

R Final Project : Breast Cancer Classification :: Notebook 1.1

Utpal Mishra - 20207425
26 December 2020

Import Libraries

```
require(dplyr)

## Loading required package: dplyr

## Warning: package 'dplyr' was built under R version 3.6.3

##
## Attaching package: 'dplyr'

## The following objects are masked from 'package:stats':
##
##   filter, lag

## The following objects are masked from 'package:base':
##
##   intersect, setdiff, setequal, union

require(repr)

## Loading required package: repr

## Warning: package 'repr' was built under R version 3.6.3

library(corrplot)

## Warning: package 'corrplot' was built under R version 3.6.3

## corrplot 0.84 loaded

library(gplots)

## Warning: package 'gplots' was built under R version 3.6.3

##
## Attaching package: 'gplots'

## The following object is masked from 'package:stats':
##
##   lowess

library(psych)

## Warning: package 'psych' was built under R version 3.6.3

library(fitdistrplus)

## Warning: package 'fitdistrplus' was built under R version 3.6.3

## Loading required package: MASS

##
## Attaching package: 'MASS'

## The following object is masked from 'package:dplyr':
##
##   select

## Loading required package: survival

library(tidyverse)

## Warning: package 'tidyverse' was built under R version 3.6.3

## -- Attaching packages ----- tidyverse 1.3.0 --

## v ggplot2 3.3.2   v purrr   0.3.4
## v tibble  3.0.4   v stringr 1.4.0
## v tidyr   1.1.2   v forcats 0.4.0
## v readr   1.3.1

## Warning: package 'ggplot2' was built under R version 3.6.3

## Warning: package 'tibble' was built under R version 3.6.3

## Warning: package 'tidyr' was built under R version 3.6.3

## Warning: package 'purrr' was built under R version 3.6.3

## -- Conflicts ----- tidyverse_conflicts() --
## x ggplot2::%>%() masks psych::%>%()
## x ggplot2::alpha() masks psych::alpha()
## x dplyr::filter() masks stats::filter()
## x dplyr::lag() masks stats::lag()
## x MASS::select() masks dplyr::select()

library(corpcor)
library("ggplot2", lib.loc=~R/win-library/3.6")
library("Ggally", lib.loc=~R/win-library/3.6")

## Warning: package 'Ggally' was built under R version 3.6.3

## Registered 53 method overwritten by 'Ggally':
##   method from
##   --.gs    ggplot2

cat("IMPORTED LIBRARIES!!!")

## IMPORTED LIBRARIES!!!
```

Import Breast Cancer Data

```
library(readxl) #reading data using the function read.csv() from the library readxl

data <- read.csv("E:/UCD/Lectures/Semester 1/Data Programming with R/Final Project/breast-cancer-wisconsin.wdbc.csv")
data <- data[c(-1)]
head(data) #view(data) #fix(data) #display first 5 rows of the data

##   diagnosis..M.malignant..B.benign. radius..nuca. texture..nuca.
## 1                M                17.99         10.38
## 2                M                20.57         17.77
## 3                M                19.69         21.25
## 4                M                11.42         20.38
## 5                M                20.29         14.34
## 6                M                12.45         15.70
##  perimeter..nuca. area..nuca. smoothness..nuca. compactness..nuca.
## 1         122.80      1081.0      0.11840      0.27550
## 2         132.90      1326.0      0.08474      0.07864
## 3         130.00      1283.0      0.10960      0.15990
## 4          77.58       386.1      0.14250      0.28390
## 5         135.10      1297.0      0.10630      0.13200
## 6          82.57       477.1      0.12780      0.17000
##  concavity..nuca. concave.points..nuca. symmetry..nuca.
## 1          0.3801      0.14710      0.2419
## 2          0.0869      0.07017      0.2812
## 3          0.1974      0.12790      0.2869
## 4          0.2414      0.10520      0.2597
## 5          0.1000      0.10430      0.1009
## 6          0.1578      0.08089      0.2087
##  fractal.dimension..nuca. radius..nucB. texture..nucB. perimeter..nucB.
## 1          0.07871      1.0950      0.9053      8.589
## 2          0.05607      0.5435      0.7339      3.398
## 3          0.05999      0.7456      0.7869      4.585
## 4          0.09744      0.4956      1.1560      3.445
## 5          0.05883      0.7572      0.7813      5.438
## 6          0.07613      0.3345      0.8902      2.217
##  area..nucB. smoothness..nucB. compactness..nucB. concavity..nucB.
## 1         153.40      0.006399      0.04904      0.05373
## 2          74.88      0.005225      0.01308      0.01860
## 3          94.03      0.006150      0.04006      0.03832
## 4          27.23      0.009110      0.07458      0.05661
## 5          94.44      0.011490      0.02461      0.05688
## 6          27.19      0.007510      0.03345      0.03672
##  concave.points..nucB. symmetry..nucB. fractal.dimension..nucB. radius..nucc.
## 1          0.01587      0.03003      0.006193      25.38
## 2          0.01340      0.01389      0.003532      24.99
## 3          0.02058      0.02250      0.004571      23.57
## 4          0.01867      0.02963      0.009200      14.91
## 5          0.01885      0.07756      0.005115      22.54
## 6          0.01137      0.02165      0.005082      15.47
##  texture..nucC. perimeter..nucC. area..nucc. smoothness..nucc.
## 1          17.33      104.60      2010.0      0.1622
## 2          23.41      158.80      1956.0      0.1238
## 3          25.53      152.50      1709.0      0.1444
## 4          26.50       98.87      567.7      0.2098
## 5          16.67      152.20      1575.0      0.1374
## 6          23.75      103.40      741.6      0.1793
##  compactness..nucc. concavity..nucc. concave.points..nucc. symmetry..nucc.
## 1          0.6656      0.7119      0.2654      0.4601
## 2          0.1866      0.2416      0.1860      0.2750
## 3          0.4245      0.4504      0.2430      0.3613
## 4          0.8663      0.6869      0.2575      0.6638
## 5          0.2050      0.4000      0.1625      0.2364
## 6          0.5249      0.5355      0.1741      0.3985
##  fractal.dimension..nucc.
## 1          0.11890
## 2          0.08902
## 3          0.00754
## 4          0.17300
## 5          0.07678
## 6          0.12440
```

Statistical values about Data

```
summary(data) # summary of the data with IQR
```

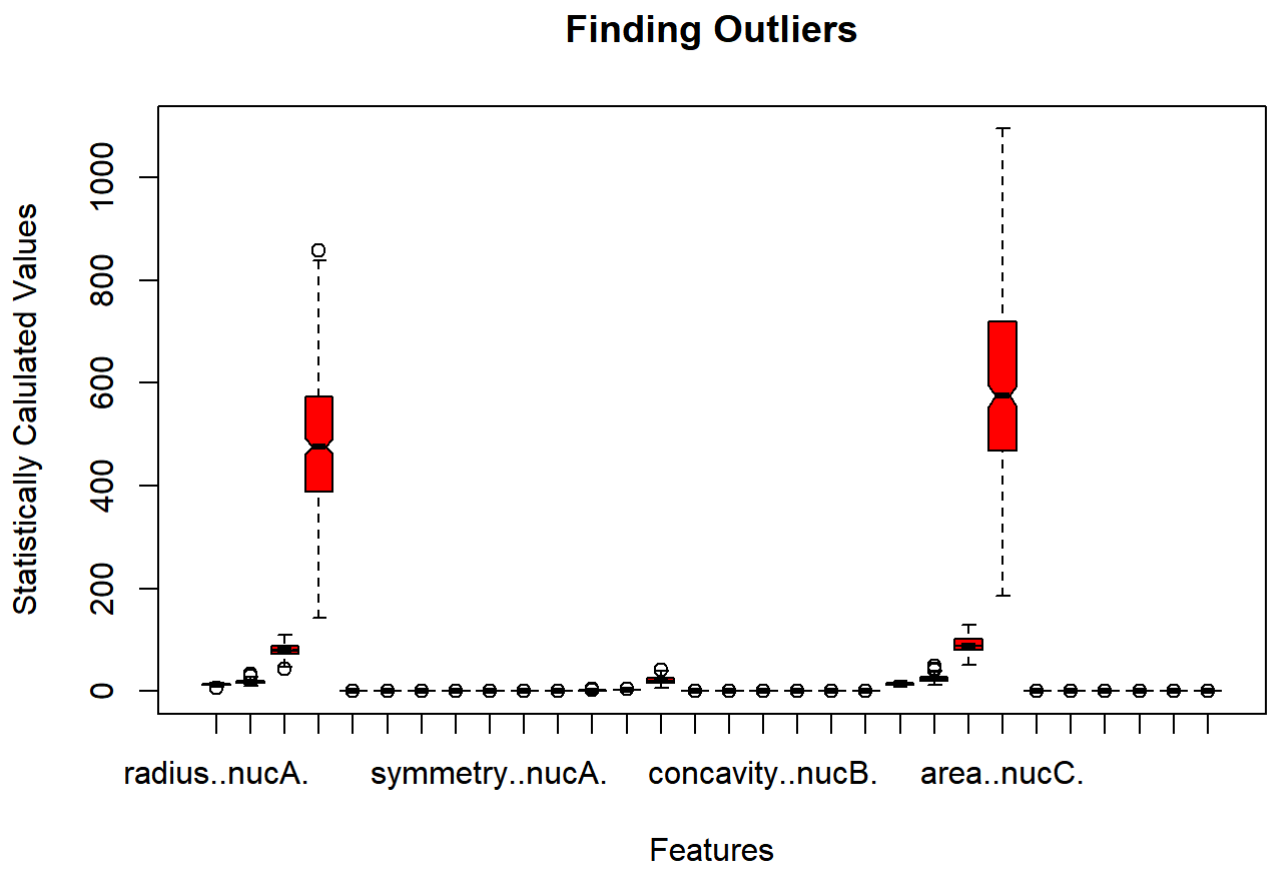


```
for(i in c(1:13)){
  outliers <- boxplot(data$sarea..nucA., plot=FALSE)$out #fetching out the outliers from the boxplot
  x <- data #making a dummy dataset
  x <- x[-which(x$sarea..nucA. %in% outliers), ] #remove outliers from the data
  #boxplot(x[c(-1)], col = "red")
  data <- x #update original data without outliers

  outliers <- boxplot(data$sarea..nucB., plot=FALSE)$out
  x <- data #making a dummy dataset
  x <- x[-which(x$sarea..nucB. %in% outliers), ] #remove outliers from the data
  #boxplot(x[c(-1)], col = "red")
  data <- x #update original data without outliers

  outliers <- boxplot(data$sarea..nucC., plot=FALSE)$out
  x <- data #making a dummy dataset
  x <- x[-which(x$sarea..nucC. %in% outliers), ] #remove outliers from the data
  #boxplot(x[c(-1)], col = "red")
  data <- x #update original data without outliers
}

boxplot(data[c(-1)], col = "red", main = "Finding Outliers", notch = TRUE, xlab = "Features", ylab = "Statistically Calulate
d Values") #using boxplot to witness the data without outliers
```



