

R programming

Introduction:-

- R language is a open source program used for performing statistical operations
- It is a command line driven program and execute one command at a time.
- It performs statistical operations and generates R data analysis reports in graphical or text formats.

History of R Language:-

- John chambers and colleagues developed R at bell laboratories.
- R is an implementation of the S programming language and combines with lexical scoping semantics inspired by Scheme.

R-Studio:-

- R studio is an integrated development environment (IDE) for R language. R studio is a code editor and development environment, with some nice features that make code development in R easy and fun.

features of R studio:-

- code highlighting
- automatic bracket matching
- code completion
- it is free to use.

components of R studio:

- 1) Source = Multiple lines of code can be entered here.
- 2) Console = All the interactive work of R programming is performed in this window.
- 3) Files = where the user can browse folders and files.
- 4) plots = where R displays the User's plots.
- 5) workspace and history = shows overview of workspace, history of the commands.
- 6) Packages = this is where the user can view a list of all the installed packages.
- 7) Help = this is where you can browse the built-in Help system of R

→ X Sepal.Length Sepal.Width Petal.Length Petal.Width Species

1	1	5.1	3.5	1.4	0.2	setosa
2	2	4.9	3.0	1.4	0.2	setosa
3	3	4.7	3.2	1.3	0.2	setosa

149 149 6.2 3.4 5.4 2.3 = virginica
50 150 5.9 3.0 5.1 1.8 virginica

→ "x" "Sepal.Length" "Sepal.Width" "Petal.Length" "Petal.Width"
"species"

→ 'data.frame': 150 obs of 6 variables

```
$x : int 1 2 3 4 . . .  
$Sepal.Length: num 5.1 4.9 4.7 4.6 . . .  
$Sepal.Width: num 3.5 3.3 3.2 3.1 . . .  
$Petal.Length: num 1.4 1.4 1.3 1.5 . . .  
$Petal.Width: num 0.2 0.2 0.2 0.2 . . .  
$Species : factor w/3 levels "setosa", "Versicolor", "virginica".
```

→ [D] 4.3

→ [177.9]

$\rightarrow [1] 4.3 \text{ } 7.9$

→ [1] 5.843333

Date :

Q) Load the 'iris.csv' file and display the names and type of each column. Find statistics such as min, max, range, mean, median, variance, standard deviation for each column of data.

> iris <- read.csv("c:/Users/Admin/Desktop/iris.csv",
header=TRUE, sep=",")

> print(iris)

> colnames(iris)

> str(iris)

> print(min(iris\$Sepal.Length))

> print(max(iris\$Sepal.Length))

> print(range(iris\$Sepal.Length))

> print(mean(iris\$Sepal.Length))

- $\rightarrow [1] 5.8$ (iii) $\frac{1}{2} \ln(1 + \sqrt{1 + 4x})$
- $\rightarrow [1] 0.6856935$ (iii) $\frac{1}{2} \ln(1 + \sqrt{1 + 4x})$
- $\rightarrow [1] 0.8280661$ (iii) $\frac{1}{2} \ln(1 + \sqrt{1 + 4x})$
- $\rightarrow [1] 2$ (iii) $\frac{1}{2} \ln(1 + \sqrt{1 + 4x})$
- $\rightarrow [1] 4.4$ (iii) $\frac{1}{2} \ln(1 + \sqrt{1 + 4x})$
- $\rightarrow [1] 2.0 \quad 4.4$ (iii) $\frac{1}{2} \ln(1 + \sqrt{1 + 4x})$
- $\rightarrow [1] 3.057333$ (iii) $\frac{1}{2} \ln(1 + \sqrt{1 + 4x})$
- $\rightarrow [1] 3$ (iii) $\frac{1}{2} \ln(1 + \sqrt{1 + 4x})$
- $\rightarrow [1] 0.1899794$ (iii) $\frac{1}{2} \ln(1 + \sqrt{1 + 4x})$
- $\rightarrow [1] 0.4358663$ (iii) $\frac{1}{2} \ln(1 + \sqrt{1 + 4x})$

> print(median(iris \$ Sepal.Length))	(5.0)
> print(var(iris \$ Sepal.Length))	(4.35)
> print(sd(iris \$ Sepal.Length))	(2.435)
> print(min(iris \$ Sepal.Width))	(1.0)
> print(max(iris \$ Sepal.Width))	(4.3)
> print(range(iris \$ Sepal.Width))	(2.8)
> print(mean(iris \$ Sepal.Width))	(3.73)
> print(median(iris \$ Sepal.Width))	(3.7)
> print(var(iris \$ Sepal.Width))	(1.11)
> print(sd(iris \$ Sepal.Width))	(1.05)

$\rightarrow [1] 1$

$\rightarrow [1] 6.9$

$\rightarrow [1] 1.0 \quad 6.9$

$\rightarrow [1] 3.758$

$\rightarrow [1] 4.35$

$\rightarrow [1] 3.116278$

$\rightarrow [1] 1.765278$

$\rightarrow [1] 0.1$

$\rightarrow [1] 2.5$

$\rightarrow [1] 0.1 \quad 2.5$

$\rightarrow [1] 1.199333$

$\rightarrow [1] 1.3$

$\rightarrow [1] 0.5810063$

$\rightarrow [1] 0.7622377$

Date :

> print(min(iris\$Petal.Length))		
> print(max(iris\$Petal.Length))		
> print(range(iris\$Petal.Length))		
> print(mean(iris\$Petal.Length))		
> print(median(iris\$Petal.Length))		
> print(var(iris\$Petal.Length))		
> print(sd(iris\$Petal.Length))		
> print(min(iris\$Petal.Width))		
> print(max(iris\$Petal.Width))		
> print(range(iris\$Petal.Width))		
> print(mean(iris\$Petal.Width))		
> print(median(iris\$Petal.Width))		
> print(var(iris\$Petal.Width))		
> print(sd(iris\$Petal.Width))		

- ```

[1] 0.2222222 0.16666667 0.1111111 0.08333333
[5] 0.19444444 0.3055556 0.08333333 0.19444444
[9] 0.02777778 0.16666667 0.30555556 0.13888889
[13] 0.0000000 0.41666666 0.38888889 0.30555556

```

→

- [0] 0.62500000 0.41666667 0.50000000 0.45833333  
 [5] 0.66666667 0.79166667 0.58333333 0.58333333  
 [9] 0.37500000 0.45833333 0.70833333 0.41666667  
 [13] 0.41666667 0.83333333 1.00000000 0.79166667

Date :

2) Write R program to normalize the variables into 0 to 1 scale Using min-max normalization.

min-max normalizations:

Min-max normalization is a normalization strategy which linearly transforms  $x$  to  $y$ .

$y = (x - \text{min}(x)) / (\text{max}(x) - \text{min}(x))$ , where  $\text{min}$  and  $\text{max}$  are the minimum and maximum values in  $x$ , where  $x$  is the set of observed values of  $x$ . It can be easily seen that when  $x = \text{min}$  then  $y = 0$  and when  $x = \text{max}$ , then  $y = 1$ .

