

DS Practicals Explanation

Practical 1: EDA

`eda_data = read.csv("file.choose()", header = TRUE, sep = ",")` → Reads data from csv
`eda_data` ⇒ Reads entire dataset from entered csv file.
`head(eda_data)` ⇒ Shows first six rows
`summary(eda_data)` ⇒ Shows mean, min, max, etc. of each column.
`str(eda_data)` ⇒ Shows data type of each column.
`eda_data[1:8]` ⇒ Shows first eight rows
`head(eda_data, 3)` ⇒ Shows first three rows
`head(eda_data, 8)` ⇒ Shows first eight rows
`tail(eda_data, 8)` ⇒ Shows last eight rows.
`eda_data[, 1:5]` ⇒ Shows first 5 columns of dataset

`newdata1 = subset(eda_data, eda_data$Education == "Grad")` ⇒ Creates new dataset, where Education column is equal to "Grad".
`newdata1` ⇒ Displays the above data set
`newdata2 = subset(eda_data, eda_data$Age == "51" & eda_data$Gender == "M")` ⇒ Creates dataset, where Age column is "51" & Gender column is "M".
`newdata2` ⇒ Displays the above data set

`a = eda_data[order(eda_data$Education)]` ⇒ Sorts Education column of eda_data in ascending order & store it in "a".

`a` ⇒ Displays the ascending values of Education Column

`a = eda_data[order(eda_data$Education, decreasing = TRUE)]` ⇒ Sorts Education Column of eda_data in descending order & store it in "a".

`a` ⇒ Displays the descending values of Education Column

`a = colSums(is.na(eda_data))` ⇒ Calculates number of missing values in each column of eda_data & stores it in "a".

`a` ⇒ Shows number of missing values.

`hist(eda_data$Age)` ⇒ Histogram of Age Column.

`bxplst(eda_data$Age)` ⇒ Boxplot of Age Column

`mean(eda_data$Age)` ⇒ Mean of Age Column

`min(eda_data$Age)` ⇒ Minimum Value of Age Column

`max(eda_data$Age)` ⇒ Maximum Value of Age Column

`median(eda_data$Age)` ⇒ Median of Age Column

`counts = table(eda_data$Education, eda_data$Gender)` ⇒ Generates a table that counts number of occurrences of each combination of Education & Gender

`counts` ⇒ Displays the table

`help(pie)` ⇒ Provides information about "pie" function

`plot(eda_dataAge, eda_dataSalary)` ⇒ Creates Scatter Ppt of Age and Salary Column

`hist(eda_data$Age, col = "green", border = "red", xlim = c(20, 70), ylim = c(0, 10))` ⇒ Creates Histogram of Age Column, with green bars & red border, with x-axis values from 20 to 70 & y-axis values from 0 to 10.

PAGE NO / /
DATE / /

PAGE NO / /
DATE / /

`plot(CatData$SalaryType, p1, col="red", main="Salary", xlim=c(0,10), ylim=c(0,80))` \Rightarrow
`bampt(CatData, main="Data Distribution by Education vs Gender")`, \Rightarrow Creates
`col=c("blue", "red"), legend, text=TRUE, xlim=c(0,10))`

`x\Rightarrow Creates Table of Education levels (Grade &
 per=round(x/sum(x)*100) \Rightarrow Calculates Percentage of Education level)`

`pie(x, labels=per, main="City Pie Chart", col=c("blue", "red")) \Rightarrow Creates`

Creates Points Plot of Salary Column, with red points & title Salary, \Rightarrow x-axis range 0 to 80

a Bar Plot

Creates Points Plot of Education levels (Grade & Percentage), \Rightarrow x-axis range 0 to 100

Creates Points Plot of Education level.

Creates Points Plot of City Pie Chart, \Rightarrow x-axis range 0 to 100

Creates Points Plot of Education level.

