

Loading and Handling Data in R

Challenges of Analytical Data Processing

Expression, Variables and Functions

Missing Values Treatment in R

Using the 'as' Operator to Change the Structure of Data,

Vectors

Matrices

Factors

List

Aggregating and Group Processing of a Variable,

Simple Analysis Using R,

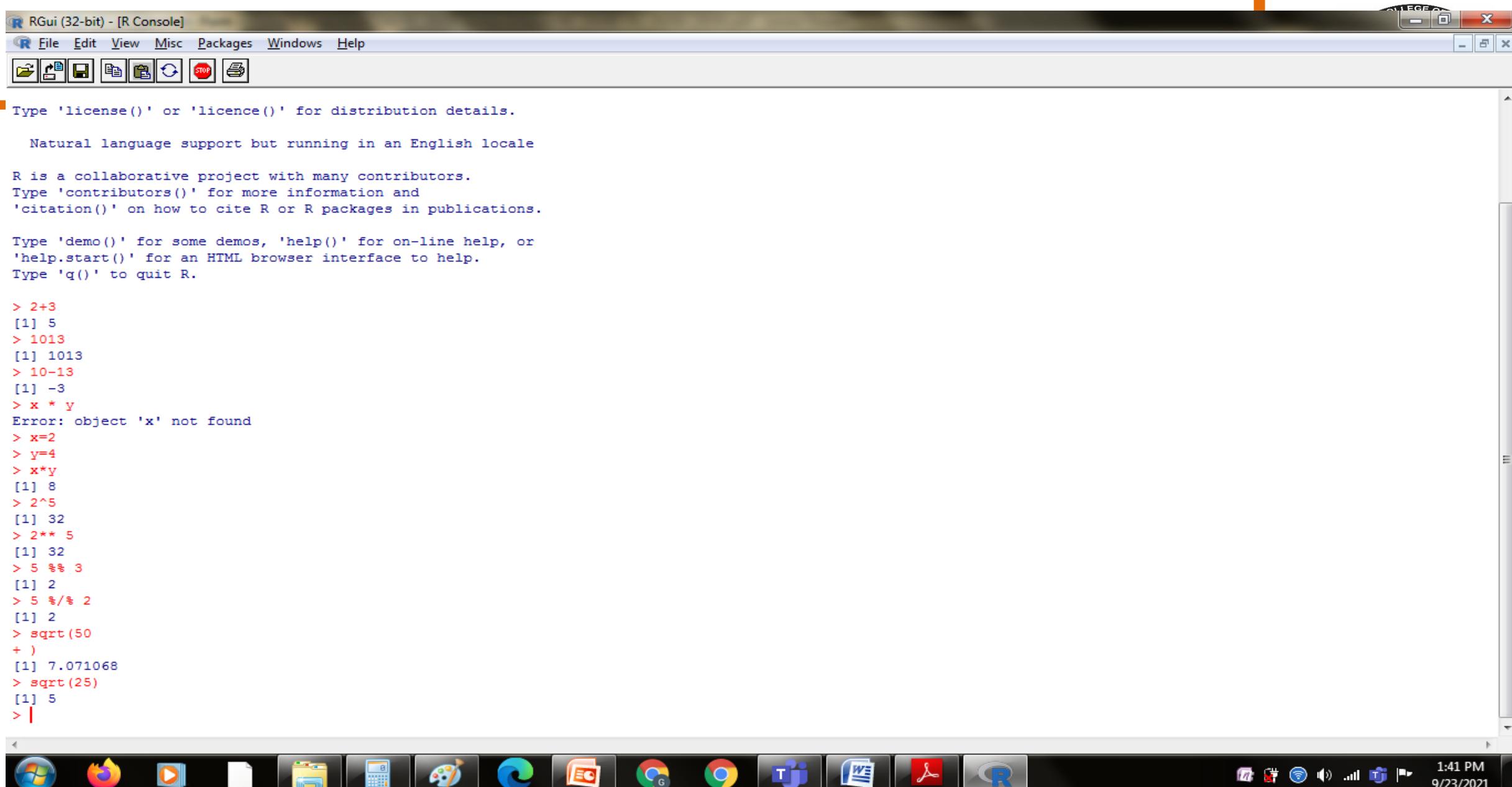
Methods for Reading Data,

Comparison of R GUIs for Data Input

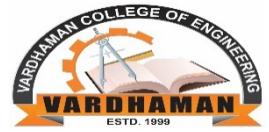
Expression, variables and functions

q Expressions

Operation	Operator	Description	Example
Addition	$x + y$	y added to x	<code>> 4 + 8 [1] 12</code>
Subtraction	$x - y$	y subtracted from x	<code>> 10 - 3 [1] 7</code>
Multiplication	$x * y$	x multiplied by y	<code>> 7 * 8 [1] 56</code>
Division	x / y	x divided by y	<code>< 8/3 [1] 2.666667</code>
Exponentiation	$x ^ y$ $x ** y$	x raised to the power y	<code>> 2 ^ 5 [1] 32</code> Or <code>> 2 ** 5 [1] 32</code>
Modulus	$x \% y$	Remainder of (x divided by y)	<code>> 5 %% 3 [1] 2</code>
Integer Division	$x \% / \% y$	x divided by y but rounded down	<code>> 5 \% / % 2 [1] 2</code>
Computing the Square Root	<code>sqrt(x)</code>	Computing the square root of x	<code>> sqrt (25) [1] 5</code>



Logical Values



- Logical values are TRUE and FALSE or T and F. Note that these are case sensitive. The equality operator is ==.

- to Clear Screen use Ctrl

```
RGui (32-bit) - [R Console]
File Edit View Misc Packages Windows Help

> 8 < 4
[1] FALSE
> 8 < 18
[1] TRUE
> 6*2 =12
Error in 6 * 2 = 12 : target of assignment expands to non-language object
> 6 * 2 == 12
[1] TRUE
> t ==True
Error: object 'True' not found
> T == true
Error: object 'true' not found
> T == TRUE
[1] TRUE
> % creating
Error: unexpected input in " % creating "
> |
```

R GUI (32-bit) - [R Console]

R File Edit View Misc Packages Windows Help

[File, Edit, View, Misc, Packages, Windows, Help, Stop, Help]

```
> # How to create a vector X and store values from 1 to 10
> x <- c(1:10)
> x
[1] 1 2 3 4 5 6 7 8 9 10
> x[(x>7) | (x<5)]
[1] 1 2 3 4 8 9 10
> x[(x<5) | (x>7)]
Error: object 'x.7' not found
> x[(x<5) | (x>7)]
[1] 1 2 3 4 8 9 10
> |
```

Part (i) Display 'TRUE' for elements whose values are more than 7, else display 'FALSE'.

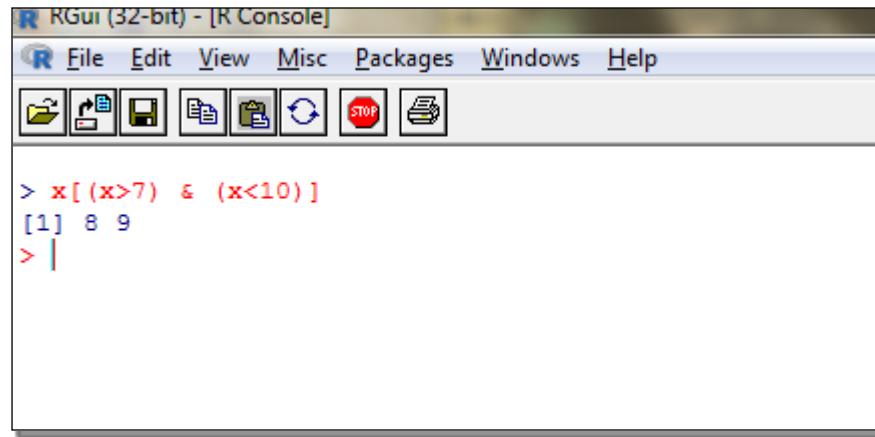
```
> x>7
```

```
[1] FALSE FALSE FALSE FALSE FALSE FALSE FALSE TRUE TRUE TRUE
```

Part (ii) Display 'TRUE' for elements whose values are less than 5, else display 'FALSE'.

```
> x<5
```

```
[1] TRUE TRUE TRUE TRUE FALSE FALSE FALSE FALSE FALSE
```



RGui (32-bit) - [R Console]

File Edit View Misc Packages Windows Help

[File, New, Open, Save, Print, Stop, Help]

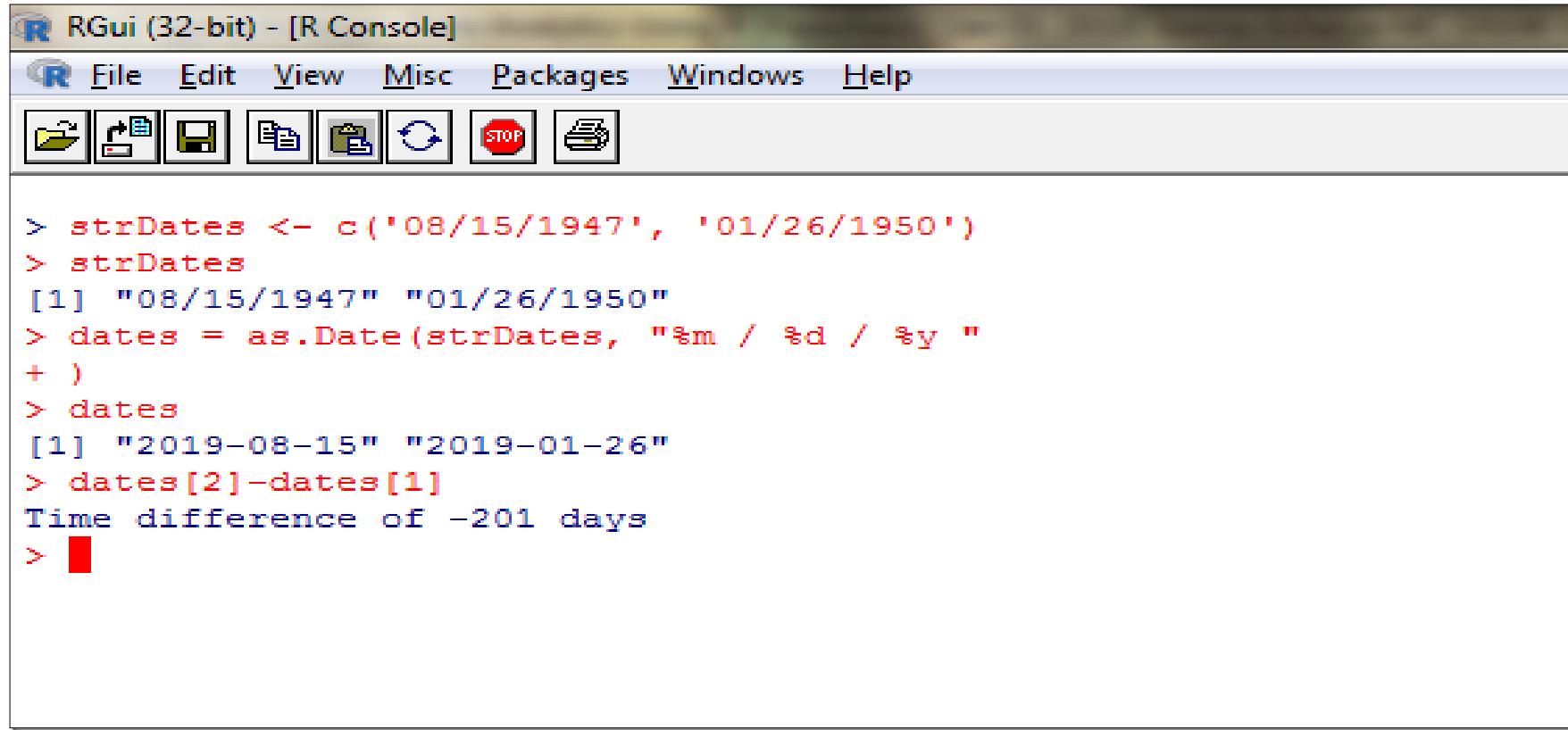
```
> x[ (x>7) & (x<10)]
[1] 8 9
> |
```

Date Activity

Dates : The default format of date is YYYY-MM-DD.

Print system's date.	Sys.Date()	[1] "2017-01-13"
Print system's time.	Sys.time()	[1] "2017-01-13 10:54:37 IST"
Print the time zone	Sys.timezone()	[1] "Asia/Calcutta"
Print today's date.	today <- Sys.Date()	today [1] "2017-01-13"
Formate Date	format (today, format = "%B %d %Y")	"January 13 2017"
Store date as a text data type.	CustomDate = "2016-01-13" class (CustomDate)	[1] "2016-01-13" [1] "character"
Convert the date stored as text data type into a date data type.	CustDate = as.Date(CustomDate) class(CustDate) CustDate	[1] "Date" [1] "2016-01-13"

Find a Difference between Two Dates



R Gui (32-bit) - [R Console]

File Edit View Misc Packages Windows Help


```
> strDates <- c('08/15/1947', '01/26/1950')
> strDates
[1] "08/15/1947" "01/26/1950"
> dates = as.Date(strDates, "%m / %d / %y "
+ )
> dates
[1] "2019-08-15" "2019-01-26"
> dates[2]-dates[1]
Time difference of -201 days
> █
```

Variables

S.N o	Variables	Explain
1	Assign a value of 50 to the variable called 'Var'.	Var <-50 or Var=5
2	Print the value in the variable, 'Var'	Var
3	Perform arithmetic operations on the variable, 'Var'.	Var + 10
4	Reassign a string value to the variable, 'Var'. > Var <- "R is a Statistical Programming Language" Print the value in the variable, 'Var'.	Var [1] "R is a Statistical Programming Language"
5	Reassign a logical value to the variable, 'Var'.	Var <- TRUE Var [1] TRUE

Loading and Handling Data in R

Challenges of Analytical Data Processing

Expression, Variables and Functions

Missing Values Treatment in R

Using the 'as' Operator to Change the Structure of Data,

Vectors

Matrices

Factors

List

Aggregating and Group Processing of a Variable,

Simple Analysis Using R,

Methods for Reading Data,

Comparison of R GUIs for Data Input

Vectors

A vector can have a list of values.

The values can be numbers, strings or logical.

All the values in a vector should be of the same data type.

A few points to remember about vectors in R are:

- Vectors are stored like arrays in C
- Vector indices begin at 1
- All vector elements must have the same mode such as integer, numeric (floating point number), character (string), logical (Boolean), complex, object, etc.