This technique is traditionally used with K-Nearest Neighbours (KNN) classifications problems.

```
> x1 <- iris $ Sepal.Length
```

```
> normalized1 = (x1 - min(x1)) / (max(x1) - min(x1))
```

```
> print(normalized1)
```

```
> x2 <- iris $ Sepal.Width
```

```
> normalized2 = (x2 - min(x2)) / (max(x2) - min(x2))
```

```
> print(normalized2)
```

30 min

-60 min

→  $\begin{bmatrix} 1 \\ 0.06779661 \\ 0.06779661 \\ 0.05084746 \\ 0.08474576 \end{bmatrix}$   
 $\begin{bmatrix} 5 \\ 0.06779661 \\ 0.11864407 \\ 0.06779661 \\ 0.08474576 \end{bmatrix}$   
 $\begin{bmatrix} 9 \\ 0.06779661 \\ 0.08474576 \\ 0.08474576 \\ 0.10169492 \end{bmatrix}$   
 $\begin{bmatrix} 13 \\ 0.06779661 \\ 0.01694915 \\ 0.03389831 \\ 0.08474576 \end{bmatrix}$

→  $\begin{bmatrix} 1 \\ 0.06779661 \\ 0.06779661 \\ 0.05084746 \\ 0.08474576 \end{bmatrix}$   
 $\begin{bmatrix} 5 \\ 0.06779661 \\ 0.11864407 \\ 0.06779661 \\ 0.08474576 \end{bmatrix}$   
 $\begin{bmatrix} 9 \\ 0.06779661 \\ 0.08474576 \\ 0.08474576 \\ 0.10169492 \end{bmatrix}$   
 $\begin{bmatrix} 13 \\ 0.06779661 \\ 0.01694915 \\ 0.03389831 \\ 0.08474576 \end{bmatrix}$

→  $\begin{bmatrix} 1 \\ 0.04166667 \\ 0.04166667 \\ 0.04166667 \\ 0.04166667 \end{bmatrix}$   
 $\begin{bmatrix} 5 \\ 0.12500000 \\ 0.08333333 \\ 0.04166667 \\ 0.04166667 \end{bmatrix}$   
 $\begin{bmatrix} 9 \\ 0.00000000 \\ 0.04166667 \\ 0.04166667 \\ 0.00000000 \end{bmatrix}$   
 $\begin{bmatrix} 13 \\ 0.416667 \\ 0.12500000 \\ 0.12500000 \\ 0.08333333 \end{bmatrix}$

→  $\begin{bmatrix} 1 \\ 0.04166667 \\ 0.04166667 \\ 0.04166667 \\ 0.04166667 \end{bmatrix}$   
 $\begin{bmatrix} 5 \\ 0.12500000 \\ 0.08333333 \\ 0.04166667 \\ 0.04166667 \end{bmatrix}$   
 $\begin{bmatrix} 9 \\ 0.00000000 \\ 0.04166667 \\ 0.04166667 \\ 0.00000000 \end{bmatrix}$   
 $\begin{bmatrix} 13 \\ 0.416667 \\ 0.12500000 \\ 0.12500000 \\ 0.08333333 \end{bmatrix}$

→  $\begin{bmatrix} 1 \\ 0.04166667 \\ 0.04166667 \\ 0.04166667 \\ 0.04166667 \end{bmatrix}$   
 $\begin{bmatrix} 5 \\ 0.12500000 \\ 0.08333333 \\ 0.04166667 \\ 0.04166667 \end{bmatrix}$   
 $\begin{bmatrix} 9 \\ 0.00000000 \\ 0.04166667 \\ 0.04166667 \\ 0.00000000 \end{bmatrix}$   
 $\begin{bmatrix} 13 \\ 0.416667 \\ 0.12500000 \\ 0.12500000 \\ 0.08333333 \end{bmatrix}$

Date :

```
> $x3 <- iris$Petal.Length
```

```
> normalized = (x3 - min(x3)) / (max(x3) - min(x3))
```

```
> print(normalized)
```

```
> $x4 <- iris$Petal.Width
```

```
> normalized = (x4 - min(x4)) / (max(x4) - min(x3))
```

```
> print(normalized)
```

Date :

- ③ Generate histograms for any one variable (sepal length / sepal width / petal length / petal width) and generate scatter plots for every pair of variables showing each species in different color.

Histogram:-

A Histogram represents the frequencies of values of a variable bucketed into ranges. Histogram is similar to bar chart but the differences is to groups the values into continuous ranges.

R creates histogram using hist() function.

Syntax:

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

'v' is a vector containing numeric values used in histogram.

'main' indicates title of the chart.

'col' is used to set border color of the bars.

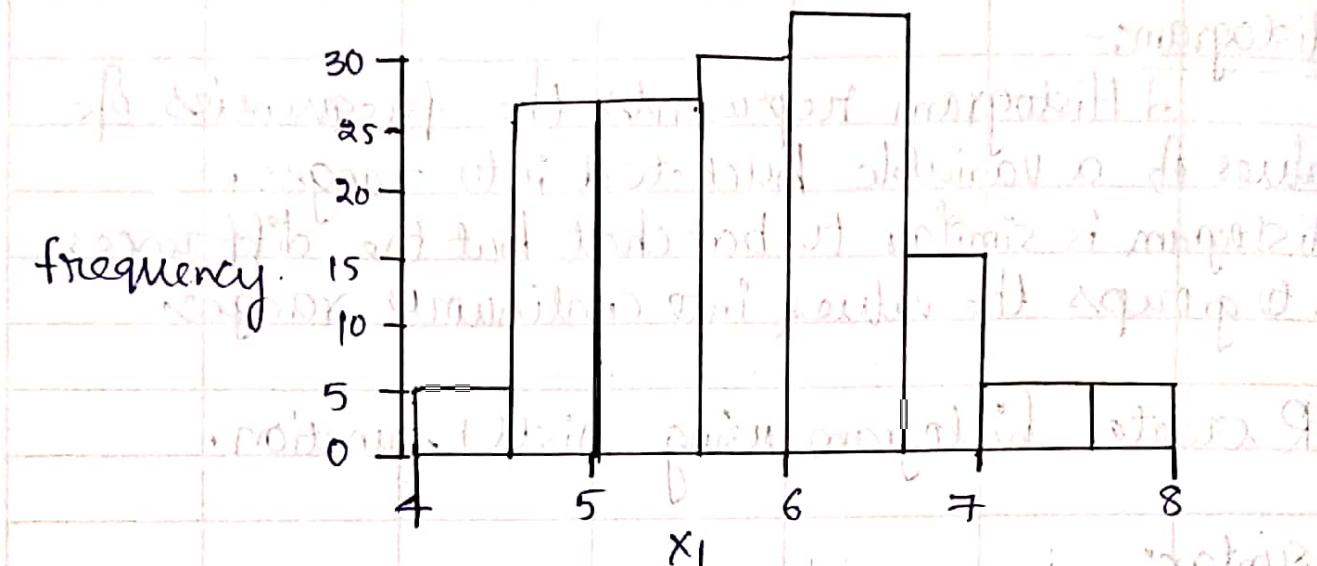
'border' is used to set border color of each bar.

'xlab' is used to give description of x-axis.

'xlim' is used to specify the range of values on the x-axis.

