

Practical File

Lab Name ... Data Analytics Lab..... Lab Code.....KIT-651.....



Estd. 2000

Name.....Harshul Gupta.....

Adm.No...2018BCI1055..... Univ. Roll No...1803211023.....

CourseB.Tech.....Branch.....CEIT.....

Sem.....6th..... Section.....A.....

ABES
Engineering College
(College Code-032)

NAAC Accredited, NBA Accredited Branches (CSE, ECE, EN & IT)

19th Km. Stone, NH-09, Ghaziabad - 201009 (UP), India
Phone : 0120-7135112, 9999889341 Fax : 0120-7135115 Website : www.abes.ac.in, Email: info@abes.ac.in





INDEX

Exp. No	Name of the Experiment	Date of Practical	Assessment by Faculty Member					
		Schedule Actual	Pre-Lab Writing work * MM:05	Implementation/ Active Participation MM:10	Graph, Results/ Output, Calc. MM: 05	In Time submission of Lab reports/ Viva Voce MM:05	Total MM:25	Signature with Date
1	Understanding of statistics on dataset	4/13/2021						
2	Handling data files in R	20/04/2021						
3	Handling Multi-dimensional Array	22/05/2021						
4	Perform mean, median, mode and standard deviation	13/04/2021						
5	Perform data processing operations	27/04/2021						
6	Perform Dimensionality Reduction	4/5/2021						
7	Perform Simple Linear Regression	24/05/2021						

Teacher's Remarks (if any):

Average Marks:

Name & Sign. Of Faculty members(s) with date

Name & Sign. Of Lab in charge

HOD

- Note:
1. *Pre-Lab writing work should include problem statement, objective , Algorithm (if applicable)and Methodology.
 2. Faculty members will check pre Lab writing work & Lab work readings/output in each class and sign with date.
 3. Please use pen with blue colour Ink only.

Practical Name: To get the input from user & perform numerical op's (Max, Min, Avg, Sum, SQRT, Round) using in R.

Objective → Understanding of statistics on dataset.

Algorithm →

- ① Take the array of number from user
- ② find minimum, maximum, average, sum and square root of each number.
- ③ Use min(), max(), avg(), sum(), sqrt(), round() functions
- ④ Display output.

Code →

```
x ← strsplit(readline(prompt = "Enter the list of numbers:"),
              ",")  

x ← as.double(x[[1]])  

print(x)  

print(c("Max:", max(x)))  

print(c("Min:", min(x)))  

print(c("Sum:", sum(x)))  

print(c("Average:", mean(x)))  

print(c("Round:", round(x)))  

print(c("SquareRoot:", sqrt(x)))
```

Practical Name: To get input from user & perform numerical operations Practical No. 1

Output →

Enter the list of numbers: 5 4.8 55 12.6

5.0 4.8 55.0 12.6

Max: 55

Min: 4.8

Sum: 77.4

Average: 19.35

Round: 5 5 55 13

SquareRoot: 2.23 2.19 7.41 3.54

Objective → Handling the data files in R programming.

Algorithm →

- ① Explore the UCI machine learning repository.
- ② Download the dataset in .CSV, .XLS and .TXT file format.
- ③ Use read() functions to import the dataset in R.
- ④ Save the file in data frame and print the header using header() function.

Code →

```
data1 <- read.csv("E:\\sample.csv")
data2 <- read.table("E:\\sample.csv", sep = ",",
                     dec = ".")
data3 <- read.xlsx("E:\\sample.xlsx",
                   sheetName = "Sheet1")
head(data1)
head(data2)
head(data3)
v <- c(1, 2, 3, 4, 5)
write.csv(v, file = "E:\\test.csv")
write.table(v, file = "E:\\test.txt", append = F,
            sep = " ", dec = ".", row.names = T,
            col.names = T)
write.xlsx(v, file = "E:\\test.xlsx",
            sheetName = "Sheet1", col.names = T,
```

Roll No: 1803211023 Date 20/04/2021 Page No.