S.N o	Vector	Example
1	Create a vector of numbers The c function (c is short for combine) creates a new vector consisting of three values, viz. 4, 7 and 8.	> c(4, 7, 8) [1] 4 7 8
2	Create a vector of string values.	> c("R", "SAS", "SPSS") [1] "R" "SAS" "SPSS"
3	Create a vector of logical values.	> c(TRUE, FALSE) [1] TRUE FALSE
4	A vector cannot hold values of different data types. Consider the example below on placing integer, string and Boolean values together in a vector.	> c(4, 8, "R", FALSE) [1] "4" "8" "R" "FALSE"
5	Declare a vector by the name, 'Project' of length 3 and store values in it. > Project <- vector(length = 3) > Project [1] <- "Finance Project" > Project [2] <- "Retail Project" > Project [3] <- "Energy Project"	Project [1] "Finance Project" "Retail Project" "Energy Project" > length (Project) [1] 3

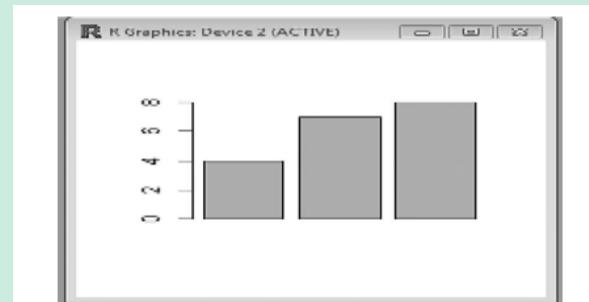
A sequence vector can be created with a start:end notation.

S.N o	Sequence Vector	Example
1	Create a sequence of numbers between 1 and 5 (both inclusive).	<pre>> 1:5 [1] 1 2 3 4 5 > seq(1:5) [1] 1 2 3 4 5</pre>
2	The default increment with seq is 1. However, it also allows the use of increments other than 1.	<pre>> seq (1, 10, 2) [1] 1 3 5 7 9 > seq (from=1, to=10, by=2) [1] 1 3 5 7 9 > seq (1, 10, by=2) [1] 1 3 5 7 9</pre>
3	seq can also generate numbers in the descending order.	<pre>> 10:1 [1] 10 9 8 7 6 5 4 3 2 1 > seq (10, 1, by=-2) [1] 10 8 6 4 2</pre>

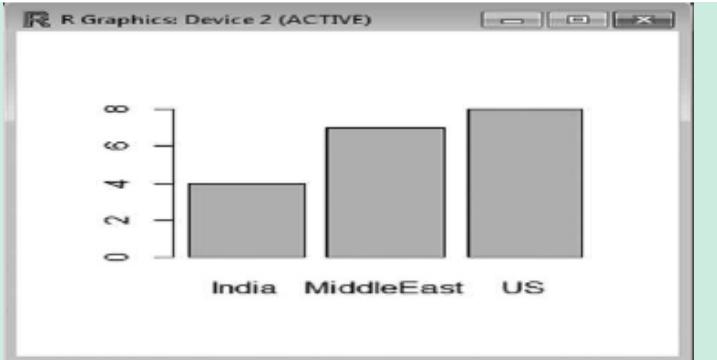
To access more than one value from the vector.

S.N o	Rep Function	Example
1	Access the first and the fifth element from the vector, 'VariableSeq'.	> VariableSeq[c(1, 5)] [1] "R" "language" "c
2	Access first to the fourth element from the vector, 'VariableSeq'.	> VariableSeq[1:4] [1] "R" "is" "a" "good programming"
3	Access the first, fourth and the fifth element from the vector, 'VariableSeq'.	> VariableSeq[c(1, 4:5)] [1] "R" "good programming" "language"
4	Retrieve all the values from the variable, 'VariableSeq'	> VariableSeq [1] "R" "is" "a" "good programming" [5] "language"

Vector Names : The names() function helps to assign names to the vector elements.

S.No	Function	Examples					
	<p>Creating a name to a vector values placeholder</p> <p>R is a programming language</p> <table> <tr> <td>1</td><td>2</td><td>3</td><td>4</td><td>5</td> </tr> </table>	1	2	3	4	5	<pre>> placeholder <- 1:5 > names(placeholder) <- c("r", "is", "a", "programming", "language")</pre>
1	2	3	4	5			
	placeholder [3]	a 3					
	<p>Plot a bar graph using the barplot function. The barplot function uses a vector's values to plot a bar chart.</p> <pre>> BarVector <- c(4, 7, 8) > barplot(BarVector)</pre>						

Vector Names : The names() function helps to assign names to the vector elements.

S.No	Function	Examples
	<p>Let us use the name function to assign names to the vector elements. These names will be used as labels in the barplot.</p> <pre>> names(BarVector) <- c("India", "MiddleEast", "US") > barplot(BarVector)</pre>	 <p>The figure shows an R graphics window titled "R Graphics: Device 2 (ACTIVE)". It displays a bar plot with three bars. The x-axis is labeled with the names "India", "MiddleEast", and "US". The y-axis has numerical tick marks at 0, 2, 4, 6, and 8. The bars are gray and have black outlines. The bar for "India" reaches approximately 4, the bar for "MiddleEast" reaches approximately 7, and the bar for "US" reaches approximately 8.</p>

S.N o	Vector Math Function	Example
1	We can run other arithmetic operations on the vector as given:	<pre>> x - 1 [1] 4 7 8 > x * 2 [1] 10 16 18 > x / 2 [1] 2.5 4.0 4.5</pre>
2	Let us practice these arithmetic operations on two vectors.	<pre>> x [1] 5 8 9 > y <- c(1, 2, 3) > y [1] 1 2 3 > x + y [1] 6 10 12</pre>

Vector Math : Let us define a vector, 'x' with three values. Let us add a scalar value (single value) to the vector. This value will get added to each vector element

S.N o	Vector Math Function	Example
1	value will get added to each vector element.	<pre>> x <- c(4, 7, 8) > x +1 [1] 5 8 9 the vector will retain its individual elements. > x [1] 4 7 8</pre>
2	If the vector needs to be updated with the new values, type the statement given below.	<pre>> x <- x + 1 > x [1] 5 8 9</pre>
3	Other Arithmetic operations	<pre>> x - 1 [1] 4 7 8 > x * 2 [1] 10 16 18 > x / 2 [1] 2.5 4.0 4.5</pre>

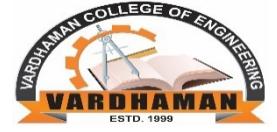
Vector Math : Let us define a vector, 'x' with three values. Let us add a scalar value (single value) to the vector. This value will get added to each vector element

S.N o	Vector Math Function	Example
1	Other arithmetic operations are:	<pre>> x - y [1] 4 6 6 > x * y [1] 5 16 27</pre>
2	Check if the two vectors are equal. The comparison takes place element by element.	<pre>> x [1] 5 8 9 > y [1] 1 2 3 > x==y [1] FALSE FALSE FALSE > x < y [1] FALSE FALSE FALSE > sin(x) [1] -0.9589243 0.9893582 0.4121185</pre>

Vector Recycling : If an operation is performed involving two vectors that requires them to be of the same length, the shorter one is recycled, i.e., repeated until it is long enough to

S.N o	Vector Recycling Function	Example
1	Add two vectors wherein one has length, 3 and the other has length, 6.	<pre data-bbox="1339 383 2018 483">> c(1, 2, 3) + c(4, 5, 6, 7, 8, 9) [1] 5 7 9 8 10 12</pre>
2	Multiply the two vectors wherein one has length, 3 and the other has length, 6.	<pre data-bbox="1339 584 2018 685">> c(1, 2, 3) * c(4, 5, 6, 7, 8, 9) [1] 4 10 18 7 16 27</pre>

Plot a Scatter Plot



The function to plot a scatter plot is 'plot'.

This function uses two vectors, i.e.

one for the x axis and another for the y axis.

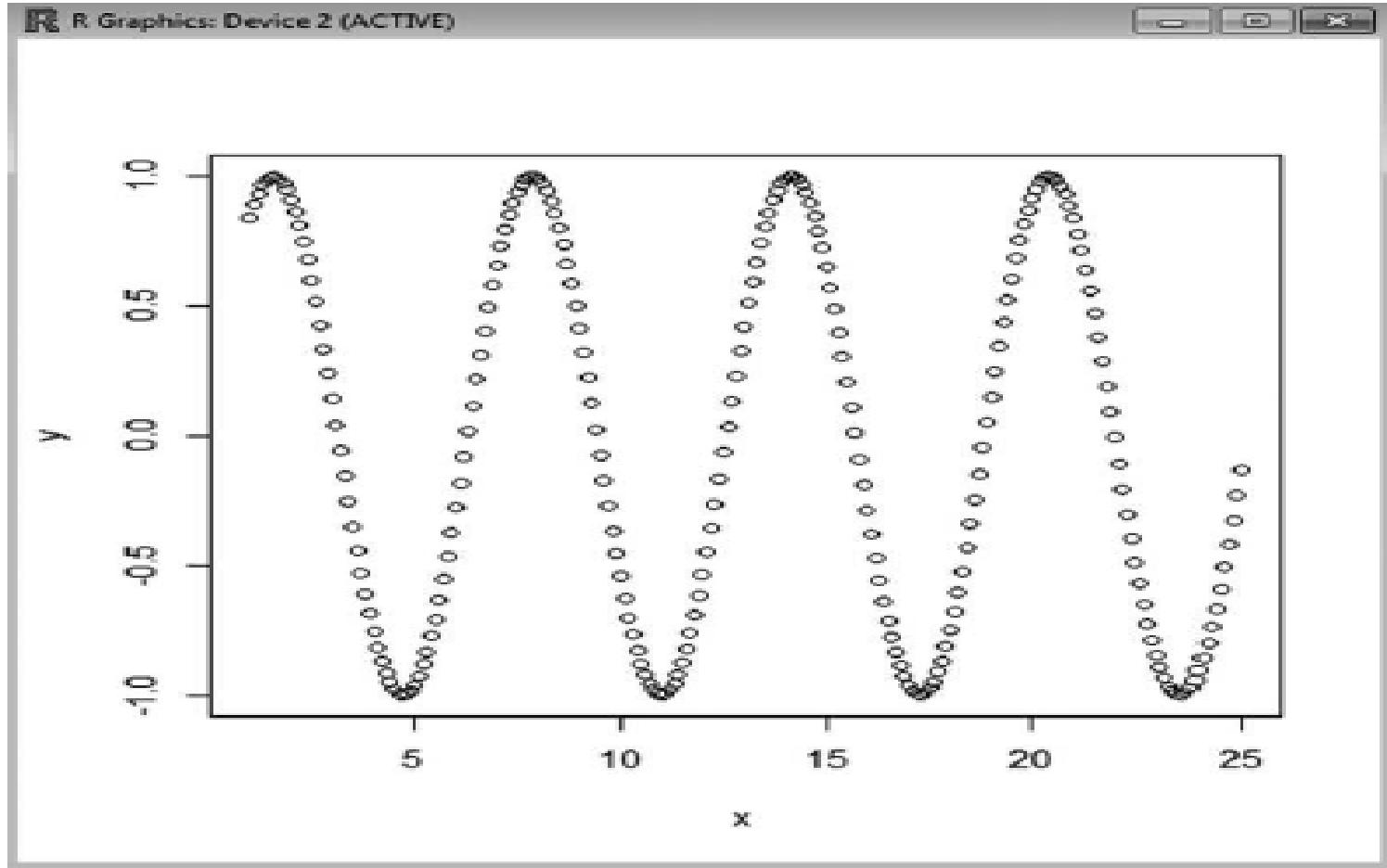
The objective is to understand the relationship between numbers and their sines.

We will use two vectors.

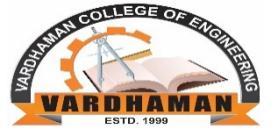
Vector, x which will have a sequence of values between 1 and 25

at an interval of 0.1 and

vector, y which stores the sines of all values held in vector, x.



Session outcomes



At the end of the session student will be able to understand:

- Introduction of R
- Install R
- Install any R package

Aggregating and group processing of a variable

R provides some functions for aggregation operation.

1. aggregate() Function: is an inbuilt function of R that aggregates data values. The function also splits data into groups after performing given statistical functions. The syntax of the aggregate() function is

aggregate(x, ...) or aggregate(x, by, FUN, ...)

where, x is an object, by argument defines the list of group elements of the specific variable of the dataset, FUN argument is a statistic function that returns a numeric value after given statistic operations and the dots ‘...’ define the other optional argument

The following example reads a table, ‘Fruit_data.csv’ into object, ‘S’. The aggregate() function computes the mean price of each type of fruit.

Here by argument is list(Fruit.Name = S\$Fruit.Name) that groups the Fruit.Name columns



RGui (64-bit)

File Edit View Misc Packages Windows Help

R Console

```
> # reading table
> S <- read.csv("Fruit_data.csv")
> S
  Fruit.Name Fruit.Price
1      Mango        80
2      Apple       100
3     Banana        40
4      Mango        70
5  Pienapple      120
6     Banana        50
7      Apple        90
8      Apple       110
9      Mango        90
> # Finding mean value of each fruit
> aggregate(S$Fruit.Price, list(Fruit.Name = S$Fruit.Name), mean)
  Fruit.Name   x
1      Apple  100
2     Banana   45
3      Mango   80
4  Pienapple 120
> |
```

The tapply() function : is also an inbuilt function of R and works in a manner similar to the function aggregate(). The function aggregates the data values into groups after performing the given statistical functions.

The syntax of the tapply () function is

tapply (x, ...) or tapply(x, INDEX, FUN, ...)

where, x is an object that defines the summary variable, INDEX argument defines the list of group elements—also called group variable, FUN argument is a statistic function that returns a numeric value after given statistic operations and the dots ‘...’ define the other optional argument.

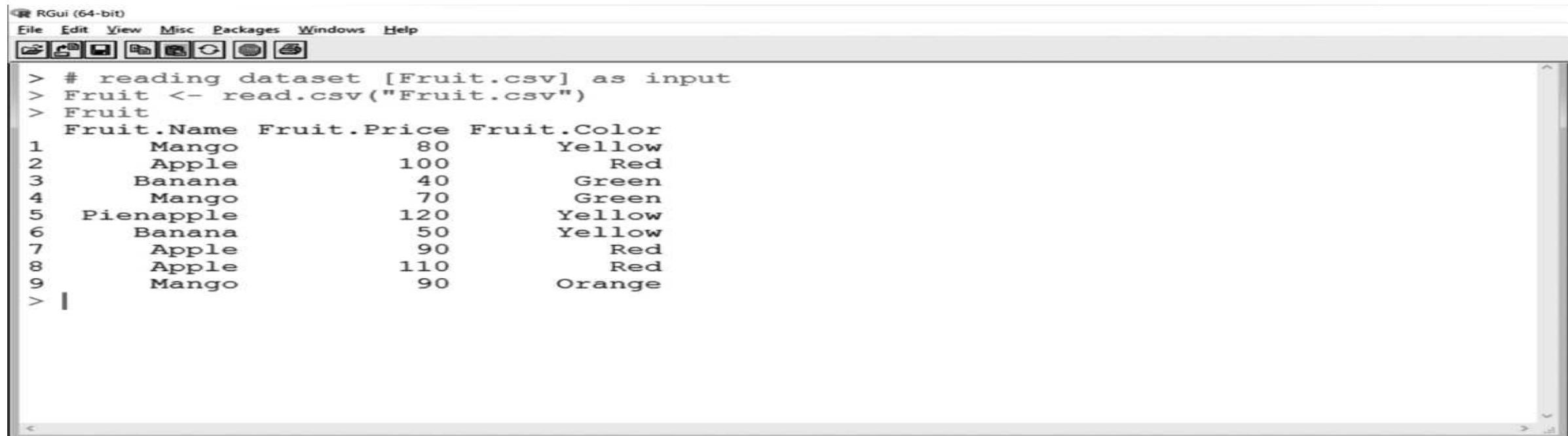
The following example reads the table, ‘Fruit_data.csv’ into object, ‘A’.
The tapply()function computes the sum and price of each type of fruit.
Here Fruit.Price is a summary variable and Fruit.Name is a grouping variable.
The FUN function is applied on the summary variable, Fruit.Price

```
RGui (64-bit)
File Edit View Misc Packages Windows Help
R Console
> A
  Fruit.Name  Fruit.Price
1      Mango        80
2      Apple       100
3     Banana        40
4      Mango        70
5  Pienapple      120
6     Banana        50
7      Apple        90
8      Apple       110
9      Mango        90
> #applying tapply()  function
> tapply(A$Fruit.Price, list(Fruit.Name = A$Fruit.Name), sum)
Fruit.Name
  Apple     Banana      Mango  Pienapple
    300       90       240       120
> tapply(A$Fruit.Price, list(Fruit.Name = A$Fruit.Name), max)
Fruit.Name
  Apple     Banana      Mango  Pienapple
    110       50       90       120
> tapply(A$Fruit.Price, list(Fruit.Name = A$Fruit.Name), min)
Fruit.Name
  Apple     Banana      Mango  Pienapple
    90       40       70       120
```

Simple analysis using R

you will learn how to read data from a dataset, perform a common operation and see the output.