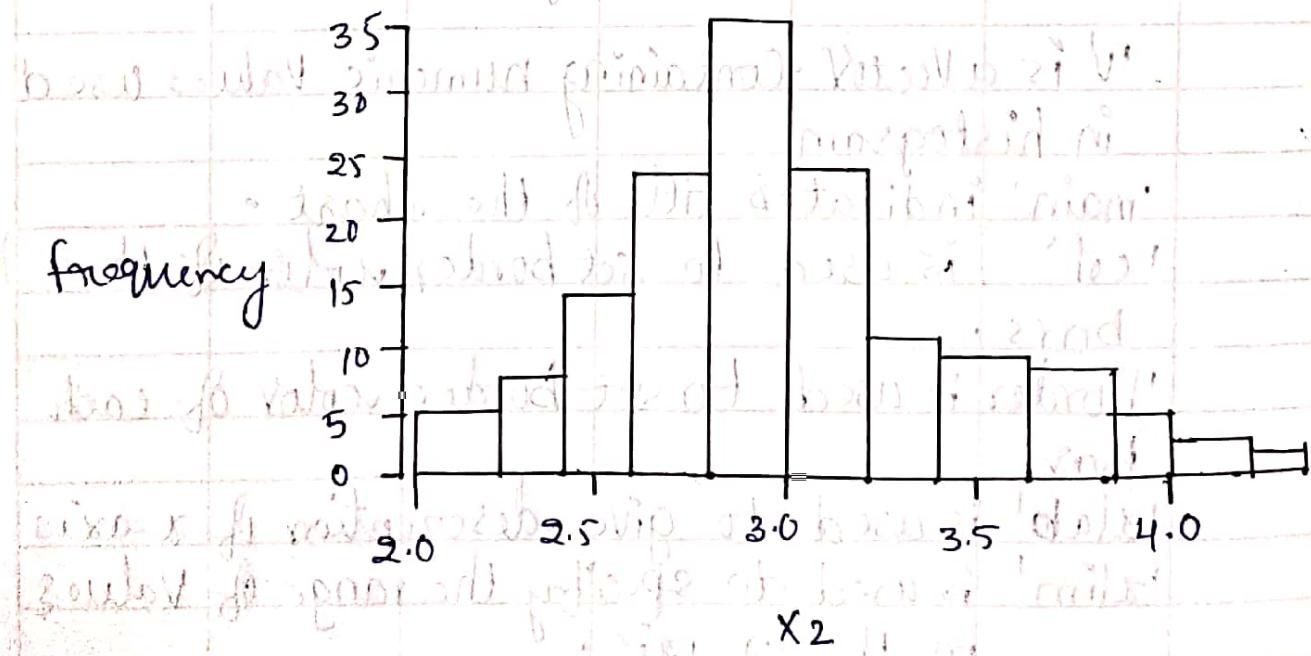
'ylim' is used to specify the range on the y-axis.

Histogram of  $X_1$  (Frequency distribution)



$\rightarrow$

Histogram of  $X_2$  (Frequency distribution)



Date :

'breaks' is used to mention the width of each bar.

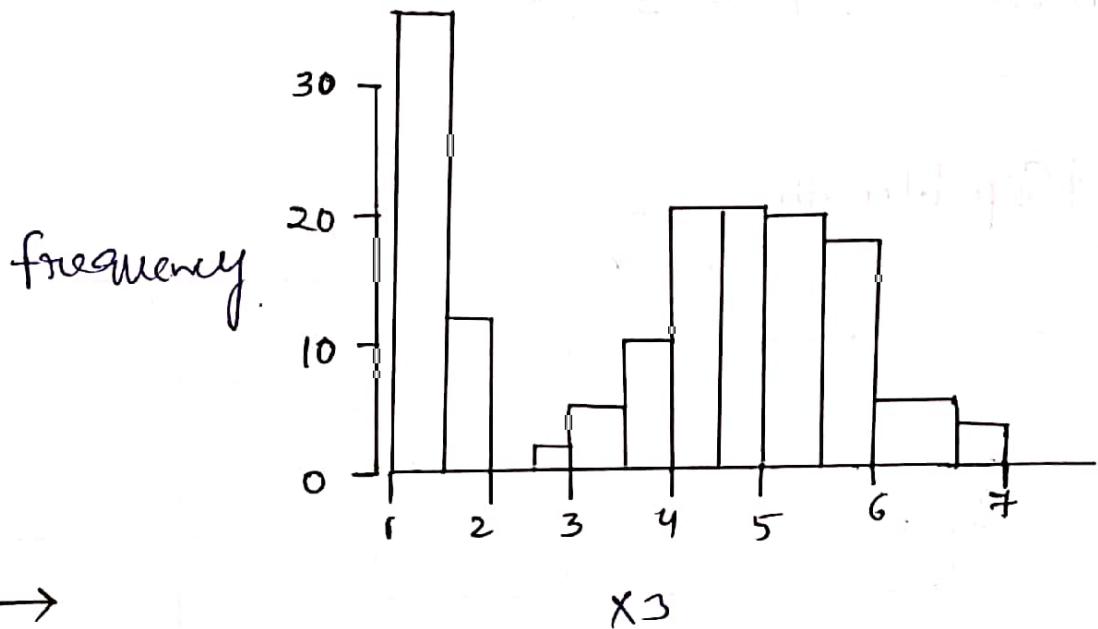
Ex:- `hist(x, main = "acme", xlab = "alpha", ylab = "beta",  
xlim = c(0, 0.4), ylim = c(0, 10), col = "blue",  
breaks = 10, border = "red")`

> `x1 <- iris$Sepal.Length`  
> `hist(x1)`

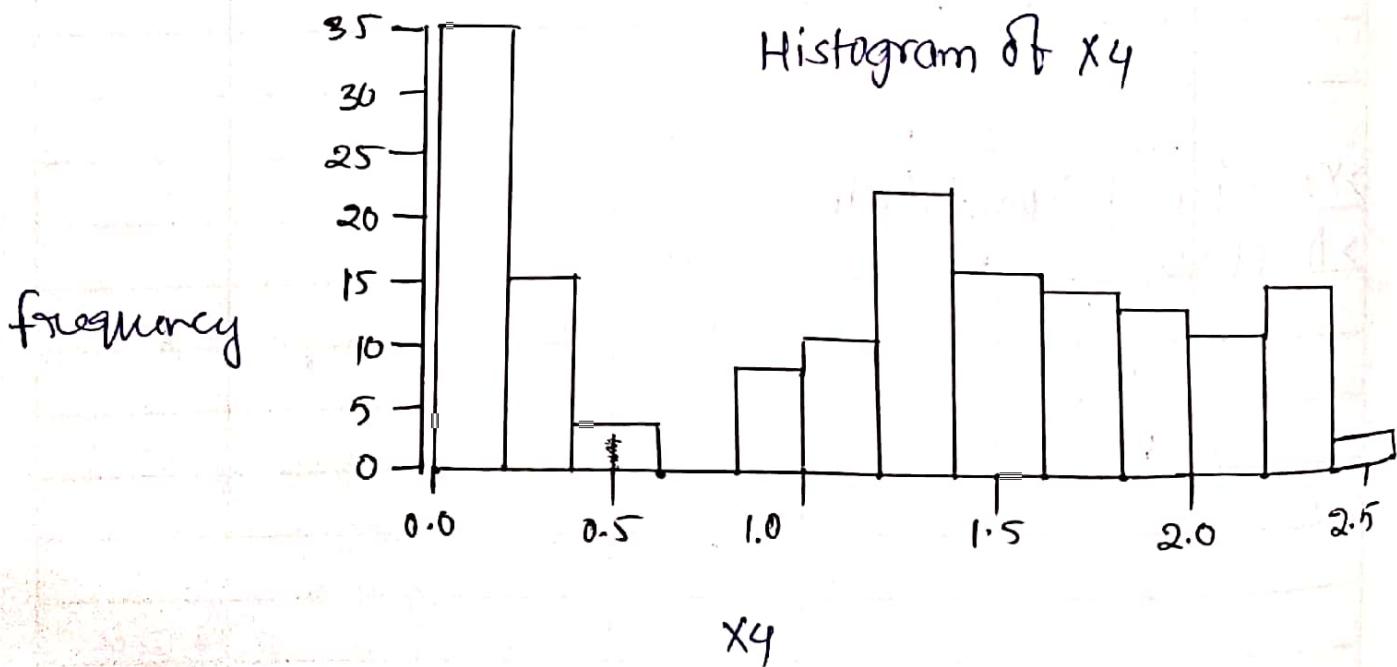
> `x2 <- iris$Sepal.Width`  
> `hist(x2)`



