

Experiment No. 1

Aim : Installation of R and RStudio

Theory :

R is a programming language and free software environment for statistical computing and graphics supported by the R Foundation for Statistical Computing. The R language is widely used among statisticians and data miners for developing statistical software and data analysis.

R provides a wide variety of statistical (linear and nonlinear modelling, classical statistical tests, time-series analysis, classification, clustering, ...) and graphical techniques, and is highly extensible. The S language is often the vehicle of choice for research in statistical methodology, and R provides an Open Source route to participation in that activity.

One of R's strengths is the ease with which well-designed publication-quality plots can be produced, including mathematical symbols and formulae where needed. Great care has been taken over the defaults for the minor design choices in graphics, but the user retains full control.

R is an integrated suite of software facilities for data manipulation, calculation and graphical display. The term "environment" is intended to characterize it as a fully planned and coherent system, rather than an incremental accretion of very specific and inflexible tools, as is frequently the case with other data analysis software.

Steps to install R :

1. Open an internet browser and go to www.r-project.org.
2. Click the "download R" link in the middle of the page under "Getting Started."
3. Select a CRAN location (a mirror site) and click the corresponding link.
4. Click on the "Download R for Windows" link at the top of the page.
5. Click on the "install R for the first time" link at the top of the page.
6. Click "Download R for Windows" and save the executable file somewhere on your computer. Run the .exe file and follow the installation instructions.
7. Now that R is installed, you need to download and install RStudio.

Steps to install RStudio :

1. Go to www.rstudio.com and click on the "Download RStudio" button.
2. Click on "Download RStudio Desktop."
3. Click on the version recommended for your system, or the latest Windows version, and save the executable file. Run the .exe file and follow the installation instructions.

Output :

The screenshot shows a terminal window with a green and blue floral background. The terminal output is as follows:

```
> sudo apt install r-base libappparm1
[sudo] password for akuma:
Reading package lists... Done
Reading state information... Done
libappparm1 is already the newest version (3.4.4-1ubuntu1).
r-base is already the newest version (3.4.4-1ubuntu1).
libappparm1 is already the newest version (2.12-4ubuntu5.1).
0 upgraded, 0 newly installed, 0 to remove and 619 not upgraded.
> R --version
R version 3.4.4 (2018-03-15) -- "Someone to Lean On"
Copyright (C) 2018 The R Foundation for Statistical Computing
Platform: x86_64-pc-linux-gnu (64-bit)

R is free software and comes with ABSOLUTELY NO WARRANTY.
You are welcome to redistribute it under the terms of the
GNU General Public License, either version 2 or 3.
For more information about these matters see
http://www.gnu.org/licenses/.

> cd ~/Downloads/
> wget -c http://download2.rstudio.org/rstudio-server-0.97.336-i386.deb -O rstudio.deb
--2020-03-31 22:24:25.  http://download2.rstudio.org/rstudio-server-0.97.336-i386.deb
Resolving download2.rstudio.org (download2.rstudio.org)... 13.35.214.118, 13.35.214.128, 13.35.214.74, ...
Connecting to download2.rstudio.org (download2.rstudio.org)[13.35.214.118]:80... connected.
HTTP request sent, awaiting response... 200 OK
Length: 17840264 (17M) [application/x-deb]
Saving to: 'rstudio.deb'

rstudio.deb          100%[=====] 17.01M  418KB/s   in 33s
2020-03-31 22:25:00 (521 KB/s) - 'rstudio.deb' saved [17840264/17840264]

> sudo dpkg -i rstudio.deb
Selecting previously unselected package rstudio-server:i386.
(Reading database ... 7%
```

Conclusion : Thus we have installed R and Rstudio.

Experiment No. 2

Aim : Basic Functionality of R - variables, basic arithmetic and help commands.

Theory :

A variable provides us with named storage that our programs can manipulate. A variable in R can store an atomic vector, group of atomic vectors or a combination of many R objects. A valid variable name consists of letters, numbers and the dot or underline characters. The variable name starts with a letter or the dot not followed by a number.

We can name variables, vectors, matrices, functions, and subroutines. You are not restricted by the names you use, as long as they are not reserved for existing functions (like the function *sum()*) in the R package you are using. R is case sensitive, so the variable *x* is different from the variable *X*.

The *help()* function and *? help* operator 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. To access documentation for the standard *lm* (linear model) function, for example, enter the command *help(lm)* or *help("lm")*, or *?lm* or *?"lm"* (i.e., the quotes are optional).

To access help for a function in a package that's *not* currently loaded, specify in addition the name of the package: For example, to obtain documentation for the *rlm()* (robust linear model) function in the MASS package, *help(rlm, package="MASS")*.

Output :

Code :

```
# BASIC ARITHMETIC OPERATORS
2-5 # subtraction
6/3 # division
3+2*5 # note order of operations exists
(3+2)*5 # if you need, force operations using
# redundant parentheses
4^3 # raise to a power is ^
exp(4) # e^4 = 54.598 (give or take)
log(2.742) # natural log of 2.74
log10(1000) # common log of 1000
pi # 3.14159...
```

```
# VARIABLE NAMES
x <- 5
x # 1 way to print out contents of a variable
```

```
var1 <- 7/2
print(var1) # another way to print contents
valid.variable.name <- 18.6
# you can even have long variable names
# if you are of that kind of weirdness
valid.variable.name
```

```
a <- 3 # This stores the value 3 into variable a
a
sqrt(a) # square root of a (which is 3 at this moment)
a^4 # a raised to the fourth power
log(a) # natural log of a
log10(a) # common log of a
exp(a) # e to the power a
tan(a) # tangent of a
b <- pi # pie, the number approx 3.14159
b
```

VECTORS

```
x <- c(2,3,5,1,4,4) # create a row vector with those 6 numbers
# and then call the vector x
x
sum(x) # sums the elements of vector x
mean(x) # finds the mean of vector x
sd(x) # finds standard deviation of vector x
median(x) # finds the median of x
sqrt(x) # finds square root of every element of x
x^2 # finds the square of every element of x
seq(1,10) # makes a vector 10 long, from 1 to 10
seq(1,10, 2) # makes a vector from 1 to 10, starting with 1, then skipping 1 ending up with
only odds from 1 to 10
seq(1:10) # same as seq(1,10)
seq(1,10,by=2) # same as seq(1,10,2)
y <- c(1:7) # create a row vector with elements 1 through 6
y
z <- 1:7 # same as y--use when you increment by 1
z
w <- c(1:12,0,-6) # use the c( when you increment by 1 followed by more numbers not in
sequence
w
```

```
# SOME OPERATIONS WITH NUMERICAL VECTORS AND LOGICAL VECTORS
x <- c(-5,0,7,-6,14,27) # data vector x
x[1] # prints first element of x
x[length(x)] # prints last element of x
i<-3
```

```

x[i] # the ith entry if 1<=i<=n, NA if i>n, all but the ith if -n<=i<=-1, an error if i < -n, and an
empty vector if i=0
x[c(2,3)] # the second and 3rd entries
x[-c(2,3)] # all but the second and third entries
x[i] <- 5 # assign a value of 5 to the first entry; also can use x[i] = 5
x[c(1,4)] <- c(2,3) # assign values to the first and fourth entries
y
x[indices] <- y # assign to those indices indicated by the values
# of INDICES: if y is not long enough, values
# are recycled; if y is too long, just its initial
# values are used and a warning is issued
x
x < 3 # vector with length n of TRUE or FALSE # depending if x[i] < 3 which(x<3) which
indices correspond to the TRUE values of x<3
x[x<3] # the x values when x<3 is TRUE -- same as x[which(x<3)]

```

```

# INTRO TO GRAPHING
# input a vector x with data
x <- c(2,4,4,6,6,5,5,7,3,7,3,8,9,7,9,6,4,3,4,4,6,2,2,1,2,4,6,6,8)
# input vector y of consecutive numbers
# 1 through 29
y <- c(1:29)
# make a histogram of x
hist(x)
# make a scatter plot of y vs x
plot(x,y)
# make a time plot of x, connecting dots
plot(x, type="b")
# make a histogram of a sampling of 20 from
# the binomial distribution B(12,.4)
y <- rbinom(20, 12, .4)
hist(y)
# make another binomial histogram, but using
a
# sample of 200
y <- rbinom(200, 12, .4)
hist(y)
# make a boxplot of x data
boxplot(x)
# make side by side boxplots of x and y data
boxplot(x,y)

```

```

#LOGICAL OPERATORS
x <- 1:5 # x is a row vector = [1,2,3,4,5]
test <- (x <4) # test is a row vector (logical <)
# = [TRUE TRUE TRUE FALSE FALSE]

```

```

test <- (x==3) # test is a row vector (logical =)
# = [FALSE FALSE TRUE FALSE FALSE]
test <- (x>1 & x<4) # test is a row vector (logical AND)
# = [FALSE TRUE TRUE FALSE FALSE]
test <- (x>4 | x==2) # test is a row vector (logical OR)
# = [FALSE TRUE FALSE FALSE TRUE]
test<- (x != 4) # test is a row vector (logical NOT =)
# = [TRUE TRUE TRUE FALSE TRUE]
test <- (x %in% c(2,4)) # test is a row vector
# testing whether the entry is
# a 2 or a 4
# = [FALSE TRUE FALSE TRUE FALSE]

```

```

# THE NORMAL DISTRIBUTION
pnorm(1.43) # finds P(Z < 1.43)
pnorm(129,100,16) # finds P(X < 129) where x comes from
# normal with mean 100 and std dev 16
qnorm(.994) # finds the .994 quantile of the standard normal
qnorm(.95,100,16) # finds the .95 quantile of N(100,16) distribution

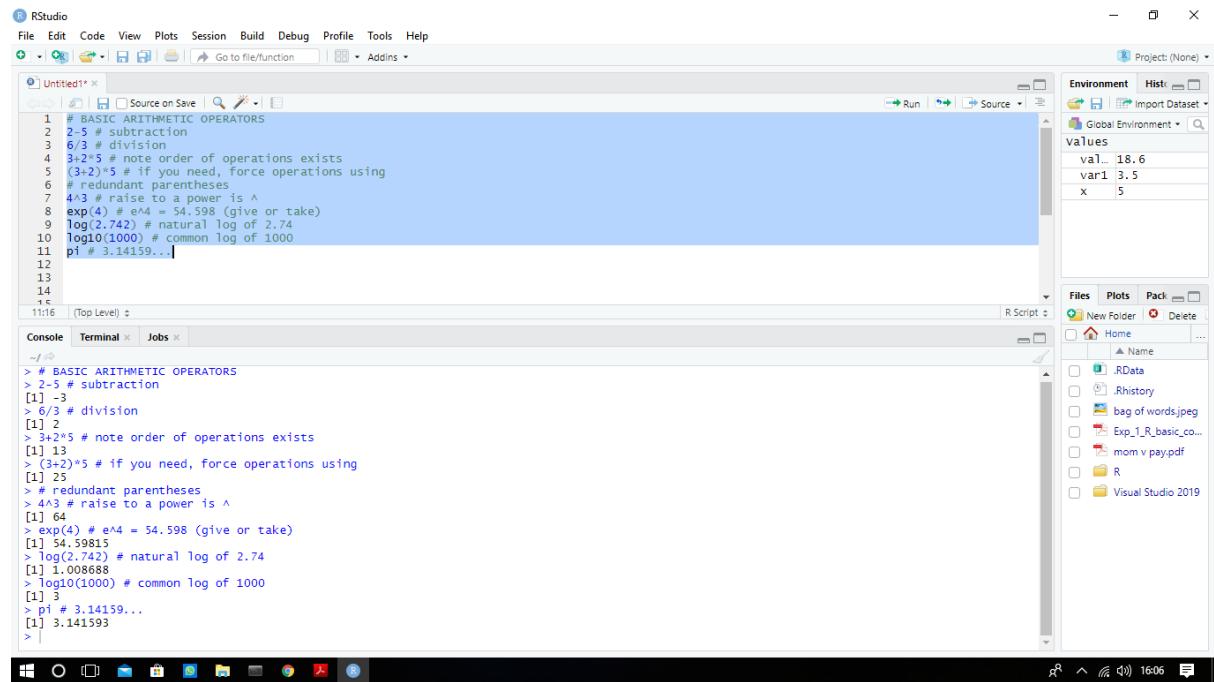
```

```

# THINGS TO KNOW
x <- runif(30,0,1) # creates a vector 30 long
x # of random values from U(0,1)
y <- c() # before you can use a vector in a loop
for (i in 1:10) { # you must "save space" for it by
  y[i] <- rbinom(1,3,.2) # typing the y <- c() command
}
y
s <- c(1:5) # for R you can use <- or =
s # for a variable assignment
t = c(2:6)
t
help(mean)

```

Screenshots :



RStudio

File Edit Code View Plots Session Build Debug Profile Tools Help

Untitled1

```
1 # BASIC ARITHMETIC OPERATORS
2 2-5 # subtraction
3 6/3 # division
4 3+2*5 # note order of operations exists
5 (3-2)^5 # if you need, force operations using
6 # redundant parentheses
7 4^3 # raise to a power is ^
8 exp(4) # e^4 = 54.598 (give or take)
9 log(2.742) # natural log of 2.74
10 log10(1000) # common log of 1000
11 pi # 3.14159...
12
13
14
15
16
```

Console Terminal Jobs

Environment Hist... Project: (None)

values

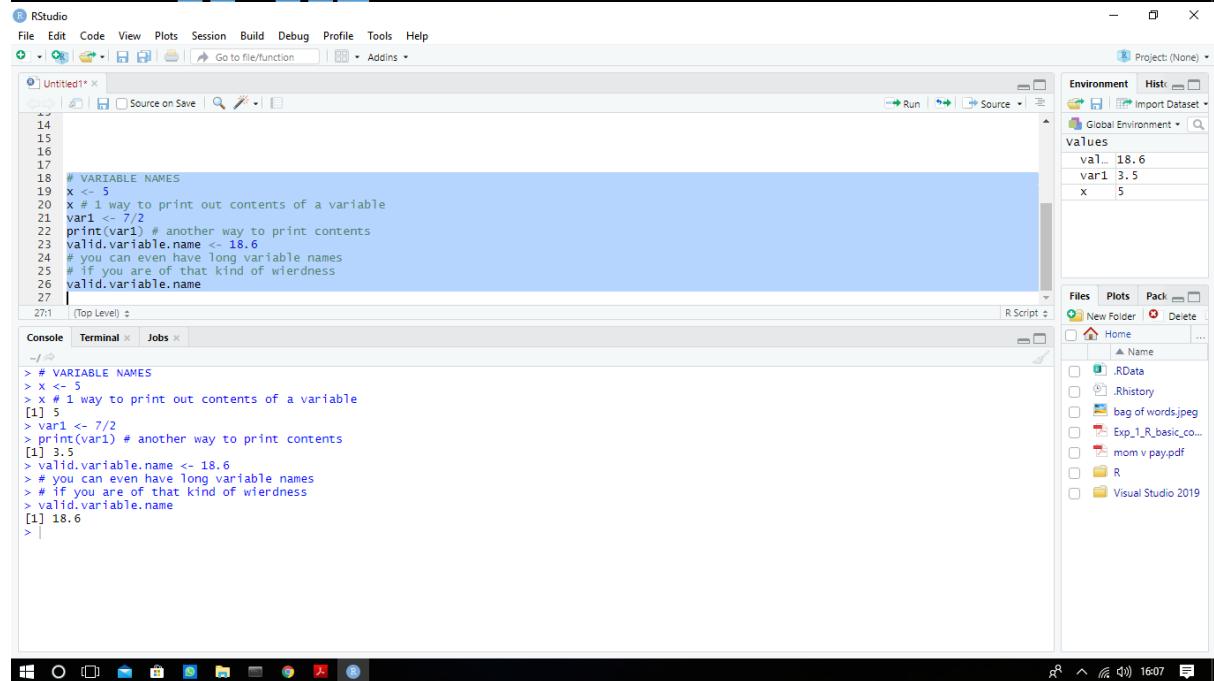
val_	18.6
var1	3.5
x	5

Files Plots Pack

New Folder Delete

Home RData Rhistory bag of words.jpeg Exp_1_R_basic.co... mom v pay.pdf R Visual Studio 2019

1606



RStudio

File Edit Code View Plots Session Build Debug Profile Tools Help

Untitled1

```
1.4
1.5
1.6
1.7
18 # VARIABLE NAMES
19 x <- 5
20 x # 1 way to print out contents of a variable
21 var1 <- 7/2
22 print(var1) # another way to print contents
23 valid.variable.name <- 18.6
24 # you can even have long variable names
25 # if you are of that kind of weirdness
26 valid.variable.name
27 | (Top Level) :|
```

Console Terminal Jobs

Environment Hist... Project: (None)

values

val_	18.6
var1	3.5
x	5

Files Plots Pack

New Folder Delete

Home RData Rhistory bag of words.jpeg Exp_1_R_basic.co... mom v pay.pdf R Visual Studio 2019

1607

RStudio

File Edit Code View Plots Session Build Debug Profile Tools Help

Untitled1* | Source on Save | Run | Source | Environment | Hist | Global Environment | Values | Files | Plots | Pack | Project: (None)

```

28
29
30 a <- 3 # This stores the value 3 into variable a
31 a
32 sqrt(a) # square root of a (which is 3 at this moment)
33 a^4 # a raised to the fourth power
34 log(a) # natural log of a
35 log10(a) # common log of a
36 exp(a) # e to the power a
37 tan(a) # tangent of a
38 b <- pi # pie, the number approx 3.14159
39 b
40
41
42
39:2 (Top Level) R Script

```

Console Terminal Jobs

```

> a <- 3 # This stores the value 3 into variable a
> a
[1] 3
> sqrt(a) # square root of a (which is 3 at this moment)
[1] 1.732051
> a^4 # a raised to the fourth power
[1] 81
> log(a) # natural log of a
[1] 1.098612
> log10(a) # common log of a
[1] 0.4771213
> exp(a) # e to the power a
[1] 20.08553
> tan(a) # tangent of a
[1] -0.1425465
> b <- pi # pie, the number approx 3.14159
> b
[1] 3.141593
>

```

16:08

RStudio

File Edit Code View Plots Session Build Debug Profile Tools Help

Untitled1* | Source on Save | Run | Source | Environment | Hist | Global Environment | Values | Files | Plots | Pack | Project: (None)

```

44 # and then call the vector x
45 x
46 sum(x) # sums the elements of vector x
47 mean(x) # finds the mean of vector x
48 sd(x) # finds standard deviation of vector x
49 median(x) # finds the median of x
50 sqrt(x) # finds square root of every element of x
51 x<2 # finds the square of every element from 1 to 10
52 seq(1,10) # makes a vector 10 long, from 1 to 10
53 seq(1,10, 2) # makes a vector from 1 to 10, starting with 1 ending up with only odds from 1 to 10
54 seq(1:10) # same as seq(1,10)
55 seq(1,10,by=2) # same as seq(1,10,2)
56 y <- c(1:7) # create a row vector with elements 1 through 6
57 y
58
62:1 (Top Level) R Script

```

Console Terminal Jobs

```

> x
[1] 1.414214 1.732051 2.236068 1.000000 2.000000 2.000000
> x<2 # finds the square of every element of x
[1] 4 9 25 1 16 16
> seq(1,10) # makes a vector 10 long, from 1 to 10
[1] 1 2 3 4 5 6 7 8 9 10
> seq(1,10, 2) # makes a vector from 1 to 10, starting with 1, then skipping 1 ending up with only odds from 1 to 10
[1] 1 3 5 7 9
> seq(1:10) # same as seq(1,10)
[1] 1 2 3 4 5 6 7 8 9 10
> seq(1,10,by=2) # same as seq(1,10,2)
[1] 1 3 5 7 9
> y <- c(1:7) # create a row vector with elements 1 through 6
> y
[1] 1 2 3 4 5 6 7
> z <- 1:7 # same as y--use when you increment by 1
> z
[1] 1 2 3 4 5 6 7
> w <- c(1:12,0,-6) # use the c() when you increment by 1 followed by more numbers not in sequence
> w
[1] 1 2 3 4 5 6 7 8 9 10 11 12 0 -6
>

```

16:12

RStudio

File Edit Code View Plots Session Build Debug Profile Tools Help

Untitled1* Source on Save Go to file/function Addins

```

73 x[c(1,4)] # all but the second and third entries
74 x[1] <- 5 # assign a value of 5 to the first entry; also can use x[1] = 5
75 x[c(1,4)] <- c(2,3) # assign values to the first and fourth entries
76 x[indices] <- y # assign to those indices indicated by the values
77 # of INDICES: if y is not long enough, values
78 # are recycled; if y is too long, just its initial
79 # values are used and a warning is issued
80 x
81 x < 3 # vector with length n of TRUE or FALSE # depending if x[i] < 3 which(x<3) which indices correspond to the TRUE values of x<3
82 x[x<3] # the x values when x<3 is TRUE -- same as x[which(x<3)]
83
84
85
86
86:1 (Top Level) R Script

```

Console Terminal Jobs

```

> x[c(1,4)] # all but the second and third entries
[1] 7
> x[c(2,3)] # the second and 3rd entries
[1] 0 7
> x[-c(2,3)] # all but the second and third entries
[1] -5 -6 14 27
> x[1] <- 5 # assign a value of 5 to the first entry; also can use x[1] = 5
> x[c(1,4)] <- c(2,3) # assign values to the first and fourth entries
> y
[1] 1 2 3 4 5 6 7
> x[indices] <- y # assign to those indices indicated by the values
Warning message:
In x[indices] <- y :
  number of items to replace is not a multiple of replacement length
> # of INDICES: if y is not long enough, values
> # are recycled; if y is too long, just its initial
> # values are used and a warning is issued
> x
[1] 1 0 5 3 14 27 2
> x < 3 # vector with length n of TRUE or FALSE # depending if x[i] < 3 which(x<3) which indices correspond to the TRUE values of x<3
[1] TRUE FALSE FALSE FALSE FALSE FALSE TRUE

```

Project (None)

Environment History

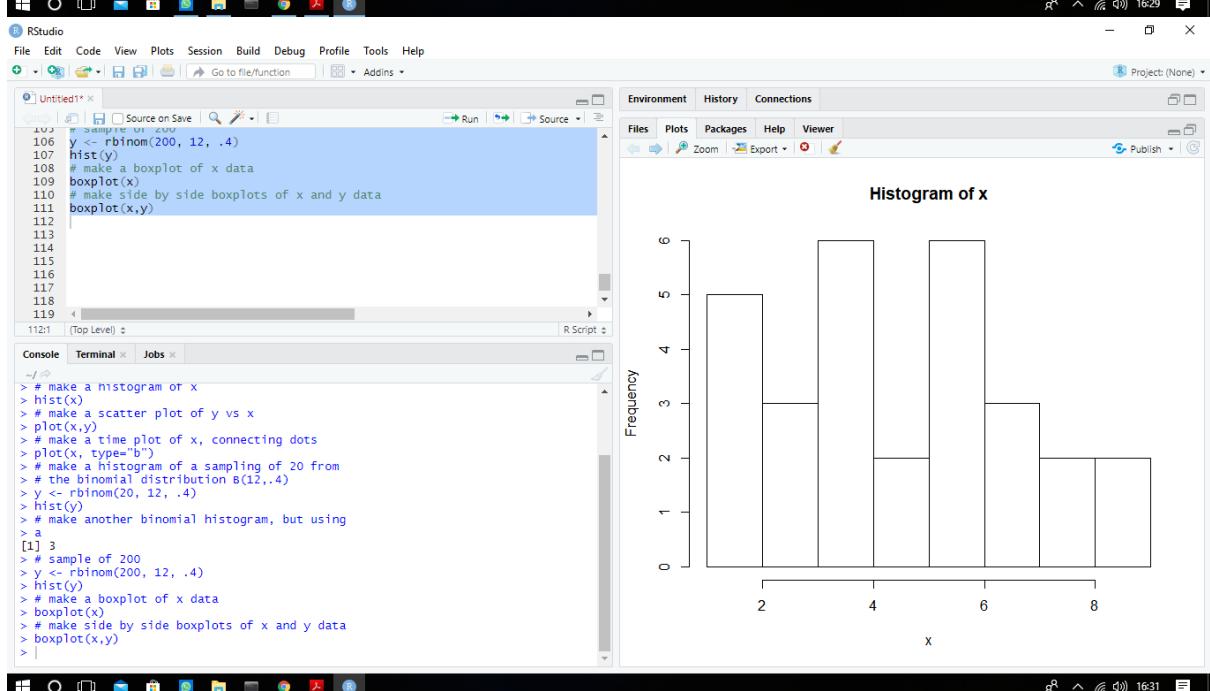
Global Environment b 3.1415926... i 3 ind_ num [1:2]... val_ 18.6 var1 3.5 w num [1:14]... x num [1:7]... y int [1:7]... z int [1:7]...

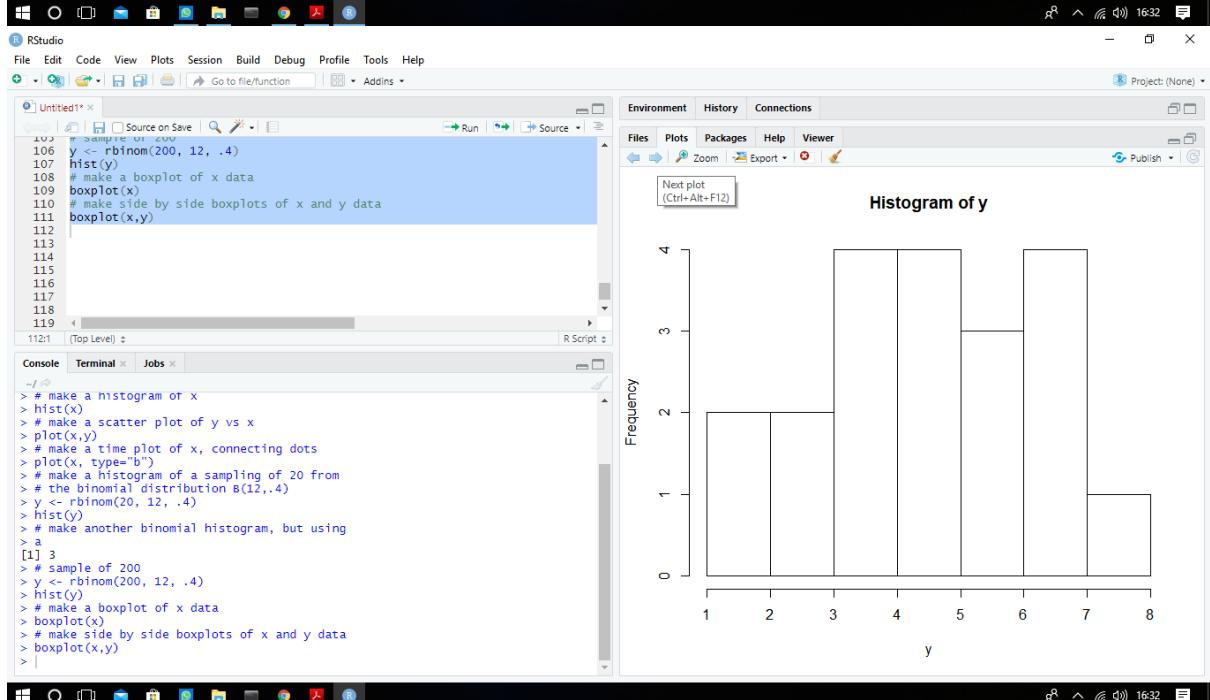
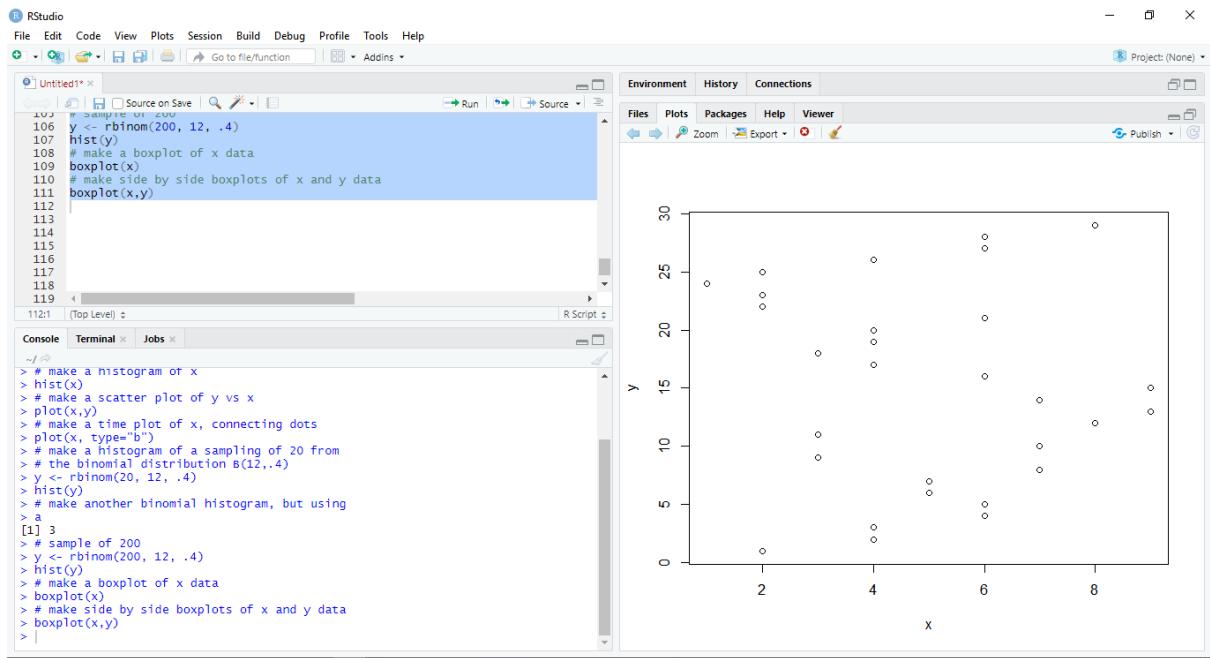
Files Plots Pack

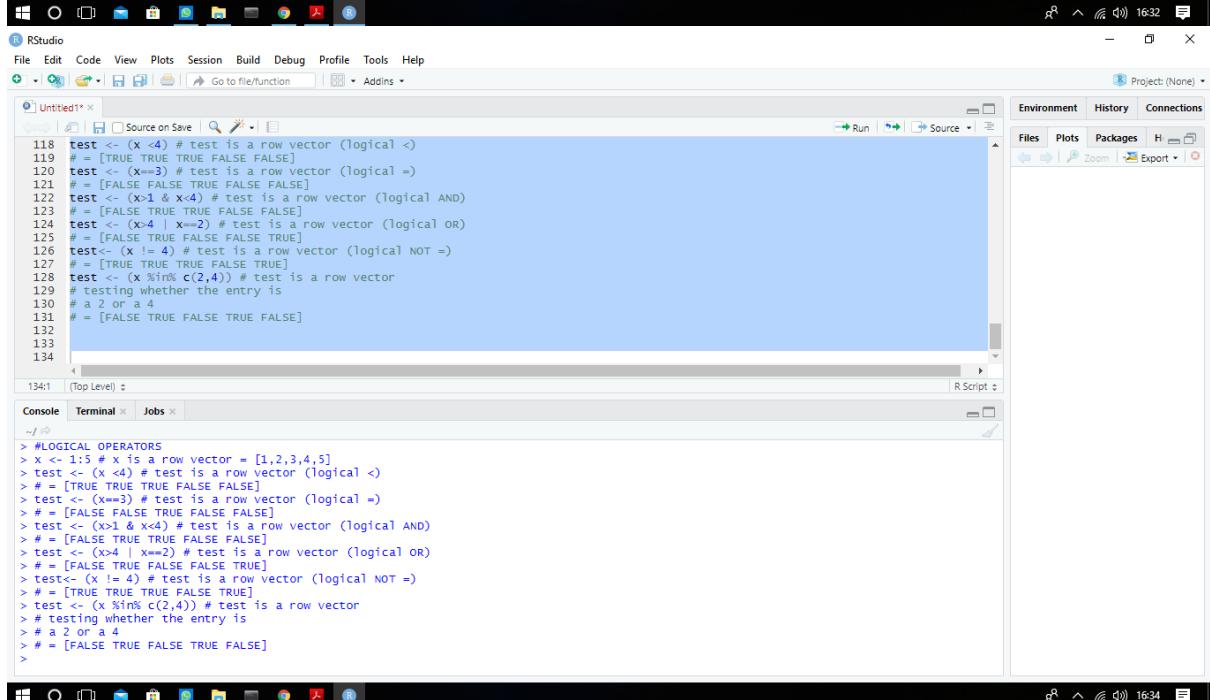
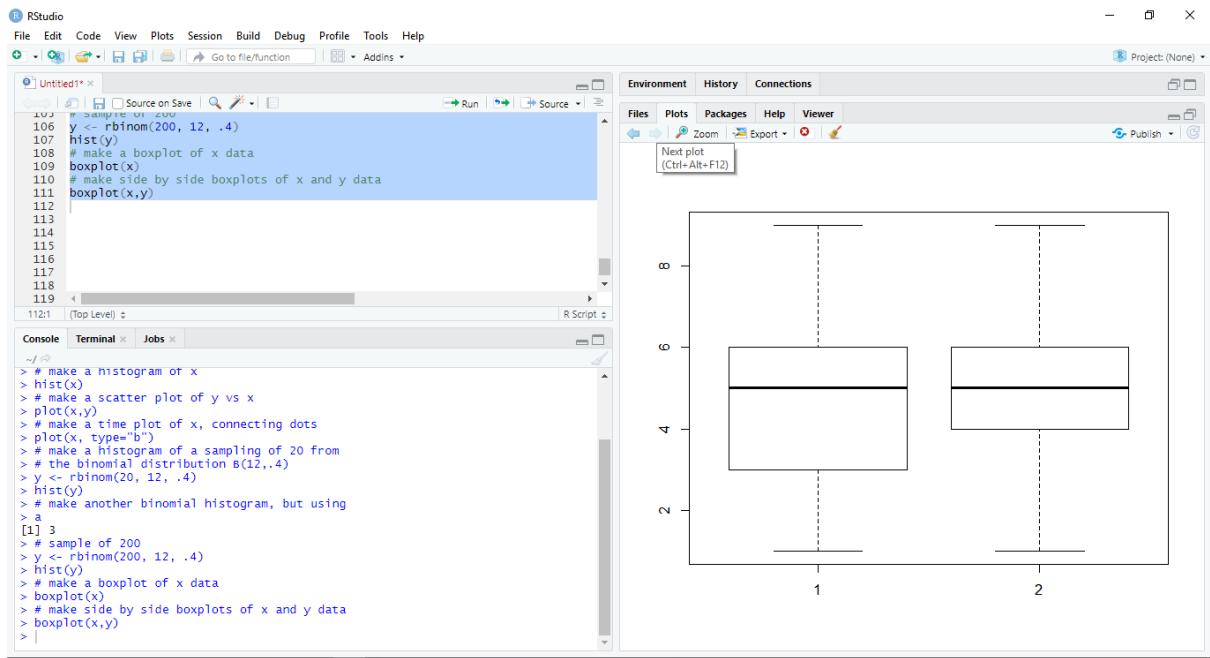
New Folder Home ...

RData Rhistory bag of words.jpeg Exp_1_R_basic.co... mom v.pdf R Visual Studio 2019

Windows Taskbar 1629







RStudio

File Edit Code View Plots Session Build Debug Profile Tools Help

Untitled1* Source on Save Run Source

```
136  
137  
138  
139  
140  
141  
142  
143  
144  
145  
146 # THE NORMAL DISTRIBUTION  
147 pnorm(1.43) # finds P(Z < 1.43)  
148 pnorm(129,100,16) # finds P(X < 129) where x comes from  
# normal with mean 100 and std dev 16  
149 qnorm(.994) # finds the .994 quantile of the standard normal  
150 qnorm(.95,100,16) # finds the .95 quantile of N(100,16) distribution  
151  
152
```

(Top Level) R Script

Console Terminal Jobs

```
> # THE NORMAL DISTRIBUTION  
> pnorm(1.43) # finds P(Z < 1.43)  
[1] 0.932633  
> pnorm(129,100,16) # finds P(X < 129) where x comes from  
# normal with mean 100 and std dev 16  
[1] 0.9650455  
> qnorm(.994) # finds the .994 quantile of the standard normal  
[1] 2.512144  
> qnorm(.95,100,16) # finds the .95 quantile of N(100,16) distribution  
[1] 126.3177  
>
```

RStudio

File Edit Code View Plots Session Build Debug Profile Tools Help

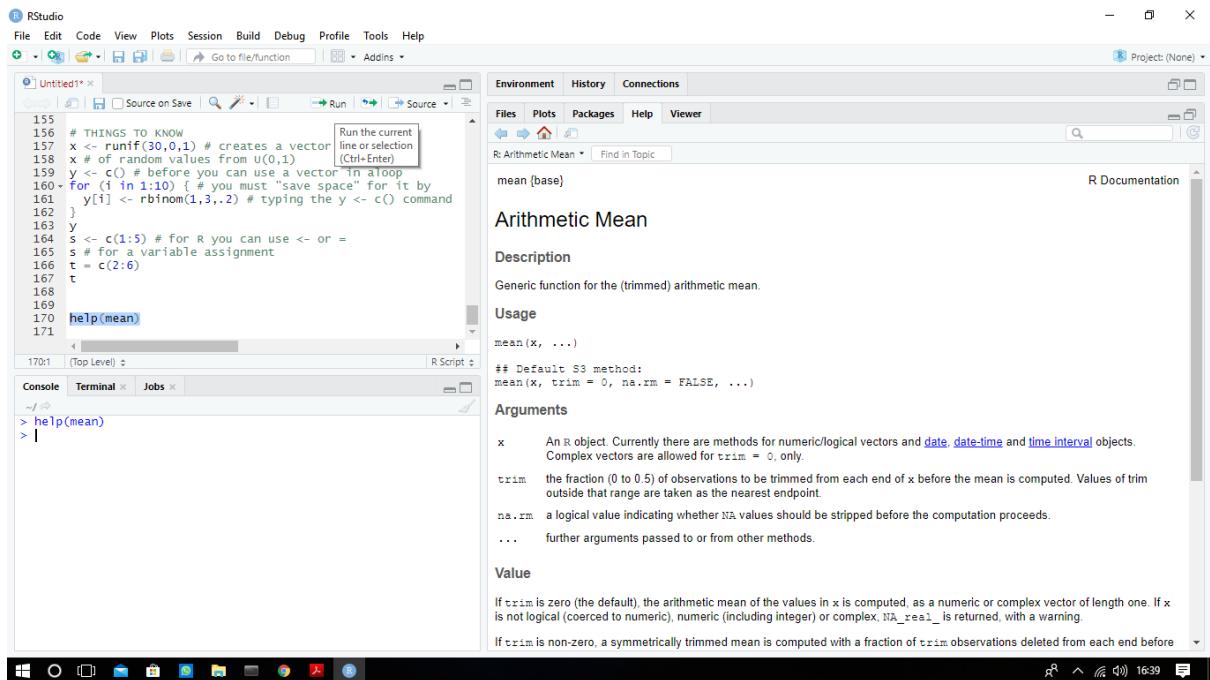
Untitled1* Source on Save Run Source

```
151 qnorm(.95,100,16) # finds the .95 quantile of N(100,16) distribution  
152  
153  
154  
155  
156 # THINGS TO KNOW  
157 x <- runif(30,0,1) # creates a vector 30 long  
158 y <- c() # before you can use a vector in a loop  
159 for (i in 1:10) { # you must "save space" for it by  
160 y[i] <- rbinom(1,3,,2) # typing the y <- c() command  
161 }  
162  
163 y  
164 s <- c(1:5) # for R you can use <- or =  
165 s # for a variable assignment  
166 t = c(2:6)  
167 t
```

(Top Level) R Script

Console Terminal Jobs

```
> x # of random values from U(0,1)  
[1] 0.661798792 0.040347270 0.546019772 0.804877741 0.634598822 0.884237377 0.791882051 0.103549371 0.247004411 0.646603940  
[1] 0.627884216 0.671641702 0.658571128 0.785585448 0.287223652 0.286631284 0.921992715 0.505599710 0.009688218 0.411355618  
[21] 0.716804039 0.361962770 0.021335714 0.206812853 0.436679048 0.588968598 0.317600365 0.890287028 0.035083916 0.119840272  
> y <- c() # before you can use a vector in a loop  
> for (i in 1:10) { # you must "save space" for it by  
+ y[i] <- rbinom(1,3,,2) # typing the y <- c() command  
+ }  
+ y  
[1] 0 1 1 1 1 2 0 1 1 1  
> s <- c(1:5) # for R you can use <- or =  
> s # for a variable assignment  
[1] 1 2 3 4 5  
> t = c(2:6)  
[1] 2 3 4 5 6  
>
```



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

Experiment No. 3

Aim : Basic Data types and data structures in R - vector ,matrices,list and dataframes.