Practical Name: To perform data import/export Practical No. 2

From names = T, append = f)

Output ->

X

1 1

2 2

3 3

4 4

5 5

X

1 1

2 2

3 3

4 4

5 5

X

1 1

2 2

3 3

4 4

5 5

Name: To get input matrix from user and perform..... Practical No. 3

Matrix addition, subtraction, multiplication,
inverse, transpose and division

Objective → Handling of multidimensional array

Algorithm →

- ① Take the input from the user in 2-D array of size (3x3).
- ② Perform addition, subtraction, multiplication,
inverse, transpose and division operation in R.
- ③ Display the output

Code →

```
print("Enter the elements of first matrix: ")
```

```
d1 ← scan(nmax=9)
```

```
print("Enter the elements of second matrix: ")
```

```
d2 ← scan(nmax=9)
```

```
rowname = c("row1", "row2", "row3")
```

```
colname = c("col1", "col2", "col3")
```

```
M1 ← matrix(d1, nrow=3, byrow=TRUE,
```

```
dimnames = list(rowname, colname))
```

```
print(M1)
```

```
M2 ← matrix(d2, nrow=3, byrow=TRUE,
```

```
dimnames = list(rowname, colname))
```

```
print(M2)
```

```

print("Addition of two matrix:")
print(M1+M2)
print("Subtraction of two matrix:")
print(M1-M2)
print("Multiplication of two matrix:")
print(M1 %*% M2)
print("Inverse of matrix:")
print(solve(M1))
print("Transpose of matrix:")
print(t(M1))
print("Division of two matrix:")
print(M1/M2)

```

Output ->

Enter the elements of first matrix:

20	30	23
14	21	28
18	19	17

Enter the elements of second matrix:

15	12	28
19	25	30
21	23	11

Addition of two matrix:

row1	col1	col2	col3
row1	35	42	51
row2	33	46	58
row3	39	42	28

Name : Practical No. 3

Subtraction of two matrix:

	Col1	Col2	Col3
Row1	5	18	-5
Row2	-5	-4	-2
Row3	-3	-4	6

Multiplication of two matrix:

	Col1	Col2	Col3
Row1	1353	1519	1713
Row2	1197	1337	1330
Row3	988	1087	1261

Inverse of matrix:

	Col1	Col2	Col3
Row1	-0.09	-0.03	1.87
Row2	0.13	-0.03	-1.25
Row3	-0.05	0.08	-3.73

Transpose of matrix:

	Row1	Row2	Row3
Col1	20	14	18
Col2	30	21	19
Col3	23	28	17

Division of matrix:

	Col1	Col2	Col3
Row1	1.33	2.50	0.82
Row2	0.73	0.84	0.93
Row3	0.85	0.82	1.54

Practical Name: To perform statistical operations (Mean, Median, Practical No. 4

Mode and Standard Deviation) using R

Objective → To find out the statistics of the dataset.

Algorithm →

- ① Apply Mean, Median and Standard Deviation formulas to find statistics of given dataset.
- ② Apply mean(), median(), mode() and sd() function.
- ③ Find all the statistics parameter on data frame calculated for &
- ④ Display the output.

Code →

```
data1 ← read.csv("E:\Sample.csv")
head(data1)
mode ← function(v) {
  unq ← unique(v)
  unq[which.max(tabulate(match(v, unq)))]
}

print(c("Mean:", mean(data1[["x"]])))
print(c("Median:", median(data1[["x"]])))
print(c("Mode:", mode(data1[["x"]])))
print(c("Standard deviation", sd(data1[["x"]])))
```

Name: To perform Statistical operations Practical No. 4

Output →

X

1 23

2 56

3 78

4 23

5 56

6 79

Mean: 53.875

Median: 56

Mode: 23

Standard deviation: 28.3167

Practical Name: To perform data pre-processing operations Practical No. 5

- (i) Handling missing data
- (ii) Min-Max normalization

Objective \Rightarrow Data Cleaning processing

Algorithm \rightarrow

- ① Find missing values in data using `is.na()` function.
- ② Fill missing values using mean/median/mode
- ③ Normalize the data using min-max normalization.

Code \rightarrow

```
v <- c(3, NA, 5, 6, 8, NA, 9, 12, 10, NA, 21, 25)
```

```
a <- v
```

```
b <- v
```

```
print(is.na(v))
```

```
a[is.na(a)] <- mean(a, na.rm=T)
```

```
b[is.na(b)] <- median(b, na.rm=T)
```

```
print("Using mean")
```

```
print(a)
```

```
print("Using median")
```

```
print(b)
```

```
d <- a
```

```
for(i in 1:length(a))
```

```
d[i] = (d[i] - min(d)) / (max(d) - min(d))
```

```
print("Normalised data")
```

```
print(d)
```

Output →

FALSE TRUE FALSE FALSE FALSE TRUE FALSE FALSE FALSE
TRUE FALSE FALSE

Using mean

3 11 5 6 8 11 9 12 10 11 21 25

Using median

3 9 5 6 8 9 9 12 10 9 21 25

Normalised data

0.00 0.36 0.09 0.13 0.22 0.36 0.27

0.40 0.31 0.36 0.81 1.00

Name: To perform dimensionality reduction
operation using PCA Practical No. 6

Objective → To get the precise and important data.