As can be compared from the above two boxplots, the outliers for the columns nucA and nucC are removed in the later one with change in the y-scale from the multiple iterations

Standardizing the Data

```
tail(data[c(-1)]) #original data

## radius..nucA texture..nucA perimeter..nucA area..nucA smoothness..nucA
## 559 14.59 22.68 96.39 657.1 0.08473
## 560 11.51 23.93 74.52 483.5 0.09261
## 561 14.85 27.15 91.38 698.4 0.09929
## 562 11.20 29.37 78.67 386.0 0.07449
## 563 15.22 38.62 103.48 716.9 0.10488
## 569 7.76 24.54 47.92 181.0 0.05263
## compactness..nucA concavity..nucA concave.points..nucA symmetry..nucA
## 559 0.11380 0.10290 0.03736 0.1454
## 560 0.10210 0.11128 0.04105 0.1388
## 561 0.11268 0.04462 0.04384 0.1537
## 562 0.03558 0.00000 0.00000 0.1060
## 563 0.20870 0.25500 0.09429 0.2128
## 569 0.04362 0.00000 0.00000 0.1587
## fractal.dimension..nucA radius..nucB texture..nucB perimeter..nucB
## 559 0.06147 0.2254 1.106 2.224
## 560 0.06570 0.2388 2.904 1.936
## 561 0.06171 0.3645 1.492 2.888
## 562 0.05502 0.3141 3.896 2.041
## 563 0.07152 0.2602 1.205 2.362
## 569 0.05884 0.3857 1.428 2.548
## area..nucB smoothness..nucB compactness..nucB concavity..nucB
## 559 19.54 0.004242 0.046390 0.06578
## 560 16.97 0.006200 0.029820 0.05738
## 561 29.84 0.007256 0.026780 0.02071
## 562 22.81 0.007594 0.008878 0.00000
## 563 22.65 0.004625 0.048440 0.07359
## 569 19.15 0.007189 0.004660 0.00000
## concave.points..nucB symmetry..nucB fractal.dimension..nucB
## 559 0.01686 0.01638 0.004406
## 560 0.01267 0.01488 0.004738
## 561 0.01626 0.02080 0.005304
## 562 0.00000 0.01989 0.001773
## 563 0.01608 0.02137 0.006142
## 569 0.00000 0.02676 0.002783
## radius..nucC texture..nucC perimeter..nucC area..nucC smoothness..nucC
## 559 15.480 27.27 105.90 733.5 0.10260
## 560 12.480 37.16 82.28 474.2 0.12980
## 561 15.300 33.17 106.20 706.7 0.12410
## 562 11.500 38.38 75.19 439.6 0.09267
## 563 17.520 42.79 128.70 915.0 0.14170
## 569 9.456 30.37 59.16 268.6 0.08996
## compactness..nucC concavity..nucC concave.points..nucC symmetry..nucC
## 559 0.13170 0.13662 0.11650 0.2258
## 560 0.25170 0.3630 0.09653 0.2112
## 561 0.22640 0.1326 0.10480 0.2250
## 562 0.00494 0.00000 0.1566
## 563 0.79170 1.1700 0.23560 0.4089
## 569 0.06444 0.00000 0.00000 0.2871
## fractal.dimension..nucC
## 559 0.00004
## 560 0.08732
## 561 0.08321
## 562 0.05905
## 563 0.14090
## 569 0.07039

data[c(-1)] = as.data.frame(scale(data[c(-1)])) #scaling the data
tail(data[c(-1)])

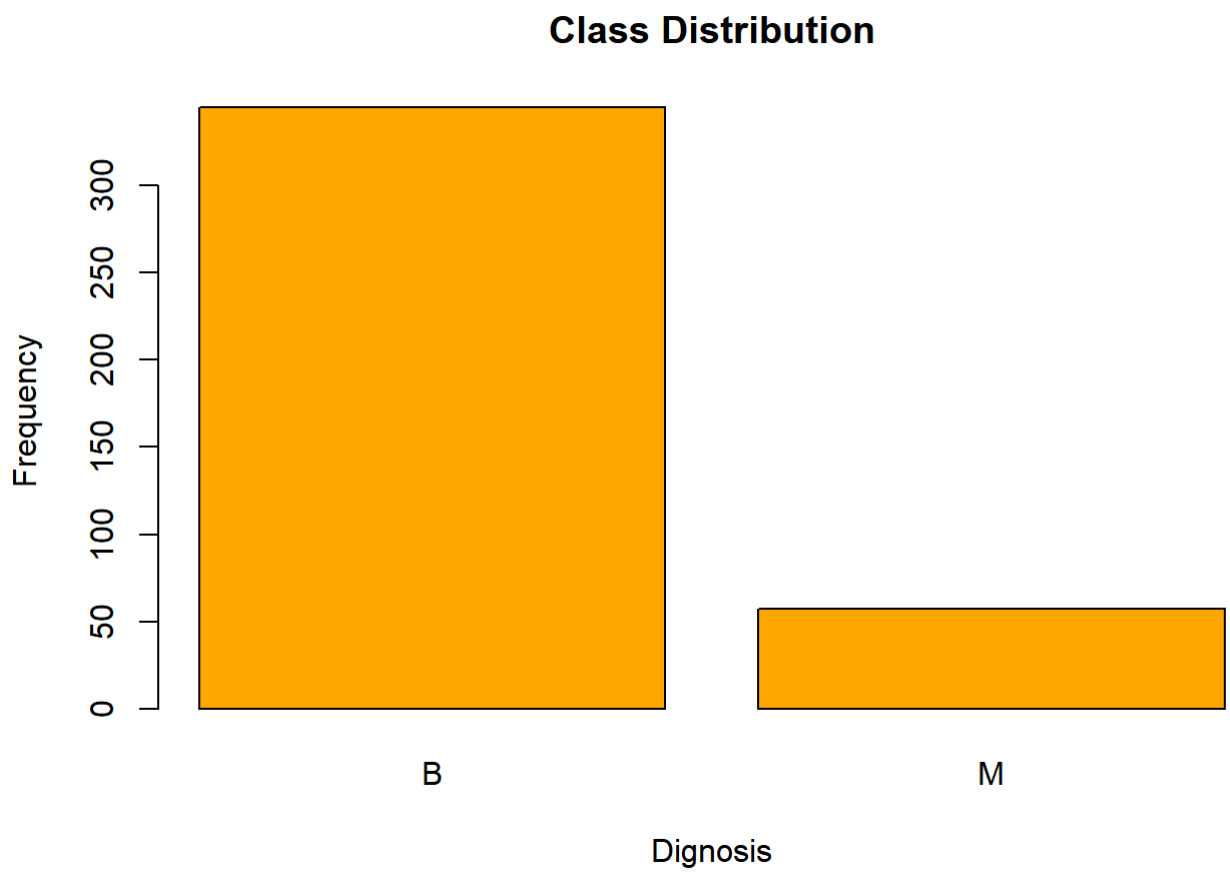
## radius..nucA texture..nucA perimeter..nucA area..nucA smoothness..nucA
## 559 1.2074246 1.044616 1.3320394 1.2613611 -0.66118207
## 560 -0.4776464 1.350437 -0.4344414 -0.5578477 -0.09080665
## 561 0.0119001 2.138231 0.9274411 0.8546216 0.37778904
## 562 -0.6472477 2.681368 -0.7450602 -0.6833846 -1.39188070
## 563 1.5520982 2.987189 1.8981540 1.6903386 0.77095869
## 569 -2.5202750 1.409677 -2.5820001 -2.1539596 -2.95175101
## compactness..nucA concavity..nucA concave.points..nucA symmetry..nucA
## 559 1.0532034 0.8868547 0.3143427 -1.2279767
## 560 0.3240996 1.0467332 0.4868179 -1.4872965
## 561 0.3718533 -0.2357619 0.5790330 -0.9020624
## 562 -1.2454793 -1.0952533 -1.4319109 -2.7760373
## 563 2.8393898 3.8166762 2.9753227 1.4202283
## 569 -1.0557707 -1.0952533 -1.4319109 -0.7054081
## fractal.dimension..nucA radius..nucB texture..nucB perimeter..nucB
## 559 -0.2747282 -0.5792937 -0.14577404 0.34864287
## 560 0.3390652 -0.4325470 3.08738365 -0.07790180
## 561 -0.2399030 0.9440246 0.5450244 1.33206532
## 562 -1.2106140 0.3920818 4.07318122 0.07700928
## 563 1.1835752 -0.1981903 0.02884528 0.55302886
## 569 -0.6563538 1.1761910 0.43028969 0.82850563
## area..nucB smoothness..nucB compactness..nucB concavity..nucB
## 559 -0.2228602 -0.95795413 1.4397753 1.501510
## 560 -0.5682078 0.39381635 0.4399588 1.179023
## 561 1.1612679 0.07800404 0.2565284 -0.228788
## 562 0.2165737 0.18617243 -0.8236596 -1.023873
## 563 0.1050729 -0.8272624 1.5634701 1.001347
## 569 -0.2752585 0.04793516 -1.0781693 -1.023873
## concave.points..nucB symmetry..nucB fractal.dimension..nucB
## 559 1.1883397 -0.53089675 0.3399485
## 560 0.5228143 -0.73236658 0.4780061
## 561 1.2276628 0.06276771 0.7182380
## 562 -1.9693001 -0.05045732 -0.7692234
## 563 1.1922720 0.13912625 1.0712521
## 569 1.9693001 0.86127453 -0.3457530
## radius..nucC texture..nucC perimeter..nucC area..nucC smoothness..nucC
## 559 0.7741599 0.4561861 1.0136417 0.7468115 -1.16107491
## 560 -0.5979053 2.1132404 -0.5104709 -0.6510854 0.02343475
## 561 0.6018360 1.4447220 0.6450415 0.6020089 -0.22470070
## 562 -0.8540242 2.3842456 -0.9679629 -0.8378647 -1.59358803
## 563 1.7071643 3.0565382 2.4848428 1.7254394 0.54165773
## 569 -1.9809471 0.9755863 -2.0623204 -1.7590778 -1.71152352
## compactness..nucC concavity..nucC concave.points..nucC symmetry..nucC
## 559 0.67109766 0.8449898 0.4597465 -0.98501486
## 560 0.22169563 0.8277076 0.1802202 -1.23834065
## 561 0.04704439 -0.4160805 0.3456950 -0.99809572
## 562 -1.11030584 -1.1327380 -1.7512527 -2.18570070
## 563 3.93235458 5.1860536 2.9628777 2.19196817
## 569 -1.06507843 -1.1327380 -1.7512527 0.07860644
## fractal.dimension..nucC
## 559 -0.1716873034
## 560 0.2212809939
## 561 -0.0005733607
## 562 -1.3047099079
## 563 3.1134044789
## 569 -0.6925862140
```