Theory :

Vectors are the most basic R data objects and there are six types of atomic vectors. They are logical, integer, double, complex, character and raw. Elements of a Vector are accessed using indexing. The [] brackets are used for indexing. Indexing starts with position 1. Giving a negative value in the index drops that element from the result. TRUE, FALSE or 0 and 1 can also be used for indexing.

Matrices are the R objects in which the elements are arranged in a two-dimensional rectangular layout. They contain elements of the same atomic types. Though we can create a matrix containing only characters or only logical values, they are not of much use. We use matrices containing numeric elements to be used in mathematical calculations.

A Matrix is created using the `matrix()` function. The basic syntax for creating a matrix in R is –

```
matrix(data, nrow, ncol, byrow, dimnames)
```

Following is the description of the parameters used –

- data is the input vector which becomes the data elements of the matrix.
- nrow is the number of rows to be created.
- ncol is the number of columns to be created.
- byrow is a logical clue. If TRUE then the input vector elements are arranged by row.
- dimname is the names assigned to the rows and columns.

Lists are the R objects which contain elements of different types like – numbers, strings, vectors and another list inside it. A list can also contain a matrix or a function as its elements. List is created using `list()` function.

R list is the object which contains elements of different types – like *strings, numbers, vectors and another list* inside it. R list can also contain a matrix or a function as its elements. The list is created using the `list()` function in R. In other words, a list is a generic vector containing other objects.

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

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

- Each column should contain the same number of data items.

Output :

Code :

```

x <- "dataset"
typeof(x)
attributes(x)
y <- 1:10
y
typeof(y)
length(y)
z <- as.numeric(y)
z
typeof(z)
vector()
vector("character", length = 5) # a vector of mode 'character' with 5 elements
character(5) # the same thing, but using the constructor directly
numeric(5) # a numeric vector with 5 elements
logical(5) # a logical vector with 5 elements
x <- c(1, 2, 3)
x1 <- c(1L, 2L, 3L)
y <- c(TRUE, TRUE, FALSE, FALSE)
z <- c("Sarah", "Tracy", "Jon")
typeof(z)
length(z)
class(z)
str(z)
z <- c(z, "Annette")
z
z <- c("Greg", z)
z
series <- 1:10
seq(10)
seq(from = 1, to = 10, by = 0.1)
x <- c(0.5, NA, 0.7)
x
x <- c(TRUE, FALSE, NA)
x
x <- c("a", NA, "c", "d", "e")
x
x <- c(1+5i, 2-3i, NA)
x

```

```

x <- c("a", NA, "c", "d", NA)
y <- c("a", "b", "c", "d", "e")
is.na(x)
is.na(y)
anyNA(x)
anyNA(y)
1/0
0/0
xx <- c(1.7, "a")
xx <- c(TRUE, 2)
xx <- c("a", TRUE)
as.numeric("1") #You can also control how vectors are coerced explicitly using the
as.<class_name>() functions:
as.character(1:2)
length(1:10) #to get length of list
nchar("Software Carpentry") #to get length of a character string
m <- matrix(nrow = 2, ncol = 2)
m
dim(m)
m <- matrix(c(1:3))
class(m)
typeof(m)
FOURS <- matrix(
  c(4, 4, 4, 4),
  nrow = 2,
  ncol = 2)
typeof(FOURS[1])
typeof(FOURS)
m <- matrix(1:6, nrow = 2, ncol = 3)
m <- 1:10 #another way to make a matrix
dim(m) <- c(2, 5)
x <- 1:3
y <- 10:12
cbind(x, y)
rbind(x, y)
mdat <- matrix(c(1, 2, 3, 11, 12, 13),
               nrow = 2,
               ncol = 3,
               byrow = TRUE)
mdat
mdat[2, 3]
x <- list(1, "a", TRUE, 1+4i)
x
x <- vector("list", length = 5) # empty list

```

```
length(x)
x[[1]]
#Vectors can be coerced to lists as follows:
```

```
x <- 1:10
x <- as.list(x)
length(x)
xlist <- list(a = "Karthik Ram", b = 1:10, data = head(iris)) #Elements of a list can be named
(i.e. lists can have the names attribute)
xlist
names(xlist)

dat <- data.frame(id = letters[1:10], x = 1:10, y = 11:20) #To create data frames by hand
dat
is.list(dat) #data frame is a special list
class(dat)
dat[1, 3]
#As data frames are also lists, it is possible to refer to columns (which are elements of such
list) using the list notation, i.e. either double square brackets or a $
dat[["y"]]
dat$y
```

The screenshot shows the RStudio interface. In the top right, the help documentation for the 'mean' function is displayed under the 'R Arithmetic Mean' section. The 'Description' and 'Usage' sections are visible. In the bottom right, the 'Arguments' section is expanded, showing that 'x' is described as an R object. The main area shows a console session with various R code snippets demonstrating different ways to create vectors and calculate their means.

```
> /<
> x <- "dataset"
> typeof(x)
[1] "character"
> attributes(x)
NULL
> y <- 1:10
> y
[1] 1 2 3 4 5 6 7 8 9 10
> typeof(y)
[1] "integer"
> length(y)
[1] 10
> z <- as.numeric(y)
> z
[1] 1 2 3 4 5 6 7 8 9 10
> typeof(z)
[1] "double"
> vector()
logical(0)
> logical("character", length = 5) # a vector of mode 'character' with 5 elements
[1] 
> character(5) # the same thing, but using the constructor directly
[1] 
> numeric(5) # a numeric vector with 5 elements
[1] 0.0 0.0 0.0
> logical(5) # a logical vector with 5 elements
[1] FALSE FALSE FALSE FALSE FALSE
> x <- c(1, 2, 3)
> x1 <- c(1L, 2L, 3L)
> y <- c(TRUE, TRUE, FALSE, FALSE)
> z <- c("Sarah", "Tracy", "Jon")
> typeof(z)
[1] "character"
> length(z)
[1]
> class(z)
[1] "character"
> str(z)
```

This screenshot is similar to the one above, showing the RStudio interface with the 'mean' function documentation open. The 'Arguments' section for 'x' is also expanded, stating it is an R object. The console session shows more examples of creating vectors and calculating their means, including the use of 'seq' and 'c' functions.

```
> /<
> str(z)
chr [1:3] "Sarah" "Tracy" "Jon"
> z <- c(z, "Annette")
> z
[1] "Sarah" "Tracy" "Jon" "Annette"
> z <- c("Greg", z)
> z
[1] "Greg" "Sarah" "Tracy" "Jon" "Annette"
> series <- 1:10
> seq(10)
[1] 1 2 3 4 5 6 7 8 9 10
> seq(from = 1, to = 10, by = 0.1)
[1] 1.0 1.1 1.2 1.3 1.4 1.5 1.6 1.7 1.8 1.9 2.0 2.1 2.2 2.3 2.4 2.5 2.6 2.7 2.8 2.9 3.0 3.1 3.2 3.3 3.4 3.5
[27] 3.6 3.7 3.8 3.9 4.0 4.1 4.2 4.3 4.4 4.5 4.6 4.7 4.8 4.9 5.0 5.1 5.2 5.3 5.4 5.5 5.6 5.7 5.8 5.9 6.0 6.1
[51] 6.2 6.3 6.4 6.5 6.6 6.7 6.8 6.9 7.0 7.1 7.2 7.3 7.4 7.5 7.6 7.7 7.8 7.9 8.0 8.1 8.2 8.3 8.4 8.5 8.6 8.7
[75] 8.8 8.9 9.0 9.1 9.2 9.3 9.4 9.5 9.6 9.7 9.8 9.9 10.0
> x <- c(0.5, NA, 0.7)
> x
[1] 0.5 NA 0.7
> x <- c(TRUE, FALSE, NA)
> x
[1] TRUE FALSE NA
> x <- c("a", NA, "c", "d", "e")
> x
[1] "a" NA "c" "d" "e"
> x <- c(1+5i, 2-3i, NA)
> x
[1] 1+5i 2-3i NA
> x <- c("a", NA, "c", "d", NA)
> y <- c("a", "b", "c", "d", "e")
> is.na(x)
[1] FALSE TRUE FALSE FALSE TRUE
> is.na(y)
[1] FALSE FALSE FALSE FALSE FALSE
> any(is.na(x))
[1] TRUE
> any(is.na(y))
[1] FALSE
```

RStudio

File Edit Code View Plots Session Build Debug Profile Tools Help

Project (None)

Source

Console Terminal Jobs

```
> f<-function()
> anyNA(y)
[1] FALSE
> 1/0
[1] Inf
> 0/0
[1] NaN
> xx <- c(1.7, "a")
> xx <- c(TRUE, 2)
> xx <- c("a", TRUE)
> as.numeric("a") # You can also control how vectors are coerced explicitly using the as.<class.name>() functions:
[1] 1
> as.character(1:2)
[1] "1" "2"
> length(1:10) #to get length of list
[1] 10
> nchar("software Carpentry") #to get length of a character string
[1] 18
> m <- matrix(nrow = 2, ncol = 2)
> m
     [,1] [,2]
[1,] NA NA
[2,] NA NA
> dim(m)
[1] 2 2
> n <- matrix(c(1:3))
> class(n)
[1] "matrix"
> typeof(n)
[1] "integer"
> FOURS <- matrix(
+   c(4, 4, 4, 4),
+   nrow = 2,
+   ncol = 2)
> typeof(FOURS[1])
[1] "double"
> typeof(FOURS)
[1] "double"
> n <- matrix(1:6, nrow = 2, ncol = 3)
```

Environment History Code

Files Plots Pack

R Arithmetic Mean

(base) Documentation

Arithmetic Mean

Description

Generic function for the (trimmed) arithmetic mean.

Usage

```
mean(x, ...)
```

Default S3 method

```
mean(x, trim = 0,
```

Arguments

x An R object. Currently there are methods for numeric/logical vectors and `date`, `date-time` and `time` `interval` objects.

RStudio

File Edit Code View Plots Session Build Debug Profile Tools Help

Project (None)

Source

Console Terminal Jobs

```
> f<-function()
> m <- matrix(1:6, nrow = 2, ncol = 3)
> m <- 1:10 #another way to make a matrix
> dim(m) <- c(2, 5)
> x <- 1:3
> y <- 10:12
> cbind(x, y)
      x   y
[1,] 1 10
[2,] 2 11
[3,] 3 12
> rbind(x, y)
     [,1] [,2] [,3]
x    1    2    3
y    10   11   12
> mdat <- matrix(c(1, 2, 3, 11, 12, 13),
+                  nrow = 2,
+                  ncol = 3,
+                  byrow = TRUE)
> mdat
     [,1] [,2] [,3]
[1,] 1    2    3
[2,] 11   12   13
> mdat[2, 3]
[1] 13
> x <- list(1, "a", TRUE, 1+4i)
> x
[[1]]
[1] 1
[[2]]
[1] "a"
[[3]]
[1] TRUE
[[4]]
[1] 1+4i
```

Environment History Code

Files Plots Pack

R Arithmetic Mean

(base) Documentation

Arithmetic Mean

Description

Generic function for the (trimmed) arithmetic mean.

Usage

```
mean(x, ...)
```

Default S3 method

```
mean(x, trim = 0,
```

Arguments

x An R object. Currently there are methods for numeric/logical vectors and `date`, `date-time` and `time` `interval` objects.

The screenshot shows the RStudio interface. In the top right, the help documentation for the 'mean' function is displayed under the 'R Arithmetic Mean' section. The main area shows a console session where a list 'x' is created and then converted into a data frame 'dat'. The 'mean' function is then used on the 'Sepal.Length' column of 'dat'.

```
> x <- vector("list", length = 5) # empty list
> length(x)
[1] 5
> x[[1]]
NULL
> #vectors can be coerced to lists as follows:
>
> x <- 1:10
> x <- as.list(x)
> length(x)
[1] 10
> xlist <- list(a = "karthik Ram", b = 1:10, data = head(iris)) #Elements of a list can be named (i.e. Lists can have the names attribute)
> xlist
$ a
[1] "karthik Ram"
$b
[1] 1 2 3 4 5 6 7 8 9 10
$data
  Sepal.Length Sepal.Width Petal.Length Petal.Width Species
1          5.1         3.5          1.4         0.2   setosa
2          4.9         3.0          1.4         0.2   setosa
3          4.7         3.2          1.3         0.2   setosa
4          4.6         3.1          1.5         0.2   setosa
5          5.0         3.6          1.4         0.2   setosa
6          5.4         3.9          1.7         0.4   setosa
> names(xlist)
[1] "a"      "b"      "data"
>
> dat <- data.frame(id = letters[1:10], x = 1:10, y = 11:20) #To create data frames by hand
> dat
  id x y
1  a 1 11
2  b 2 12
3  c 3 13
4  d 4 14
5  e 5 15
6  f 6 16
7  g 7 17
8  h 8 18
9  i 9 19
10 j 10 20
> is.list(dat) #data frame is a special list
[1] TRUE
> class(dat)
[1] "data.frame"
> dat[, 3]
[1] 11 12 13 14 15 16 17 18 19 20
> dat$y
[1] 11 12 13 14 15 16 17 18 19 20
>
```

This screenshot shows a second console session in RStudio. It demonstrates creating a list 'xlist' and then creating a data frame 'dat' from it. The 'mean' function is applied to the 'Sepal.Length' column of 'dat'. The help documentation for 'mean' is visible on the right side of the interface.

```
> xlist
  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
> names(xlist)
[1] "a"      "b"      "data"
>
> dat <- data.frame(id = letters[1:10], x = 1:10, y = 11:20) #To create data frames by hand
> dat
  id x y
1  a 1 11
2  b 2 12
3  c 3 13
4  d 4 14
5  e 5 15
6  f 6 16
7  g 7 17
8  h 8 18
9  i 9 19
10 j 10 20
> is.list(dat) #data frame is a special list
[1] TRUE
> class(dat)
[1] "data.frame"
> dat[, 3]
[1] 11 12 13 14 15 16 17 18 19 20
> dat$y
[1] 11 12 13 14 15 16 17 18 19 20
>
```

Conclusion : Thus, we have studied basic Data types and data structures in R - vector ,matrices,list and dataframes.

Experiment No. 4

Aim : Programming constructs in R - loops ,conditional executions

Theory :

A loop statement allows us to execute a statement or group of statements multiple times and the following is the general form of a loop statement in most of the programming languages .

1. Repeat Loop : The Repeat loop executes the same code again and again until a stop condition is met.
2. While Loop : The Repeat loop executes the same code again and again until a stop condition is met.
3. For Loop : A For loop is a repetition control structure that allows you to efficiently write a loop that needs to execute a specific number of times.

A conditional expression or a conditional statement is a programming construct where a decision is made to execute some code based on a boolean (true or false) condition. A more commonly used term for conditional expression in programming is an 'if-else' condition. In plain English, this is stated as 'if this test is true then do this operation; otherwise do this different operation'.

Output :

Code :

```
#IF STATEMENTS
values <- 1:10
if (sample(values,1) <= 10)
  print(paste(values, "is less than or equal to 10"))

#if else
val1 = 10                      #Creating our first variable val1
val2 = 5                        #Creating second variable val2
if (val1 > val2){               #Executing Conditional Statement based on the comparison
  print("Value 1 is greater than Value 2")
} else if (val1 < val2){
  print("Value 1 is less than Value 2")

}

#For loop
values <- c(1,2,3,4,5)
for(id in 1:5){
```

```
    print(values[id]^values[id])
}
```

```
#Nested for loop
mat <- matrix(1:10, 2)
for (id1 in seq(nrow(mat))) {
  for (id2 in seq(ncol(mat))) {
    print(mat[id1, id2])
  }
}
mat
```

```
val = 2.987
while(val <= 4.987) {
  val = val + 0.987
  print(c(val, val-2, val-1))
}
```

```
val <- 5
repeat {
  print(val)
  val <- val+1
  if (val == 10){
    break
  }
}
```

```
val <- 5
repeat {
  print(val)
  val <- val+1
  if (val == 10){
    break
  }
}
values = 1:10
for (id in values){
  if (id == 2){
    break
  }
  print(id)
}
x = 1: 4
for (i in x) {
```

```

if (i == 2) {
  next
}
print(i)
}

check <- function(x) {
  if (x > 0) {
    result <- "Positive"
  } else if (x < 0) {
    result <- "Negative"
  } else {
    result <- "Zero"
  }
  return(result)
}
check(1)
check(0)
check(10)

```

Screenshots

The screenshot shows the RStudio interface with the following details:

- Console Tab:** Displays the R code and its output. The code includes an if-statement, a nested loop, and a matrix operation. The output shows the results of these operations.
- Environment Tab:** Shows the current environment variables, including `val1` and `val2` with their assigned values (10 and 5 respectively), and a matrix `mat` with dimensions 10x2.
- File Explorer:** Shows a project structure with files like `a.R`, `b.R`, `E..`, `R`, `V..`, and `d..`.
- Task View:** Shows various available packages and tools.
- System Tray:** Shows standard Windows icons for network, battery, and system status.

```

RStudio
File Edit Code View Plots Session Build Debug Profile Tools Help
Go to file/function Addins
Project (None)
Source
Console Terminal Jobs
~/r
> #IF STATEMENTS
> values <- 1:10
> if (sample(values,1) <= 10)
+   print(paste(values, "is less than or equal to 10"))
[1] "1 is less than or equal to 10" "2 is less than or equal to 10" "3 is less than or equal to 10" "4 is less than or equal to 10"
[5] "5 is less than or equal to 10" "6 is less than or equal to 10" "7 is less than or equal to 10" "8 is less than or equal to 10"
[9] "9 is less than or equal to 10" "10 is less than or equal to 10"
>
> #IF ELSE
> val1 = 10
> val2 = 5
#Creating our first variable val1
#Creating second variable val2
> if (val1 > val2){
+   print("Value 1 is greater than value 2")
+ } else if (val1 < val2){
+   print("Value 1 is less than value 2")
+
[1] "Value 1 is greater than value 2"
>
> #For loop
> values <- c(1,2,3,4,5)
> for(id in 1:5){
+   print(values[id]^values[id])
+
[1] 1
[1] 4
[1] 27
[1] 256
[1] 3125
>
> #nested for loop
> mat <- matrix(1:10, 2)
> for (id1 in seq(nrow(mat))) {
+   for (id2 in seq(ncol(mat))) {
+     print(mat[id1, id2])
+
[1] 1
[1] 1

```

RStudio

File Edit Code View Plots Session Build Debug Profile Tools Help

Project: (None)

Source

Console Terminal Jobs

```

> mat <- matrix(1:10, 2)
> for (id1 in seq(nrow(mat))) {
+   for (id2 in seq(ncol(mat))) {
+     print(mat[id1, id2])
+   }
+ }
[1] 1
[1] 3
[1] 5
[1] 7
[1] 9
[1] 2
[1] 4
[1] 6
[1] 8
[1] 10
> mat
[,1] [,2] [,3] [,4] [,5]
[1,]    1    3    5    7    9
[2,]    2    4    6    8   10
>
> val = 2.987
> while(val <= 4.987) {
+   val = val + 0.987
+   print(c(val, val-2, val-1))
+ }
[1] 3.974 1.974 2.974
[1] 4.961 2.961 3.961
[1] 5.948 3.948 4.948
>
> val <= 5
> repeat {
+   print(val)
+   val <- val+1
+   if (val == 10){
+     break
+   }
+ }

```

Windows Taskbar: RStudio, File Explorer, Google Chrome, etc.

RStudio

File Edit Code View Plots Session Build Debug Profile Tools Help

Project: (None)

Source

Console Terminal Jobs

```

> val <- 5
> repeat {
+   print(val)
+   val <- val+1
+   if (val == 10){
+     break
+   }
+ }
[1] 5
[1] 6
[1] 7
[1] 8
[1] 9
> val <- 5
> repeat {
+   print(val)
+   val <- val+1
+   if (val == 10){
+     break
+   }
+ }
[1] 5
[1] 6
[1] 7
[1] 8
[1] 9
> values = 1:10
> for (id in values){
+   if (id == 2){
+     break
+   }
+   print(id)
+ }
[1] 1
>
> x = 1: 4

```

Windows Taskbar: RStudio, File Explorer, Google Chrome, etc.

The screenshot shows the RStudio IDE interface. The main area is a code editor with tabs for 'Source', 'Terminal', and 'Jobs'. The 'Source' tab is active, displaying the following R code:

```
~/r
> values = 1:10
> for (id in values){
+   if (id == 2){
+     break
+   }
+   print(id)
+ }
[1] 1
>
> x = 1: 4
> for (i in x) {
+   if (i == 2) {
+     next
+   }
+   print(i)
+ }
[1] 1
[1] 3
[1] 4
>
> check <- function(x) {
+   if (x > 0) {
+     result <- "Positive"
+   } else if (x < 0) {
+     result <- "Negative"
+   } else {
+     result <- "zero"
+   }
+   return(result)
+ }
> check(1)
[1] "Positive"
> check(0)
[1] "zero"
> check(10)
[1] "Positive"
>
```

To the right of the code editor is the 'Environment' pane, which displays a file tree under 'Files'. The tree includes 'New Folder', 'Home', and several files and folders such as 'b...', 'E...', 'R', 'V...', 'A...', and 'd...'. The status bar at the bottom right shows the time as 17:26.

Conclusion : Thus we have studied the programming constructs in 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").

Output:

Code:

```
install.packages("expss")
library(expss)
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 the 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(a22, other ~ copy) %into% h22
    recode(b22, other ~ copy) %into% p22
    c1r = c1
  }) %>%
  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(a22, other ~ copy) %into% p22
recode(b22, other ~ copy) %into% h22
recode(c1, 1 ~ 2, 2 ~ 1, other ~ copy) %into% c1r
}) %>%
compute({
  # recode age by groups
  age_cat = recode(s2a, lo %thru% 25 ~ 1, lo %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
")

overall_liking_scale = num_lab("
  1 Extremely poor
  2 Very poor
  3 Quite poor
  4 Neither good, nor poor
  5 Quite good
  6 Very good
  7 Excellent
")

w = apply_labels(w,
  c1r = "Preferences",
  c1r = num_lab(
    1 VSX123
    2 SDF456
    3 Hard to say
  ),
  age_cat = "Age",
  age_cat = c("18 - 25" = 1, "26 - 35" = 2),
  h1_1 = "Likes. VSX123",
)

```

```

p1_1 = "Likes. SDF456",
h1_1 = codeframe_likes,
p1_1 = codeframe_likes,

h_likes = "Number of likes. VSX123",
p_likes = "Number of likes. SDF456",

h22 = "Overall quality. VSX123",
p22 = "Overall quality. SDF456",
h22 = overall_liking_scale,
p22 = overall_liking_scale
)
# 'tab_mis_val(3)' remove 'hard to say' from vector
w %>% tab_cols(total(), age_cat) %>%
  tab_cells(c1r) %>%
  tab_mis_val(3) %>%
  tab_stat_cases() %>%
  tab_last_sig_cases() %>%
  tab_pivot()

# lets specify repeated parts of table creation chains
banner = w %>% tab_cols(total(), age_cat, c1r)
# column percent with significance
tab_cpct_sig = . %>% tab_stat_cpct() %>%
  tab_last_sig_cpct(sig_labels = paste0("<b>", LETTERS, "</b>"))

# means with significance
tab_means_sig = . %>% tab_stat_mean_sd_n(labels = c("<b><u>Mean</u></b>", "sd", "N"))
%>%
  tab_last_sig_means(
    sig_labels = paste0("<b>", LETTERS, "</b>"),
    keep = "means")

# Preferences
banner %>%
  tab_cells(c1r) %>%
  tab_cpct_sig() %>%
  tab_pivot()

# Overall liking
banner %>%
  tab_cells(h22) %>%
  tab_means_sig() %>%
  tab_cpct_sig() %>%
  tab_cells(p22) %>%
  tab_means_sig() %>%
  tab_cpct_sig() %>%
  tab_pivot()

# Likes

```

```

banner %>%
  tab_cells(h_likes) %>%
  tab_means_sig() %>%
  tab_cells(mrset(h1_1 %to% h1_6)) %>%
  tab_cpct_sig() %>%
  tab_cells(p_likes) %>%
  tab_means_sig() %>%
  tab_cells(mrset(p1_1 %to% p1_6)) %>%
  tab_cpct_sig() %>%
  tab_pivot()

# below more complicated table where we compare likes side by side
# Likes - side by side comparison
w %>%
  tab_cols(total(label = "#Total| |"), c1r) %>%
  tab_cells(list(unvr(mrset(h1_1 %to% h1_6)))) %>%
  tab_stat_cpct(label = var_lab(h1_1)) %>%
  tab_cells(list(unvr(mrset(p1_1 %to% p1_6)))) %>%
  tab_stat_cpct(label = var_lab(p1_1)) %>%
  tab_pivot(stat_position = "inside_columns")
write_labelled_csv(w, file = "product_test.csv")

library(exps)
w = text_to_columns("
  a b c
  2 15 50
  1 70 80
  3 30 40
  2 30 40
")
w$d = ifelse(w$b>60, 1, 0)
w = compute(w, {
  d = ifelse(b>60, 1, 0)
  e = 42
  abc_sum = sum_row(a, b, c)
  abc_mean = mean_row(a, b, c)
})
count_if(1, w)
calculate(w, count_if(1, a, b, c))

w$d = count_row_if(gt(1), w)
w = compute(w, {
  d = count_row_if(gt(1), a, b, c)
})
count_col_if(le(1), w$a)
sum(w, na.rm = TRUE)
w$d = mean_row(w)
#or
w = compute(w, {
  d = mean_row(a, b, c)
})

```

```

})
sum_col(w$a)
sum_if(gt(40), w)
#or
calculate(w, sum_if(gt(40), a, b, c))
w$d = sum_row_if(lt(40), w)
#or
w = compute(w, {
  d = sum_row_if(lt(40), a, b, c)
})
mean_col_if(lt(3), w$a, data = w$b)
#or
calculate(w, mean_col_if(lt(3), a, data = sheet(b, c)))

dict = text_to_columns(
  x y
  1 apples
  2 oranges
  3 peaches
)
w$d = vlookup(w$a, dict, 2)
#or
w$d = vlookup(w$a, dict, "y")

```

Screenshots:

The screenshot shows the RStudio interface with the following details:

- MenuBar:** File, Edit, Code, View, Plots, Session, Build, Debug, Profile, Tools, Help.
- Console Tab:** Shows the R code being run and its output. The output includes:
 - The downloaded source packages are in C:\Users\shail\AppData\Local\Temp\Rtmpw0LFYB\downloaded_packages
 - library(exps)
 - 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 - vsxl23 sample, 'p' variables - 'SDFA56' 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(a22, other ~ copy) %into% h22
 - recode(b22, other ~ copy) %into% p22
 - cir ~ c1
 - }) %>%
 - 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(a22, other ~ copy) %into% p22
 - recode(b22, other ~ copy) %into% h22
 - recode(c1, 1 ~ 2, 2 ~ 1, other ~ copy) %into% cir
 - }) %>%
 - 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)
 - })
- Environment Tab:** Shows the current environment with various objects listed.
- Taskbar:** Shows the Windows taskbar with multiple pinned icons.
- System Tray:** Shows the system tray with icons for battery, signal, and time (17:47).

RStudio

```

File Edit Code View Plots Session Build Debug Profile Tools Help
Project (None)
Untitled1* Source on Save Run Source
72     h1_1 = codeframe_likes,
73     p1_1 = codeframe_likes,
74
75     h_likes = "number of likes. VSX123",
76     p_likes = "number of likes. SDF456",
77
78     h2_1 = "overall quality. VSX123",
79     p2_1 = "overall quality. SDF456",
80     h2_2 = overall_liking_scale,
81     p2_2 = overall_liking_scale
82 )
83 # `tab_mis_val()` remove "hard to say" from vector
84 w %>% tab_cols(total(), age_cat) %>%
85   tab_cells(cir) %>%
86   tab_mis_val() %>%
87   tab_stat_cases() %>%
88   tab_last_sig_cases() %>%
89   tab_pivot()
89:14 (Top Level) s

```

Console Terminal Jobs

```

> # `tab_mis_val()` remove "hard to say" from vector
> w %>% tab_cols(total(), age_cat) %>%
+   tab_cells(cir) %>%
+   tab_mis_val() %>%
+   tab_stat_cases() %>%
+   tab_last_sig_cases() %>%
+   tab_pivot()
#>
#> | | #Total | 18 - 25 | 26 - 35 |
#> Preferences |-----|-----|-----|
#> | VSX123 | 94.0 | 46.0 | 48.0 |
#> | SDF456 | 50.0 | 22.0 | 28.0 |
#> | Hard to say | 1.0 | 0.0 | 0.0 |
#> | #chi-squared p-value | <0.05 | (warn.) |
#> | #Total cases | 144.0 | 68.0 | 76.0 |
#>

```

RStudio

```

File Edit Code View Plots Session Build Debug Profile Tools Help
Project (None)
Untitled1* Source on Save Run Source
90
91 # lets specify repeated parts of table creation chains
92 banner = w %>% tab_cols(total(), age_cat, cir)
93 # column percent with significance
94 tab_cptc_sig = . %>% tab_stat_cptc() %>%
95   tab_last_sig_cptc(sig_labels = paste0("<b>", LETTERS, "</b>"))
96
97 # means with significance
98 tab_means_sig = . %>% tab_stat_mean_sd_n(labels = c("<b><u>Mean</u></b>", "sd", "N")) %>%
99   tab_last_sig_means(
100   sig_labels = paste0("<b>", LETTERS, "</b>"),
101   keep = "means")
102
103 # Preferences
104 banner %>%
105   tab_cells(cir) %>%
106   tab_cptc_sig() %>%
107   tab_pivot()
107:15 (Top Level) s

```

Console Terminal Jobs

```

#> keep = "means"
#> # Preferences
#> banner %>%
#>   tab_cells(cir) %>%
#>   tab_cptc_sig() %>%
#>   tab_pivot()
#>
#> | | #Total | 18 - 25 | 26 - 35 | Preferences | VSX123 | SDF456 | Hard to say |
#> |-----|-----|-----|-----|-----|-----|-----|-----|
#> Preferences |-----|-----|-----|-----|-----|-----|-----|
#> | VSX123 | 62.7 | 65.7 | 60.0 | 100.0 | 100.0 | 100.0 | 6 |
#> | SDF456 | 33.3 | 31.4 | 35.0 | 100.0 | 100.0 | 100.0 | 6 |
#> | Hard to say | 4.0 | 2.9 | 5.0 | 100.0 | 100.0 | 100.0 | 6 |
#> | #Total cases | 150 | 70 | 80 | 94 | 50 | 100 | 6 |
#>

```

The screenshot shows the RStudio interface with the following components:

- Top Bar:** RStudio, File, Edit, Code, View, Plots, Session, Build, Debug, Profile, Tools, Help.
- Left Sidebar:** Project (None), Environment, Files, New File, Home.
- Script Editor:** An R script titled "Untitled1" containing code for creating a pivot table. The code includes functions like `tab_cells`, `tab_pct`, `tab_sig`, and `tab_pivot`.
- Console:** Displays the R command history and the resulting pivot table output.
- Environment:** Shows the current objects in memory, including the pivot table.

```
90
91 # lets specify repeated parts of table creation chains
92 banner = w %>% tab_cols(total(), age_cat, clr)
93 # column percent with significance
94 tab_pct_sig = . %>% tab_stat_pct() %>%
95   tab_last_sig_pct(sig_labels = paste0("<b>", LETTERS, "</b>"))
96
97 # means with significance
98 tab_mean_sig = . %>% tab_stat_mean_sd_n(tables = c("<b><u>Mean</u></b>", "sd", "n")) %>%
99   tab_last_mean()
100 sig_labels = paste0("<b>", LETTERS, "</b>"),
101   keep = "means")
102
103 # Preferences
104 banner %>%
105   tab_cells(clr) %>%
106   tab_pct_sig() %>%
107   tab_pivot()
```

107:15 [Top Level] 2

Console Terminal Jobs

```
* keep = "means")
>
> # Preferences
> banner %>%
+   tab_cells(clr) %>%
+   tab_pct_sig() %>%
+   tab_pivot()
```

	#Total	Age	Preferences			
	18 - 25	26 - 35	A	B	C	D
Preferences	vSX123	62.7	65.7	60.0	100.0	
	SDF456	33.3	34.3	35.0		100.0
	Hard to say	4.0	2.9	5.0		100.0
#Total cases	150	70	80	94	50	6

File Edit Code View Plots Session Build Debug Profile Tools Help

Untitled

118 #_Likes
119 banner %>
120 tab_cells(h1_likes) %>
121 tab_means_sig() %>
122 tab_cells(mset(h1_1, h1_2, h1_3)) %>
123 tab_collapse_sig() %>
124 tab_collapse() %>
125

126 (top-level)

Console Terminal Jobs

number of Likes, vsx123 Likes, vsx123

	#total	Age	Preferences	SOF456	Hard to say
Liked everything	2.0	2.0	2.2	1.9	2.2
Disliked everything	3.3	1.4	5.0	4.3	2.0
Chocolate	34.0	38.6	30.0	35.1	32.0
Appearance	29.3	21.4	36.2 A	25.5	38.0
Taste	32.0	38.6	26.2	23.4	48.0 A
Stuffing	27.3	20.0	31.8	28.7	26.0
Nuts	66.7	72.9	61.3	69.1	60.0
Consistency	12.0	4.3	18.8 A	8.5	14.0
other				50.0 A	B
Hard to answer					
Total cases	150	70	80	94	50
Liked everything	2.0	2.0	2.2	2.0	2.0

Number of Likes, SOF456 Likes, SOF456

	#total	Age	Preferences	SOF456	Hard to say
Liked everything	2.0	2.0	2.2	2.0	2.0
Disliked everything	1.3	1.4	1.2	2.1	1.1
Chocolate	32.0	27.1	36.2	29.8	34.0
Appearance	32.0	35.7	28.7	34.0	30.0
Taste	39.3	42.9	36.2	36.2	44.0
Stuffing	27.3	24.3	30.0	31.9	20.0
Nuts	61.3	60.0	62.5	58.5	68.0
Consistency	10.0	5.7	13.8	11.7	6.0
other	0.7		1.2	1.1	
Hard to answer					
Total cases	150	70	80	94	50
Liked everything	2.0	2.0	2.2	2.0	2.0

RStudio

File Edit Code View Plots Session Build Debug Profile Tools Help

Untitled1.R

Source on Save | Run | Source

```

130 # below more complicated table where we compare Likes side by side
131 # Likes - side by side comparison
132 w %>%
133   tab_cols(total(label = "#total | "), c1r) %>%
134   tab_cells(list(unvr(mrset(hd_1 %to% hd_6)))) %>%
135   tab_stat_cpct(label = var_lab(hd_1)) %>%
136   tab_pivot(stat_position = "inside_columns")
137 (Top Level) 1

```

Console Terminal Jobs

```

> w
6.0
> # below more complicated table where we compare Likes side by side
> # Likes - side by side comparison
> w %>%
+   tab_cols(total(label = "#total | "), c1r) %>%
+   tab_cells(list(unvr(mrset(hd_1 %to% hd_6)))) %>%
+   tab_stat_cpct(label = var_lab(hd_1)) %>%
+   tab_pivot(stat_position = "inside_columns")

```

	#Total	Preferences	SDF456	Hard to say
	VSX123	Likes. SDF456	Likes. VSX123	Likes. SDF456
Liked everything	3.3	1.3	4.3	2.1
Disliked everything	34.0	32.0	35.1	28.8
Chocolate	29.3	32.0	25.5	34.0
Appearance	32.0	39.3	23.4	36.2
Taste	27.3	27.3	28.7	31.9
Stuffing	66.7	61.3	69.1	58.5
Nuts	12.0	10.0	8.5	11.7
Consistency	0.7	0.7	1.1	1.1
other	150.0	150.0	94.0	94.0
Hard to answer			50	50
#Total cases				6.0

RStudio

File Edit Code View Plots Session Build Debug Profile Tools Help

Untitled1.R

Source on Save | Run | Source

```

133 e = 42
134 abc_sum = sum_row(a, b, c)
135 abc_mean = mean_row(a, b, c)
136 })
137 count_if(1, w)
138 calculate(w, count_if(1, a, b, c))
139
140 (Top Level) 1

```

Console Terminal Jobs

```

> library(exps)
> w = text_to_columns(""
+   " a b c
+   2 15 50
+   1 70 80
+   3 30 40
+   2 30 40
+   ")
> w$b = ifelse(w$b>60, 1, 0)
> w = compute(w, {
+   d = ifelse(b>60, 1, 0)
+   e = 42
+   abc_sum = sum_row(a, b, c)
+   abc_mean = mean_row(a, b, c)
+ })
> count_if(1, w)
[1] 2
> calculate(w, count_if(1, a, b, c))
[1] 1
>

```

RStudio

File Edit Code View Plots Session Build Debug Profile Tools Help

Untitled1* Project (None)

```
159 w$id = count_row_if(gt(1), w)
160 w = compute(w, {
161   d = count_row_if(gt(1), a, b, c)
162 })
163 count_col_if(lt(1), w$ia)
164
165 (Top Level) 1
```

Console Terminal Jobs

```
> library(exps)
> w = text_to_columns("a b c
+ 1 15 50
+ 1 70 80
+ 3 30 40
+ 2 30 40
+ ")
> w$d = ifelse(w$b>60, 1, 0)
> w = compute(w, {
+   d = ifelse(b>60, 1, 0)
+   e = 42
+   abc_sum = sum_row(a, b, c)
+   abc_mean = mean_row(a, b, c)
+ })
> count_if(l, w)
[1] 2
> calculate(w, count_if(l, a, b, c))
[1] 2
> w$d = count_row_if(gt(1), w)
> w = compute(w, {
+   d = count_row_if(gt(1), a, b, c)
+ })
> count_col_if(lt(1), w$ia)
[1] 1
>
```

RStudio

File Edit Code View Plots Session Build Debug Profile Tools Help

Untitled1* Project (None)

```
166 w$io = mean_row(w)
167 #or
168 w = compute(w, {
169   d = mean_row(a, b, c)
170 })
171 sum_col(w$ia)
172 sum_if(gt(40), w)
173
174 calculate(w, sum_if(gt(40), a, b, c))
175 w$d = sum_row_if(lt(40), w)
184 (Top Level) 1
```

Console Terminal Jobs

```
> sum(w, na.rm = TRUE)
[1] 1026
> w$d = mean_row(w)
> #or
> w = compute(w, {
+   d = mean_row(a, b, c)
+ })
> sum_col(w$ia)
[1] 8
> sum_if(gt(40), w)
[1] 831.6667
> #or
> calculate(w, sum_if(gt(40), a, b, c))
[1] 200
> w$d = sum_row_if(lt(40), w)
> #or
> w = compute(w, {
+   d = sum_row_if(lt(40), a, b, c)
+ })
> mean_col_if(lt(3), w$ia, data = w$ib)
[1] 38.33333
> #or
> calculate(w, mean_col_if(lt(3), a, data = sheet(b, c)))
  b           c
38.33333 56.66667
>
```

The screenshot shows the RStudio interface. The top menu bar includes File, Edit, Code, View, Plots, Session, Build, Debug, Profile, Tools, Help, and Project (None). The main area has tabs for Untitled1 and a Go to file/function search bar. Below is a code editor with the following R script:

```
171 sum_col(w$a)
172 sum_if(gt(40), w)
173 #or
174 calculate(w, sum_if(gt(40), a, b, c))
175 w$ld = sum_row_if(lt(40), w)
176 #or
177 w = compute(w,
178   d = sum_row_if(lt(40), a, b, c)
179 )
180 mean_col_if(lt(3), w$ia, data = w$b)
181 #or
182 calculate(w, mean_col_if(lt(3), a, data = sheet(b, c)))
183
184 dict = text_to_columns(
185   x y
186   1 apples
187   2 oranges
188   3 peaches
189 )
190 w$ld = vlookup(w$ia, dict, 2)
191 #or
192 w$ld = vlookup(w$ia, dict, "y")
193
194
195
```

The status bar at the bottom indicates (Top Level) x. The bottom left shows a console tab with the following output:

```
> dict = text_to_columns(
+   x y
+   1 apples
+   2 oranges
+   3 peaches
+ )
> w$ld = vlookup(w$ia, dict, 2)
> #or
> w$ld = vlookup(w$ia, dict, "y")
> |
```

Experiment No. 6

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

Theory :

Data manipulation operations such as subset, group, update, join etc., are all inherently related. Keeping these related operations together allows for:

- concise and consistent syntax irrespective of the set of operations you would like to perform to achieve your end goal.
- performing analysis fluidly without the cognitive burden of having to map each operation to a particular function from a potentially huge set of functions available before performing the analysis.
- automatically optimising operations internally, and very effectively, by knowing precisely the data required for each operation, leading to very fast and memory efficient code.

data.table is an R package that provides an enhanced version of data.frames, which are the standard data structure for storing data in base R.

dplyr is a grammar of data manipulation, providing a consistent set of verbs that help you solve the most common data manipulation challenges:

- mutate() adds new variables that are functions of existing variables
- select() picks variables based on their names.
- filter() picks cases based on their values.
- summarise() reduces multiple values down to a single summary.
- arrange() changes the ordering of the rows.