Histogram of  $x_3$



$x_3$

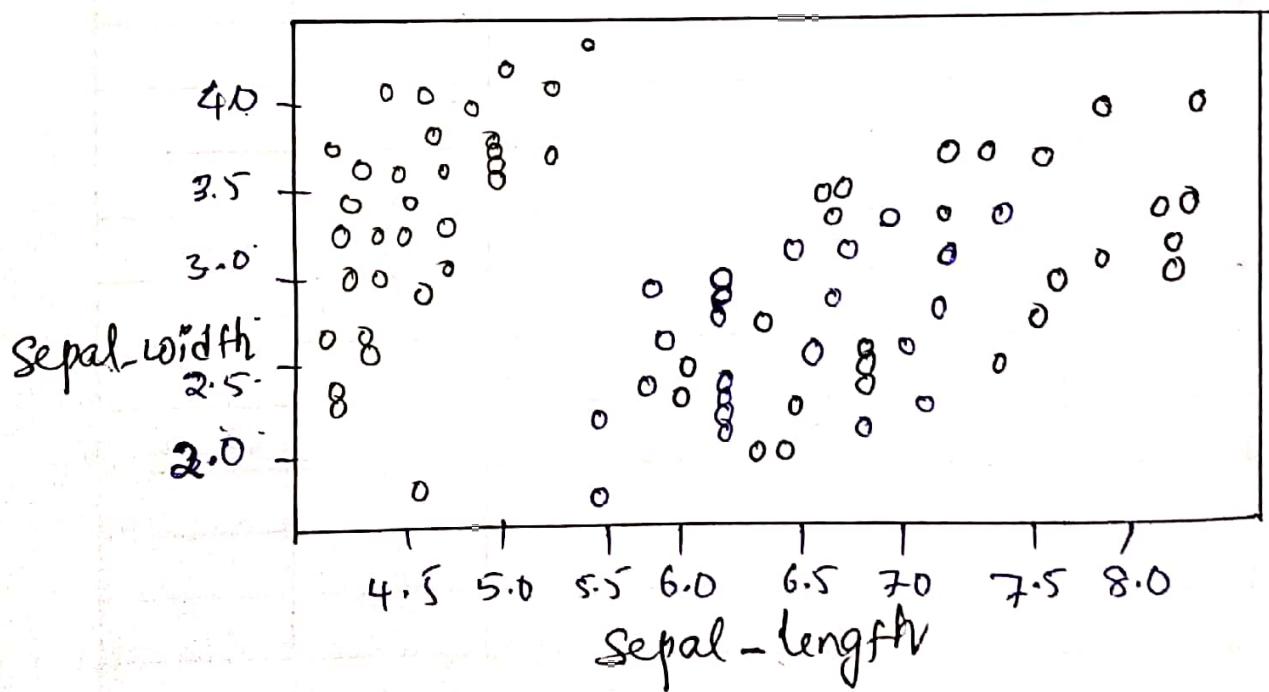


>  $x_3 \leftarrow \text{iris} \$ \text{Petal.Length}$   
> hist( $x_3$ )

>  $x_4 \leftarrow \text{iris} \$ \text{Petal.Width}$   
> hist( $x_4$ )



Sepal-length Vs Sepal-width



Date :

Scatter plots :-

Scatter plots show many points plotted in the cartesian plane. Each point represents the values of two variables. One variable is chosen in the horizontal axis and another in the vertical axis.

Syntax :-

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

'x' is the data set whose values are the horizontal coordinates,

'y' is the data set whose values are vertical coordinates

'main' is title of the graph

'xlab' is the label in the horizontal axis

'ylab' is the label in the vertical axis

'xlim' is the limits of the values of x

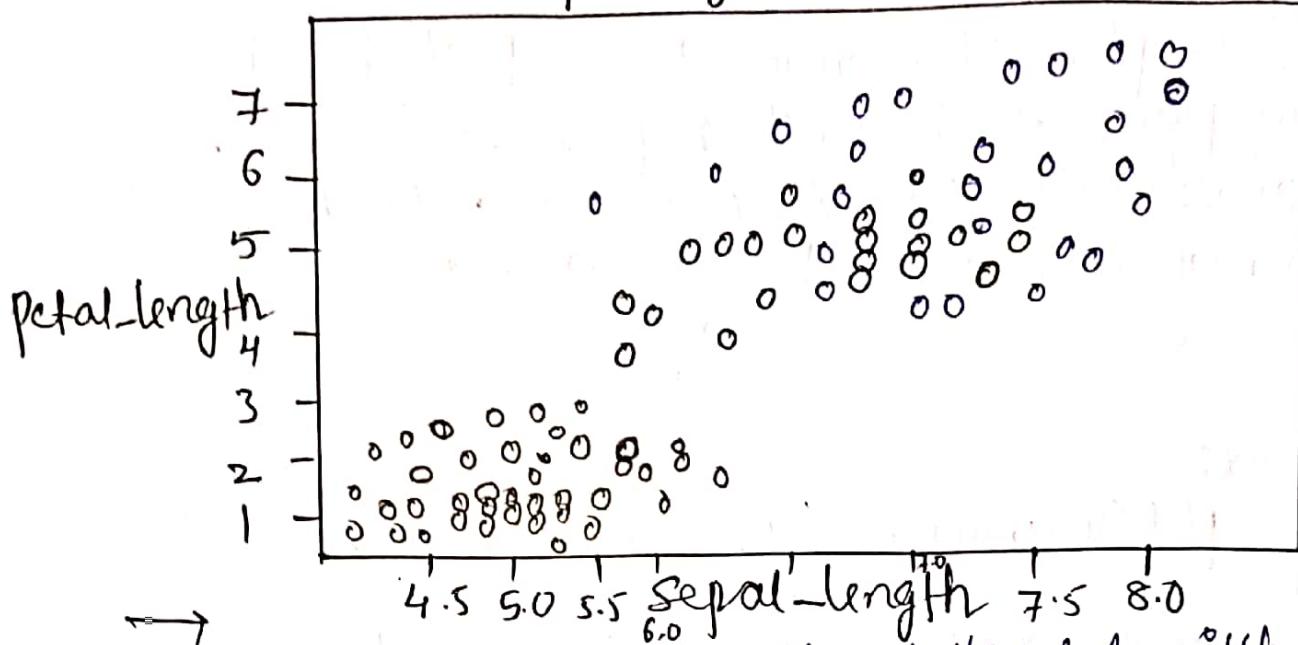
'ylim' is the limits of the values of y

'axes' indicates whether both axes should be drawn on the plot.

```
> plot(x=x1, y=x2, col=c("red", "blue", "green")) [iris$Species],
 xlab="Sepal.Length", ylab="Sepal.Width",
 main="Sepal.Length vs Sepal.Width")
```

