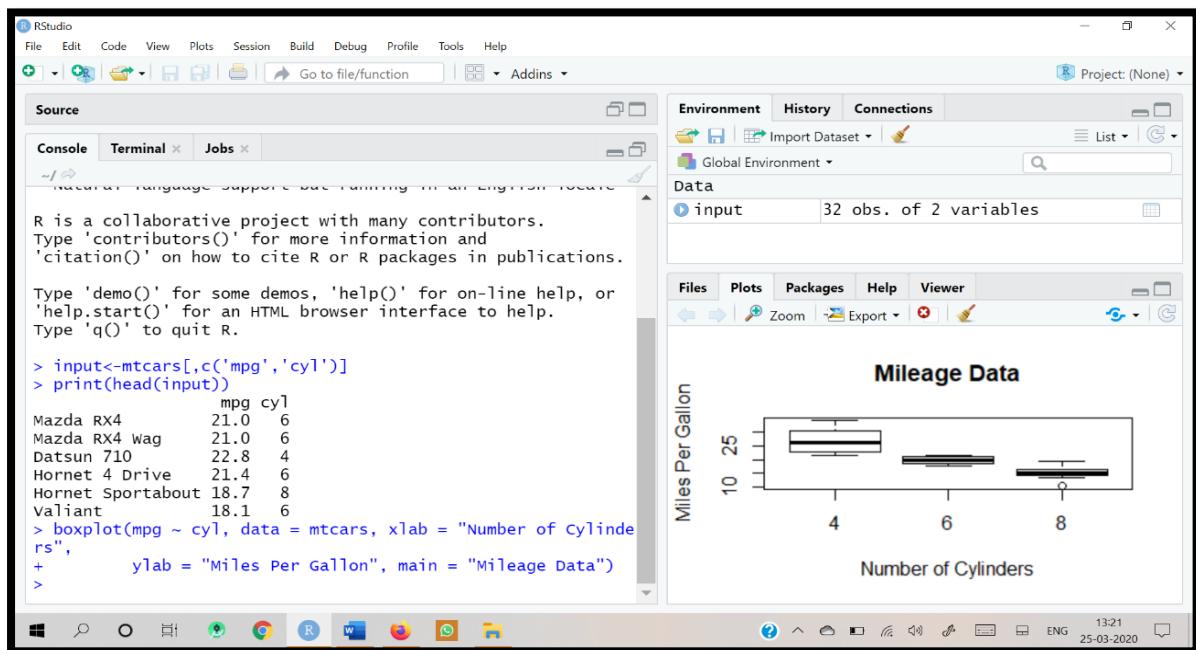
```
boxplot(x, data, notch, varwidth, names, main)
```

Following is the description of the parameters used –

- **x** is a vector or a formula.

- **data** is the data frame.
- **notch** is a logical value. Set as TRUE to draw a notch.
- **varwidth** is a logical value. Set as true to draw width of the box proportionate to the sample size.
- **names** are the group labels which will be printed under each boxplot.
- **main** is used to give a title to the graph.

## Output:



## H. Scatter-plot in R:

Scatterplots show many points plotted in the Cartesian plane. Each point represents the values of two variables. One variable is chosen in the horizontal axis and another in the vertical axis. The simple scatterplot is created using the **plot ()** function.

### Syntax:

The basic syntax for creating scatterplot in R is –

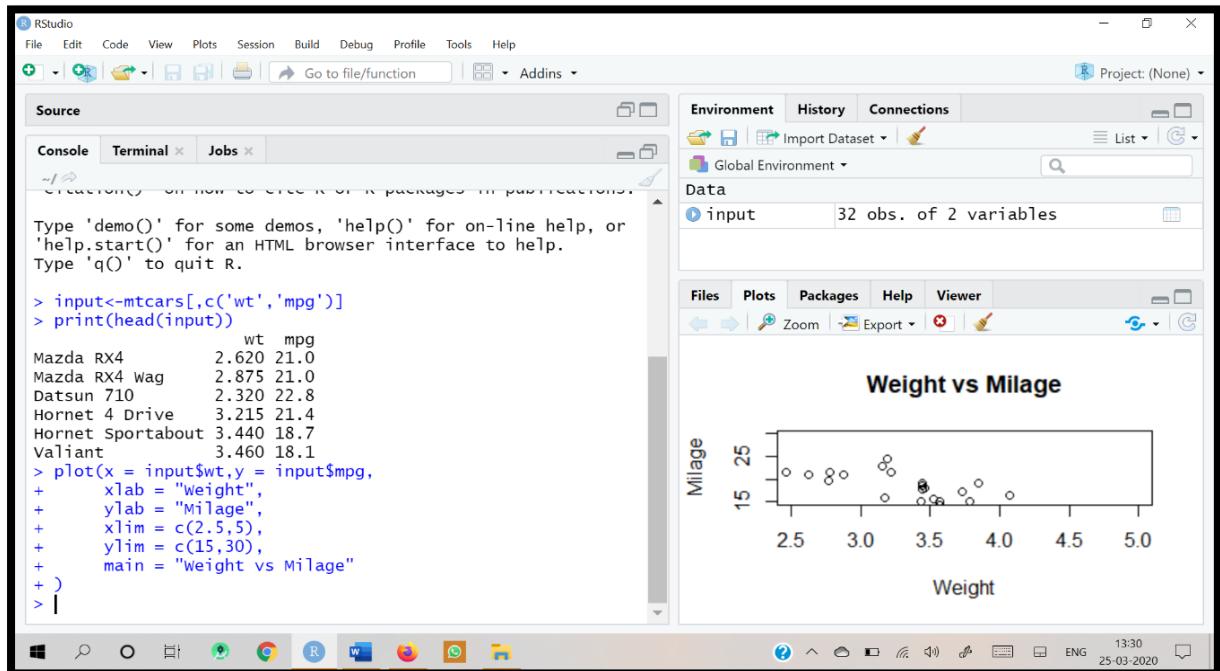
```
plot (x, y, main, xlab, ylab, xlim, ylim, axes)
```

Following is the description of the parameters used –

- **x** is the data set whose values are the horizontal coordinates.
- **y** is the data set whose values are the vertical coordinates.

- **main** is the title of the graph.
- **xlab** is the label in the horizontal axis.
- **ylab** is the label in the vertical axis.
- **xlim** is the limits of the values of x used for plotting.
- **ylim** is the limits of the values of y used for plotting.
- **axes** indicate whether both axes should be drawn on the plot.

## Output:



## Conclusion:

Thus, we have studied exploratory data analysis in R.

## Experiment no:8

**Aim:** Basics of Advanced graphics in R.

### **Theory:**

Graphics are important for conveying important features of the data. R includes at least three graphical systems, the standard graphics package, the lattice package for Trellis graphs. R has good graphical capabilities but there are some alternatives like gun plot.

R provides the usual range of standard statistical plots, including scatterplots, boxplots, histograms, bar plots, pie charts, and basic 3D plots.

They can be used to examine

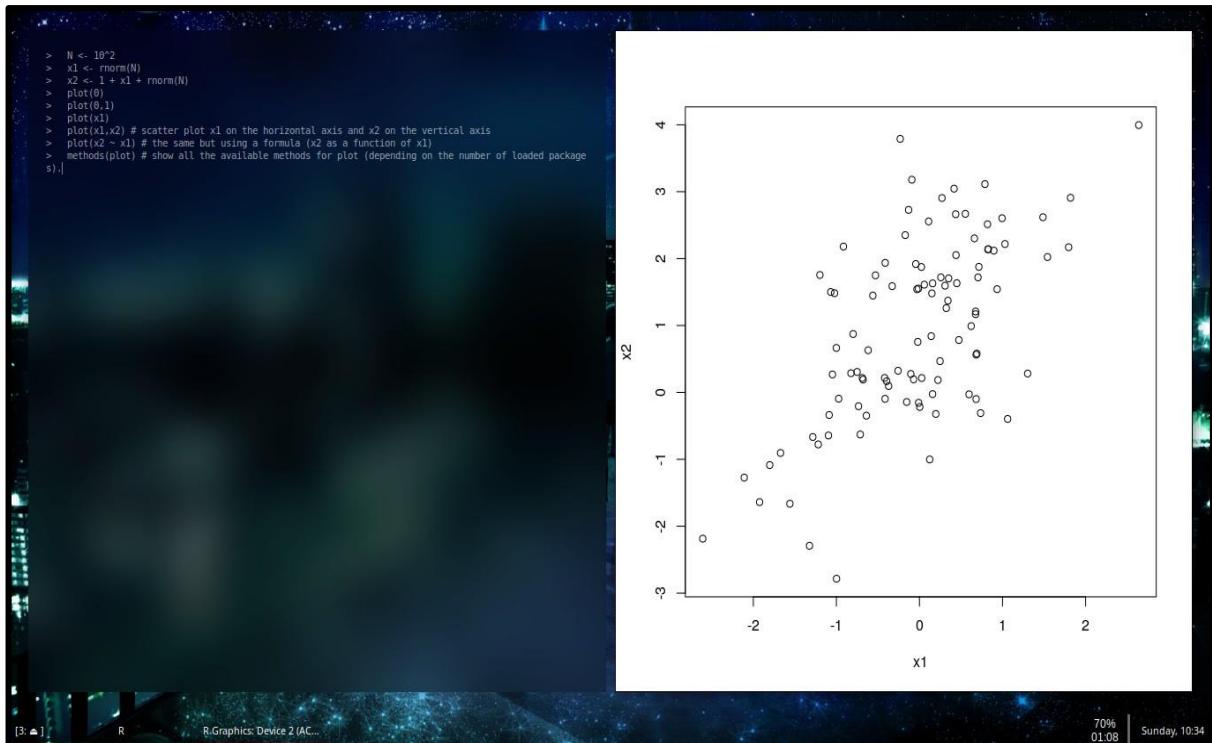
- Marginal distributions
- Relationships between variables
- Summary of very large data

Important complement to many statistical and computational techniques.

### **Graphics in R:**

- `plot ()` is the main function for graphics. The arguments can be a single point such as 0 or c (.3,.7), a single vector, a pair of vectors or many other R objects.
- `par ()` is another important function which defines the default settings for plots.
- There are many other plot functions which are specific to some tasks such as `hist ()`, `boxplot ()`, etc. Most of them take the same arguments as the `plot ()` function.

```
N <- 10^2
x1 <- rnorm(N)
x2 <- 1 + x1 + rnorm(N)
plot(0)
plot(0,1)
plot(x1)
plot(x1, x2) # scatter plot x1 on the horizontal axis and x2 on the vertical axis
plot(x2 ~ x1) # the same but using a formula (x2 as a function of x1)
```



## R Histograms

In this article, you'll learn to use `hist()` function to create histograms in R programming with the help of numerous examples.

Histogram can be created using the `hist()` function in R programming language. This function takes in a vector of values for which the histogram is plotted.

Let us use the built-in dataset `air quality` which has Daily air quality measurements in New York, May to September 1973.-R documentation.

`str(airquality)`

'data. frame':153 obs. of 6 variables: \$ Ozone: int 41 36 12 18 NA 28 23 19 8

NA ...\$

Solar.R: int 190 118 149 313 NA NA 299 99 19 194 ...\$

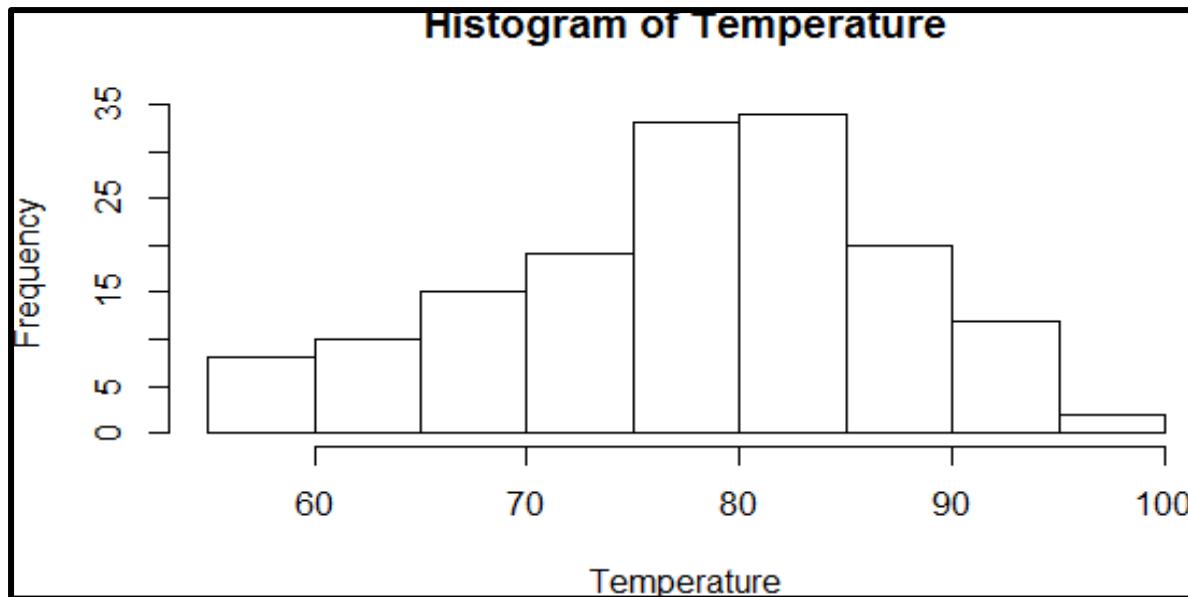
Wind: num 7.4 8 12.6 11.5 14.3 14.9 8.6 13.8 20.1 8.6 ...\$

Temp: int 67 72 74 62 56 66 65 59 61 69 ...\$

Month: int 5 5 5 5 5 5 5 5 5 ...\$

Day: int 1 2 3 4 5 6 7 8 9 10 ...

`Temperature <- airquality$Temp`hist (Temperature)



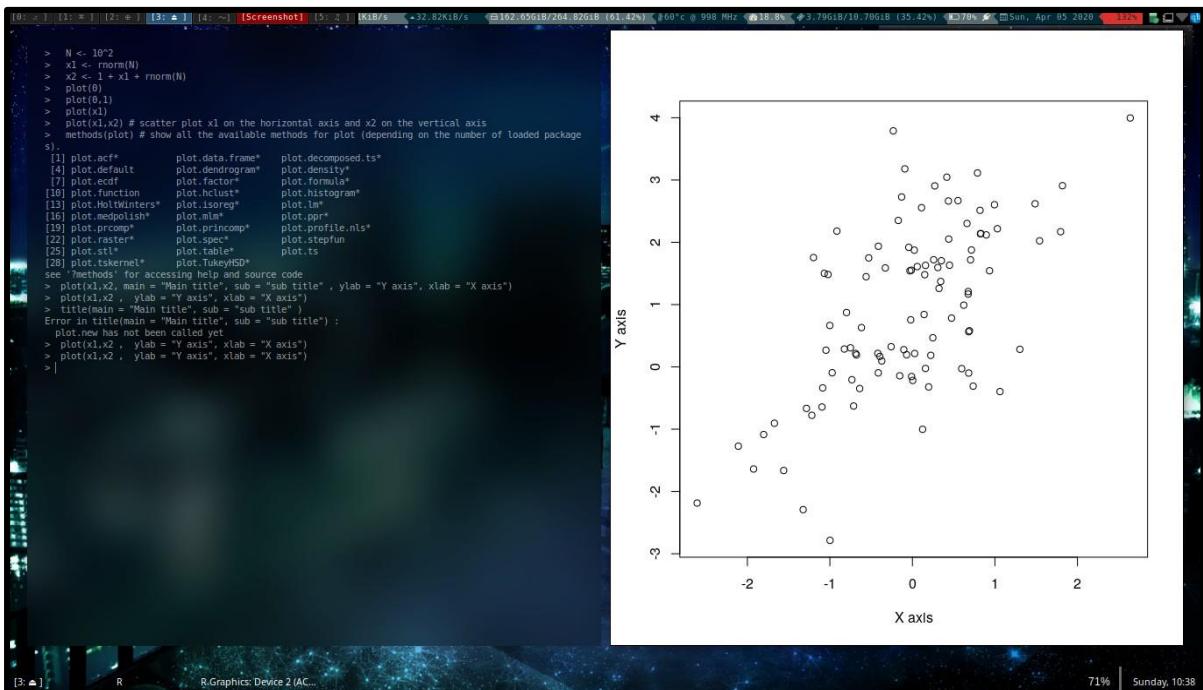
## Titles, legends and annotations

Titles main gives the main title, sub the subtitle. They can be passed as argument of the plot () function or using the title () function. xlab the name of the x axis and ylab the name of the y axis.