These all combine naturally with group_by() which allows you to perform any operation “by group”. dplyr is designed to abstract over how the data is stored. That means as well as working with local data frames, you can also work with remote database tables, using exactly the same R code. Install the dplyr package then read vignette("databases", package = "dbplyr").

Output :

Code :

```
rm(list=ls())
N_id=10000
N_tr=200
T_tr=80
set.seed(1)
N=rpois(N_id,N_tr)
N_traj=rpois(sum(N),T_tr)
```

```

origin_lat=runif(N_id,-5,5)
origin_lon=runif(N_id,-5,5)
lat=lon=Traj_Id=rep(NA,sum(N_traj))
Pers_Id=rep(NA,length(N_traj))
s=1
for(i in 1:N_id){Pers_Id[s:(s+N[i])]=i;s=s+N[i]}
s=1
for(i in 1:length(N_traj))
{
  lat[s:(s+N_traj[i])]=origin_lat[Pers_Id[i]]+cumsum(c(0,rnorm(N_traj[i]-1,0,SD=.2)));
  lon[s:(s+N_traj[i])]=origin_lon[Pers_Id[i]]+cumsum(c(0,rnorm(N_traj[i]-1,0,SD=.2)));
  s=s+N_traj[i]
}

Pers_Id=rep(NA,sum(N_traj))
N0=cumsum(c(0,N))
s=1
for(i in
1:N_id){Pers_Id[s:(s+sum(N_traj[(N0[i]+1):N0[i+1]]))]=i;s=s+sum(N_traj[(N0[i]+1):N0[i+1
]])}
s=1
for(i in 1:sum(N)){Traj_Id[s:(s+N_traj[i])]=i;s=s+N_traj[i]}

df=data.frame(Pers_Id,Traj_Id,lat,lon)
library(dplyr)
ldf=data_frame(Pers_Id,Traj_Id,lat,lon)

library(data.table)
dt=data.table(Pers_Id,Traj_Id,lat,lon)

object.size(df)

rm(list=ls())
library(data.table)
system.time( load("dt_json_2.RData") )

dropbox_ldf <- "https://dl.dropboxusercontent.com/s/l7rolinwojn37e2/ldf_json_2.RData"
dropbox_df <- "https://dl.dropboxusercontent.com/s/wzr19v9pyl7ah0j/df_json_2.RData"
dropbox_dt <- "https://dl.dropboxusercontent.com/s/nlwi1mmsxr5f23k/dt_json_2.RData"

source_https <- function(loc,...)
{
  curl=getCurlHandle()
  curlSetOpt(cookiejar="cookies.txt", useragent="Mozilla/5.0", followlocation=TRUE)
  tmp <- getURLContent(loc, .opts=list(ssl.verifypeer=FALSE),curl=curl)
  fch <- source(tmp)
  fch
}

system.time(load("df_json_2.RData"))

```

```

system.time(n <- nrow(df))
system.time(df$first<-c(1,df$Traj_Id[2:n]!=df$Traj_Id[1:(n-1)]))
system.time(df$last<-c(df$Traj_Id[2:n]!=df$Traj_Id[1:(n-1)],1))
object.size(df)

lat_0=0
lon_0=0
system.time(df$test <-(lat_0-df$lat)^2+(lon_0-df$lon)^2) <=1)&(df$last==1)

f$first <- NULL
df$last <- NULL
object.size(df)
system.time(df$test<-(lat_0-df$lat)^2+(lon_0-df$lon)^2) <=1&(c(df$Traj_Id[2:n]!=
df$Traj_Id[1:(n-1)],1)==1)
system.time(list_Traj <- unique(df$Traj_Id[df$test==TRUE]))
system.time(base <- df[(df$Traj_Id %in% list_Traj)& (c(1,df$Traj_Id[2:n]!=df$Traj_Id[1:(n-
1)])==1),c("lat","lon")])
head(base)

X <- base[,c("lon","lat")]
library(KernSmooth)
kde2d <- bkde2D(X, bandwidth=c(bw.ucv(X[,1]),bw.ucv(X[,2])),gridsize = c(251L, 251L))
image(x=kde2d$x1, y=kde2d$x2,z=kde2d$fhat,col=rev(heat.colors(100)))
contour(x=kde2d$x1, y=kde2d$x2,z=kde2d$fhat, add=TRUE)

rm(list=ls())
library(data.table)
system.time( load("dt_json_2.RData") )

system.time( setkey(dt,Traj_Id) )
system.time(depart <-dt[!duplicated(Traj_Id)])
system.time( depart <- dt[J(unique(Traj_Id)), mult = "first"])
system.time( arrivee <- dt[J(unique(Traj_Id)), mult = "last"] )
lat_0=0
lon_0=0
system.time( arrivee[,dist:=(lat-lat_0)^2+(lon-lon_0)^2] )

system.time( fin <- subset(arrivee,dist <= 1) )
system.time( fin[,Pers_Id:=NULL] )
system.time( fin[,lat:=NULL] )
system.time( fin[,lon:=NULL] )
system.time( setkey(fin, Traj_Id) )
system.time( base <<- merge(fin,depart,all.x=TRUE) )
system.time( base <<- merge(fin,depart,all.x=TRUE) )
head(base)
rm(list=ls())
library(dplyr)
library(data.table)
system.time( load("ldf_json_2.RData") )
system.time( ldepart <<- ldf %>% group_by(Traj_Id) %>%

```

```

+ summarise(first_lat=head(lat,1),
           first_lon=head(lon,1)) )
system.time( larrive <- ldf %>% group_by(Traj_Id) %>%
  + summarise(last_lat=tail(lat,1),last_lon=tail(lon,1)) )
lat_0=0
lon_0=0
system.time( system.time( larrive <-> mutate(larrive,dist=(last_lat-lat_0)^2+(last_lon-
lon_0)^2) ) )
system.time( lfin <- filter(larrive,dist<=1) )
system.time( lbase <- left_join(lfin,ldepart) )
head(lbase)

```

Screenshots :

The screenshot shows the RStudio interface. The top menu bar includes File, Edit, Code, View, Plots, Session, Build, Debug, Profile, Tools, Help, and a Project tab. Below the menu is a toolbar with various icons. The main area has two tabs: 'Untitled1' (active) and 'Console'. The 'Untitled1' tab contains R code for generating trajectories and calculating distances. The 'Console' tab shows the execution of this code in the R environment.

```

8 N_traj=rpois(sum(N),T_tr)
9
10 origin_lat=runiform(N_id,-5,5)
11 origin_lon=runiform(N_id,-5,5)
12
13 lat=lon=Traj_Id=rep(NA,sum(N_traj))
14 pers_id=rep(NA,length(N_traj))
15 s=1
16 for(i in 1:N_id){pers_id[s:(s+N[i])]=i;s=s+N[i]}
17 s=1
18 for(i in 1:length(N_traj)){
19   lat[s:(s+N_traj[i])]=origin_lat[pers_id[i]]+cumsum(c(0,rnorm(N_traj[i]-1,0,sd=.2)))
20   lon[s:(s+N_traj[i])]=origin_lon[pers_id[i]]+cumsum(c(0,rnorm(N_traj[i]-1,0,sd=.2)))
21   s=s+N_traj[i]
}

```

```
1 #> set.seed(1)
2 N=rpois(N_id,N_traj)
3 N_traj=rpois(sum(N),T_mtr)
4
5 origin_lat=rnorm(N_id,-5,5)
6 origin_lon=rnorm(N_id,-5,5)
7
8 Lat=Lat=Traj_Id=rep(NA,sum(N_traj))
9 Pers_Id=rep(NA,length(N_traj))
10 s=1
11 for(i in 1:N_id){Pers_Id[s:(s+N[i])]=i;s=s+N[i]}
12 s=1
13 for(i in 1:length(N_traj)){
14   Lat[s:(s+N_traj[i])]=origin_lat[Pers_Id[i]]+cumsum(c(0,rnorm(N_traj[i]-1,0,sd=.2)));
15   Lon[s:(s+N_traj[i])]=origin_lon[Pers_Id[i]]+cumsum(c(0,rnorm(N_traj[i]-1,0,sd=.2)));
16   s=s+N_traj[i]
17 }
18 Pers_Id=rep(NA,sum(N_traj))
19 N0=cumsum(c(0,N))
20 s=1
21 for(i in 1:N_id){Pers_Id[s:(s+sum(N_traj[(N0[i]+1):N0[i+1]))]]=i;s=s+sum(N_traj[(N0[i]+1):N0[i+1]]);}
22 for(i in 1:sum(N)){Traj_Id[s:(s+N_traj[i])]=i;s=s+N_traj[i]}
23
24 df=data.frame(Pers_Id,Traj_Id,Lat,lon)
25
26 library(dplyr)
27 tdf=data.frame(Pers_Id,Traj_Id,Lat,lon)
28
29 library(data.table)
30 dt=data.table(Pers_Id,Traj_Id,Lat,lon)
31
32 object.size(df)
33
```

The screenshot shows an RStudio interface with an R script file named "Untitled.R" open. The code generates a trajectory dataset. It starts by setting a seed, defining parameters, and initializing variables. It then loops through individuals (N_id) and generates their trajectories (Lat and Lon) based on their origin coordinates and a random walk model. Finally, it creates a data frame (df) and a data.table (dt) from the generated data.

Conclusion : Thus, we have studied working with large datasets with dplyr and data.table.

Experiment No : 7

Aim : EDA with R - Range, Summary, Mean, Variance, Median, Standard Deviation, Histogram, Box plot, Scatter Graph.

Theory:

Exploratory data analysis or “EDA” is a critical first step in analyzing the data from an experiment. Here are the main reasons we use EDA:

- detection of mistakes
- checking of assumptions
- preliminary selection of appropriate models
- determining relationships among the explanatory variables, and
- assessing the direction and rough size of relationships between explanatory and outcome variables.

Loosely speaking, any method of looking at data that does not include formal statistical modeling and inference falls under the term exploratory data analysis.

Types of EDA:

Exploratory data analysis is generally cross-classified in two ways. First, each method is either non-graphical or graphical. And second, each method is either univariate or multivariate (usually just bivariate).

Non-graphical methods generally involve calculation of summary statistics, while graphical methods obviously summarize the data in a diagrammatic or pictorial way. Univariate methods look at one variable (data column) at a time, while multivariate methods look at two or more variables at a time to explore relationships. Usually our multivariate EDA will be bivariate (looking at exactly two variables), but occasionally it will involve three or more variables.

Univariate non-graphical EDA:

The data that come from making a particular measurement on all of the subjects in a sample represent our observations for a single characteristic such as age, gender, speed at a task, or response to a stimulus. We should think of these measurements as representing a “sample distribution” of the variable, which in turn more or less represents the “population distribution” of the variable. The usual goal of univariate non-graphical EDA is to better appreciate the “sample distribution” and also to make some tentative conclusions about what population distribution(s) is/are compatible with the sample distribution. Outlier detection is also a part of this analysis.

1. Categorical data
2. Quantitative data

Several statistics are commonly used as a measure of the spread of a distribution, including variance, standard deviation, and interquartile range. Spread is an indicator of how far away from the center we are still likely to find data values. The variance is a standard measure of spread. It is calculated for a list of numbers, e.g., the n observations of a particular measurement labeled x_1 through x_n , based on the n sample deviations (or just “deviations”). A third measure of spread is the interquartile range. To define IQR, we first need to define the concepts of quartiles. The quartiles of a population or a sample are the three values which divide the distribution or observed data into even fourths.

Skewness is a measure of asymmetry. Kurtosis is a more subtle measure of peakedness compared to a Gaussian distribution.

Univariate graphical EDA:

If we are focusing on data from observation of a single variable on n subjects, i.e., a sample of size n, then in addition to looking at the various sample statistics discussed in the previous section, we also need to look graphically at the distribution of the sample. Non-graphical and graphical methods complement each other.

1. Histogram
2. Stem and leaf graph - A stem and leaf plot shows all data values and the shape of the distribution.
3. Boxplots - Boxplots show robust measures of location and spread as well as providing information about symmetry and outliers.
4. Quantile-normal plots - Quantile-Normal plots allow detection of non-normality and diagnosis of skewness and kurtosis.

Multivariate non-graphical EDA:

Multivariate non-graphical EDA techniques generally show the relationship between two or more variables in the form of either cross-tabulation or statistics.

Multivariate graphical EDA:

There are few useful techniques for graphical EDA of two categorical random variables. The only one used commonly is a grouped barplot with each group representing one level of one of the variables and each bar within a group representing the levels of the other variable.

Output:

Code:

```
install.packages('aqp', dep=TRUE)
install.packages("soilDB")
```

```

library(aqp)
library(soilDB)

# Load from the the loakercreek dataset
data("loafercreek")

# Construct generalized horizon designations
n <- c("A",
       "BAt",
       "Bt1",
       "Bt2",
       "Cr",
       "R")
# REGEX rules
p <- c("A",
       "BA|AB",
       "Bt|Bw",
       "Bt3|Bt4|2B|C",
       "Cr",
       "R")

# Compute genhz labels and add to loafercreek dataset
loafercreek$genhz <- generalize.hz(loafercreek$hzname, n, p)

# Extract the horizon table
h <- horizons(loafercreek)

# Examine the matching of pairing of the genhz label to the hznames
table(h$genhz, h$hzname)
View(h)
vars <- c("genhz", "clay", "total_frags_pct", "phfield", "effclass")
summary(h[, vars])

# just for factors
levels(h$genhz)
# for characters and factors
sort(unique(h$hzname))

h$hzname <- ifelse(h$hzname == "BT", "Bt", h$hzname)

# or

h$hzname[h$hzname == "BT"] <- "Bt"

# or as a last resort we could manually edit the spreadsheet in R

edit(h)

# first remove missing values and create a new vector

```

```

clay <- na.exclude(h$clay)
mean(clay)

mean(h$clay, na.rm = TRUE)
median(clay)
sort(table(round(h$clay)), decreasing = TRUE)[1] # sort and select the 1st value, which will
be the mode
table(h$genhz)
table(h$genhz, h$texcl)

# append the table with row and column sums
addmargins(table(h$genhz, h$texcl))

# calculate the proportions relative to the rows, margin = 1 calculates for rows, margin = 2
calculates for columns, margin = NULL calculates for total observations
round(prop.table(table(h$genhz, h$texture_class), margin = 1) * 100)
knitr::kable(addmargins(table(h$genhz, h$texcl)))

aggregate(clay ~ genhz, data = h, mean)
aggregate(clay ~ genhz, data = h, median)
# or we could use the summary() function to get both the mean and median
aggregate(clay ~ genhz, data = h, summary)

var(h$clay, na.rm=TRUE)
sd(h$clay, na.rm = TRUE)
cv <- sd(clay) / mean(clay) * 100
cv

quantile(clay)
range(clay)
diff(range(clay))
IQR(clay)

install.packages("ggplot2")
library(ggplot2)

# bar plot
ggplot(h, aes(x = texcl)) +
  geom_bar()
ggplot(h, aes(x = clay)) +
  geom_histogram(bins = nclass.Sturges(h$clay))
ggplot(h, (aes(x = genhz, y = clay))) +
  geom_boxplot()

h$hzdepm <- (h$hzdepb + h$hzdept) / 2 # Compute the middle horizon depth

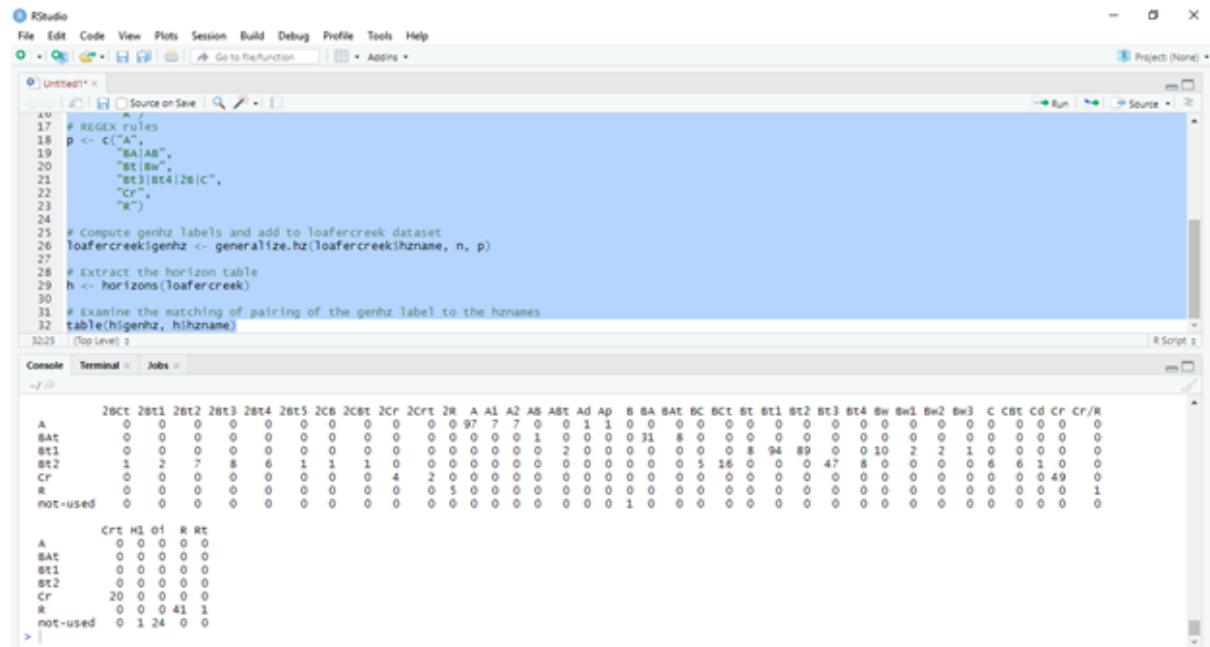
vars <- c("hzdepm", "clay", "sand", "total_frags_pct", "phfield")

round(cor(h[, vars], use = "complete.obs"), 2)

```

```
# scatter plot
ggplot(h, aes(x = clay, y = hzdepm)) +
  geom_point() +
  ylim(100, 0)
```

Screenshots:



```

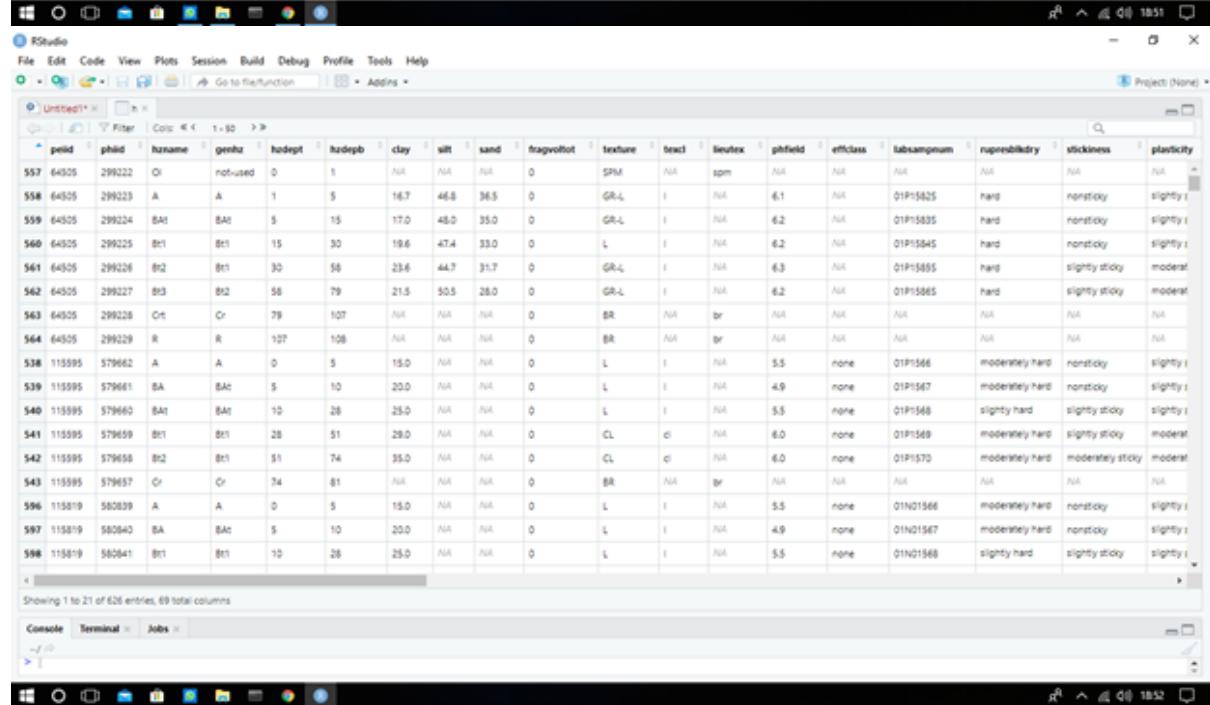
x.v
# REGEX rules
p <- c("A",
      "BA|A",
      "Bt|Bt",
      "Bt3)Bt4|2B|C",
      "Cr",
      "R")
24
25 # Compute genhz labels and add to loafercreek dataset
26 loafercreekgenhz <- generalize.hz(loafercreek$hzname, n, p)
27
28 # Extract the horizon table
29 h <- horizons(loafercreek)
30
31 # Examine the matching of pairing of the genhz label to the hznames
32 table(hgenhz, hzname)
32-25 (Top Level) s

```

Console Terminal Jobs

	2Bt1	2Bt2	2Bt3	2Bt4	2Bt5	2CB	2Cbt	2Cr	2Crt	2R	A	A1	A2	A8	ABt	Ad	Ap	B	BA	BAT	BC	Bct	Bt	Bt1	Bt2	Bt3	Bt4	Bw	Bw1	Bw2	Bw3	C	CBt	Cd	Cr	R/R						
A	0	0	0	0	0	0	0	0	0	0.97	7	7	0	0	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
BAE	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
Bt1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
Bt2	1	2	7	8	6	1	1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
Cr	0	0	0	0	0	0	0	0	0	0	4	2	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
R	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
not-used	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	
	Crt	Hl	Gf	R	Rt																																					
A	0	0	0	0	0																																					
BAE	0	0	0	0	0																																					
Bt1	0	0	0	0	0																																					
Bt2	0	0	0	0	0																																					
Cr	20	0	0	0	0																																					
R	0	0	0	0	41	1																																				
not-used	0	1	24	0	0																																					

> |



peid	phid	hzname	genhz	hdept	hzdepb	clay	silt	sand	fragvoltot	texture	feach	feutex	phifield	effclass	labsamplenum	resistiblity	thickness	plasticity	
557	64905	299222	Ol	not-used	0	1	NA	NA	NA	0	SPM	NA	SPM	NA	NA	NA	NA	NA	NA
558	64905	299223	A	A	1	5	16.7	46.8	36.5	0	GR-L	I	NA	6.1	NA	01P15825	hard	nonsticky	slightly i
559	64905	299224	BAE	BAE	5	15	17.0	48.0	35.0	0	GR-L	I	NA	6.2	NA	01P15835	hard	nonsticky	slightly i
560	64905	299225	Bt1	Bt1	15	30	19.6	47.4	33.0	0	L	I	NA	6.2	NA	01P15845	hard	nonsticky	slightly i
561	64905	299226	Bt2	Bt1	30	56	23.6	44.7	31.7	0	GR-L	I	NA	6.3	NA	01P15855	hard	slightly sticky	moderat
562	64905	299227	Bt3	Bt2	58	79	21.5	50.5	28.0	0	GR-L	I	NA	6.2	NA	01P15865	hard	slightly sticky	moderat
563	64905	299228	Or	Or	79	107	NA	NA	NA	0	BR	NA	BR	NA	NA	NA	NA	NA	NA
564	64905	299229	R	R	107	108	NA	NA	NA	0	BR	NA	BR	NA	NA	NA	NA	NA	NA
538	115995	579662	A	A	0	5	15.0	NA	NA	0	L	I	NA	5.5	none	01P1566	moderately hard	nonsticky	slightly i
539	115995	579661	BA	BAE	5	10	20.0	NA	NA	0	L	I	NA	4.9	none	01P1567	moderately hard	nonsticky	slightly i
540	115995	579660	Bt1	Bt1	10	28	25.0	NA	NA	0	L	I	NA	5.5	none	01P1568	slightly hard	slightly sticky	slightly i
541	115995	579659	Bt1	Bt1	28	51	29.0	NA	NA	0	CL	CI	NA	6.0	none	01P1569	moderately hard	slightly sticky	moderat
542	115995	579658	Bt2	Bt1	51	74	35.0	NA	NA	0	CL	CI	NA	6.0	none	01P1570	moderately hard	moderately sticky	moderat
543	115995	579657	Or	Or	74	81	NA	NA	NA	0	BR	NA	BR	NA	NA	NA	NA	NA	NA
546	115995	580839	A	A	0	5	15.0	NA	NA	0	L	I	NA	5.5	none	01N01566	moderately hard	nonsticky	slightly i
547	115995	580640	BA	BAE	5	10	20.0	NA	NA	0	L	I	NA	4.9	none	01N01567	moderately hard	nonsticky	slightly i
548	115995	580641	Bt1	Bt1	10	28	25.0	NA	NA	0	L	I	NA	5.5	none	01N01568	slightly hard	slightly sticky	slightly i

Showing 1 to 21 of 626 entries, 69 total columns

Console Terminal Jobs

> |

RStudio

```

File Edit Code View Plots Session Build Debug Profile Tools Help
Project (None)
Untitled1.R Source on Save Run Source
13 "Bt1"
14 "Bt2"
15 "Cr"
16 "R")
17 # REGEX rules
18 p <- c("A",
19 "BA|AB",
20 "Bt|Bw",
21 "Bt3|Bt4|2B|C",
22 "Cr",
23 "R")
24
25 # Compute genhz labels and add to loafercreek dataset
26 loafercreekgenhz <- generalize.hz(loafercreek$hzname, n, p)
27
28 # Extract the horizon table
29 h <- horizons(loafercreek)
30
31 # Examine the matching of pairing of the genhz label to the hznames
32 table(h$genhz, hzname)
33 view(h)
34 vars <- c("genhz", "clay", "total frags_pct", "phfield", "effclass")
35 summary(h[, vars])
35:19 (Top Level) s

```

Console Terminal Jobs

```

> view(h)
> vars <- c("genhz", "clay", "total frags_pct", "phfield", "effclass")
> summary(h[, vars])
   genhz      clay    total frags_pct    phfield    effclass
A:113 Min. :10.00  Min. :0.00  Min. :14.00  very slight: 0
BAT :40  1st Qu.:18.00  1st Qu.:0.00  1st Qu.:16.00  slight : 0
Bt1 :208 Median :22.00  median :5.00  Median :16.30  strong : 0
Bt2 :116 Mean  :23.67  mean  :14.18  Mean  :16.18  violent : 0
Cr  : 5  3rd Qu.:28.00  3rd Qu.:20.00  3rd Qu.:16.50  none   : 86
R   :48  Max. :60.00  Max. :95.00  Max. :7.00   NA's   :540
not-used:26  NA's  :173  NA's  :381

```

RStudio

```

File Edit Code View Plots Session Build Debug Profile Tools Help
Project (None)
Untitled1.R Source on Save Run Source
28 # Extract the horizon table
29 h <- horizons(loafercreek)
30
31 # examine the matching of pairing of the genhz label to the hznames
32 table(h$genhz, hzname)
33 view(h)
34 vars <- c("genhz", "clay", "total frags_pct", "phfield", "effclass")
35 summary(h[, vars])
36
37 # just for factors
38 levels(h$genhz)
39 # for characters and factors
40 sort(unique(h$hzname))
41
42 h$hzname <- ifelse(h$hzname == "BT", "Bt", h$hzname)
43
44 # or
45
46 h$hzname[h$hzname == "BT"] <- "Bt"
47
48 # or as a last resort we could manually edit the spreadsheet in R
49
50 edit(h)
50:19 (Top Level) s

```

Data Editor

row.names	pelid	phid	hzname	genhz	hzdept	hzdegb	clay	silt
1	557	64505	299222	O1	not-used	0	1	NA
2	558	64505	299223	A	A	1	5	16.7
3	559	64505	299224	BAt	BAt	5	15	49
4	560	64505	299225	Bt1	Bt1	15	30	19.4
5	561	64505	299226	Bt2	Bt1	30	58	23.6
6	562	64505	299227	Bt3	Bt2	58	79	21.5
7	563	64505	299228	Crt	Crt	79	107	NA
8	564	64505	299229	R	R	107	108	NA
9	538	115595	579662	A	A	0	5	15
10	539	115595	579661	NA	NA	5	10	NA
11	540	115595	579660	BAt	BAt	10	28	25
12	541	115595	579659	Bt1	Bt1	28	51	29
13	542	115595	579658	Bt2	Bt1	51	74	35
14	543	115595	579657	Cr	Cr	74	81	NA
15	596	115819	580839	A	A	0	5	15
16	597	115819	580840	BA	BAt	5	10	NA
17	598	115819	580841	Bt1	Bt1	10	28	25
18	599	115819	580842	Bt2	Bt1	28	51	29
19	600	115819	580843	Bt3	Bt2	51	74	35

Console Terminal Jobs

```

> levels(h$genhz)
[1] "A"      "BAt"    "Bt1"    "Bt2"    "Cr"     "R"     "not-used"
> sort(unique(h$hzname))
[1] "2Bt"    "2Bt1"   "2Bt2"   "2Bt3"   "2Bt4"   "2Bt5"   "2C8"    "2C8t"   "2Cr"    "2Crt"   "2R"    "A"      "A1"    "A2"    "AB"    "ABt"   "Ad"    "Ap"    "B"      "BA"    "BAT"   "BC"
[21] "Bt"     "Bt1"    "Bt2"    "Bt3"    "Bt4"    "Bw"    "Bw1"   "Bw2"    "Bw3"    "C"     "C8t"   "Cd"    "Cr"    "Cr/R"  "Crt"   "H1"    "O1"    "R"      "RA"    "Rt"
> # or
>

```

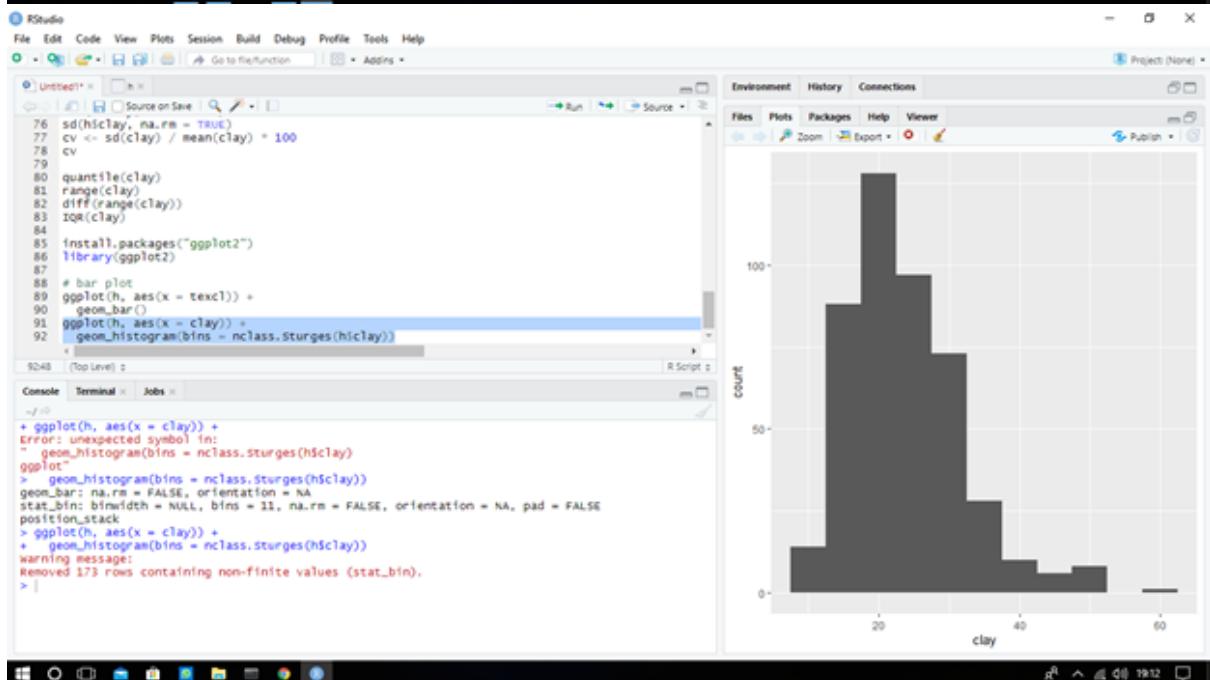
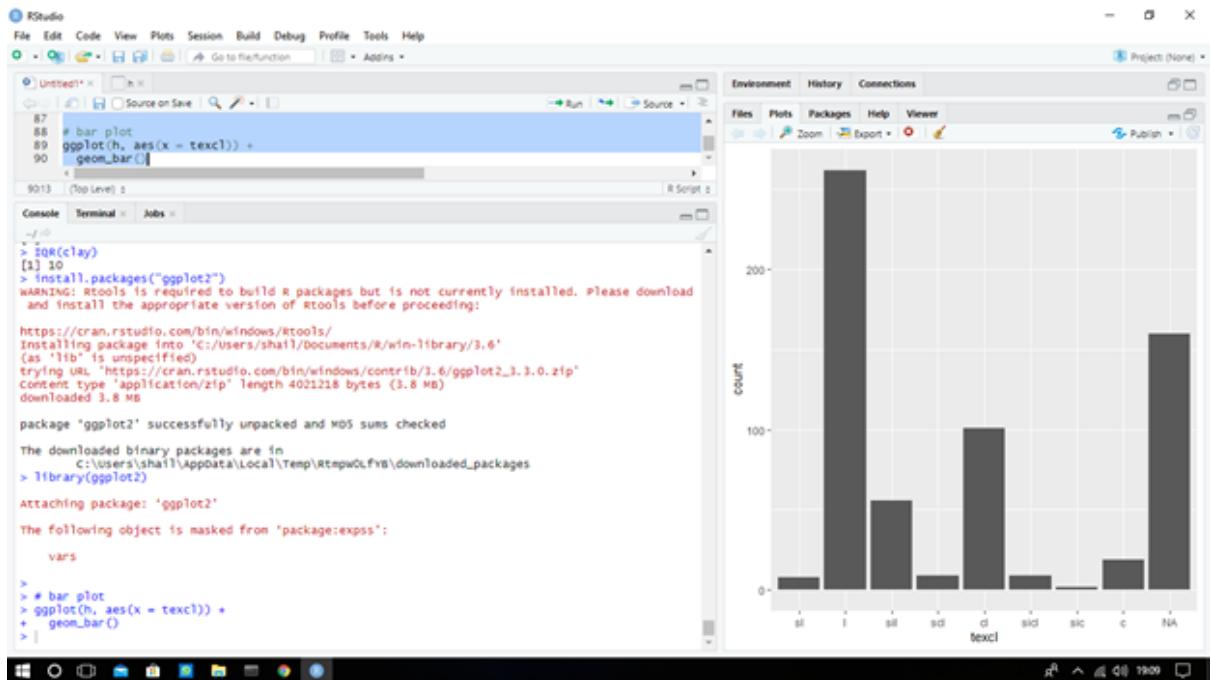

The screenshot shows two instances of RStudio running side-by-side. Both instances have the same session history and environment. The top instance has a blue header bar, while the bottom one has a white header bar. The code being run is as follows:

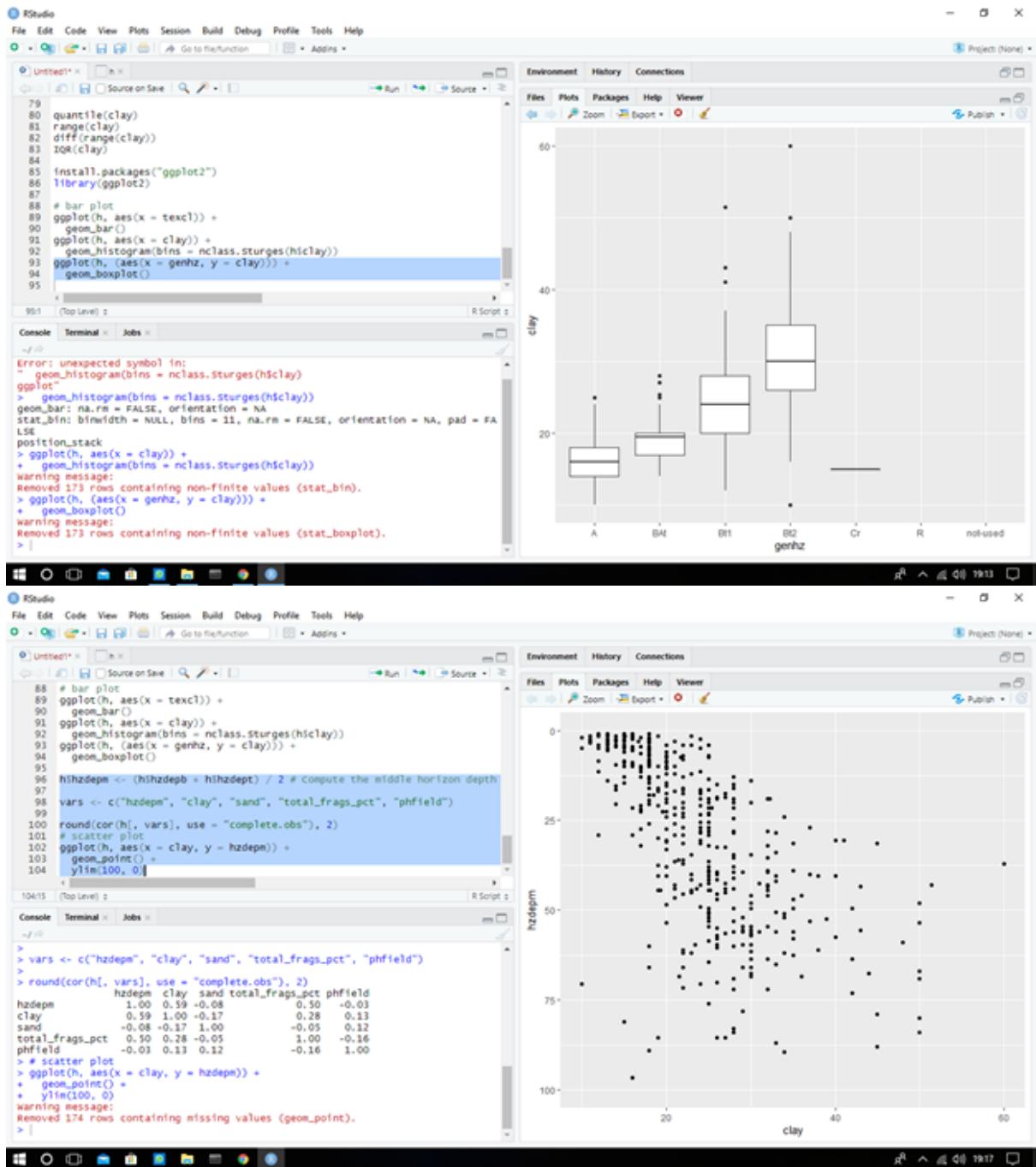
```
r> # aggregate(clay ~ genhz, data = h, mean)
genhz <- clay
1 A 16.23113
2 BAT 19.51889
3 BT1 24.14221
4 BT2 31.35045
5 CR 15.00000
> # or we could use the summary() function to get both the mean and median
> aggregate(clay ~ genhz, data = h, summary)
genhz clay.min. clay.1st qu. clay.Median clay.mean. clay.3rd qu. clay.Max.
1 A 10.00000 14.00000 16.23113 18.00000 25.00000
2 BAT 14.00000 17.00000 19.50000 19.51889 20.00000 28.00000
3 BT1 12.00000 20.00000 24.00000 24.14221 28.00000 51.40000
4 BT2 10.00000 26.00000 30.00000 31.35045 35.00000 60.00000
5 CR 15.00000 15.00000 15.00000 15.00000 15.00000 15.00000
>
> var(hiclay, na.rm=TRUE)
[1] 64.89187
> sd(hiclay, na.rm = TRUE)
[1] 8.055549
> cv <- sd(clay) / mean(clay) * 100
> cv
[1] 34.03087
>
```

The bottom instance continues the session with:

```
79 | quantile(clay)
80 | range(clay)
81 | diff(range(clay))
82 |
83 |
791 | (Top Level) :
```

Then it runs the same aggregate and variance calculations again, followed by quantile, range, and diff range calculations.