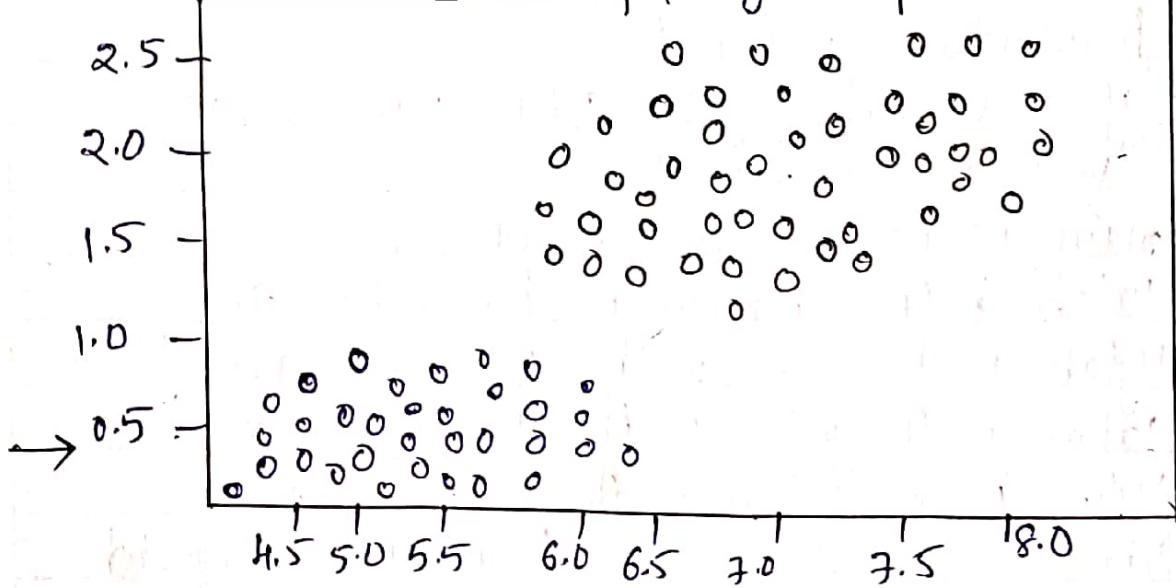
→

### sepal length Vs petal-length



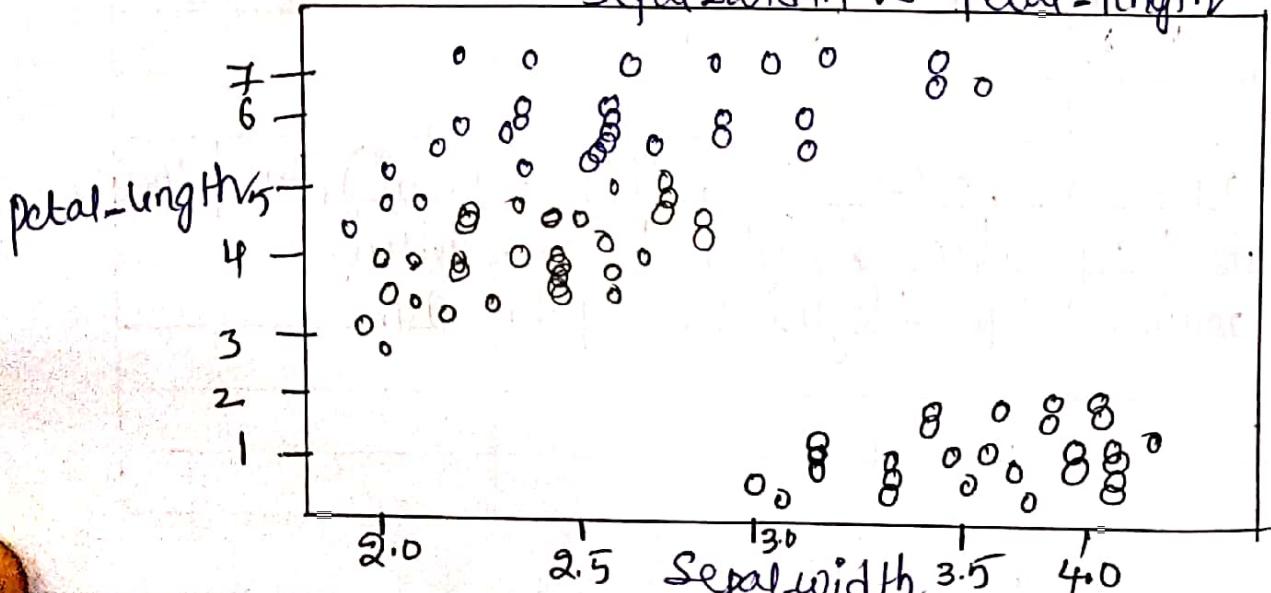
→

### sepal length Vs petal-width



→

### Sepal width Vs Petal-length

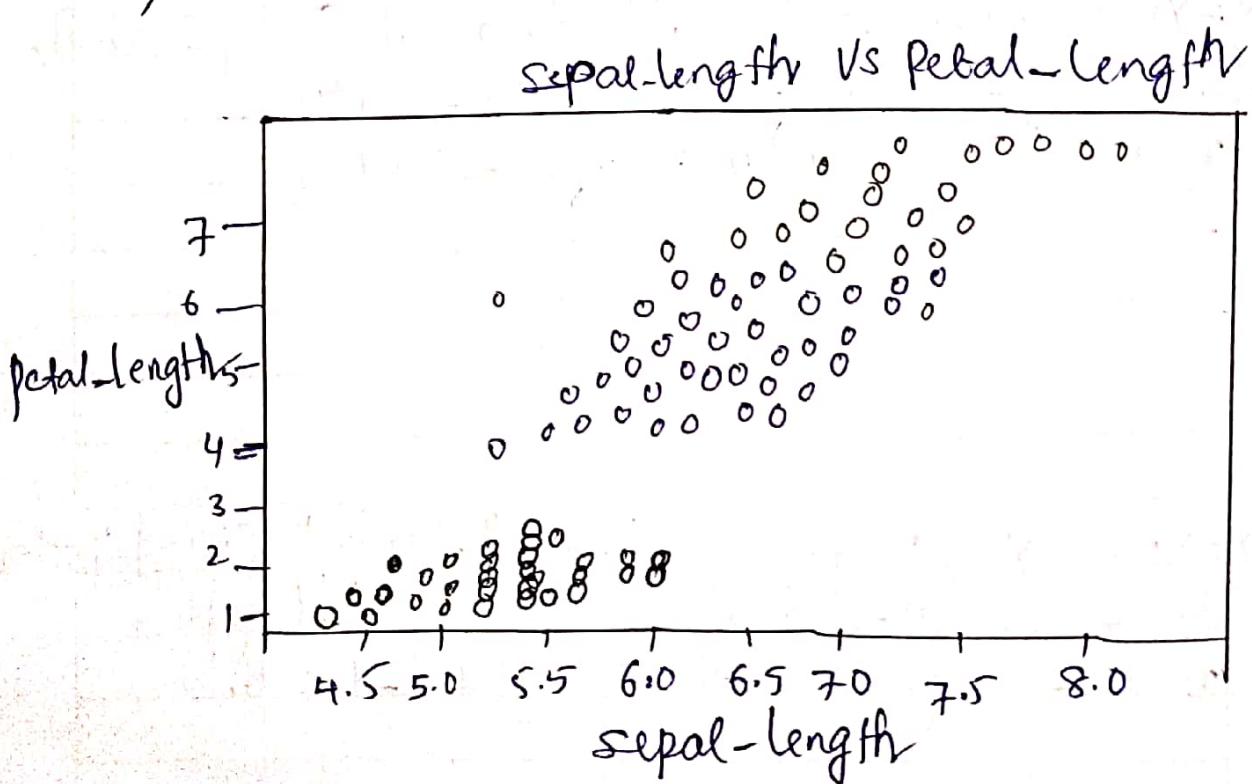
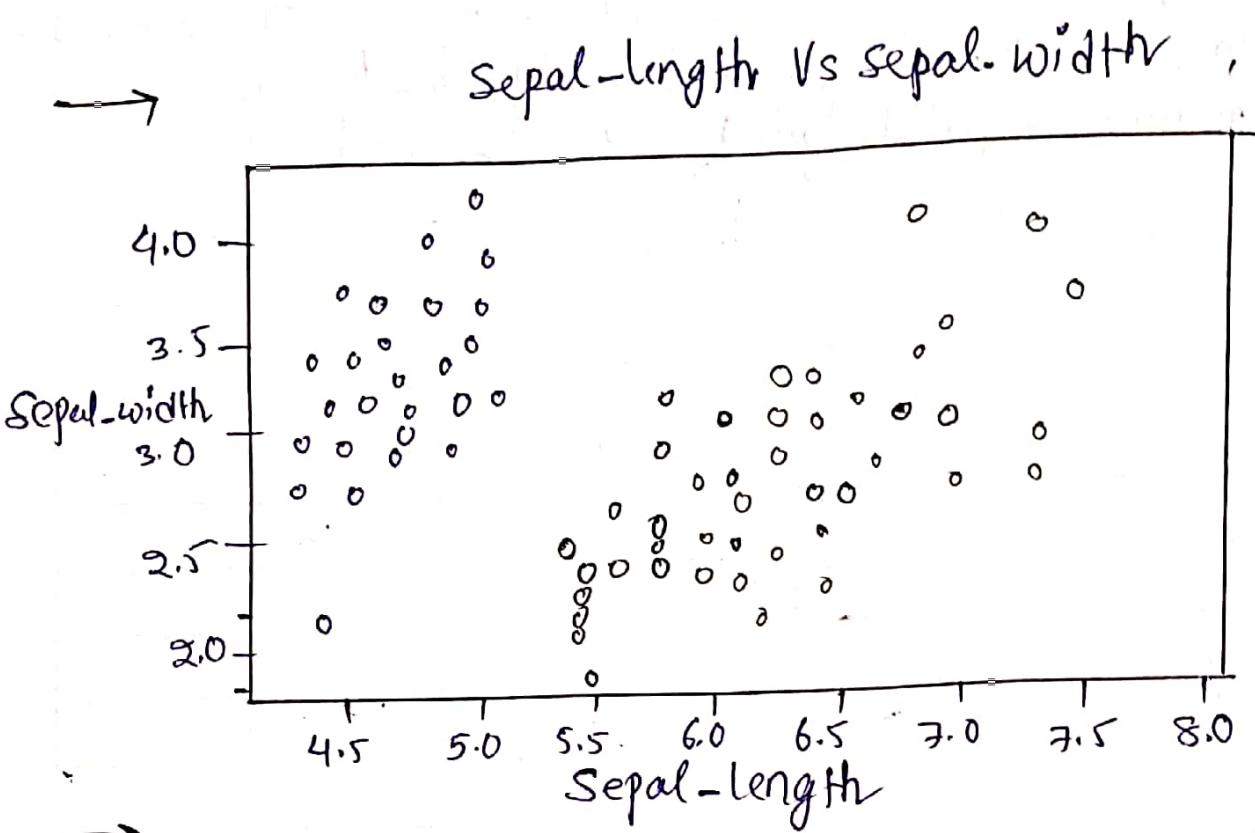