Plot (x1, x2, main = "Main title", sub = "sub title", ylab = "Y axis", xlab = "X axis")

plot (x1, x2, ylab = "Y axis", xlab = "X axis")

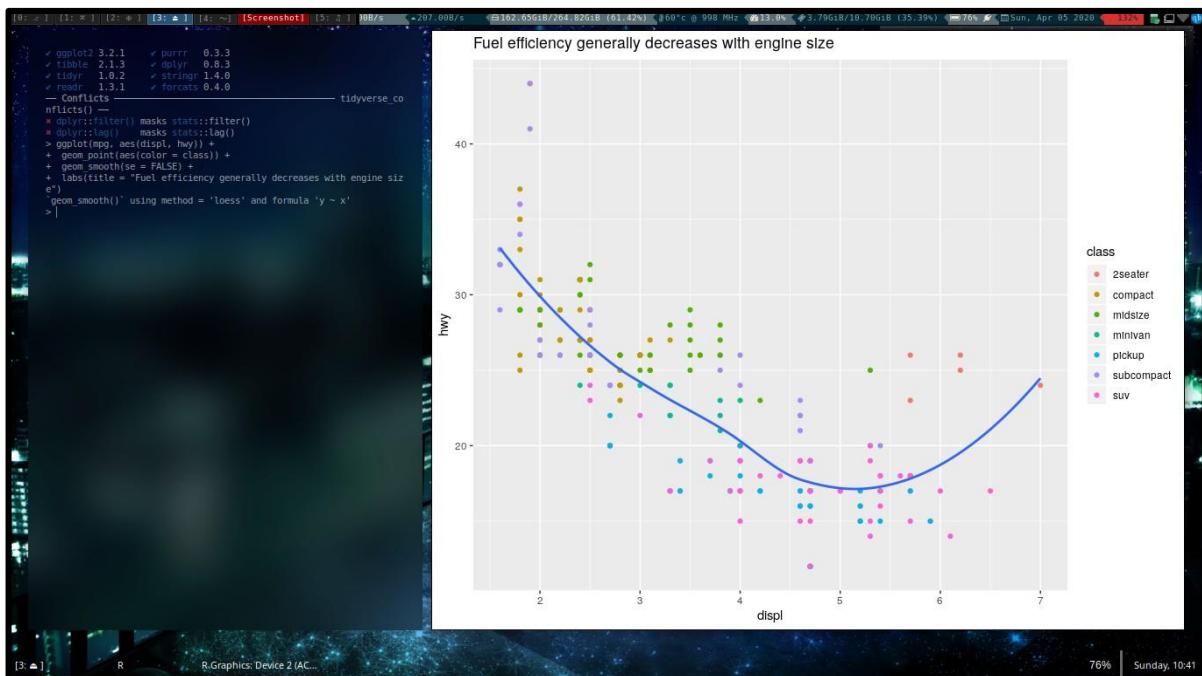
title (main = "Main title", sub = "sub title")



## Label

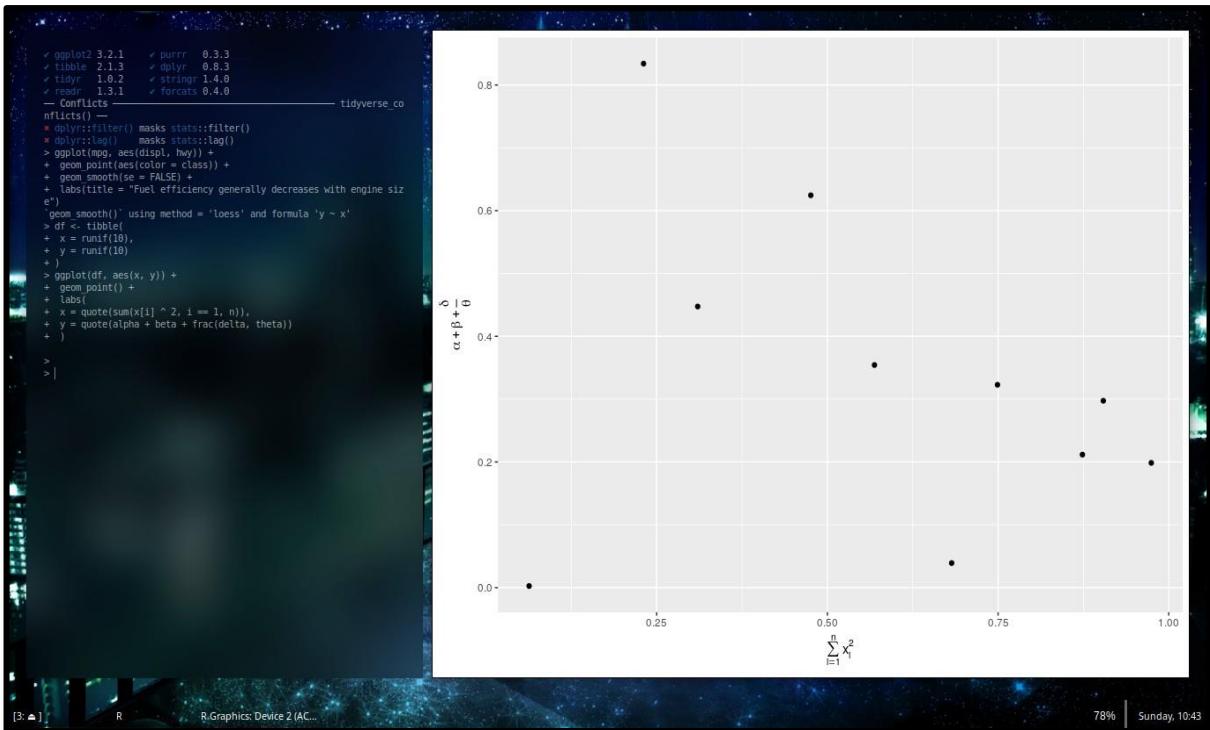
The easiest place to start when turning an exploratory graphic into an expository graphic is with good labels. You add labels with the labs () function. This example adds a plot title:

```
ggplot(mpg, aes(displ, hwy)) +  
  geom_point(aes(color = class)) +  
  geom_smooth(se = FALSE) +  
  labs(title = "Fuel efficiency generally decreases with engine size")
```



It's possible to use mathematical equations instead of text strings. Just switch "" out for quote () and read about the available options in? plotmath:

```
df <- tibble (  
  x = runif (10),  
  y = runif (10)  
)  
ggplot (df, aes (x, y)) +  
  geom_point () +  
  labs (  
    x = quote(sum(x[i] ^ 2, i == 1, n)),  
    y = quote (alpha + beta + frac (delta, theta))  
)
```



## Annotations

In addition to labelling major components of your plot, it's often useful to label individual observations or groups of observations. The first tool you have at your disposal is `geom_text()`. `geom_text()` is similar to `geom_point()`, but it has an additional aesthetic: `label`. This makes it possible to add textual labels to your plots.

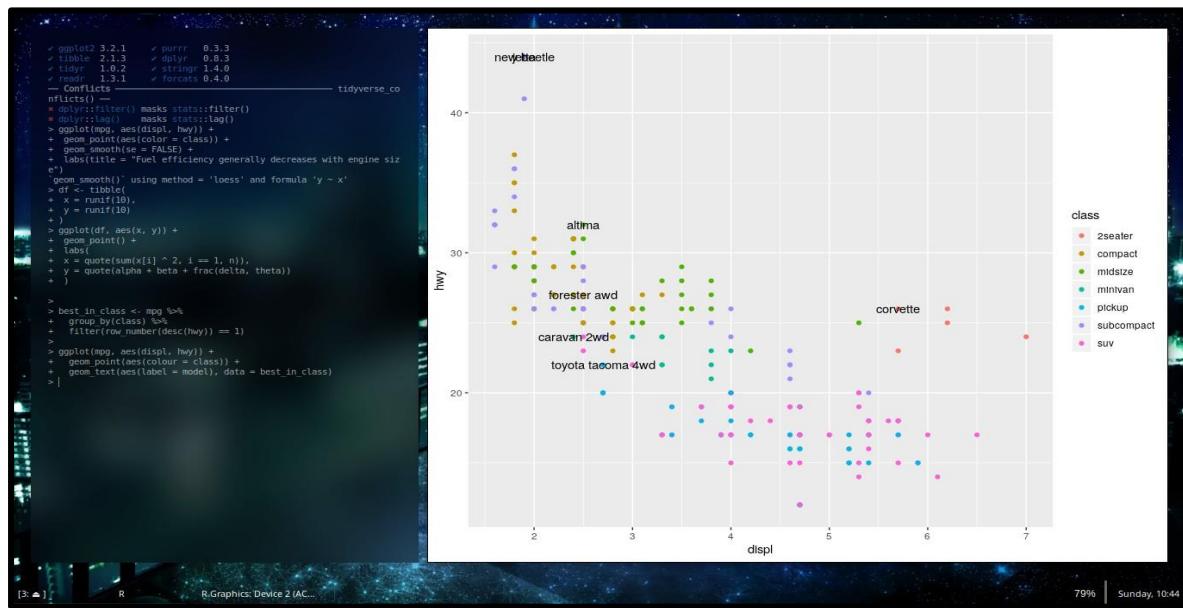
There are two possible sources of labels. First, you might have a tibble that provides labels. The plot below isn't terribly useful, but it illustrates a useful approach: pull out the most efficient car in each class with `dplyr`, and then label it on the plot:

```

best_in_class <- mpg %>%
  group_by(class) %>%
  filter(row_number(desc(hwy)) == 1)

ggplot (mpg, aes(displ, hwy)) +
  geom_point (aes(colour = class)) +
  geom_text (aes(label = model), data = best_in_class)

```



This is hard to read because the labels overlap with each other, and with the points. We can make things a little better by switching to `geom_label()` which draws a rectangle behind the text. We also use the `nudge_y` parameter to move the labels slightly above the corresponding points:

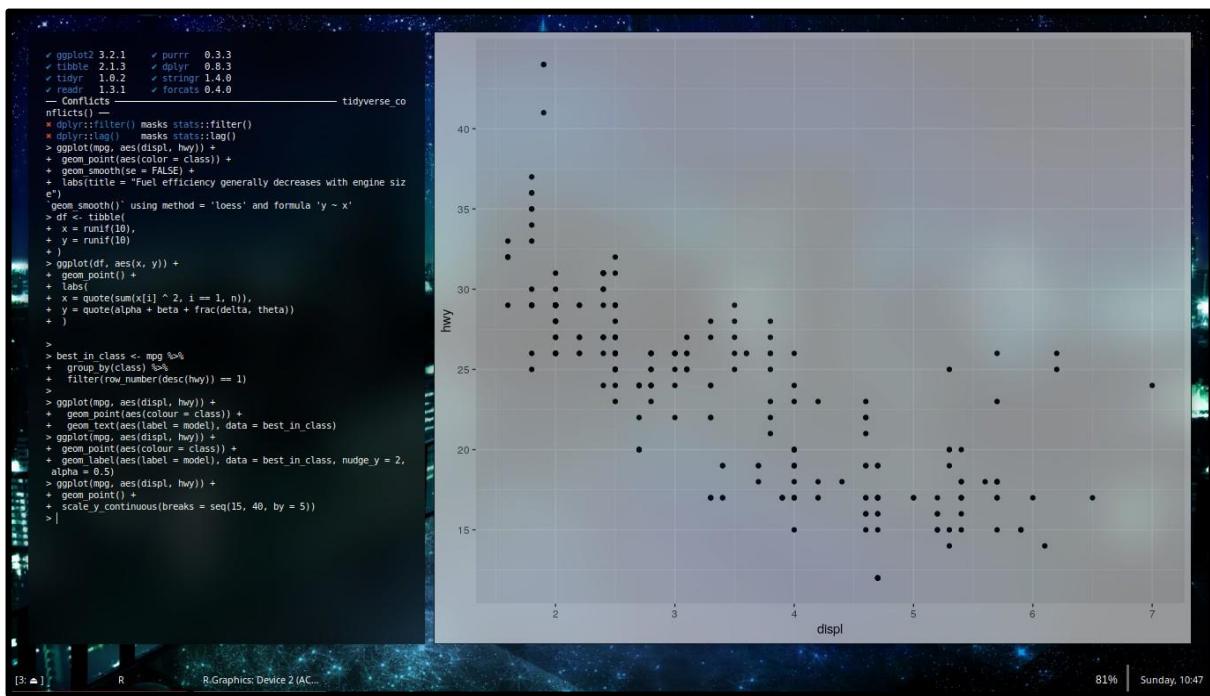
```
ggplot (mpg, aes(displ, hwy)) +
  geom_point (aes(colour = class)) +
  geom_label (aes(label = model), data = best_in_class, nudge_y = 2, alpha = 0.5)
```



## Axis ticks and legend keys

There are two primary arguments that affect the appearance of the ticks on the axes and the keys on the legend: breaks and labels. Breaks controls the position of the ticks, or the values associated with the keys. Labels controls the text label associated with each tick/key. The most common use of breaks is to override the default choice:

```
ggplot(mpg, aes(displ, hwy)) +  
  geom_point() +  
  scale_y_continuous(breaks = seq(15, 40, by = 5))
```

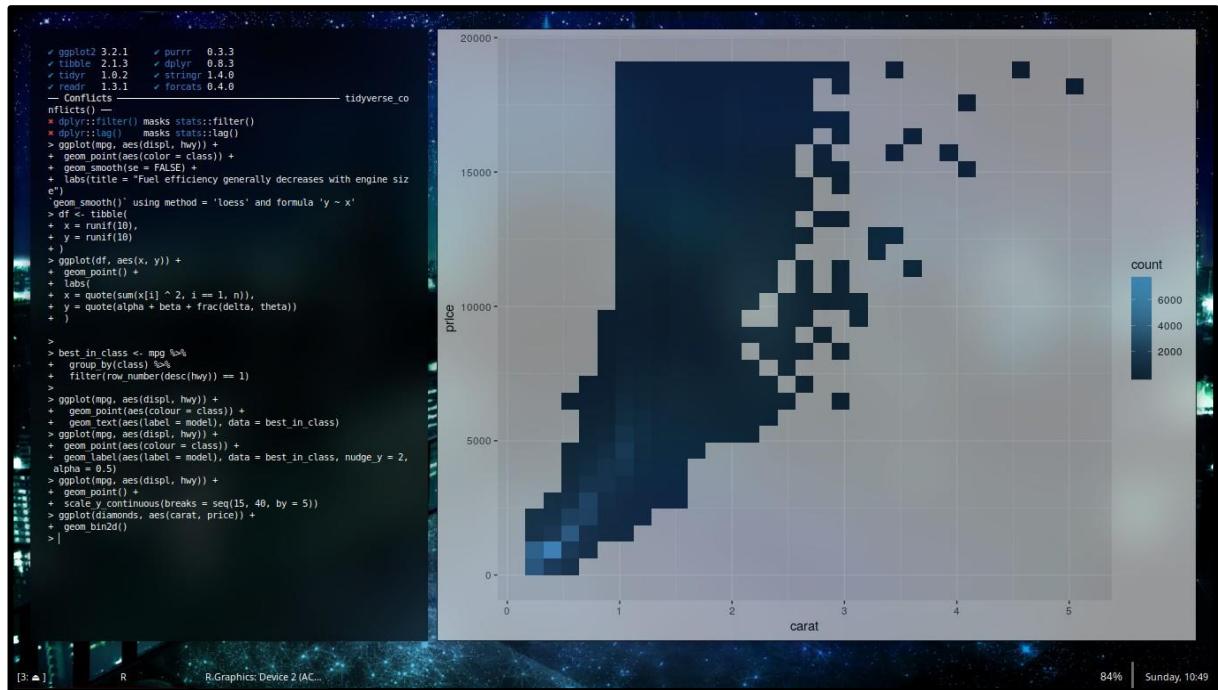


## Replacing a scale

Instead of just tweaking the details a little, you can instead replace the scale altogether. There are two types of scales you're mostly likely to want to switch out: continuous position scales and colour scales. Fortunately, the same principles apply to all the other aesthetics, so once you've mastered position and colour, you'll be able to quickly pick up other scale replacements.

It's very useful to plot transformations of your variable. For example, as we've seen in diamond prices it's easier to see the precise relationship between carat and price if we log transform them:

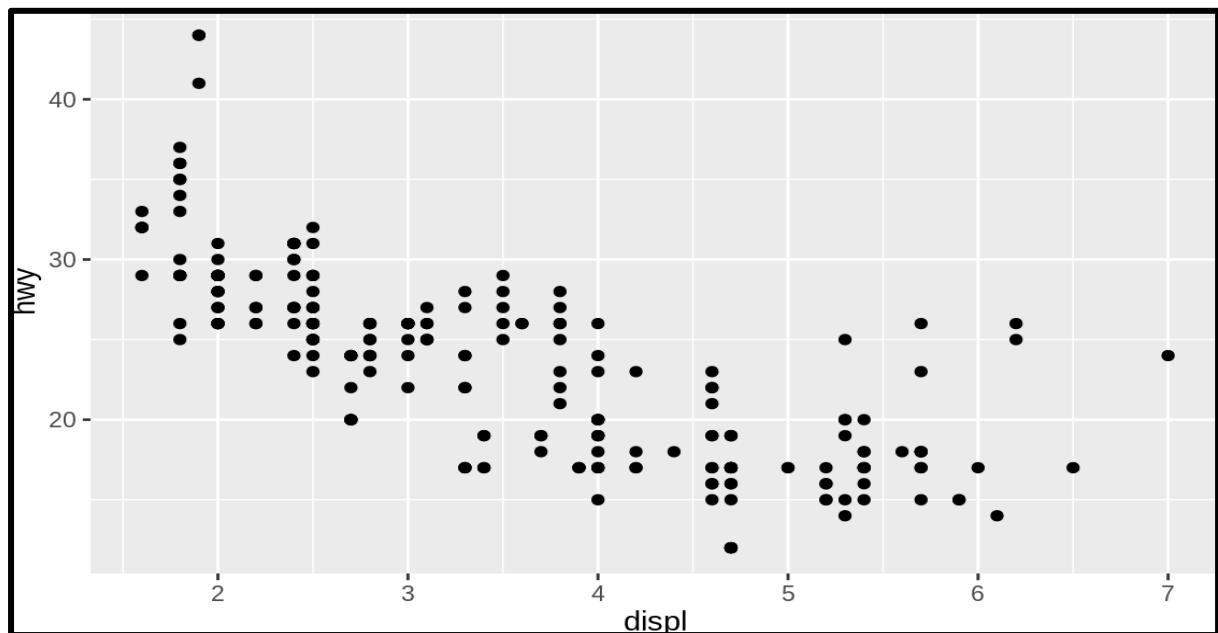
```
ggplot(diamonds, aes(carat, price)) +
  geom_bin2d()
```