Conclusion: Thus, we have studied EDA in R.

Experiment No. 8

Aim : Basic and Advanced graphics in R

Theory :

One of the main reasons data analysts turn to R is for its strong graphic capabilities. These include density plots (histograms and kernel density plots), dot plots, bar charts (simple, stacked, grouped), line charts, pie charts (simple, annotated, 3D), box plots (simple, notched, violin plots, bagplots) and Scatter Plots (simple, with fit lines, scatterplot matrices, high density plots, and 3D plots).

The R base function plot() can be used to create graphs. Some of them are :

- A scatter plot can be created using the function plot(x, y). The function lm() will be used to fit linear models between y and x. A regression line will be added on the plot using the function abline(), which takes the output of lm() as an argument. You can also add a smoothing line using the function loess().
- 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.
- A histogram represents the frequencies of values of a variable bucketed into ranges. Histogram is similar to a 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 a histogram using hist() function. This function takes a vector as an input and uses some more parameters to plot histograms.
- A line chart is a graph that connects a series of points by drawing line segments between them. These points are ordered in one of their coordinate (usually the x-coordinate) values. Line charts are usually used in identifying the trends in data. The plot() function in R is used to create the line graph.
- A bar chart represents data in rectangular bars with length of the bar proportional to the value of the variable. R uses the function barplot() to create bar charts. R can draw both vertical and Horizontal bars in the bar chart. In the bar chart each of the bars can be given different colors.
- A pie-chart is a representation of values as slices of a circle with different colors. The slices are labeled and the numbers corresponding to each slice is also represented in the chart. In R the pie chart is created using the pie() function which takes positive numbers as a vector input.

Output :

Code :

```

City = read.csv("http://course1.winona.edu/bdeppa/DSCI%20415/Data/City77.csv")
# Copy city names to row names, i.e. observation labels
row.names(City) = City$City
# Remove city name column from the data frame
City = City[,-1]
names(City)

hist(City$medinc,prob=T,nclass=20,col="pink",main="Histogram of Median
Income",xlab="Median Income")
lines(density(City$medinc),lty=2,col="red",lwd=2)

require(ggplot2)

a = ggplot(City,aes(medinc))
a + geom_histogram(binwidth=1000,aes(y=..density..),fill="pink",color="black") +
  geom_density(linetype=2) +
  xlab("Median Income") +
  ggtitle("Median Income for Top 77 Cities")

plot(City$poverty,City$infmort,xlab="Percent Below Poverty Level",ylab="Infant Mortality
Rate")
lines(lowess(City$poverty,City$infmort),lty=2,col="pink",lwd=2)
abline(lm(infmort~poverty,data=City),lwd=2)
# identify(City$poverty,City$infmort,labels=row.names(City))
title(main="Infant Morality vs. Poverty Level")

a = ggplot(City,aes(poverty,infmort))
a + geom_point() + geom_smooth() +
  geom_text(aes(label=row.names(City)),size=2) +
  xlab("Poverty") + ylab("Infant Mortality") +
  ggtitle("Infant Mortality vs. Poverty Rate")

Olives = read.csv("http://course1.winona.edu/bdeppa/DSCI%20415/Data/Olives.txt")
names(Olives)

par(mfrow=c(3,2)) # sets a plotting region for with 3 rows and 2 columns
boxplot(split(Olives$Oleic,Olives$Area.Name),xlab="Area",ylab="Oleic Acid")
boxplot(split(Olives$Linoleic,Olives$Area.Name),xlab="Area",ylab="Linoleic Acid")
boxplot(split(Olives$Stearic,Olives$Area.Name),xlab="Area",ylab="Stearic Acid")
boxplot(split(Olives$Linolenic,Olives$Area.Name),xlab="Area",ylab="Linolenic Acid")
boxplot(split(Olives$Eicosanoic,Olives$Area.Name),xlab="Area",ylab="Eicosanoic Acid")
boxplot(split(Olives$Eicosenoic,Olives$Area.Name),xlab="Area",ylab="Eicosenoic Acid")

require(s20x)

boxqq(Oleic~Region.Name,data=Olives)

a = ggplot(Olives,aes(Area.Name,Oleic))
a + geom_boxplot(fill="lightblue") + xlab("Area Name") +
  ylab("Oleic Acid") + ggtitle("Oleic Acid vs. Growing Area")

```

```

require(violinmplot)

violinmplot(Area.Name~Oleic,xlab="Oleic Acid",ylab="Area
Name",data=Olives,horizontal=T)

a = ggplot(Olives,aes(Area.Name,Oleic))
a + geom_violin(fill="lightblue") + xlab("Area Name") +
ylab("Oleic Acid") + ggtitle("Oleic Acid vs. Growing Area")

data("UCBAdmissions")
UCBAdmissions

temp = data.frame(UCBAdmissions)
temp

DeptCount = margin.table(UCBAdmissions,3)
DeptCount

par(mfrow=c(1,2))
pie(DeptCount)
barplot(DeptCount,xlab="Department",ylab="Frequency")

par(mfrow=c(1,2))
DeptAdmit = margin.table(UCBAdmissions,c(1,3))
DeptAdmit

barplot(DeptAdmit,xlab="Department",ylab="Frequency")
barplot(DeptAdmit,xlab="Department",ylab="Frequency",beside=TRUE)

GenderAdmit = margin.table(UCBAdmissions,c(1,2))
GenderAdmit

pie(GenderAdmit[,1],main="Males")
pie(GenderAdmit[,2],main="Females")

barplot(GenderAdmit,xlab="Gender",ylab="Frequency")
barplot(GenderAdmit,xlab="Gender",ylab="Frequency",beside=T)

par(mfrow=c(1,2))
mosaicplot(~Gender+Admit,data=UCBAdmissions,color=T)
mosaicplot(~Admit+Gender,data=UCBAdmissions,color=T)

par(mfrow=c(1,2))
mosaicplot(~Dept+Admit,data=UCBAdmissions,color=T)
mosaicplot(~Admit+Dept,data=UCBAdmissions,color=T)

mosaicplot(~Dept+Gender+Admit,data=UCBAdmissions,color=T)

```

```

mosaicplot(~Dept+Admit+Gender,data=UCBAdmissions,color=T)

mosaicplot(~Gender+Dept+Admit,data=UCBAdmissions,color=T)

mosaicplot(~Admit+Dept+Gender,data=UCBAdmissions,color=T)

names(Olives)

olive.mat = Olives[,c(7,6,5,3)]
pairs(olive.mat)

panel.cor = function (x, y, digits = 2, prefix = "", cex.cor)
{
  usr <- par("usr")
  on.exit(par(usr))
  par(usr = c(0, 1, 0, 1))
  r <- abs(cor(x, y))
  txt <- format(c(r, 0.123456789), digits = digits)[1]
  txt <- paste(prefix, txt, sep = "")
  if (missing(cex.cor))
    cex <- 0.8/strwidth(txt)
  text(0.5, 0.5, txt, cex = cex * r)
}

panel.trendsd = function (x, y, f = 0.5)
{
  par(err = -1)
  xs <- sort(x, index = T)
  x <- xs$x
  ix <- xs$ix
  y <- y[ix]
  trend <- lowess(x, y, f)
  e2 <- (y - trend$y)^2
  scatter <- lowess(x, e2)
  uplim <- trend$y + sqrt(abs(scatter$y))
  lowlim <- trend$y - sqrt(abs(scatter$y))
  points(x, y, pch = 1)
  lines(trend, col = "Blue")
  lines(scatter$x, uplim, lty = 2, col = "Red")
  lines(scatter$x, lowlim, lty = 2, col = "Red")
  invisible()
}

panel.hist = function (x, ...)
{
  usr <- par("usr")
  on.exit(par(usr))
  par(usr = c(usr[1:2], 0, 1.5))
  h <- hist(x, plot = FALSE)
}

```

```

breaks <- h$breaks
nB <- length(breaks)
y <- h$counts
y <- y/max(y)
rect(breaks[-nB], 0, breaks[-1], y, col = "cyan", ...)
}

pairs.trendsd = function(data,...) {
  pairs(data,lower.panel=panel.cor,upper.panel=panel.trendsd,
    diag.panel=panel.hist,...)}

require(lattice)
pairs.trendsd(olive.mat)

Boston = read.csv("http://course1.winona.edu/bdeppa/DSCI%20415/Data/Boston.txt")
names(Boston)

require(corrgram)
corrgram(Boston[,-c(1:4)],lower.panel=panel.pts,upper.panel=panel.pie)
corrgram(Boston[,-c(1:4)],lower.panel=panel.shade,upper.panel=panel.pie)

require(corrplot)
Boston.corr = cor(Boston[,-c(1:4)]) # find pairwise correlations between all variables.
options(digits=2) # reduce the number of significant digits shown.
Boston.corr

corrplot(Boston.corr)

corrplot(Boston.corr,method="ellipse")

corrplot(Boston.corr,order="FPC")

corrplot(Boston.corr,order="hclust")

data(iris)
names(iris)

library(lattice)
xyplot(Sepal.Length~Petal.Length|Species,data=iris)

xyplot(Sepal.Length~Petal.Length|Species,layout=c(2,2),data=iris)

SepWid = equal.count(iris$Sepal.Width)
plot(SepWid)

print(SepWid)

```

```

xyplot(Sepal.Length~Petal.Length|SepWid,data=iris)

xyplot(Sepal.Length~Petal.Length|SepWid,groups=Species,layout=c(3,2),
       auto.key=list(columns=3),main="Sepal Length * Petal Length | Sepal
Width",data=iris)

splom(~iris[,1:4],groups=Species,auto.key=list(columns=3),data=iris)

bwplot(Species~Sepal.Width,data=iris,xlab="Sepal Width (mm)",main="Boxplots of Iris
Sepal Width")

stripplot(Species ~ jitter(Sepal.Width), data = iris, xlab="Sepal Width (mm)",main="Sepal
Width Across Species")

stripplot(Species ~ jitter(Sepal.Width), data = iris, aspect = 1,
          jitter=T,xlab="Sepal Width (mm)",main="Sepal Width Across Species")

coplot(Sepal.Width~Sepal.Length|Species,data=iris)

coplot(Sepal.Width~Sepal.Length|Petal.Width,number=4,overlap=.2,data=iris)

coplot(Sepal.Width~Sepal.Length|Petal.Width*Petal.Length,number=c(3,3),
       overlap=.5,col=as.numeric(iris$Species),pch=as.numeric(iris$Species)+1,data=iris)

Ozdata = read.csv("http://course1.winona.edu/bdeppa/DSCI%20415/Data/Ozdata.csv")
names(Ozdata)

coplot(upoz~safb|inbh*v500,number=c(4,4),panel=panel.smooth,data=Ozdata)

coplot(upoz~safb|inbh*v500,number=c(4,4),overlap=.25,panel=function(x,y,...)
       panel.smooth(x,y,span=.6,...),data=Ozdata)

library(ash)
Swiss = read.csv("http://course1.winona.edu/bdeppa/DSCI%20415/Data/Swiss.csv")
names(Swiss)

d1 = ash1(bin1(Swiss$diagon,nbin=50),3)
hist(Swiss$diagon,nclass=20,prob=T,col="pink",main="Histogram of Bill Diagonal (mm)")
lines(d1)

diagrt.bin <- bin2(cbind(Swiss$diagon,Swiss$right),nbin=c(50,50))
diagrt.1 <- ash2(diagrt.bin,m=c(5,5))
persp(diagrt.1,xlab="Diagonal Length",ylab="Right Height",zlab="",cex=.5,theta=-
45,phi=30,shade=1,col="pink")

image(diagrt.1,xlab="Bill Diagonal",ylab="Right Height")

```

```

contour(diagrt.1,xlab="Bill Diagonal",ylab="Right Height",add=T)
points(Swiss$diagon,Swiss$right,pch=as.character(Swiss$genu),cex=.4)

Genuine <- Swiss$genu
Genuine[Genuine==0] <- "Forged"
Genuine[Genuine==1] <- "Real"
xyplot(Swiss$diagon~Swiss$bottom|Genuine,groups=Genuine)

splom(~Swiss[,1:6],groups=Genuine,auto.key=list(columns=2))

pairs.image <- function(x) {
  pairs(x,panel=function(x,y) {
    foo <- bin2(cbind(x,y),nbins=c(75,75))
    foo <- ash2(foo,m=c(6,6))
    image(foo,add=T,xlab="",ylab="",col=topo.colors(1000))
    points(x,y,pch=".")
  })
}
pairs.image(Swiss[,1:6])

pairs.persp <- function(x) {
  par(bg="sky blue")
  pairs(x,panel=function(x,y) {
    foo <- bin2(cbind(x,y),nbins=c(75,75))
    foo <- ash2(foo,m=c(8,8))
    par(new=T)
    persp(foo,xlab="",ylab="",theta=-45,phi=35,col="red",
          shade=.75,box=F,scale=F,border=NA,expand=.9)
  })
  par(new=F,bg="white")
}
pairs.persp(Swiss[,1:3])

NHL = read.csv("http://course1.winona.edu/bdeppa/DSCI%20415/Data/NHL.csv")
names(NHL)

plot(NHL$GF,NHL$GA,xlab="Goals For",ylab="Goals Against",main="Plot of GA vs. GF
with Symbols = Wins")
symbols(NHL$GF,NHL$GA,circles=NHL$W,add=T,inches=.25) # circles are too big if
inches is not specified!
symbols(NHL$GF,NHL$GA,circles=NHL$W,add=T,bg="lightblue",inches=.25) # bg
option colors circles light blue
text(NHL$GF,NHL$GA,labels=NHL$TEAM,cex=.6)

kodiak.crab = read.csv("http://course1.winona.edu/bdeppa/DSCI%20415/Data/kodiak.csv")
survey.crab = read.csv("http://course1.winona.edu/bdeppa/DSCI%20415/Data/survey.csv")
attach(kodiak.crab) # This file contains latitude and longitude of the island coastline.
attach(survey.crab) # This file contains crab counts by type, year, and location (lat, long)

names(survey.crab)

```

```

table(year)

crab.plot <- function(lat,long,y,year,yname=deparse(substitute(y))){
  symbols(long[year==year],lat[year==year],circles=y[year==year],inches=.1,bg="blue",xlab=
  "Longitude",ylab="Latitude",
  main=paste("Circles Proportional to",yname))
  title(sub= paste("Year =",year))
  points(kodiak.crab$longit,kodiak.crab$latit,pch=".")
}

nr.pp = nrec/npots
crab.plot(lat,long,y=nr.pp,year=83,yname="Number of Recruits per Pot")

detach(survey.crab)
detach(kodiak.crab) # if you attach, be sure to detach when finished!

attach(City)      # This command readies the data frame City for analysis.

names(City)       # The names command lists the names of the variables in a data set.

View(City)

city.mat <-cbind(popdens,pctblack,pcthisp,medinc,pct.pa,poverty,infmort,unempl)
stars(city.mat, key.loc = c(15,1.25),main = "Starplots for U.S.
Cities",label=row.names(City),cex=.7)

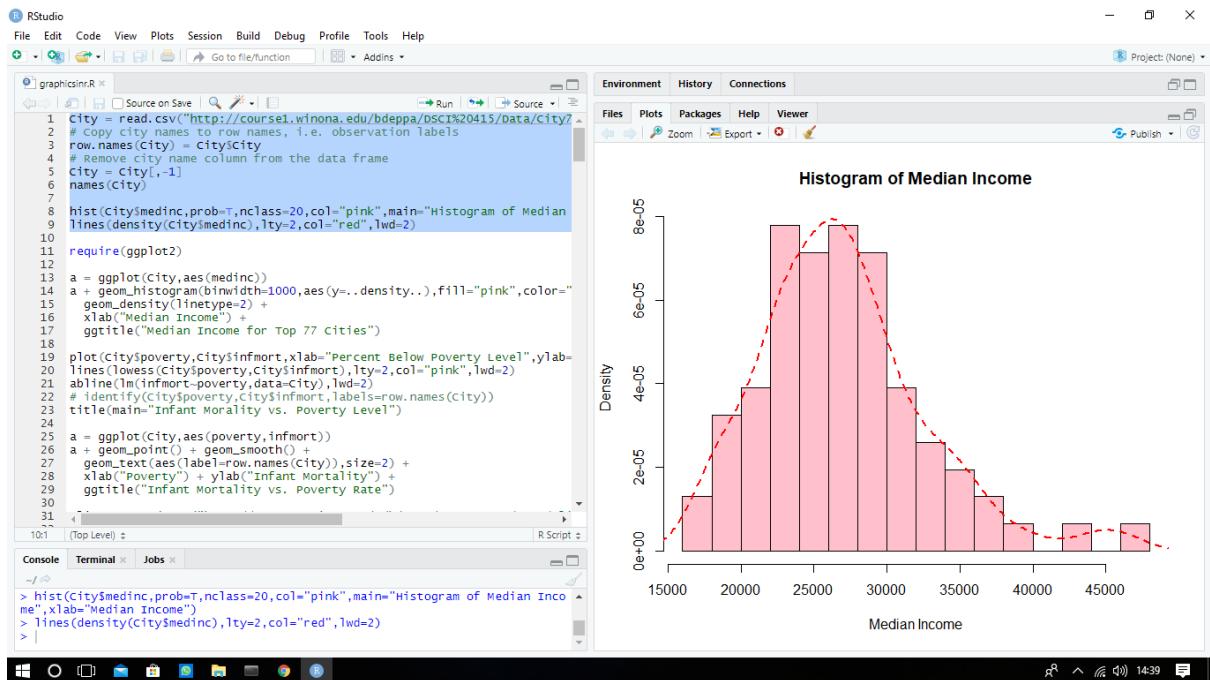
palette(rainbow(12, s = 0.6, v = 0.75))
stars(city.mat,len=.80,key.loc = c(15,1.25),main = "Starplots for U.S. Cities",draw.segments
= TRUE,label=row.names(City),cex=.7)

palette("default")

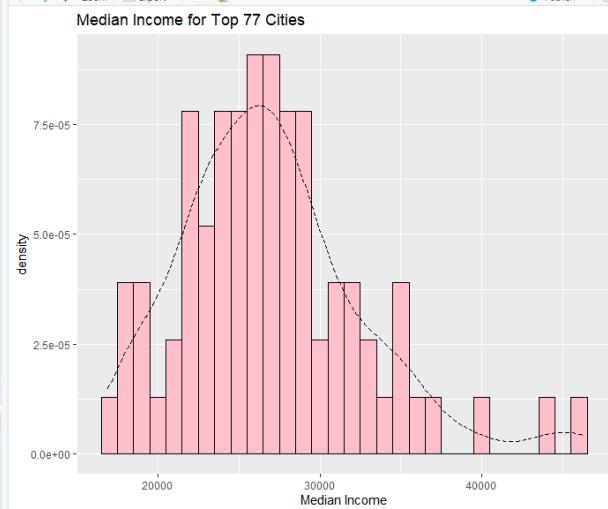
library(MASS)

parcoord(city.mat)
names(Olives)
olive.mat = Olives[,3:9]
palette(rainbow(12,s=0.6,v=0.75))
parcoord(olive.mat,col=as.numeric(Olives$Area.Name))
Screenshots :

```



The figure shows a histogram of Median Income for the top 77 cities. The x-axis is labeled "Median Income" and ranges from approximately 20,000 to 45,000. The y-axis is labeled "density" and ranges from 0.0e+00 to 7.5e-05. The histogram bars are pink, and a smooth density curve is overlaid in black. The distribution is roughly bell-shaped, peaking around 28,000.



The figure shows a screenshot of the RStudio interface. The left pane displays an R script titled "graphics.r" with code for reading a CSV file, creating a histogram, and plotting infant mortality against poverty level. The right pane shows a scatter plot titled "Infant Morality vs. Poverty Level" with a linear regression line and a dashed trend line.

```
City = read.csv("http://course.winona.edu/bdeppa/DSCI320415/data/city7")
# Copy city names to row names, i.e. observation labels
row.names(City) = City$city
# Remove city name column from the data frame
City = City[, -1]
names(City)
hist(City$medinc, prob=T, nclass=20, col="pink", main="Histogram of Median Income")
lines(density(City$medinc), lty=2, col="red", lwd=2)
require(ggplot2)
a = ggplot(City, aes(medinc))
a + geom_histogram(binwidth=1000, aes(y=..density..), fill="pink", color="black")
geom_density(linetype=2) +
  xlab("Median Income") +
  ylab("Density") +
  ggtitle("Median Income for Top 77 Cities")
plot(City$poverty, City$infmort, xlab="Percent Below Poverty Level", ylab="Infant Mortality Rate")
lines(lm(City$poverty, City$infmort), lty=2, col="pink", lwd=2)
abline(lm(infmort~poverty, data=City), lwd=2)
# identify(City$poverty, City$infmort, labels=row.names(City))
title(main="Infant Mortality vs. Poverty Level")
a = ggplot(City, aes(poverty, infmort))
a + geom_point() + geom_smooth()
geom_text(aes(label=row.names(City)), size=2) +
  xlab("Poverty") + ylab("Infant Mortality") +
  ggtitle("Infant Mortality vs. Poverty Rate")

```

RStudio

File Edit Code View Plots Session Build Debug Profile Tools Help

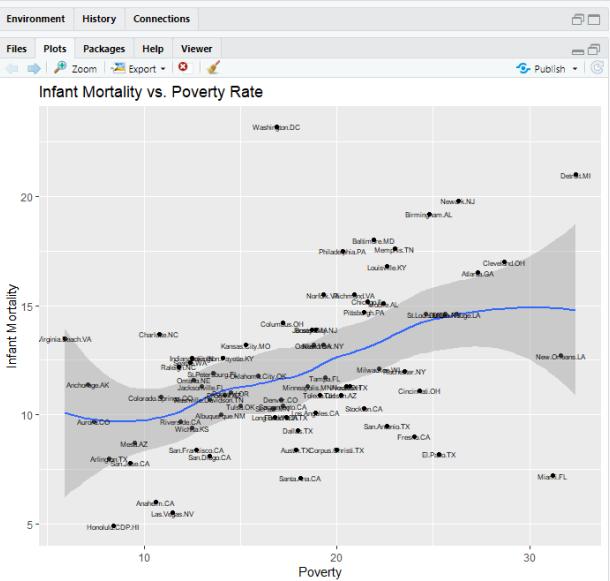
graphicsinst.R

```
23 title(main="Infant Mortality vs. Poverty Level")
24
25 a = ggplot(city,aes(poverty,infmort))
26 a + geom_point() + geom_smooth()
27 geom_text(aes(label=rownames(city)),size=2) +
28 xlab("Poverty") + ylab("Infant Mortality") +
29 ggtitle("Infant Mortality vs. Poverty Rate")
30
31 Olives = read.csv("http://course1.winona.edu/bdeppa/DSCI2041S/data/oil13names.csv")
32
33 par(mfrow=c(3,2)) # sets a plotting region for with 3 rows and 2 columns
34 boxplot(split(Olives$oleic,Olives$Area.Name),xlab="Area",ylab="Oleic Acid")
35 boxplot(split(Olives$Linoleic,Olives$Area.Name),xlab="Area",ylab="Linoleic Acid")
36 boxplot(split(Olives$Stearic,Olives$Area.Name),xlab="Area",ylab="Stearic Acid")
37 boxplot(split(Olives$Linolenic,Olives$Area.Name),xlab="Area",ylab="Linolenic Acid")
38 boxplot(split(Olives$Eicosanoic,Olives$Area.Name),xlab="Area",ylab="Eicosanoic Acid")
39 boxplot(split(Olives$Eicosenoic,Olives$Area.Name),xlab="Area",ylab="Eicosenoic Acid")
40 boxplot(split(Olives$Eicosenoic,Olives$Area.Name),xlab="Area",ylab="Eicosenoic Acid")
41
42 require(s20x)
43
44 boxqq(oleic~Region.Name,data=Olives)
45
46 a = ggplot(Olives,aes(Area.Name,oleic))
47 a + geom_boxplot(fill="lightblue") + xlab("Area Name") +
48 ylab("Oleic Acid") + ggtitle("Oleic Acid vs. Growing Area")
49
50 require(violinomplot)
51
52 violinomplot(Area.Name~oleic,xlab="oleic Acid",ylab="Area Name",data=Olives)
53
54
```

41:1 (Top Level) ↵ R Script ↵

Console Terminal ↵ Jobs ↵

```
oil("acid")
> boxplot(split(Olives$Eicosanoic,Olives$Area.Name),xlab="Area",ylab="Eicosanoic Acid")
> |
```



RStudio

File Edit Code View Plots Session Build Debug Profile Tools Help

graphcsm.R

Go to file/function Addins

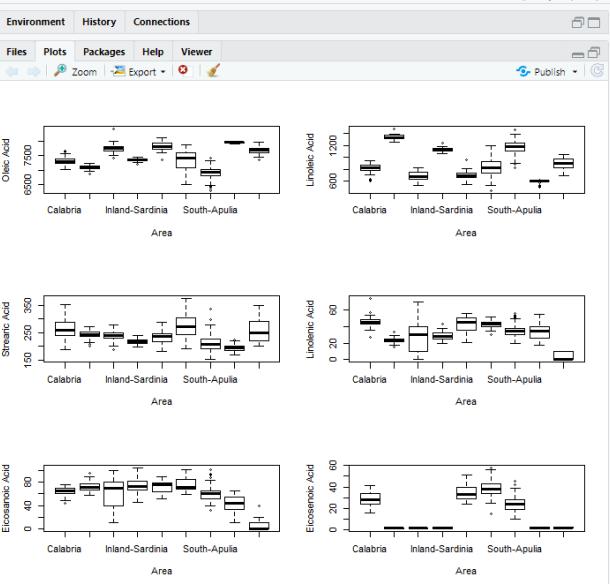
```
23 title(main="Infant Mortality vs. Poverty Level")
24
25 a = ggplot(city,aes(poverty,infmort))
26 a + geom_point() + geom_smooth()
27 geom_text(aes(label=rownames(city)),size=2) +
28 xlab("Poverty") + ylab("Infant Mortality") +
29 ggtitle("Infant Mortality vs. Poverty Rate")
30
31 olives = read.csv("http://course1.winona.edu/bdeppa/DSCI20415/Data/olives.csv")
32 names(olives)
33
34 par(mfrow=c(3,2)) # sets a plotting region for with 3 rows and 2 columns
35 boxplot(split(olives$oleic,olives$Area.Name),xlab="Area",ylab="Oleic Acid")
36 boxplot(split(olives$linoleic,olives$Area.Name),xlab="Area",ylab="Linoleic Acid")
37 boxplot(split(olives$stearic,olives$Area.Name),xlab="Area",ylab="Stearic Acid")
38 boxplot(split(olives$linolenic,olives$Area.Name),xlab="Area",ylab="Linolenic Acid")
39 boxplot(split(olives$Eicosanoic,olives$Area.Name),xlab="Area",ylab="Eicosanoic Acid")
40 boxplot(split(olives$Eicosenoic,olives$Area.Name),xlab="Area",ylab="Eicosenoic Acid")
41 |
42 require(s20x)
43
44 boxqq(oleic~Region.Name,data=olives)
45
46 a = ggplot(olives,aes(Area.Name,oleic))
47 a + geom_boxplot(fill="lightblue") + xlab("Area Name") +
48 ylab("Oleic Acid") + ggtitle("Oleic Acid vs. Growing Area")
49
50 require(violinplot)
51
52 violinplot(Area.Name~oleic,xlab="oleic Acid",ylab="Area Name",data=olives)
53 |
```

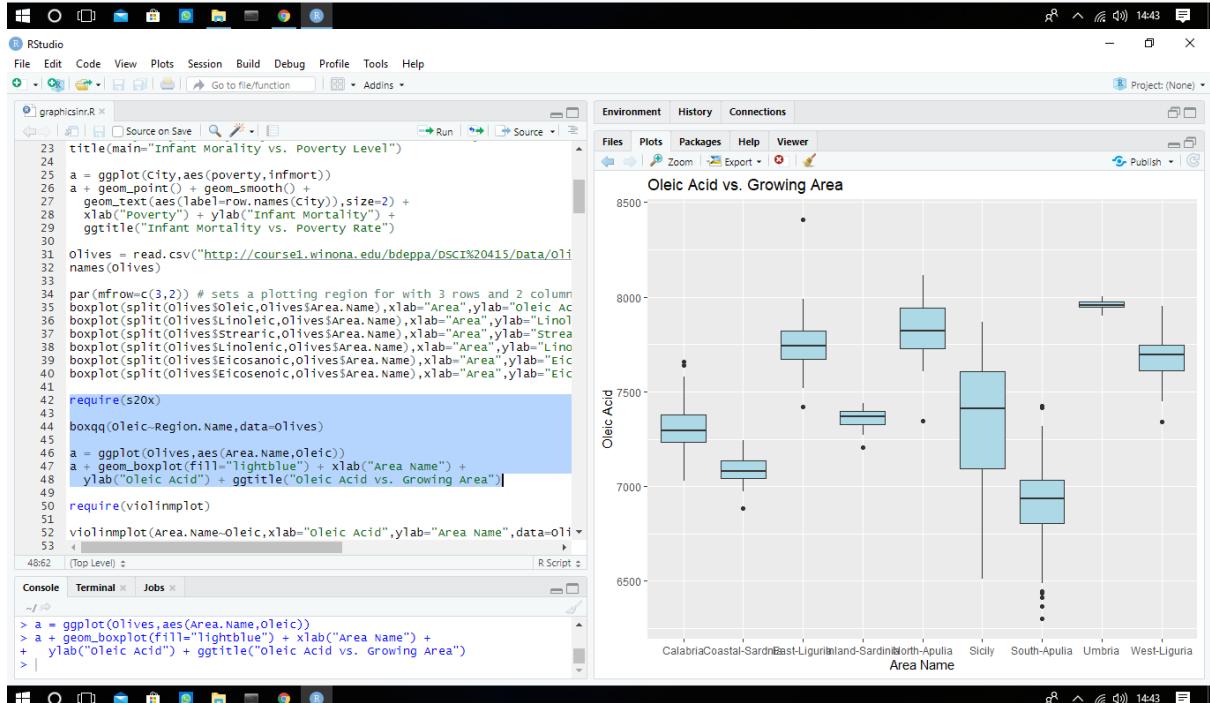
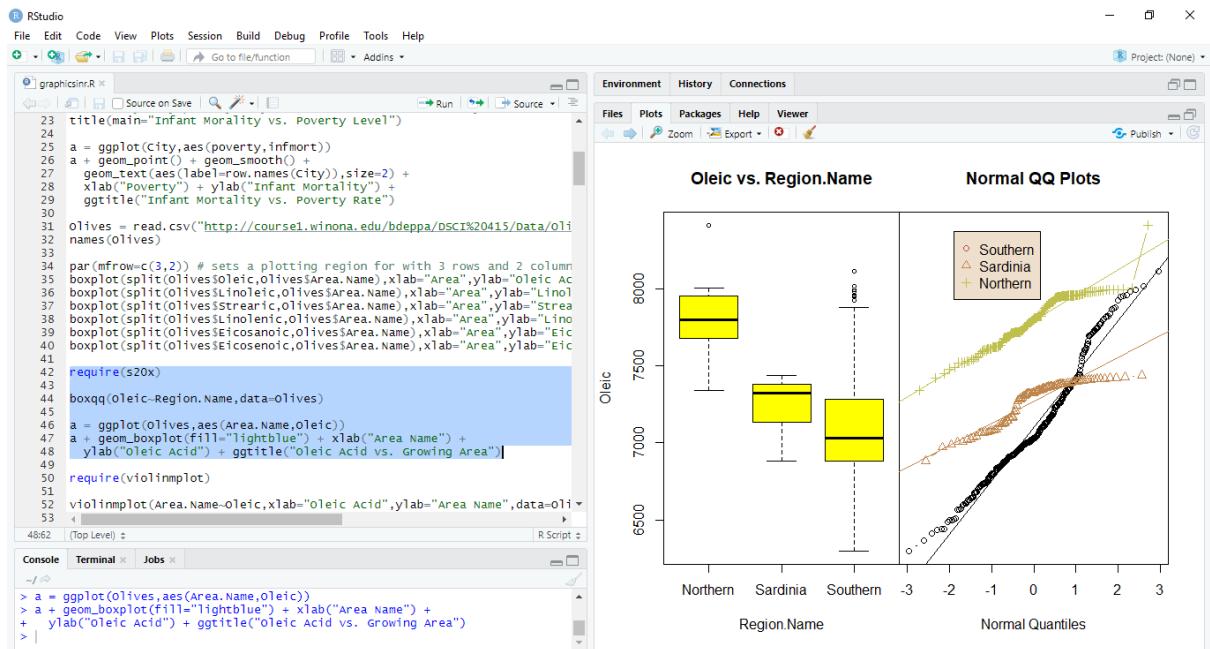
41:1 (Top Level) ↵ R Script ↵

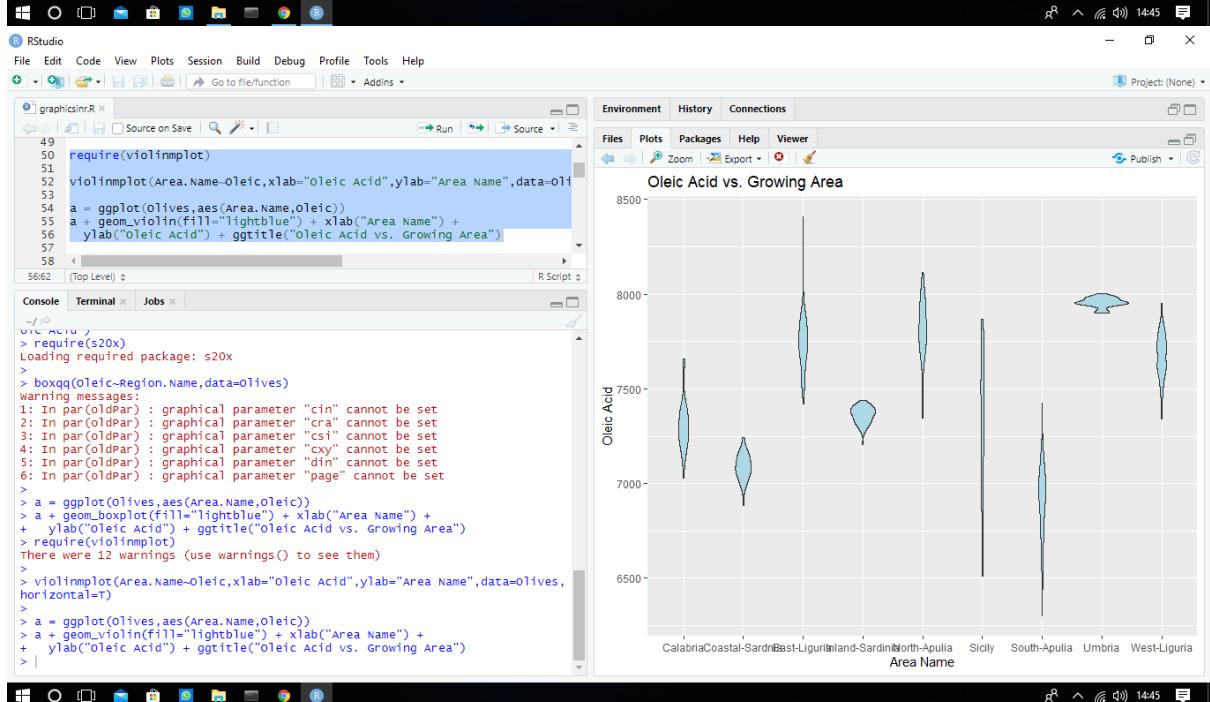
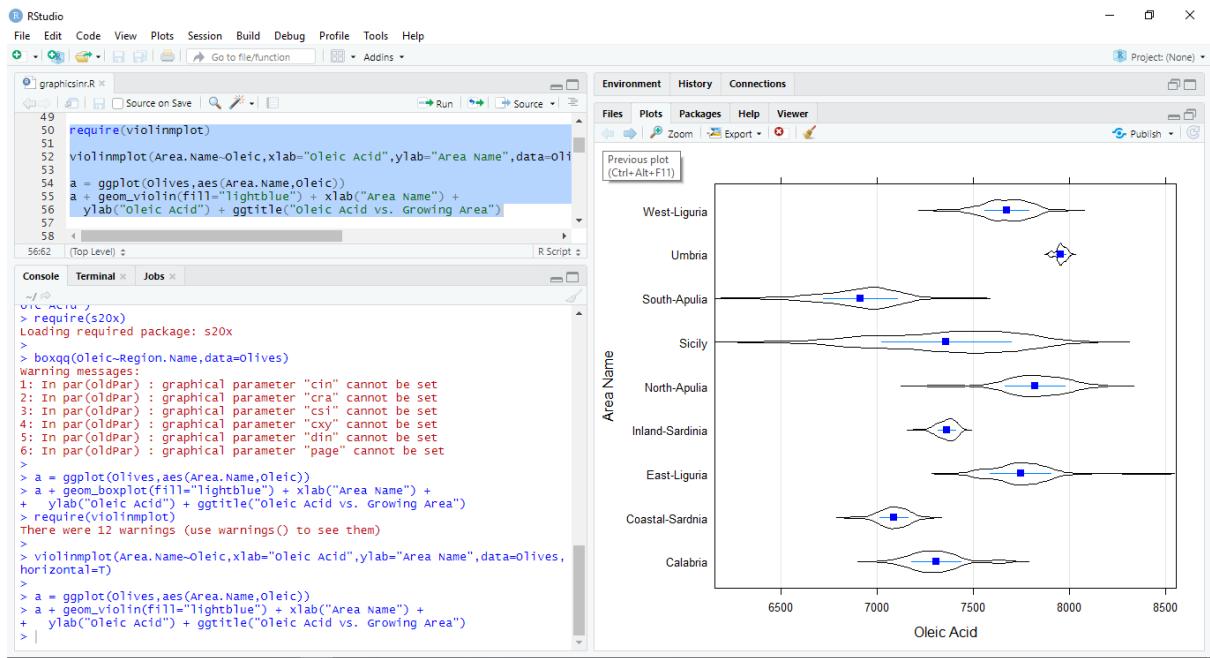
Console Terminal Jobs ↵

