

Practical File

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



Estd. 2000

Name.....Tanuj Yadav.....

Adm.No...2018BCI1026..... Univ. Roll No...1803211050.....

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 -1

Objective : To understand the statistics on data base.

Algorithm:

1. Take a array of numbers from user.
2. find out minimum , maximum , average sum and Square root of each number
3. use min() , max() , avg() , sum() , sqrt() , round() function
4. Display out put

Code:

```

x ← scan()
print(x)
# maximum of an array
max = x[1]
for [i in x]
{
    if (i > max)
        max = i
}
print(paste0("without function = " max))
print(paste0("without function = " max))

# minimum of an array
min = x[1]
for [i in x]
{
    if (i < min)
        min = i
}
print(paste0("without function = " min))

```

`print (paste 0 ("with function = ", min(x)))`

Average of an array
 $j = 0$

`counter = 0`

`for (i in x)`

`{ j = j + i`

`counter = counter + 1`

`print (paste 0 ("without function = ", sum))`

`print (paste 0 ("with function = ", sum(x)))`

Square root of a number

`num = x[3]`

`print (paste 0 ("without function = ", num^0.5))`

`print (sqrt (x[3]))`

Round off

`print (round (x[4]))`

Output

1 2 3 4 5 6 ← Input

[1, 2, 3, 4, 5, 6]

without function = 6

with function = 6

without function = 1

with function = 1

without function = 3.5

with function = 3.5

without function = 2.1

with function = 2.1

without function = 1.73205080756888

with function = 1.732051

Practical - 2

Objective : To handle the data in R program.

1. Explore the UCI machine learning repository
2. Download the dataset in CSV, XLS, .txt file format
3. Use read() function to import the dataset
4. Save the file in data frame & print the head using head() function.

Code

```

data <- read.csv("file=~/home/RTHK/Desktop/prac2.csv", header=T)
head(data 1)
v<- c(1,2,3,4,5,6,7)
write.csv(v, file = "~/home/utsav/Desktop/prac2")
data 2<- read.table (.file.choose())
head(data 2)
write.table(v,file = "~/home/RTHK/Desktop/prac2",
append=F,sep="/",dec=",",row.names
=T,col.names=T)
data 3<- read.xlsx("~/home/utsav/Desktop/prac2.xlsx",
sheetname="Sheet")
head(data 3)
write.xlsx(v,file = "~/home/utsav/Desktop/text.xlsx",
append=F,sep="/",dec=",",row.names
=T,col.names=T)

```

Output

	Name	Roll No
1	A	1
2	B	2

Date / /

Saathi

3	C	3
4	D	4
5	E	5
6	F	6

walte.csv(v,file = "/home/lustav/Desktop/lab.csv")

- [1] Testing .data for reading .text file in R
<views> (Or o.length row.names)

	Name	Branch
1	A	CS IT
2	B	CS IT
3	C	IT
4	D	EC
5	E	ME

Practical-3

Objectives: To understand the handling of multidimensional array

To get the input matrix from user & perform matrix addition, subtraction, multiplication

Inverse, transpose & division operation

using vector concept in R

Algorithm:

1. Take a input from user in 2-D arrays of size 3×3 , i.e.

Matrix 1: 20 30 23	14 21 28	18 19 17	Matrix 2: 15, 12, 22
15, 12, 22	19, 23, 30	21, 23, 11	

2. Arrays :-
3. Perform addition, subtraction, multiplication
Inverse, transpose & division operation on R
4. Display output

Code: \rightarrow `Ave1 <- matrix(c(20,30,23,14,21,28,19,18,19,17), nrow = 3, by.col = TRUE)`

`nrow = 3, by.col = TRUE`
`print(Ave1)`

`Ave2 <- matrix(c(15,12,22,19,25,30,21,23,11), nrow = 3, by.col = TRUE)`
`print(Ave2)`

`Addition <- Ave1 + Ave2`
`print(Addition)`

`Subt Ave1 - Ave2`
`print(Subt)`

`Multiplication <- Ave1 %.*% Ave2`
`print(Multiplication)`

`Inverse <- Inv(Ave1)`

point (Inverse)

Trans $\leftarrow t(Aw_1 w_2)$

point (Trans)

Div $\leftarrow Aw_1 + Aw_2$

point (DPV)

output \div

20	30	23
14	21	28
18	19	17

15	12	28
19	25	30
21	23	11

5	18	-5
-5	-4	-2
-3	-4	6

1353	1519	1713
1197	1337	1330
988	1082	1261

-0.09191176	0.18970588	-0.05882353
-0.03834034	-0.03886555	0.08403361
0.18750000	-0.02500000	0.00000000

15	19	21
12	25	23
20	20	11

Saath!