Input: is the first step in any processing, including analytical data processing. Here, the input is dataset, 'Fruit'. For reading the dataset into R, use `read.table()` or `read.csv()` function.

A screenshot of the RGui (64-bit) interface. The menu bar includes File, Edit, View, Misc, Packages, Windows, and Help. Below the menu is a toolbar with various icons. The main window shows R code and its output. The code reads a CSV file named 'Fruit.csv' and prints the resulting data frame. The data frame has columns: Fruit.Name, Fruit.Price, and Fruit.Color. The output shows 9 rows of data.

```
RGui (64-bit)
File Edit View Misc Packages Windows Help
File Edit View Misc Packages Windows Help
> # reading dataset [Fruit.csv] as input
> Fruit <- read.csv("Fruit.csv")
> Fruit
  Fruit.Name Fruit.Price Fruit.Color
1   Mango        80      Yellow
2   Apple       100       Red
3  Banana        40      Green
4   Mango        70      Green
5 Pienapple     120      Yellow
6  Banana        50      Yellow
7   Apple        90       Red
8   Apple       110       Red
9   Mango        90     Orange
>
```

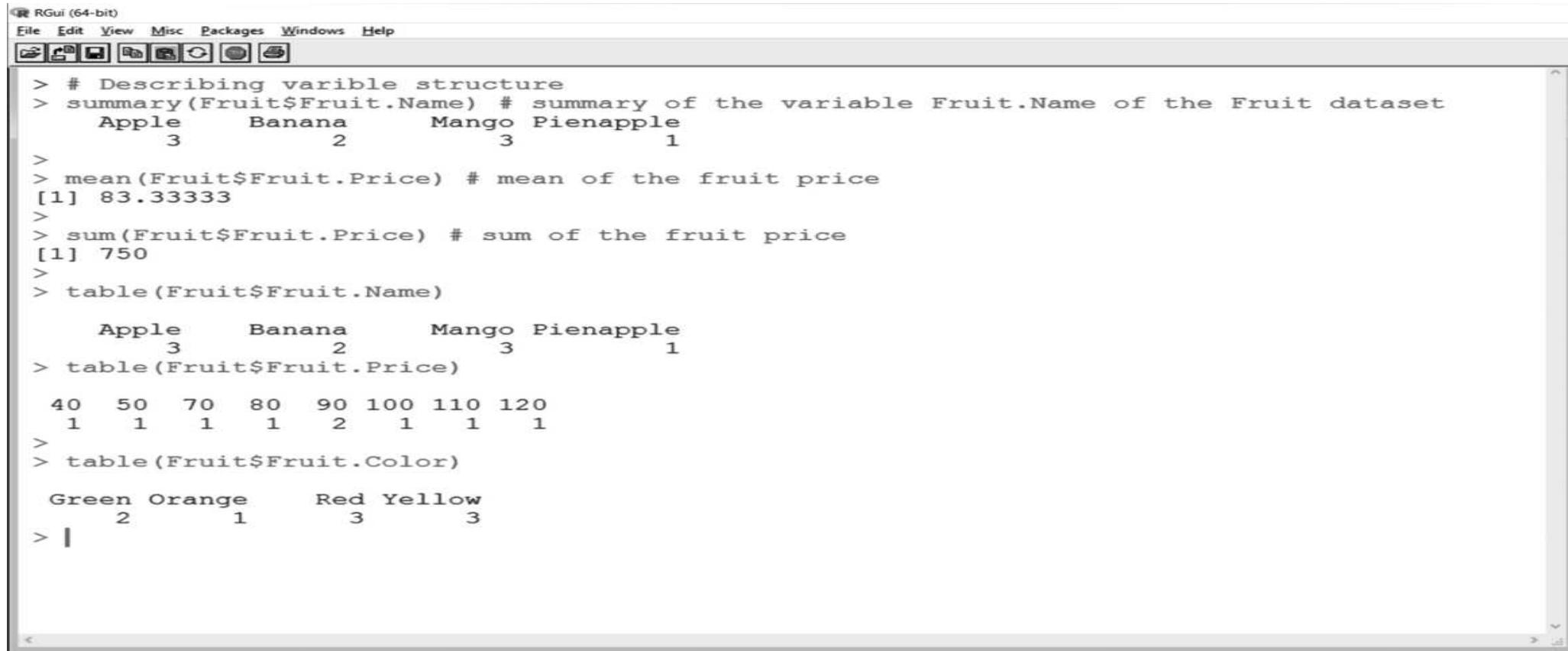
Describe Data Structure

After reading the dataset into the R workspace, the dataset can be described using different functions like `names()`, `str()`, `summary()`, `head()` and `tail()`.

```
RGui (64-bit)
File Edit View Misc Packages Windows Help
[1] > # Describing data structure
[2] > names(Fruit) # variables of the dataset
[3] [1] "Fruit.Name" "Fruit.Price" "Fruit.Color"
[4] >
[5] > str(Fruit) # structure of the dataset
[6] 'data.frame': 9 obs. of 3 variables:
[7] $ Fruit.Name : Factor w/ 4 levels "Apple", "Banana", ... : 3 1 2 3 4 2 1 1 3
[8] $ Fruit.Price: int 80 100 40 70 120 50 90 110 90
[9] $ Fruit.Color: Factor w/ 4 levels "Green", "Orange", ... : 4 3 1 1 4 4 3 3 2
[10]>
[11]> head(Fruit,3) #reading top 3 rows of the dataset
[12] Fruit.Name Fruit.Price Fruit.Color
[13] 1      Mango        80       Yellow
[14] 2      Apple         100        Red
[15] 3     Banana         40       Green
[16]>
[17]> tail(Fruit,3) #reading bottom 3 rows of the dataset
[18] Fruit.Name Fruit.Price Fruit.Color
[19] 7      Apple         90        Red
[20] 8      Apple         110        Red
[21] 9      Mango         90       Orange
[22]>
[23]> summary(Fruit) #summary of the dataset
[24]   Fruit.Name   Fruit.Price   Fruit.Color
[25] Apple       : 3    Min.   : 40.00  Green  : 2
[26] Banana      : 2  1st Qu.: 70.00  Orange: 1
[27] Mango       : 3  Median  : 90.00  Red    : 3
[28] Pienapple   : 1  Mean    : 83.33  Yellow: 3
[29]                   3rd Qu.:100.00
```

Describe Variable Structure

After describing the dataset, you can also describe the variables of the dataset using different functions. For describing the variables and performing operations on them, many functions are available.



The screenshot shows the RGui (64-bit) interface with the following R code:

```
RGui (64-bit)
File Edit View Misc Packages Windows Help
[Icons]
> # Describing variable structure
> summary(Fruit$Fruit.Name) # summary of the variable Fruit.Name of the Fruit dataset
  Apple     Banana      Mango   Pienapple
    3         2          3        1
>
> mean(Fruit$Fruit.Price) # mean of the fruit price
[1] 83.33333
>
> sum(Fruit$Fruit.Price) # sum of the fruit price
[1] 750
>
> table(Fruit$Fruit.Name)
  Apple     Banana      Mango   Pienapple
    3         2          3        1
> table(Fruit$Fruit.Price)
  40   50   70   80   90  100  110  120
  1     1     1     1     2     1     1     1
>
> table(Fruit$Fruit.Color)
  Green Orange      Red Yellow
    2       1       3       3
> |
```

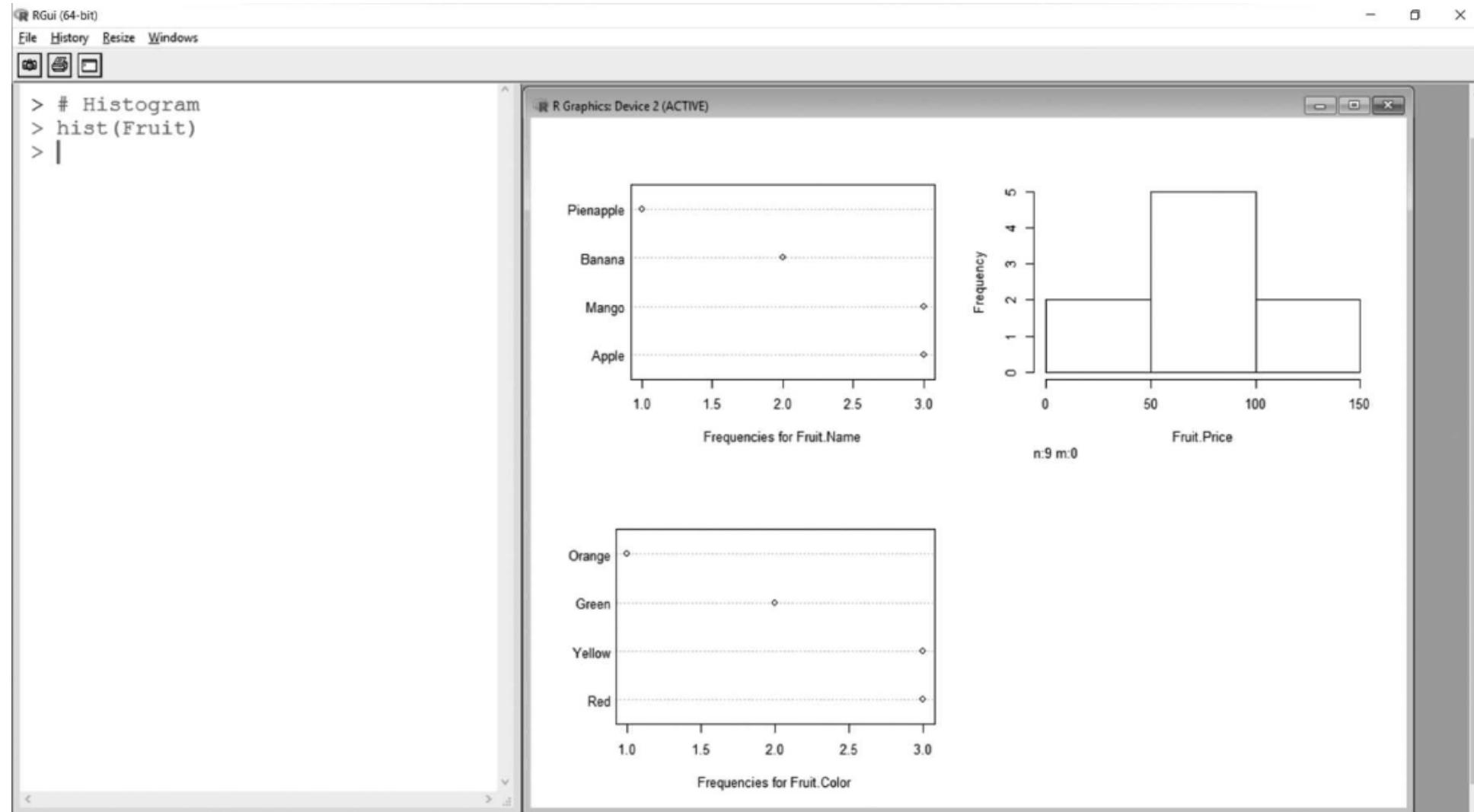
describes the histogram of the ‘Fruit’ dataset using the hist() function. A histogram is a graphical display of data that uses many bars of different heights.

The complete syntax for **hist()** function is:

```
hist(x, breaks = 'Sturges', freq = NULL, probability = !freq, include.lowest = TRUE, right = TRUE,  
density = NULL, angle = 45, col = NULL, border = NULL, main = paste('Histogram of ', xname),  
xlim = range(breaks), ylim = NULL, xlab = xname, ylab, axes = TRUE, plot = TRUE, labels = FALSE,  
nclass = NULL, warn.unused = TRUE, ...)
```

x is the vector for which a histogram is required.

freq is a logical value. If TRUE, the histogram graphic is a representation of frequencies, the counts component of the result. If FALSE, the probability densities and component density are plotted. main, xlab, ylab are arguments to title. plot is a logical value. If TRUE (default), a histogram is plotted, else a list of breaks and counts is returned



A box and whisker plot summarises the group values into boxes.