Date :

```
> plot(x=x1, y=x3, col=c("red", "blue", "green")) [iris$Species],
 xlab = "Sepal-length", ylab = "petal-length",
 main = "Sepal-length VS petal-length")
```

```
> plot(x=x1, y=x4, col = ("red", "blue", "green")) [iris$Species],
 xlab = "Sepal-length", ylab = "petal-width",
 main = "Sepal-length VS petal-width")
```

```
> plot(x=x2, y=x3, col = c("red", "blue", "green")) [iris$Species],
 xlab = "Sepal-width", ylab = "petal-length",
 main = "Sepal-width VS petal-length")
```



Date :

```
> plot(x=x2, y=x4, col=c("red", "blue", "green")) [iris$Species]
 xlab = "Sepal_width", ylab = "petal_width",
 main = "sepal_width vs petal_width")
```

```
> plot(x=x3, y=x4, col=c("red", "blue", "green")) [iris$Species],
 xlab = "petal_length", ylab = "petal_width",
 main = "petal_length vs petal_width")
```

Date :

- ④ Generate box plots for each of the numeric attributes. Identify the attribute with the highest variance.

Boxplots :-

Boxplots are a measure of how well distributed is the data in a dataset. It divides the dataset 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.

Syntax :-

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

'x' is a vector or a formula.

'data' is the data frame.

'notch' is a logical value. Set as TRUE to draw a notch.

'varwidth' is a logical value. Set as true to draw width of the box proportionate to the sample size.

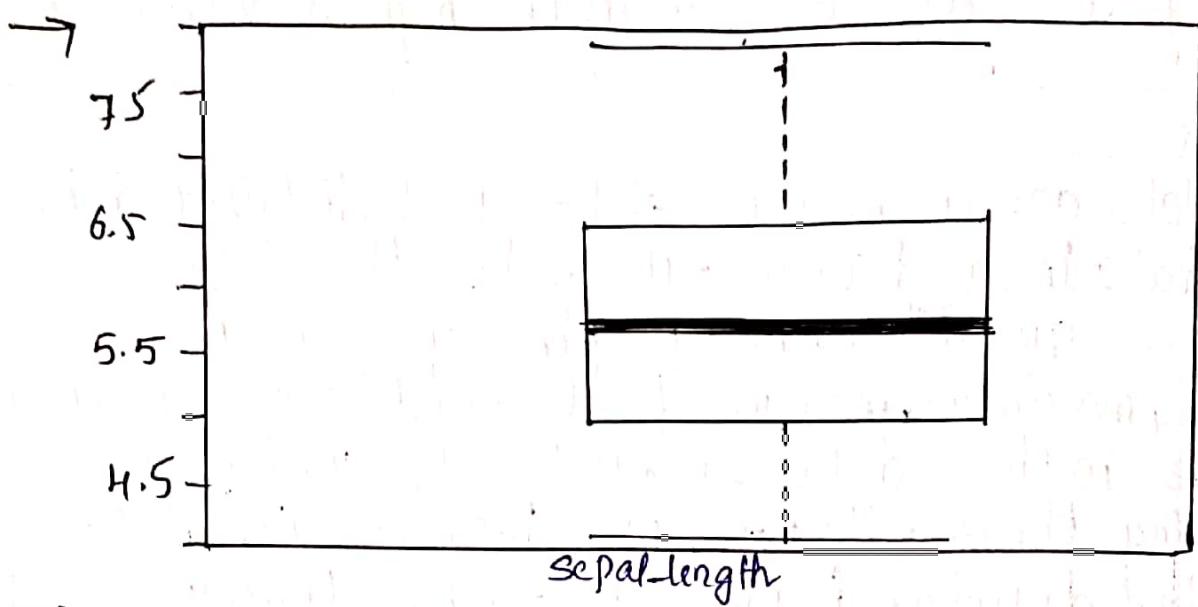
'names' are the group labels which will be printed under each boxplot.

'main' is used to give a title to the graph.

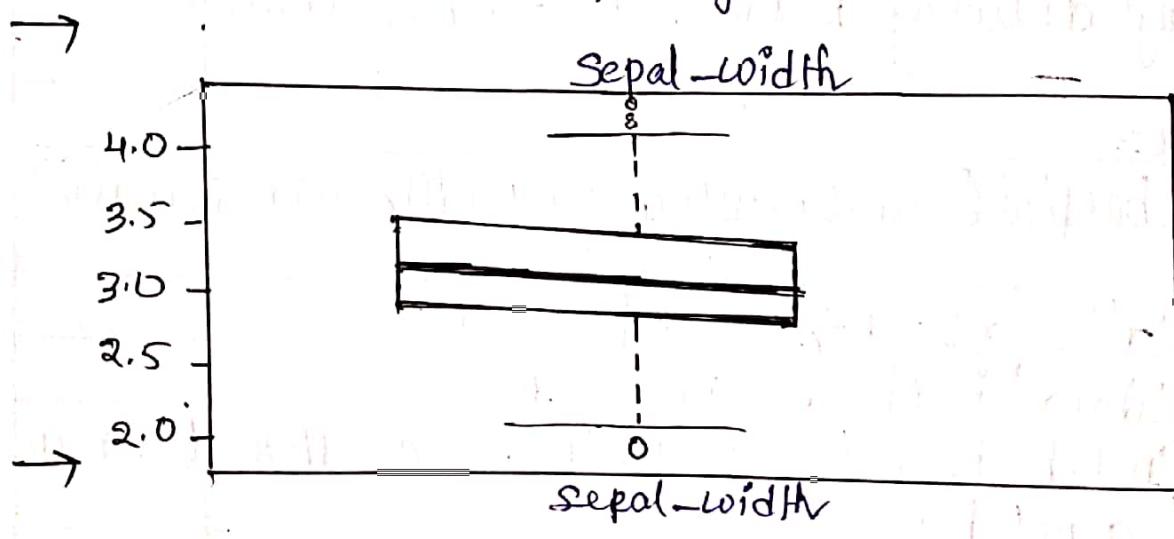
```
> x1 <- iris$Sepal.Length
```

```
> boxplot(x1, xlab =
```