Resampling

```
prop.table(table(data$diagnosis..M.malignant..B.benign.)) #frequency table for diagnosis into Malignant and Benign

## B M
## 0.858209 0.141791

barplot(table(data$diagnosis..M.malignant..B.benign.), col = "orange", xlab = "Dignosis", ylab = "Frequency", main = "Class
Distribution") #frequency plot for diagnosis into Malignant and Benign
```



```
#install.packages("ROSE")
require(ROSE) #using rose library for over-sampling the data

## Loading required package: ROSE

## Warning: package 'ROSE' was built under R version 3.6.3

## Loaded ROSE 0.0-3

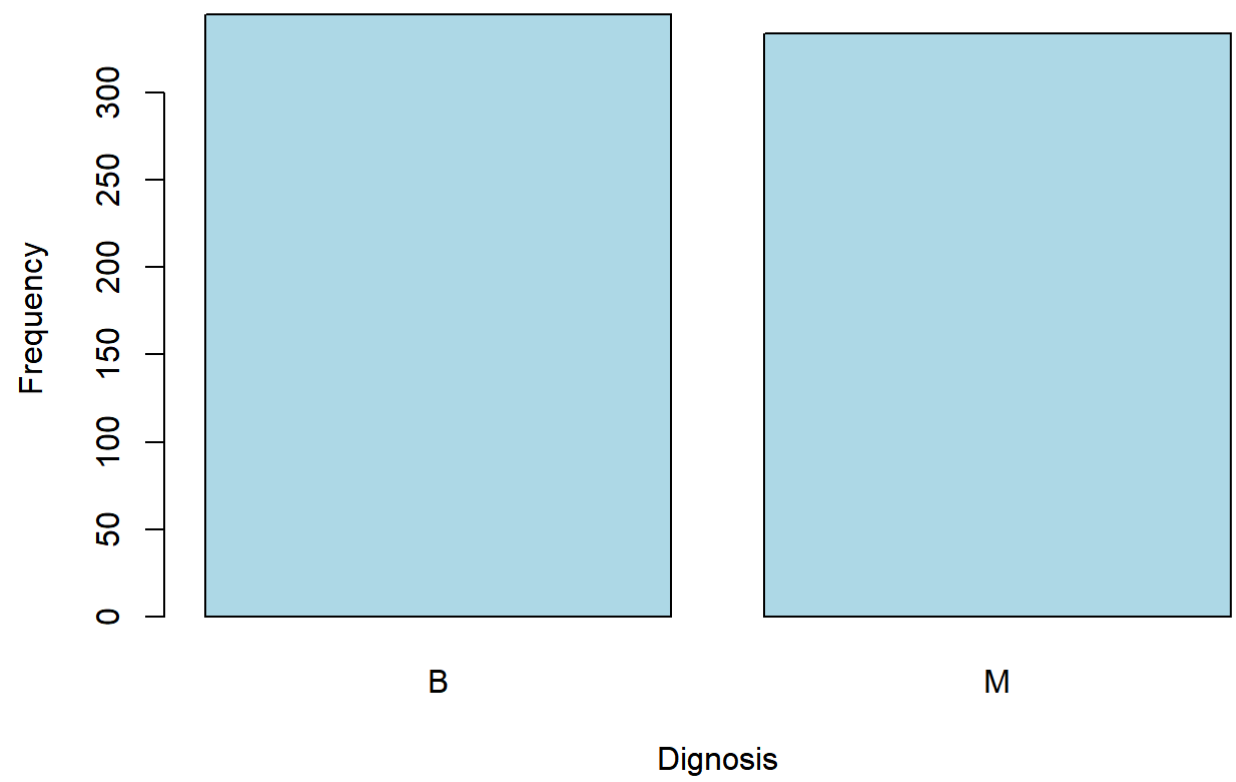
data = ovun.sample(diagnosis..M.malignant..B.benign. ~ ., data = data, method = "over", seed = 1)$data

prop.table(table(ovun.sample(diagnosis..M.malignant..B.benign. ~ ., data = data, method = "over", seed = 1)$data$diagnosis..
M.malignant..B.benign.)) #frequency table for diagnosis into Malignant and Benign

## B M
## 0.5081001 0.4918999

barplot(table(ovun.sample(diagnosis..M.malignant..B.benign. ~ ., data = data, method = "over", seed = 1)$data$diagnosis..M.m
alignant..B.benign.), col = "lightblue", xlab = "Dignosis", ylab = "Frequency", main = "Over-sampled Class Distribution") #f
requency plot for diagnosis into Malignant and Benign
```


Over-sampled Class Distribution



Building Classification Model

Spliting the Data into Training and Testing Data

```
library(caTools)

data[c(-1)] = scale(data[c(-1)])

#data$diagnosis..M.malignant..B.benign. = factor(data$diagnosis..M.malignant..B.benign., levels = c(0, 1))
sample.split(data$diagnosis..M.malignant..B.benign., SplitRatio = 0.80) -> split_data

subset(data, split_data == TRUE) -> train_data
subset(data, split_data == FALSE) -> test_data
```