Date ____ / ____ / ____

1.33333333	2.5000000	0.8214285
0.7368421	0.8400000	0.9333332
0.8571429	0.8260870	1.5454545

Objectives : To find the statistics of data set

Algorithm

- * Apply mean, median & standard deviation formula to find statistics of given data set
- * Apply mean(), median(), mode() & sd() function
- * Find all the parameters of data if you will calc.
- * Display the output

Code: data \leftarrow read.csv ("E:/II Sample.csv")
head(data)

```
mode  $\leftarrow$  function(v) {
    unique  $\leftarrow$  unique(v)
    unique [which.max(tabulate(match(v, unique)))]
}
print(CC["mean":= mean(data)[c("n")]]))
print(CC["median":= median(data)[c("n")]]))
print(CC["mode":= mode(data)[c("n")]]))
print(CC["standard deviation":= sd(data)[c("n")]]))
```

Output \rightarrow

12 3

25 6

37 8

42 3

55 6

43 50 67 9

mean: 53.375

Median: 56

Mode: 23

Standard deviation: 20.3167

Objective: To understand data cleaning process.

Algorithm

1. find missing value in data using `is.na()` function
2. fill missing value using mean | median | mode
3. Normalize the data using min-max normalization

Code

```
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(b)
```

```
d <- a
```

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

```
{d[j] = [(d[i] - min(a)) / (max(a) - min(a))]}
```

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

```
print(d)
```

Output →

~~false~~

FALSE TRUE FALSE FALSE FALSE TRUE FALSE

FALSE FALSE TRUE FALSE FALSE

using mean

3, 11, 56, 8, 11, 9, 12, 10, 11 21, 25

using median

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

Date / /

Period

Saathi

Normalized data

0.00, 0.36 0.09 0.13, 0.22, 0.36, 0.21, 0.40
0.31, 0.36, 0.81, 1.00

Normalized data for each period is as follows:

Period	Normalized Data
1	0.00, 0.36
2	0.09
3	0.13, 0.22
4	0.36
5	0.21
6	0.40
7	0.31
8	0.36
9	0.81, 1.00

Practical - 6

Objective : To get the precise & importa

Algorithm :

1. Pre process the dataset
2. Built piechart function for exploratory data analysis on variance
3. Standardize the data by using scale & apply "pcacomp" function
4. Choose the principal components with highest variance
5. Visualization of data in the new reduced dimension

Code :

```

data1 <- read.csv ("E:/11 Sample csv," header=F)
n <- as.matrix(data1)
print ("original data")
print(n)
a <- scale(n)
c <- pcacomp(a) $ rotation[1:2]
print ("Required principle components")
print(c)
y <- n %*% 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 principle components

 PC_1 and PC_2

$$\begin{matrix} -0.7071 & 0.0000 \\ -0.6660 & 0.3333 \\ -0.2357 & -0.9428 \end{matrix}$$

Reduced dimension data

PC_1	PC_2
-124.33	-64.85
-130.71	1.71
-96.56	36.56
-103.63	-64.85
-48.28	-18.28

Objective To perform Simple Linear regression with

- Performing the regression technique on given dataset

Algorithm

1. Load the data & make sure the data meets the assumption :- * Normality \rightarrow Correlation
* Linearity
2. The formula for linear regression is
 $y = a_0 + a_1 x$ $y \rightarrow$ predicted value ; a_0 - value by of y when x is 0 $a_1 \rightarrow$ regression coefficient
 $x \rightarrow$ independent variable
3. Display the regression line.

Code:

```

M2 <- Matrix(c(95, 85, 85, 95, 80, 70, 70, 70, 65, 60, 70)
  .byrow = TRUE)
meanx = mean(M2[,1])
meany = mean(M2[,2])
print ("mean of x:")
print (meanx)
print ("mean of y:")
print (meany)
ain = 0
old = 0
for (i in 1:length(M2[,1])) {
  ain = ain + ((M2[,1][i] - meanx) * (M2[,2][i] - meany))
  old = old + (M2[,1][i] - meanx)^2
}
ain / old

```

Date / /

Saathie

$$a_1 = \text{ain} / \text{ald}$$

$a_0 = \text{mean } y - (a_1 * \text{mean } x)$

point ("Regression Line")

point (at " $y =$ ", a_0 , " x ", a_1 , " x ")

output

mean of x :

78

mean of y

77

Regression line;

$$y = 26.78082 + 0.6458356x$$

Objective → Perform the clustering technique to categorise iris data in groups.

Algorithm →

- 1) Load the iris dataset
- 2) Specify the number of clusters to be created
- 3) Select randomly K objects from the dataset as the initial cluster center or means.
- 4) Assign each observation to their closest centroid based on the Euclidean distance b/w 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 data points in 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 cluster formed.

Code → library (datasets)

df ← iris

set.seed(25)

rows ← sample (nrow(df))

df ← df [rows,]