PAGE NO.	
DATE	/ /

PAGE NO.	
DATE	/ /

Practical 2 :- Hypothesis Testing in R for Statistical Inference

* 2.1:- Independent t-test :-

data1 = read.csv("file.choose()", sep = "", header = T) \Rightarrow Reads CSV file.

t.test(data1\$buk, data1\$Germany, alternative = "two.sided", var.equal = FALSE) \Rightarrow Performs two-tailed t-test; alternative = "two.sided" \Rightarrow Alternative Hypothesis; var.equal = FALSE \Rightarrow variance \neq equal

* 2.2:- Paired t-test :-

data2 = read.csv("file.choose()", sep = "", header = T) \Rightarrow Reads CSV file

t.test(data2\$before, data2\$after, alternative = "greater", paired = T) \Rightarrow Performs Paired t-test

t.test(data2\$before, data2\$after, alternative = "less", paired = T) \Rightarrow Performs Paired t-test

t.test(data2\$before, data2\$after, alternative = "two.sided", paired = T) \Rightarrow Performs Paired t-test

* 2.3:- Test for Correlation :-

data3 = read.csv("file.choose()", sep = "", header = T) \Rightarrow Reads CSV file

cor.test(data3\$mpg, data3\$wt, alternative = "two.sided", method = "pearson") \Rightarrow Perform Pearson Correlation Test.

* 2.4:- One Sample t-test :-

data4 = read.csv("file.choose()", sep = "", header = T) \Rightarrow Reads CSV file

t.test(data4\$b, alternative = "greater", mu = 100) \Rightarrow Performs One Sample t-test

Results & Inference

Independent t-test

Hypothesis: $H_0: \mu_b = \mu_G$; $H_a: \mu_b \neq \mu_G$ (two-tailed)

Performing two-tailed t-test; alternative = "two.sided" \Rightarrow Alternative Hypothesis; var.equal = FALSE \Rightarrow variance \neq equal

Hypothesis: $H_0: \mu_b - \mu_G = 0$; $H_a: \mu_b - \mu_G \neq 0$ (two-tailed)

Hypothesis: $H_0: \mu_b - \mu_G \geq 0$; $H_a: \mu_b - \mu_G < 0$ (less)

Hypothesis: $H_0: \mu_b - \mu_G \leq 0$; $H_a: \mu_b - \mu_G > 0$ (greater)

Performing one-tailed t-test; alternative = "less" \Rightarrow Alternative Hypothesis

Performing one-tailed t-test; alternative = "greater" \Rightarrow Alternative Hypothesis

Hypothesis: $H_0: \mu_b = \mu_G$; $H_a: \mu_b \neq \mu_G$ (two-tailed)

Performing two-tailed t-test; alternative = "two.sided" \Rightarrow Alternative Hypothesis

Hypothesis: $H_0: \mu_b = \mu_G$; $H_a: \mu_b \neq \mu_G$ (two-tailed)

Performing two-tailed t-test; alternative = "two.sided" \Rightarrow Alternative Hypothesis

Hypothesis: $H_0: \mu_b = \mu_G$; $H_a: \mu_b \neq \mu_G$ (two-tailed)

Performing two-tailed t-test; alternative = "two.sided" \Rightarrow Alternative Hypothesis

Hypothesis: $H_0: \mu_b = \mu_G$; $H_a: \mu_b \neq \mu_G$ (two-tailed)

Performing two-tailed t-test; alternative = "two.sided" \Rightarrow Alternative Hypothesis

Hypothesis: $H_0: \mu_b = \mu_G$; $H_a: \mu_b \neq \mu_G$ (two-tailed)

Performing two-tailed t-test; alternative = "two.sided" \Rightarrow Alternative Hypothesis

Hypothesis: $H_0: \mu_b = \mu_G$; $H_a: \mu_b \neq \mu_G$ (two-tailed)

Performing two-tailed t-test; alternative = "two.sided" \Rightarrow Alternative Hypothesis

Hypothesis: $H_0: \mu_b = \mu_G$; $H_a: \mu_b \neq \mu_G$ (two-tailed)

Performing two-tailed t-test; alternative = "two.sided" \Rightarrow Alternative Hypothesis

Hypothesis: $H_0: \mu_b = \mu_G$; $H_a: \mu_b \neq \mu_G$ (two-tailed)

Performing two-tailed t-test; alternative = "two.sided" \Rightarrow Alternative Hypothesis

Hypothesis: $H_0: \mu_b = \mu_G$; $H_a: \mu_b \neq \mu_G$ (two-tailed)

Performing two-tailed t-test; alternative = "two.sided" \Rightarrow Alternative Hypothesis