Decision Tree

Fitting Model

```
library(rpart) #using rpart function to build a decision tree classification model

rpart(diagnosis..M.malignant..B.benign. ~., data = train_data) -> dtmodel #fitting the model
summary(dtmodel) #model summary
```

```
## Call:
## rpart(formula = diagnosis..M.malignant..B.benign. ~., data = train_data)
## n= 543
##
##          CP nsplit  rel error   xerror     xstd
## 1  0.8314607      0  1.00000000  1.03745318  0.04362842
## 2  0.03745318      1  0.16853933  0.16853933  0.02486078
## 3  0.03183521      2  0.13108614  0.15738337  0.02331483
## 4  0.02247191      4  0.06701573  0.10466891  0.01930863
## 5  0.01000000      5  0.04494382  0.07490637  0.01643822
##
## Variable importance
## concave.points..nucC.      concave.points..nucA.      concavity..nucC.
##              18                  15                  14
## concavity..nucA.      compactness..nucC.      compactness..nucA.
##              14                  12                  12
## radius..nucC.      perimeter..nucC.      area..nucC.
##              2                  2                  2
## area..nucA.      perimeter..nucA.      radius..nucA.
##              2                  2                  2
## texture..nucC. fractal.dimension..nucC.      smoothness..nucA.
##              1                  1                  1
## texture..nucA.
##              1
##
## Node number 1: 543 observations,      complexity param=0.8314607
## predicted class=B      expected loss=0.4917127      P(node) =1
## class counts:      276      267
## probabilities: 0.508 0.492
## left son=2 (251 obs) right son=3 (292 obs)
## Primary splits:
## concave.points..nucC. < -0.1665203      to the left,      improve=190.6126, (0 missing)
## concave.points..nucA. < 0.1343971      to the left,      improve=170.9446, (0 missing)
## concavity..nucC. < -0.4698279      to the left,      improve=161.7647, (0 missing)
## perimeter..nucC. < 0.1631721      to the left,      improve=161.7817, (0 missing)
## concavity..nucA. < -0.2312098      to the left,      improve=153.0918, (0 missing)
## Surrogate splits:
## concave.points..nucA. < -0.2748458      to the left,      agree=0.917, adj=0.821, (0 split)
## concavity..nucC. < -0.2262298      to the left,      agree=0.908, adj=0.801, (0 split)
## concavity..nucA. < -0.3020997      to the left,      agree=0.895, adj=0.773, (0 split)
## compactness..nucC. < -0.197674      to the left,      agree=0.877, adj=0.733, (0 split)
## compactness..nucA. < -0.2867455      to the left,      agree=0.869, adj=0.717, (0 split)
##
## Node number 2: 251 observations,      complexity param=0.02247191
## predicted class=B      expected loss=0.03984064      P(node) =0.4622468
## class counts:      241      10
## probabilities: 0.960 0.040
## left son=4 (243 obs) right son=5 (8 obs)
## Primary splits:
## radius..nucA. < 1.303419      to the left,      improve=11.52726, (0 missing)
## perimeter..nucA. < 1.129559      to the left,      improve=11.52726, (0 missing)
## area..nucA. < 1.43081      to the left,      improve=11.52726, (0 missing)
## smoothness..nucB. < -1.483797      to the right,      improve=11.52726, (0 missing)
## radius..nucC. < 1.210413      to the left,      improve=11.52726, (0 missing)
## Surrogate splits:
## perimeter..nucA. < 1.129559      to the left,      agree=1, adj=1, (0 split)
## area..nucA. < 1.43081      to the left,      agree=1, adj=1, (0 split)
## radius..nucC. < 1.210413      to the left,      agree=1, adj=1, (0 split)
## perimeter..nucC. < 0.9525673      to the left,      agree=1, adj=1, (0 split)
## area..nucC. < 1.333808      to the left,      agree=1, adj=1, (0 split)
##
## Node number 3: 292 observations,      complexity param=0.03745318
## predicted class=M      expected loss=0.119863      P(node) =0.5377532
## class counts:      35      257
## probabilities: 0.120 0.880
## left son=6 (10 obs) right son=7 (282 obs)
## Primary splits:
## texture..nucC. < -1.096192      to the left,      improve=16.04221, (0 missing)
## concave.points..nucA. < 0.1343971      to the left,      improve=14.39048, (0 missing)
## texture..nucA. < -1.068879      to the left,      improve=14.38697, (0 missing)
## radius..nucC. < -0.084517485      to the left,      improve=14.23838, (0 missing)
## perimeter..nucC. < -0.6000321      to the left,      improve=12.74339, (0 missing)
## Surrogate splits:
## texture..nucA. < -1.068879      to the left,      agree=0.983, adj=0.5, (0 split)
## perimeter..nucB. < -1.347062      to the left,      agree=0.973, adj=0.2, (0 split)
##
## Node number 4: 243 observations
## predicted class=B      expected loss=0.01234568      P(node) =0.4475138
## class counts:      240      3
## probabilities: 0.988 0.012
##
## Node number 5: 8 observations
## predicted class=M      expected loss=0.125      P(node) =0.01473297
## class counts:      1      7
## probabilities: 0.125 0.875
##
## Node number 6: 10 observations
## predicted class=B      expected loss=0      P(node) =0.01841621
## class counts:      10      0
## probabilities: 1.000 0.000
##
## Node number 7: 282 observations,      complexity param=0.03183521
## predicted class=M      expected loss=0.08865248      P(node) =0.519337
## class counts:      25      257
## probabilities: 0.090 0.910
## left son=14 (34 obs) right son=15 (248 obs)
## Primary splits:
## radius..nucC. < -0.004517485      to the left,      improve=15.021360, (0 missing)
## perimeter..nucC. < -0.6000321      to the left,      improve=13.670600, (0 missing)
## area..nucC. < 0.01442689      to the left,      improve=12.172330, (0 missing)
## concave.points..nucA. < 0.1343971      to the left,      improve=10.035930, (0 missing)
## concave.points..nucC. < 0.233253      to the left,      improve= 7.572433, (0 missing)
## Surrogate splits:
## perimeter..nucC. < -0.07098562      to the left,      agree=0.989, adj=0.912, (0 split)
## area..nucC. < 0.01442689      to the left,      agree=0.979, adj=0.824, (0 split)
## area..nucA. < -0.6703464      to the left,      agree=0.947, adj=0.559, (0 split)
## radius..nucA. < -0.6659743      to the left,      agree=0.943, adj=0.529, (0 split)
## perimeter..nucA. < -0.5931788      to the left,      agree=0.943, adj=0.529, (0 split)
##
## Node number 14: 34 observations,      complexity param=0.03183521
## predicted class=B      expected loss=0.47050802      P(node) =0.0626151
## class counts:      18      16
## probabilities: 0.529 0.471
## left son=28 (17 obs) right son=29 (17 obs)
## Primary splits:
## concave.points..nucA. < 0.4028951      to the left,      improve=15.058820, (0 missing)
## concave.points..nucC. < 0.3325954      to the left,      improve=10.541100, (0 missing)
## smoothness..nucA. < 0.6708758      to the left,      improve= 9.322129, (0 missing)
## concavity..nucA. < 0.2170118      to the left,      improve= 9.322129, (0 missing)
## concavity..nucC. < 0.3597701      to the left,      improve= 8.213904, (0 missing)
## Surrogate splits:
## concave.points..nucC. < 0.3325954      to the left,      agree=0.912, adj=0.824, (0 split)
## concavity..nucA. < 0.2170118      to the left,      agree=0.882, adj=0.765, (0 split)
## concavity..nucC. < 0.3597701      to the left,      agree=0.853, adj=0.706, (0 split)
## smoothness..nucA. < 0.6708758      to the left,      agree=0.824, adj=0.647, (0 split)
## fractal.dimension..nucC. < 0.1733008      to the left,      agree=0.824, adj=0.647, (0 split)
##
## Node number 15: 248 observations
## predicted class=M      expected loss=0.02822581      P(node) =0.4567219
## class counts:      7      241
## probabilities: 0.028 0.972
##
## Node number 28: 17 observations
## predicted class=B      expected loss=0      P(node) =0.03130755
## class counts:      17      0
## probabilities: 1.000 0.000
##
## Node number 29: 17 observations
## predicted class=M      expected loss=0.05882353      P(node) =0.03130755
## class counts:      1      16
## probabilities: 0.059 0.941
```

Predictions

```
library(caret) #using caret to make model predictions
```

```
## Warning: package 'caret' was built under R version 3.6.3
```

```
## Loading required package: lattice
```

```
##
## Attaching package: 'caret'
```

```
## The following object is masked from 'package:purrr':
## lift
```

```
## The following object is masked from 'package:survival':
## cluster
```

```
predict(dtmodel, test_data, type = "class") -> dtresult
#table(test_data$diagnosis..M.malignant..B.benign., dtresult)
```

Confusion Matrix

```
library(caret)
```

```
confusionMatrix(table(test_data$diagnosis..M.malignant..B.benign., dtresult)) #the maximum accuracy of the model is 94.12
```

```
## Confusion Matrix and Statistics
##
##      dresult
##      0  M
##      B 63  6
##      M  2 65
##
##              Accuracy : 0.9412
##              95% CI   : (0.8874, 0.9743)
##      No Information Rate : 0.5221
##      P-Value [Acc > NIR] : <2e-16
##
##              Kappa : 0.8824
##
##      Mcnemar's Test P-Value : 0.2888
##
##              Sensitivity : 0.9692
##              Specificity : 0.9155
##              Pos Pred Value : 0.9130
##              Neg Pred Value : 0.9703
##              Prevalence : 0.4779
##              Detection Rate : 0.4632
##              Detection Prevalence : 0.5074
##              Balanced Accuracy : 0.9424
##
##              'Positive' Class : B
##
```