```

print("head of iris dataset")
head(df)
df <- df[1:2:4]

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[, -s]
smdf <- smdf[, -c]
print("centroid 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+1] > dist) {
  df[i,m+1] <- j
  df[i,m+2] <- dist
}

sm <- aggregate (n = 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 → Read of Iris dataset

sepal length	sepal width	petal length	petal width
6.1	2.6	5.6	1.4
5.1	3.3	1.7	0.5
5.3	3.0	1.5	0.2
6.1	2.8	4.0	1.3
6.5	3.0	5.8	2.2

Centroid selected →

sepal length	sepal width	petal width	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 centroids →

4.291667	3.208333	6.25833	2.233333
5.388636	3.132955	5.998809	0.8357812
6.296000	2.888000	4.846000	1.044000

Objective → To perform KNN classification for the given dataset

Algorithm →

- ① Load the Dataset.
- ② Calculate the Euclidean measure b/w the data points.
- ③ Arrange the Euclidean distance in ascending order.
- ④ Initialize K and take the first K distance from the list.
- ⑤ Calculate the No. of datapoints belonging to i^{th} class among K points.
- ⑥ If $k_i > k_j \text{ & } i \neq j$, put x in class i .

Code →

```
df <- read.csv(file = "C:/users/Coba.csv", header = TRUE)
head(df)

eu-dist <- function(x, df) {
  df[["distance"]] <- NULL
  for (i in 1: nrow(df)) {
    d <- 0
    for (j in 1: (ncol(df)-2)) {
      d <- d + (x[j] - df[i, j])^2
    }
    df[i, "distance"] <- sqrt(d)
  }
  return(df)
}
```

```
        }  
        return(df)
```

```
KNN <- function(dt,f,k,x) {  
  df <- ec-dist(x,df)  
  df <- df[order(df[ncol(df)]),]  
  m <- 0  
  l <- 0  
  for ( i in 1:k) {  
    if (df[i,3] == "M") {  
      m = m + 1  
    }  
    else {  
      l = l + 1  
    }  
  }  
}
```

```
class- function(m,n) { if (m>1) { print ("m-1")}  
else { print ("L") } }  
print (" class is: ")  
class(m,l)  
}
```

```
print("New data to be predicted: ")  
weight <- readline (prompt= "Enter weight")  
height <- readline (prompt= "Enter height")
```

$x \leftarrow c$ (as integer('height'), as integer('weight'))
KNN(df, s, x)

Outcome → New data point to predict

Enter width : 68 "

Enter height: 168

class is: c

Practical Name → To Perform Market Basket analysis
Using Association Rules

Objective → Perform the association rules for Analysis

Algorithm →

- I Start with itemset containing just a single item such as {apple}, {pear}
- II Determine the support for itemsets. Keep the itemsets that meet your minimum support threshold and remove the itemset that do not.
- III Using the itemsets you have kept from Step II
- IV Repeat Step I & II until there are no newer itemsets

Code →

```
df ← readCSV(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], split = ","))
  a ← unlist(union(a, s))
}
for (i in 1 = length(a)) {
  df[a[i]] ← NA
}
```

Date _____ / _____ / _____

```
for (i in 1:nrow(df)) {
    s <- unlist(strsplit(df[i,], split = " ", ))
    for (j in 1:length(a)) {
        if (a[j] %in% s) {
            df[i, a[j]] <- 1
        } else {
            df[i, a[j]] <- 0
        }
    }
}
```

```
df <- df[, 3:ncol(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 (columns(data.frame(Sdf[b[i]] * Sdf[c(i)])) < 2)

|| (gapply(a[j], b[i]))) {

}

else {

b <- gppmll(b, part(b[i], a[j], sep = ":"))

Sdf[part(b[i], b[j], sep = ":"),]

= Sdf[b[i]] * Sdf[b[j]]

}

k <- k + 1

}

Date — / — /

```

        print ("Rule Formed")
        for ( i in 1:length(b) ) {
            p <- colsum ( diag(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

Rule formed:

1 → 2

3 → 3

2 → 3

5 → 3

31 → 2

23 → 2

53 → 2

52 → 3

235 → 2