Hypothesis: $H_0: \mu_b = \mu_G$; $H_a: \mu_b \neq \mu_G$ (two-tailed)

Performing two-tailed t-test; alternative = "two.sided" \Rightarrow Alternative Hypothesis

Hypothesis: $H_0: \mu_b = \mu_G$; $H_a: \mu_b \neq \mu_G$ (two-tailed)

Performing two-tailed t-test; alternative = "two.sided" \Rightarrow Alternative Hypothesis

Practical 4:- Anova

* Ftest :-

`ftest = read.csv(file.choose(), sep = "", header = T)` \Rightarrow Reads CSV file
`var(ftest, ftest[,1], ftest[,2], alternative = "two.sided")` \Rightarrow Performs ftest

* One Way Anova :-

`data1 = read.csv(file.choose(), sep = "", header = T)` \Rightarrow Reads CSV file
`names(data1)` \Rightarrow Returns names of columns in data1 i.e. dept, exp, sal
`summary(data1)` \Rightarrow Gives Max, min, mean, etc of each column in data1 i.e. dept, exp, sal
`head(data1)` \Rightarrow Displays first 6 rows of data1 i.e. dept, exp, sal
`anov = aov(formula = sal ~ dept, data = data1)` \Rightarrow Performs Anova, with sal as response variable & dept as predictor variable.
`summary(anov)` \Rightarrow Displays Summary of Anova results.

* Two Way Anova :-

`data2 = read.csv(file.choose(), sep = "", header = T)` \Rightarrow Reads CSV file
`names(data2)` \Rightarrow Returns names of columns in data2.
`summary(data2)` \Rightarrow Gives Max, min, mean, etc of each column in data2.
`head(data2)` \Rightarrow Displays first 6 rows of data2 i.e. dept, exp, sal
`anov = aov(formula = sal ~ dept + exp + dept * exp, data = data2)` \Rightarrow Performs ANOVA, dependent variable \rightarrow sal, independent variables \rightarrow Dept & exp
`summary(anov)` \Rightarrow Displays Summary of Anova results.

previous work with R language

Fourth step which is analysis of variance which can be done by ftest or anova function. In R, ftest is used for one way analysis of variance and anova is used for two way analysis of variance. Both these functions return a list which contains the following information:-

- 1. Sum of squares (SS) :- It is the sum of squared differences between the observed values and their respective group means.
- 2. Degrees of freedom (df) :- It is the number of independent pieces of information that go into the estimate of a parameter.
- 3. Mean square (MS) :- It is calculated by dividing the sum of squares by the degrees of freedom.
- 4. F-value :- It is the ratio of the mean square between groups to the mean square within groups.
- 5. P-value :- It is the probability of observing a test statistic as extreme as the one calculated, assuming that the null hypothesis is true.

ANOVA table for one way analysis of variance:-

Source of Variation	Sum of Squares (SS)	Degrees of Freedom (df)	Mean Square (MS)	F-value	P-value
Dept	SS_Dept	df_Dept	MS_Dept	F_Dept	P_Dept
Exp	SS_Exp	df_Exp	MS_Exp	F_Exp	P_Exp
Error	SS_Error	df_Error	MS_Error	-	-
Total	SS_Total	df_Total	-	-	-

ANOVA table for two way analysis of variance:-

Source of Variation	Sum of Squares (SS)	Degrees of Freedom (df)	Mean Square (MS)	F-value	P-value
Dept	SS_Dept	df_Dept	MS_Dept	F_Dept	P_Dept
Exp	SS_Exp	df_Exp	MS_Exp	F_Exp	P_Exp
Dept * Exp	SS_Dept_Exp	df_Dept_Exp	MS_Dept_Exp	F_Dept_Exp	P_Dept_Exp
Error	SS_Error	df_Error	MS_Error	-	-
Total	SS_Total	df_Total	-	-	-

Practical 5: Time Series Analysis

`AirPassengers` \Rightarrow Loads the AirPassengers dataset

`view(AirPassengers)` \Rightarrow Opens new window to view dataset.