Tree Model

```
#install.packages("party")
library(party)

## Warning: package 'party' was built under R version 3.6.3

## Loading required package: grid

## Loading required package: mvtnorm

## Warning: package 'mvtnorm' was built under R version 3.6.3

## Loading required package: modeltools

## Warning: package 'modeltools' was built under R version 3.6.3

## Loading required package: stats4

## Loading required package: strucchange

## Warning: package 'strucchange' was built under R version 3.6.3

## Loading required package: zoo

##
## Attaching package: 'zoo'

## The following objects are masked from 'package:base':
##
##      as.Date, as.Date.numeric

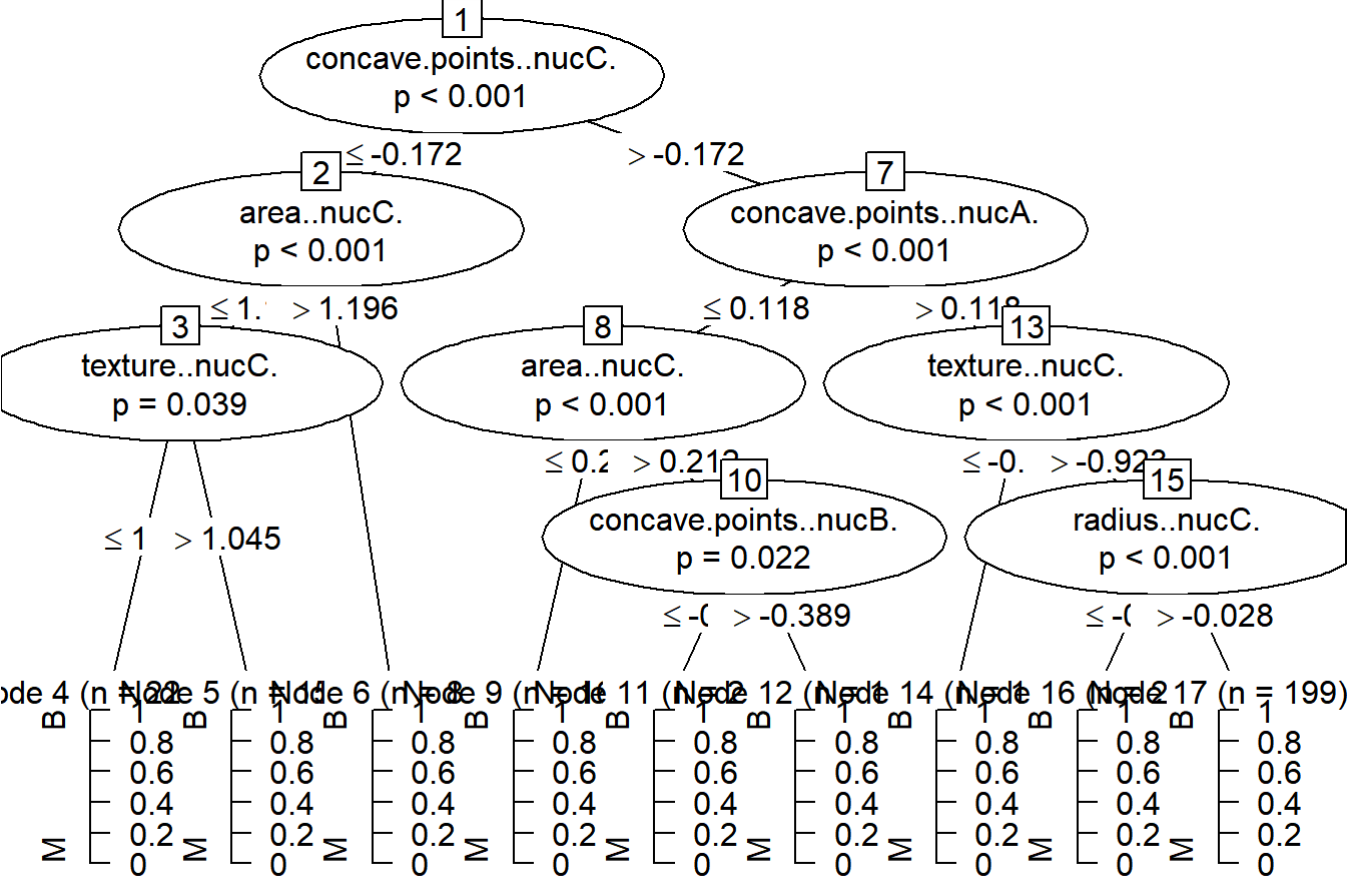
## Loading required package: sandwich

## Warning: package 'sandwich' was built under R version 3.6.3

##
## Attaching package: 'strucchange'

## The following object is masked from 'package:stringr':
##
##      boundary

plot(ctree(diagnosis..M.malignant..B.benign. ~., data = train_data)) #tree model
```



Random Forest

Fitting Model

```
#install.packages("randomForest")
library(randomForest) #using randomForest function to build a random forest classification model

## randomForest 4.6-14

## Type rfNews() to see new features/changes/bug fixes.

##
## Attaching package: 'randomForest'

## The following object is masked from 'package:ggplot2':
##
##      margin

## The following object is masked from 'package:psych':
##
##      outlier

## The following object is masked from 'package:dplyr':
##
##      combine

randomForest(formula = diagnosis..M.malignant..B.benign. ~., data = train_data) -> rfmodel #fitting the model
summary(rfmodel) #model summary

##              Length Class Mode
## call              3 -none- call
## type              1 -none- character
## predicted         543 factor numeric
## err.rate          1500 -none- numeric
## confusion          6 -none- numeric
## votes            1086 matrix numeric
## oob.times         543 -none- numeric
## classes           2 -none- character
## importance        30 -none- numeric
## importanceSD       0 -none- NULL
## localImportance    0 -none- NULL
## proximity          0 -none- NULL
## ntree             1 -none- numeric
## mtry              1 -none- numeric
## forest            14 -none- list
## y                 543 factor numeric
## test              0 -none- NULL
## inbag             0 -none- NULL
## terms             3 terms call
```

Predictions

```
predict(rfmodel, test_data, type = "class") -> rfresult #using caret to make model predictions
#table(test_data$diagnosis..M.malignant..B.benign., rfresult)
```

Confusion Matrix

```
confusionMatrix(table(test_data$diagnosis..M.malignant..B.benign., rfresult)) #the maximum accuracy of the model is 97.79

## Confusion Matrix and Statistics
##
##      rfresult
##      0  M
##      B 66  3
##      M  0 67
##
##              Accuracy : 0.9779
##              95% CI   : (0.9369, 0.9954)
##      No Information Rate : 0.5147
##      P-Value [Acc > NIR] : <2e-16
##
##              Kappa : 0.9559
##
##      Mcnemar's Test P-Value : 0.2482
##
##              Sensitivity : 1.0000
##              Specificity : 0.9571
##              Pos Pred Value : 0.9565
##              Neg Pred Value : 1.0000
##              Prevalence : 0.4853
##              Detection Rate : 0.4853
##              Detection Prevalence : 0.5074
##              Balanced Accuracy : 0.9786
##
##              'Positive' Class : B
##
```

Support Vector Machine

```
#install.packages("e1071")
library(e1071) #using library e1071 to build a SVM classification model

## Warning: package 'e1071' was built under R version 3.6.3

Fitting Model

svm(diagnosis..M.malignant..B.benign. ~., data = train_data, type = 'C-classification', kernel = 'linear') -> svmmodel #fitting the model
summary(svmmodel) #model summary
```



```
## Call:
## svm(formula = diagnosis..M.malignant..B.benign, ~., data = train_data,
## type = "C-classification", kernel = "linear")
##
## Parameters:
## SVM-Type: C-classification
## SVM-Kernel: linear
## cost: 1
##
## Number of Support Vectors: 49
##
## ( 24 25 )
##
##
## Number of Classes: 2
##
## Levels:
## B M
```

Predictions

```
predict(svmmodel, test_data, type = "class") -> svmresult #using caret to make model predictions
#table(test_data$diagnosis..M.malignant..B.benign., svmresult)
```