Algorithm →

- ① Pre-process the dataset
- ② Built boxplot function for EDA on variance.
- ③ Standardize the data by using scale and apply "percomp" function.
- ④ Choose the principal components with highest variance
- ⑤ Visualize Data in new reduced dimension

Code →

```
data1 ← read.csv("E:\\sample.csv", header=F)
x ← as.matrix(data1)
print("Original Data")
print(x)
a ← scale(x)
c ← prcomp(a)$rotation[,1:2]
print("Required Principal Components")
print(c)
y ← x %*% c
print("Reduced Dimension Data")
print(y)
```

Output →

Original Matrix

 $v_1 \ v_2 \ v_3$

90 60 90

90 90 30

60 60 60

60 60 90

30 30 30

Required Principal Components

PC1 PC2

-0.7071 0.0000

-0.6666 0.3333

-0.2357 -0.9428

Reduced Dimension Data

PC1 PC2

-124.85 -64.85

-130.71 1.71

-96.56 -36.56

-103.63 -64.85

-48.28 -18.28

Name : To perform Simple linear regression with Practical No. 7
Roll No: B03211023 Date 24/05/2021 Page No.

Objective → performing the regression technique on the given dataset

Algorithm -

- (1) Load the data
- (2) Make sure the data meet the assumptions:
 - Normality
 - Correlation
 - Linearity

(3) The formula for linear regression is:

$$y = a_0 + a_1 x$$

$y \rightarrow$ predicted value

$a_0 \rightarrow$ the value of y when x is 0

$a_1 \rightarrow$ regression coefficient

$x \rightarrow$ independent variable

- (4) Display the regression line

Code →

```
M2 ← matrix(c(95, 85, 85, 95, 80, 70, 70, 70, 65, 60, 20),  
nrow = 5, byrow = TRUE)  
mean_x = mean(M2[, 1])  
mean_y = mean(M2[, 2])  
print("Mean of x:")  
print(mean_x)  
print("Mean of y:")  
print(mean_y)
```

$$a1n = 0$$

$$a1d = 0$$

for (i in 1: length(m2[,1])) {

$$a1n = a1n + (m2[,1][i] - \text{mean}_x) * (m2[,2][i] - \text{mean}_y)$$

$$a1d = a1d + (m2[,1][i] - \text{mean}_x)^2 / 2$$

}

$$a1 = a1n / a1d$$

$$a0 = \text{mean}_y - (a1 * \text{mean}_x)$$

print ("Regression line: ")

print (cat("y = ", a0, "+", a1, "x"))

Output →

Mean of x:

78

Mean of y:

77

Regression line:

$$y = 26.78082 + 0.6488356 x$$

Practical Name: To perform K-Means clustering operation on Iris dataset Practical No. 8

Objective → Perform the clustering technique to categorize Iris data in groups.

Algorithm →

- ① Load the Iris dataset
- ② Specify the number of clusters (k) to be created
- ③ Select randomly k -objects from the dataset as the initial cluster centers or means.
- ④ Assign each observation to their closest centroid, based on the Euclidean distance between the object and the centroid.
- ⑤ For each of the k clusters update the cluster centroid by calculating the new mean values of all the datapoints in the cluster.
- ⑥ Iteratively minimize the total within sum of square i.e., iterate steps 3 & 4 until cluster can't be updated further.
- ⑦ Print the clusters formed.

