

Assignment 2

Siddharth Harsukh Pansuria | CWID: 20005837

2/21/2022

Clearing environmental variables

```
rm(list = ls())
```

Loading breast-cancer-wisconsin.csv and replacing ‘?’ with NA

```
data <- read.csv("breast-cancer-wisconsin.csv", na.strings = '?')
```

Viewing the top 10 rows of the loaded data

```
head(data, n = 10)
```

```
##      Sample F1 F2 F3 F4 F5 F6 F7 F8 F9 Class
## 1  1000025  5  1  1  1  2  1  3  1  1     2
## 2  1002945  5  4  4  5  7 10  3  2  1     2
## 3  1015425  3  1  1  1  2  2  3  1  1     2
## 4  1016277  6  8  8  1  3  4  3  7  1     2
## 5  1017023  4  1  1  3  2  1  3  1  1     2
## 6  1017122  8 10 10  8  7 10  9  7  1     4
## 7  1018099  1  1  1  1  2 10  3  1  1     2
## 8  1018561  2  1  2  1  2  1  3  1  1     2
## 9  1033078  2  1  1  1  2  1  1  1  5     2
## 10 1033078  4  2  1  1  2  1  2  1  1     2
```

Deleting rows with missing values

```
data <- na.omit(data)
```

Removing the first column since it only contains ID

```
data <- data[-1]
```

Setting column ‘Class’ as a factor of benign and malignant

```
data$Class <- factor(data$Class, levels = c(2,4), labels = c("benign","malignant"))
```

Viewing the top 10 rows of the loaded data

```
head(data, n = 10)
```

```
##   F1 F2 F3 F4 F5 F6 F7 F8 F9      Class
## 1  5  1  1  1  2  1  3  1  1      benign
## 2  5  4  4  5  7 10  3  2  1      benign
## 3  3  1  1  1  2  2  3  1  1      benign
## 4  6  8  8  1  3  4  3  7  1      benign
## 5  4  1  1  3  2  1  3  1  1      benign
## 6  8 10 10  8  7 10  9  7  1 malignant
## 7  1  1  1  1  2 10  3  1  1      benign
## 8  2  1  2  1  2  1  3  1  1      benign
## 9  2  1  1  1  2  1  1  1  5      benign
## 10 4  2  1  1  2  1  2  1  1      benign
```

Creating a function for min max normalization

```
min_max_norm <- function(x,minX,maxX)
{
  ans <- (x - minX) / (maxX - minX)
  return(ans)
}
```

Getting number of columns with the ‘Class’ column

```
target <- ncol(data)-1
```

Normalizing the data

```
normalized_data <- data.frame(matrix(NA,nrow = nrow(data), ncol = ncol(data)))
colnames(normalized_data) <- c(colnames(data))

for (i in c(1:target)) {
  normalized_data[,i] <- min_max_norm(data[,i], min(data[,i]), max(data[,i]))
}

normalized_data[,ncol(data)] <- data[,ncol(data)]
```

Viewing the top 10 rows of the loaded data

```
head(normalized_data, n = 10)

##          F1          F2          F3          F4          F5          F6          F7
## 1  0.4444444  0.0000000  0.0000000  0.0000000  0.1111111  0.0000000  0.2222222
## 2  0.4444444  0.3333333  0.3333333  0.4444444  0.6666667  1.0000000  0.2222222
## 3  0.2222222  0.0000000  0.0000000  0.0000000  0.1111111  0.1111111  0.2222222
## 4  0.5555556  0.7777778  0.7777778  0.0000000  0.2222222  0.3333333  0.2222222
## 5  0.3333333  0.0000000  0.0000000  0.2222222  0.1111111  0.0000000  0.2222222
## 6  0.7777778  1.0000000  1.0000000  0.7777778  0.6666667  1.0000000  0.8888889
## 7  0.0000000  0.0000000  0.0000000  0.0000000  0.1111111  1.0000000  0.2222222
## 8  0.1111111  0.0000000  0.1111111  0.0000000  0.1111111  0.0000000  0.2222222
## 9  0.1111111  0.0000000  0.0000000  0.0000000  0.1111111  0.0000000  0.0000000
## 10 0.3333333  0.1111111  0.0000000  0.0000000  0.1111111  0.0000000  0.1111111

##          F8          F9      Class
## 1  0.0000000  0.0000000  benign
## 2  0.1111111  0.0000000  benign
## 3  0.0000000  0.0000000  benign
## 4  0.6666667  0.0000000  benign
## 5  0.0000000  0.0000000  benign
## 6  0.6666667  0.0000000 malignant
## 7  0.0000000  0.0000000  benign
## 8  0.0000000  0.0000000  benign
## 9  0.0000000  0.4444444  benign
## 10 0.0000000  0.0000000 benign
```

Splitting the dataset into training and test data

```
id <- sort(sample(nrow(normalized_data), as.integer(.70*nrow(normalized_data)))))

train_data <- normalized_data[id,]
test_data <- normalized_data[-id,]
```

Installing knn

```
install.packages("kknn", repos = "http://cran.us.r-project.org")

## Installing package into 'C:/Users/cookie_/Documents/R/win-library/4.1'
## (as 'lib' is unspecified)

## package 'kknn' successfully unpacked and MD5 sums checked

## Warning: cannot remove prior installation of package 'kknn'

## Warning in file.copy(savedcopy, lib, recursive = TRUE): problem copying C:
## \Users\cookie_\Documents\R\win-library\4.1\00LOCK\kknn\libs\x64\kknn.dll to C:
## \Users\cookie_\Documents\R\win-library\4.1\kknn\libs\x64\kknn.dll: Permission
## denied
```

```
## Warning: restored 'kknn'

##
## The downloaded binary packages are in
##   C:\Users\cookie\AppData\Local\Temp\Rtmpelj4HO\downloaded_packages

library(kknn)
```

Implementing the KNN model for k = 3

```
prediction_k3 <- kknn(formula=Class~., train_data, test_data[,-11], k=3, kernel ="rectangular")
fit_k3 <- fitted(prediction_k3)
```

Implementing the KNN model for k = 5

```
prediction_k5 <- kknn(formula=Class~., train_data, test_data[,-11], k=5, kernel ="rectangular")
fit_k5 <- fitted(prediction_k5)
```

Implementing the KNN model for k = 10

```
prediction_k10 <- kknn(formula=Class~., train_data, test_data[,-11], k=10, kernel ="rectangular")
fit_k10 <- fitted(prediction_k10)
```

Frequency table for KNN mode where k = 3

```
table(Actual=test_data$Class, Fitted=fit_k3)

##           Fitted
## Actual      benign malignant
##   benign        133       2
##   malignant      3       67
```

Frequency table for KNN mode where k = 5

```
table(Actual=test_data$Class, Fitted=fit_k5)

##           Fitted
## Actual      benign malignant
##   benign        133       2
##   malignant      2       68
```

Frequency table for KNN mode where k = 10

```
table(Actual=test_data$Class, Fitted=fit_k10)

##           Fitted
## Actual      benign malignant
##   benign        133       2
##   malignant      2       68
```