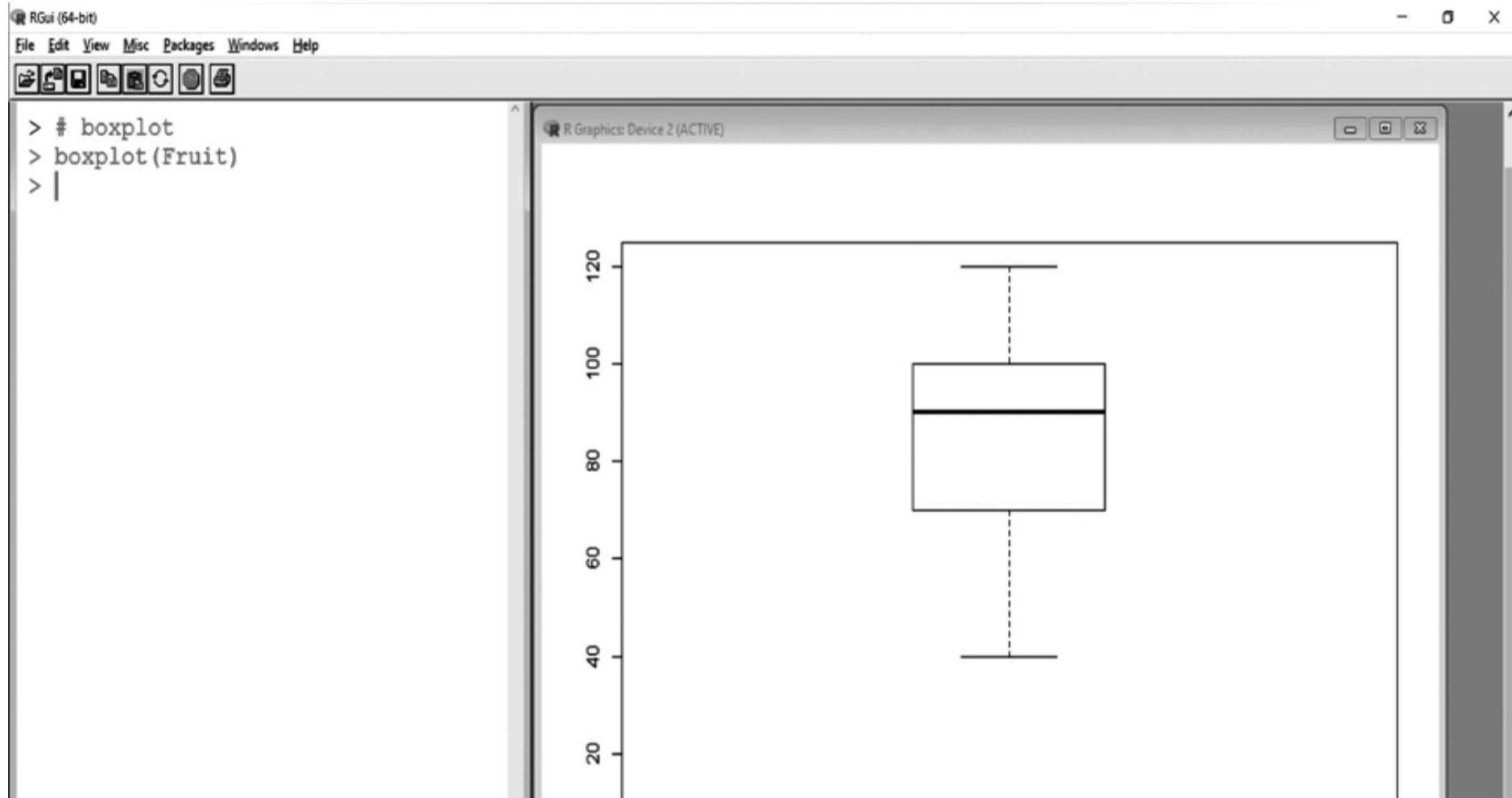
The syntax for boxplot() function is:

```
boxplot(x, ..., range = 1.5, width = NULL, varwidth = FALSE, notch = FALSE, outline = TRUE, names, plot = TRUE, border = par('fg'), col = NULL, log = "", pars = list(boxwex = 0.8, staplewex = 0.5, outwex = 0.5), horizontal = FALSE, add = FALSE, at = NULL)
```

where x is a numeric vector or a single list containing such vectors.

outline - If outline is not true, the outliers are not drawn.

range - This determines how far the plot whiskers extend out from the box.



Methods for reading data

CSV and Spreadsheets: Comma separated value (CSV) files and spreadsheets are used for storing small size data. R has an inbuilt function facility through which analysts can read both types of files.

Reading CSV Files

A CSV file uses **.csv extension** and stores data in a table structure format in any plain text. The following function reads data from a CSV file:

read.csv('filename')

where,

filename is the name of the CSV file that needs to be imported.

The `read.table()` function can also read data from CSV files. The syntax of the function is

`read.table('filename', header=TRUE, sep=',',...)`

where,

`filename` argument defines the path of the file to be read, `header` argument contains logical values `TRUE` and `FALSE` for defining whether the file has header names on the first line or not, `sep` argument defines the character used for separating each column of the file and the dots `'...'` define the other optional arguments.

RGui (64-bit)

File Edit View Misc Packages Windows Help

[Icons for File, Edit, View, Misc, Packages, Windows, Help]

R Console

```
> # Reading CSV file using .csv function()
> read.csv("Hardware.csv")
  SN      HName HPrice HCompany
1 1      Mouse   1000    Sony
2 2  Touchpad   2000    Dell
3 3 Motherboard 10000    HTC
4 4     Speaker    500     LG
>
>
> # Reading CSV file using read.table function()
> read.table("Hardware.csv", header=TRUE, sep=", ")
  SN      HName HPrice HCompany
1 1      Mouse   1000    Sony
2 2  Touchpad   2000    Dell
3 3 Motherboard 10000    HTC
4 4     Speaker    500     LG
> |
```

Reading Spreadsheets: A spreadsheet is a table that stores data in rows and columns. Many applications are available for creating a spreadsheet. Microsoft Excel is the most popular for creating an Excel file.

An Excel file **uses .xlsx** extension and stores data in a spreadsheet. In R, different **packages are available such as gdata, xlsx, etc.**, that provide functions for reading Excel files.

Importing such packages is necessary before using any inbuilt function of any package. The **read.xlsx() is an inbuilt function of ‘xlsx’ package for reading Excel files.** The **syntax of the read.xlsx() function is read.xlsx(‘filename’,...)**

filename argument defines the path of the file to be read and the dots ‘...’ define the other optional arguments.

Softdrink.xlsx - Microsoft Excel

FILE HOME INSERT PAGE LAYOUT FORMULAS DATA REVIEW VIEW Foxit Reader PDF TEAM

Cut Copy Paste Format Painter

Font Alignment Number

A15

	A	B	C	D	E	F	G	H	I
1	Software	drinkName	Price [per leter]						
2	Coca-cola		100						
3	Fruite		80						
4	Pepsi		70						
5	Maza		50						
6									
7									
8									
9									

RGui (64-bit)

File Edit View Misc Packages Windows Help

R Console

```
> # Reading spreadsheet [Excel files] after converting into [.csv file]
> read.csv("Softdrink.csv")
Software.drinkName Price...per.leter. X
1 Coca-cola 100 NA
2 Fruite 80 NA
3 Pepsi 70 NA
4 Maza 50 NA
>
```

Reading Data from Packages

library() Function : loads packages into the R workspace. It is compulsory to import the package before reading the available dataset of that package. The syntax of the library() function is:

library(packagename)

Where, packagename argument is the name of the package to be read.

data() Function: lists all the available datasets of the loaded package into the R workspace. For loading a new dataset into the loaded packages, users need to pass the name of the new dataset into data() function. The syntax of the data() function is: **data(datasetname)**

Where, datasetname argument is the name of the dataset to be read.

RGui (64-bit)

File Edit Windows



R Console

```
> # Loading a package
> library(Matrix)
> data()
> Orange
Tree   age circumference
1      1     118              30
2      1     484              58
3      1     664              87
4      1    1004             115
5      1    1231             120
6      1    1372             142
7      1    1582             145
8      2     118              33
9      2     484              69
10     2     664             111
11     2    1004             156
12     2    1231             172
13     2    1372             203
14     2    1582             203
15     3     118              30
16     3     484              51
17     3     664              75
18     3    1004             108
19     3    1231             115
20     3    1372             139
21     3    1582             140
```

R data sets

Data sets in package 'datasets':

AirPassengers	Monthly Airline Passenger Numbers 1949–1960
BJsales	Sales Data with Leading Indicator
BJsales.lead (BJsales)	Sales Data with Leading Indicator
BOD	Biochemical Oxygen Demand
CO2	Carbon Dioxide Uptake in Grass Plants
ChickWeight	Weight versus age of chicks on different diets
DNase	Elisa assay of DNase
EuStockMarkets	Daily Closing Prices of Major European Stock Indices, 1991–1998
Formaldehyde	Determination of Formaldehyde
HairEyeColor	Hair and Eye Color of Statistics Students
Harman23.cor	Harman Example 2.3
Harman74.cor	Harman Example 7.4
Indometh	Pharmacokinetics of Indomethacin
InsectSprays	Effectiveness of Insect Sprays
JohnsonJohnson	Quarterly Earnings per Johnson & Johnson \$
LakeHuron	Level of Lake Huron 1875–1972
LifeCycleSavings	Intercountry Life-Cycle Savings Data
Loblolly	Growth of Loblolly pine trees
Nile	Flow of the River Nile
Orange	Growth of Orange Trees
OrchardSprays	Potency of Orchard Sprays
PlantGrowth	Results from an Experiment on Plant Growth

Reading Data from Web/APIs

Nowadays most business organisations are using the Internet and cloud services for storing data. This online dataset is directly accessible through packages and application programming interfaces (APIs).

Different packages are available in R for reading from online datasets.

Packages	Description	Download Link
RCurl	The package permits download of files from the web server and post forms.	https://cran.r-project.org/web/packages/RCurl/index.html
Google Prediction API	It allows uploading of data to Google storage and then training them for Google Prediction API.	http://code.google.com/p/r-google-predictionapi-v121
Infochimps	The package provides the functions for accessing all API.	http://api.infochimps.com
HttpRequest	The package reads the web data with the help of an HTTP request protocol and implements the GET, POST request.	https://cran.r-project.org/web/packages/httpRequest/index.html
WDI	The package reads all World Bank data.	http://cransprojectorg/web/packages/WDI/index.html
XML	The package reads and creates an XML and HTML document with the help of an HTTP or FTP protocol.	http://cransprojectorg/web/packages/XML/index.html
Quantmod	The package reads finance data from Yahoo finance.	http://crans-projectorg/web/packages/quantmod/index.html
ScrapeR	The package reads online data.	http://crans-projectorg/web/packages/scrapeR/index.html

The following example illustrates web scraping.

Ø Web scraping extracts data from any webpage of a website.

Ø Here package ‘**RCurl**’ is used for web scraping . At first, the package, ‘**RCurl**’ is imported into the workspace and then **getURL()** function of the package, ‘**RCurl**’ takes the required webpage.

Ø Now **htmlTreeParse()** function parses the content of the webpage.

Reading a JSON (Java Script Object Notation) Document

Step 1: Install rjson package.

```
> install.packages("rjson")
```

Installing package into ‘C:/Users/seema_acharya/Documents/R/winlibrary / 3.2’(as ‘lib’ is unspecified) trying URL ‘https://cran.hafro.is/bin/windows/contrib/3.2/rjson_0.2.15.zip’

Content type ‘application/zip’ length 493614 bytes (482 KB) downloaded 482 KB package ‘rjson’ successfully unpacked and MD5 sums checked

RStudio

File Edit Code View Plots Session Build Debug Tools Help

Console

```
> # loading RCurl Package
> library(RCurl)
>
> #passing the URL for which web data required
> wd <- getURL("https://www.techopedia.com/definition/5212/web-scraping", ssl.verifypeer = FALSE)
>
> # Now parsing the webdata
> wd_parsed <- htmlTreeParse(wd)
>
> # Displaying the content of webpage
> wd_parsed
```

Environment History

Import Dataset List

Global Environment

Values

answer3	11.25
jan09	"<!DOCTYPE HTML PUBLIC ...
jan09_pa...	Large XMLDocumentContent...
p	num [1:6] 2 2.2 2.4 2.6...
p1	num [1:6] 2 2.2 2.4 2.6...
s	num [1:6] 5 7 9 11 NA NA
wd	"\r\n<!DOCTYPE html>\r\..."
wd_parsed	Large XMT.DocumentContent...

Files Plots Packages Help Viewer

New Folder Delete Rename More

Home

Name	Size	Modified
.RData	2.7KB	Jul 3, 2016, 10:54 AM
Rhistory	919 B	Jul 3, 2016, 10:54 AM
2-extra.docx	19.6 KB	Mar 3, 2016, 2:12 PM
25-extra.docx	37.7 KB	Feb 25, 2016, 11:08 PM
Bluetooth.Folder		
CoffeeCup Software		
Custom Office Templates		
Datasheet.accdt	381 KB	May 28, 2015, 11:50 PM

Step 2: Input data.

Store the data given below in a text file ('D:/Jsondoc.json'). Ensure that the file is saved with an extension of .json

```
{  
    'EMPID': ['1001', '2001', '3001', '4001', '5001', '6001', '7001', '8001'  
    ],  
    'Name': ['Ricky', 'Danny', 'Michelle', 'Ryan', 'Gerry', 'Nonita', 'Sim  
    on', 'Gallop'],  
    'Dept': ['IT', 'Operations', 'IT', 'HR', 'Finance', 'IT', 'Operations'  
, 'Finance']  
}
```

A JSON document begins and ends with a curly brace ({}). A JSON document is a set of key value pairs. Each key:value pair is delimited using , as a delimiter

Step 3: Read the JSON file, ‘d:/Jsondoc.json’.

```
> output <- fromJSON(file = "d:/Jsondoc.json")  
> output  
$EMPID  
[1] "1001" "2001" "3001" "4001" "5001" "6001" "7001" "8001"  
$Name  
[1] "Ricky" "Danny" "Michelle" "Ryan" "Gerry" "Nonita"  
[7] "Simon" "Gallop"  
$Dept  
[1] "IT" "Operations" "IT" "HR" "Finance"  
[6] "IT" "Operations" "Finance"
```

Step 4: Convert JSON to a data frame.

```
> JSONDataFrame <- as.data.frame(output)
```

Display the content of the data frame, 'output'.

```
> JSONDataFrame
```

EMPID Name Dept

1 1001 Ricky IT

2 2001 Danny Operations

3 3001 Michelle IT

4 4001 Ryan HR

5 5001 Gerry Finance

6 6001 Nonita IT

7 7001 Simon Operations

8 8001 Gallop Finance

Reading an XML File

Step 1: Install an XML package.

```
> install.packages("XML")
```

Installing package into 'C:/Users/seema_acharya/Documents/R/winlibrary/ 3.2'(as 'lib' is unspecified)

trying URL 'https://cran.hafro.is/bin/windows/contrib/3.2/XML_3.98-1.3.zip'
Content type 'application/zip' length 4299803 bytes (4.1 MB)

downloaded 4.1 M package 'XML' successfully unpacked and MD5 sums checked

Step 2: Input data.

Store the data below in a text file (XMLFile.xml in the D: drive). Ensure that the file is saved with an extension of .xml.

```
<RECORDS>
  <EMPLOYEE>
    <EMPID>1001</EMPID>
    <EMPNAME>Merrilyn</EMPNAME>
    <SKILLS>MongoDB</SKILLS>
    <DEPT>Computer Science</DEPT>
  </EMPLOYEE>

  <EMPLOYEE>
    <EMPID>1002</EMPID>
    <EMPNAME>Ramya</EMPNAME>
    <SKILLS>People Management</SKILLS>
    <DEPT>Human Resources</DEPT>
  </EMPLOYEE>

  <EMPLOYEE>
    <EMPID>1003</EMPID>
    <EMPNAME>Fedora</EMPNAME>
    <SKILLS>Recruitment</SKILLS>
    <DEPT>Human Resources</DEPT>
  </EMPLOYEE>
</RECORDS>
```

Reading an XML File: The xml file is read in R using the function `xmlParse()`. It is stored as a list in R.