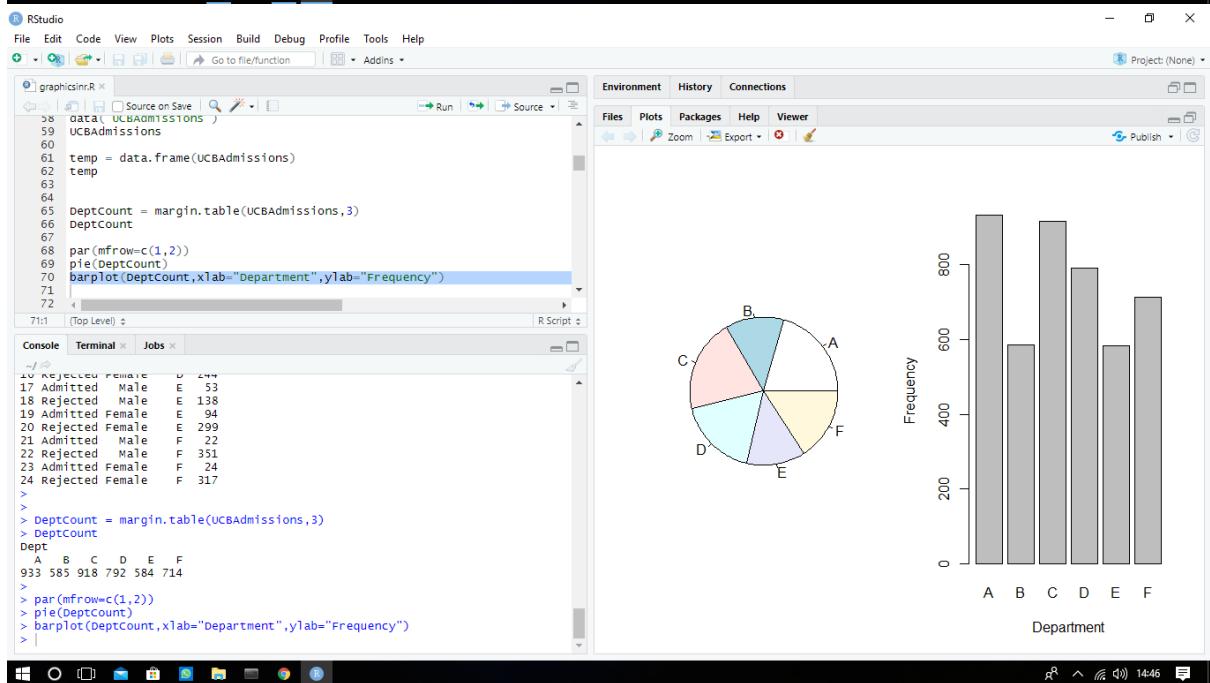
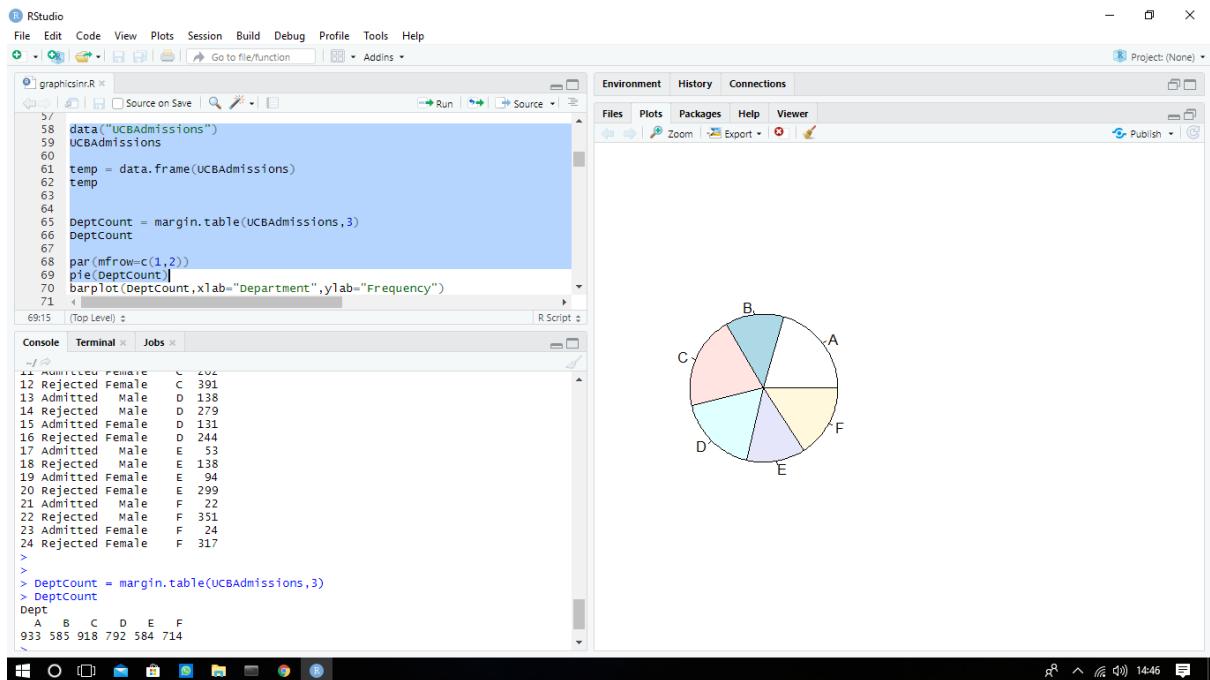
~/ -/

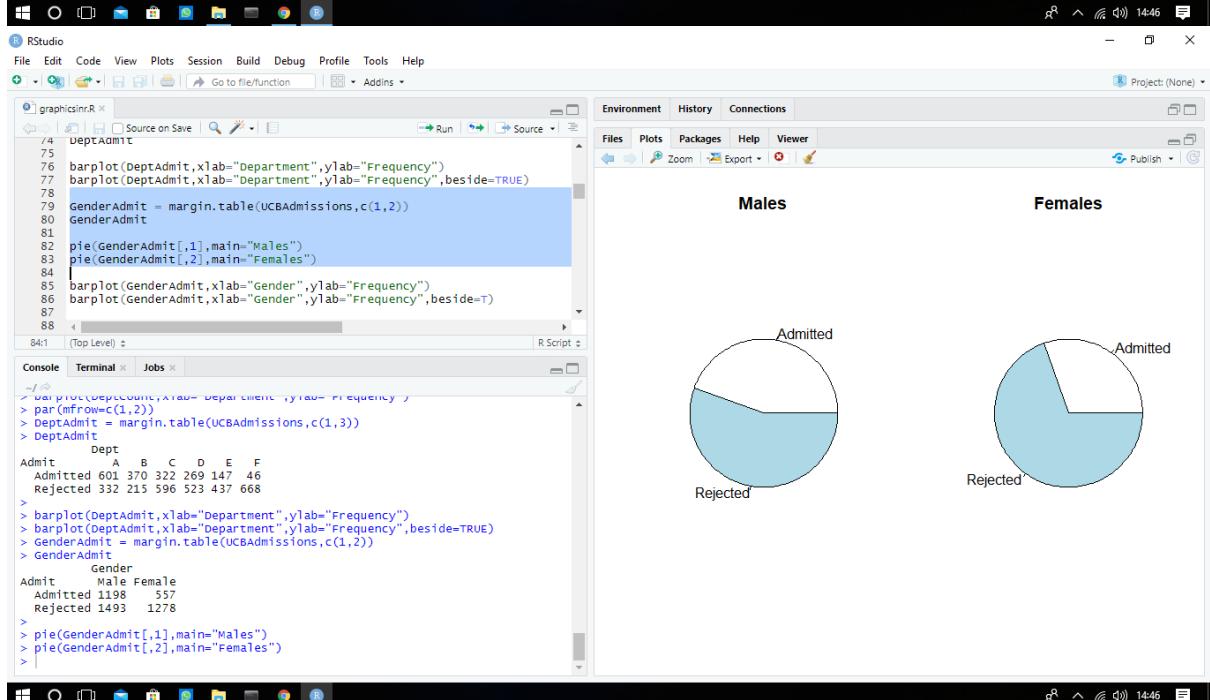
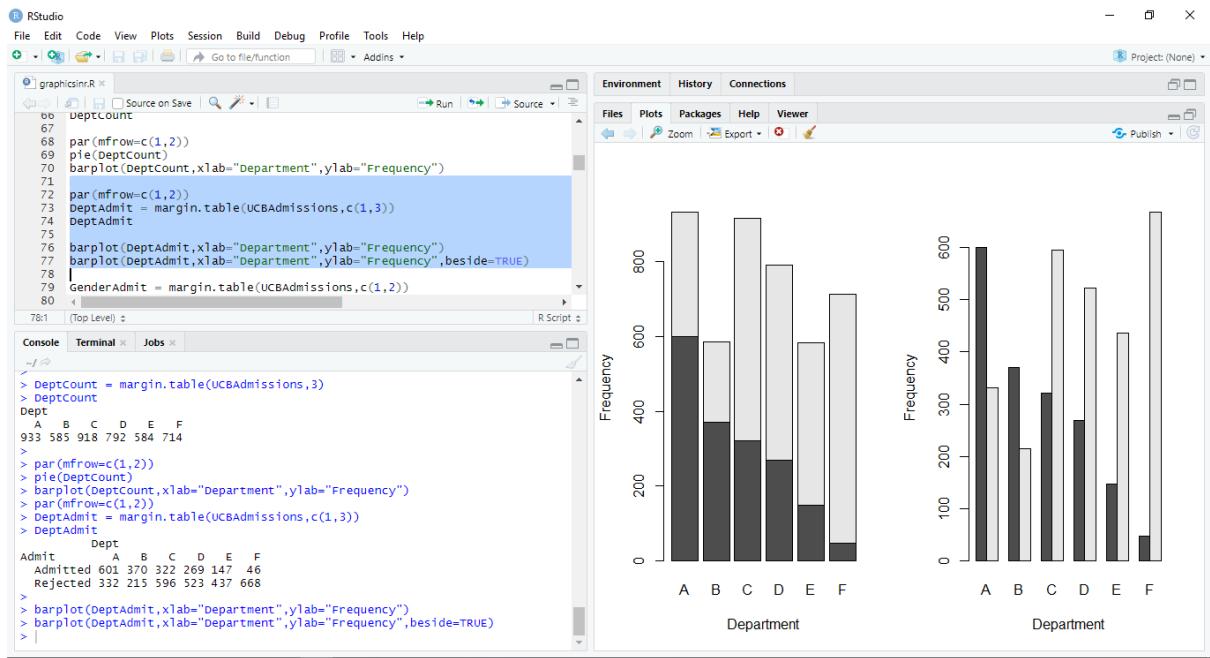
```
oleic Acid")
> boxplot(split(olives$Eicosenoic,olives$Area.Name),xlab="Area",ylab="Eicosenoic Acid")
> |
```

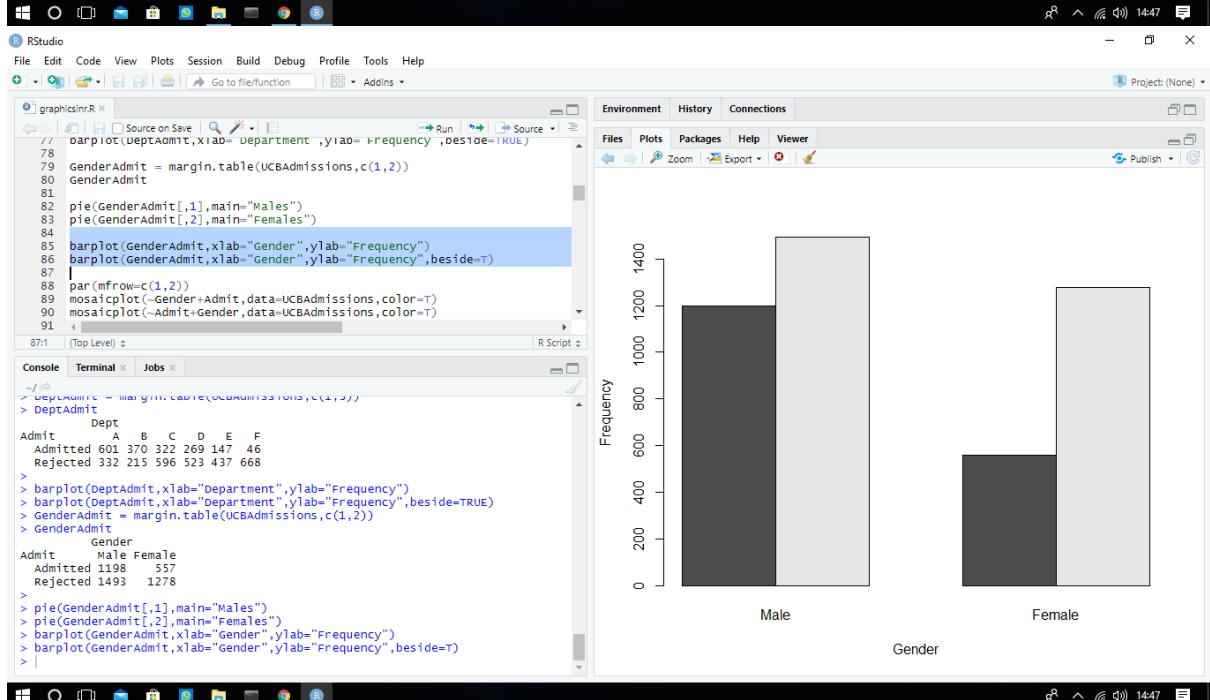
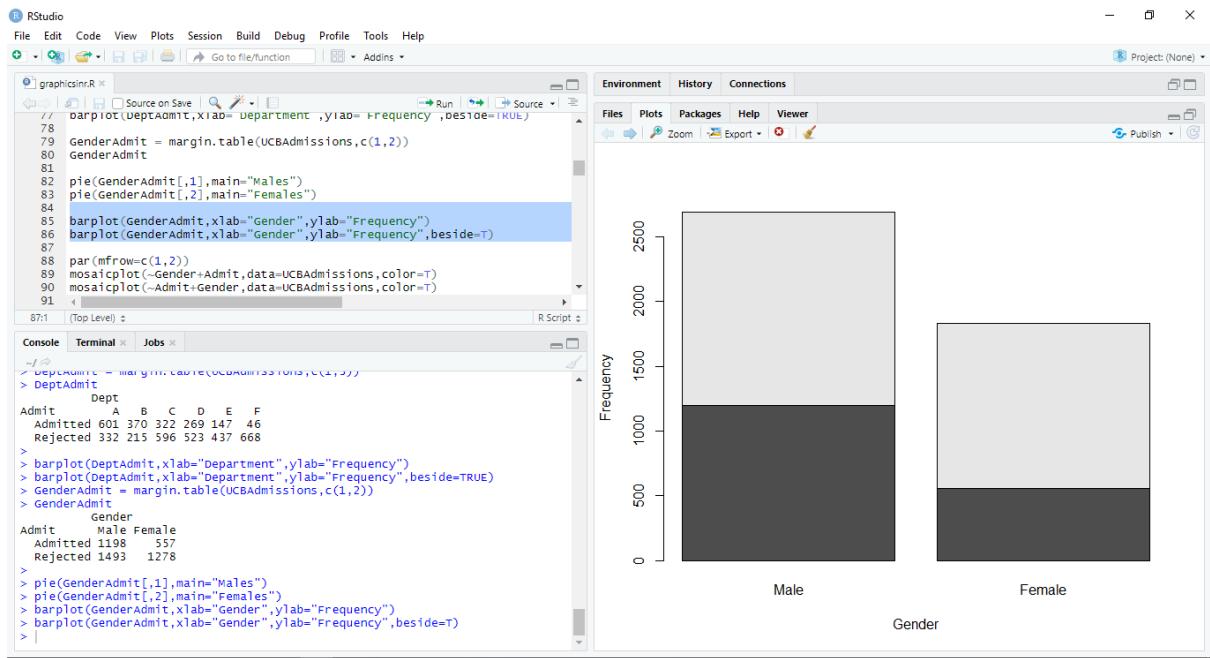


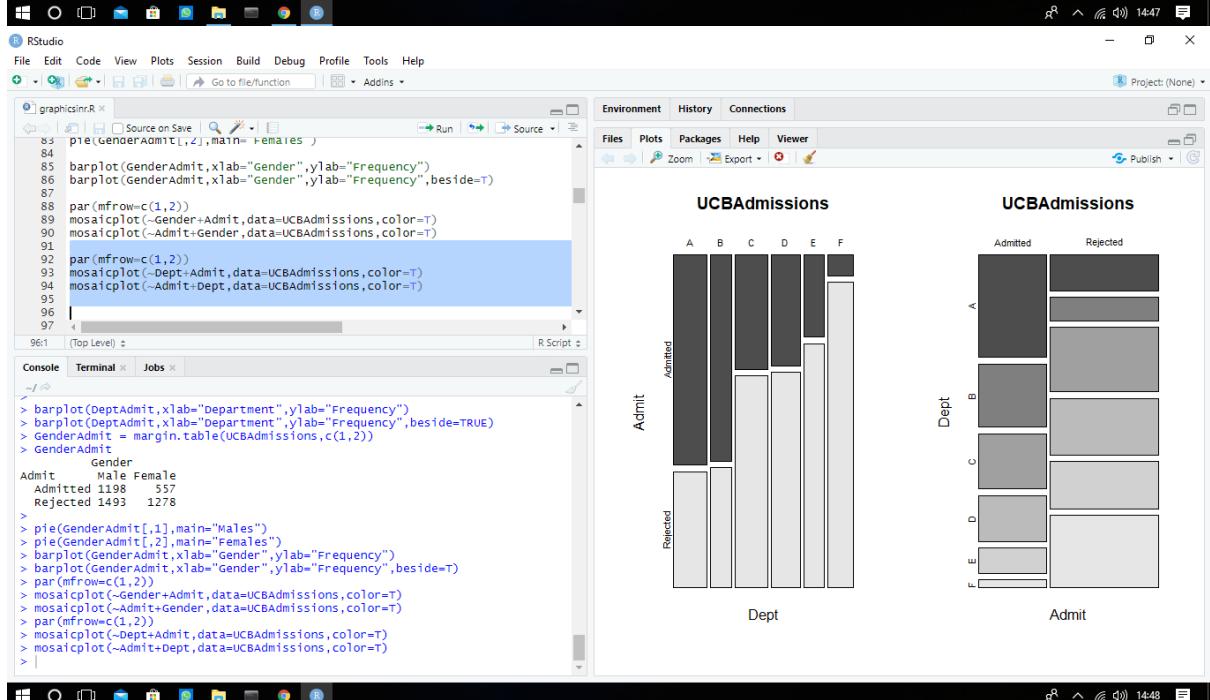
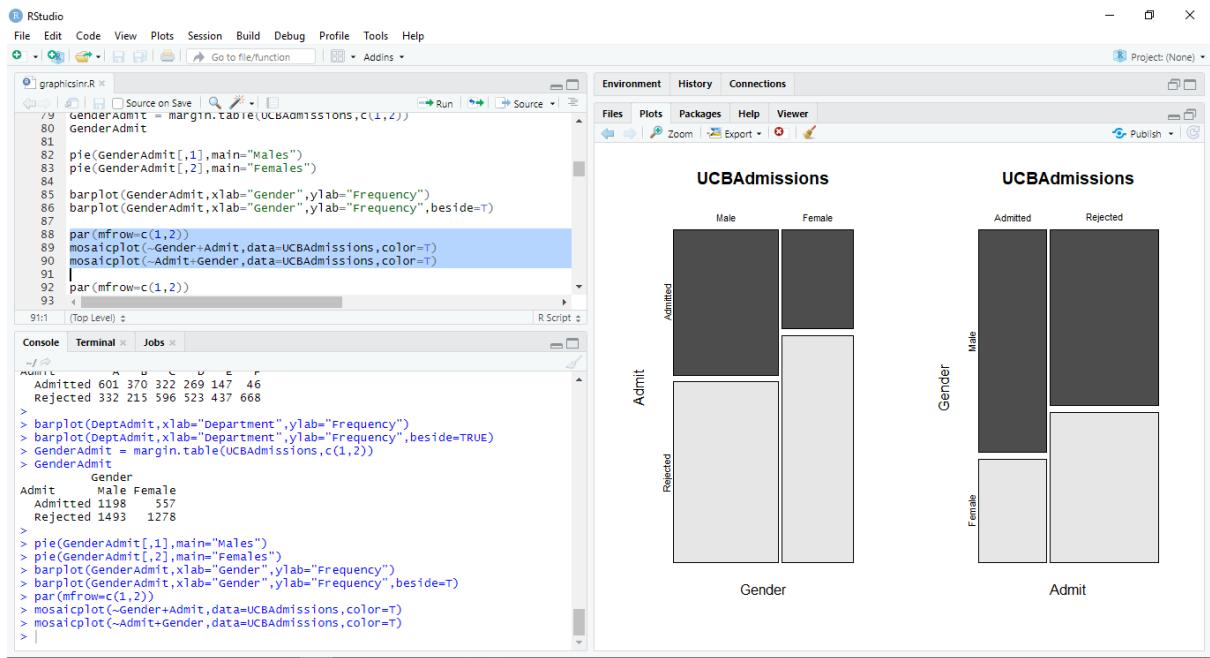


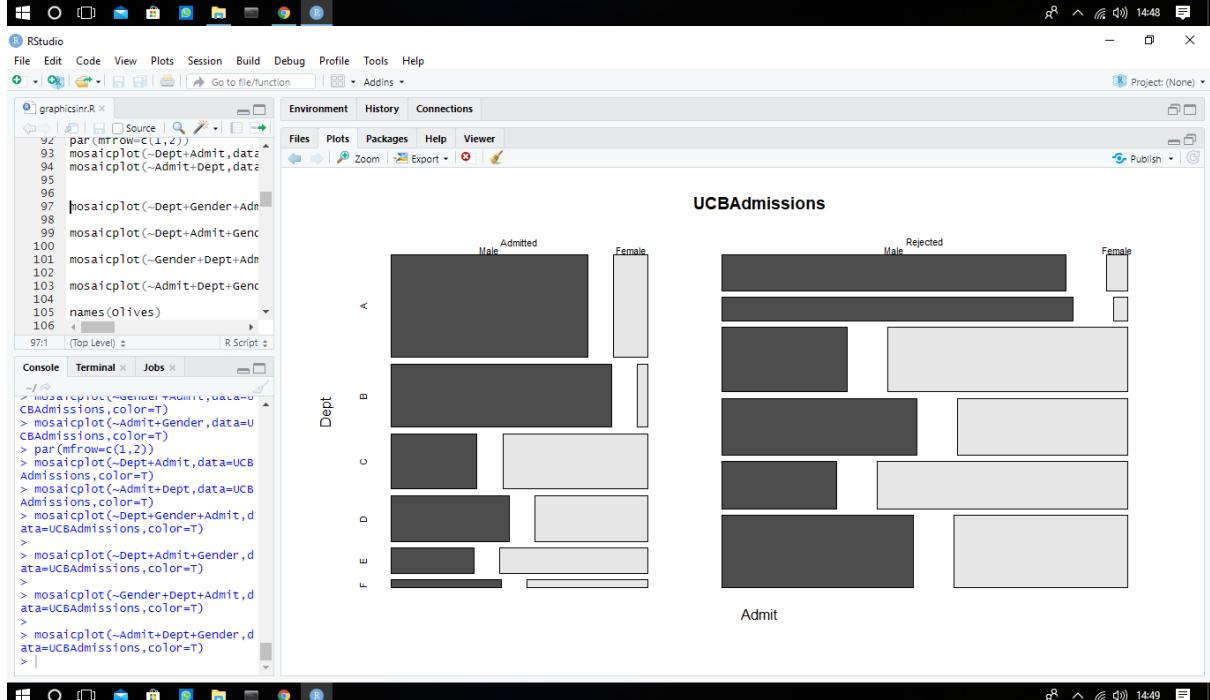
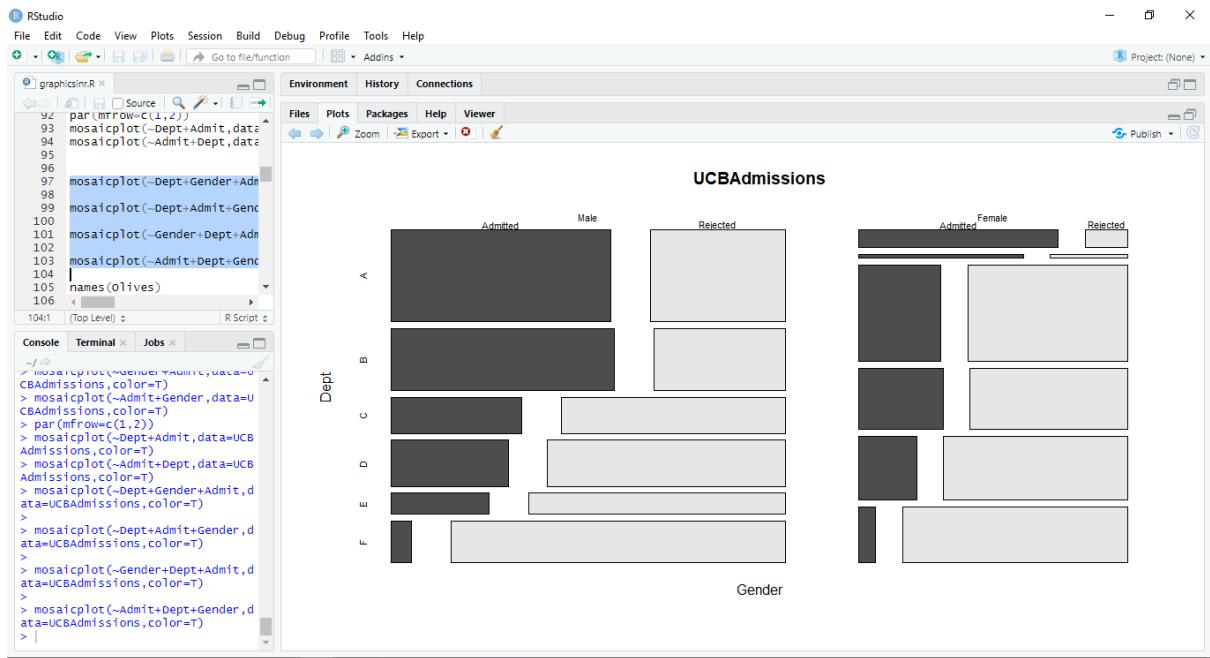


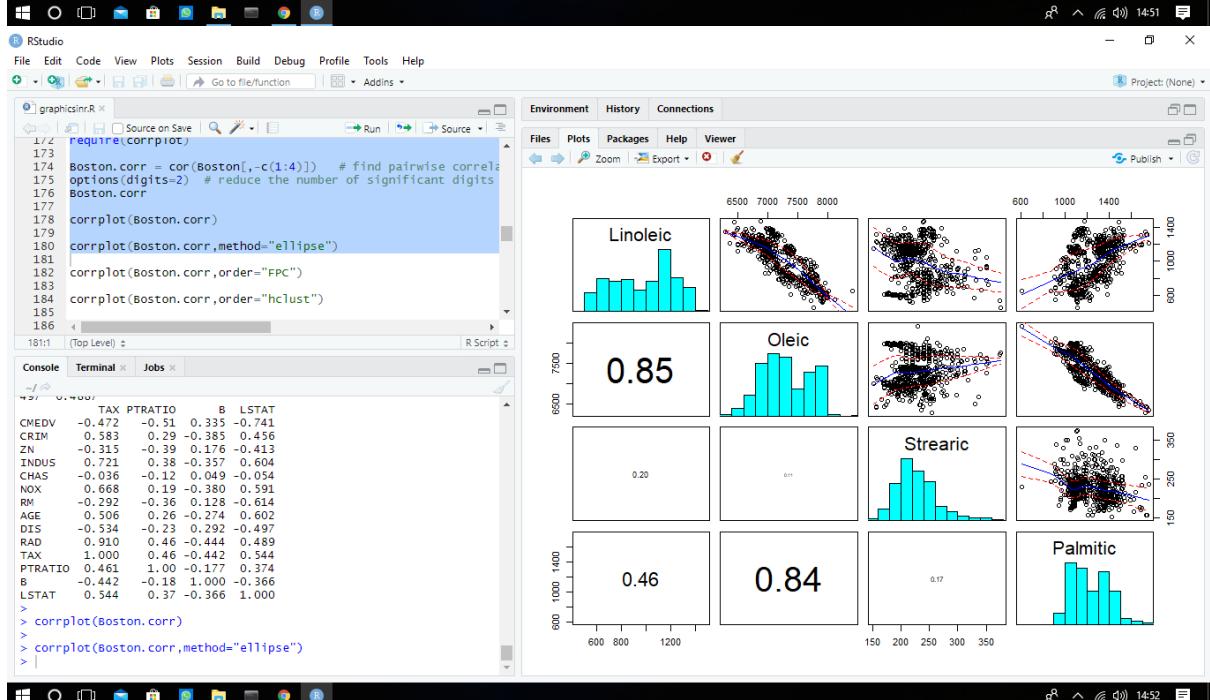
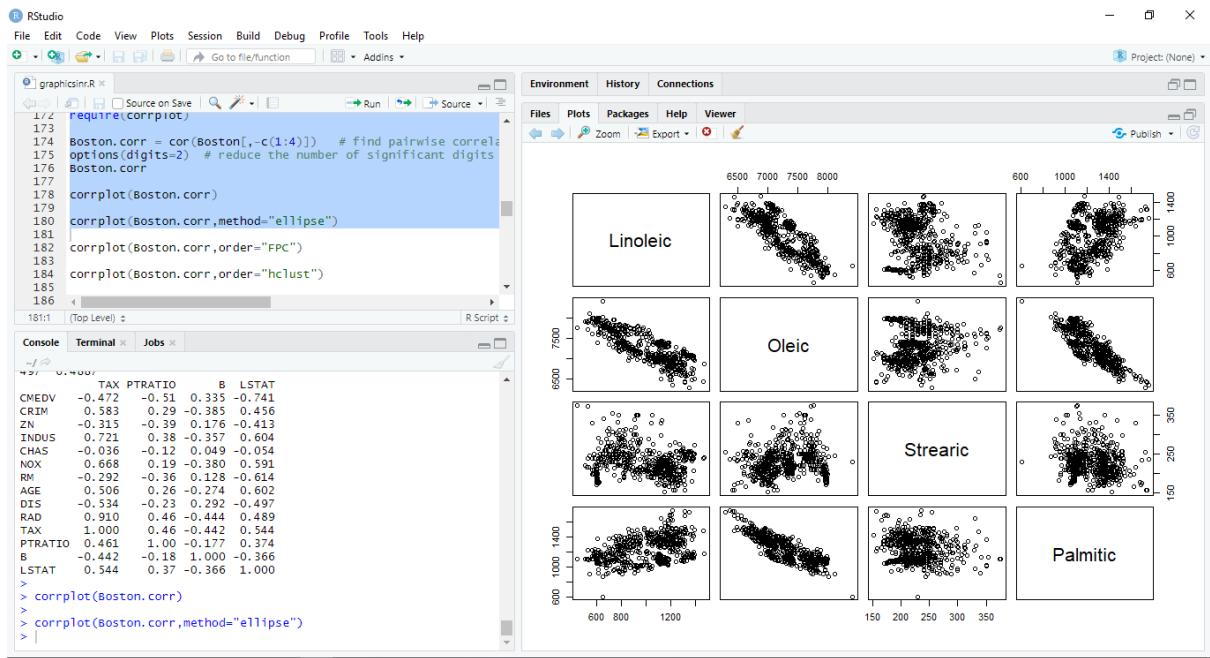


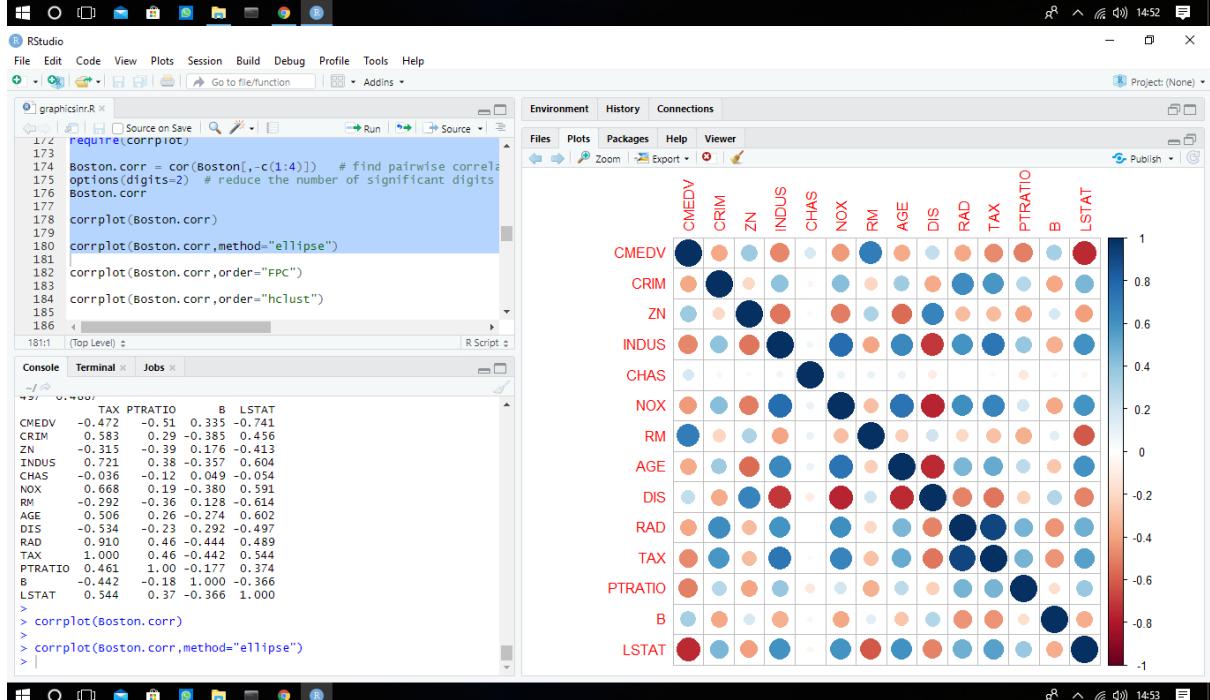
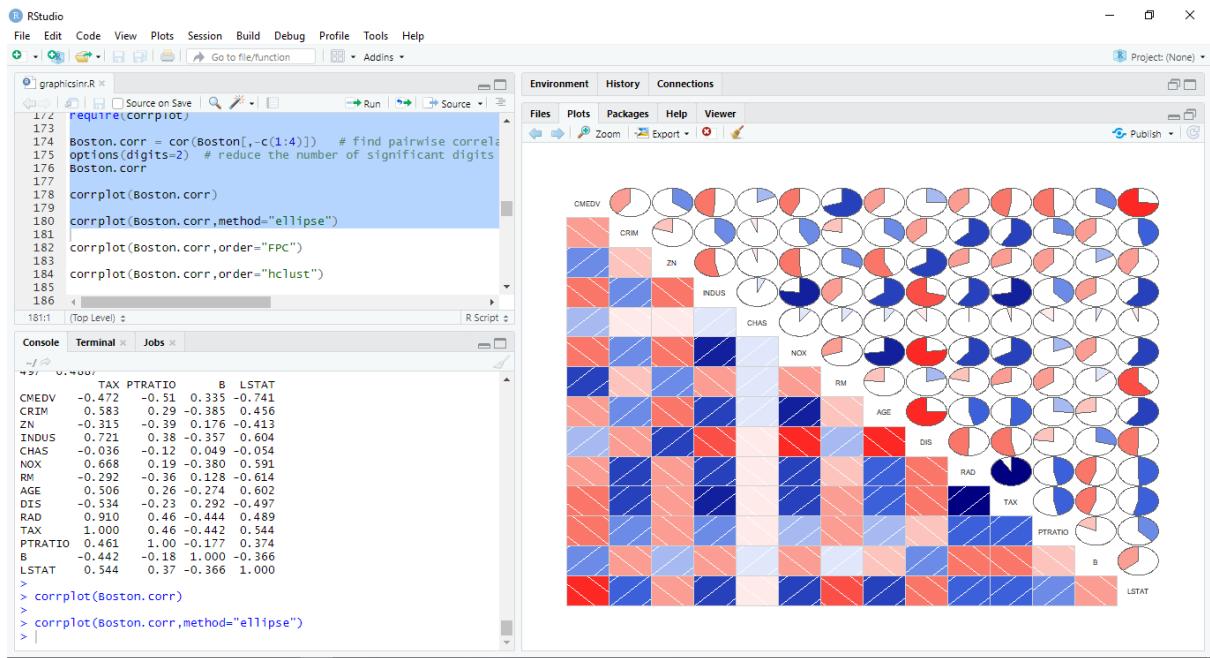