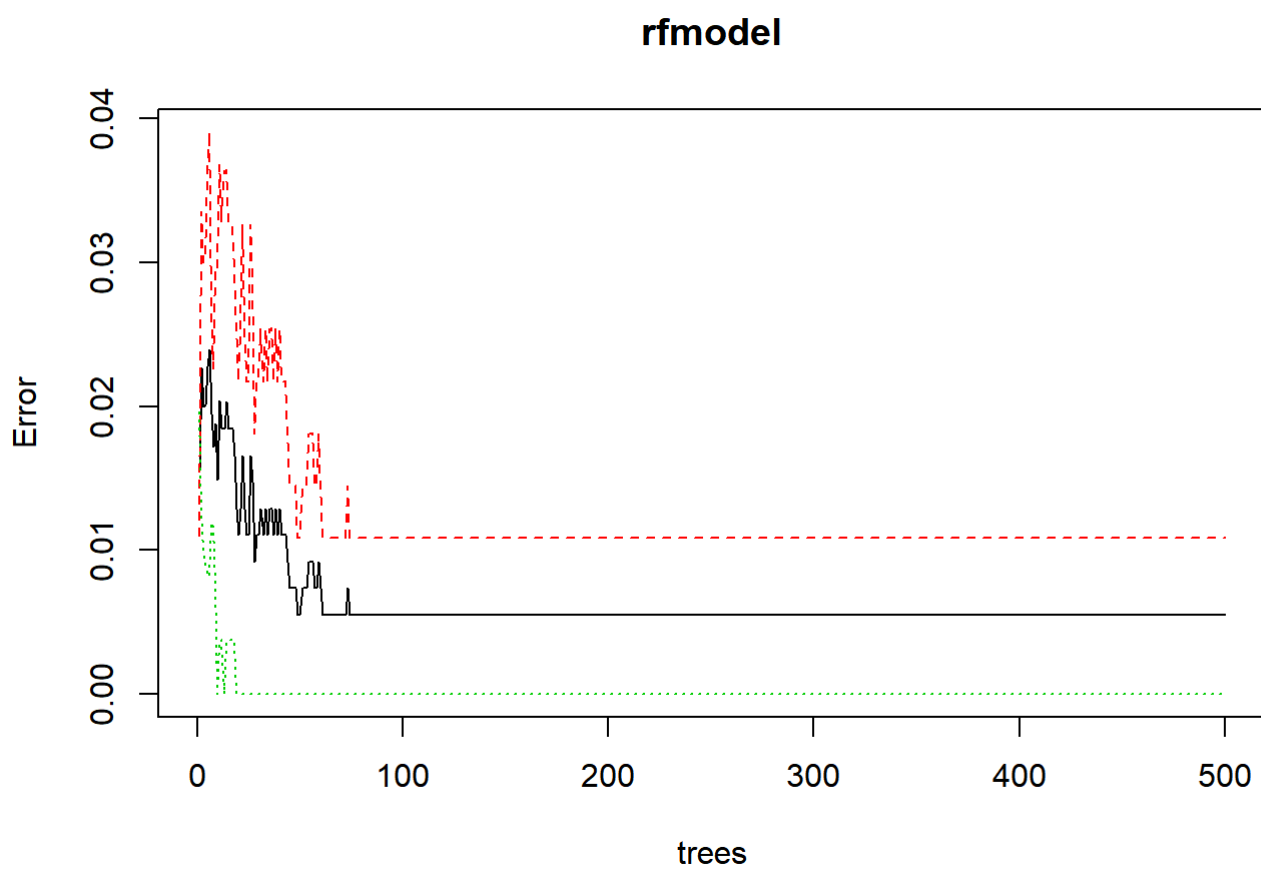
Confusion Matrix

```
confusionMatrix(table(test_data$diagnosis..M.malignant..B.benign., svmresult)) #the maximum accuracy of the model is 95.59
```

```
## Confusion Matrix and Statistics
##
##      svmresult
##      B  M
## B 67  2
## M  4 63
##
##              Accuracy : 0.9559
##              95% CI   : (0.9064, 0.9836)
## No Information Rate : 0.5221
## P-Value [Acc > NRI] : <2e-16
##
##              Kappa   : 0.9117
##
## Mcnemar's Test P-Value : 0.6831
##
## Sensitivity : 0.9437
## Specificity : 0.9692
## Pos Pred Value : 0.9718
## Neg Pred Value : 0.9403
## Prevalence : 0.5221
## Detection Rate : 0.4926
## Detection Prevalence : 0.5074
## Balanced Accuracy : 0.9564
##
## 'Positive' Class : B
##
```

Error vs Model Plot

```
plot(rfmodel)
```



Naive Bayes

```
#install.packages('e1071')
#library(e1071) #using library e1071 to build a Naive Bayes classification model
```

Fitting Model

```
naiveBayes(diagnosis..M.malignant..B.benign, ~., data = train_data, laplace = 1) -> nbmodel #fitting the model
summary(nbmodel) #model summary
```

```
##              Length Class Mode
## apriori      2      table  numeric
## tables      30     -none- list
## levels       2     -none- character
## isnumeric    30     -none- logical
## call         4     -none- call
```

Predictions

```
predict(nbmodel, test_data, type = "class") -> nbresult #using caret to make model predictions
#table(test_data$diagnosis..M.malignant..B.benign., nbresult)
```

Confusion Matrix

```
confusionMatrix(table(test_data$diagnosis..M.malignant..B.benign., nbresult)) #the maximum accuracy of the model is 90.44
```

```
## Confusion Matrix and Statistics
##
##      nbresult
##      B  M
## B 62  7
## M  6 61
##
##              Accuracy : 0.9044
##              95% CI   : (0.8421, 0.9481)
## No Information Rate : 0.5
## P-Value [Acc > NRI] : <2e-16
##
##              Kappa   : 0.8088
##
## Mcnemar's Test P-Value : 1
##
## Sensitivity : 0.9118
## Specificity : 0.8971
## Pos Pred Value : 0.8986
## Neg Pred Value : 0.9184
## Prevalence : 0.5000
## Detection Rate : 0.4559
## Detection Prevalence : 0.5074
## Balanced Accuracy : 0.9044
##
## 'Positive' Class : B
##
```

KNN

```
# library(class) #using library class to build a KNN model
#
# knn(train, test, cl = train$diagnosis..M.malignant..B.benign., k=5) -> knnmodel #fitting the model
# confusionMatrix(table(test$diagnosis..M.malignant..B.benign., knnmodel)) #the maximum accuracy of the model is 98.75
```

Neural Network: Model 1

```
#install.packages('neuralnet')
library(neuralnet) #using library neuralnet to build a neural network classification model
```

```
## Warning: package 'neuralnet' was built under R version 3.6.3
```

```
##
## Attaching package: 'neuralnet'
```

```
## The following object is masked from 'package:dplyr':
##
## compute
```

```
train = train_data #creating dummy training data
test = test_data #creating dummy testing data
```

Categorical Encoding

```
train$diagnosis..M.malignant..B.benign. <- ifelse(train$diagnosis..M.malignant..B.benign. %in% c("B", "B"), 0, 1) #encoding
the categorical/ response variable in training data
tail(train)
```



```
## Confusion Matrix and Statistics
##
##      prediction
## actual 0 1
##      0 67 2
##      1 2 65
##
##              Accuracy : 0.9706
##              95% CI   : (0.9264, 0.9919)
##      No Information Rate : 0.5074
##      P-Value [Acc > NIR] : <2e-16
##
##              Kappa : 0.9412
##
##      Mcnemar's Test P-Value : 1
##
##              Sensitivity : 0.9710
##              Specificity : 0.9701
##      Pos Pred Value : 0.9710
##      Neg Pred Value : 0.9701
##      Prevalence : 0.5074
##      Detection Rate : 0.4926
##      Detection Prevalence : 0.5074
##      Balanced Accuracy : 0.9706
##
##      'Positive' Class : 0
```

Neural Network: Model 2

Fitting Model

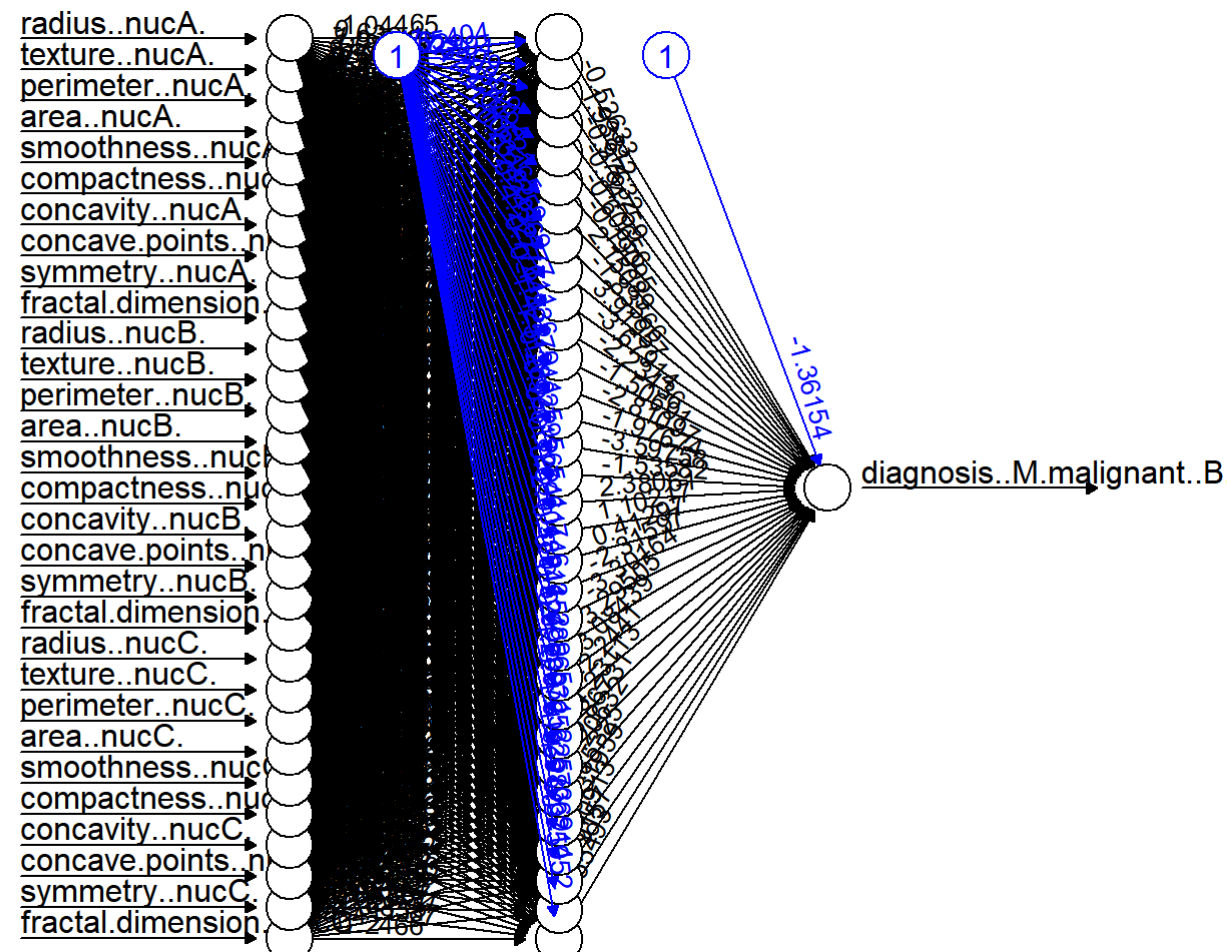
```
neuralnet(diagnosis..M.malignant..B.benign. ~., data = train, threshold = 0.03, hidden = 32, err.fct = "ce", linear.output = FALSE, lifesign = "full",
act.fct = "logistic",rep = 1, algorithm = "backprop", learningrate = 0.003, stepmax = 100000) -> nnmodel
```

```
## hidden: 32   thresh: 0.03   rep: 1/1   steps:   1000 min thresh: 0.217965433376419
##                                                    2000 min thresh: 0.0962128267962868
##                                                    3000 min thresh: 0.0604581594587903
##                                                    4000 min thresh: 0.0437216015083215
##                                                    5000 min thresh: 0.0340804896648614
##                                                    5611 error: 0.18605 time: 14.07 secs
```

```
summary(nnmodel) ##model summary
```

```
##              Length Class      Mode
## call          13 -none-      call
## response       543 -none-     numeric
## covariate     16298 -none-     numeric
## model.list      2 -none-      list
## err.fct        1 -none-     function
## act.fct        1 -none-     function
## linear.output   1 -none-     logical
## data           31 data.frame list
## exclude        0 -none-      NULL
## net.result      1 -none-      list
## weights        1 -none-      list
## generalized.weights 1 -none-  list
## startweights   1 -none-      list
## result.matrix  1028 -none-    numeric
```

```
plot(nnmodel, rep = 1) ##network architecture
```