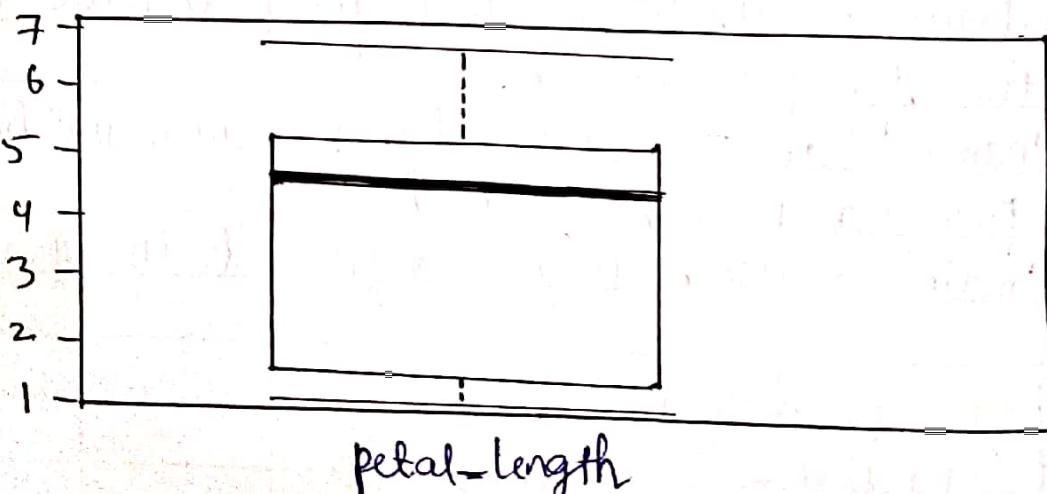
sepal-length



sepal-length



Petal-length



petal-length

Date : \_\_\_\_\_

```
x2 <- iris$Sepal.Width
x3 <- iris$Sepal.Length
x4 <- iris$Petal.Length

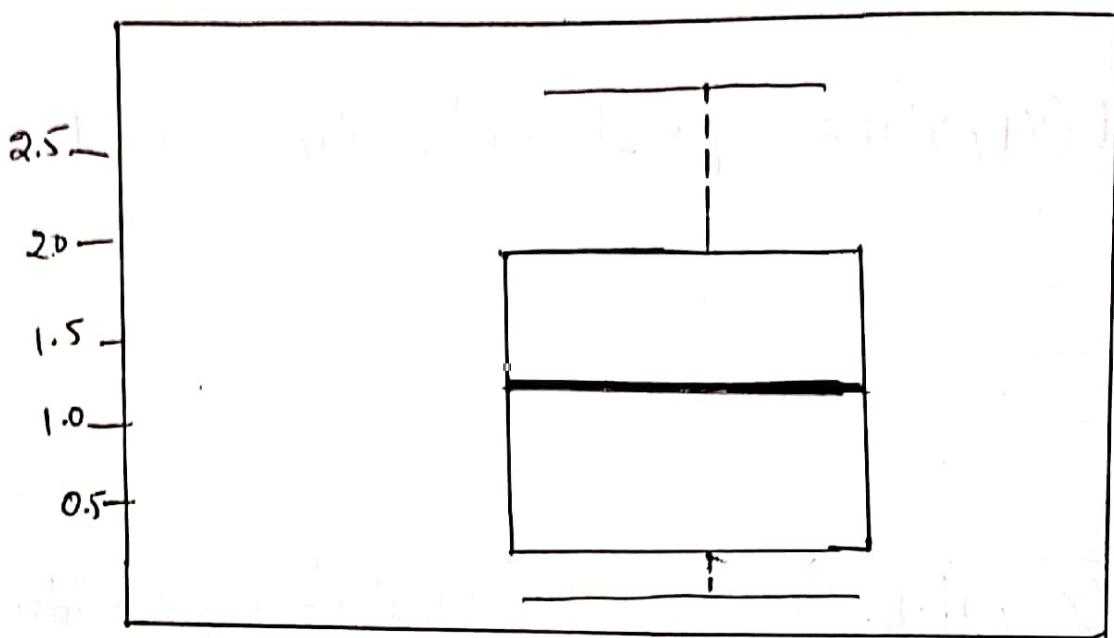
> boxplot(x1, xlab = "sepal.length", main = "sepal.length")

> boxplot(x2, xlab = "sepal.width", main = "sepal.width")

> boxplot(x3, xlab = "petal.length", main = "petal.length")
```



Petal-Width



petal-width



[1] 3.0116278

Date : \_\_\_\_\_

```
> boxplot(x4, -1) lab = "Petal-width", main = "Petal-width"

> Var1 = Var(x1)
> Var2 = Var(x2)
> Var3 = Var(x3)
> Var4 = Var(x4)

> max1 = max(c(Var1, Var2, Var3, Var4))

> print(max1)
```

→ [1] "red" "green" "yellow"

→ [1] "1" "5" "6" "11" "8" "11" "a"

→ [1] 5

→ [1] "2" "3" "5" "aa" "bb" "cc" "dd" "ee"

→ [1] 2 5 9 15 10

Date :

⑤ Study of homogeneous and heterogeneous iris structures such as Vector, matrix, array, list, data frame in R

### Vector :-

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

```
> apple <- c("red", "green", "yellow")
> print(apple)
```

```
> num <- c(15, 6, 1, 8, 11, "a")
> print(num)
```

```
> length(c("aa", "bb", "cc", "dd", "ee"))
```

```
> n = c(2, 3, 5)
> s = c("aa", "bb", "cc", "dd", "ee")
> c(n, s)
```

```
> a = c(1, 3, 5, 7)
> b = c(1, 2, 4, 8, 9)
```

```
> a + b
```

$\rightarrow$  [1,1] 2 [1,2] 4 [1,3] 3  
[2,1] 1 [2,2] 5 [2,3] 7

$\rightarrow$  [1,1] 2 [1,2] 3 [1,3] 5  
[2,1] 4 [2,2] 1 [2,3] 7

$\rightarrow$  [1] 7

$\rightarrow$  [1] 4 1 7

$\rightarrow$  [1] 5 7

$\rightarrow$  [1,1] 2 [1,2] 5  
[2,1] 4 [2,2] 7

$\rightarrow$  [1,1] 2 [1,2] 4  
[2,1] 3 [2,2] 1  
[3,1] 5 [3,2] 7

## Matrices :-

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

#matrix

$A = \text{matrix}(C$

+ c(2, 4, 3, 1, 5, 7), # data elements

+ nrow = 2, # number of rows

+ ncol = 3, # number of columns

+ byrow = TRUE) # fill matrix by rows

> A = matrix(c(2, 4, 3, 1, 5, 7), nrow = 2, ncol = 3, byrow = TRUE)

> A = matrix(c(2, 4, 3, 1, 5, 7), nrow = 2, ncol = 3)

> A[2, 3]

> A[2, ]

> A[, 3]

> A[, c(1, 3)]

> t(A)

→

, , 1

|       | [1]      | [2]      | [3]      |
|-------|----------|----------|----------|
| [1, ] | "green"  | "yellow" | "green"  |
| [2, ] | "yellow" | "green"  | "yellow" |
| [3, ] | "green"  | "yellow" | "green"  |

, , 2

|       | [1]      | [2]      | [3]      |
|-------|----------|----------|----------|
| [1, ] | "yellow" | "green"  | "yellow" |
| [2, ] | "green"  | "yellow" | "green"  |
| [3, ] | "yellow" | "green"  | "yellow" |

→

[1]

|   | gender | height | weight | Age |
|---|--------|--------|--------|-----|
| 1 | male   | 152.0  | 81     | 42  |
| 2 | male   | 171.5  | 93     | 38  |
| 3 | Female | 165.0  | 78     | 26  |

## Arrays:

While matrices were confined to two dimensions, arrays can be of any number of dimensions. The array function takes a dim attribute which creates the required number of dimensions of

```
a <- array(c("green", "yellow"), dim = c(3, 3, 2))
```

## Data frames:

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

Data frames are created using the data.frame() function.