Step 1: Begin by loading the required packages.

```
> library("XML")
```

Warning message:

package ‘XML’ was built under R version 3.2.3

```
> library ("methods")
```

```
> output <- xmlParse(file = "d:/XMLFile.xml")
```

```
> print(output)
```

```
<?xml version="1.0"?>
<RECORDS>
    <EMPLOYEE>
        <EMPID>1001</EMPID>
        <EMPNAME>Merrilyn</EMPNAME>
        <SKILLS>MongoDB</SKILLS>
        <DEPT>ComputerScience</DEPT>
    </EMPLOYEE>
    <EMPLOYEE>
        <EMPID>1002</EMPID>
        <EMPNAME>Ramya</EMPNAME>
        <SKILLS>PeopleManagement</SKILLS>
        <DEPT>HumanResources</DEPT>
    </EMPLOYEE>
    <EMPLOYEE>
        <EMPID>1003</EMPID>
        <EMPNAME>Fedora</EMPNAME>
        <SKILLS>Recruitment</SKILLS>
        <DEPT>HumanResources</DEPT>
    </EMPLOYEE>
</RECORDS>
```

Step 2: Extract the root node from the XML file.

```
> rootnode <- xmlRoot(output)
```

Find the number of nodes in the root.

```
> rootsize <- xmlSize(rootnode)
```

```
> rootsize
```

```
[1] 3
```

Let us display the details of the first node.

```
> print (rootnode[1])
```

```
$EMPLOYEE
```

```
<EMPLOYEE>
```

```
<EMPID>1001</EMPID>
```

```
<EMPNAME>Merrilyn</EMPNAME>
```

```
<SKILLS>MongoDB</SKILLS>
```

```
<DEPT>ComputerScience</DEPT>
```

```
</EMPLOYEE>
```

```
attr(“class”)
```

```
[1] “XMLInternalNodeList” “XMLNodeList”
```

Let us display the details of the first element of the first node.

```
> print(rootnode[[1]][[1]])
```

```
<EMPID>1001</EMPID>
```

Let us display the details of the third element of the first node.

```
> print(rootnode[[1]][[3]])
```

```
<SKILLS>MongoDB</SKILLS>
```

Next, display the details of the third element of the second node.

```
> print(rootnode[[2]][[3]])
```

```
<SKILLS>PeopleManagement</SKILLS>
```

We can also display the value of 2nd element of the first node.

```
> output <-xmlValue(rootnode[[1]][[2]])
```

```
> output
```

```
[1] "Merrilyn"
```

Step 3: Convert the input xml file to a data frame using the `xmlToDataFrame` function.

```
> xmldataframe <- xmlToDataFrame("d:/XMLFile.xml")
```

Display the output of the data frame.

```
> xmldataframe
```

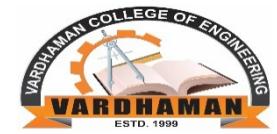
EMPID	EMPNAME	SKILLS	DEPT
-------	---------	--------	------

1	1001	Merrilyn	MongoDB ComputerScience
---	------	----------	-------------------------

2	1002	Ramya	PeopleMananement HumanResources
---	------	-------	---------------------------------

3	1003	Fedora	Recruitment HumanResources
---	------	--------	----------------------------

Comparison of R GUIs for data input



- ØR is mainly used for statistical analytical data processing.
- ØAnalytical data processing needs a large dataset that is stored in a tabular form. Sometimes it is difficult to use inbuilt functions of R for doing such analytical data processing operations in R console.
- ØHence, to overcome this problem, GUI is developed for R. Graphical user interface is a graphical medium through which users interact with the language or perform operations.
- ØDifferent GUIs are available for data input in R.
- ØEach GUI has its own features.

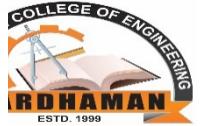
GUI Name	Description	Download Weblink
RCommander (Rcmdr)	<ul style="list-style-type: none"> RCommander was developed by John Fox and licensed under the GNU public license. It comes with many plug-ins and has a very simple interface. Users can install it like other packages of R within language. 	http://socserv.mcmaster.ca/jfox/Misc/Rcmdr/ Or https://cran.r-project.org/web/packages/Rcmdr/index.html
Rattle	<ul style="list-style-type: none"> Dr. Graham Williams developed the Rattle GUI package written in R. Data mining operation is the main application area of Rattle. It offers statistical analysis, validation, testing and other operations. 	http://rattle.togaware.com/ Or http://rattle.togaware.com/rattle-install-mswindows.html
RKWard	<ul style="list-style-type: none"> RKWard community developed the RKWard package. It provides a transparent front end and supports different features for doing analytical operations in R. It supports different platforms, such as Windows, Linux, BSD, and OS X. 	https://rkward.kde.org/ Or http://download.kde.org/stable/rkward/0.6.5/win32/install_rkward_0.6.5.exe
JGR (Java GUI for R)	<ul style="list-style-type: none"> Markus Helbig, Simon Urbanek, and Ian Fellows developed JGR. JGR is a universal GUI for R that supports cross platform. Users can use it as a replacement for the default R GUI on Windows. 	http://www.rforge.net/JGR/ Or https://cran.r-project.org/web/packages/JGR/
Deducer	<ul style="list-style-type: none"> Deducer is another simple GUI that has a menu system for doing common data operations, analytical processing and other operations. It is mainly designed to use it with the Java-based R Console [JGR]. 	Main.DeducerManual">http://www.deducer.org/pmwiki/pmwiki.php?n>Main.DeducerManual Or Main.DownloadingAndInstallingDeducer">http://www.deducer.org/pmwiki/index.php?n>Main.DownloadingAndInstallingDeducer

Using R With Databases

RODBC is a package of languages that interacts with a database. Michael Lapsley and Brian Ripley developed this package.

RODBC helps in accessing databases such as MS Access and Microsoft SQL Server through an ODBC interface. Its package has many inbuilt functions for performing database operations on the database.

| Major functions of RODBC



Function	Description
<code>odbcConnect(dsn, uid= '', pwd= '')</code> where, dsn is domain name server, uid is the user ID and pwd is the password.	The function opens a connection to an ODBC database.
<code>sqlFetch(sqltable)</code> where, sqltable is name of the SQL table.	The function reads a table from an ODBC database into a data frame.
<code>sqlQuery(query)</code> where, query is the SQL query.	The function takes a query, sends to an ODBC database and returns its result.
<code>sqlSave(dataframe, tablename= 'sqltable')</code> where, data frame defines the data frame object and tablename argument is the name of the table.	The function writes or updates a data frame to a table in the ODBC database.
<code>sqlDrop(sqltable)</code> where, sqltable is the name of the SQL table.	The function removes a table from the ODBC database.
<code>odbcClose()</code>	The function closes the open connection.

Here is a sample code where package RODBC is used for reading data from a database.

```
># importing package
> library(RODBC)
> connect1 <- odbcConnect(dsn = 'servername', uid= '    ', pwd= '    ')
#Open connection
> query1 <- 'Select * from lib.table where...'
> Demodb <- sqlQuery(connect1, query1, errors = TRUE)
> odbcClose(connection) #Close the connection
```

Using MySQL and R

- Ø MySQL is an open source SQL database system.
- Ø It is a small-sized popular database that is available for free download.
- Ø For accessing MySQL database, users need to install the MySQL database system on their computers.
- Ø MySQL database can be downloaded and installed from its official website.
- Ø R also provides a package, ‘RMySQL’ used for accessing the database from the MySQL database. Like other packages, RMySQL has many inbuilt functions for interacting with a database.

Major functions of RMySQL

Function	Description
<pre>dbConnect (MySQL(), uid= '', pwd= '', dbname = '',...)</pre> <p>where, MySQL() is MySQL driver, uid is the user ID, pwd is the password and dbname is the database name.</p>	The function opens a connection to the MySQL database.
<pre>dbDisconnect (connectionname)</pre> <p>where, Connectionname defines the name of the connection.</p>	The function closes the open connection.
<pre>dbSendQuery (connectionname, sql)</pre> <p>where, connectionname defines the name of the connection.</p>	The function runs the SQL queries of the open connection.
<pre>dbListTables (connectionname)</pre> <p>where, connectionname defines the name of the connection.</p>	The function lists the tables of the database of the open connection.
<pre>dbWriteTable (connectionname, name = 'table name', value = data.frame.name)</pre> <p>where, connectionname defines the name of the connection.</p>	The function creates the table and alternatively writes or updates a data frame in the database.

A sample code to illustrate the use of RMySQL for reading data from a database is given below.

```
># importing package  
>library(RMySQL)  
#Open connection 'connectm'  
> connectm <- odbcConnect(MySQL(), uid= "", pwd= "", dbname = "", host = "")  
> querym <- 'Select * from lib.table where...'  
> Demom<- dbSendQuery(connectm, querym)  
#Close the connection 'connect'  
>dbDisconnect(connectm)
```

Using PostgreSQL and R

- Ø PostgreSQL is an open source and customisable SQL database system.
- Ø After MySQL, PostgreSQL database is used for business analytical processing.
- Ø For accessing the PostgreSQL database, users need to install the PostgreSQL database system on their computer system.
- Ø Please note that it requires a server.
- Ø Users can get a server on rent, download and install the MySQL database from its official website.
- Ø R has a package, ‘RPostgreSQL’ that is used for accessing the database from the PostgreSQL database.
- Ø Like other packages, RPostgreSQL has many inbuilt functions for interacting with its database.

Major functions of the RPostgreSQL

<i>Function</i>	<i>Description</i>
<pre>dbConnect (driverobject, uid= '', pwd= '', dbname = '',...)</pre> <p>where, driverobject is an object of database driver, uid is the user ID, pwd is the password and dbname is the database name.</p>	The function opens a connection to an RPostgreSQL database.
<pre>dbDisconnect (connectionname)</pre> <p>where, Connectionname defines the name of the connection.</p>	The function closes the open connection.

Using SQLite and R

- Ø SQLite is a server-less, self-contained, transactional and zero-configuration SQL database system.
- Ø It is an embedded SQL database engine that does not require any server, due to which it is called a serverless database.
- Ø The database also supports all business analytical data processing.
- Ø R has an RSQLite package that is used for accessing a database from the SQLite database.
- Ø The RSQLite5 has many inbuilt functions for working with the database.
- Ø Like other packages used for accessing a database, as explained in the previous sections, users can use the same methods—dbConnect() and dbDisconnect() for opening and closing the connection from the SQLite database, respectively.
- Ø The only difference here is that users have to pass the SQLite database driver object in the dbConnect() function.

Using JasperDB and R

- Ø JasperDB is another open source database system integrated with R.
- Ø It was developed by the Jaspersoft community. It provides many business intelligence tools for analytical business processing.
- Ø A Java library interface is used between JasperDB and R. It is called ‘RevoConnectR for JasperReports Server’.
- Ø The dashboard of the JasperReports Server provides many features through which R charts, an output of the RevoDeploy R, etc., are easily accessible.
- Ø Like other packages, JasperDB has a package or web service framework called ‘RevoDeployR’ developed by Revolution Analytics.
- Ø RevoDeploy R6 provides a set of web services with security features, scripts, APIs and libraries in a single server.
- Ø It easily integrates with the dynamic R-based computations into web applications

Using Pentaho and R

- Ø Pentaho is one of the most famous companies in the data integration field that develops different products and provides services for big data deployment and business analytics.
- Ø The company provides different open source-based and enterprise-class platforms.
- Ø Pentaho Data Integration (PDI) is one of the products of Pentaho7 used for accessing database and Ø analytical data processing. It prepares and integrates data for creating a perfect picture Ø of any business.
- Ø The tool provides accurate and analytics-ready data reports to the end users, eliminates the coding complexity and uses big data in one place.
- Ø R Script Executor is one of the inbuilt tools of the PDI tool for establishing a relationship between R and Pentaho Data Integration.
- Ø Through R Script Executor, users can access data and perform analytical data operations. If users have R in their system already, then they just need to install PDI from its official website.
- Ø The users need to configure environment variables, Spoon, DI Server, and Cluster nodes as well.
- Ø Although users can try PDI and transform a database using R Script Executor, PDI is a paid tool for doing analytical data integration operation.

Thank You!

Topic: Introduction to R Language

Session no: 2

Course name: DATA ANALYTICS USING R

Semester: VII

Regulation : R18

Presented by

Loading and Handling Data in R

Challenges of Analytical Data Processing

Expression, Variables and Functions

Missing Values Treatment in R

Using the 'as' Operator to Change the Structure of Data,

Vectors

Matrices

Factors

List

Aggregating and Group Processing of a Variable,

Simple Analysis Using R,

Methods for Reading Data,

Comparison of R GUIs for Data Input

List : List is similar to C Struct.

S.No	Command	Example
1	Create a list in R. To create a list, 'emp' having three elements, 'EmpName', 'EmpUnit' and 'EmpSal'.	<pre>Øemp <- list ("EmpName="Alex", EmpUnit = "IT", EmpSal = 55000) > emp \$EmpName [1] "Alex" \$EmpUnit [1] "IT"</pre>
2	Emp list	<pre>EmpList <- list("Alex", "IT", 55000) EmpList [[1]] [1] "Alex" [[2]] [1] "IT" [[3]] [1] 55000</pre>

A list has elements. The elements in a list can have names, which are referred to as tags.

Elements can also have values.

S.No	Command	Example
1	Retrieve the names of the elements in the list 'emp'.	> names(emp) [1] "EmpName" "EmpUnit" "EmpSal"
2	Retrieve the values of the elements in the list 'emp'.	> unlist(emp) EmpName EmpUnit EmpSal "Alex" "IT" "55000"
3	Retrieve the value of the element 'EmpName' in the list 'emp'.	> unlist(emp["EmpName"]) EmpName "Alex"
4	Retrieve the value of the element 'EmpName' in the list 'emp'.	> emp[["EmpName"]] [1] "Alex"

List : List is similar to C Struct.

Add/Delete Element to or from a List	
1	<p>Before adding an element to the list 'emp', let us verify what elements exist in the list.</p>
2	<p>Add an element with the name 'EmpDesg' and value 'Software Engineer' to the list, 'emp'. <code>> emp\$EmpDesg = "Software Engineer"</code></p>

List : List is similar to C Struct.