## Saving your plots

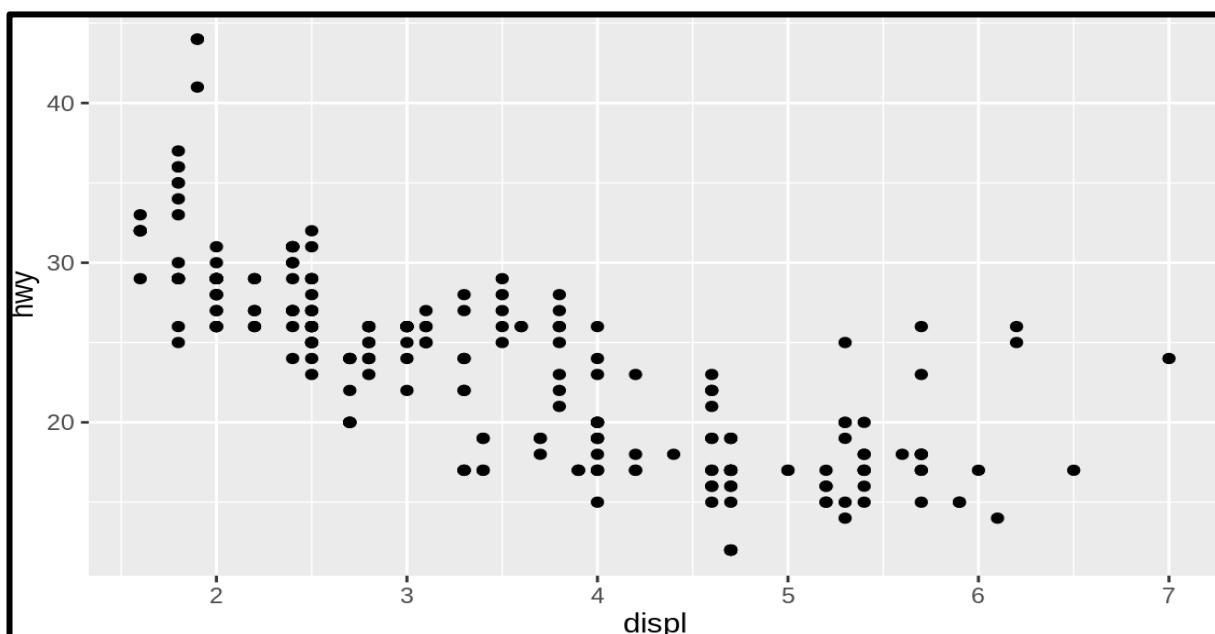
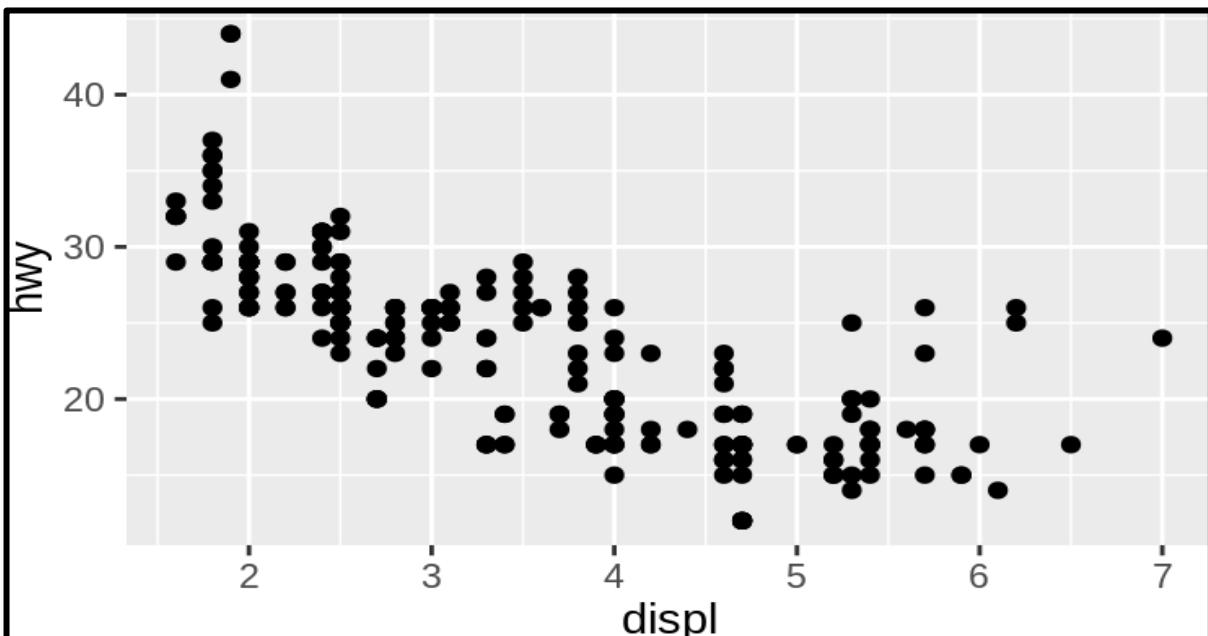
There are two main ways to get your plots out of R and into your final write-up:  
`ggsave()` and `knitr::ggsave()` will save the most recent plot to disk:

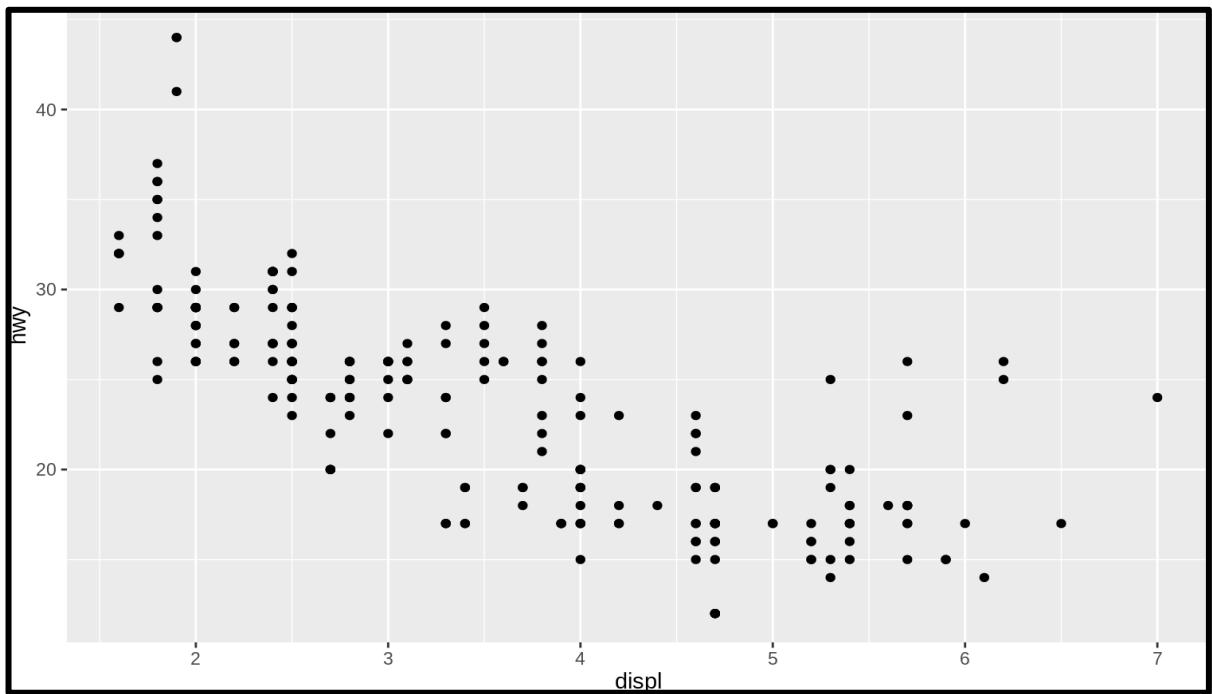
```
ggplot(mpg, aes(displ, hwy)) + geom_point()
```



## Figure sizing

The biggest challenge of graphics in R Markdown is getting your figures the right size and shape. There are five main options that control figure sizing: `fig.width`, `fig.height`, `fig.asp`, `out.width` and `out.height`. Image sizing is challenging because there are two sizes (the size of the figure created by R and the size at which it is inserted in the output document), and multiple ways of specifying the size (i.e., height, width, and aspect ratio: pick two of three).





**Conclusion:** Thus, we have studied basics of advanced graphics in R.

## Experiment no:9

**Aim:** Basic of correlation, simple and multiple regression in R.

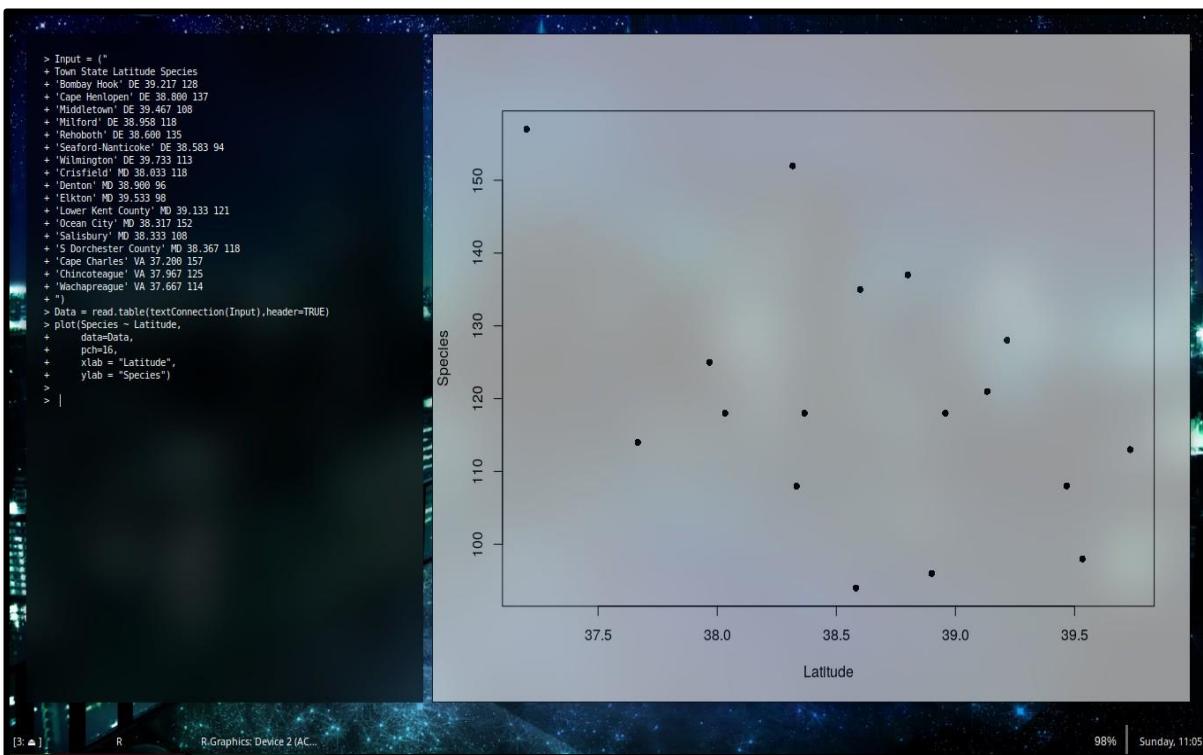
### Theory:

Correlation and linear regression each explore the relationship between two quantitative variables. Both are very common analyses.

Correlation determines if one variable varies systematically as another variable changes. It does not specify that one variable is the dependent variable and the other is the independent variable. Often, it is useful to look at which variables are correlated to others in a data set, and it is especially useful to see which variables correlate to a particular variable of interest.

So as to generate the simple plot for the data:

```
plot(Species ~ Latitude,  
      data=Data,  
      pch=16,  
      xlab = "Latitude",  
      ylab = "Species")
```



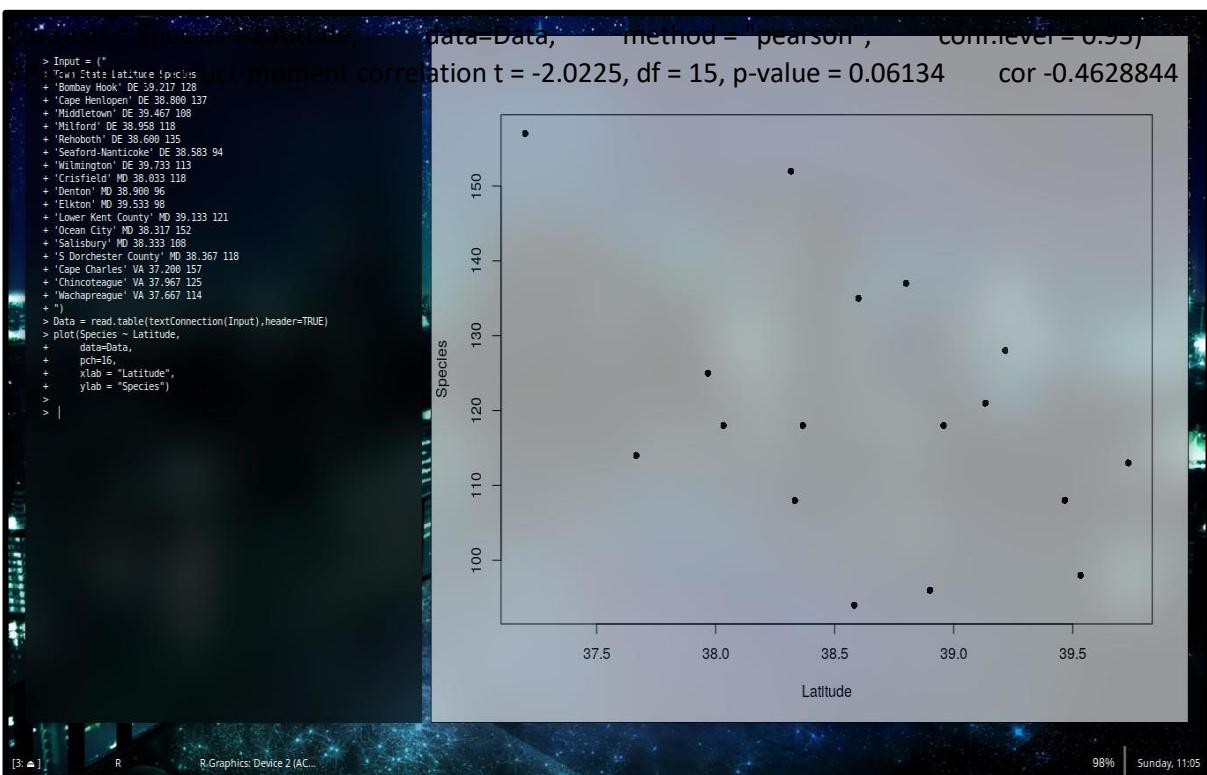
## Correlation

Correlation can be performed with the `cor.test` function in the native stats package. It can perform Pearson, Kendall, and Spearman correlation procedures. Methods for multiple correlation of several variables simultaneously are discussed in the Multiple regression chapter.

### Pearson correlation

Pearson correlation is the most common form of correlation. It is a parametric test, and assumes that the data are linearly related and that the residuals are normally distributed.

```
cor.test(~Species + Latitude,  
        data=Data,  
        method = "pearson",  
        conf.level = 0.95)  
  
#Pearson's product-moment correlation  
t = -2.0225, df = 15, p-value = 0.06134  
cor -0.4628844
```



## Kendall correlation

Kendall rank correlation is a non-parametric test that does not assume a distribution of the data or that the data are linearly related. It ranks the data to determine the degree of correlation.

```
cor.test (~ Species + Latitude,  
         data=Data,  
         method = "kendall",  
         continuity = FALSE,  
         conf. level = 0.95)
```