Results

```
nnresults <- compute(nnmodel, test_data)
results <- data.frame(actual = test$diagnosis..M.malignant..B.benign., prediction = nnresults$net.result)
```

```
head(results)
```

```
##      actual prediction
## 5          0 4.462708e-10
## 12         0 1.937442e-10
## 16         0 4.749167e-11
## 17         0 6.620901e-08
## 23         0 1.671339e-07
## 26         0 2.386366e-01
```

```
predict(nnmodel, test_data, type = "class") -> nnresult ##using caret to make model predictions
#table(test_data$diagnosis..M.malignant..B.benign., nnresult)
```

Confusion Matrix

```
#confusionMatrix(table(test_data$diagnosis..M.malignant..B.benign., nnresult))
roundedresults <- sapply(results,round,digits = 0)
roundedresultsdata = data.frame(roundedresults)
attach(roundedresultsdata)
```

```
## The following objects are masked from roundedresultsdata (pos = 3):
##
##      actual, prediction
```

```
table(actual, prediction)
```

```
##      prediction
## actual 0 1
##      0 68 1
##      1 0 67
```

Model Statistics

```
confusionMatrix(table(actual, prediction)) ##the maximum accuracy of the model is 98.75
```

```
## Confusion Matrix and Statistics
##
##      prediction
## actual 0 1
##      0 68 1
##      1 0 67
##
##              Accuracy : 0.9926
##              95% CI   : (0.9597, 0.9998)
##      No Information Rate : 0.5
##      P-Value [Acc > NIR] : <2e-16
##
##              Kappa : 0.9853
##
##      Mcnemar's Test P-Value : 1
##
##              Sensitivity : 1.0000
##              Specificity : 0.9853
##      Pos Pred Value : 0.9855
##      Neg Pred Value : 1.0000
##      Prevalence : 0.5000
##      Detection Rate : 0.5000
##      Detection Prevalence : 0.5074
##      Balanced Accuracy : 0.9926
##
##      'Positive' Class : 0
```

Hybrid Models

Decision Tree and Random Forest

```
confusionMatrix(table(round((ifelse(dresult %in% c("B", "B"), 0, 1) + ifelse(rresult %in% c("B", "B"), 0, 1))/2), test$diagnosis..M.malignant..B.benign.)) ##averaged ensemble model with maximum accuracy 96.32
```

```
## Confusion Matrix and Statistics
##
##      prediction
## actual 0 1
##      0 66 2
##      1 3 65
##
##              Accuracy : 0.9632
##              95% CI   : (0.9163, 0.988)
##      No Information Rate : 0.5074
##      P-Value [Acc > NIR] : <2e-16
##
##              Kappa : 0.9265
##
##      Mcnemar's Test P-Value : 1
##
##              Sensitivity : 0.9565
##              Specificity : 0.9701
##      Pos Pred Value : 0.9706
##      Neg Pred Value : 0.9559
##      Prevalence : 0.5074
##      Detection Rate : 0.4853
##      Detection Prevalence : 0.5000
##      Balanced Accuracy : 0.9633
##
##      'Positive' Class : 0
```

Decision Tree and SVM

```
confusionMatrix(table(round((ifelse(dresult %in% c("B", "B"), 0, 1) + ifelse(svrresult %in% c("B", "B"), 0, 1))/2), test$diagnosis..M.malignant..B.benign.)) ##average d ensemble model with maximum accuracy 96.32
```

```
## Confusion Matrix and Statistics
##
##      prediction
## actual 0 1
##      0 68 4
##      1 1 63
##
##              Accuracy : 0.9632
##              95% CI   : (0.9163, 0.988)
##      No Information Rate : 0.5074
##      P-Value [Acc > NIR] : <2e-16
##
##              Kappa : 0.9264
##
##      Mcnemar's Test P-Value : 0.3711
##
##              Sensitivity : 0.9855
##              Specificity : 0.9403
##      Pos Pred Value : 0.9444
##      Neg Pred Value : 0.9844
##      Prevalence : 0.5074
##      Detection Rate : 0.5000
##      Detection Prevalence : 0.5294
##      Balanced Accuracy : 0.9629
##
##      'Positive' Class : 0
```

Random Forest and SVM


```
confusionMatrix(table(round((ifelse(rfresult %in% c("B", "B"), 0, 1) +
                               ifelse(svmresult %in% c("B", "B"), 0, 1))/2), test$diagnosis..M.malignant..B.benign.)) #averaged
ensemble model with maximum accuracy 96.32

## Confusion Matrix and Statistics
##
##      0  1
## 0 68  4
## 1  1 63
##
##      Accuracy : 0.9632
##      95% CI   : (0.9163, 0.988)
## No Information Rate : 0.5074
## P-Value [Acc > NIR] : <2e-16
##
##      Kappa : 0.9264
##
## Mcnemar's Test P-Value : 0.3711
##
##      Sensitivity : 0.9855
##      Specificity : 0.9403
##      Pos Pred Value : 0.9444
##      Neg Pred Value : 0.9844
##      Prevalence : 0.5074
##      Detection Rate : 0.5000
##      Detection Prevalence : 0.5294
##      Balanced Accuracy : 0.9629
##
##      'Positive' Class : 0
##
```

Random Forest and Naive Bayes

```
confusionMatrix(table(round((ifelse(rfresult %in% c("B", "B"), 0, 1) +
                               ifelse(nbrresult %in% c("B", "B"), 0, 1))/2), test$diagnosis..M.malignant..B.benign.)) #averaged
ensemble model with maximum accuracy 94.12

## Confusion Matrix and Statistics
##
##      0  1
## 0 67  6
## 1  2 61
##
##      Accuracy : 0.9412
##      95% CI   : (0.8874, 0.9743)
## No Information Rate : 0.5074
## P-Value [Acc > NIR] : <2e-16
##
##      Kappa : 0.8822
##
## Mcnemar's Test P-Value : 0.2888
##
##      Sensitivity : 0.9710
##      Specificity : 0.9104
##      Pos Pred Value : 0.9178
##      Neg Pred Value : 0.9683
##      Prevalence : 0.5074
##      Detection Rate : 0.4926
##      Detection Prevalence : 0.5368
##      Balanced Accuracy : 0.9407
##
##      'Positive' Class : 0
##
```

Random Forest and KNN

```
# confusionMatrix(table(round((ifelse(rfresult %in% c("B", "B"), 0, 1))*1.00 +
                               (ifelse(knnmodel %in% c("B", "B"), 0, 1))*1.00))/2), test$diagnosis..M.malignant..B.benign.)) #
averaged ensemble model with maximum accuracy 97.79
```

Random Forest and Neural Network

```
confusionMatrix(table(round((ifelse(rfresult %in% c("B", "B"), 0, 1))*0.90 +
                               (ifelse(nnresult %in% c("B", "B"), 0, 1))*0.90))/2), test$diagnosis..M.malignant..B.benign.)) #av
eraged ensemble model with maximum accuracy 97.79

## Confusion Matrix and Statistics
##
##      0  1
## 0 66  0
## 1  3 67
##
##      Accuracy : 0.9779
##      95% CI   : (0.9369, 0.9954)
## No Information Rate : 0.5074
## P-Value [Acc > NIR] : <2e-16
##
##      Kappa : 0.9559
##
## Mcnemar's Test P-Value : 0.2482
##
##      Sensitivity : 0.9565
##      Specificity : 1.0000
##      Pos Pred Value : 1.0000
##      Neg Pred Value : 0.9571
##      Prevalence : 0.5074
##      Detection Rate : 0.4853
##      Detection Prevalence : 0.4853
##      Balanced Accuracy : 0.9783
##
##      'Positive' Class : 0
##
```