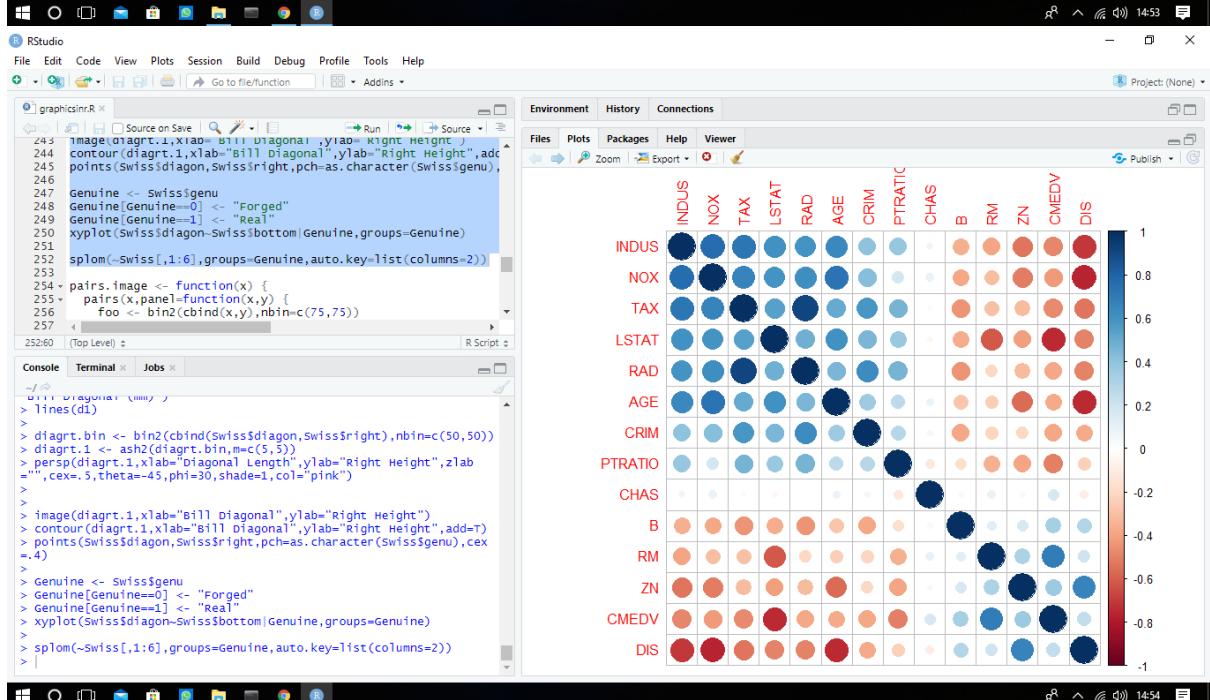
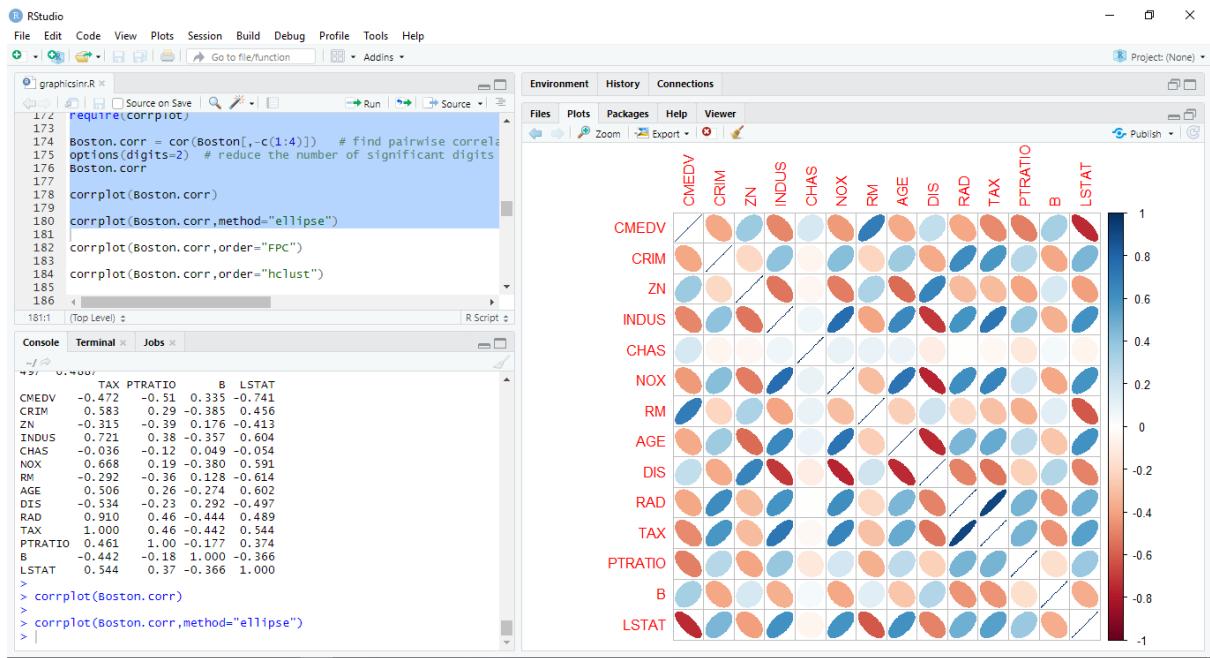


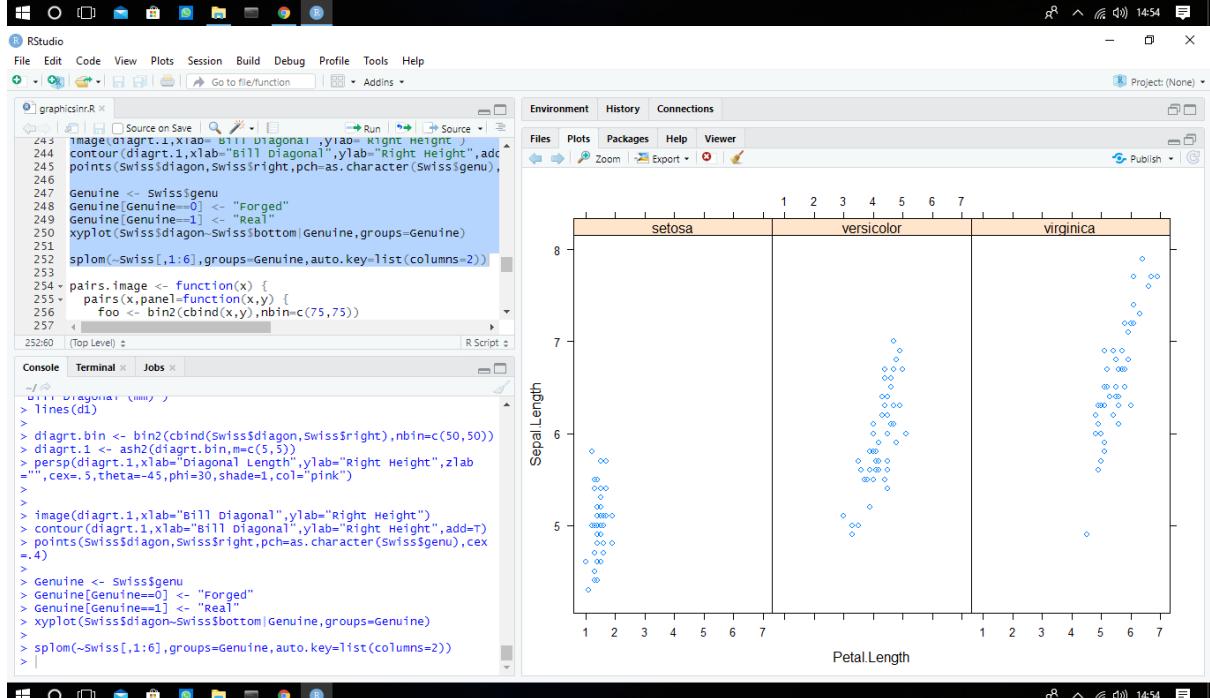
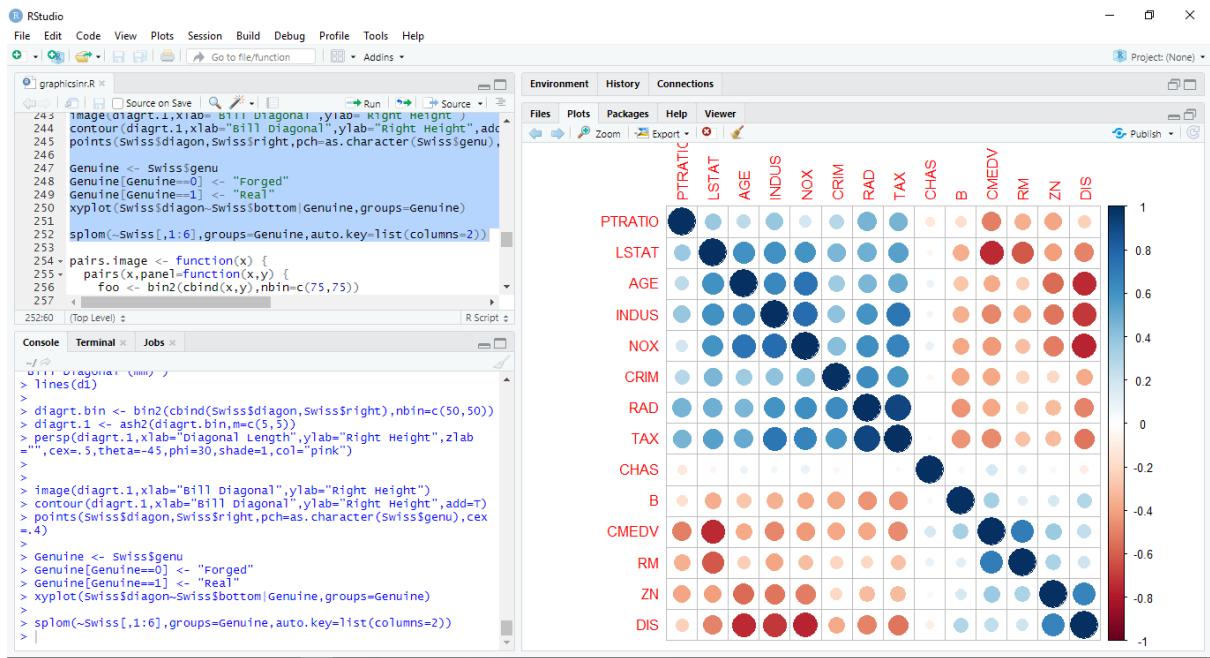


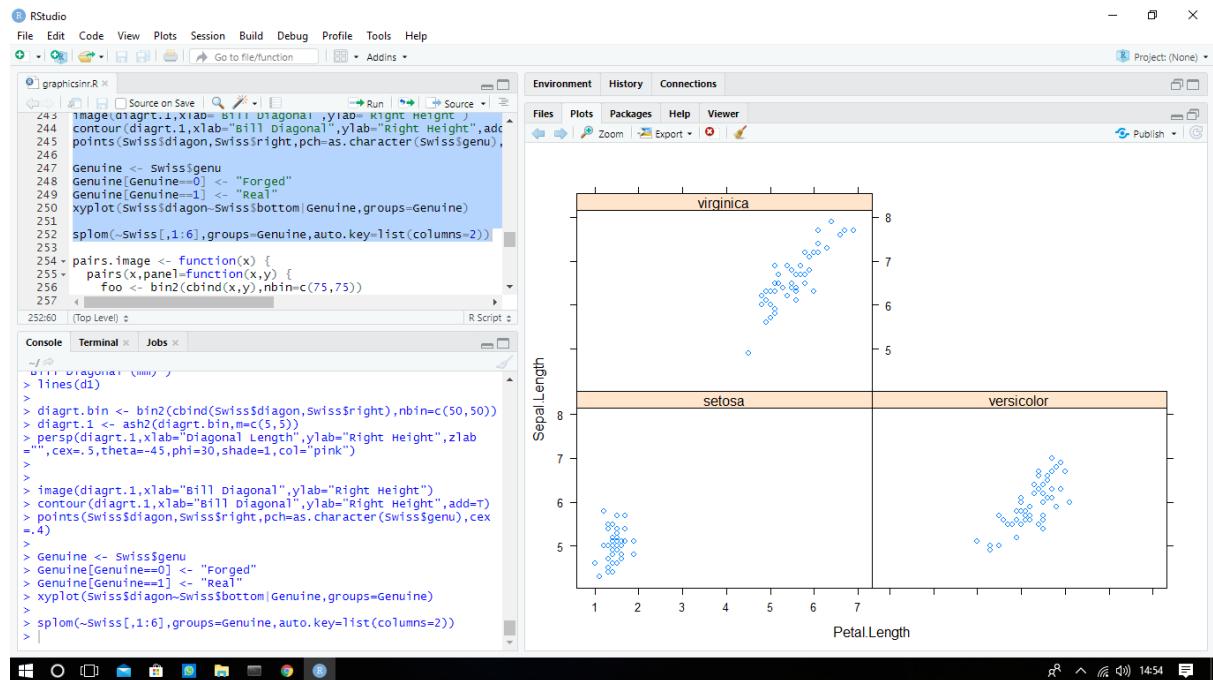


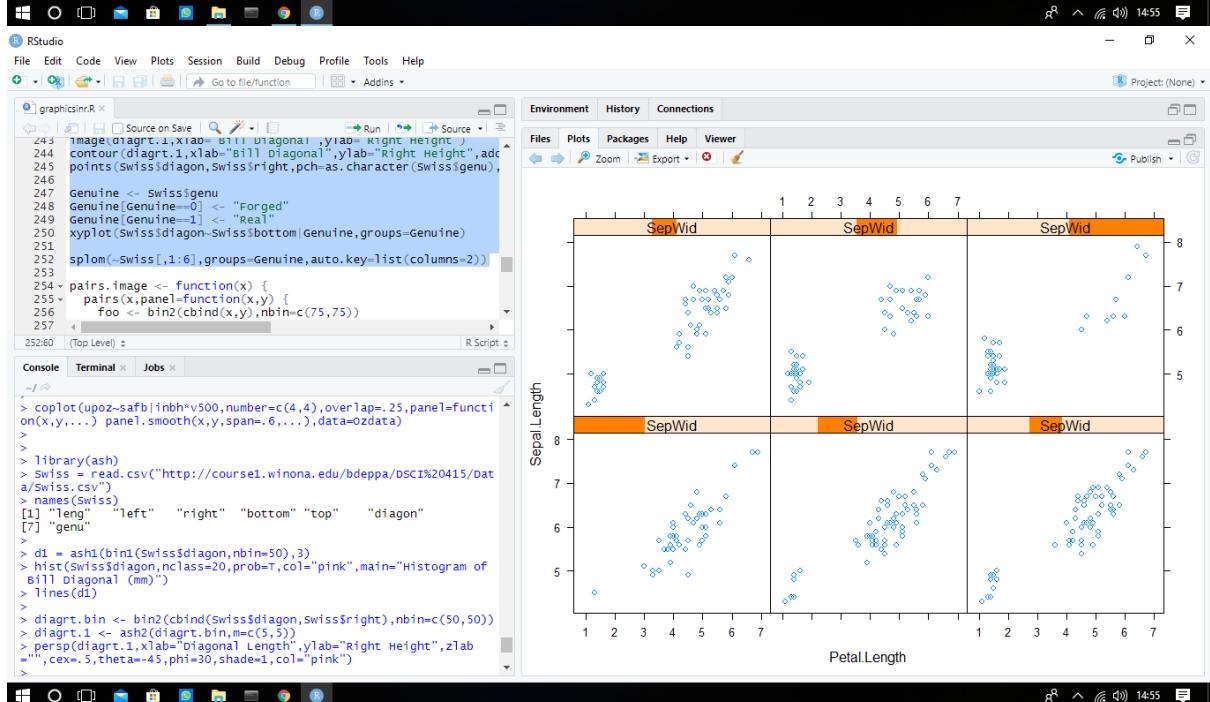
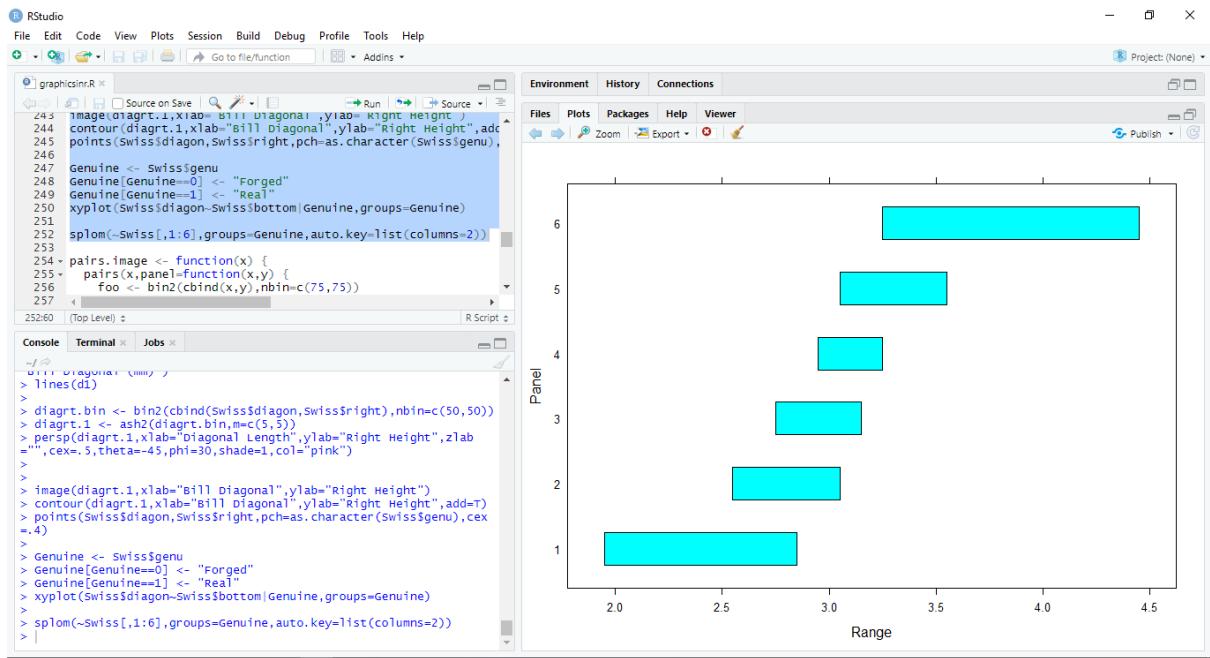


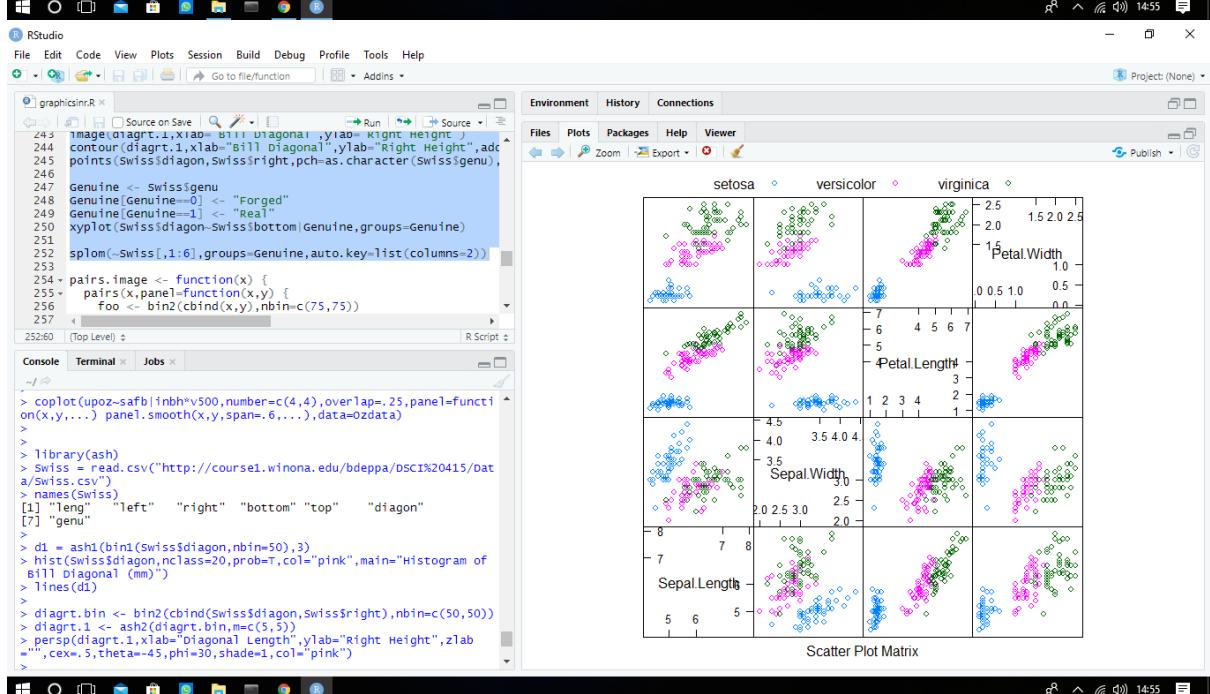
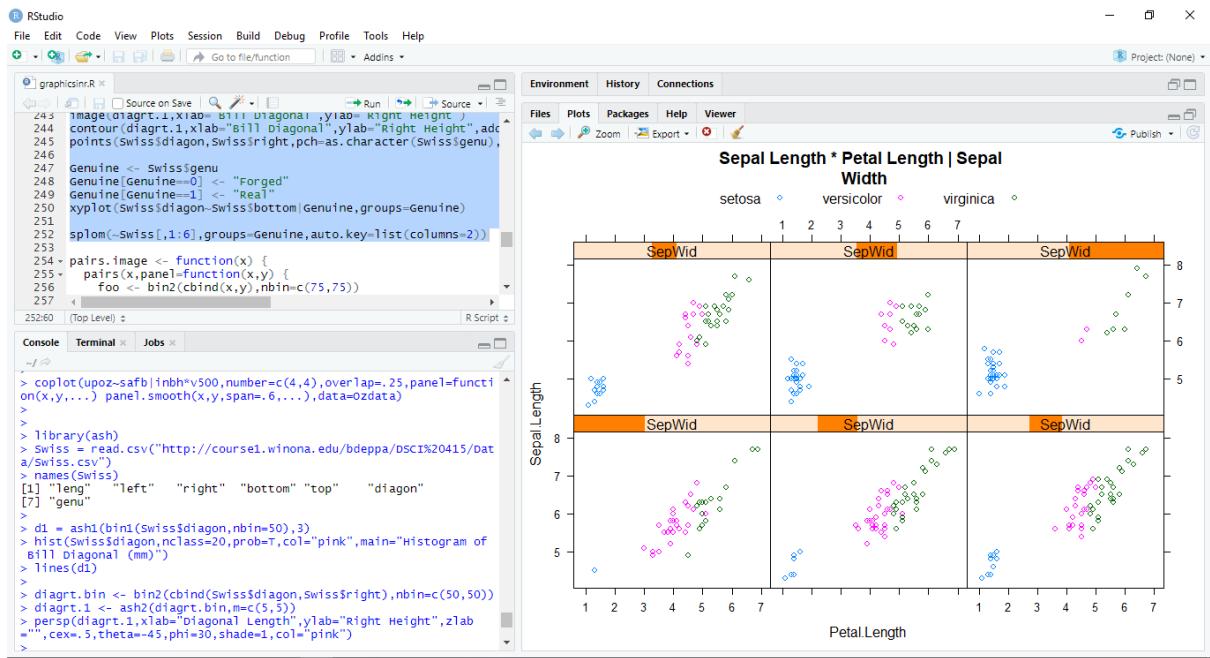


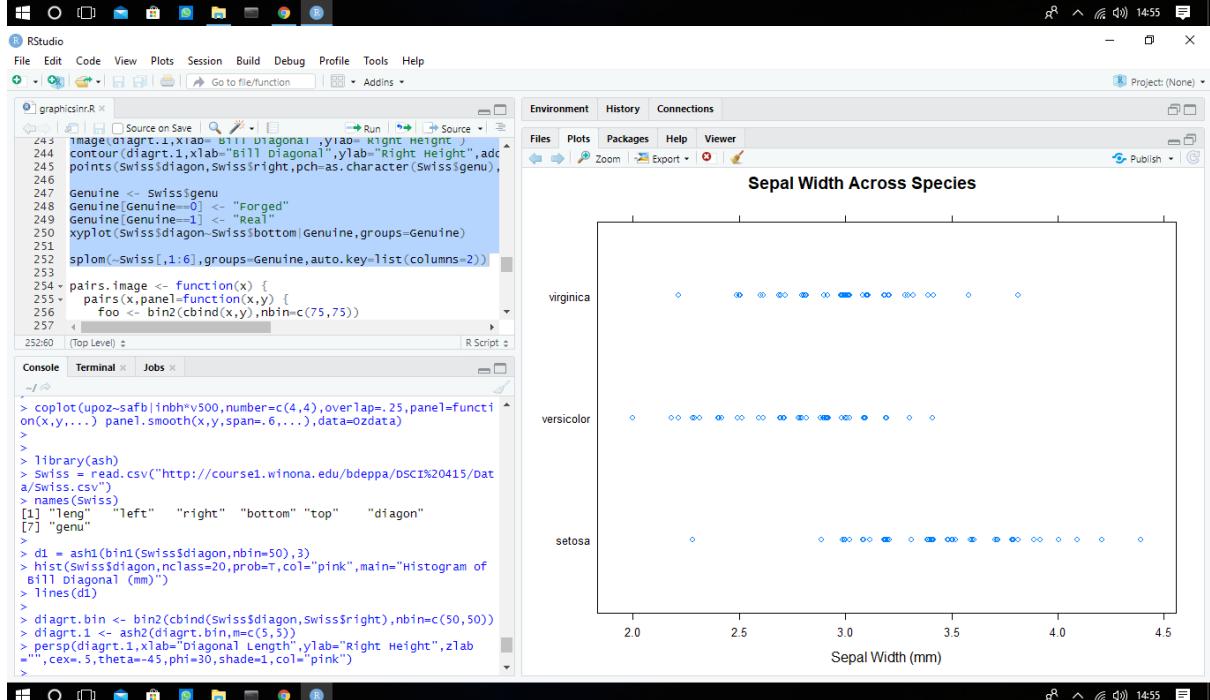
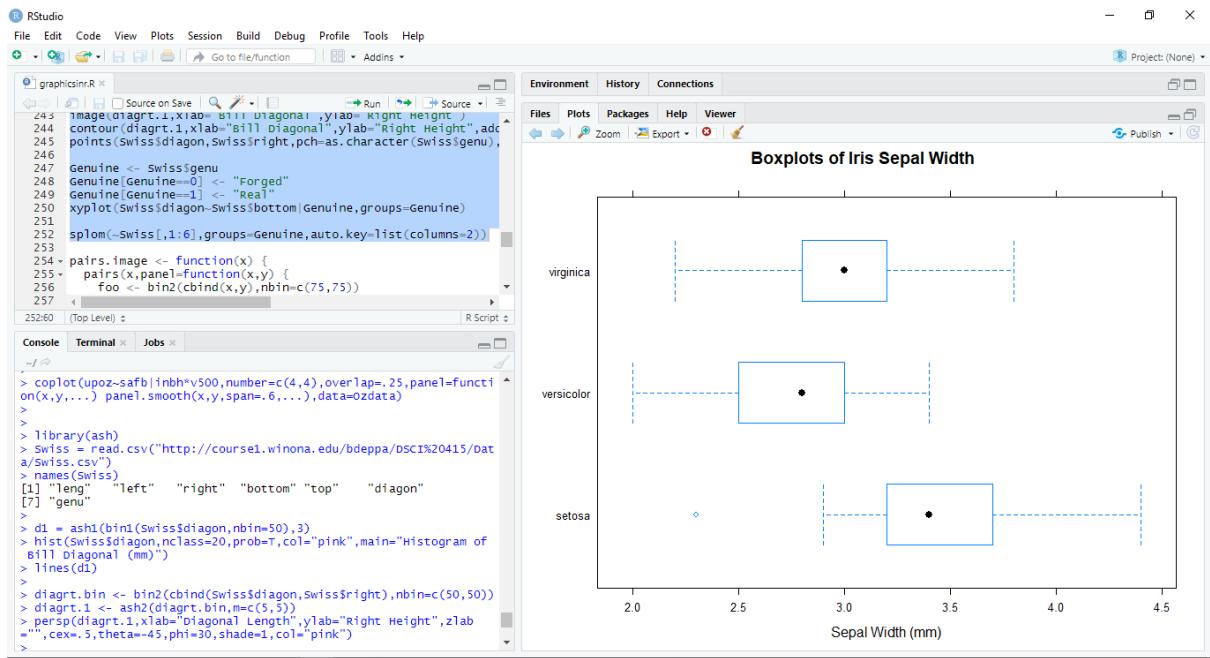


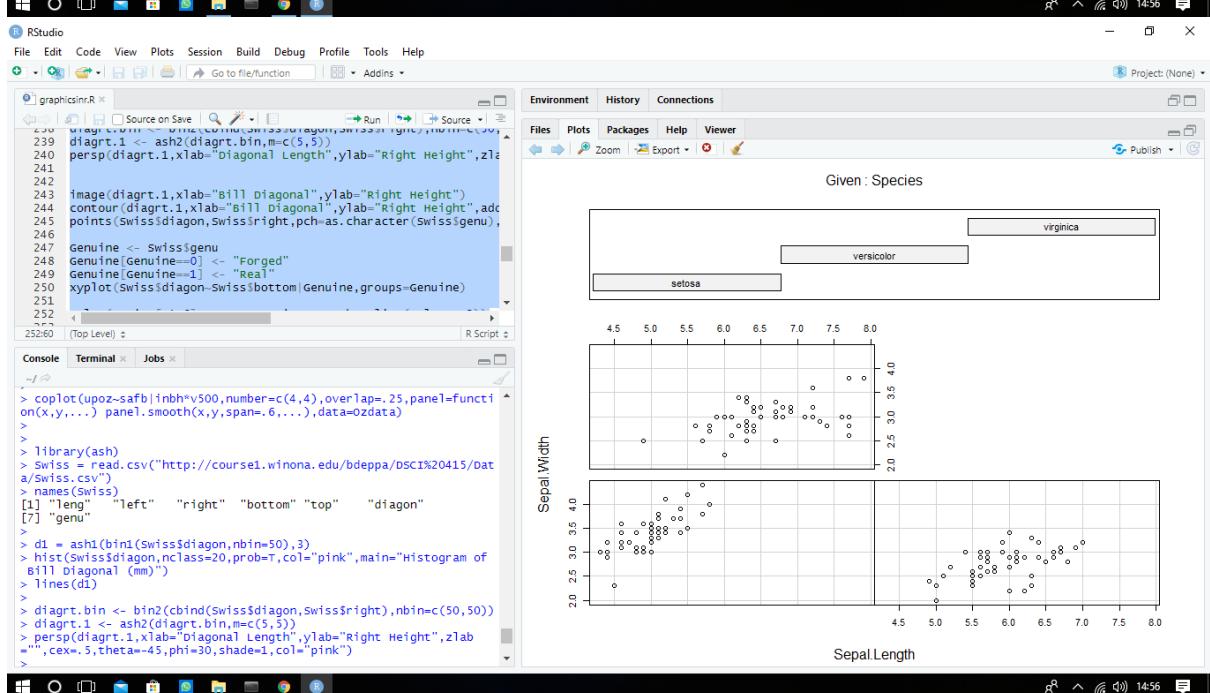
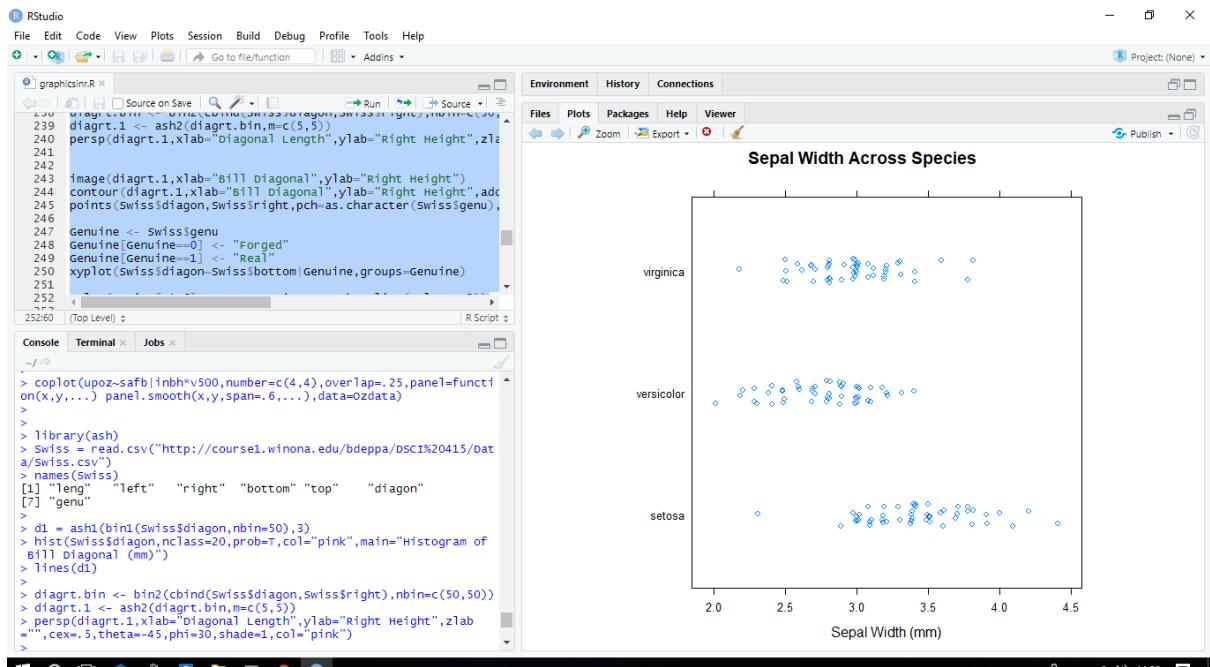


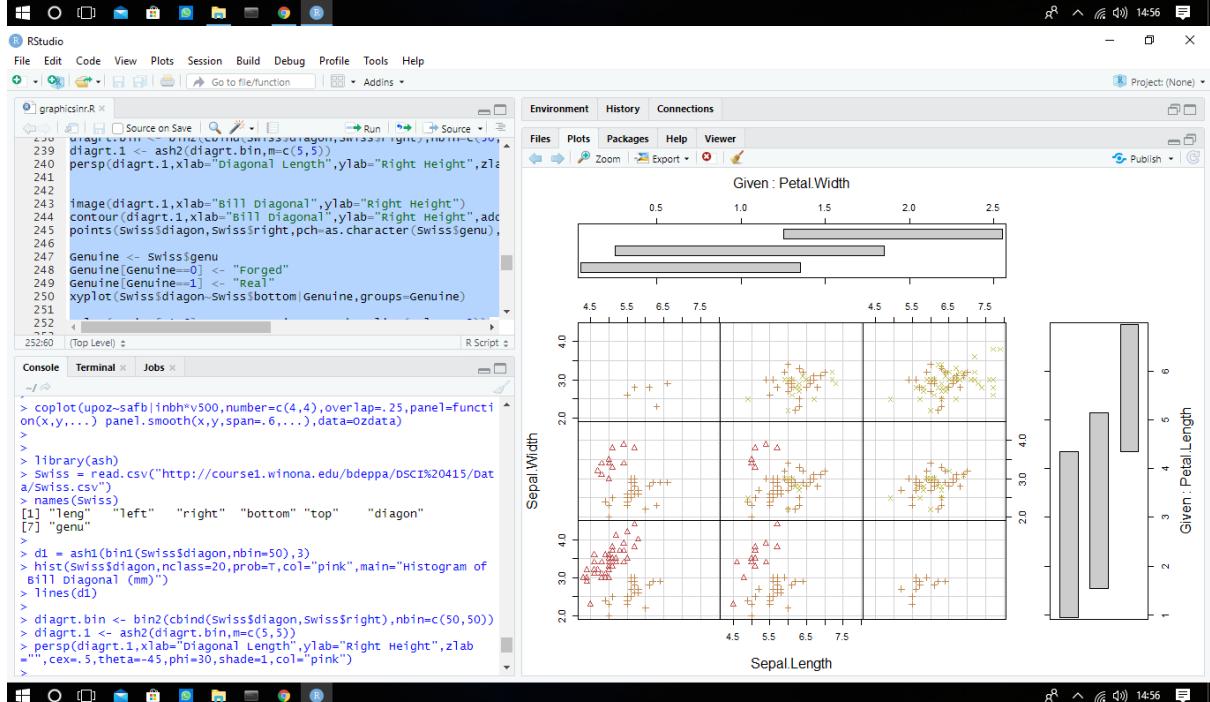
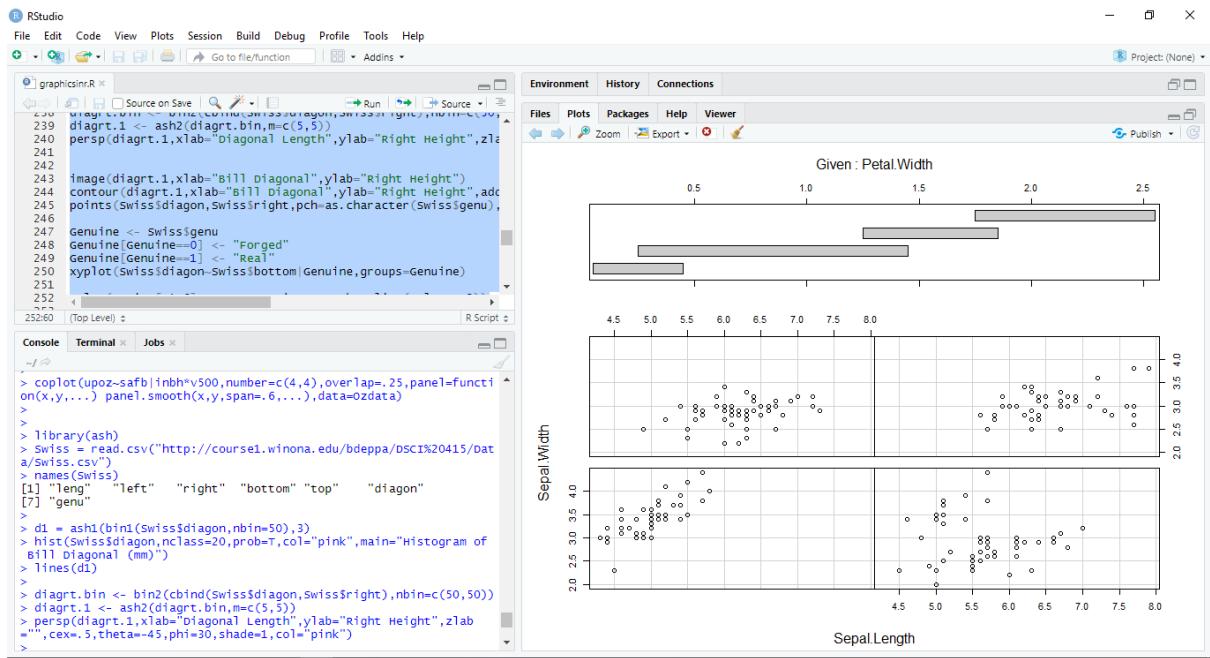