Kendall's rank correlation tau  
z = -1.3234, p-value = 0.1857  
tau -0.2388326

## Spearman correlation

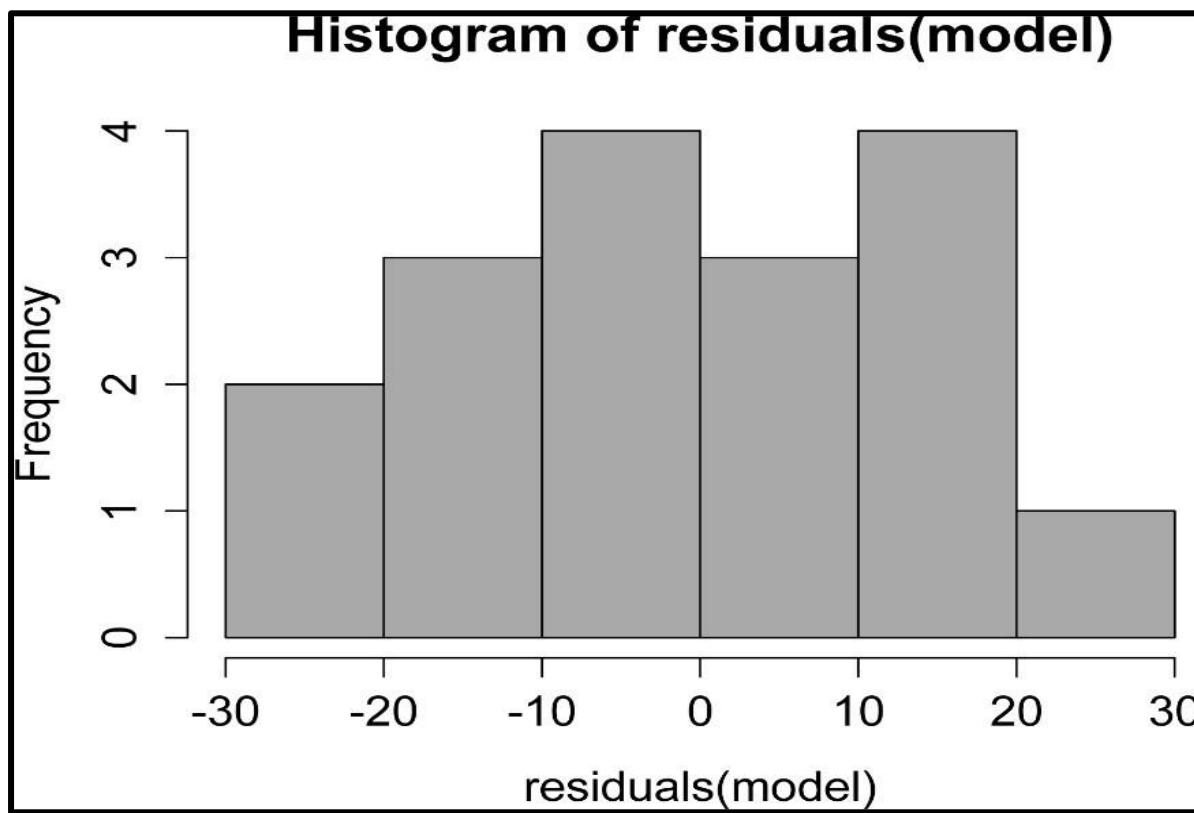
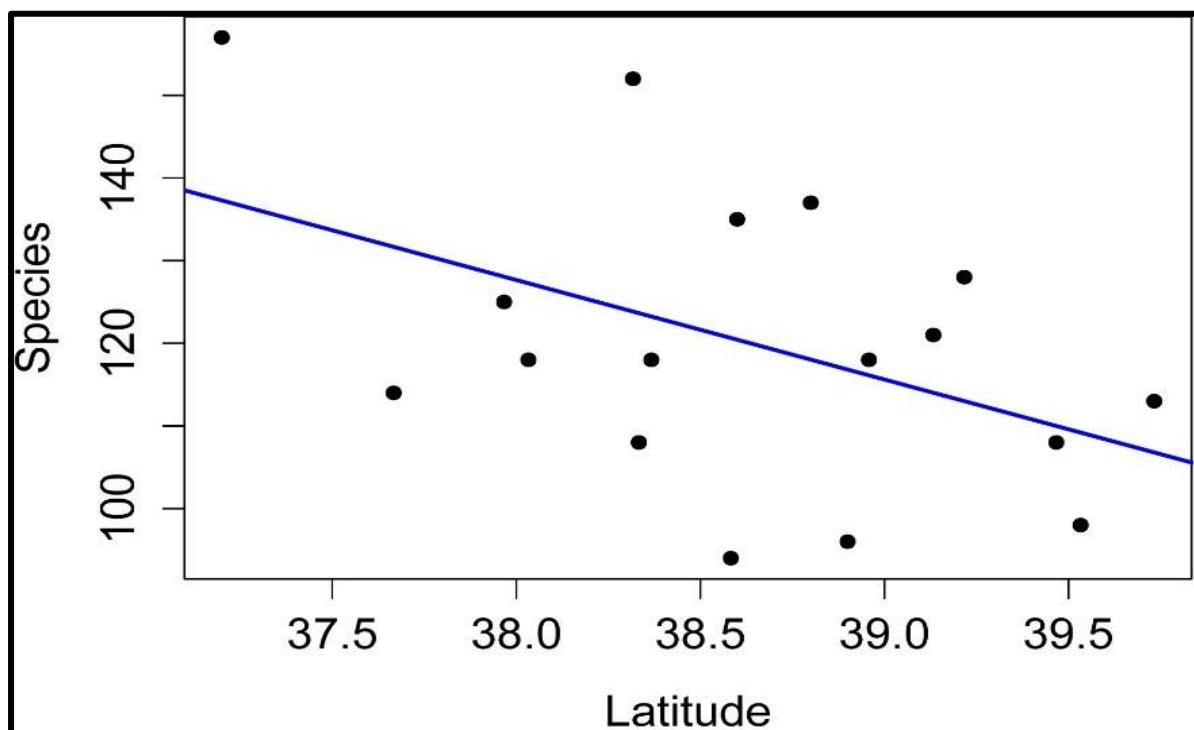
Spearman rank correlation is a non-parametric test that does not assume a distribution of the data or that the data are linearly related. It ranks the data to determine the degree of correlation, and is appropriate for ordinal measurements.

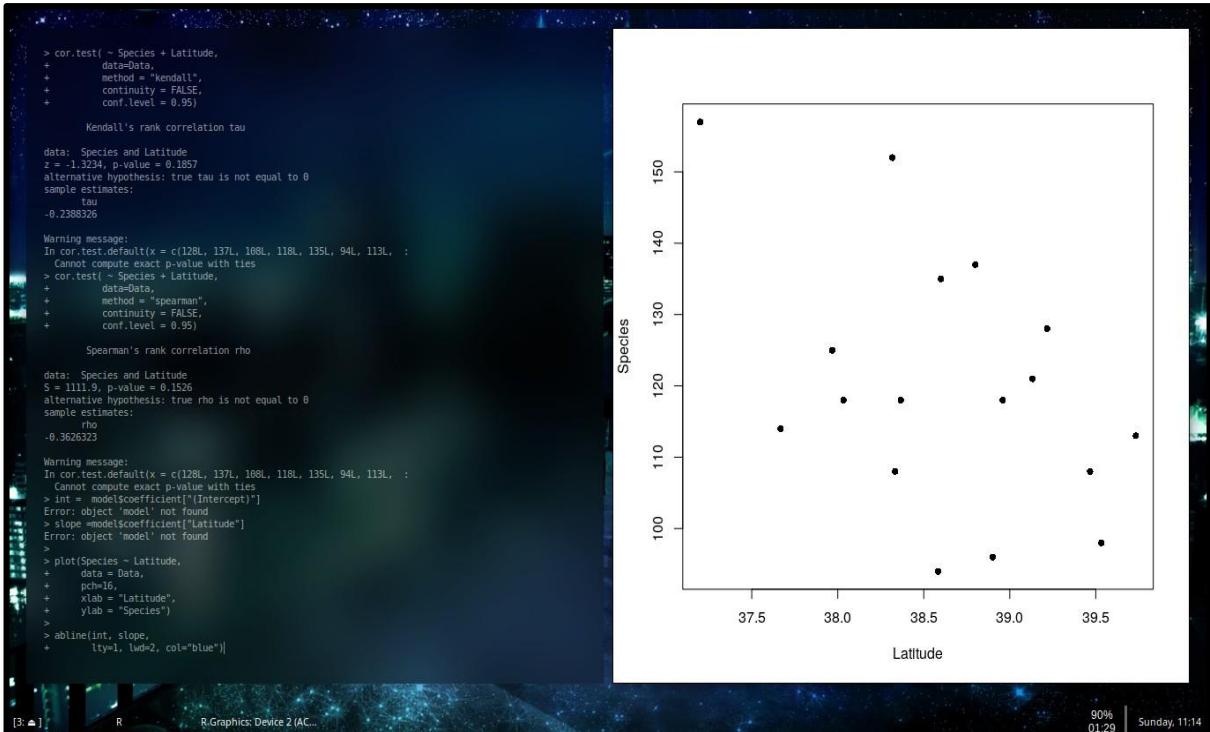
```
cor.test (~ Species + Latitude,  
         data=Data,  
         method = "spearman",  
         continuity = FALSE,  
         conf. level = 0.95)
```

Spearman's rank correlation rho  
S = 1111.908, p-value = 0.1526  
rho -0.3626323

## Linear regression

Linear regression can be performed with the lm function in the native stats package. A robust regression can be performed with the lmrob function in the robust base package.





**Conclusion:** Thus, we have successfully studied Basics of correlation and regression in R.

## Experiment no:10

**Aim:** Clustering in R: K-means

### **Theory:**

K Means Clustering is an unsupervised learning algorithm that tries to cluster data based on their similarity. Unsupervised learning means that there is no outcome to be predicted, and the algorithm just tries to find patterns in the data. In k means clustering, we have to specify the number of clusters we want the data to be grouped into. The algorithm randomly assigns each observation to a cluster, and finds the centroid of each cluster. Then, the algorithm iterates through two steps:

- Reassign data points to the cluster whose centroid is closest.
- Calculate new centroid of each cluster.

These two steps are repeated till the within cluster variation cannot be reduced any further. The within cluster variation is calculated as the sum of the Euclidean distance between the data points and their respective cluster centroids.

### **Exploring the data**

The iris dataset contains data about sepal length, sepal width, petal length, and petal width of flowers of different species.

Let us see what it looks like:

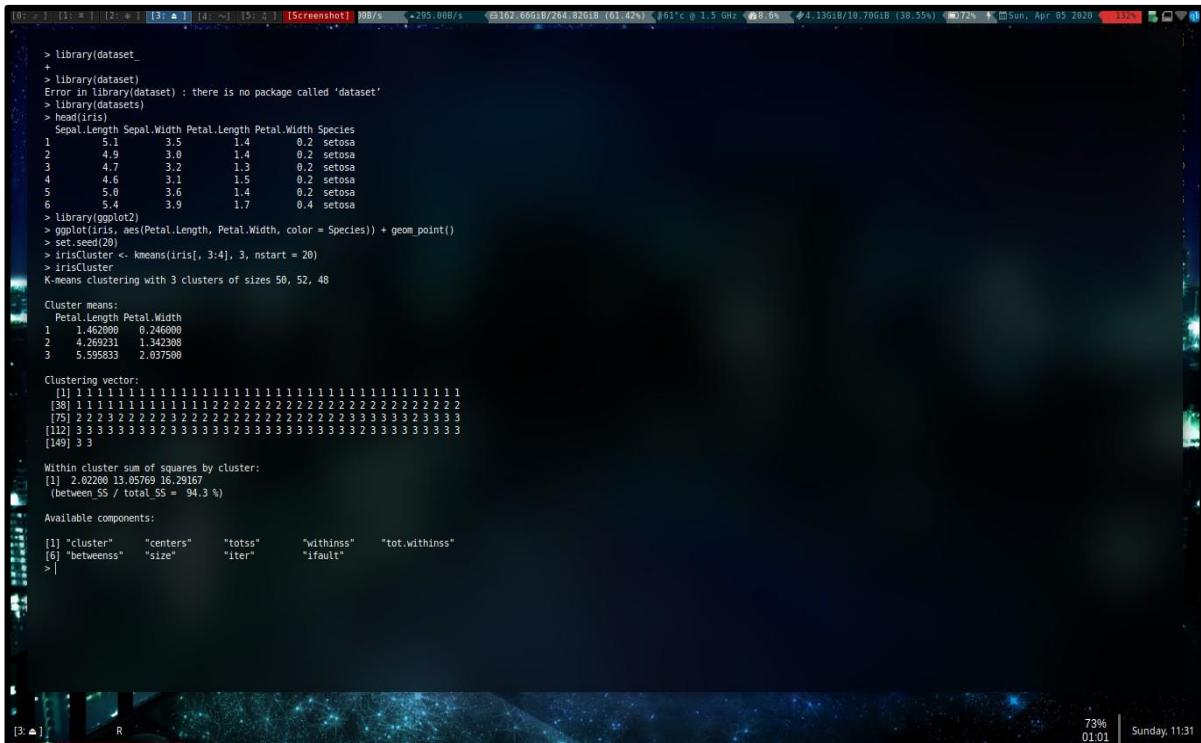
```
library(datasets)
head(iris)

  Sepal.Length Sepal.Width Petal.Length Petal.Width Species
1          5.1         3.5          1.4         0.2   setosa
2          4.9         3.0          1.4         0.2   setosa
3          4.7         3.2          1.3         0.2   setosa
4          4.6         3.1          1.5         0.2   setosa
5          5.0         3.6          1.4         0.2   setosa
6          5.4         3.9          1.7         0.4   setosa
```

```
library(ggplot2)
ggplot(iris, aes(Petal.Length, Petal.Width, color = Species)) + geom_point()
```



Since we know that there are 3 species involved, we ask the algorithm to group the data into 3 clusters, and since the starting assignments are random, we specify nstart = 20 . This means that R will try 20 different random starting assignments and then select the one with the lowest within cluster variation. We can see the cluster centroids, the clusters that each data point was assigned to, and the within cluster variation.



A screenshot of a terminal window on a Mac OS X desktop. The terminal is showing R code being run. The code performs the following steps:

- Loads the 'iris' dataset (using head(iris) to view it).
- Creates a scatter plot of Sepal.Length vs Sepal.Width, colored by Species.
- Runs k-means clustering with 3 clusters and nstart = 20.
- Prints the cluster means for Petal.Length and Petal.Width.
- Prints the clustering vector, which maps each of the 150 data points to a cluster index (1, 2, or 3).
- Prints the within cluster sum of squares by cluster.
- Shows the available components of the k-means object.

The R console shows the results of these operations, including the cluster assignments and the final within-cluster sum of squares.

Let us compare the clusters with the species.

```
table(irisCluster$cluster, iris$Species)
    setosa versicolor virginica
1      0          2       44
2      0         48        6
3     50          0        0
```

As we can see, the data belonging to the setosa species got grouped into cluster 3, versicolor into cluster 2, and virginica into cluster 1. The algorithm wrongly classified two data points belonging to versicolor and six data points belonging to virginica.

We can also plot the data to see the clusters:

```
irisCluster$cluster <- as.factor(irisCluster$cluster)
ggplot(iris, aes(Petal.Length, Petal.Width, color = iris$cluster)) +
  geom_point()
```



**Conclusion:** Thus, we have successfully studied K-means clustering on iris dataset in R.

## Experiment no:11

**Aim:** Classification on large and noisy dataset with R - logistic regression and naive Bayes.

### Theory:

#### Logic Regression:

Logistic regression is yet another technique borrowed by machine learning from the field of statistics. It's a powerful statistical way of modelling a binomial outcome with one or more explanatory variables. It measures the relationship between the categorical dependent variable and one or more independent variables by estimating probabilities using a logistic function, which is the cumulative logistic distribution.

Logistic regression is an instance of classification technique that you can use to predict a qualitative response. More specifically, logistic regression models the probability that belongs to a particular category.

That means that, if you are trying to do gender classification, where the response falls into one of the two categories, male or female, you'll use logistic regression models to estimate the probability that belongs to a particular category.

#### Naïve Bayes Classifier

$$P(c|x) = \frac{P(x|c)P(c)}{P(x)}$$

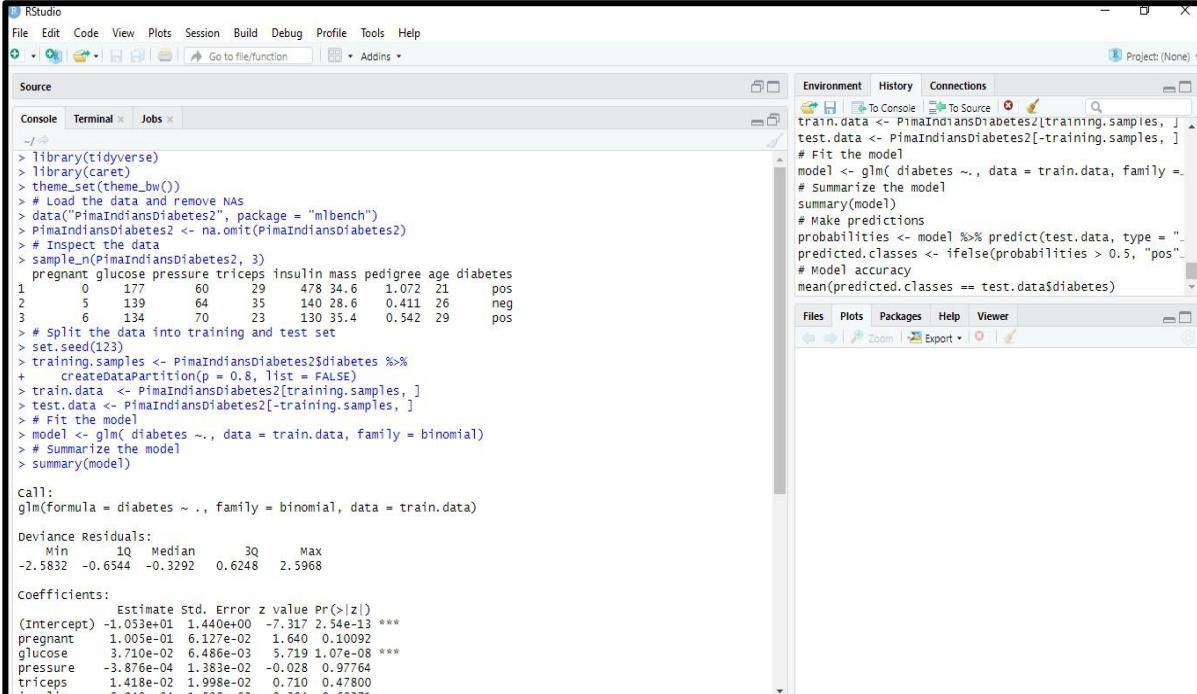
**Likelihood**                    **Class Prior Probability**  
**Posterior Probability**            **Predictor Prior Probability**

$$P(c|X) = P(x_1|c) \times P(x_2|c) \times \dots \times P(x_n|c) \times P(c)$$

The Naïve Bayes classifier is a simple probabilistic classifier which is based on Bayes theorem but with strong assumptions regarding independence. Historically, this technique became popular with applications in email filtering, spam detection, and document categorization. Although it is often

outperformed by other techniques, and despite the naïve design and oversimplified assumptions, this classifier can perform well in many complex real-world problems. And since it is a resource efficient algorithm that is fast and scales well, it is definitely a machine learning algorithm to have in your toolkit.