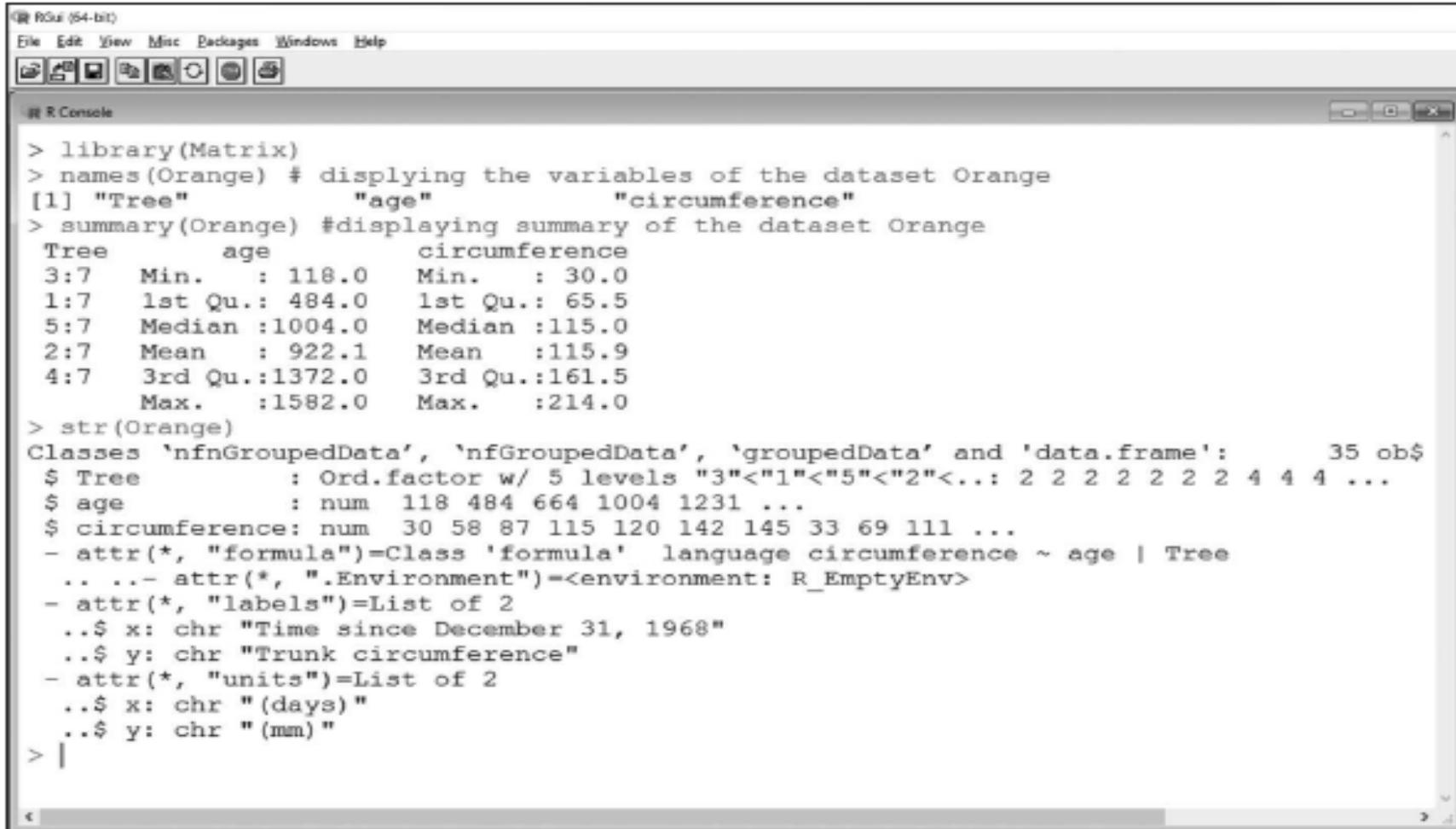
Add/Delete Element to or from a List	
1	<p>Delete an element with the name 'EmpUnit' and value 'IT' from the list, 'emp'.</p> <pre>> emp\$EmpUnit <- NULL</pre>
2	<p>Determine the number of elements in the list, 'emp'</p> <p>length() function can be used to determine the number of elements present in the list. The list, 'emp' has three elements as shown:</p>

	<p>Recursive List: A recursive list means a list within a list.</p> <p>Create a list within a list.</p> <p>Let us begin with two lists, 'emp' and 'emp1'.</p> <p>The elements in both the lists are as shown below.</p> <p>We would like to combine both the lists into a single list called 'EmpList'.</p> <pre>> EmpList <- list(emp, emp1)</pre>	<pre>> emp \$EmpName [1] "Alex" 78 Data Analytics using R \$EmpSal [1] 55000 \$EmpDesg [1] "Software Engineer" > emp1 \$EmpUnit [1] "IT" \$EmpCity [1] "Los Angeles"</pre>

Task

TABLE 3.4 Functions for exploring a dataset

Functions	Function Arguments	Description
<code>names(dataset)</code>	<ul style="list-style-type: none"> Dataset argument contains the name of the dataset. 	The function displays the variables of the given dataset.
<code>summary(dataset)</code>	<ul style="list-style-type: none"> Dataset argument contains the name of the dataset. 	The function displays the summary of the given dataset.
<code>str(dataset)</code>	<ul style="list-style-type: none"> Dataset argument contains the name of the dataset. 	The function displays the structure of the given dataset.
<code>head(dataset, n)</code>	<ul style="list-style-type: none"> Dataset argument contains the name of the dataset. <i>n</i> is a numeric value to display the number of top rows. 	The function displays the top rows according to the value of <i>n</i> . If the value of <i>n</i> is not provided in the function then by default the function displays the top 6 rows of the dataset.
<code>tail(dataset, n)</code>	<ul style="list-style-type: none"> Dataset argument contains the name of the dataset. <i>n</i> is a numeric value to display the number of bottom rows. 	The function displays the top rows according to the value of <i>n</i> . If the value of <i>n</i> is not provided in the function then by default the function displays the bottom 6 rows of the dataset.
<code>class(dataset)</code>	<ul style="list-style-type: none"> Dataset argument contains the name of the dataset. 	The function displays the class of the dataset.
<code>dim(dataset)</code>	<ul style="list-style-type: none"> Dataset argument contains the name of the dataset. 	The function returns the dimension of the dataset which implies the total number of rows and columns of the dataset.
<code>table(dataset\$variable names)</code>	<ul style="list-style-type: none"> Dataset argument contains the name of the dataset. Variable name contains the name of the variable names. 	The function returns the number of categorical values after counting them.

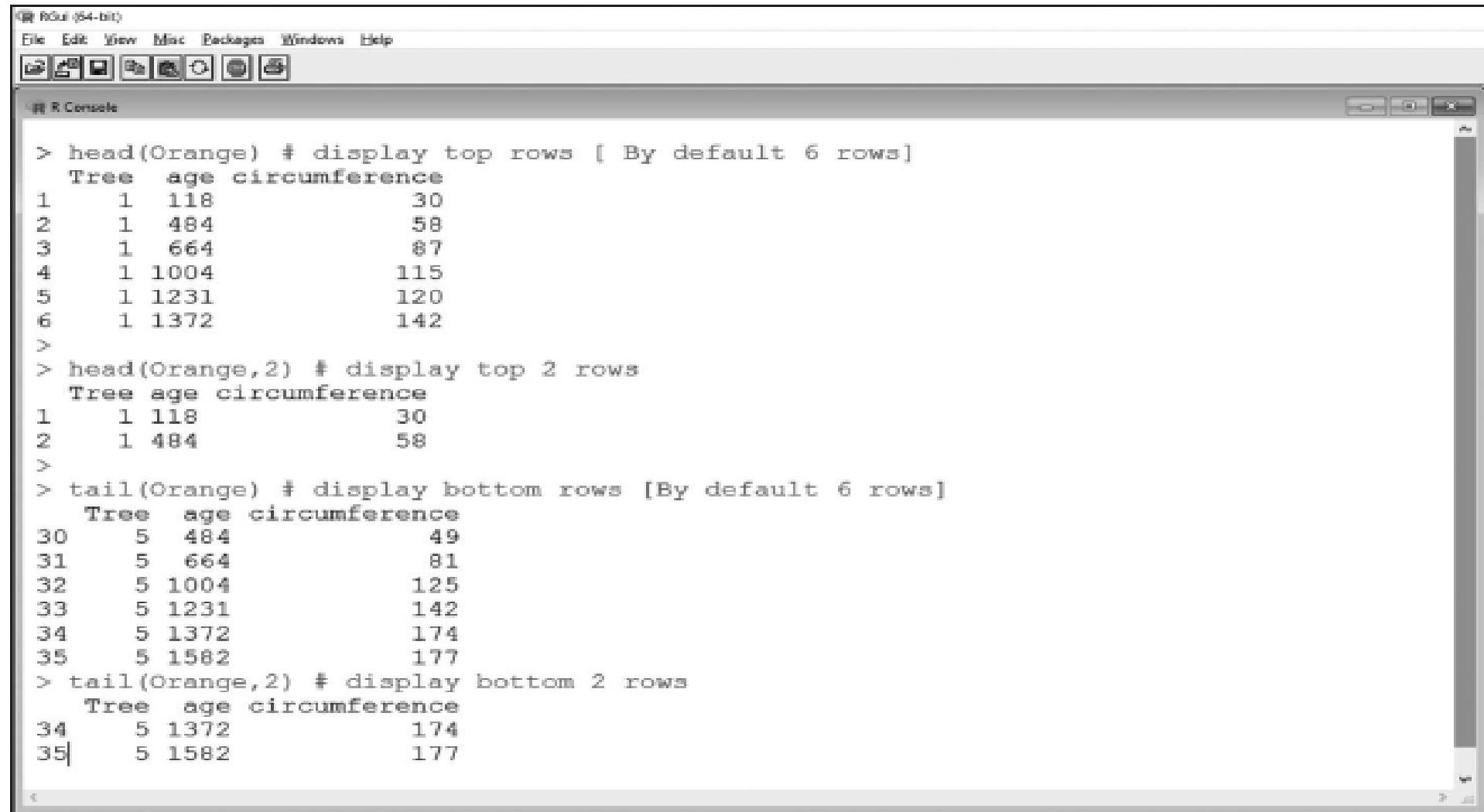


The screenshot shows the RStudio interface with the R Console tab active. The console window displays the following R code and its output:

```
RGui (64-bit)
File Edit View Misc Packages Windows Help
R Console

> library(Matrix)
> names(Orange) # displaying the variables of the dataset Orange
[1] "Tree"          "age"           "circumference"
> summary(Orange) #displaying summary of the dataset Orange
Tree      age      circumference
3:7   Min. :118.0  Min. : 30.0
1:7   1st Qu.:484.0  1st Qu.: 65.5
5:7   Median :1004.0 Median :115.0
2:7   Mean   :922.1  Mean  :115.9
4:7   3rd Qu.:1372.0 3rd Qu.:161.5
      Max. :1582.0  Max. :214.0
> str(Orange)
Classes 'nfnGroupedData', 'nfGroupedData', 'groupedData' and 'data.frame': 35 obs$ 
 $ Tree       : Ord.factor w/ 5 levels "3" < "1" < "5" < "2" < ...: 2 2 2 2 2 2 2 4 4 4 ...
 $ age        : num  118 484 664 1004 1231 ...
 $ circumference: num  30 58 87 115 120 142 145 33 69 111 ...
 - attr(*, "formula")=Class 'formula' language circumference ~ age | Tree
 .. . . . attr(*, ".Environment")=<environment: R_EmptyEnv>
 - attr(*, "labels")=List of 2
 .. $ x: chr "Time since December 31, 1968"
 .. $ y: chr "Trunk circumference"
 - attr(*, "units")=List of 2
 .. $ x: chr "(days)"
 .. $ y: chr "(mm)"
> |
```

FIGURE 3.13 Exploring a dataset using `names()`, `summary()` and `str()` functions

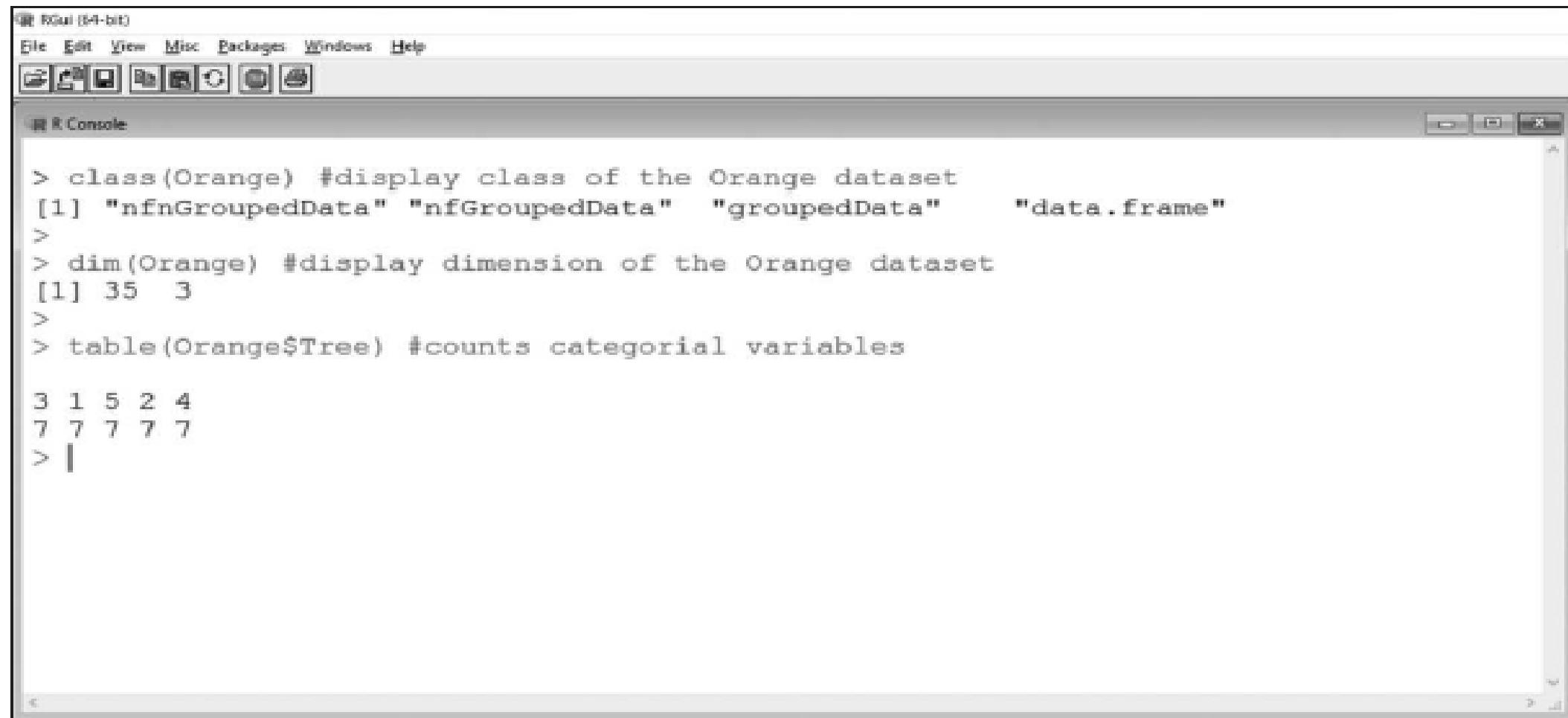


RGui (64-bit)
File Edit View Mac Packages Windows Help

R Console

```
> head(Orange) # display top rows [ By default 6 rows]
  Tree age circumference
  1    1    118          30
  2    1    484          58
  3    1    664          87
  4    1   1004         115
  5    1   1231         120
  6    1   1372         142
>
> head(Orange,2) # display top 2 rows
  Tree age circumference
  1    1    118          30
  2    1    484          58
>
> tail(Orange) # display bottom rows [By default 6 rows]
  Tree age circumference
 30   5    484          49
 31   5    664          81
 32   5   1004         125
 33   5   1231         142
 34   5   1372         174
 35   5   1582         177
> tail(Orange,2) # display bottom 2 rows
  Tree age circumference
 34   5   1372         174
 35   5   1582         177
```

FIGURE 3.14 Exploring a dataset using `head()` and `tail()` functions



RGui (64-bit)

File Edit View Misc Packages Windows Help

R Console

```
> class(Orange) #display class of the Orange dataset
[1] "nfNGroupedData" "nfGroupedData"   "groupedData"      "data.frame"
>
> dim(Orange) #display dimension of the Orange dataset
[1] 35  3
>
> table(Orange$Tree) #counts categorial variables

 3 1 5 2 4
7 7 7 7 7
> |
```