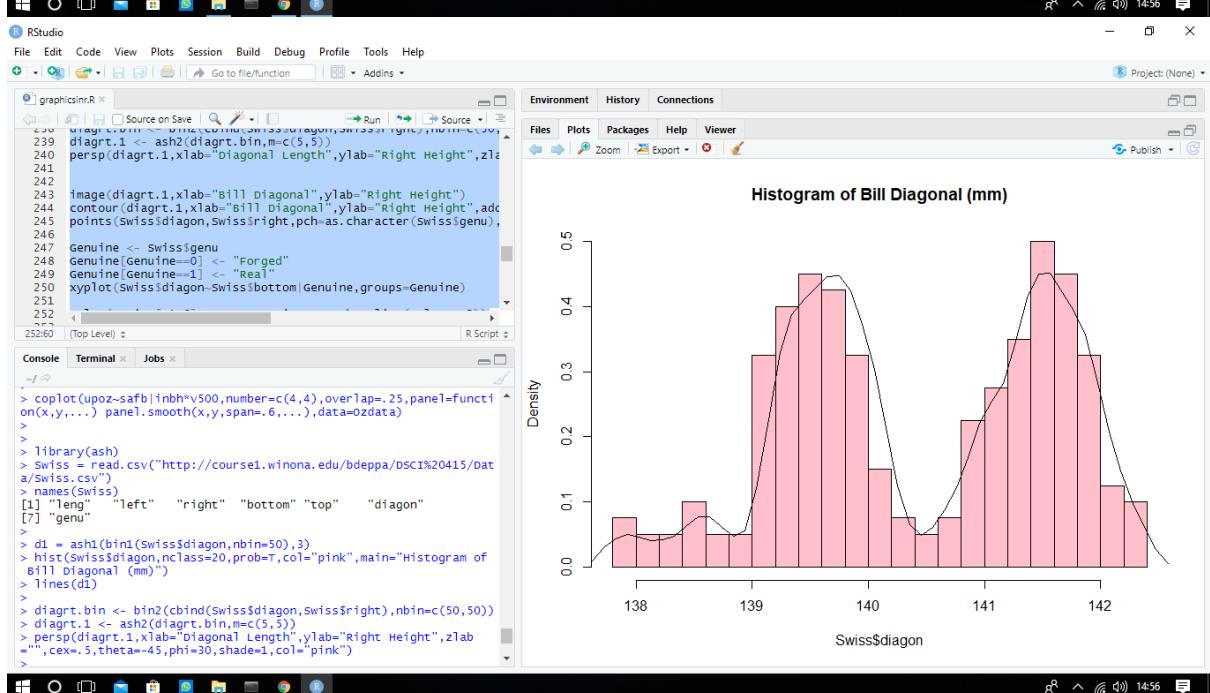
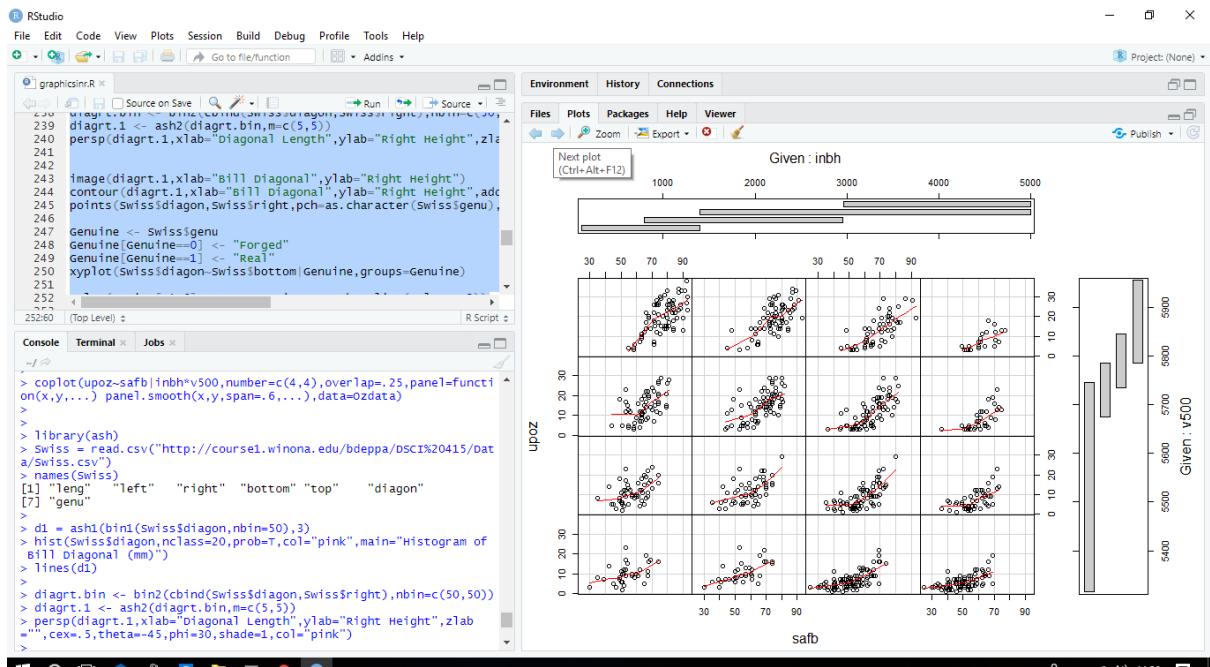


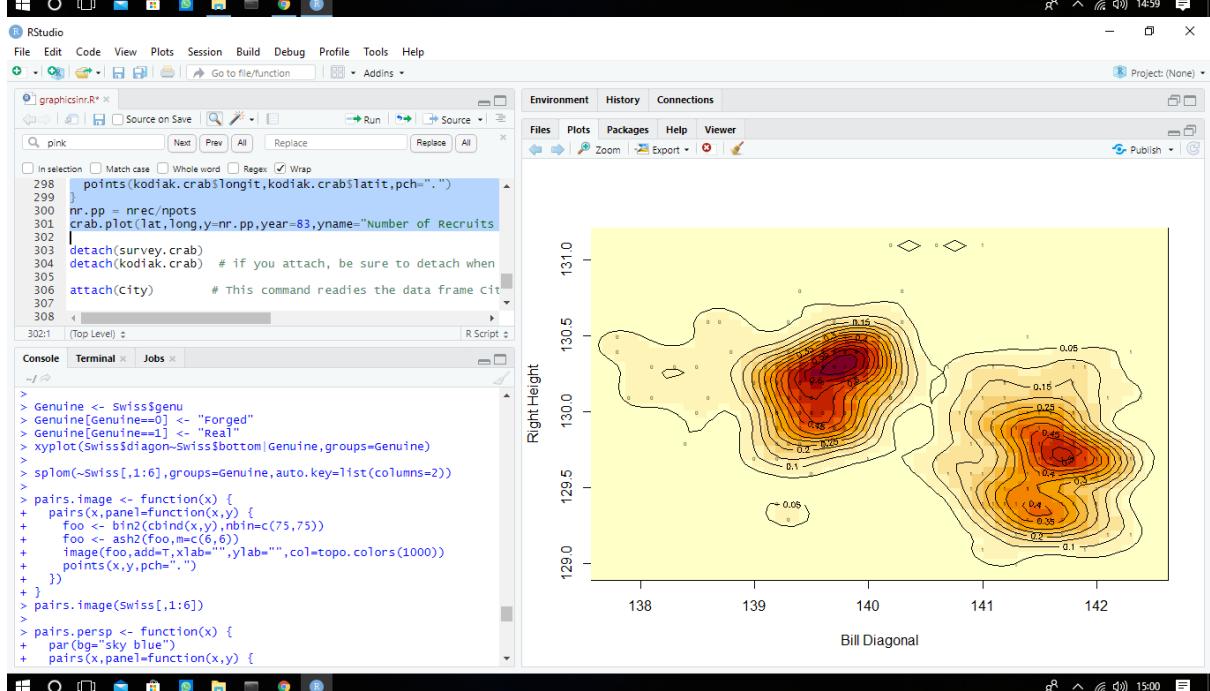
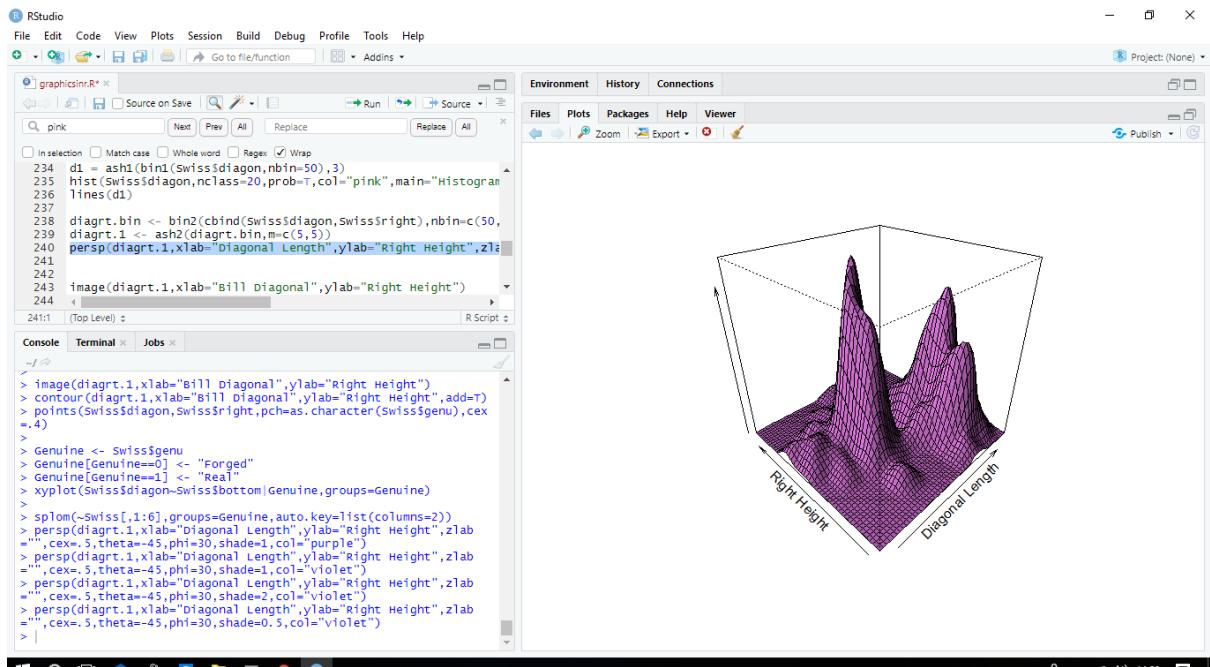


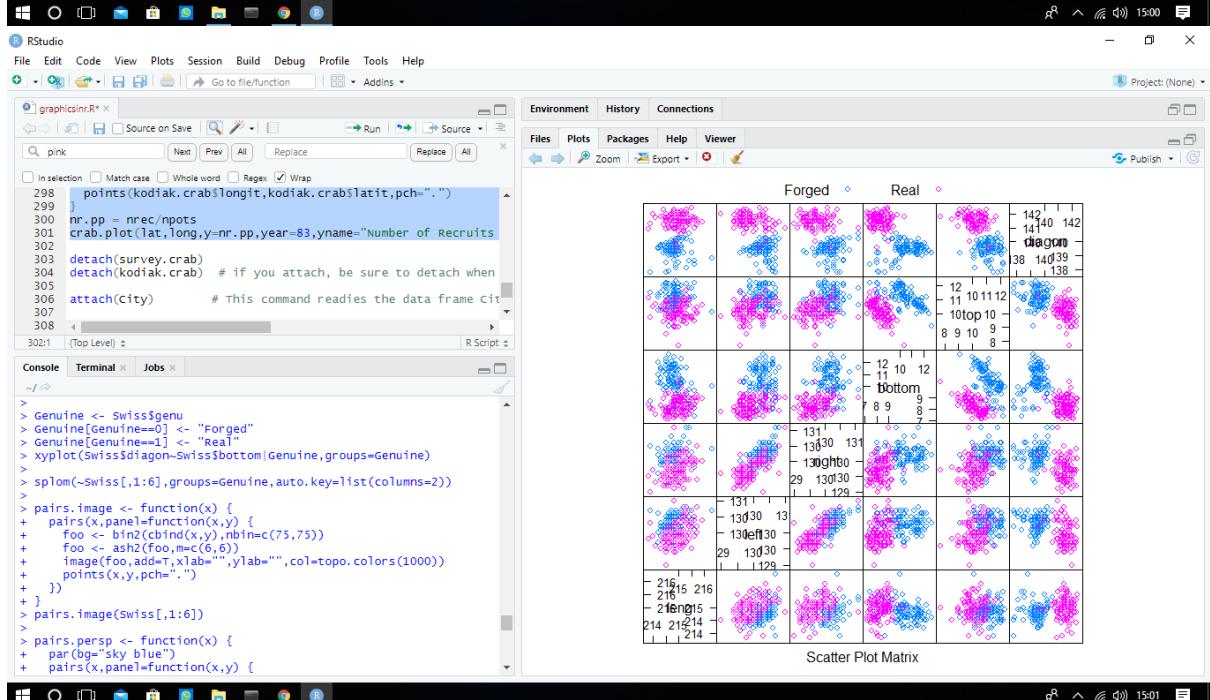
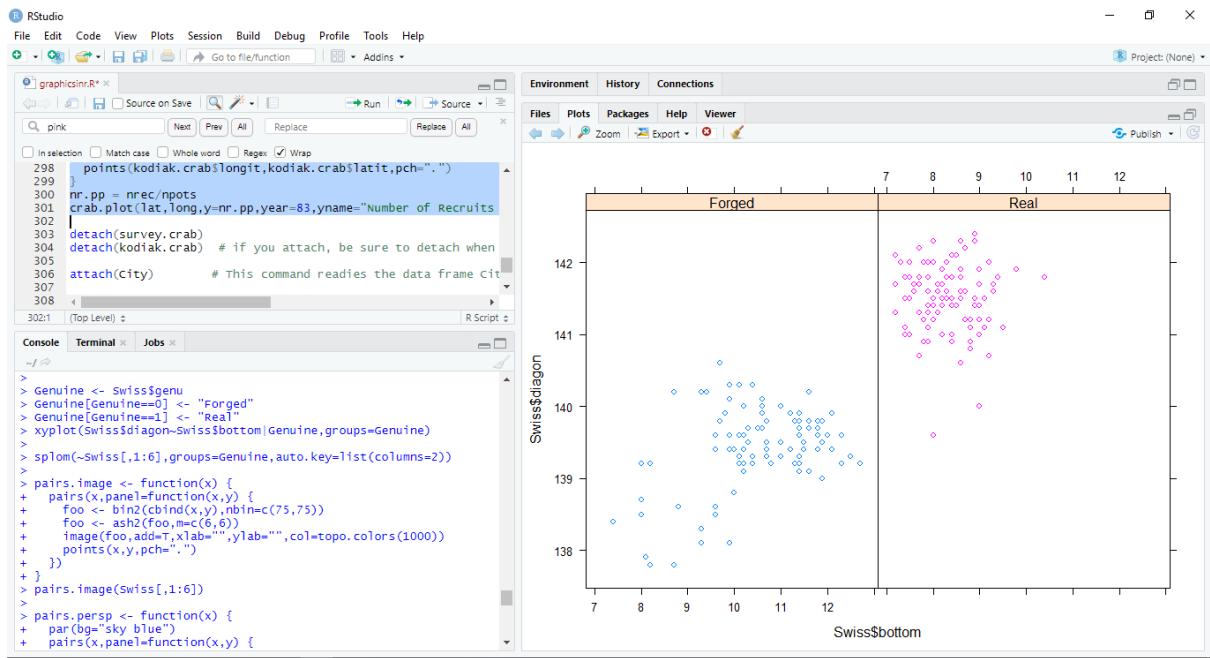


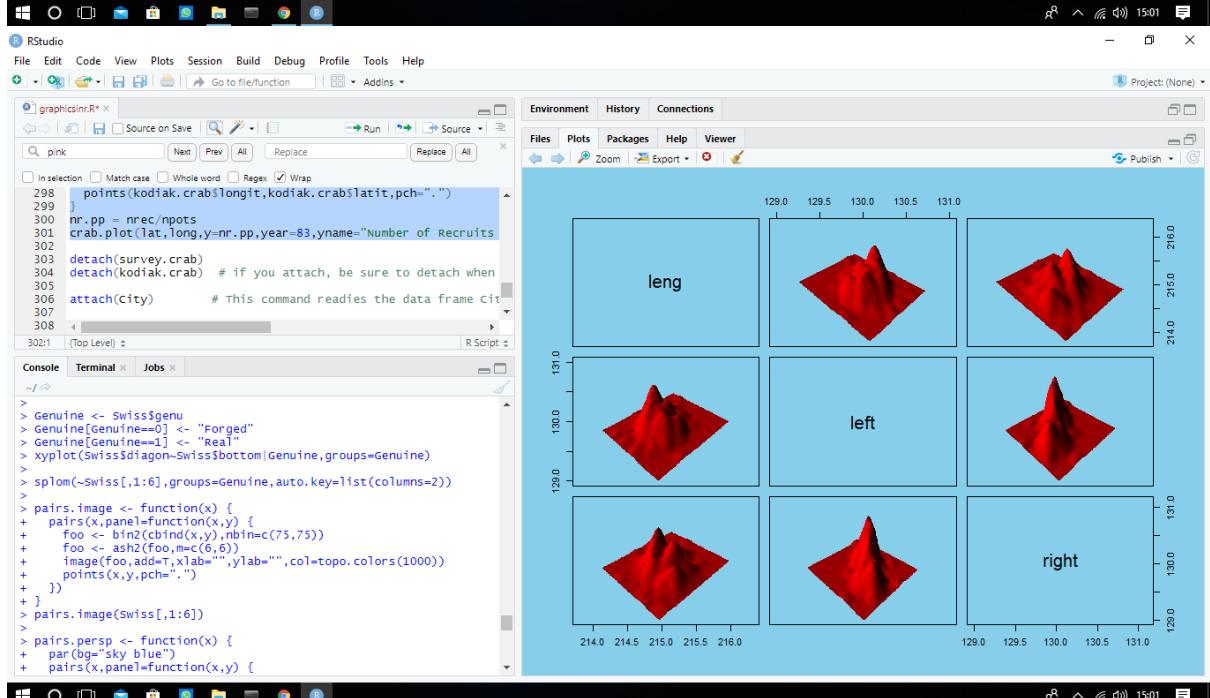
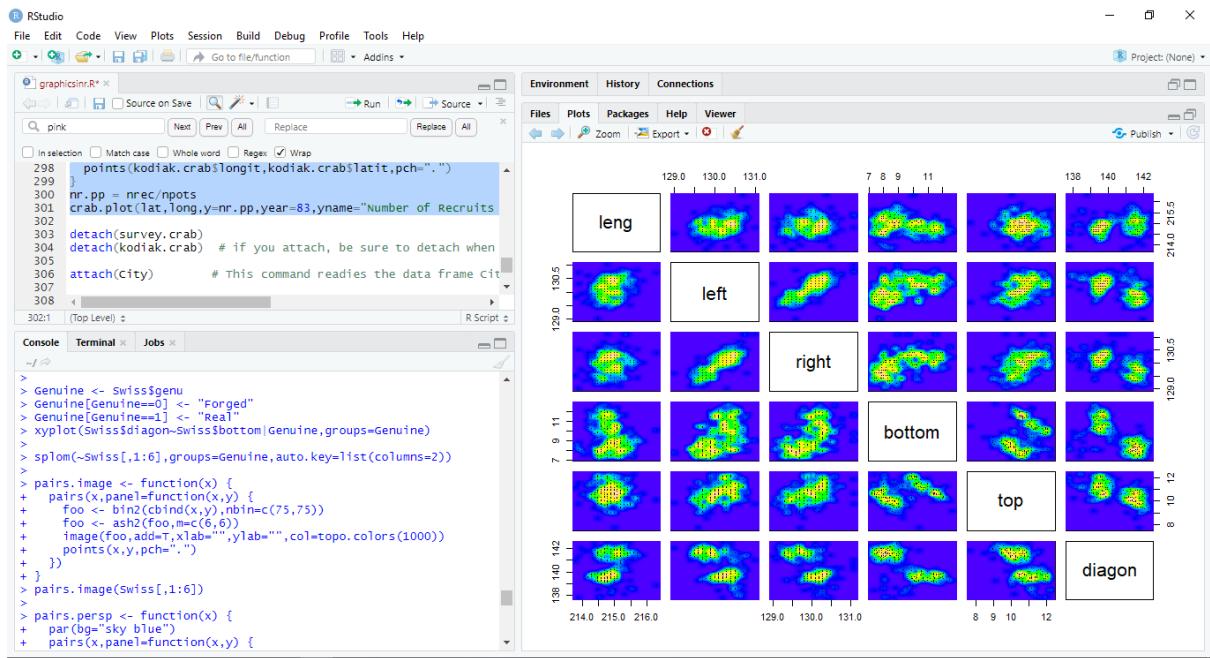


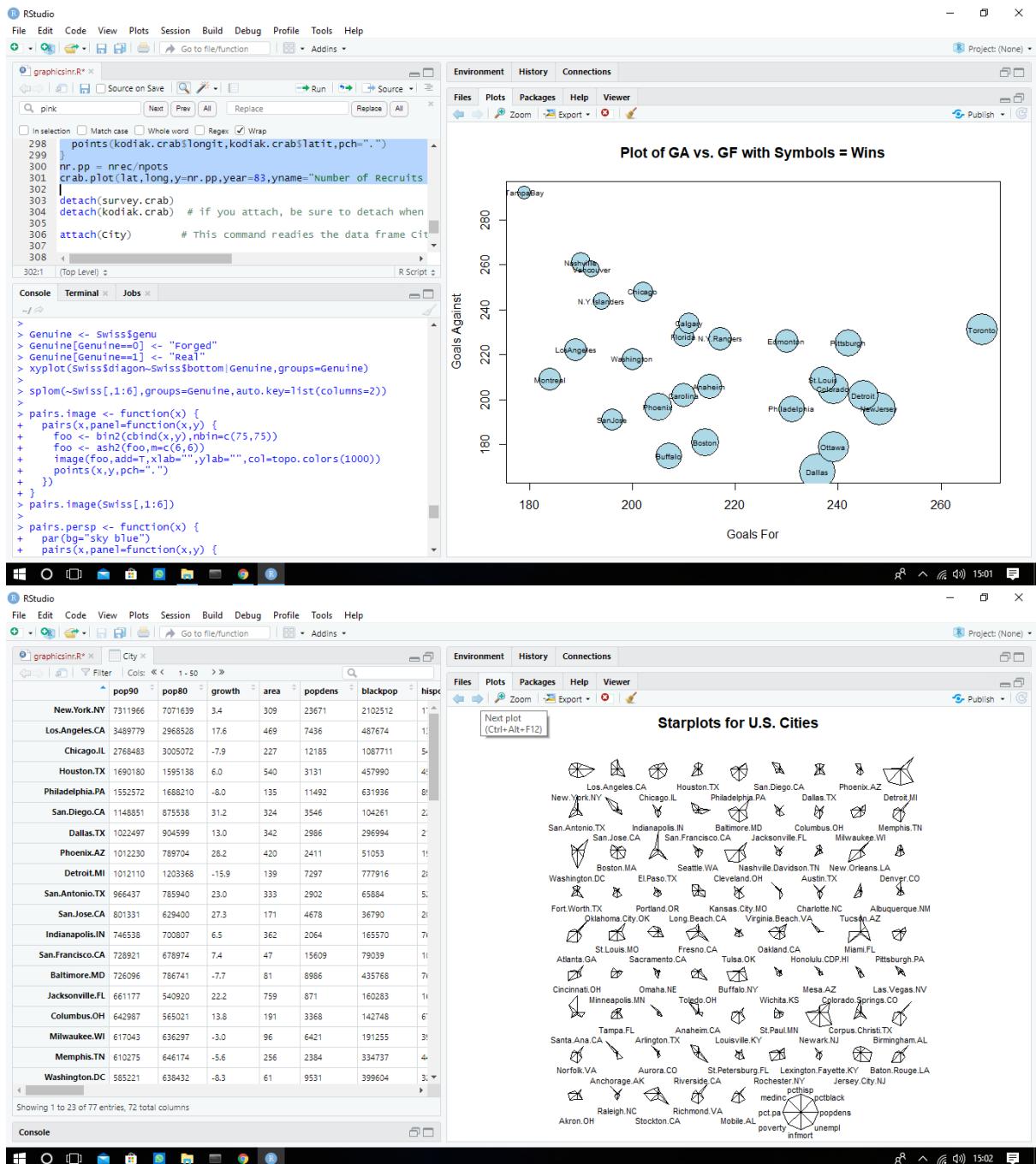


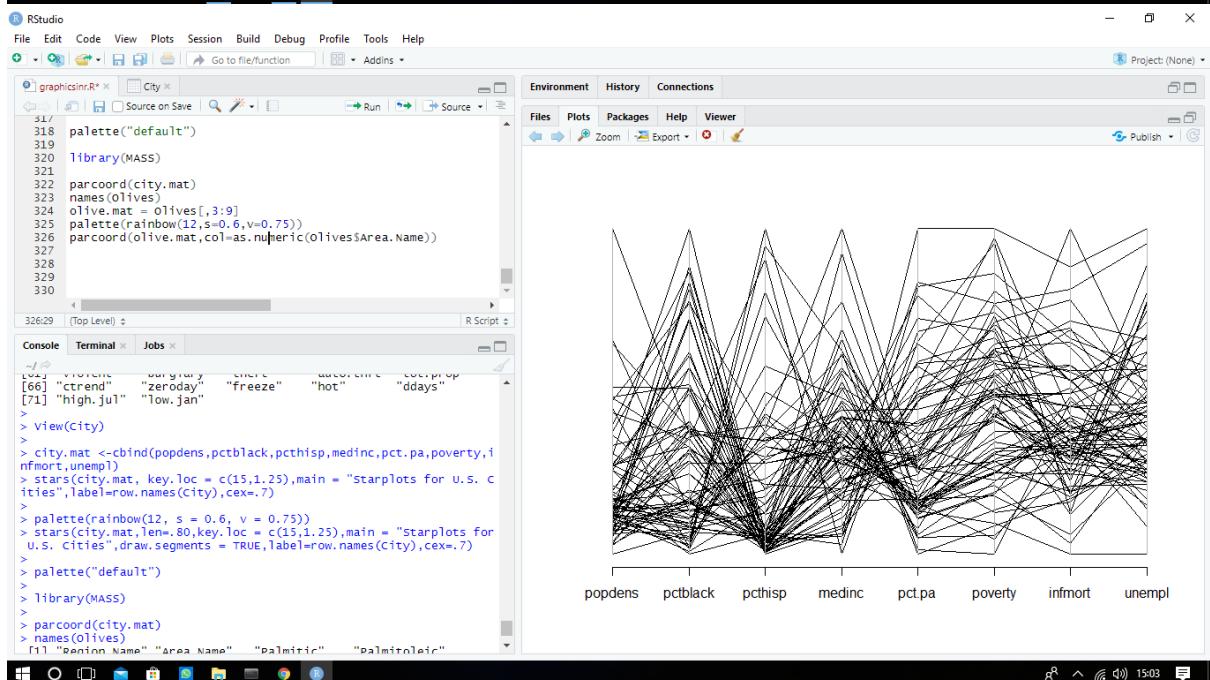
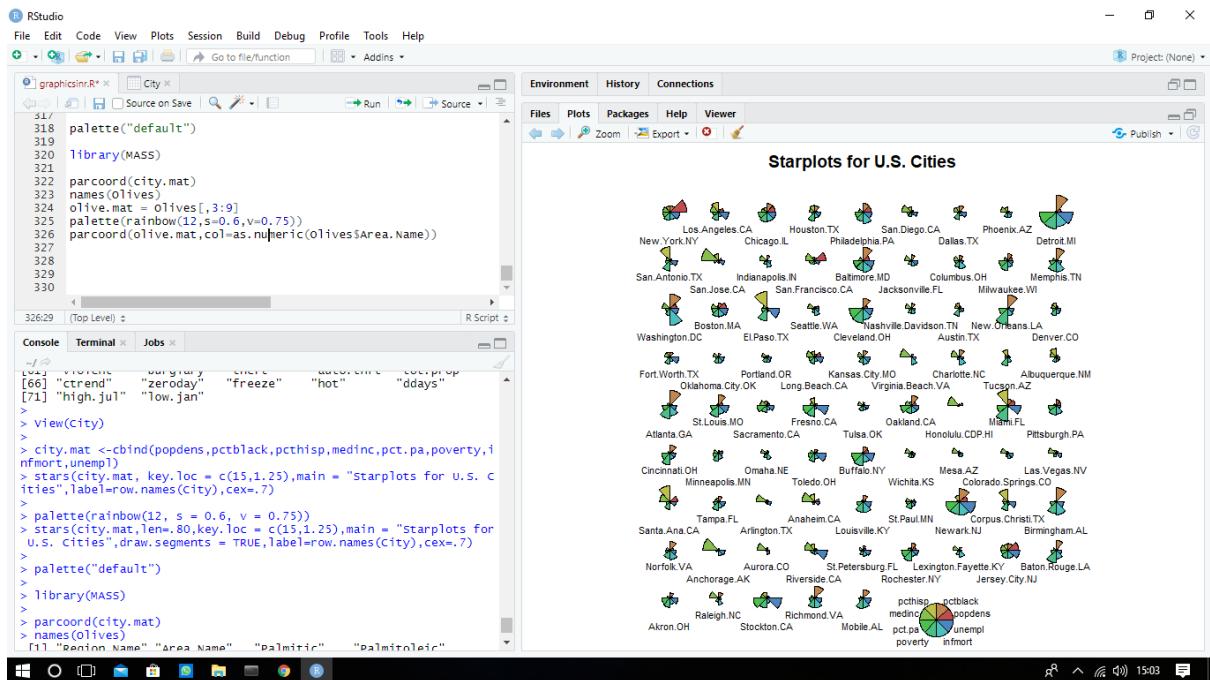


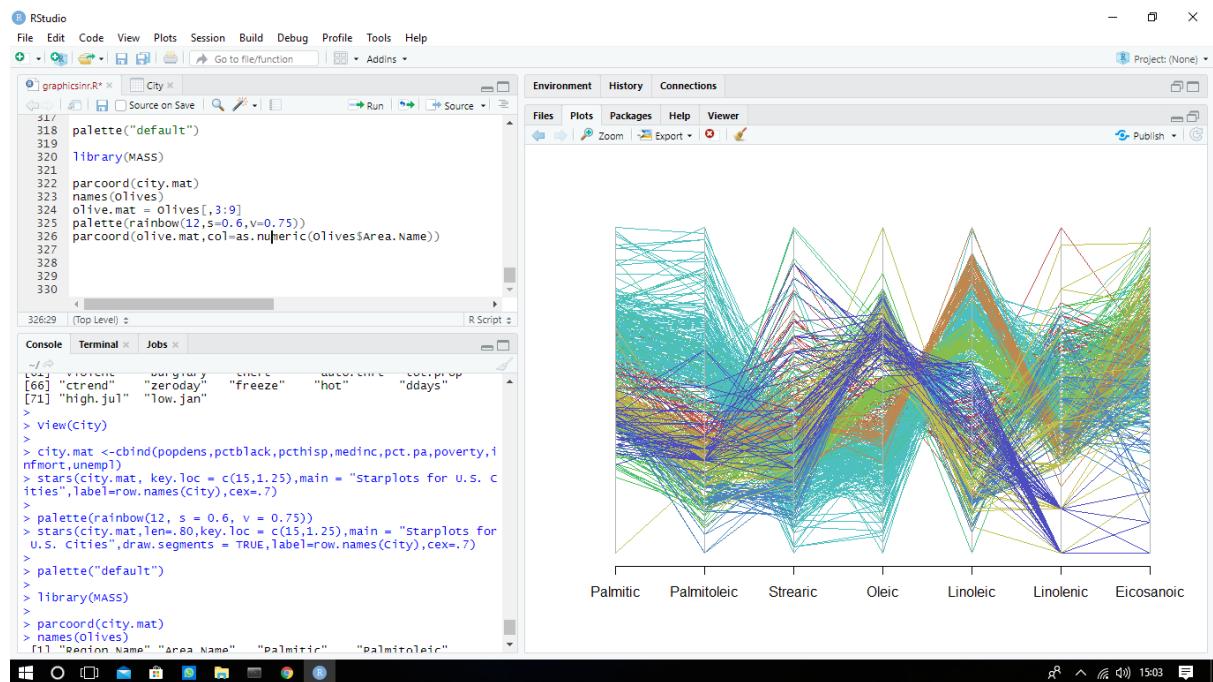


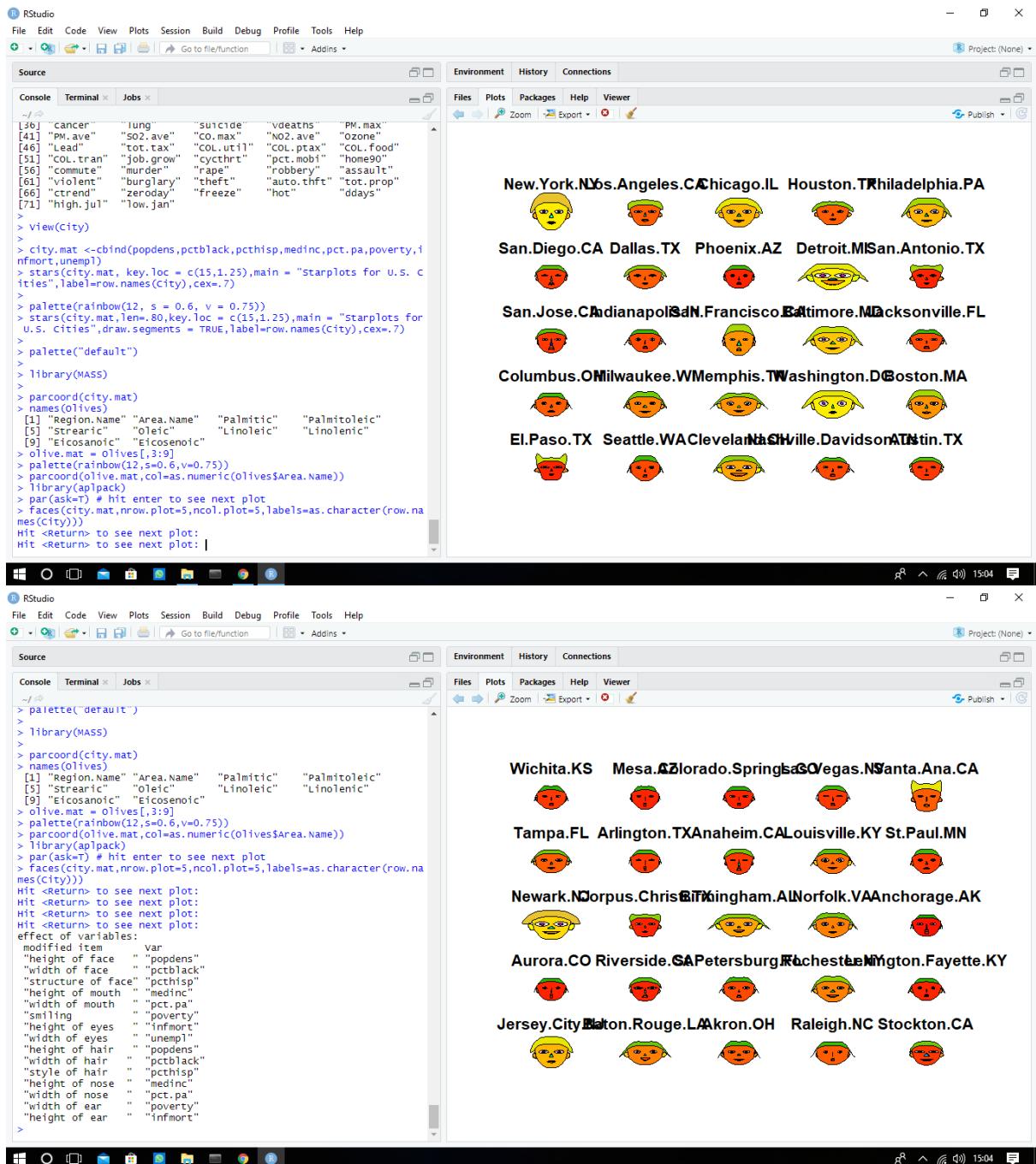












Conclusion : Thus we have studied basic and advanced graphics in R.

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.

The three forms of correlation presented here are Pearson, Kendall, and Spearman. The test determining the p -value for Pearson correlation is a parametric test that assumes that data are bivariate normal. Kendall and Spearman correlation use nonparametric tests.

In contrast, linear regression specifies one variable as the independent variable and another as the dependent variable. The resultant model relates the variables with a linear relationship.

The tests associated with linear regression are parametric and assume normality, homoscedasticity, and independence of residuals, as well as a linear relationship between the two variables.

Appropriate data

- For Pearson correlation, two interval/ratio variables. Together the data in the variables are bivariate normal. The relationship between the two variables is linear. Outliers can detrimentally affect results.
- For Kendall correlation, two variables of interval/ratio or ordinal type.
- For Spearman correlation, two variables of interval/ratio or ordinal type.
- For linear regression, two interval/ratio variables. The relationship between the two variables is linear. Residuals are normal, independent, and homoscedastic. Outliers can affect the results unless robust methods are used.

Hypotheses

- For correlation, null hypothesis: The correlation coefficient (r , τ , or ρ) for the sampled population is zero. Or, there is no correlation between the two variables.
- For linear regression, null hypothesis: The slope of the fit line for the sampled population is zero. Or, there is no linear relationship between the two variables.

Interpretation

- For correlation, reporting significant results as “Variable A was significantly correlated to Variable B” is acceptable. Alternatively, “A significant correlation between Variable A and Variable B was found.” Or, “Variables A, B, and C were significantly correlated.”

- For linear regression, reporting significant results as “Variable A was significantly linearly related to Variable B” is acceptable. Alternatively, “A significant linear regression between Variable A and Variable B was found.”