## Output:



```

RStudio
File Edit Code View Plots Session Build Debug Profile Tools Help
Go to file/function Addins Project: (None)

Source
Console Terminal Jobs
~/
> library(tidyverse)
> library(caret)
> theme_set(theme_bw())
> # Load the data and remove NAs
> data("PimaIndiansDiabetes")
> PimaIndiansDiabetes2 <- na.omit(PimaIndiansDiabetes)
> # Inspect the data
> sample_n(PimaIndiansDiabetes2, 3)
#> # Split the data into training and test set
> set.seed(123)
> training.samples <- PimaIndiansDiabetes2$diabetes %>%
+   createDataPartition(p = 0.8, list = FALSE)
> train.data <- PimaIndiansDiabetes2[training.samples, ]
> test.data <- PimaIndiansDiabetes2[-training.samples, ]
> # Fit the model
> model <- glm(diabetes ~ ., data = train.data, family = binomial)
> # Summarize the model
> summary(model)

Call:
glm(formula = diabetes ~ ., family = binomial, data = train.data)

Deviance Residuals:
    Min      1Q  Median      3Q     Max 
-2.5832 -0.6544 -0.3292  0.6248  2.5968 

Coefficients:
            Estimate Std. Error z value Pr(>|z|)    
(Intercept) -1.053e+01  1.440e+00 -7.317 2.54e-13 ***
pregnant    1.005e-01  6.127e-02  1.640  0.10092    
glucose     3.710e-02  6.486e-03  5.719 1.07e-08 ***
pressure   -3.876e-04  1.383e-02 -0.028  0.97764    
triceps    1.418e-02  1.998e-02  0.710  0.47800    
insulin     5.940e-04  1.508e-03  0.394  0.69371    
mass        7.997e-02  3.180e-02  2.515  0.01190 *  
pedigree   1.329e+00  4.823e-01  2.736  0.00385 ** 
age         2.718e-02  2.020e-02  1.346  0.17840    
---
signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

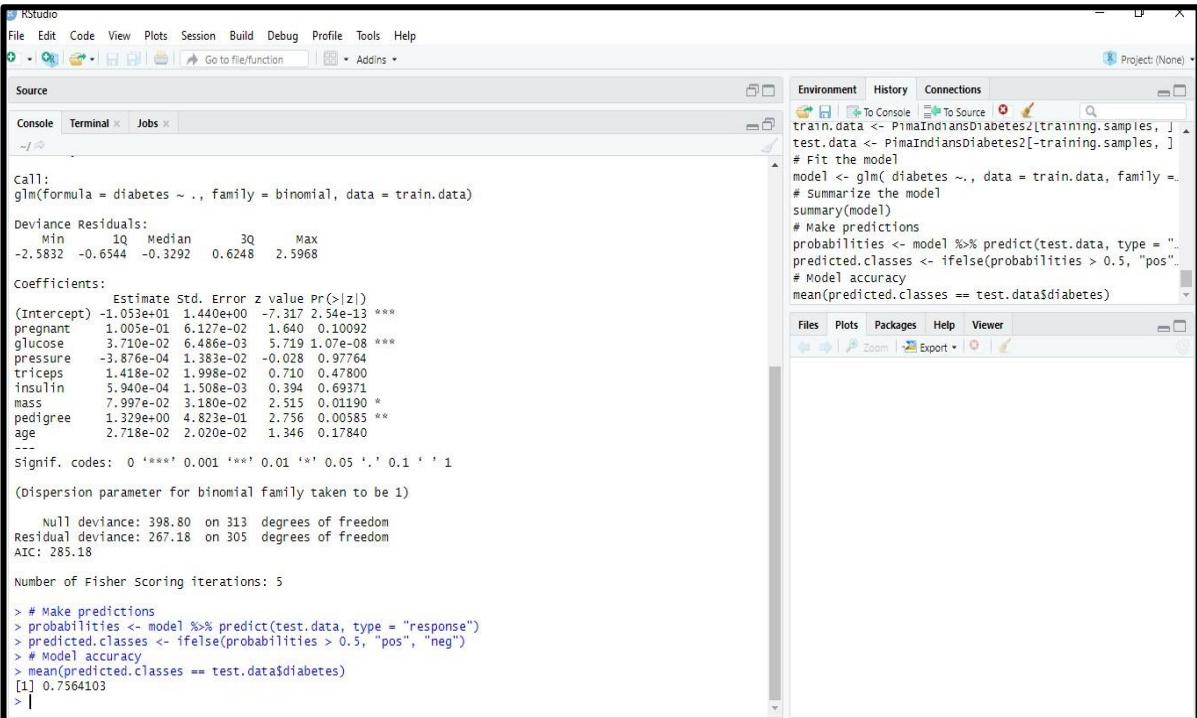
Dispersion parameter for binomial family taken to be 1

Null deviance: 398.80 on 313 degrees of freedom
Residual deviance: 267.18 on 305 degrees of freedom
AIC: 285.18

Number of Fisher scoring iterations: 5

> # Make predictions
> probabilities <- model %>% predict(test.data, type = "response")
> predicted.classes <- ifelse(probabilities > 0.5, "pos", "neg")
> # Model accuracy
> mean(predicted.classes == test.data$diabetes)
[1] 0.7564103
>

```



```

RStudio
File Edit Code View Plots Session Build Debug Profile Tools Help
Go to file/function Addins Project: (None)

Source
Console Terminal Jobs
~/
Call:
glm(formula = diabetes ~ ., family = binomial, data = train.data)

Deviance Residuals:
    Min      1Q  Median      3Q     Max 
-2.5832 -0.6544 -0.3292  0.6248  2.5968 

Coefficients:
            Estimate Std. Error z value Pr(>|z|)    
(Intercept) -1.053e+01  1.440e+00 -7.317 2.54e-13 ***
pregnant    1.005e-01  6.127e-02  1.640  0.10092    
glucose     3.710e-02  6.486e-03  5.719 1.07e-08 ***
pressure   -3.876e-04  1.383e-02 -0.028  0.97764    
triceps    1.418e-02  1.998e-02  0.710  0.47800    
insulin     5.940e-04  1.508e-03  0.394  0.69371    
mass        7.997e-02  3.180e-02  2.515  0.01190 *  
pedigree   1.329e+00  4.823e-01  2.736  0.00385 ** 
age         2.718e-02  2.020e-02  1.346  0.17840    
---
signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Dispersion parameter for binomial family taken to be 1

Null deviance: 398.80 on 313 degrees of freedom
Residual deviance: 267.18 on 305 degrees of freedom
AIC: 285.18

Number of Fisher scoring iterations: 5

> # Make predictions
> probabilities <- model %>% predict(test.data, type = "response")
> predicted.classes <- ifelse(probabilities > 0.5, "pos", "neg")
> # Model accuracy
> mean(predicted.classes == test.data$diabetes)
[1] 0.7564103
>

```

RStudio

File Edit Code View Plots Session Build Debug Profile Tools Help

Console Terminal Jobs

```

~/
> # Make predictions
> probabilities <- model %>% predict(test.data, type = "response")
> predicted.classes <- ifelse(probabilities > 0.5, "pos", "neg")
> # Model accuracy
> mean(predicted.classes == test.data$diabetes)
[1] 0.7564103
> model <- glm(diabetes ~ glucose, data = train.data, family = binomial)
> summary(model)$coef
  Estimate Std. Error z value Pr(>|z|)
(Intercept) -6.15882009 0.700096646 -8.797100 1.403974e-18
glucose     0.04327234 0.005341133 8.101716 5.418949e-16
> newdata <- data.frame(glucose = c(20, 180))
> probabilities <- model %>% predict(newdata, type = "response")
> predicted.classes <- ifelse(probabilities > 0.5, "pos", "neg")
> predicted.classes
  1   2 
"neg" "pos"
> train.data %>%
+   mutate(prob = ifelse(diabetes == "pos", 1, 0)) %>%
+   ggplot(aes(glucose, prob)) +
+   geom_point(alpha = 0.2) +
+   geom_smooth(method = "glm", method.args = list(family =
+     binomial))
> +
+   labs(
+     title = "Logistic Regression Model",
+     x = "Plasma Glucose Concentration",
+     y = "Probability of being diabete-pos"
+   )
> geom_smooth()^ using formula 'y ~ x'
> model <- glm(diabetes ~ glucose + mass + pregnant,
+               data = train.data, family = binomial)
> summary(model)$coef
  Estimate Std. Error z value Pr(>|z|)
(Intercept) -9.32369818 1.125997285 -8.280391 1.227711e-16
glucose     0.03886154 0.005404219 7.190962 6.433636e-13
mass        0.09458458 0.023529905 4.019760 5.823738e-05
pregnant    0.14466661 0.045125/29 3.20585/ 1.346611e-03
> |

```

Files Plots Packages Help Viewer

**Logistic Regression Model**

Probability of being diabetic-pos

Plasma Glucose Concentration

7:49 PM 3/16/2020

RStudio

File Edit Code View Plots Session Build Debug Profile Tools Help

Console Terminal Jobs

```

~/
> mass      0.09458458 0.023529905 4.019760 5.823738e-05
> pregnant  0.14466661 0.045125729 3.205857 1.346611e-03
> model <- glm(diabetes ~., data = train.data, family = binomial)
> summary(model)$coef
  Estimate Std. Error z value Pr(>|z|)
(Intercept) -0.153400e+01 1.439679266 -7.31690975 2.537464e-13
pregnant    1.005031e-01 0.061266974 1.64041157 1.009196e-01
glucose     3.709621e-02 0.006486093 5.71934633 1.069346e-08
pressure   -3.875933e-04 0.013826185 -0.02803328 9.776356e-01
triceps    1.417771e-02 0.019981885 0.70952823 4.779967e-01
insulin     5.939876e-04 0.001508231 0.39383055 6.937061e-01
mass       7.997447e-02 0.031798907 2.51500698 1.190300e-02
pedigree   1.329149e+00 0.482291020 2.75590704 5.852963e-03
age        2.718224e-02 0.020199295 1.34570257 1.783985e-01
> coef(model)
(Intercept) pregnant      glucose      pressure      triceps      insulin      mass
-1.053400e+01 1.005031e-01 3.709621e-02 -3.875933e-04 1.417771e-02 5.939876e-04 7.997447e-02
pedigree      age
1.329149e+00 2.718224e-02
> summary(model)$coef
  Estimate Std. Error z value Pr(>|z|)
(Intercept) -0.153400e+01 1.439679266 -7.31690975 2.537464e-13
pregnant    1.005031e-01 0.061266974 1.64041157 1.009196e-01
glucose     3.709621e-02 0.006486093 5.71934633 1.069346e-08
pressure   -3.875933e-04 0.013826185 -0.02803328 9.776356e-01
triceps    1.417771e-02 0.019981885 0.70952823 4.779967e-01
insulin     5.939876e-04 0.001508231 0.39383055 6.937061e-01
mass       7.997447e-02 0.031798907 2.51500698 1.190300e-02
pedigree   1.329149e+00 0.482291020 2.75590704 5.852963e-03
age        2.718224e-02 0.020199295 1.34570257 1.783985e-01
> model <- glm(diabetes ~ pregnant + glucose + pressure + mass + pedigree,
+               data = train.data, family = binomial)
> probabilities <- model %>% predict(test.data, type = "response")
> head(probabilities)
  19    21    32    55    64    71 
0.1352603 0.5127526 0.6795376 0.7517408 0.2734867 0.1648174
> |

```

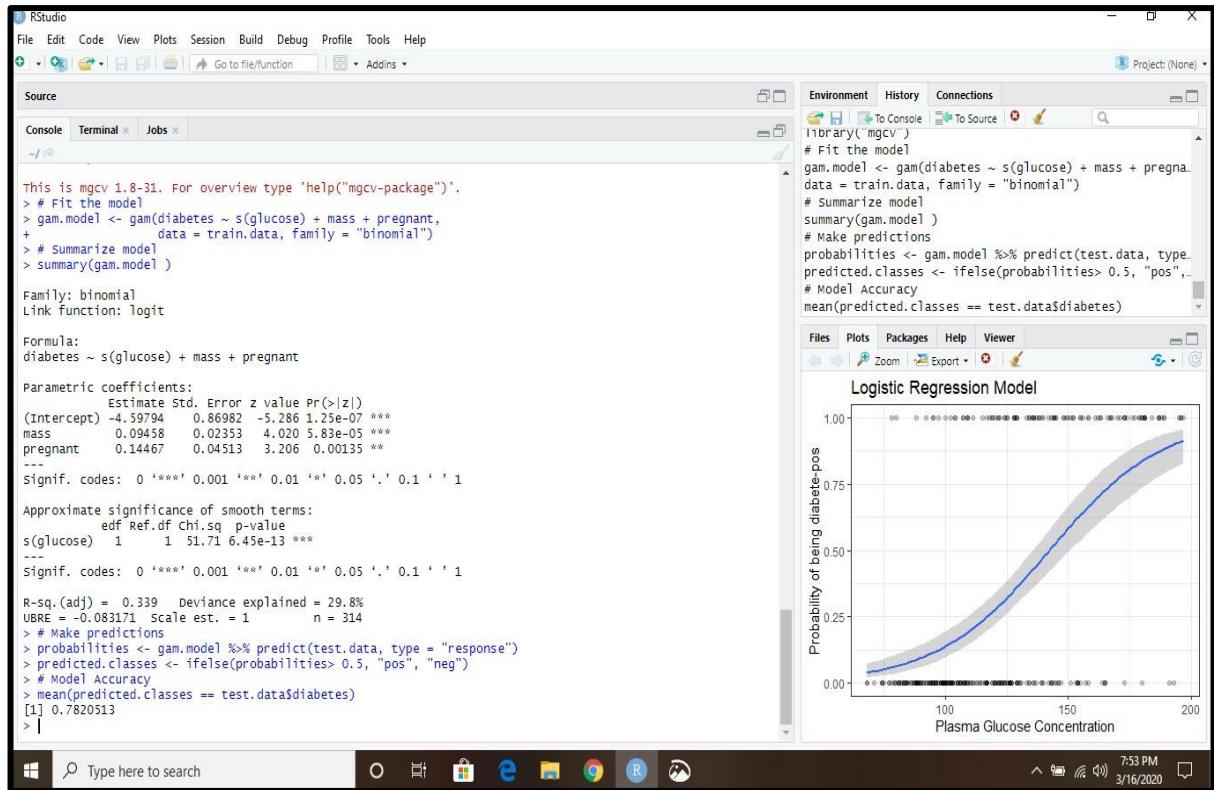
Files Plots Packages Help Viewer

**Logistic Regression Model**

Probability of being diabetic-pos

Plasma Glucose Concentration

7:51 PM 3/16/2020



**Conclusion:** Thus, we have studied Classification on large and noisy dataset with R.

## Experiment no:12

### Project: Mini-Project for R-Lab

#### Project-Name: Salary Prediction

#### Dataset:

age	workclass	fnlwgt	education	education-num	marital-status	occupation	relationship	race	sex	capital-gain	capital-loss	house-per-native	cou-income	native-country
39	State-gov	77516	Bachelors	13	Never-married	Adm-clerical	Not-in-family	White	Male	2174	0	40	United-States	<=50K
50	Self-employed	83311	Bachelors	13	Married-civ-spouse	Exe-manag-sell-relate	Unmarried	White	Male	0	0	13	United-States	<=50K
38	Private	215646	HS-grad	9	Divorced	Handlers-cleaners	Not-in-family	White	Male	0	0	40	United-States	<=50K
53	Private	234721	11th	7	Married-civ-spouse	Husband	Black	Male	0	0	40	United-States	<=50K	
28	Private	338409	Bachelors	13	Married-civ-spouse	Prof-specialty	Wife	Black	Female	0	0	40	Cuba	<=50K
37	Private	284582	Masters	14	Married-civ-spouse	Exec-managerial	Wife	White	Female	0	0	40	United-States	<=50K
49	Private	160187	9th	5	Married-civ-spouse	Other-service	Not-in-family	Black	Female	0	0	16	Jamaica	<=50K
52	Self-employed	209642	HS-grad	9	Married-civ-spouse	Exe-manag-sell-relate	Unmarried	White	Male	0	0	45	United-States	>50K
31	Private	45781	Masters	14	Never-married	Prof-specialty	Not-in-family	White	Female	14084	0	50	United-States	>50K
42	Private	159449	Bachelors	13	Married-civ-spouse	Husband	White	Male	5178	0	40	United-States	>50K	
37	Private	280464	Some-college	10	Married-civ-spouse	Husband	Black	Male	0	0	80	United-States	>50K	
30	State-gov	141297	Bachelors	13	Married-civ-spouse	Prof-specialty	Husband	Asian-Pac Islander	Male	0	0	40	India	>50K
23	Private	122272	Bachelors	13	Never-married	Adm-clerical	Own-child	White	Female	0	0	30	United-States	<=50K
32	Private	205019	Assoc-accel	12	Never-married	Sales	Not-in-family	Black	Male	0	0	50	United-States	<=50K
40	Private	121772	Assoc-voc	11	Married-civ-spouse	Craft-repair	Husband	Asian-Pac Islander	Male	0	0	40	?	>50K
34	Private	245487	7th-8th	4	Married-civ-spouse	Transportation-moving	Husband	American-Indian	Male	0	0	45	Mexico	<=50K
25	Self-employed	176756	HS-grad	9	Never-married	Farming-fishing	Own-child	White	Male	0	0	35	United-States	<=50K
32	Private	186824	HS-grad	9	Never-married	Machine-guardian	Unmarried	White	Male	0	0	40	United-States	<=50K
38	Private	28887	11th	7	Married-civ-spouse	Sales	Husband	White	Male	0	0	50	United-States	<=50K
43	Self-employed	292175	Masters	14	Divorced	Exe-manag-sell-relate	Unmarried	White	Female	0	0	45	United-States	>50K
40	Private	102534	Doctorate	16	Married-civ-spouse	Prof-specialty	Husband	White	Male	0	0	60	United-States	<=50K

age	workclass	fnlwgt	education	education-num	marital-status	occupation	relationship	race	sex	capital-gain	capital-loss	house-per-native	cou-income	native-country
38	Private	28887	11th	7	Married-civ-spouse	Sales	Husband	White	Male	0	0	50	United-States	<=50K
43	Self-employed	292175	Masters	14	Divorced	Exe-manag-sell-relate	Unmarried	White	Female	0	0	45	United-States	>50K
40	Private	193524	Doctorate	16	Married-civ-spouse	Prof-specialty	Husband	White	Male	0	0	60	United-States	>50K
54	Private	302146	HS-grad	9	Separated	Other-service	Unmarried	Black	Female	0	0	20	United-States	<=50K
35	Federal-gov	76845	9th	5	Married-civ-spouse	Farming-fishing	Husband	Black	Male	0	0	40	United-States	<=50K
43	Private	117037	11th	7	Married-civ-spouse	Transportation-moving	Husband	White	Male	0	2042	40	United-States	<=50K
59	Private	109015	HS-grad	9	Divorced	Tech-support	Unmarried	White	Female	0	0	40	United-States	<=50K
56	Local-government	216851	Bachelors	13	Married-civ-spouse	Tech-support	Husband	White	Male	0	0	40	United-States	>50K
19	Private	166294	HS-grad	9	Never-married	Craft-repair	Own-child	White	Male	0	0	40	United-States	<=50K
54	?	180211	Some-college	10	Married-civ-spouse	?	Husband	Asian-Pac Islander	Male	0	0	60	South	>50K
39	Private	367260	HS-grad	9	Divorced	Exe-manag-sell-relate	Not-in-family	White	Male	0	0	80	United-States	<=50K
49	Private	193366	HS-grad	9	Married-civ-spouse	Craft-repair	Husband	White	Male	0	0	40	United-States	<=50K
23	Local-government	190709	Assoc-accel	12	Never-married	Protective-serv	Not-in-family	White	Male	0	0	52	United-States	<=50K
20	Private	266015	Some-college	10	Never-married	Sales	Own-child	Black	Male	0	0	44	United-States	<=50K
45	Private	386940	Bachelors	13	Divorced	Exe-manag-sell-relate	Own-child	White	Male	0	1408	40	United-States	<=50K
30	Federal-gov	59951	Some-college	10	Married-civ-spouse	Adm-clerical	Own-child	White	Male	0	0	40	United-States	<=50K
22	State-government	311512	Some-college	10	Married-civ-spouse	Other-service	Husband	Black	Male	0	0	15	United-States	<=50K
48	Private	242406	11th	7	Never-married	Machine-guardian	Unmarried	White	Male	0	0	40	Puerto-Rico	<=50K
21	Private	197200	Some-college	10	Never-married	Machine-guardian	Own-child	White	Male	0	0	40	United-States	<=50K
39	Private	544091	HS-grad	9	Married-civ-spouse	Adm-clerical	Wife	White	Female	0	0	25	United-States	<=50K
40	Private	84154	Some-college	10	Married-civ-spouse	Sales	Husband	White	Male	0	0	38	?	>50K
40	Self-employed	365877	Assoc-accel	13	Married-civ-spouse	Prof-specialty	Husband	White	Male	0	0	40	United-States	<=50K

adult - Excel

File Home Insert Draw Page Layout Formulas Data Review View Help Tell me what you want to do Asmita Wagh All Share