Code →

```
library(datasets)
df <- iris
set.seed(25)
nrows <- sample(nrow(df))
df <- df[nrows, ]
```

Practical Name: K-Means Clustering Practical No. 8

print ("Head of this dataset:")

head(df)

df <- df[, 2:4]

n <- nrow(df)

cols <- names(df)

df[["Cluster"]] <- NA

df[["Dist"]] <- 999

m <- length(cols)

n-c <- 3

s <- sample(1:n, n-c)

smdf <- df[0,]

for (i in s) {

smdf <- rbind(smdf, df[i,])

{}

s <- nrow(smdf)

smdf <- smdf[, -5]

smdf <- smdf[, -5]

print ("Centroids selected :")

smdf

a <- FALSE

while (a == FALSE) {

for (i in 1:n) {

for (j in 1:s) {

dist <- 0

for (k in 1:m) {

dist <- dist + abs(df[i, k] - smdf[j, k])

{}

```
if (df[i, m+2] > dist) {
```

```
  df[i, m+1] <- i
```

```
  df[i, m+2] <- dist
```

{

{

~~sm = sm[1:m]~~

④

```
sm <- aggregate(x = df[, 1:5],  
by = list(df$cluster),  
FUN = mean)
```

```
sm <- sm[-1]
```

```
a <- identical(sm, smdf)
```

```
smdf <- sm
```

```
print("Number of clusters found: ")
```

```
df %>% count(cluster)
```

```
print("final centroids")
```

```
print(smdf)
```

Output =

Head of iris Dataset :

Sepal.length	Sepal.Width	Petal.length	Petal.Width	Species
5.1	2.6	5.6	1.4	Virginica
5.1	3.3	1.7	0.5	Setosa
5.3	3.7	1.5	0.2	Setosa
6.1	2.8	4.0	1.3	Versicolor
6.5	3.0	5.8	2.2	Virginica

Centroid Selected:

Sepal.length	Sepal.Width	Petal.length	Petal.Width
7.2	3.6	6.1	2.5
6.4	2.8	5.6	2.1
7.2	3.0	5.8	1.6

Number of clusters found:

Cluster	n
1	12
2	88
3	50

Final Centroid:

Sepal.length	Sepal.Width	Petal.length	Petal.Width
7.291667	3.208233	6.25833	2.233333
5.388636	3.132955	2.798864	0.8056818
6.296000	2.888000	4.846000	1.644000

Practical Name: To perform market basket analysis using
Association Rules (Apriori) Practical No. 10

Objective → Perform the association rule for analysis

Algorithm →

- (1) Start with itemsets containing just a single item, such as {apple} & {pear}.
- (2) Determine the support for itemsets. Keep the itemsets that meet your minimum support threshold, and remove itemsets that do not.
- (3) Using the itemsets you have kept from Step 1, generate all the possible itemset configuration.
- (4) Repeat Step 1 & 2 until there are no newer itemsets.

Code →

```

df ← read.csv(file = "C:/Users/purchase.csv", header=TRUE)
head(df)
a ← unlist(strsplit(df[, 2], split = ","))
for (j in 2:nrow(df)) {
  s ← unlist(strsplit(df[j, 2], split = ","))
  a ← unlist(cunion(a, s))
}
for (i in 1:length(a)) {
  df[a[i]] ← NA
}
  
```

```

for(i in 1:length(df)) {
  s <- unlist(cbind(split(df[,1], split = " ", ))[1])
  for(j in 1:length(a)) {
    if(a[j] %in% s) {
      df[i, a[j]] <- 1
    } else {
      df[i, a[j]] <- 0
    }
  }
  sdj <- df[, 3 == row(df)]
  k <- 0
  m <- length(a)
  b <- a
  while(k < m) {
    for(j in 1:length(a)) {
      s <- length(b)
      for(i in 1:s) {
        if((colSums(data.frame(sdj[b[i], ] * sdj[a[j], ])) < 2) ||
           (grep(a[j], b[i]))) {
          ??
        } else {
          b <- append(b, paste(b[i], a[j], sep = " "))
          sdj[paste(b[i], b[j], sep = " ")] = sdj[b[i]] +
            sdj[b[j]]
        }
      }
      k <- k + 1
    }
  }
}

```

```

for (i in 1:length(b))
print("Rules Formed")
for (i in 1: length(b)) {
  p <- colnames(edf(b[i]))
  if (p >= 2) {
    print(paste(b[i], "→", p))
  }
}
  
```

Output

TID Items

100	1,3,4
200	2,3,5
300	1,2,3,5
400	2,5

Rules Formed:

- 1 → 2
- 3 → 3
- 2 → 3
- 5 → 3
- 31 → 2
- 23 → 2
- 53 ← 2
- 52 → 3
- 235 → 2