PROJECT REPORT



AMERICAN INTERNATIONAL UNIVERSITY BANGLADESH

Submitted By

Name: AHMED IMTIAZ

ID: 20-42933-1

Section: C

Course: INTRODUCTION TO DATA SCIENCE

Semester: SUMMER 2022-2023

Dataset Description

The dataset is chosen from the Kaggle website. The National Institute of Diabetes and Digestive and Kidney Diseases is the source of this dataset. The goal is to determine whether a patient has diabetes based on diagnostic parameters. Here all the patients in this dataset are female. There are 9 attributes and 768 instances in the dataset. The details of the attributes of this dataset are provided below.

- > Pregnancies: Number of times pregnant
- ➤ Glucose: Plasma glucose concentration 2 hours in an oral glucose tolerance test
- ➤ Blood Pressure: Diastolic blood pressure (mm Hg)
- ➤ Skin Thickness: Triceps skin fold thickness (mm)
- > Insulin: 2-Hour serum insulin (mu U/ml)
- ➤ BMI: Body mass index (weight in kg/ (height in m) ^2)
- > Diabetes Pedigree Function: Diabetes pedigree function
- > Age: Age (years)
- > Outcome: Class variable (0 or 1)

Data Preparation

Code Segment: -

```
data <- read.csv("D:/INTRODUCTION TO DATA SCIENCE/Final Term/Project/diabetes.csv")
str(data)
summary(data)</pre>
```

Output Segment: -

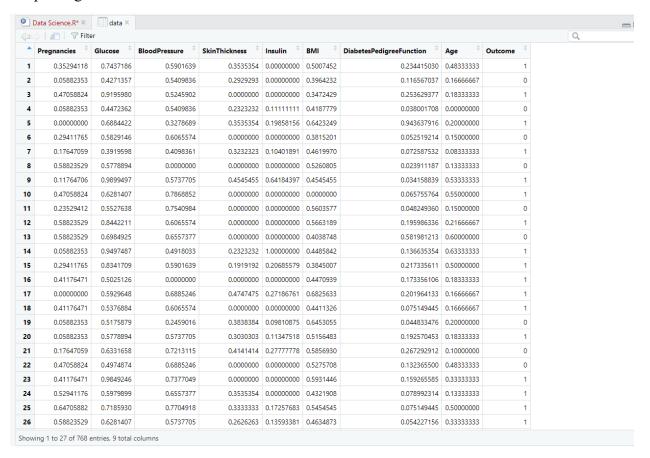
```
> str(data)
'data.frame':
                768 obs. of 9 variables:
                            : num 0.3529 0.0588 0.4706 0.0588 0 ...
$ Pregnancies
                            : num 0.744 0.427 0.92 0.447 0.688 ...
 $ Glucose
                                   0.59 0.541 0.525 0.541 0.328 ...
 $ BloodPressure
                            : num
 $ SkinThickness
                                   0.354 0.293 0 0.232 0.354 ...
                            : num
 $ Insulin
                            : num
                                   0 0 0 0.111 0.199
                                  0.501 0.396 0.347 0.419 0.642 ...
0.234 0.117 0.254 0.038 0.944 ...
0.483 0.167 0.183 0 0.2 ...
 $ BMI
                            : num
 $ DiabetesPedigreeFunction: num
                            $ Outcome
> summary(data)
                                                       SkinThickness
 Pregnancies
                      Glucose
                                     BloodPressure
                                                                            Insulin
                                                                                                 BMI
Min. :0.00000
1st Qu.:0.05882
                   Min. :0.0000
1st Qu.:0.4975
                                                                        Min. :0.00000
1st Qu.:0.00000
                                     Min. :0.0000
                                                       Min. :0.0000
                                                                                           Min.
                                                                                                  :0.0000
                                     1st Qu.:0.5082
                                                       1st Qu.:0.0000
                                                                                           1st Qu.:0.4069
                   Median :0.5879
                                                                         Median :0.03605
 Median :0.17647
                                     Median :0.5902
                                                       Median :0.2323
                                                                                           Median :0.4769
 Mean :0.22618
                   Mean :0.6075
                                     Mean :0.5664
                                                       Mean :0.2074
                                                                        Mean :0.09433
                                                                                           Mean :0.4768
 3rd Qu.:0.35294
                   3rd Qu.:0.7048
                                     3rd Qu.:0.6557
                                                       3rd Qu.:0.3232
                                                                        3rd Qu.:0.15041
                                                                                           3rd Qu.:0.5455
                                            :1.0000
                                                                               :1.00000
 Max.
       :1.00000
                   Max.
                          :1.0000
                                     Max.
                                                      Max.
                                                              :1.0000
                                                                        Max.
                                                                                           Max.
                                                                                                  :1.0000
                          Age
Min. '
DiabetesPedigreeFunction
                                               Outcome
       :0.00000
                                  :0.0000
                                            Min.
                                                   :0.000
 Min.
 1st Qu.:0.07077
                           1st Qu.:0.0500
                                            1st Qu.:0.000
 Median :0.12575
                           Median :0.1333
                                             Median :0.000
 Mean
       :0.16818
                           Mean
                                 :0.2040
                                            Mean :0.349
 3rd Qu.: 0.23409
                           3rd Qu.: 0.3333
                                             3rd Qu.:1.000
        :1.00000
                                  :1.0000
                                                    :1.000
                           Max.
                                            Max.
 Max.
```

Description: - Here, str code is used to check the structure of the dataset, and in the output, we can see that there is no character-related data; all data is in the numeric or number format. On the other hand, the summary code is used to check if there is a null value in the dataset. In the output, we can see no null value in the dataset.