Clipboard Font Alignment Number Styles Cells Editing

A1 : fx age

41	48	Self-emp-	265477	Assoc-acc	12	Married-c	Prof-spec Husband	White	Male	K	0	0	40	United-St <=50K
42	31	Private	507875	9th	5	Married-c	Machine- Husband	White	Male	L	0	0	43	United-St <=50K
43	53	Self-emp-	88506	Bachelors	13	Married-c	Prof-spec Husband	White	Male	M	0	0	40	United-St <=50K
44	24	Private	172987	Bachelors	13	Married-c	Tech-supj Husband	White	Male	N	0	0	50	United-St <=50K
45	49	Private	94638	HS-grad	9	Separatec	Adm-cleri Unmarrie	White	Female	O	0	0	40	United-St <=50K
46	25	Private	289980	HS-grad	9	Never-ma	Handlers- Not-in-far	White	Male	P	0	0	35	United-St <=50K
47	57	Federal-g	337895	Bachelors	13	Married-c	Prof-spec Husband	Black	Male	Q	0	0	40	United-St >50K
48	53	Private	144361	HS-grad	9	Married-c	Machine- Husband	White	Male	R	0	0	38	United-St <=50K
49	44	Private	128354	Masters	14	Divorced	Exec-man Unmarrie	White	Female	S	0	0	40	United-St <=50K
50	41	State-gov	101603	Assoc-voc	11	Married-c	Craft-rep; Husband	White	Male		0	0	40	United-St <=50K
51	29	Private	271466	Assoc-voc	11	Never-ma	Prof-spec Not-in-far	White	Male		0	0	43	United-St <=50K
52	25	Private	32275	Some-coll	10	Married-c	Exec-man Wife	Other	Female		0	0	40	United-St <=50K
53	18	Private	226956	HS-grad	9	Never-ma	Other-ser Own-chilc	White	Female		0	0	30	? <=50K
54	47	Private	51835	Prof-scho	15	Married-c	Prof-spec Wife	White	Female		0	1902	60	Honduras >50K
55	50	Federal-g	251585	Bachelors	13	Divorced	Exec-man Not-in-far	White	Male		0	0	55	United-St >50K
56	47	Self-emp-	109832	HS-grad	9	Divorced	Exec-man Not-in-far	White	Male		0	0	60	United-St <=50K
57	43	Private	237993	Some-coll	10	Married-c	Tech-supj Husband	White	Male		0	0	40	United-St >50K
58	46	Private	216666	5th-6th	3	Married-c	Machine- Husband	White	Male		0	0	40	Mexico <=50K
59	35	Private	56352	Assoc-voc	11	Married-c	Other-ser Husband	White	Male		0	0	40	Puerto-Ri <=50K
60	41	Private	147372	HS-grad	9	Married-c	Adm-cleri Husband	White	Male		0	0	48	United-St <=50K
61	30	Private	188146	HS-grad	9	Married-c	Machine- Husband	White	Male		5013	0	40	United-St <=50K
62	30	Private	59496	Bachelors	13	Married-c	Sales Husband	White	Male		2407	0	40	United-St <=50K
63	32	?	293936	7th-8th	4	Married-s ?	Not-in-far	White	Male		0	0	40	? <=50K
64	48	Private	149640	HS-grad	9	Married-c	Transport Husband	White	Male		0	0	40	United-St <=50K
65	42	Private	116632	Doctorate	16	Married-c	Prof-spec Husband	White	Male		0	0	45	United-St >50K
66	29	Private	105598	Some-coll	10	Divorced	Tech-supj Not-in-far	White	Male		0	0	58	United-St <=50K
67	36	Private	155537	HS-grad	9	Married-c	Craft-rep Husband	White	Male		0	0	40	United-St <=50K
68	28	Private	183175	Some-coll	10	Divorced	Adm-cleri Not-in-far	White	Female		0	0	40	United-St <=50K
69	53	Private	169846	HS-grad	9	Married-c	Adm-cleri Wife	White	Female		0	0	40	United-St >50K
70	49	Self-emp-	191681	Some-coll	10	Married-c	Exec-man Husband	White	Male		0	0	50	United-St <=50K
71	25	?	200681	Some-coll	10	Never-ma ?	Own-chilc	White	Male		0	0	40	United-St <=50K
72	19	Private	101509	Some-coll	10	Never-ma	Prof-spec Own-chilc	White	Male		0	0	32	United-St <=50K
73	31	Private	309974	Bachelors	13	Separatec	Sales Own-chilc	Black	Female		0	0	40	United-St <=50K
74	29	Self-emp-	162298	Bachelors	13	Married-c	Sales Husband	White	Male		0	0	70	United-St >50K
75	23	Private	211678	Some-coll	10	Never-ma	Machine- Not-in-far	White	Male		0	0	40	United-St <=50K
76	79	Private	124744	Some-coll	10	Married-c	Prof-spec Other-rel	White	Male		0	0	20	United-St <=50K
77	27	Private	213921	HS-grad	9	Never-ma	Other-ser Own-chilc	White	Male		0	0	40	Mexico <=50K
78	40	Private	32214	Assoc-acc	12	Married-c	Adm-cleri Husband	White	Male		0	0	40	United-St <=50K
79	67	?	212759	10th	6	Married-c ?	Husband	White	Male		0	0	2	United-St <=50K
80	18	Private	309634	11th	7	Never-ma	Other-ser Own-chilc	White	Female		0	0	22	United-St <=50K
81	31	Local-gov	125027	7th-8th	4	Married-c	Farming-f Husband	White	Male		0	0	40	United-St <=50K

adult - Excel

File Home Insert Draw Page Layout Formulas Data Review View Help Tell me what you want to do Asmita Wagh All Share

Clipboard Font Alignment Number Styles Cells Editing

A1 : fx age

60	41	Private	147372	HS-grad	9	Married-c	Adm-cleri Husband	White	Male	K	0	0	48	United-St <=50K
61	30	Private	188146	HS-grad	9	Married-c	Machine- Husband	White	Male	L	5013	0	40	United-St <=50K
62	30	Private	59496	Bachelors	13	Married-c	Sales Husband	White	Male	M	2407	0	40	United-St <=50K
63	32	?	293936	7th-8th	4	Married-s ?	Not-in-far	White	Male	N	0	0	40	? <=50K
64	48	Private	149640	HS-grad	9	Married-c	Transport Husband	White	Male	O	0	0	40	United-St <=50K
65	42	Private	116632	Doctorate	16	Married-c	Prof-spec Husband	White	Male	P	0	0	45	United-St >50K
66	29	Private	105598	Some-coll	10	Divorced	Tech-supj Not-in-far	White	Male	Q	0	0	58	United-St <=50K
67	36	Private	155537	HS-grad	9	Married-c	Craft-rep Husband	White	Male	R	0	0	40	United-St <=50K
68	28	Private	183175	Some-coll	10	Divorced	Adm-cleri Not-in-far	White	Female	S	0	0	40	United-St <=50K
69	53	Private	169846	HS-grad	9	Married-c	Adm-cleri Wife	White	Female		0	0	40	United-St >50K
70	49	Self-emp-	191681	Some-coll	10	Married-c	Exec-man Husband	White	Male		0	0	50	United-St <=50K
71	25	?	200681	Some-coll	10	Never-ma ?	Own-chilc	White	Male		0	0	40	United-St <=50K
72	19	Private	101509	Some-coll	10	Never-ma	Prof-spec Own-chilc	White	Male		0	0	32	United-St <=50K
73	31	Private	309974	Bachelors	13	Separatec	Sales Own-chilc	Black	Female		0	0	40	United-St <=50K
74	29	Self-emp-	162298	Bachelors	13	Married-c	Sales Husband	White	Male		0	0	70	United-St >50K
75	23	Private	211678	Some-coll	10	Never-ma	Machine- Not-in-far	White	Male		0	0	40	United-St <=50K
76	79	Private	124744	Some-coll	10	Married-c	Prof-spec Other-rel	White	Male		0	0	20	United-St <=50K
77	27	Private	213921	HS-grad	9	Never-ma	Other-ser Own-chilc	White	Male		0	0	40	Mexico <=50K
78	40	Private	32214	Assoc-acc	12	Married-c	Adm-cleri Husband	White	Male		0	0	40	United-St <=50K
79	67	?	212759	10th	6	Married-c ?	Husband	White	Male		0	0	2	United-St <=50K
80	18	Private	309634	11th	7	Never-ma	Other-ser Own-chilc	White	Female		0	0	22	United-St <=50K
81	31	Local-gov	125027	7th-8th	4	Married-c	Farming-f Husband	White	Male		0	0	40	United-St <=50K

adult - Excel

Clipboard

Font

Alignment

Number

Styles

Cells

Editing

AutoSum

Fill

Sort & Find & Filter

Clear

File Home Insert Draw Page Layout Formulas Data Review View Help Tell me what you want to do

A1

age

18 Private 309634 11th 7 Never-ma Other-ser Own-chil White Female 0 0 22 United-St <=50K  
31 Local-gov 125927 7th-8th 4 Married-c Farming-f Husband White Male 0 0 40 United-St <=50K  
18 Private 446839 HS-grad 9 Never-ma Sales Not-in-far White Male 0 0 30 United-St <=50K  
52 Private 276515 Bachelors 13 Married-c Other-ser Husband White Male 0 0 40 Cuba <=50K  
46 Private 51618 HS-grad 9 Married-c Other-ser Wife White Female 0 0 40 United-St <=50K  
59 Private 159937 HS-grad 9 Married-c Sales Husband White Male 0 0 48 United-St <=50K  
44 Private 343591 HS-grad 9 Divorced Craft-rep: Not-in-far White Female 14344 0 40 United-St >50K  
53 Private 346253 HS-grad 9 Divorced Sales Own-chik White Female 0 0 35 United-St <=50K  
49 Local-gov 268234 HS-grad 9 Married-c Protective Husband White Male 0 0 40 United-St >50K  
33 Private 202051 Masters 14 Married-c Prof-spec Husband White Male 0 0 50 United-St <=50K  
30 Private 54334 9th 5 Never-ma Sales Not-in-far White Male 0 0 40 United-St <=50K  
43 Federal-g 410867 Doctorate 16 Never-ma Prof-spec Not-in-far White Female 0 0 50 United-St >50K  
57 Private 249977 Assoc-voc 11 Married-c Prof-spec Husband White Male 0 0 40 United-St <=50K  
37 Private 286730 Some-coll 10 Divorced Craft-rep: Unmarrie White Female 0 0 40 United-St <=50K  
28 Private 212563 Some-coll 10 Divorced Machine- Unmarrie Black Female 0 0 25 United-St <=50K  
30 Private 117747 HS-grad 9 Married-c Sales Wife Asian-Pac Female 0 1573 35 ? <=50K  
34 Local-gov 226296 Bachelors 13 Married-c Protective Husband White Male 0 0 40 United-St >50K  
29 Local-gov 115585 Some-coll 10 Never-ma Handlers- Not-in-far White Male 0 0 50 United-St <=50K  
48 Self-emp- 191277 Doctorate 16 Married-c Prof-spec Husband White Male 0 1902 60 United-St >50K  
99 37 Private 202683 Some-coll 10 Married-c Sales Husband White Male 0 0 48 United-St >50K  
48 Private 171095 Assoc-acc 12 Divorced Exec-man Unmarrie White Female 0 0 40 England <=50K  
22 End-of-cou 204000 HS-grad 9 Never-ma Other-ser Own-chil Black Male 0 0 40 United-St <=50K  
adult

Clipboard

Font

Alignment

Number

Styles

Cells

Editing

AutoSum

Fill

Sort & Find & Filter

Clear

File Home Insert Draw Page Layout Formulas Data Review View Help Tell me what you want to do

A1

age

48 Private 171095 Assoc-acc 12 Divorced Exec-man Unmarrie White Female 0 0 40 England <=50K  
32 Federal-g 249409 HS-grad 9 Never-ma Other-ser Own-chik Black Male 0 0 40 United-St <=50K  
76 Private 124191 Masters 14 Married-c Exec-man Husband White Male 0 0 40 United-St >50K  
44 Private 198282 Bachelors 13 Married-c Exec-man Husband White Male 15024 0 60 United-St >50K  
47 Self-emp- 149116 Masters 14 Never-ma Prof-spec Not-in-far White Female 0 0 50 United-St <=50K  
20 Private 188300 Some-coll 10 Never-ma Tech-supj Own-chik White Female 0 0 40 United-St <=50K  
29 Private 103432 HS-grad 9 Never-ma Craft-rep: Not-in-far White Male 0 0 40 United-St <=50K  
32 Self-emp- 317660 HS-grad 9 Married-c Craft-rep: Husband White Male 7688 0 40 United-St >50K  
17 ? 304873 10th 6 Never-ma ? Own-chik White Female 34095 0 32 United-St <=50K  
30 Private 194901 11th 7 Never-ma Handlers- Own-chik White Male 0 0 40 United-St <=50K  
31 Local-gov 189265 HS-grad 9 Never-ma Adm-cleri Not-in-far White Female 0 0 40 United-St <=50K  
42 Private 124692 HS-grad 9 Married-c Handlers- Husband White Male 0 0 40 United-St <=50K  
24 Private 432376 Bachelors 13 Never-ma Sales Other-rel. White Male 0 0 40 United-St <=50K  
38 Private 65324 Prof-scho 15 Married-c Prof-spec Husband White Male 0 0 40 United-St >50K  
56 Self-emp- 335605 HS-grad 9 Married-c Other-ser Husband White Male 0 1887 50 Canada >50K  
28 Private 377869 Some-coll 10 Married-c Sales Wife White Female 4064 0 25 United-St <=50K  
36 Private 102864 HS-grad 9 Never-ma Machine- Own-chik White Female 0 0 40 United-St <=50K  
53 Private 95647 9th 5 Married-c Handlers- Husband White Male 0 0 50 United-St <=50K  
56 Self-emp- 303090 Some-coll 10 Married-c Sales Husband White Male 0 0 50 United-St <=50K  
49 Local-gov 197371 Assoc-voc 11 Married-c Craft-rep: Husband Black Male 0 0 40 United-St >50K  
55 Private 247552 Some-coll 10 Married-c Sales Husband White Male 0 0 56 United-St <=50K  
22 Director 103622 HS-grad 9 Never-ma Craft-rep: Not-in-far White Male 0 0 40 United-St <=50K  
adult

adult - Excel

Clipboard

Font

Alignment

Number

Styles

Cells

Editing

AutoSum

Fill

Sort & Find & Filter

Clear

File Home Insert Draw Page Layout Formulas Data Review View Help Tell me what you want to do

A1

age

48 Private 171095 Assoc-acc 12 Divorced Exec-man Unmarrie White Female 0 0 40 England <=50K  
32 Federal-g 249409 HS-grad 9 Never-ma Other-ser Own-chik Black Male 0 0 40 United-St <=50K  
76 Private 124191 Masters 14 Married-c Exec-man Husband White Male 0 0 40 United-St >50K  
44 Private 198282 Bachelors 13 Married-c Exec-man Husband White Male 15024 0 60 United-St >50K  
47 Self-emp- 149116 Masters 14 Never-ma Prof-spec Not-in-far White Female 0 0 50 United-St <=50K  
20 Private 188300 Some-coll 10 Never-ma Tech-supj Own-chik White Female 0 0 40 United-St <=50K  
29 Private 103432 HS-grad 9 Never-ma Craft-rep: Not-in-far White Male 0 0 40 United-St <=50K  
32 Self-emp- 317660 HS-grad 9 Married-c Craft-rep: Husband White Male 7688 0 40 United-St >50K  
17 ? 304873 10th 6 Never-ma ? Own-chik White Female 34095 0 32 United-St <=50K  
30 Private 194901 11th 7 Never-ma Handlers- Own-chik White Male 0 0 40 United-St <=50K  
110 31 Local-gov 189265 HS-grad 9 Never-ma Adm-cleri Not-in-far White Female 0 0 40 United-St <=50K  
42 Private 124692 HS-grad 9 Married-c Handlers- Husband White Male 0 0 40 United-St <=50K  
24 Private 432376 Bachelors 13 Never-ma Sales Other-rel. White Male 0 0 40 United-St <=50K  
38 Private 65324 Prof-scho 15 Married-c Prof-spec Husband White Male 0 0 40 United-St >50K  
56 Self-emp- 335605 HS-grad 9 Married-c Other-ser Husband White Male 0 1887 50 Canada >50K  
28 Private 377869 Some-coll 10 Married-c Sales Wife White Female 4064 0 25 United-St <=50K  
36 Private 102864 HS-grad 9 Never-ma Machine- Own-chik White Female 0 0 40 United-St <=50K  
53 Private 95647 9th 5 Married-c Handlers- Husband White Male 0 0 50 United-St <=50K  
56 Self-emp- 303090 Some-coll 10 Married-c Sales Husband White Male 0 0 50 United-St <=50K  
49 Local-gov 197371 Assoc-voc 11 Married-c Craft-rep: Husband Black Male 0 0 40 United-St >50K  
55 Private 247552 Some-coll 10 Married-c Sales Husband White Male 0 0 56 United-St <=50K  
22 Director 103622 HS-grad 9 Never-ma Craft-rep: Not-in-far White Male 0 0 40 United-St <=50K  
adult

Clipboard

Font

Alignment

Number

Styles

Cells

Editing

AutoSum

Fill

Sort & Find & Filter

Clear

File Home Insert Draw Page Layout Formulas Data Review View Help Tell me what you want to do

A1

age

adult - Excel

The screenshot displays the first 141 rows of the adult dataset in Microsoft Excel. The columns are labeled A through S. The data includes various demographic and socioeconomic information. The Excel ribbon is visible at the top, and the formula bar shows 'age'.