FIGURE 3.15 Exploring a dataset using `class()`, `dim()` and `table()` functions

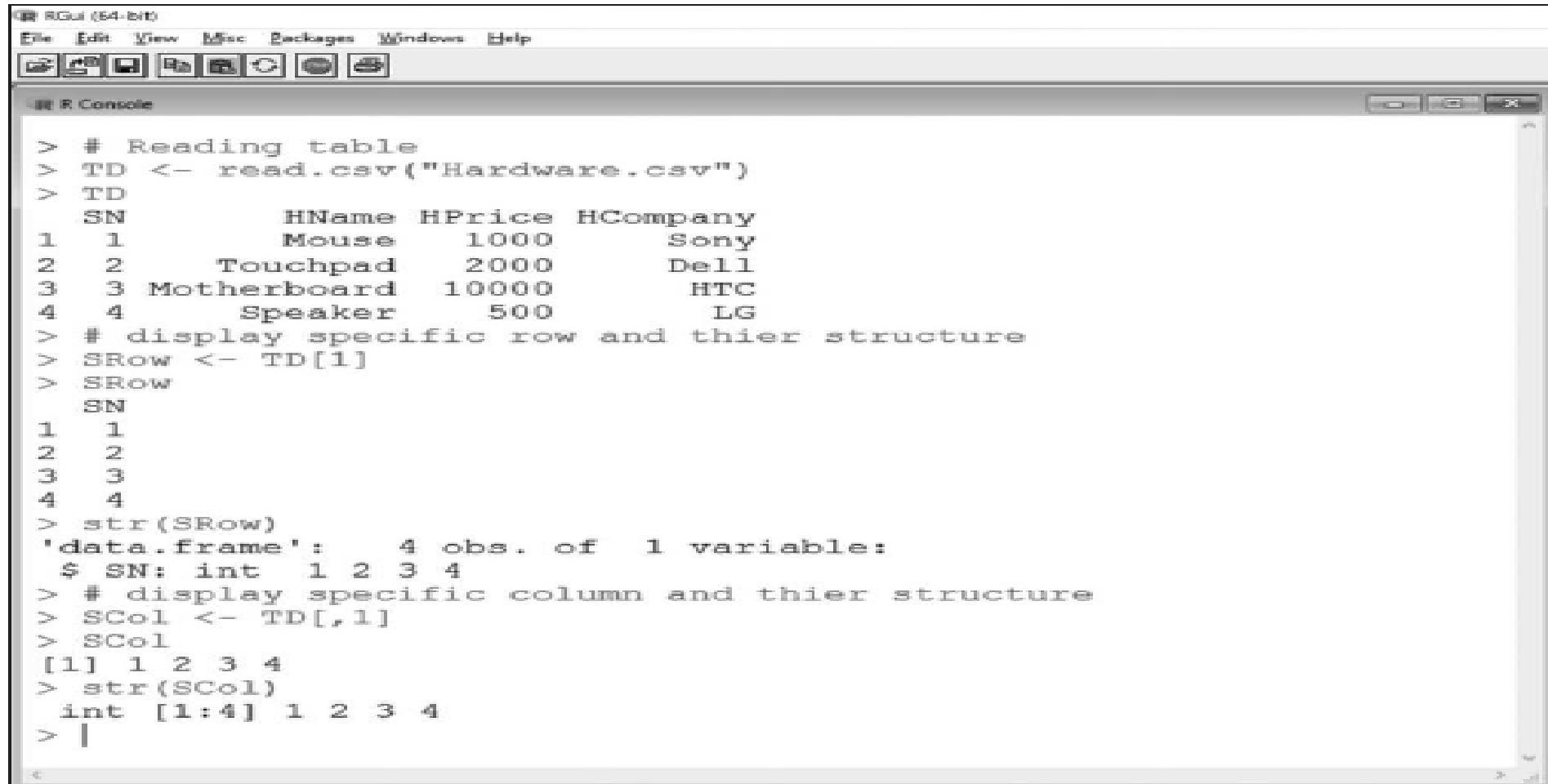
- Analytical data processing sometimes may require specific rows and columns of a dataset.

TABLE 3.5 Commands for accessing specific rows and columns of a dataset

<i>Commands</i>	<i>Command Arguments</i>	<i>Description</i>
Tablename [n]	<i>n</i> is a numeric value.	The command displays the rows according to the given value of argument <i>n</i> of the table.
Tablename [, n]	<i>n</i> is a numeric value.	The command displays the columns according to the given value of argument <i>n</i> of the table.

The following example reads a table, 'Hardware.csv' into object, 'TD' on the R workspace. The TD[1] and TD[, 1] commands displays rows and columns (Figure 3.16).

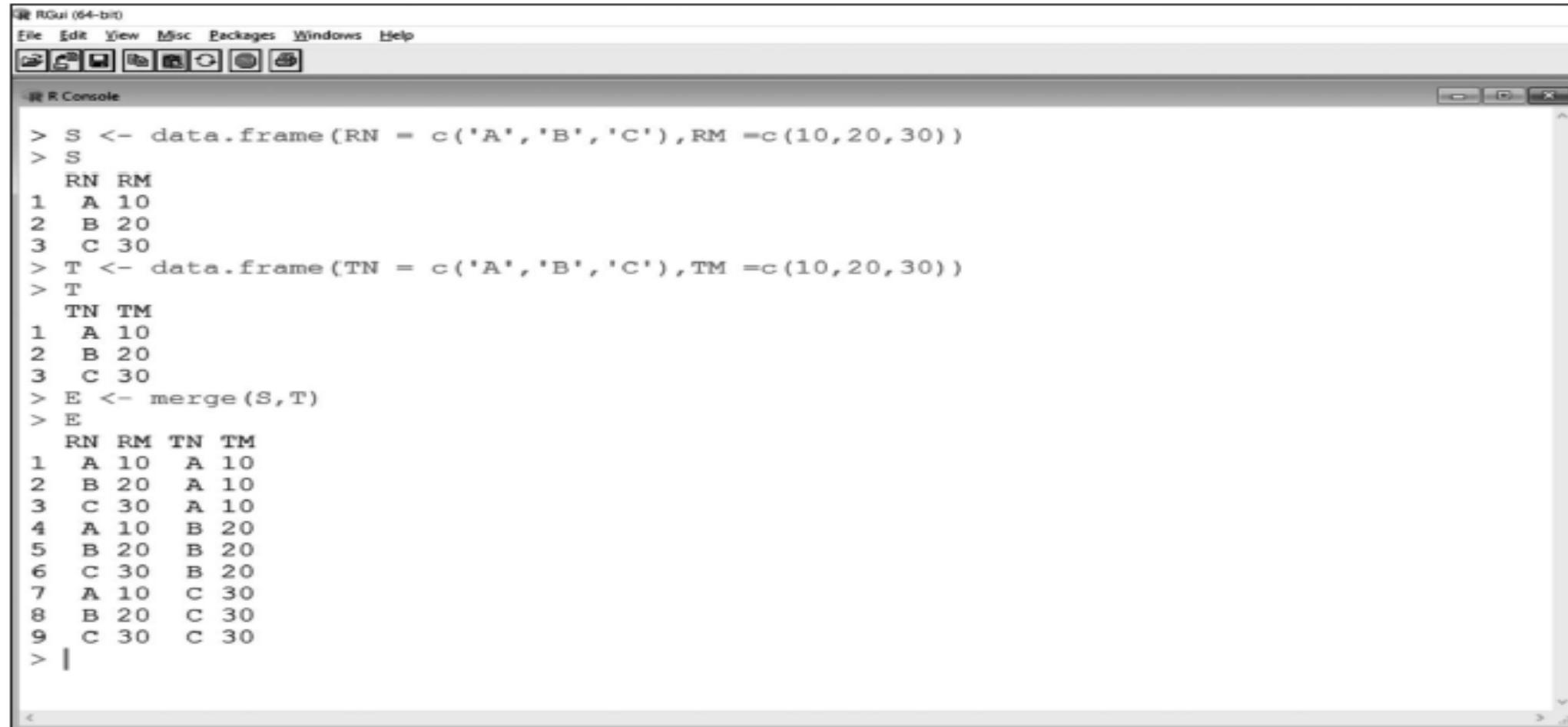
Merging Data



```
RGui (64-bit)
File Edit View Misc Packages Windows Help
R Console

> # Reading table
> TD <- read.csv("Hardware.csv")
> TD
  SN      HName HPrice HCompany
1  1       Mouse   1000    Sony
2  2     Touchpad   2000    Dell
3  3 Motherboard 10000    HTC
4  4     Speaker    500     LG
> # display specific row and thier structure
> SRow <- TD[1]
> SRow
  SN
1  1
2  2
3  3
4  4
> str(SRow)
'data.frame': 4 obs. of 1 variable:
 $ SN: int 1 2 3 4
> # display specific column and thier structure
> SCol <- TD[,1]
> SCol
[1] 1 2 3 4
> str(SCol)
  int [1:4] 1 2 3 4
> |
```

SC 2.16: Conditional manipulation of a dataset



RGui (64-bit)

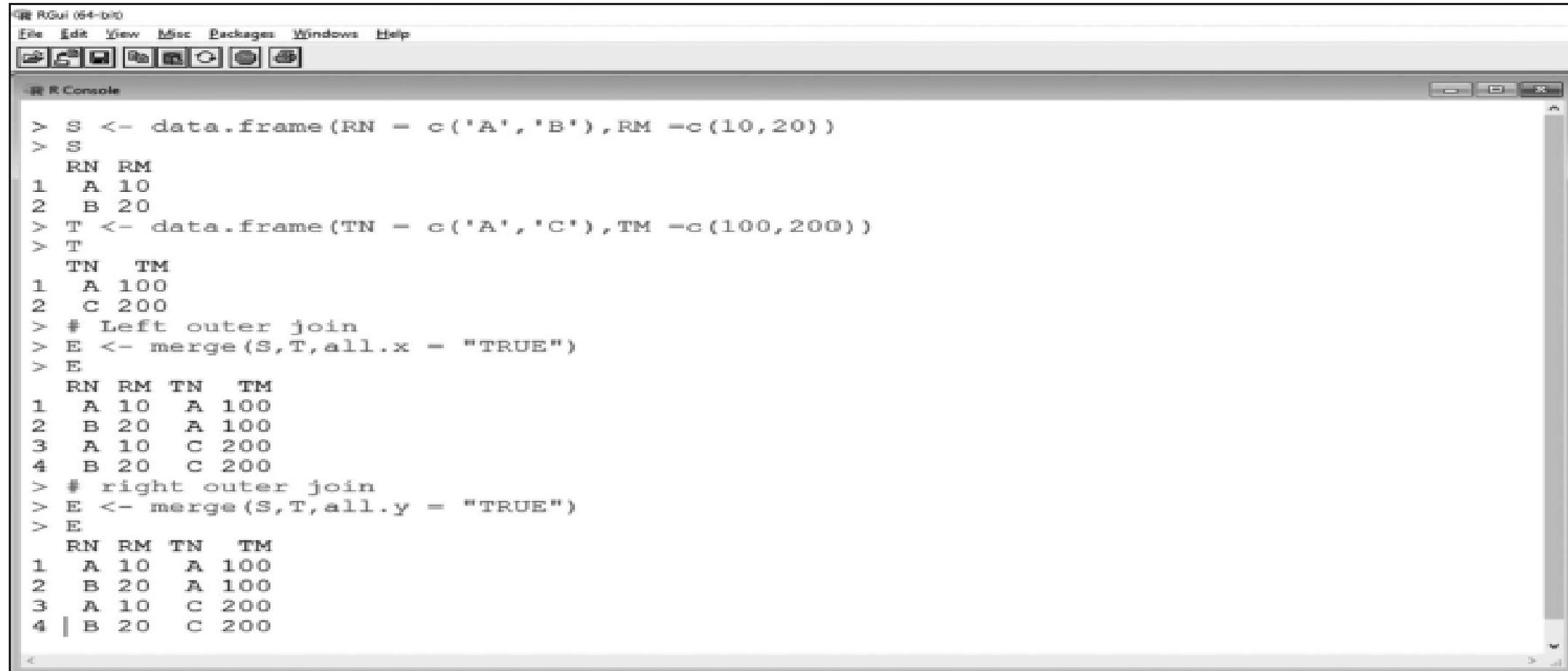
File Edit View Misc Packages Windows Help

R Console

```
> S <- data.frame(RN = c('A', 'B', 'C'), RM =c(10,20,30))
> S
  RN RM
1  A 10
2  B 20
3  C 30
> T <- data.frame(TN = c('A', 'B', 'C'), TM =c(10,20,30))
> T
  TN TM
1  A 10
2  B 20
3  C 30
> E <- merge(S, T)
> E
  RN RM TN TM
1  A 10  A 10
2  B 20  A 10
3  C 30  A 10
4  A 10  B 20
5  B 20  B 20
6  C 30  B 20
7  A 10  C 30
8  B 20  C 30
9  C 30  C 30
> |
```

FIGURE 3.17 Merging data

FIGURE 3.18 Merging data



The screenshot shows an RStudio interface with the title bar "RStudio 1.4.1090". The menu bar includes "File", "Edit", "View", "Misc", "Packages", "Windows", and "Help". Below the menu is a toolbar with various icons. The main area is titled "R Console". The console output shows the following R code and its execution results:

```
RStudio (64-bit)
File Edit View Misc Packages Windows Help
RStudio 1.4.1090

> S <- data.frame(RN = c("A", "B"), RM =c(10, 20) )
> S
RN RM
1 A 10
2 B 20
> T <- data.frame(TN = c("A", "C"), TM =c(100, 200) )
> T
TN TM
1 A 100
2 C 200
# Left outer join
> E <- merge(S,T,all.x = "TRUE")
> E
RN RM TN TM
1 A 10 A 100
2 B 20 A 100
3 A 10 C 200
4 B 20 C 200
# right outer join
> E <- merge(S,T,all.y = "TRUE")
> E
RN RM TN TM
1 A 10 A 100
2 B 20 A 100
3 A 10 C 200
4 | B 20 C 200
```

FIGURE 3.18 Merging data using join condition

Loading and Handling Data in R

Challenges of Analytical Data Processing

Expression, Variables and Functions

Missing Values Treatment in R

Using the 'as' Operator to Change the Structure of Data,

Vectors

Matrices

Factors

List

Aggregating and Group Processing of a Variable,

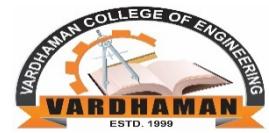
Simple Analysis Using R,

Methods for Reading Data,

Comparison of R GUIs for Data Input

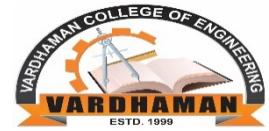
Objective 1	
1	<p>Let us create a matrix which is 3 rows by 4 columns and set all its elements to 1.</p> <pre>> matrix (1, 3, 4) [, 1] [, 2] [, 3] [, 4] [1,] 1 1 1 1 [2,] 1 1 1 1 [3,] 1 1 1 1</pre>
2	<p>Use a vector to create an array, 3 rows high and 3 columns wide.</p> <p>Begin by creating a vector that has elements from 10 to 90 with an interval of 10.</p> <p>Validate by printing the value of vector a.</p> <p>Call the matrix function with vector, 'a' the number of rows and the number of columns.</p> <pre>0a <- seq(10, 90, by = 10) > a [1] 10 20 30 40 50 60 70 80 90 > matrix (a, 3, 3) [, 1] [, 2] [, 3] [1,] 10 40 70 [2,] 20 50 80 [3,] 30 60 90</pre>

Matrices : Matrices Are Nothing But Two-dimensional Arrays.