SVM and Naive Bayes

```
confusionMatrix(table(round((ifelse(dtrresult %in% c("B", "B"), 0, 1) +
                               ifelse(nbrresult %in% c("B", "B"), 0, 1))/2), test$diagnosis..M.malignant..B.benign.)) #averaged
ensemble model with maximum accuracy 92.65

## Confusion Matrix and Statistics
##
##      0  1
## 0 65  6
## 1  4 61
##
##      Accuracy : 0.9265
##      95% CI   : (0.8689, 0.9642)
## No Information Rate : 0.5074
## P-Value [Acc > NIR] : <2e-16
##
##      Kappa : 0.8528
##
## Mcnemar's Test P-Value : 0.7518
##
##      Sensitivity : 0.9420
##      Specificity : 0.9104
##      Pos Pred Value : 0.9155
##      Neg Pred Value : 0.9385
##      Prevalence : 0.5074
##      Detection Rate : 0.4779
##      Detection Prevalence : 0.5221
##      Balanced Accuracy : 0.9262
##
##      'Positive' Class : 0
##
```

SVM and KNN

```
# confusionMatrix(table(round((ifelse(svmresult %in% c("B", "B"), 0, 1))*0.80 +
                               (ifelse(knnmodel %in% c("B", "B"), 0, 1))*0.90))/2), test$diagnosis..M.malignant..B.benign.)) #
averaged ensemble model with maximum accuracy 95.59
```

SVM and Neural Network

```
confusionMatrix(table(round((ifelse(svmresult %in% c("B", "B"), 0, 1) +
                               ifelse(nnresult %in% c("B", "B"), 0, 1))/2), test$diagnosis..M.malignant..B.benign.)) #averaged
ensemble model with maximum accuracy 95.59

## Confusion Matrix and Statistics
##
##      0  1
## 0 67  4
## 1  2 63
##
##      Accuracy : 0.9559
##      95% CI   : (0.9064, 0.9836)
## No Information Rate : 0.5074
## P-Value [Acc > NIR] : <2e-16
##
##      Kappa : 0.9117
##
## Mcnemar's Test P-Value : 0.6831
##
##      Sensitivity : 0.9710
##      Specificity : 0.9403
##      Pos Pred Value : 0.9437
##      Neg Pred Value : 0.9692
##      Prevalence : 0.5074
##      Detection Rate : 0.4926
##      Detection Prevalence : 0.5221
##      Balanced Accuracy : 0.9557
##
##      'Positive' Class : 0
##
```

Naive Bayes and Neural Network

```
confusionMatrix(table(round((ifelse(nbrresult %in% c("B", "B"), 0, 1) +
                               ifelse(nnresult %in% c("B", "B"), 0, 1))/2), test$diagnosis..M.malignant..B.benign.)) #averaged
ensemble model with maximum accuracy 90.44

## Confusion Matrix and Statistics
##
##      0  1
## 0 62  6
## 1  7 61
##
##      Accuracy : 0.9044
##      95% CI   : (0.8421, 0.9481)
## No Information Rate : 0.5074
## P-Value [Acc > NIR] : <2e-16
##
##      Kappa : 0.8088
##
## Mcnemar's Test P-Value : 1
##
##      Sensitivity : 0.8986
##      Specificity : 0.9104
##      Pos Pred Value : 0.9118
##      Neg Pred Value : 0.8971
##      Prevalence : 0.5074
##      Detection Rate : 0.4559
##      Detection Prevalence : 0.5000
##      Balanced Accuracy : 0.9045
##
##      'Positive' Class : 0
##
```

Random Forest, SVM and Neural Network

```
confusionMatrix(table(round((ifelse(rfresult %in% c("B", "B"), 0, 1))*0.90 +
                               ifelse(svmresult %in% c("B", "B"), 0, 1))*0.85 +
                               ifelse(nnresult %in% c("B", "B"), 0, 1))*0.90))/3), test$diagnosis..M.malignant..B.benign.)) #av
eraged ensemble model with maximum accuracy 97.00
```

```
## Confusion Matrix and Statistics
##
##      0  1
## 0 65  0
## 1  4 67
##
##      Accuracy : 0.9786
##      95% CI : (0.9264, 0.9919)
##      No Information Rate : 0.5874
##      P-Value [Acc > NIR] : <2e-16
##
##      Kappa : 0.9412
##
##      Mcnemar's Test P-Value : 0.1336
##
##      Sensitivity : 0.9420
##      Specificity : 1.0000
##      Pos Pred Value : 1.0000
##      Neg Pred Value : 0.9437
##      Prevalence : 0.5874
##      Detection Rate : 0.4779
##      Detection Prevalence : 0.4779
##      Balanced Accuracy : 0.9718
##
##      'Positive' Class : 0
```

Random Forest, SVM and KNN

```
# confusionMatrix(table(round((ifelse(rfresult %in% c("B", "B"), 0, 1)*0.99+
#                               (ifelse(svmresult %in% c("B", "B"), 0, 1)*0.99 +
#                               (ifelse(knmodel %in% c("B", "B"), 0, 1)*0.99))/3), test$diagnosis..M.malignant..B.benign.)) #
averaged ensemble model with maximum accuracy 97.06
```

Ensemble Model: Random Forest, SVM -> Neural Network

Creating Sample Datasets

```
rftrain <- train #creating dummy training data for random forest algorithm
rftest  <- test  #creating dummy training data for random forest algorithm

svmtrain <- train #creating dummy training data for svm algorithm
svmtest  <- test  #creating dummy testing data for svm algorithm

ensembletrain <- train #creating dummy training data for stacked ensemble model
ensembletest  <- test  #creating dummy testing data for stacked ensemble model
```

Prediction for training data using Random Forest and SVM

```
rftrain$diagnosis..M.malignant..B.benign. <- ifelse(predict(rfmodel, train_data, type = "class") %in% c("B", "B"), 0, 1) #en
coding the categorical/ response variable in training data for random forest
svmtrain$diagnosis..M.malignant..B.benign. <- ifelse(predict(svmmodel, train_data, type = "class") %in% c("B", "B"), 0, 1) #
encoding the categorical/ response variable in training data for svm

ensembletrain$diagnosis..M.malignant..B.benign. <- round((rftrain$diagnosis..M.malignant..B.benign. + svmtrain$diagnosis..M.m
alignant..B.benign.)/2) #encoding the categorical/ response variable in training data for stacked ensemble model
```

Predction for testing data using Random Forest and SVM

```
rftest$diagnosis..M.malignant..B.benign. <- ifelse(rfresult %in% c("B", "B"), 0, 1) #encoding the categorical/ response vari
able in testing data for random forest
svmtest$diagnosis..M.malignant..B.benign. <- ifelse(svmresult %in% c("B", "B"), 0, 1) #encoding the categorical/ response va
riable in testing data for svm

ensembletest$diagnosis..M.malignant..B.benign. <- round((rftest$diagnosis..M.malignant..B.benign. + svmtest$diagnosis..M.mal
ignant..B.benign.)/2) #encoding the categorical/ response variable in testing data for stacked ensemble model
```

Training the Neural Network

```
neuralnet(diagnosis..M.malignant..B.benign. ~., data = ensembletrain, threshold = 0.03, hidden = 32, err.fct = "ce", linear.
output = FALSE, lifesign = 'full',
act.fct = "logistic",rep = 1, algorithm = "backprop", learningrate = 0.003, stepmax = 100000) -> ensemblemodel #fitting th
e model

## hidden: 32   thresh: 0.03   rep: 1/1   steps:   1000 min thresh: 0.192942924087557
##                                     2000 min thresh: 0.0902704276126855
##                                     3000 min thresh: 0.0575713627933113
##                                     4000 min thresh: 0.0418529489415114
##                                     5000 min thresh: 0.0327856433197237
##                                     5408 error: 0.18613   time: 14.5 secs
```



Model Results

```
ensbleresults <- compute(ensemblemodel, ensembletest)
ensbleresults <- data.frame(actual = ensembletest$diagnosis..M.malignant..B.benign.,
                             prediction = ensbleresults$net.result)
head(ensbleresults)

## actual prediction
## 5      0 4.270698e-11
## 12     0 4.867729e-09
## 16     0 1.470873e-09
## 17     0 5.565320e-10
## 23     0 2.397216e-07
## 26     0 9.533038e-01
```

Prediction

```
predict(ensemblemodel, ensembletest, type = "class") -> ensembleresult #using caret to make model predictions

#confusionMatrix(table(test_data$diagnosis..M.malignant..B.benign., nresult))
roundedresults <- sapply(ensbleresults,round,digits = 0)
roundedresultsdata = data.frame(roundedresults)
attach(roundedresultsdata)

## The following objects are masked from roundedresultsdata (pos = 3):
##
##      actual, prediction

## The following objects are masked from roundedresultsdata (pos = 4):
##
##      actual, prediction

#table(actual, prediction)

confusionMatrix(table(actual, prediction)) #the maximum accuracy of the model is 97.79

## Confusion Matrix and Statistics
##
##      prediction
## actual 0  1
##      0 69  3
##      1  0 64
##
##      Accuracy : 0.9779
##      95% CI : (0.9369, 0.9954)
##      No Information Rate : 0.5874
##      P-Value [Acc > NIR] : <2e-16
##
##      Kappa : 0.9558
##
##      Mcnemar's Test P-Value : 0.2482
##
##      Sensitivity : 1.0000
##      Specificity : 0.9552
##      Pos Pred Value : 0.9583
##      Neg Pred Value : 1.0000
##      Prevalence : 0.5874
##      Detection Rate : 0.5874
##      Detection Prevalence : 0.5294
##      Balanced Accuracy : 0.9776
##
##      'Positive' Class : 0
```