ID	Education	Age	Marital Status	Work Class	Race	Gender	Capital Gain	Capital Loss	Hours per Week	Country	
121	22 Private	102632 HS-grad	9 Never-ma	Craft-rep; Not-in-far	White	Male	0	0	41	United-States <=50K	
122	21 Private	199915 Some-col	10 Never-ma	Other-ser Own-chilk	White	Female	0	0	40	United-States <=50K	
123	40 Private	118853 Bachelors	13 Married-c	Exec-man Husband	White	Male	0	0	60	United-States <=50K	
124	30 Private	77143 Bachelors	13 Never-ma	Exec-man Own-chilk	Black	Male	0	0	40	Germany <=50K	
125	29 State-gov	267989 Bachelors	13 Married-c	Prof-spec Husband	White	Male	0	0	50	United-States >50K	
126	19 Private	301605 Some-col	10 Never-ma	Other-ser Own-chilk	Black	Male	0	0	35	United-States <=50K	
127	47 Private	287828 Bachelors	13 Married-c	Exec-man Wife	White	Female	0	0	40	United-States >50K	
128	20 Private	111697 Some-col	10 Never-ma	Adm-cleri Own-chilk	White	Female	0	1719	28	United-States <=50K	
129	31 Private	114937 Assoc-acc	12 Married-c	Adm-cleri Husband	White	Male	0	0	40	United-States >50K	
130	35 ?	129305 HS-grad	9 Married-c ?	Husband	White	Male	0	0	40	United-States <=50K	
131	39 Private	365739 Some-col	10 Divorced	Craft-rep; Not-in-far	White	Male	0	0	40	United-States <=50K	
132	28 Private	69621 Assoc-acc	12 Never-ma	Sales	Not-in-far	White	Female	0	0	60	United-States <=50K
133	24 Private	43323 HS-grad	9 Never-ma	Other-ser	Not-in-far	White	Female	0	1762	40	United-States <=50K
134	38 Self-emp-	120985 HS-grad	9 Married-c	Craft-rep; Husband	White	Male	4386	0	35	United-States <=50K	
135	37 Private	254202 Bachelors	13 Married-c	Sales	Husband	White	Male	0	0	50	United-States <=50K
136	46 Private	146195 Assoc-acc	12 Divorced	Tech-supj	Not-in-far	Black	Female	0	0	36	United-States <=50K
137	38 Federal-g	125933 Masters	14 Married-c	Prof-spec Husband	White	Male	0	0	40	Iran >50K	
138	43 Self-emp-	56920 HS-grad	9 Married-c	Craft-rep; Husband	White	Male	0	0	60	United-States <=50K	
139	27 Private	163127 Assoc-voc	11 Married-c	Adm-cleri Wife	White	Female	0	0	35	United-States <=50K	
140	20 Private	34310 Some-col	10 Never-ma	Sales	Own-chilk	White	Male	0	0	20	United-States <=50K
141	49 Private	81973 Some-col	10 Married-c	Craft-rep; Husband	Asian-Pac	Male	0	0	40	United-States >50K	
142	61 Self-emp-	66614 HS-grad	9 Married-c	Craft-rep; Husband	White	Male	0	0	40	United-States <=50K	
143	27 Private	232782 Some-col	10 Never-ma	Sales	Own-chilk	White	Female	0	0	40	United-States <=50K
144	19 Private	316868 Some-col	10 Never-ma	Other-ser	Own-chilk	White	Male	0	0	30	Mexico <=50K
145	45 Private	196584 Assoc-voc	11 Never-ma	Prof-spec	Not-in-far	White	Female	0	1564	40	United-States >50K
146	70 Private	105376 Some-col	10 Never-ma	Tech-supj	Other-rel	White	Male	0	0	40	United-States <=50K
147	31 Private	185814 HS-grad	9 Never-ma	Transport Unmarrie	Black	Female	0	0	30	United-States <=50K	
148	22 Private	175374 Some-col	10 Married-c	Other-ser	Husband	White	Male	0	0	24	United-States <=50K
149	36 Private	108293 HS-grad	9 Widowed	Other-ser	Unmarrie	White	Female	0	0	24	United-States <=50K
150	64 Private	181232 11th	7 Married-c	Craft-rep; Husband	White	Male	0	2179	40	United-States <=50K	
151	43 ?	174662 Some-col	10 Divorced ?	Not-in-far	White	Female	0	0	40	United-States <=50K	
152	47 Local-gov	186009 Some-col	10 Divorced	Adm-cleri	Unmarrie	White	Female	0	0	38	Mexico <=50K
153	34 Private	198183 HS-grad	9 Never-ma	Adm-cleri	Not-in-far	White	Female	0	0	40	United-States <=50K
154	33 Private	163003 Bachelors	13 Never-ma	Exec-man	Other-rel	Asian-Pac	Female	0	0	40	Philippine <=50K
155	21 Private	296158 HS-grad	9 Never-ma	Craft-rep;	Own-chilk	White	Male	0	0	35	United-States <=50K
156	52 ?	252903 HS-grad	9 Divorced ?	Not-in-far	White	Male	0	0	45	United-States >50K	
157	48 Private	187715 HS-grad	9 Married-c	Craft-rep; Husband	White	Male	0	0	46	United-States <=50K	
158	23 Private	214542 Bachelors	13 Never-ma	Handlers-	Not-in-far	White	Male	0	0	40	United-States <=50K
159	71 Self-emp-	494223 Some-col	10 Separates	Sales	Unmarrie	Black	Male	0	1816	2	United-States <=50K
160	29 Private	191535 HS-grad	9 Divorced	Craft-rep; Not-in-far	White	Male	0	0	60	United-States <=50K	
161	42 Private	728456 Bachelors	12 Separates	Other-ser	Other-rel	Black	Male	0	0	50	United-States <=50K

adult - Excel

The screenshot displays the last 161 rows of the adult dataset in Microsoft Excel. The columns are labeled A through S. The data continues from row 140 to 300, showing a mix of marital statuses, work classes, and education levels. The Excel ribbon is visible at the top, and the formula bar shows 'adult'.

ID	Education	Age	Marital Status	Work Class	Race	Gender	Capital Gain	Capital Loss	Hours per Week	Country	
140	20 Private	34310 Some-col	10 Never-ma	Sales	Own-chilk	White	Male	0	0	20	United-States <=50K
141	49 Private	81973 Some-col	10 Married-c	Craft-rep; Husband	Asian-Pac	Male	0	0	40	United-States >50K	
142	61 Self-emp-	66614 HS-grad	9 Married-c	Craft-rep; Husband	White	Male	0	0	40	United-States <=50K	
143	27 Private	232782 Some-col	10 Never-ma	Sales	Own-chilk	White	Female	0	0	40	United-States <=50K
144	19 Private	316868 Some-col	10 Never-ma	Other-ser	Own-chilk	White	Male	0	0	30	Mexico <=50K
145	45 Private	196584 Assoc-voc	11 Never-ma	Prof-spec	Not-in-far	White	Female	0	1564	40	United-States >50K
146	70 Private	105376 Some-col	10 Never-ma	Tech-supj	Other-rel	White	Male	0	0	40	United-States <=50K
147	31 Private	185814 HS-grad	9 Never-ma	Transport	Unmarrie	Black	Female	0	0	30	United-States <=50K
148	22 Private	175374 Some-col	10 Married-c	Other-ser	Husband	White	Male	0	0	24	United-States <=50K
149	36 Private	108293 HS-grad	9 Widowed	Other-ser	Unmarrie	White	Female	0	0	24	United-States <=50K
150	64 Private	181232 11th	7 Married-c	Craft-rep; Husband	White	Male	0	2179	40	United-States <=50K	
151	43 ?	174662 Some-col	10 Divorced ?	Not-in-far	White	Female	0	0	40	United-States <=50K	
152	47 Local-gov	186009 Some-col	10 Divorced	Adm-cleri	Unmarrie	White	Female	0	0	38	Mexico <=50K
153	34 Private	198183 HS-grad	9 Never-ma	Adm-cleri	Not-in-far	White	Female	0	0	40	United-States <=50K
154	33 Private	163003 Bachelors	13 Never-ma	Exec-man	Other-rel	Asian-Pac	Female	0	0	40	Philippine <=50K
155	21 Private	296158 HS-grad	9 Never-ma	Craft-rep;	Own-chilk	White	Male	0	0	35	United-States <=50K
156	52 ?	252903 HS-grad	9 Divorced ?	Not-in-far	White	Male	0	0	45	United-States >50K	
157	48 Private	187715 HS-grad	9 Married-c	Craft-rep; Husband	White	Male	0	0	46	United-States <=50K	
158	23 Private	214542 Bachelors	13 Never-ma	Handlers-	Not-in-far	White	Male	0	0	40	United-States <=50K
159	71 Self-emp-	494223 Some-col	10 Separates	Sales	Unmarrie	Black	Male	0	1816	2	United-States <=50K
160	29 Private	191535 HS-grad	9 Divorced	Craft-rep; Not-in-far	White	Male	0	0	60	United-States <=50K	
161	42 Private	728456 Bachelors	12 Separates	Other-ser	Other-rel	Black	Male	0	0	50	United-States <=50K

## Theory:

The data was extracted from 1994 Census bureau database. Please view the adult.csv file for the dataset, this dataset is what is being used in the entire study of ours.

income <- read.csv ('adult.csv', na.strings = c ('', '?')) ..... To assign the csv data into variable income, the na. strings are basically used to specify the content that is unknown in the dataset or left blank.

supply (income, function(x) sum(is.na(x))) ..... basically, counts the data that is unavailable.

After finding it out, I realized that numbers are a bit annoying to comprehend and select the reliable points, henceforth we went ahead and plotted the graphs of it to see how they are.

```
library(Amelia)
missmap(income, main = "Missing values vs observed")
table (complete.cases (income))`
```

This basically allows us to see the plot of missing and observed data.

Check Rplots.pdf for the plot

Since the dataset has already been cleaned running it up again, will not show any false values.

```
library(Amelia)
Loading required package: foreign
##
## Amelia II: Multiple Imputation
## (Version 1.6.4, built: 2012-12-17)
## Copyright (C) 2005-2020 James Honaker, Gary King and Matthew Blackwell
## Refer to http://gking.harvard.edu/amelia/ for more information
##
> missmap(income, main = "Missing values vs observed")
> table (complete.cases (income))

TRUE
32561
```

So, we go ahead and make a boxplot of everything with respect to income level.

```
library(gridExtra)
p1 <- ggplot(aes(x=income, y=age), data = income) + geom_boxplot() +
  ggtitle('Age vs. Income Level')
p2 <- ggplot(aes(x=income, y=education.num), data = income) + geom_boxplot() +
  ggtitle('Years of Education vs. Income Level')
str(income)
p3 <- ggplot(aes(x=income, y=house.per.week), data = income) + geom_boxplot() + ggtitle('Hours Per week vs. Income Level')
p4 <- ggplot(aes(x=income, y=capital.gain), data=income) + geom_boxplot() +
  ggtitle('Capital Gain vs. Income Level')
p5 <- ggplot(aes(x=income, y=capital.loss), data=income) + geom_boxplot() +
  ggtitle('Capital Loss vs. Income Level')
p6 <- ggplot(aes(x=income, y=fnlwgt), data=income) + geom_boxplot() +
  ggtitle('Final Weight vs. Income Level')
grid.arrange(p1, p2, p3, p4, p5, p6, ncol=3)
income$fnlwgt <- NULL
```

“Age”, “Years of education” and “hours per week” all show significant variations with income level. Therefore, they are kept for the regression analysis. “Final Weight” does not show any variation with income level, therefore, it has been excluded from the analysis. It’s hard to see whether “Capital gain” and “Capital loss” have variation with Income level from the above plot, so we shall keep them for now.

```
library(dplyr)
by_workclass <- income %>% group_by(workclass, income) %>% summarise(n=n())
by_education <- income %>% group_by(education, income) %>% summarise(n=n())
by_education$education <- ordered(by_education$education,
                                     levels = c('Preschool', '1st-4th', '5th-6th', '7th-8th', '9th', '10th', '11th', '12th'))
by_marital <- income %>% group_by(marital.status, income) %>% summarise(n=n())
by_occupation <- income %>% group_by(occupation, income) %>% summarise(n=n())
by_relationship <- income %>% group_by(relationship, income) %>% summarise(n=n())
by_race <- income %>% group_by(race, income) %>% summarise(n=n())
by_sex <- income %>% group_by(sex, income) %>% summarise(n=n())
by_country <- income %>% group_by(native.country, income) %>% summarise(n=n())
p7 <- ggplot(aes(x=workclass, y=n, fill=income), data=by_workclass) + geom_bar(stat = 'identity', position = position_dodge())
p8 <- ggplot(aes(x=education, y=n, fill=income), data=by_education) + geom_bar(stat = 'identity', position = position_dodge())
p9 <- ggplot(aes(x=marital.status, y=n, fill=income), data=by_marital) + geom_bar(stat = 'identity', position=position_dodge())
p10 <- ggplot(aes(x=occupation, y=n, fill=income), data=by_occupation) + geom_bar(stat = 'identity', position=position_dodge())
p11 <- ggplot(aes(x=relationship, y=n, fill=income), data=by_relationship) + geom_bar(stat = 'identity', position=position_dodge())
p12 <- ggplot(aes(x=race, y=n, fill=income), data=by_race) + geom_bar(stat = 'identity', position = position_dodge())
p13 <- ggplot(aes(x=sex, y=n, fill=income), data=by_sex) + geom_bar(stat = 'identity', position = position_dodge()) + ggtitle('Sex vs. Income')
p14 <- ggplot(aes(x=native.country, y=n, fill=income), data=by_country) + geom_bar(stat = 'identity', position = position_dodge())
grid.arrange(p7, p8, p9, p10, ncol=2)
grid.arrange(p11,p12,p13, ncol=2)
```

Most of the data is collected from the United States, so variable “native country” does not have effect on our analysis, we shall exclude it from regression model.