`class(AirPassengers)` \Rightarrow Displays class of dataset, which "ts" for Time Series

`str(AirPassengers)` \Rightarrow Displays structure of dataset i.e Dimension & Attribute

`start(AirPassengers)` \Rightarrow Displays start date of Time Series

`end(AirPassengers)` \Rightarrow Displays end date of Time Series

`frequency(AirPassengers)` \Rightarrow Displays no. of observation per unit time

`summary(AirPassengers)` \Rightarrow Displays summary of dataset.

`plot(AirPassengers)` \Rightarrow Creates plot of Time Series Data

`abline(lreg = lm(AirPassengers ~ time(AirPassengers)))` \Rightarrow Adds regression line to plot.

`cycle(AirPassengers)` \Rightarrow Calculates position in cycle of each observation

`plot(Caggregate(AirPassengers, Freq=mean))` \Rightarrow Creates Plot of mean values of time series data

`boxplot(AirPassenger ~ cycle(AirPassengers))` \Rightarrow Creates Boxplot of time series data, grouped by cycle.

`acf(AirPassengers)` \Rightarrow Creates Auto-correlation Function Plot

`acf(Clog(AirPassengers))` \Rightarrow Creates Auto-correlation Plot, by applying log & difference in respective methods

`acf(Cdiff(Clog(AirPassengers)))` \Rightarrow log & difference in respective methods

`pacf(Cdiff(Clog(AirPassengers)))` \Rightarrow Creates Partial Auto-correlation Plot, by

`pacf(Cdiff(Clog(AirPassengers)))` \Rightarrow applying log & difference & plot it.

`fit_arima(Clog(AirPassengers), c(0, 1, 1) seasonal = list(Order = c(0, 1, 1) Period = 12))` \Rightarrow Fits an ARIMA model, with first order difference & seasonal difference of order 1

`pred = predict(fit, n.ahead = 10 * 12)` \Rightarrow Predict next 10 years i.e. 120 months.

`pred1 = round(C2.118 * pred[, pred])` \Rightarrow Provides predicted values to nearest integer.

`pred1` \Rightarrow Displays Predicted values.

`plot(AirPassengers, pred1, log = "y", lty = c(1, 3))` \Rightarrow Creates plot of

Original Time Series data & Predicted values

log is used to make variance stationary

Practical 8:- Decision Tree

titanic = read.csv("titanic.csv", header = TRUE) \Rightarrow Reads CSV file

library(rpart) \Rightarrow Loads rpart package used for building decision tree

fit = rpart(Survived ~ Pclass + Gender + Age + SibSp + Parch + Fare + Embarked, data = titanic, method = "class")

plot(fit) \Rightarrow Plots the Created Decision Tree

text(fit) \Rightarrow Adds text to the plot.

install.packages("Rattle") \Rightarrow This package provides interface for

install.packages("rpart.plot") \Rightarrow It provides additional plotting capability

install.packages("RColorBrewer") \Rightarrow It provides Color Palettes for Data View

library(Rattle)

library(rpart.plot)

library(RColorBrewer)

fancyRpartPlot(fit) \Rightarrow Creates fancy Plot using "rpart.plot" package

Prediction = predict(fit, titanic, type = "class") \Rightarrow Predicts the Survival Status of Titanic