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

→

[C1])

[C1] 2 5

[C2])

[C1] "red" "green" "yellow"

[C3])

[C1] 21.3

[C4])

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

Lists :-

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

eglist  $\leftarrow$  list(c(2,5), c("red", "green", "yellow"), 21.3, BMI), mat  
print(eglist)



[1] -0.89767388 -1.13920048 -1.38072709 -1.50149039  
[5] -1.01843718 -0.93538397 -1.50149039 -1.01843718  
[9] -1.74301699 -1.13920048 -0.93538397 -1.25996379  
[13] -1.25996379 -1.86378030 -0.05233076 -0.17301407

Date :

- ⑥ Write R program using "apply" group of functions to create and apply normalization function on each of the numeric variables/columns of iris dataset to transform them into a value around 0 with z-score normalization

Z-score:

Z-score (aka, a standard score) indicates how many standard deviations an element is from the mean. A z-score can be calculated from the following formula :

$$z = (x - \mu) / \sigma$$

where  $z$  is the z-score,

$x$  is the value of the element,

$\mu$  is the population mean

$\sigma$  is the standard deviation

```
> x1 <- iris $Sepal.Length
> zscr1 <- (x1 - mean(x1)) / sd(x1)
> zscr1
```

$\rightarrow$  [1] 1.01560199 -0.13153881 0.32731751 0.09768935  
[5] 1.24503015 1.93331463 0.78617383 0.78617383  
[9] -0.36096697 0.09788935 1.47445831 0.78617383

$\rightarrow$  [1] -1.33575163 -1.33575163 -1.39239929 -1.27910398  
[5] -1.33575163 -1.16580868 -1.33575163 -1.27910398  
[9] -1.33575163 -1.27910398 -1.27910398 -1.22245633

$\rightarrow$  [1] -1.3110521482 -1.3110521482 -1.3110521482 -1.3110521482  
[5] -1.3110521482 -1.9486667950 -1.1798594716 -1.3110521482  
[9] -1.3110521482 -1.4922448248 -1.3110521482 -1.3110521482

Date :

```
> X2 <- iris $Sepal.Width
> ZSCX2 <- (X2 - mean(X2)) / sd(X2)
> ZSCY2
```

```
> X3 <- iris $Petal.Length
> ZSCX3 <- (X3 - mean(X3)) / sd(X3)
> ZSCY3
```

```
> X4 <- iris $Petal.Width
> ZSCX4 <- (X4 - mean(X4)) / sd(X4)
> ZSCY4
```

Date :

a) Use R to apply linear regression to predict evaporation coefficient in terms of air velocity using the data given below.

|                                                  |                                                              |
|--------------------------------------------------|--------------------------------------------------------------|
| Air Velocity (cm/sec)                            | 20, 60, 100, 140, 180, 220, 260, 300<br>340, 380             |
| Evaporation Coefficient (cm <sup>3</sup> mm/sec) | 0.18, 0.37, 0.35, 0.78, 0.56, 0.75<br>1.18, 1.36, 1.17, 1.65 |

### Regression:

Regression analysis is a very widely used statistical tool to establish a relationship model between two variables. One of these variable is called predictor variable whose value is gathered through experiments. The other variable is called response variable whose value is derived from the predictor variable.

In linear Regression these two variables are related through an equation, where exponent of both these variables is 1. Mathematically a linear relationship represents a straight line when plotted as a graph. A non-linear relationship where the exponent of any variable is not equal to 1 creates a curve.

→

air velocity

evaporation coefficient

|    | air velocity | evaporation coefficient |
|----|--------------|-------------------------|
| 1  | 20           | 0.18                    |
| 2  | 60           | 0.37                    |
| 3  | 100          | 0.35                    |
| 4  | 140          | 0.78                    |
| 5  | 180          | 0.56                    |
| 6  | 220          | 0.75                    |
| 7  | 260          | 1.18                    |
| 8  | 300          | 1.36                    |
| 9  | 340          | 1.17                    |
| 10 | 380          | 1.65                    |

→

Call:

lm (formula = AV \$ air velocity ~ AV \$ evaporation coefficient)

Coefficients:

(Intercept)

2.564

AV \$ evaporation coefficient

236.450

Linear regression,  $y = ax + b$

$y$  is the response variable  
 $x$  is the predictor variable  
 $a$  and  $b$  are constants which  
are called the coefficients.

### lm() function

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

#### Syntax:

$lm(\text{formula}, \text{data})$

↳ vector on which formula will be applied.

↳ symbol presenting the relation between  $x$  and  $y$

```
> AV <- data.frame(airVelocity = c(20, 60, 100, 140, 180, 220, 260,
 300, 340, 380), evaporationCoefficient = c(0.18, 0.3),
 0.35, 0.78, 0.56, 0.35, 1.18, 1.36, 1.17, 1.65))
> print(AV)
```

```
> model <- lm(AV ~ airVelocity + evaporationCoefficient)
> print(model)
```



Call:

lm(formula = AV\$airVelocity ~ AV\$evaporationCoefficient)

Residuals:

|  | Min    | 1Q     | Median | 3Q    | Max   |
|--|--------|--------|--------|-------|-------|
|  | -46.99 | -24.88 | -17.14 | 33.74 | 60.79 |

Coefficients:

|                            | Estimated Std. Error | t value | Pr(> t )           |
|----------------------------|----------------------|---------|--------------------|
| (Intercept)                | 2.564                | 25.804  | 0.099 0.923        |
| AV\$evaporationCoefficient | 23.6.456             | 27.035  | 8.746 2.29e-05 *** |

Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 '''

Residual standard error: 39.53 on 8 degrees of freedom

Multiple R-squared: 0.9053, Adjusted R-squared: 0.8935

F-statistic: 76.49 on 1 and 8 DF, p-value: 2.28e-05

→ [1] 0.9514814

→ [1] 0.9514814

b) Analyze the significance of residual standard-error value, R-squared value, F-statistic. Find the correlation coefficient for this data and analyze the significance of the correlation value.

Residual Standard error is measure of the quality of a linear regression fit ... the Residual Standard Error is the average amount of that the response (dist) will deviate from the true regression line.

R-squared statistic provides a measure of how well the model is fitting the actual data.

F-statistic is a good indicator of whether there is a relationship between our predictor and the response variable.

Is it true that we can have a F value indicating a strong relationship that is NON LINEAR so that our RSE is high and our R squared is low

> Summary(model)

> cor(AV\$airVelocity, AV\$evaporationCoefficient)

> y <- cor(AV\$airVelocity, AV\$evaporationCoefficient)

> y

→ [1] 2.995732 4.094345 4.605170 4.941642 5.192957  
5.393628 5.560682 5.703782 5.828946  
5.940171

→ call:  
 $\text{lm}(\text{formula} = \text{x} \sim \text{AV\$ evaporation coefficient})$

Coefficients:

| (Intercept) | AV \\$evaporation coefficient |
|-------------|-------------------------------|
| 3.678       | 1.614                         |

(C) Perform a log transformation on the 'Air Velocity' column, perform linear regression again and analyze all the relevant values.

>  $x \leftarrow \log(\text{AV}[\text{air velocity}])$

>  $m \leftarrow \text{lm}(x \sim \text{AV}[\text{evaporation Coefficient}])$

>  $m$