Output:

Code:

```
install.packages("ggpubr")
library("ggpubr")
```

```
my_data <- mtcars
head(my_data, 6)
```

```
ggscatter(my_data, x = "mpg", y = "wt",
      add = "reg.line", conf.int = TRUE,
      cor.coef = TRUE, cor.method = "pearson",
      xlab = "Miles/(US) gallon", ylab = "Weight (1000 lbs)")
```

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

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

```
print(relation)
```

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

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

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

```
# The predictor vector.
```

```
x <- c(151, 174, 138, 186, 128, 136, 179, 163, 152, 131)
```

```
# The response vector.
```

```
y <- c(63, 81, 56, 91, 47, 57, 76, 72, 62, 48)
```

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

```
# Find the weight of a person with height 170.
```

```
a <- data.frame(x = 170)
result <- predict(relation,a)
print(result)
```

```

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

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

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

# Save the file.
dev.off()

input <- mtcars[,c("mpg","disp","hp","wt")]
print(head(input))
input <- mtcars[,c("mpg","disp","hp","wt")]

# Create the relationship model.
model <- lm(mpg~disp+hp+wt, data = input)

# Show the model.
print(model)

# Get the Intercept and coefficients as vector elements.
cat("## ## # The Coefficient Values ## ## ",'\n')

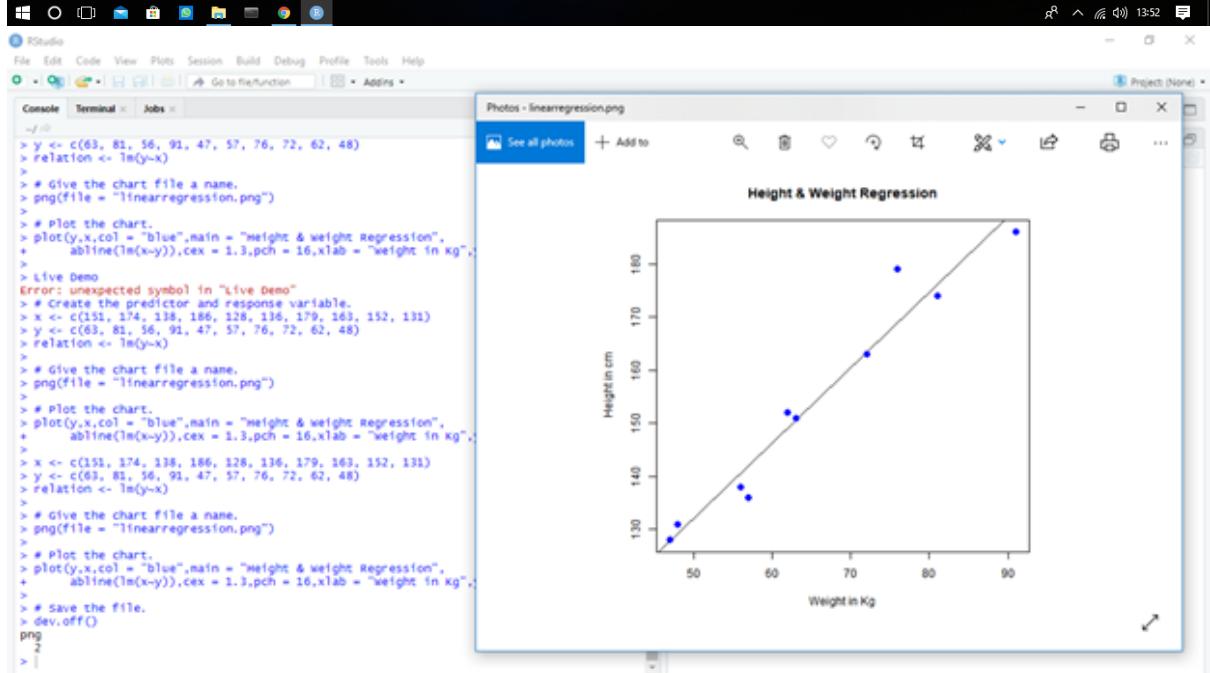
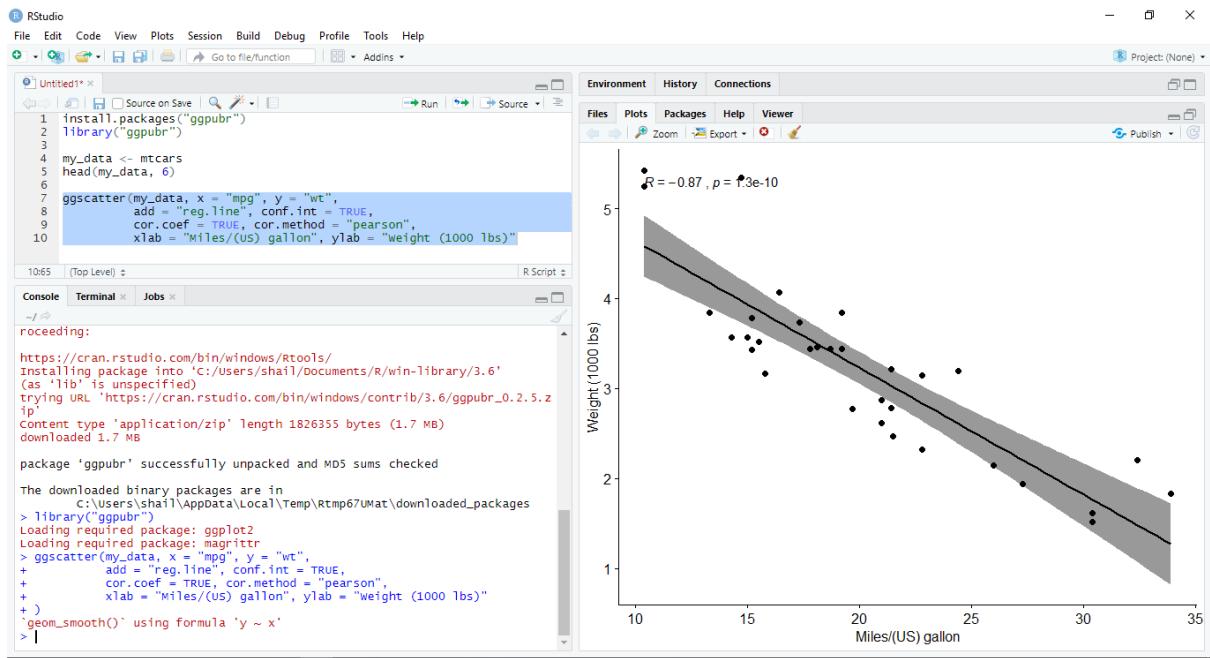
a <- coef(model)[1]
print(a)

Xdisp <- coef(model)[2]
Xhp <- coef(model)[3]
Xwt <- coef(model)[4]

print(Xdisp)
print(Xhp)
print(Xwt)

```

Screenshot:



The screenshot shows the RStudio interface with two sessions. The top session displays the following R code:

```

53 xnp <- coef(model)[3]
54 Xwt <- coef(model)[4]
55
56 print(xdisp)
57 print(xhp)
58 print(Xwt)
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
101
102
103
104
105
106
107
108
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161
162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215
216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269
270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323
324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377
378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431
432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485
486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
7010
7011

```

The bottom session shows the same R code, but the global environment viewer displays different variables and their values.

Conclusion: Thus, we have studied the basics of correlation, simple and multiple 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 the new centroid of each cluster.

These two steps are repeated until 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.

Clustering Types:

Clustering can be broadly divided into two subgroups:

- Hard clustering: in hard clustering, each data object or point either belongs to a cluster completely or not. For example in the Uber dataset, each location belongs to either one borough or the other.
- Soft clustering: in soft clustering, a data point can belong to more than one cluster with some probability or likelihood value. For example, you could identify some locations as the border points belonging to two or more boroughs.

The clustering algorithm that we are going to use is the K-means algorithm, which we can find in the package stats. The K-means algorithm accepts two parameters as input:

- The data;
- A **K** value, which is the number of groups that we want to create.

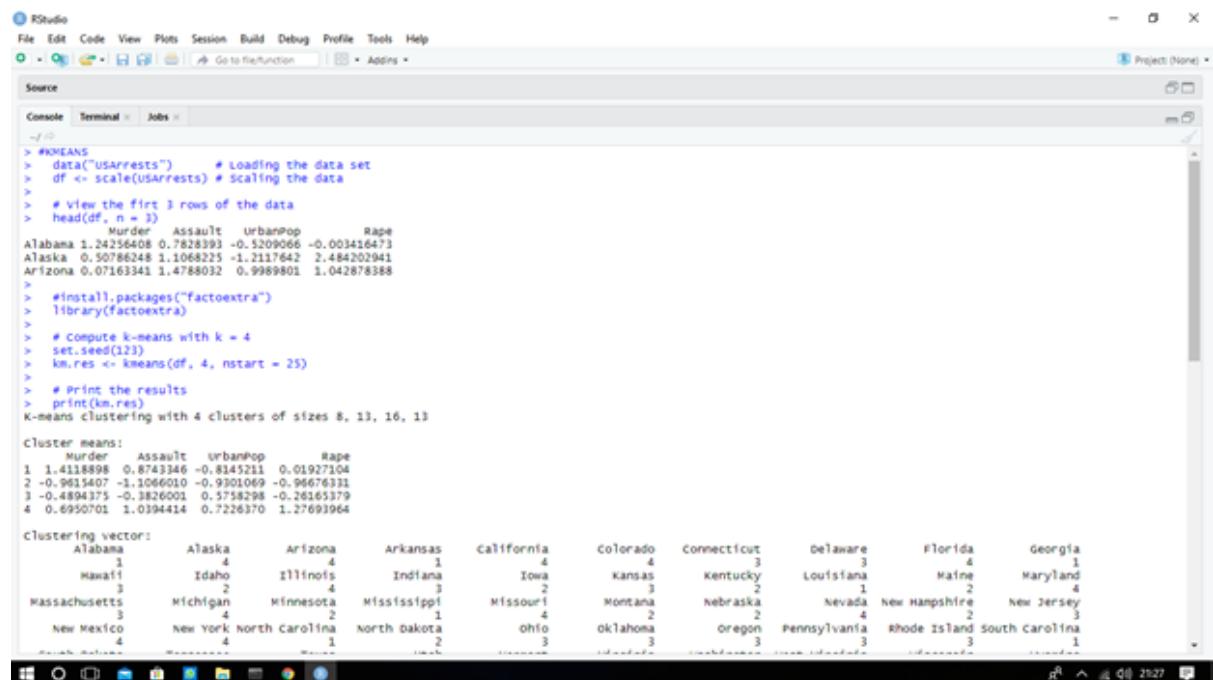
Conceptually, the K-means behaves as follows:

1. It chooses K centroids randomly;
2. Matches each point in the data (in our case, each mammal) with the closest centroid in an n-dimensional space where n is the number of features used in the clustering (in our example, 5 features — water, protein, fat, lactose, ash). After this step, each point belongs to a group.
3. Now, it recalculates the centroids as being the mean point (vector) of all other points in the group.
4. It keeps repeating the steps 2 and 3 until either when the groups are stabilized, that is, when no points are reallocated to another centroid or when it reaches the maximum number of iterations (the stats library uses 10 as default).

Output :

Code :

Screenshots :



The screenshot shows the RStudio interface with the Source tab selected. The code in the console is as follows:

```

> #KMEANS
> data("USArrests")      # Loading the data set
> df <- scale(USArrests) # scaling the data
>
> # View the first 3 rows of the data
> head(df, n = 3)
   Murder Assault UrbanPop Rape
Alabama 1.24256408 0.7828393 -0.5209066 -0.003416473
Alaska 0.50786248 1.1068225 -1.2117642 2.484202941
Arizona 0.07163341 1.4788032 0.9989801 1.042878388
>
> #install.packages("factoextra")
> library(factoextra)
>
> # Compute k-means with k = 4
> set.seed(123)
> km.res <- kmeans(df, 4, nstart = 25)
>
> # Print the results
> print(km.res)
K-means clustering with 4 clusters of sizes 8, 13, 16, 13

Cluster means:
   Murder Assault UrbanPop Rape
1 1.4118898 0.8743346 -0.8145211 0.01927104
2 -0.9815407 -1.1066010 -0.9301069 -0.96676331
3 -0.4894375 -0.8826001 0.5758298 -0.26165379
4  0.6950701 1.0394414 0.7226370 1.27693964

Clustering vector:
 Alabama    Alaska    Arizona    Arkansas    California    Colorado    Connecticut    Delaware    Florida    Georgia
          1           4           4           1           4           4           3           3           4           4           1
 Hawaii    Idaho    Illinois    Indiana    Iowa    Kansas    Kentucky    Louisiana    Maine    Maryland
          3           2           2           4           3           2           3           2           1           2           4
 Massachusetts    Michigan    Minnesota    Mississippi    Missouri    Montana    Nebraska    Nevada    New Hampshire    New Jersey
          3           4           2           1           4           4           2           2           1           2           3
 New Mexico    New York    North Carolina    North Dakota    Ohio    Oklahoma    Oregon    Pennsylvania    Rhode Island    South Carolina
          4           4           1           2           3           3           3           3           3           3           1
  
```

```
RStudio
File Edit Code View Plots Session Build Debug Profile Tools Help
Project (None)
Source
Console Terminal Jobs

Clustering vector:
Alabama    Alaska    Arizona    Arkansas    California    Colorado    Connecticut    Delaware    Florida    Georgia
1          4          4          1          4          4          3          3          4          1          1
Hawaii    Idaho    Illinois    Indiana    Iowa    Kansas    Kentucky    Louisiana    Maine    Maryland
3          2          4          3          2          3          2          1          2          2          4
Massachusetts    Michigan    Minnesota    Mississippi    Missouri    Montana    Nebraska    Nevada    New Hampshire    New Jersey
3          4          2          3          4          1          2          2          4          3          3
New Mexico    New York    North Carolina    North Dakota    Ohio    Oklahoma    Oregon    Pennsylvania    Rhode Island    South Carolina
4          4          1          1          3          3          3          2          3          1          1
South Dakota    Tennessee    Texas    Utah    Vermont    Virginia    Washington    West Virginia    Wisconsin    Wyoming
2          1          4          3          2          3          3          3          2          2          3

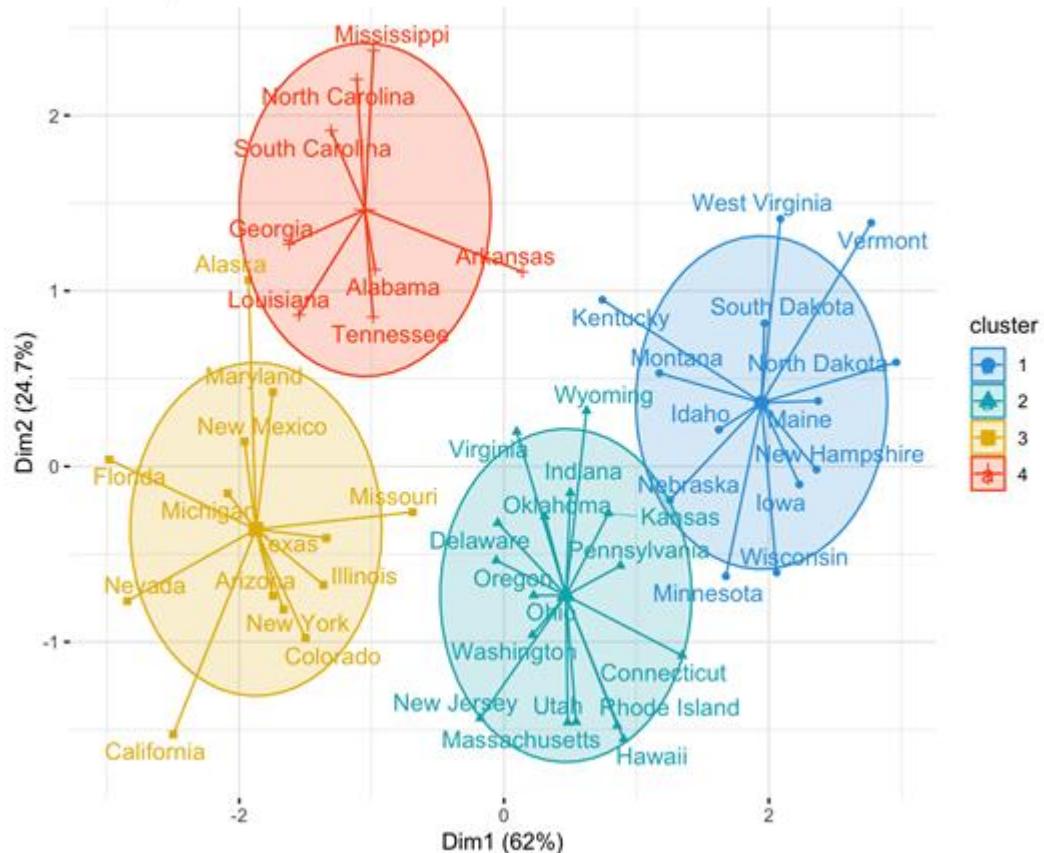
within cluster sum of squares by cluster:
[1] 8.316061 11.952463 16.212213 19.922437
(between_SS / total_SS = 71.2 %)

Available components:
[1] "cluster"    "centers"    "totss"      "withins"    "tot.withinss" "betweenss"   "size"       "iter"       "ifault"
> aggregate(usArrests, by=list(cluster=km.res$cluster), mean)
cluster Murder Assault UrbanPop Rape
1 13.9750 243.62500 53.75000 21.41250
2 3.60000 78.53846 52.07692 12.17692
3 5.65625 138.87500 73.87500 18.78125
4 4.10.85538 257.38462 76.00000 33.19231
>
> dd <- cbind(usArrests, cluster = km.res$cluster)
> head(dd)
Murder Assault UrbanPop Rape cluster
Alabama 13.2 236 58 21.2 1
Alaska 10.0 263 48 44.5 4
Arizona 8.1 294 80 31.0 4
Arkansas 8.8 190 50 19.5 1
California 9.0 276 91 40.6 4
Colorado 7.9 204 78 38.7 4
>
> dd <- cbind(usArrests, cluster = km.res$cluster)
R ^ Q 2527
```

```
RStudio
File Edit Code View Plots Session Build Debug Profile Tools Help
Project (None)
Source
Console Terminal Jobs

-> dd <- cbind(usArrests, cluster = km.res$cluster)
4 10.85538 257.38462 76.00000 33.19231
>
> dd <- cbind(usArrests, cluster = km.res$cluster)
> head(dd)
Murder Assault UrbanPop Rape cluster
Alabama 13.2 236 58 21.2 1
Alaska 10.0 263 48 44.5 4
Arizona 8.1 294 80 31.0 4
Arkansas 8.8 190 50 19.5 1
California 9.0 276 91 40.6 4
Colorado 7.9 204 78 38.7 4
>
> dd <- cbind(usArrests, cluster = km.res$cluster)
> head(dd)
Murder Assault UrbanPop Rape cluster
Alabama 13.2 236 58 21.2 1
Alaska 10.0 263 48 44.5 4
Arizona 8.1 294 80 31.0 4
Arkansas 8.8 190 50 19.5 1
California 9.0 276 91 40.6 4
Colorado 7.9 204 78 38.7 4
>
> head(km.res$cluster, 4)
Alabama Alaska Arizona Arkansas
1 1 4 4 1
>
> # Cluster size
> km.res$size
[1] 8 13 16 13
> # Cluster means
> km.res$means
Murder Assault UrbanPop Rape
1 1.4118898 0.8743346 -0.8145211 0.01927104
2 -0.9815407 -1.1066010 -0.9301069 -0.96676331
3 -0.4894375 -0.3826001 0.5758298 -0.26165379
4 0.6950701 1.0304414 0.7226370 1.27693964
>
```

Cluster plot



Conclusion : This, we have studied clustering in R - k-means

Experiment No. 11

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

Theory :

In statistics, the logistic model (or logit model) is used to model the probability of a certain class or event existing such as pass/fail, win/lose, alive/dead or healthy/sick. This can be extended to model several classes of events such as determining whether an image contains a cat, dog, lion, etc. Each object being detected in the image would be assigned a probability between 0 and 1 and the sum adding to one.

Logistic regression is a statistical model that in its basic form uses a logistic function to model a binary dependent variable, although many more complex extensions exist. In regression analysis, logistic regression (or logit regression) is estimating the parameters of a logistic model (a form of binary regression).

R makes it very easy to fit a logistic regression model. The function to be called is `glm()` and the fitting process is not so different from the one used in linear regression. In this post I am going to fit a binary logistic regression model and explain each step.

Naïve Bayes classification is a kind of simple probabilistic classification method based on Bayes' theorem with the assumption of independence between features. The model is trained on a training dataset to make predictions by `predict()` function. This article introduces two functions `naiveBayes()` and `train()` for the performance of Naïve Bayes classification.

Output:

Code:

```
training_data <- read.csv("aps_failure_training_set.csv", skip = 20, na.strings = "na")
test_data <- read.csv("aps_failure_test_set.csv", skip = 20, na.strings = "na")
```

```
dim(training_data)
dim(test_data)
```

```
library(dplyr)
library(caret)
library(caretEnsemble)
library(mice)
library(doParallel)
library(car)
```

```
install.packages("car")
glimpse(training_data)
```

```

glimpse(test_data)

summary(training_data$class)
summary(test_data$class)

options(digits = 2)
prop.table(table(training_data$class))
prop.table(table(test_data$class))

options(scipen = 999)
summary_df <- do.call(cbind, lapply(training_data[, 2:ncol(training_data)], summary))
summary_df_t <- as.data.frame(round(t(summary_df),0))
names(summary_df_t)[7] <- paste("Missing_values")
summary_df_t_2 <- summary_df_t %>%
  mutate(obs = nrow(training_data),
        Missing_prop = Missing_values / obs)
print(summary_df_t_2)

summary_df_t_2 %>% summarise(Min = mean(Min.),
                               first_Q = mean(`1st Qu.`),
                               Median = median(Median),
                               Mean = mean(Mean),
                               third_Q = mean(`3rd Qu.`),
                               Max = max(Max.),
                               mean_MV = mean(Missing_values),
                               obs = mean(obs),
                               mean_MV_perc = mean_MV / obs)

#replicate our sets
training_data_bind <- training_data
test_data_bind <- test_data
#create a new column "set" to label the observations
training_data_bind$set <- "TRAIN"
test_data_bind$set <- "TEST"
#merge them into 1 single set
full_dataset <- rbind(training_data_bind, test_data_bind)
dim(full_dataset)

set.seed(123)
imputed_full <- mice(full_dataset,
                      m=1,
                      maxit = 5,
                      method = "pmm",
                      seed = 500)

full_imputed <- complete(imputed_full, 1)
dim(full_imputed)

#subset the full_imputed_filtered dataset

```



```
trControl = my_ctrl)
```

```
model_list$glm  
model_list$nb
```

Screenshots:

STEPS:

1. Download the dataset from <https://archive.ics.uci.edu/ml/machine-learning-databases/00421/>
 2. Load the data set and View Data (str or dplyr's glimpse)

The screenshot shows the RStudio interface with a script named 'exp1.R' open. The code reads two CSV files ('apps_failure_training_set.csv' and 'apps_failure_test_set.csv') and performs dimensionality checks. It then loads several packages: dplyr, caret, caretEnsemble, mice, parallel, and car. It installs the 'car' package if it's not already installed and then loads it. The 'glimpse' function is used to inspect the 'training_data' dataset.

```
1 training_data <- read.csv("apps_failure_training_set.csv", skip = 20, na.strings = "na")
2 test_data <- read.csv("apps_failure_test_set.csv", skip = 20, na.strings = "na")
3
4 dim(training_data)
5 dim(test_data)
6
7 library(dplyr)
8 library(caret)
9 library(caretEnsemble)
10 library(mice)
11 library(parallel)
12 library(car)
13
14 #install.packages("car")
15
16 glimpse(training_data)
17
```

The 'Global Environment' sidebar on the right lists objects like 'test_16000' and 'train_60000' with their respective sizes and types.

3. Check How many neg and pos do we have in each set and check proportions

The screenshot shows the RStudio interface. The top menu bar includes File, Edit, Code, View, Plots, Session, Build, Debug, Profile, Tools, Help, and a Go to function search bar. The left sidebar has tabs for Project (None), Environment, History, Import Dataset, Global Environment, Data, and Values. The main area contains a script editor with code for loading packages like dplyr, caret, and parallel, and performing data summaries and prop.table calculations. Below the script editor is the R Script tab. The bottom left shows the Console tab with R command history. The bottom right shows the Windows taskbar with various icons.

```
explorR.R
1 library(dplyr)
2 library(caret)
3 library(caretEnsemble)
4 library(rms)
5 library(parallel)
6 library(car)
7 library(gbm)
8 library(e1071)
9 library(pROC)
10 library(rpart)
11 library(rpartOrdinal)
12 library(rpart)
13 library(rpart)
14 install.packages("car")
15 glimpse(training_data)
16 glimpse(test_data)
17
18 summary(training_data$class)
19 summary(test_data$class)
20
21 options(digits = 2)
22 prop.table(table(training_data$class))
23 prop.table(table(test_data$class))
24
25 (Top Level) :
```

Console

```
> summary(training_data$class)
  neg  pos 
59000 1000 
> summary(test_data$class)
  neg  pos 
15625 375 
> options(digits = 2)
> prop.table(table(training_data$class))

  neg  pos 
0.983 0.017 
> prop.table(table(test_data$class))

  neg  pos 
0.977 0.023 
> |
```

4. Rest of the features. Summary statistics. Understand your data

```

summary_df <- do.call(cbind, lapply(training_data,
                                     function(x) sum(x == NA, na.rm = TRUE) / ncol(training_data), summary))
names(summary_df_t)[7] <- paste("Missing_values", 0)
summary_df_t<- summary_df_t %>
  mutate(obs = nrow(training_data),
        Missing_prop = Missing_values / obs)
print(summary_df_t_2)

```

	Min	1st Qu.	Median	Mean	3rd Qu.	Max.	Missing_values	obs	Missing_prop
1	0	834	30776	59336	48668	2746564	0	60000	0.0000
2	0	0	0	1	0	204	46329	60000	0.7722
3	0	16	152	356014263	964	2330706796	3335	60000	0.0556
4	0	24	126	190621	430	8584297742	14881	60000	0.2477
5	0	0	126	0	7	21050	2500	60000	0.0417
6	0	0	0	11	0	20070	2500	60000	0.0417
7	0	0	0	222	0	3376892	671	60000	0.0112
8	0	0	0	976	0	4109372	671	60000	0.0112
9	0	0	0	8668	0	10552856	671	60000	0.0112
10	0	0	0	88591	0	63402074	671	60000	0.0112
11	0	308	3872	437097	49522	228810570	671	60000	0.0112
12	0	13834	176020	1108374	913964	179187978	671	60000	0.0112
13	0	10608	930336	1657818	1886508	94020666	671	60000	0.0112
14	0	0	119204	499310	588820	63346754	671	60000	0.0112
15	0	0	1786	35570	26690	17702522	671	60000	0.0112
16	0	0	5115	364	21398514	671	60000	0.0112	
17	0	29733	1002420	1809931	1601366	74247318	645	60000	0.0107
18	0	0	0	9017	0	16512852	629	60000	0.0105
19	0	0	0	1144	0	5629340	629	60000	0.0105
20	0	0	0	979	0	104494924	4400	60000	0.0733
21	0	0	0	59130	1204	34762578	642	60000	0.0107
22	0	0	0	93281	2364	55903508	629	60000	0.0105
23	0	73532	501865	3241027	317618	16661830	645	60000	0.0107

```

summary_df <- do.call(cbind, lapply(training_data,
                                     function(x) sum(x == NA, na.rm = TRUE) / ncol(training_data), summary))
names(summary_df_t)[7] <- paste("Missing_values", 0)
summary_df_t<- summary_df_t %>
  mutate(obs = nrow(training_data),
        Missing_prop = Missing_values / obs)
print(summary_df_t_2)

```

	Min	1st Qu.	Median	Mean	3rd Qu.	Max.	Missing_values	obs	Missing_prop
1	0	834	30776	59336	48668	2746564	0	60000	0.0000
2	0	0	0	1	0	204	46329	60000	0.7722
3	0	16	152	356014263	964	2330706796	3335	60000	0.0556
4	0	24	126	190621	430	8584297742	14881	60000	0.2477
5	0	0	126	0	7	21050	2500	60000	0.0417
6	0	0	0	11	0	20070	2500	60000	0.0417
7	0	0	0	222	0	3376892	671	60000	0.0112
8	0	0	0	976	0	4109372	671	60000	0.0112
9	0	0	0	8668	0	10552856	671	60000	0.0112
10	0	0	0	88591	0	63402074	671	60000	0.0112
11	0	308	3872	437097	49522	228810570	671	60000	0.0112
12	0	13834	176020	1108374	913964	179187978	671	60000	0.0112
13	0	10608	930336	1657818	1886508	94020666	671	60000	0.0112
14	0	0	119204	499310	588820	63346754	671	60000	0.0112
15	0	0	1786	35570	26690	17702522	671	60000	0.0112
16	0	0	5115	364	21398514	671	60000	0.0112	
17	0	29733	1002420	1809931	1601366	74247318	645	60000	0.0107
18	0	0	0	9017	0	16512852	629	60000	0.0105
19	0	0	0	1144	0	5629340	629	60000	0.0105
20	0	0	0	979	0	104494924	4400	60000	0.0733
21	0	0	0	59130	1204	34762578	642	60000	0.0107
22	0	0	0	93281	2364	55903508	629	60000	0.0105
23	0	73532	501865	3241027	317618	16661830	645	60000	0.0107

5. Our features present more than 8% missing values on average. Because it's a lot of information we don't want to lose, our approach is going to be to impute them. Impute the full dataset using mean imputation

RStudio

```

53 full_dataset <- rbind(training_data_bind, test_data_bind)
54 dim(full_dataset)
55
56 set.seed(123)
57 imputed_full <- mice(full_dataset,
58   m=1,
59   maxit = 5,
60   method = "mean",
61   seed = 500)
62
63 full_imputed <- complete(imputed_full, 1)
64 dim(full_imputed)
65
66 #subset the full_imputed filtered dataset
67 training_data_imp <- subset(training_data)
68 test_data_imp <- subset(test_data)
69 #drop the "set" column, we don't need it anymore
70 training_data_impset <- NULL
71 test_data_impset <- NULL
72 #check dimensions
73 dim(training_data_imp)
74 dim(test_data_imp)
75
76
77
78
79
80
81
82
83
84

```

Console Terminal Jobs

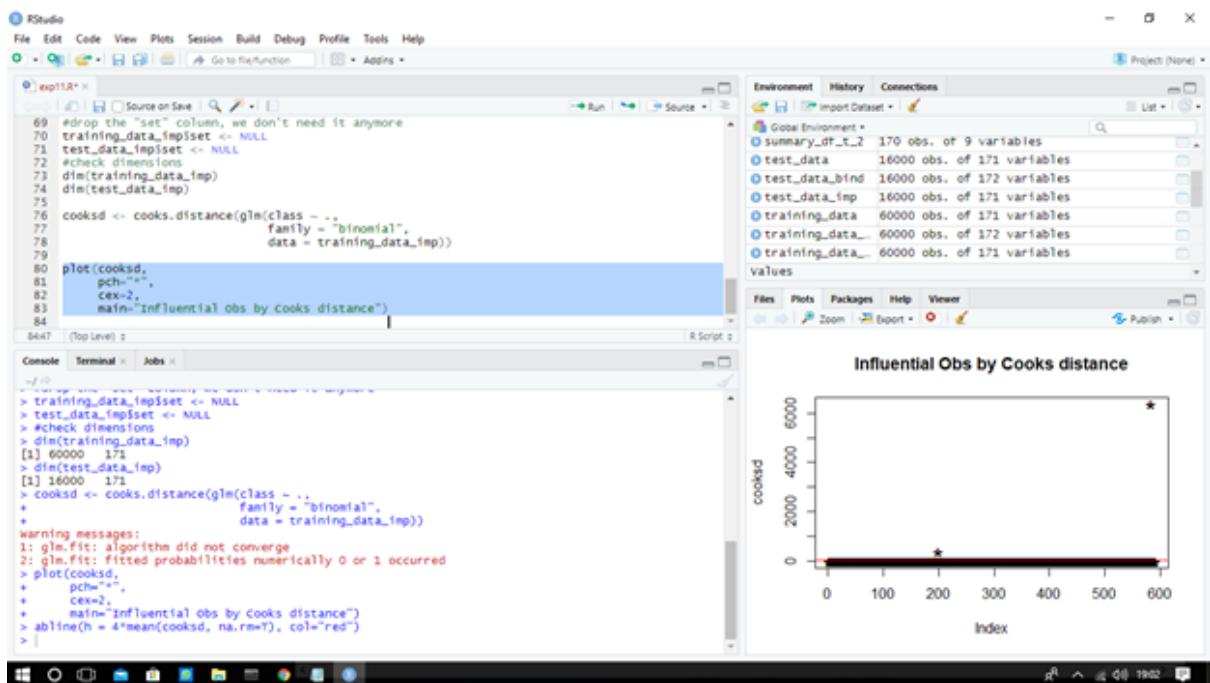
```

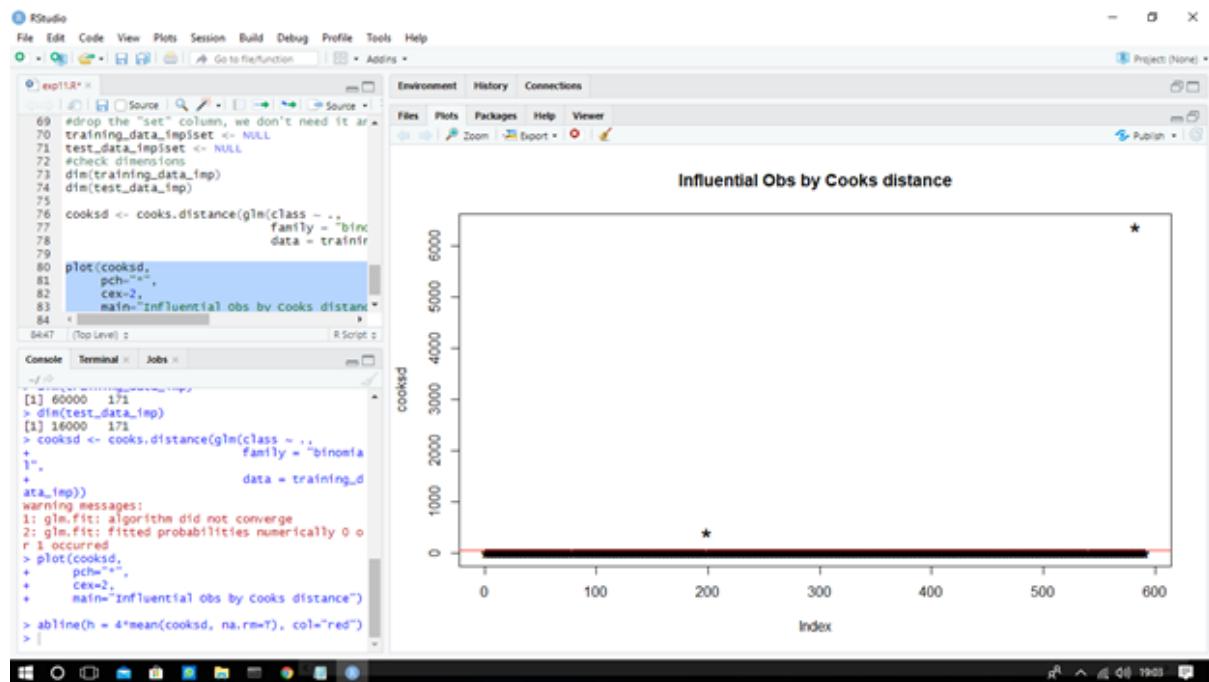
> #subset the full_imputed filtered dataset
> training_data_imp <- subset(training_data)
> test_data_imp <- subset(test_data)
> #drop the "set" column, we don't need it anymore
> training_data_impset <- NULL
> test_data_impset <- NULL
> #check dimensions
> dim(training_data_imp)
[1] 60000  171
> dim(test_data_imp)
[1] 16000  171
>

```

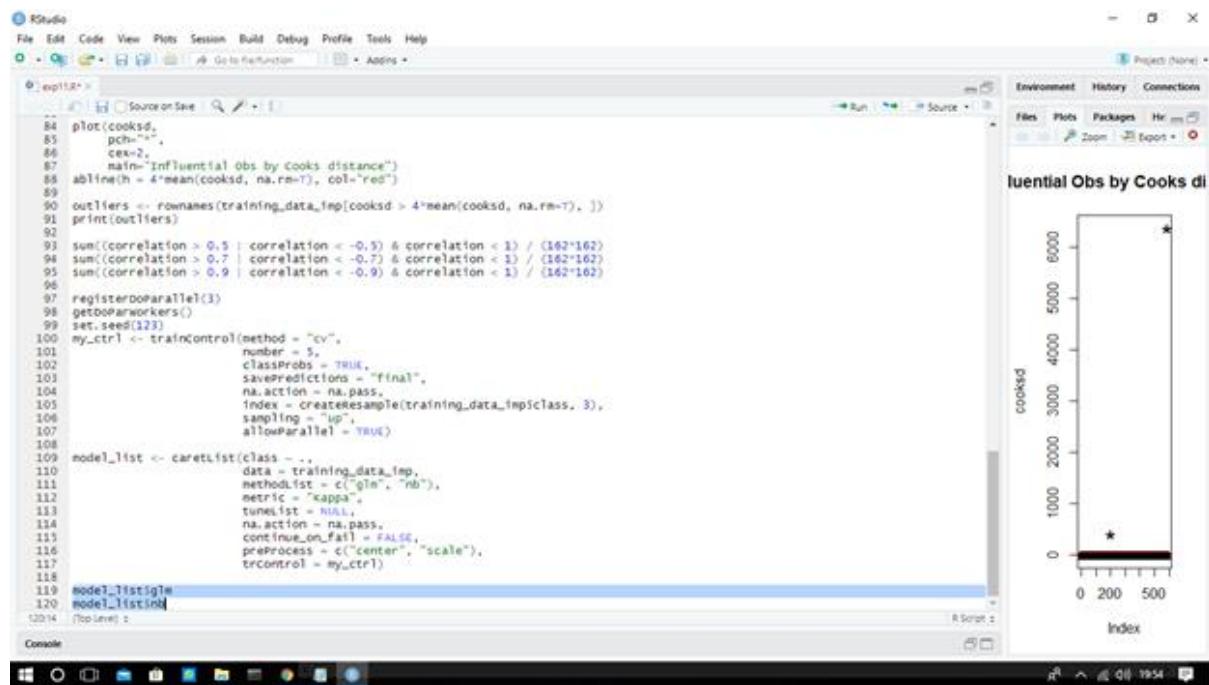
R Script Project (None)

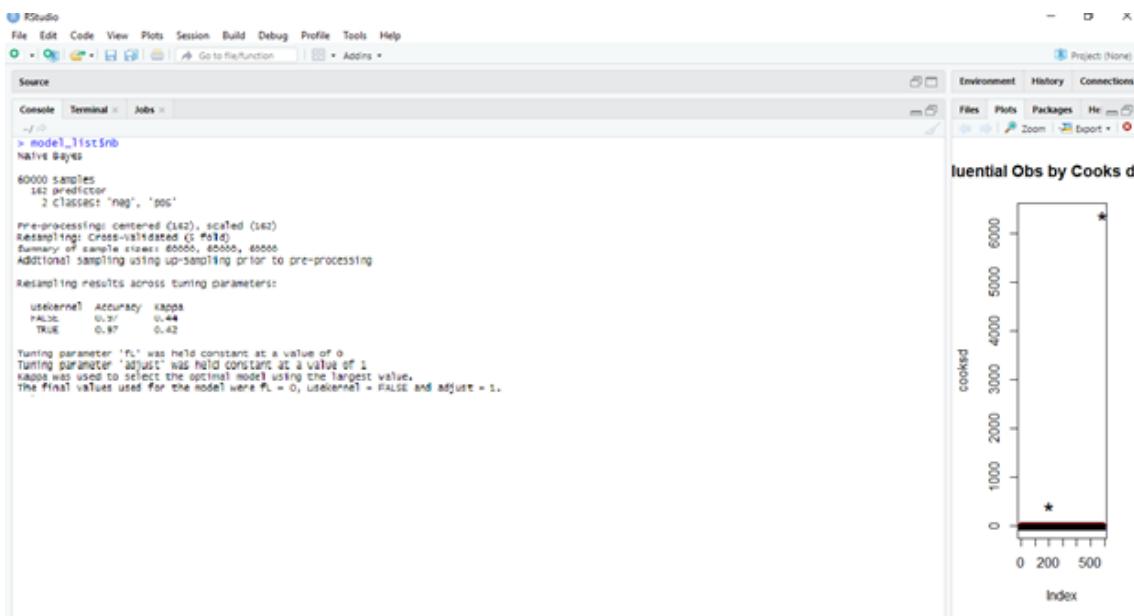
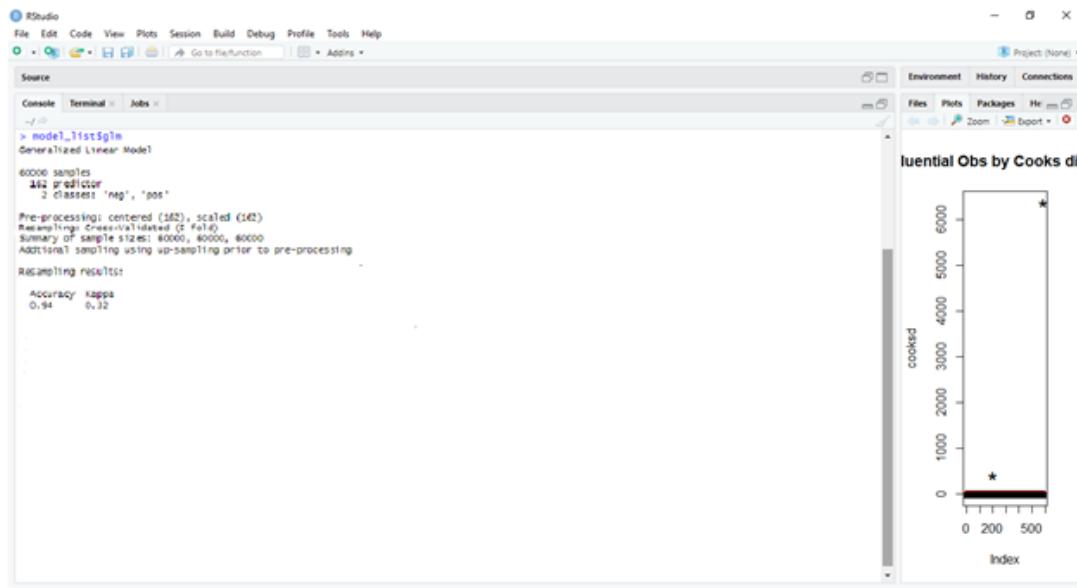
6. Check for outliers and other inconsistent data points. Box-plots. Cooks' distance. DBSCAN





7. We'll use caretEnsemble's caretList() to train both at the same time and with the same resampling .





Conclusion: Thus, we have studied classification on large and noisy dataset with R - logistic regression and naive bayes.