Prediction \Rightarrow Displays the Prediction

PAGE NO. / /
DATE / /

PAGE NO. / /
DATE / /

Decision Tree Model

Decision Tree Model is a hierarchical tree structure

Root node is called root of tree, child nodes

Further child of previous is called children of previous

Leaves of tree are called leaf nodes or terminal nodes

Decision Tree Model is a tree structure with root node, child nodes

Leaves of tree are called leaf nodes or terminal nodes

Decision Tree Model is a tree structure with root node, child nodes

Leaves of tree are called leaf nodes or terminal nodes

Decision Tree Model is a tree structure with root node, child nodes

Leaves of tree are called leaf nodes or terminal nodes

Decision Tree Model is a tree structure with root node, child nodes

Leaves of tree are called leaf nodes or terminal nodes

Decision Tree Model is a tree structure with root node, child nodes

Leaves of tree are called leaf nodes or terminal nodes

Decision Tree Model is a tree structure with root node, child nodes

Leaves of tree are called leaf nodes or terminal nodes

Decision Tree Model is a tree structure with root node, child nodes

Leaves of tree are called leaf nodes or terminal nodes

Decision Tree Model is a tree structure with root node, child nodes

Leaves of tree are called leaf nodes or terminal nodes

Decision Tree Model is a tree structure with root node, child nodes

Leaves of tree are called leaf nodes or terminal nodes

Decision Tree Model is a tree structure with root node, child nodes

Leaves of tree are called leaf nodes or terminal nodes

Decision Tree Model is a tree structure with root node, child nodes

Leaves of tree are called leaf nodes or terminal nodes

Decision Tree Model is a tree structure with root node, child nodes

Leaves of tree are called leaf nodes or terminal nodes

Decision Tree Model is a tree structure with root node, child nodes

Leaves of tree are called leaf nodes or terminal nodes

Decision Tree Model is a tree structure with root node, child nodes

Leaves of tree are called leaf nodes or terminal nodes

Decision Tree Model is a tree structure with root node, child nodes

Leaves of tree are called leaf nodes or terminal nodes

Decision Tree Model is a tree structure with root node, child nodes

Leaves of tree are called leaf nodes or terminal nodes

Decision Tree Model is a tree structure with root node, child nodes

Leaves of tree are called leaf nodes or terminal nodes

Decision Tree Model is a tree structure with root node, child nodes

Leaves of tree are called leaf nodes or terminal nodes

Decision Tree Model is a tree structure with root node, child nodes

Leaves of tree are called leaf nodes or terminal nodes

Decision Tree Model is a tree structure with root node, child nodes

Leaves of tree are called leaf nodes or terminal nodes

Decision Tree Model is a tree structure with root node, child nodes

Leaves of tree are called leaf nodes or terminal nodes

Decision Tree Model is a tree structure with root node, child nodes

Leaves of tree are called leaf nodes or terminal nodes

number of passengers for whom survival was predicted to be 1

PAGE NO. / /
DATE / /

Practical 10 :- PCA

`datairis = iris[1:4]` → Assigns first 4 columns of iris dataset to data variable.

`datairis` → Displays the assigned first 4 columns.

`Cov.data = cov(datairis)` → Calculates covariance matrix of datairis.

`Cov.data` → Displays the covariance matrix.

`Eigen.data = eigen(Cov.data)` → Calculates Eigenvectors & Eigenvalues of Covdata.

`Eigen.data` → Displays the Eigenvectors & Eigenvalues.

`PCA.data = princomp(datairis, cor = "False")` → Performs PCA using princomp.

`PCA.data` → Displays the result of PCA.

`Eigen.data$values` → Accesses Eigenvalues of Covariance Matrix.

`PCA.data$dev^2` → Accesses variance of each principal component obtained from PCA.

`PCA.data$loadings[, 1:4]` → Accesses the loadings (i.e. eigenvectors) of first Eigen-data vectors.

`Eigen.data$vector` → Accesses eigenvectors of covariance matrix.

`summary(PCA.data)` → Provides summary of PCA results.

`biplot(PCA.data)` → Creates biplot showing loadings & scores of first two

`scoresplot(PCA.data, type = "lines")` → Creates scoreplot that shows variance

`model2 = PCA.data$loadings[, 1]` → Selects the principal component for second mode.

`model2_scores = as.matrix(datairis) %*% model2` → Calculates scores of second mode by multiplying data by first principal component loading.

`library(Classe)`
`install.packages("e1071")`
`library(e1071)`

libraries for naiveBayes model.

PAGE No. / /
DATE / /

PAGE No. / /
DATE / /

mod1 = naiveBayes (iris[, 1:4], iris[, 5]) \Rightarrow fits a Naive Bayes model using mod1

first 4 columns of iris dataset as predictors & fifth column of iris dataset as response variable

mod2 = naiveBayes (model1\$scores, iris[, 5]) \Rightarrow fits a Naive Bayes model using mod2

scores of first model as predictor & fifth column of iris dataset as response variable

table (predict (mod1, iris[, 1:4]), iris[, 5]) \Rightarrow calculates accuracy of first model.

table (predict (mod2, model1\$scores), iris[, 5]) \Rightarrow calculates accuracy of second model.

round (cor (PCA\$loadings)) \Rightarrow calculates correlation between scores of principal components

Dimensionality Reduction

PCA

bioconductor package

prcomp function

PCA

first two principal components of iris