Objective 2		
1	Re-shape the vector itself into an array using the dim function. Begin by creating a vector that has elements from 10 to 90 with an interval of 10	<pre>a <- seq(10, 90, by = 10)</pre>
2	Validate by printing the value of vector, a.	<pre>> a [1] 10 20 30 40 50 60 70 80 90</pre>
3	Assign new dimensions to vector, a by passing a vector having 3 rows and 3 columns (c(3, 3)). Print the values of vector, a. You will notice that the values have shifted to form 3 rows by 3 columns. The vector is no longer one dimensional. It has been converted into a two-dimensional matrix that is 3 rows high and 3 columns wide.	<pre>> dim(a) <- c(3, 3) > a [, 1] [, 2] [, 3] [1,] 10 40 70 [2,] 20 50 80 [3,] 30 60 90</pre>

Matrices : Matrices Are Nothing But Two-dimensional Arrays.

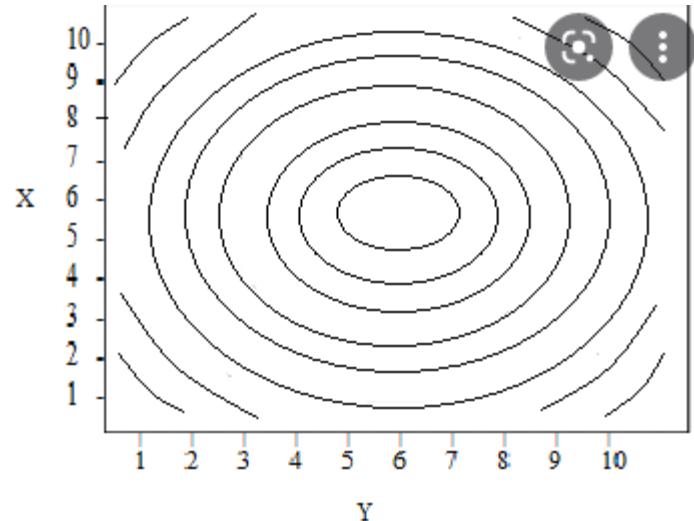


Matrix Access		
1	<p>Access the elements of a 3 *4 matrix. Create a matrix, 'mat', 3 rows high and 4 columns wide using a vector.</p>	<pre>> x <- 1:12 > x [1] 1 2 3 4 5 6 7 8 9 10 11 12 > mat <- matrix(x, 3, 4) > mat [, 1] [, 2] [, 3] [, 4] [1,] 1 4 7 10 [2,] 2 5 8 11 [3,] 3 6 9 12</pre>
2	<p>Access the element present in the second row and third column of the matrix, 'mat'.</p>	<pre>> mat [2, 3] [1] 8</pre>

Matrix Access		
1	<p>Access the third row of an existing matrix. Let us begin by printing the values of an existing matrix, 'mat'</p> <p>To access the third row of the matrix, simply provide the row number and omit the column number.</p> <p>To access the second column of the matrix, simply provide the column number and omit the row number.</p>	<pre>> mat [, 1] [, 2] [, 3] [, 4] [1,] 1 4 7 10 [2,] 2 5 8 11 [3,] 3 6 9 12</pre> <pre>> mat [3,] [1] 3 6 9 12</pre> <pre>> mat[, 2] [1] 4 5 6</pre>
2	<p>To access the second and third columns of the matrix, simply provide the column numbers and omit the row number.</p>	<pre>> mat[, 2:3] [, 1] [, 2] [1,] 4 7 [2,] 5 8 [3,] 6 9</pre>

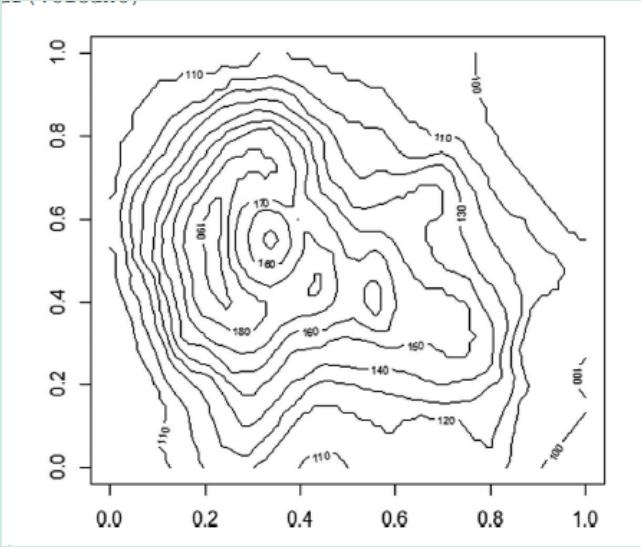
Create a contour plot

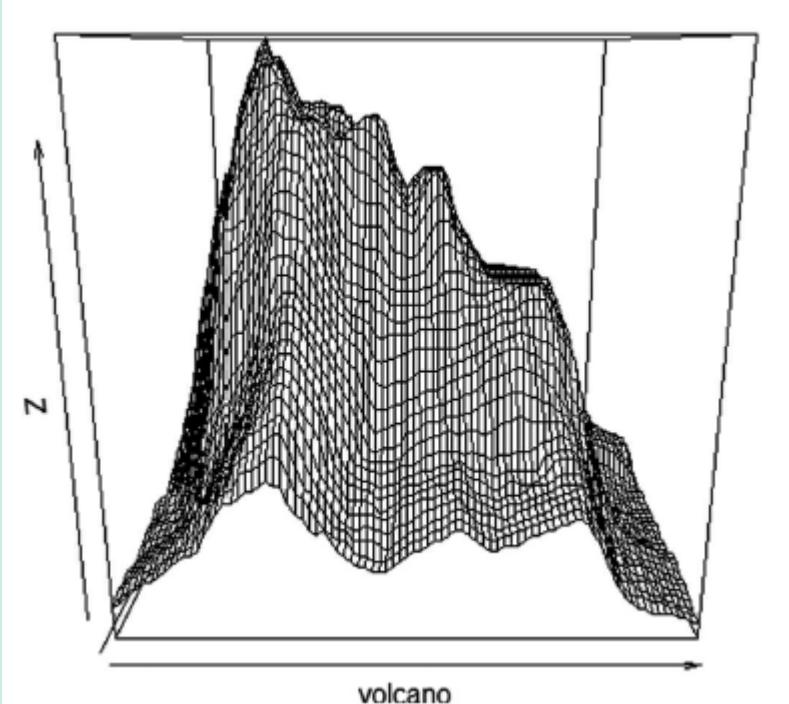
- Contour plots (sometimes called Level Plots) are **a way to show a three-dimensional surface on a two-dimensional plane.**
- It graphs two predictor variables X Y on the y-axis and a response variable Z as contours.
- These contours are sometimes called z-slices or iso-response values

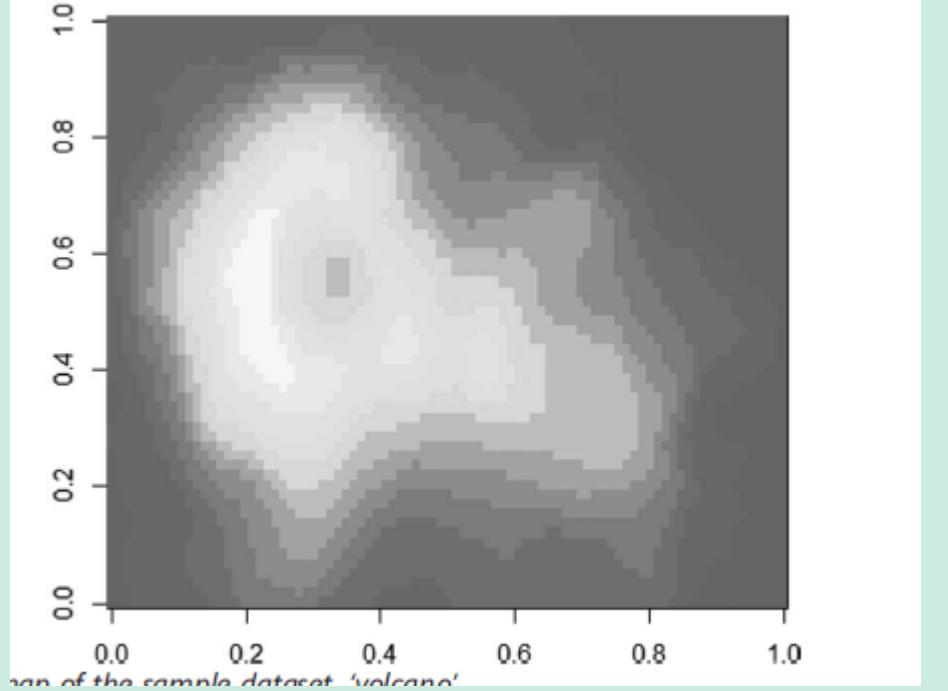


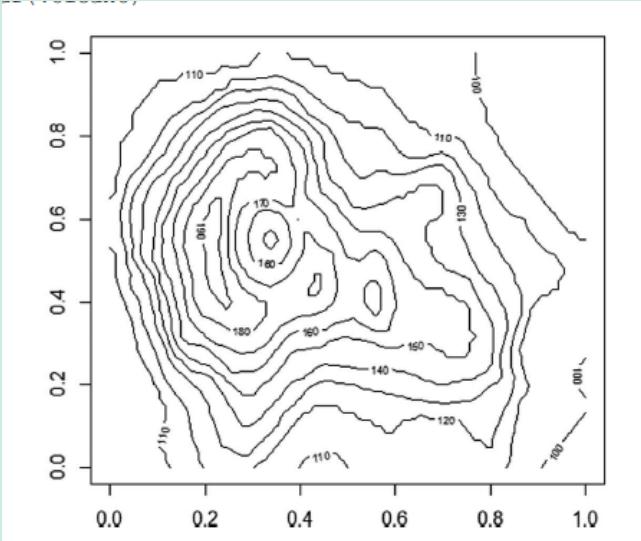
S.N O	Function	Example
	<p>Create a matrix, 'mat' which is 9 rows high and 9 columns wide and assign the value '1' to all its elements.</p>	<pre data-bbox="1462 447 2422 1422">mat <- matrix(1, 9, 9) > mat [, 1] [, 2] [, 3] [, 4] [, 5] [, 6] [, 7] [, 8] [, 9] [1,] 1 1 1 1 1 1 1 1 1 [2,] 1 1 1 1 1 1 1 1 1 [3,] 1 1 1 1 1 1 1 1 1 [4,] 1 1 1 1 1 1 1 1 1 [5,] 1 1 1 1 1 1 1 1 1 [6,] 1 1 1 1 1 1 1 1 1 [7,] 1 1 1 1 1 1 1 1 1</pre>

S.N O	Function	Example
	Create a matrix, 'mat' which is 9 rows high and 9 columns wide and assign the value '1' to all its elements.	<pre>mat <- matrix(1, 9, 9) > mat [, 1] [, 2] [, 3] [, 4] [, 5] [, 6] [, 7] [, 8] [, 9] [1,] 1 1 1 1 1 1 1 1 1 [2,] 1 1 1 1 1 1 1 1 1 [3,] 1 1 0 1 1 1 1 1 1 [4,] 1 1 1 1 1 1 1 1 1 [5,] 1 1 1 1 1 1 1 1 1 [6,] 1 1 1 1 1 1 1 1 1 [7,] 1 1 1 1 1 1 1 1 1</pre>

Dataset	Output
<p>R includes some sample data sets. One of these is 'volcano', which is a 3D map of a dormant New Zealand volcano. Create a contour map of the volcano dataset (Figure 3.7).</p> <pre data-bbox="105 630 532 673">> contour(volcano)</pre>	 <p>A contour plot showing the elevation of a dormant New Zealand volcano. The x-axis and y-axis both range from 0.0 to 1.0. The plot features a dense network of contour lines, with values labeled along several of them. The highest peak is at approximately (0.3, 0.5) with a value of 260. Other labeled values include 240, 220, 200, 180, 160, 140, 120, 110, 100, 90, 80, 70, 60, 50, 40, 30, 20, 10, and 0.</p>

Dataset	Output
<p>Let us create a 3D perspective map of the sample data set, 'volcano' (Figure 3.8). > persp(volcano)</p>	 A 3D perspective plot of the 'volcano' dataset. The vertical axis is labeled 'Z', the horizontal axis is labeled 'volcano', and the depth axis is implied by the perspective grid. The surface shows a central peak with a complex, multi-layered structure, characteristic of a volcano's profile.

Dataset	Output
<p>Create a heat map of the sample dataset, 'volcano' (Figure 3.9). > <code>image(volcano)</code></p>	 <p>A grayscale heatmap visualization of the 'volcano' dataset. The plot shows a central bright white peak surrounded by concentric rings of gray, transitioning to black at the edges. The x-axis and y-axis both range from 0.0 to 1.0, with major ticks at 0.0, 0.2, 0.4, 0.6, 0.8, and 1.0. The plot is titled "heat map of the sample dataset, 'volcano'" at the bottom.</p>

Dataset	Output
<p>R includes some sample data sets. One of these is 'volcano', which is a 3D map of a dormant New Zealand volcano. Create a contour map of the volcano dataset (Figure 3.7).</p> <pre data-bbox="105 630 532 673">> contour(volcano)</pre>	 <p>A contour plot showing the elevation of a dormant New Zealand volcano. The x-axis and y-axis both range from 0.0 to 1.0. The plot features a dense network of contour lines, with values labeled along several of them. The highest peak is at approximately (0.3, 0.5) with a value of 260. Other labeled values include 220, 200, 180, 160, 140, 120, 110, 100, 90, 80, and 70. The contours are irregular, reflecting the shape of the volcano.</p>

Factors :

```
> HouseColor <- c('red', 'green', 'blue', 'yellow', red', 'green', 'blue', 'blue')
> types <- factor(HouseColor)
> HouseColor
[1] "red" "green" "blue" "yellow" "red" "green" "blue" "blue"
> print(HouseColor)
[1] "red" "green" "blue" "yellow" "red" "green" "blue" "blue"
> print (types)
[1] red green blue yellow red green blue blue
```

Levels: blue green red yellow

Levels denotes the unique values. The above has four distinct values such as 'blue', 'green', 'red' and 'yellow'.

```
> as.integer(types)
[1] 3 2 1 4 3 2 1 1
```

The above output is explained as given below.

1 is the number assigned to blue. 2 is the number assigned to green. 3 is the number assigned to red. 4 is the number assigned to yellow.

```
> levels(types) [1] "blue" "green" "red" "yellow"
```