And all the other categorial variables seem to have reasonable variation, so they will be kept.

income\$income = as.factor(ifelse(income\$income==income\$income[1],0,1))  
basically if >50k it will be set to 1 else to 0

```
train <- income[1:22793,]
test <- income[22793:32561,]
model <- glm(income ~., family=binomial(link='logit'), data=train)
summary(model)
```

So as to fit the model and view the details.

```
Deviance Residuals:
    Min      Q1      Median      Q3      Max 
-4.9275 -0.5171 -0.1885 -0.0327  3.2694 

Coefficients: (2 not defined because of singularities)
              Estimate Std. Error z value Pr(>|z|)    
(Intercept) -8.325e+00 4.908e-01 -16.963 < 2e-16 ***
age          2.505e-02 1.932e-03 12.961 < 2e-16 ***
workclass Federal-gov 1.100e+00 1.818e-01  6.055 1.41e-09 ***
workclass Local-gov 4.509e-01 1.660e-01  2.716 0.006599 ** 
workclass Never-worked -9.225e+00 6.808e+02 -0.014 0.989189  
workclass Private   6.245e-01 1.481e-01  4.216 2.49e-05 *** 
workclass Self-emp-inc 8.046e-01 1.766e-01  4.557 5.20e-06 *** 
workclass Self-emp-not-inc 1.674e-01 1.623e-01  1.031 0.302392  
workclass State-gov 2.851e-01 1.804e-01  1.580 0.114139  
workclass Without-pay -1.279e+01 4.165e+02 -0.031 0.975497  
education 11th       3.831e-02 2.421e-01  0.158 0.874287  
education 12th       3.946e-01 3.182e-01  1.240 0.214856  
education 1st-4th    -5.194e-01 5.324e-01 -0.976 0.329244  
education 5th-6th    -4.645e-01 3.784e-01 -1.228 0.219579  
education 7th-8th    -4.160e-01 2.683e-01 -1.550 0.121113  
education 9th        -4.346e-01 3.154e-01 -1.378 0.168233  
education Assoc-acdm 1.312e+00 2.054e-01  6.388 1.68e-10 *** 
education Assoc-voc 1.307e+00 1.972e-01  6.627 3.43e-11 *** 
education Bachelors 1.885e+00 1.821e-01 10.351 < 2e-16 *** 
education Doctorate 2.840e+00 2.522e-01 11.261 < 2e-16 *** 
education HS-grad    7.689e-01 1.773e-01  4.336 1.45e-05 *** 
education Masters   2.181e+00 1.947e-01 11.204 < 2e-16 *** 

education Preschool -1.996e+01 2.167e+02 -0.092 0.926598  
education Prof-school 2.614e+00 2.353e-01 11.109 < 2e-16 *** 
education Some-college 1.062e+00 1.802e-01 5.897 3.71e-09 *** 
education.num        NA         NA         NA         NA        
marital.status Married-AF-spouse 2.687e+00 6.768e-01 3.970 7.19e-05 *** 
marital.status Married-civ-spouse 2.246e+00 3.188e-01 7.043 1.88e-12 *** 
marital.status Married-spouse-absent -2.308e-02 2.629e-01 -0.088 0.930043  
marital.status Never-married    -3.864e-01 1.039e-01 -3.718 0.000201 *** 
marital.status Separated       5.553e-02 1.931e-01 0.288 0.773704  
marital.status Widowed       1.271e-01 1.917e-01 0.663 0.507395  
occupation Adm-clerical   1.184e-02 1.163e-01 0.102 0.918946  
occupation Armed-Forces   -1.325e+01 5.210e+02 -0.025 0.979717  
occupation Craft-repair  1.318e-01 9.983e-02 1.320 0.186672  
occupation Exec-managerial 8.385e-01 1.027e-01 8.166 3.19e-16 *** 
occupation Farming-fishing -9.949e-01 1.696e-01 -5.867 4.45e-09 *** 
occupation Handlers-cleaners -6.268e-01 1.732e-01 -3.620 0.000295 *** 
occupation Machine-op-inspct -3.041e-01 1.256e-01 -2.420 0.015510 * 
occupation Other-service   -7.233e-01 1.435e-01 -5.041 4.63e-07 *** 
occupation Priv-house-serv -1.205e+01 1.174e+02 -0.103 0.918289  
occupation Prof-specialty  6.148e-01 1.103e-01 5.573 2.50e-08 *** 
occupation Protective-serv  5.988e-01 1.526e-01 3.924 8.70e-05 *** 
occupation Sales        2.675e-01 1.058e-01 2.527 0.011507 * 
occupation Tech-support   7.188e-01 1.403e-01 5.125 2.98e-07 *** 
occupation Transport-moving  NA         NA         NA         NA        
relationship Not-in-family 5.715e-01 3.148e-01 1.815 0.069470 . 
relationship Other-relative -5.502e-01 2.934e-01 -1.876 0.060718 . 
relationship Own-child     -6.178e-01 3.113e-01 -1.984 0.047209 * 
relationship Unmarried    3.774e-01 3.342e-01 1.129 0.258751  
relationship Wife        1.357e+00 1.213e-01 11.193 < 2e-16 ***
```

```

relationship Wife          1.357e+00  1.213e-01  11.193 < 2e-16 ***
race Asian-Pac-Islander  2.078e-01  2.824e-01   0.736 0.461944
race Black               3.286e-01  2.690e-01   1.222 0.221864
race Other               -5.340e-01 4.306e-01  -1.240 0.214975
race White               4.402e-01  2.565e-01   1.716 0.086104 .
sex Male                 8.320e-01  9.420e-02   8.832 < 2e-16 ***
capital.gain             3.039e-04  1.187e-05  25.598 < 2e-16 ***
capital.loss              6.484e-04  4.445e-05  14.587 < 2e-16 ***
house.per.week            2.919e-02  1.934e-03  15.094 < 2e-16 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for binomial family taken to be 1)

Null deviance: 25044  on 22792  degrees of freedom
Residual deviance: 14581  on 22736  degrees of freedom
AIC: 14695

Number of Fisher Scoring iterations: 14

Analysis of Deviance Table

Model: binomial, link: logit

Response: income

Terms added sequentially (first to last)

```

	Df	Deviance	Resid. Df	Resid. Dev	Pr(>Chi)
NULL		22792	25044		
age	1	1152.1	22791	23892 < 2.2e-16 ***	
workclass	8	558.3	22783	23334 < 2.2e-16 ***	
education	15	2551.2	22768	20783 < 2.2e-16 ***	
education.num	0	0.0	22768	20783	
marital.status	6	3687.4	22762	17096 < 2.2e-16 ***	
occupation	13	552.5	22749	16543 < 2.2e-16 ***	
relationship	5	165.6	22744	16377 < 2.2e-16 ***	
race	4	19.9	22740	16358 0.0005324 ***	
sex	1	105.8	22739	16252 < 2.2e-16 ***	
capital.gain	1	1210.5	22738	15041 < 2.2e-16 ***	
capital.loss	1	227.1	22737	14814 < 2.2e-16 ***	
house.per.week	1	233.4	22736	14581 < 2.2e-16 ***	
---					
Signif. codes:	0 '***'	0.001 '**'	0.01 '*'	0.05 '.'	0.1 ' ' 1
[1]	"Accuracy :	0.851776026205343"			

Which gives an accuracy of upto 85%.

## Code:

```

1 q()
2 income <- read.csv('adult.csv',na.strings = c("','?"))

```

```
3 str(income0
4 str(income)
5 summary(income)
6 sapply(income,function(x) sum(is.na(x)))
7 sapply(income, function(x) length(unique(x)))
8 library(Amelia)
9 missmap(income, main = "Missing values vs observed")
10 table (complete.cases (income))
11 income <- income[complete.cases(income),]
12 library(ggplot2)
13 library(gridExtra)
14 p1 <- ggplot(aes(x=income, y=age), data = income) + geom_boxplot() +
15   ggttitle('Age vs. Income Level')
16 p2 <- ggplot(aes(x=income, y=education.num), data = income) +
17   geom_boxplot() +
18 ggttitle('Years of Education vs. Income Level')
19 str(incomeo)
20 str(income)
21 p3 <- ggplot(aes(x=income, y=hours.per.week), data = income) +
22   geom_boxplot() +
23 p3 <- ggplot(aes(x=income, y=house.per.week), data = income) +
24   geom_boxplot() +
25 p3 <- ggplot(aes(x=income, y=house.per.week), data = income) +
26   geom_boxplot() + ggttitle('Hours Per week vs. Income Level')
27 p4 <- ggplot(aes(x=income, y=capital.gain), data=income) +
28   geom_boxplot()+
29 ggttitle('Capital Gain vs. Income Level')
30 p5 <- ggplot(aes(x=income, y=capital.loss), data=income) +
31   geom_boxplot()+
32 ggttitle('Capital Loss vs. Income Level')
33 p6 <- ggplot(aes(x=income, y=fnlwgt), data=income) + geom_boxplot() +
34   ggttitle('Final Weight vs. Income Level')
35 grid.arrange(p1, p2, p3, p4, p5, p6, ncol=3)
36 income$fnlwgt <- NULL
37 #cuz final weight shows no variation wrt income
38 library(dplyr)
39 by_workclass <- income %>% group_by(workclass, income) %>%
40   summarise(n=n())
41 by_education <- income %>% group_by(education, income) %>%
42   summarise(n=n())
```

```

35 by_education$education <- ordered(by_education$education,
36 levels = c('Preschool', '1st-4th', '5th-6th', '7th-8th', '9th', '10th', '11th',
37 '12th', 'HS-grad', 'Prof-school', 'Assoc-acdm', 'Assoc-voc', 'Some-college',
38 'Bachelors', 'Masters', 'Doctorate'))
39 by_marital <- income %>% group_by(marital.status, income) %>%
40 summarise(n=n())
41 by_occupation <- income %>% group_by(occupation, income) %>%
42 summarise(n=n())
43 by_relationship <- income %>% group_by(relationship, income) %>%
44 summarise(n=n())
45 by_race <- income %>% group_by(race, income) %>% summarise(n=n())
46 by_sex <- income %>% group_by(sex, income) %>% summarise(n=n())
47 by_country <- income %>% group_by(native.country, income) %>%
48 summarise(n=n())
49 p7 <- ggplot(aes(x=workclass, y=n, fill=income), data=by_workclass) +
50 geom_bar(stat = 'identity', position = position_dodge()) + ggtitle('Workclass
51 with Income Level') + theme(axis.text.x = element_text(angle = 45, hjust =
52 1))
53 p8 <- ggplot(aes(x=education, y=n, fill=income), data=by_education) +
54 geom_bar(stat = 'identity', position = position_dodge()) + ggtitle('Education
55 vs. Income Level') + coord_flip()
56 p9 <- ggplot(aes(x=marital.status, y=n, fill=income), data=by_marital) +
57 geom_bar(stat = 'identity', position=position_dodge()) + ggtitle('Marital
58 Status vs. Income Level') + theme(axis.text.x = element_text(angle = 45,
59 hjust = 1))
60 p10 <- ggplot(aes(x=occupation, y=n, fill=income), data=by_occupation) +
61 geom_bar(stat = 'identity', position=position_dodge()) + ggtitle('Occupation
62 vs. Income Level') + coord_flip()
63 p11 <- ggplot(aes(x=relationship, y=n, fill=income), data=by_relationship) +
64 geom_bar(stat = 'identity', position=position_dodge()) +
65 ggtitle('Relationship vs. Income Level') + coord_flip()
66 p12 <- ggplot(aes(x=race, y=n, fill=income), data=by_race) + geom_bar(stat
67 = 'identity', position = position_dodge()) + ggtitle('Race vs. Income Level') +
68 coord_flip()
69 p13 <- ggplot(aes(x=sex, y=n, fill=income), data=by_sex) + geom_bar(stat =
70 'identity', position = position_dodge()) + ggtitle('Sex vs. Income Level')
71 p14 <- ggplot(aes(x=native.country, y=n, fill=income), data=by_country) +
72 geom_bar(stat = 'identity', position = position_dodge()) + ggtitle('Native
73 Country vs. Income Level') + coord_flip()
74 grid.arrange(p7, p8, p9, p10, ncol=2)

```

```
52 #categorical variable exploration and plotting
53 income$native.country <- NULL
54 #for simplification and saving myself from headache I prefer to not conduct
  my study as per the countries
55 income$income =
  as.factor(ifelse(income$income==income$income[1],0,1))
56 #basically if >50k it will be set to 1 else to 0
57 summary(income)
58 str(income)
59 len(income)
60 count(income)
61 train <- income(22793)
62 train <- income[1:22793,]
63 test <- income[22793:32561,]
64 model <- glm(income ~.,family=binomial(link='logit'),data=train)
65 summary(model)
66 anova(model,test="Chisq")
67 fitted.results <- predict(model,newdata=test,type='response')
68 fitted.results <- ifelse(fitted.results > 0.5,1,0)
69 misClasificError <- mean(fitted.results != test$income)
70 print(paste('Accuracy :',1-misClasificError))
71 library(ROCR)
72 install.packages(ROCR)
73 R CMD INSTALL ROCR_1.0-1.tar.gz
74 install.packages(gplot)
75 q()
76 getwd()
77 ls()
78 savehistory(file="SalaryPrediction")
```

## Output:

```
91  ***R  
92  train <- income[1:22793,]  
93  test <- income[22793:32561]  
94  model <- glm(income ~ ., family=binomial(link='logit'), data=train)  
95  summary(model)  
96  ...  
97  So as to fit the model and view the details.  
98  ***R  
99  Deviance Residuals  
100 Min 1Q Median 3Q Max  
101 -4.9275 -0.5171 0.1885 -0.0327 3.2694  
102  
103 Coefficients: (2 not defined because of singularities)  
104 Estimate Std. Error z value Pr(>|z|)  
105 (Intercept) -8.325e+08 4.908e-01 16.963 <2e-16 ***  
106 age 2.505e-02 1.932e-03 12.961 <2e-16 ***  
107 workclass Federal_gov 1.106e+00 1.818e-01 6.055 1.41e-09 ***  
108 workclass Local_gov 4.599e-01 1.669e-01 2.716 0.006599 **  
109 workclass Never_worked 9.335e-01 2.481e-01 3.778 0.000177 ***  
110 workclass Private 6.245e-01 1.481e-01 4.216 2.49e-05 ***  
111 workclass Self_emp_inc 8.046e-01 1.766e-01 4.557 2.50e-06 ***  
112 workclass Self_emp_not_inc 1.674e-01 1.623e-01 1.031 0.302392  
113 workclass State_gov 2.851e-01 1.804e-01 1.589 0.114139  
114 workclass Without_pay 1.739e-01 1.588e-01 1.089 0.277777  
115 education 12th 3.631e-02 2.421e-01 0.158 0.874287  
116 education 13th 3.946e-01 3.182e-01 1.240 0.214856  
117 education 1st_4th -5.194e-01 5.324e-01 -0.976 0.329244  
118 education 5th_6th -4.645e-01 3.784e-01 -1.228 0.219579  
119 education 9th -4.150e-01 2.936e-01 -1.393 0.121213  
120 education Assoc_acdm 4.346e-01 2.159e-01 1.978 0.123233  
121 education Assoc_voc 1.312e+00 2.054e-01 6.388 1.68e-10 ***  
122 education Bachelors 1.307e+00 1.972e-01 6.627 4.34e-11 ***  
123 education Doctorate 1.685e+00 1.821e-01 10.351 <2e-16 ***  
124 education Highスク 2.848e+00 2.522e-01 11.261 <2e-16 ***  
125 education Masters 7.699e+00 2.522e-01 30.336 <2e-16 ***  
126 education Preschool 2.181e+00 1.947e-01 11.204 <2e-16 ***  
127 education Prof_school 1.996e+01 2.167e+02 0.092 0.926598  
128 education Some_college 2.614e+00 2.352e-01 11.109 <2e-16 ***  
129 education扭 NA NA NA NA NA ***  
130 maritalstatus Married_AF_spouse 2.687e+00 6.769e-01 3.955 7.19e-05 ***  
131 maritalstatus Married_civ_spouse 2.246e+00 3.189e-01 7.043 1.88e-12 ***  
132 maritalstatus Married_spouse_absent 2.308e-02 2.629e-01 -0.688 0.530043  
133 maritalstatus Never_married -3.664e-01 1.039e-01 -3.718 0.000201 ***  
134 maritalstatus Separated 5.553e-02 1.931e-01 0.288 0.773704  
135 married扭 NA NA NA NA NA ***  
136 occupation Adm Clerical 1.177e-01 1.588e-01 0.735 0.467575  
137 occupation Armed_Forces 1.184e-02 1.163e-01 0.102 0.919846  
138 occupation Craft_repair 1.325e+01 5.210e+02 0.025 0.979717  
139 occupation Exec_managerial 8.385e-01 1.027e-01 8.166 3.19e-16 ***  
140 occupation Farming_Fishing 9.949e-01 1.699e-01 -5.867 4.45e-09 ***  
141 occupation Handlers_cleaners -6.266e-01 1.732e-01 -3.620 0.000293 ***  
142  
143  
144  
145  
146  
147  
148  
149  
150 race Asian_Pac_Islander 2.078e-01 2.824e-01 0.736 0.461944  
151 race Black 3.286e-01 2.696e-01 1.222 0.221864  
152 race Other 5.340e-01 4.306e-01 1.240 0.214975  
153 sex Male 8.402e-01 2.359e-01 1.104 0.066104 *  
154 sex Female 8.330e-01 2.426e-01 1.182 0.066104 *  
155 capital_gain 3.039e-01 1.187e-05 25.598 <2e-16 ***  
156 capital_loss 6.404e-04 4.445e-05 14.587 <2e-16 ***  
157 house_per.week 2.919e-02 1.934e-03 15.094 <2e-16 ***  
158 ***  
159 Signif. codes: 0 ***** 0.001 *** 0.01 ** 0.05 . 0.1 ***  
160 (Dispersion parameter for binomial family taken to be 1)  
161 Null deviance 25044 on 22792 degrees of freedom  
162 Residual deviance 14581 on 22736 degrees of freedom  
163 AIC: 14695  
164 Number of Fisher Scoring iterations: 14  
165 Analysis of Deviance Table  
166 Model: binomial, link: logit  
167 Response: income  
168 Terms added sequentially (first to last)  
169  
170 capital.loss 1 227.1 22737 14814 <2.e-16 ***  
171 house_per.week 1 233.4 22736 14581 <2.e-16 ***  
172  
173 Signif. codes: 0 ***** 0.001 *** 0.01 ** 0.05 . 0.1 ***  
174 [1] "Accuracy : 0.851776026205343"  
175  
176 Which gives an accuracy of upto 85%.  
177  
178 # Conclusions:  
179 Interpreting the results of the logistic regression model:  
180  
181 "Age", "Hours per week", "sex", "capital gain" and "capital loss" are the most statistically significant variables. Their lowest p-values suggesting a strong association with the probability of wage>50K from the data.  
182 OK for the results.  
183 "Workclass", "education", "marital status", "occupation" and "relationship" are all across the table.  
184 e. so cannot be eliminated from the model.  
185 Which basically implies the role of each in earning and salary.  
186 "Race" category is not statistically significant and can be eliminated from the model.  
187  
188 *This project is mainly a subset from susanl2010 's Census Income!  
189  
190 It has been adapted for simplification and understanding at beginner level and for our need of mini project  
191  
192 [END]
```

Sun, 05 Apr 19:05 92% akuma@shiro

## Conclusion:

Interpreting the results of the logistic regression model:

- “Age”, “Hours per week”, “sex”, “capital gain” and “capital loss” are the most statistically significant variables. Their lowest p-values suggesting a strong association with the probability of wage>50K from the data.
- “Workclass”, “education”, “marital status”, “occupation” and “relationship” are all across the table. so, cannot be eliminated from the model.
- Which basically implies the role of each in earning and salary.
- “Race” category is not statistically significant and can be eliminated from the model.