Code Segment: -

```
normalize <- function(x) {
  return ((x - min(x)) / (max(x) - min(x)))
}
data <- as.data.frame(lapply(data, normalize))</pre>
```

Output Segment: -



Description: - At first, a function is created by writing code normalize function (x); the data values are given then normalization is done.

Pearson's Correlation Coefficient

Code Segment: -

```
install.packages("corrplot")
library(corrplot)

correlation_matrix <- cor(data, method = "pearson")
correlation_with_outcome <- correlation_matrix[, "outcome"]
print(correlation_with_outcome)</pre>
```

Output Segment: -

•	Pregnancies [‡]	Glucose [‡]	BloodPressure [‡]	SkinThickness [‡]	Insulin [‡]	вмі 🗦	DiabetesPedigreeFunction [‡]	Age [‡]	Outcome [‡]
Pregnancies	1.00000000	0.12945867	0.14128198	-0.08167177	-0.07353461	0.01768309	-0.03352267	0.54434123	0.22189815
Glucose	0.12945867	1.00000000	0.15258959	0.05732789	0.33135711	0.22107107	0.13733730	0.26351432	0.46658140
BloodPressure	0.14128198	0.15258959	1.00000000	0.20737054	0.08893338	0.28180529	0.04126495	0.23952795	0.06506836
SkinThickness	-0.08167177	0.05732789	0.20737054	1.00000000	0.43678257	0.39257320	0.18392757	-0.11397026	0.07475223
Insulin	-0.07353461	0.33135711	0.08893338	0.43678257	1.00000000	0.19785906	0.18507093	-0.04216295	0.13054795
ВМІ	0.01768309	0.22107107	0.28180529	0.39257320	0.19785906	1.00000000	0.14064695	0.03624187	0.29269466
DiabetesPedigreeFunction	-0.03352267	0.13733730	0.04126495	0.18392757	0.18507093	0.14064695	1.00000000	0.03356131	0.17384407
Age	0.54434123	0.26351432	0.23952795	-0.11397026	-0.04216295	0.03624187	0.03356131	1.00000000	0.23835598
Outcome	0.22189815	0.46658140	0.06506836	0.07475223	0.13054795	0.29269466	0.17384407	0.23835598	1.00000000

Description: - Here, a built-in library is installed first, then a correlation matrix is created, and the Pearson method is used here; the outcome is the target values. Correlation is calculated against the targeted value with all other values.

KNN Find Accuracy 10-Fold Cross Validation

Code Segment: -

```
install.packages(caret)
24 library(caret)
26 • if (!"Outcome" %in% colnames(data)) {
      stop("'Outcome' column not found in the dataset.")
28 - }
29
30 trainIndex <- createDataPartition(data$Outcome, p = 0.7, list = FALSE)</pre>
31 training_data <- data[trainIndex, ]</pre>
32 testing_data <- data[-trainIndex, ]</pre>
33
34 k_values <- expand.grid(k = c(5))
35
36 knn_model <- train(</pre>
37
      Outcome ~ .,
38
      data = training_data,
      method = "knn",
preProcess = c("center", "scale"),
39
40
      trControl = trainControl(method = "cv", number = 10),
41
42
      tuneGrid = k_values
43
44
45
    print(knn_model)
46
47
    predictions <- predict(knn_model, newdata = testing_data)</pre>
48
49
    accuracy <- sum(predictions == testing_data$Outcome) / nrow(testing_data)</pre>
50 cat("Accuracy with optimal k:", accuracy, "\n")
51
```

Output Segment: -

```
k-Nearest Neighbors

538 samples
8 predictor

Pre-processing: centered (8), scaled (8)
Resampling: Cross-Validated (10 fold)
Summary of sample sizes: 484, 484, 484, 484, 484, ...
Resampling results:

RMSE Rsquared MAE
0.4322101 0.2118599 0.3169672

Tuning parameter 'k' was held constant at a value of 5

> # Predict on the testing set using the optimal model
> predictions <- predict(knn_model, newdata = testing_data)
> # Calculate the accuracy
> accuracy <- sum(predictions == testing_data$Outcome) / nrow(testing_data)
> cat("Accuracy with optimal k: ", accuracy, "\n")
Accuracy with optimal k: 0.3173913
```

Description: - At first, the outcome column is checked, then data is split in the ratio of train 70% and test 30%, then the model knn is applied, then the accuracy is found out by using 10-fold cross-validation.

Confusion Matrix

Code Segment: -

```
install.packages ("Metrics")
library(Metrics)

pred_values <- factor(data$Outcome)
actual_values <- factor(data$Outcome)

cf <- caret::confusionMatrix(data=pred_values,reference=actual_values)
print(cf)

predicted <- c(data$Outcome)
actual <- c(data$Outcome)

Metrics::precision(predicted, actual)

Metrics::recall(predicted, actual)
```

Output Segment: -

```
Confusion Matrix and Statistics
           Reference
Prediction
          on 0 1
0 500 0
          1 0 268
    Accuracy : 1
95% CI : (0.9952, 1)
No Information Rate : 0.651
    P-Value [Acc > NIR] : < 2.2e-16
                    Kappa: 1
 Mcnemar's Test P-Value : NA
              Sensitivity: 1.000
              Specificity: 1.000
          Pos Pred Value : 1.000
          Neg Pred Value : 1.000
   Prevalence : 0.651
Detection Rate : 0.651
Detection Prevalence : 0.651
      Balanced Accuracy: 1.000
        'Positive' Class : 0
> predicted <- c(data$Outcome)</pre>
> actual <- c(data$Outcome)</pre>
> Metrics::precision(predicted, actual)
> Metrics::recall(predicted, actual)
[1] 1
```

Description: - First of all, the metrics library is installed, then the predicted value and the actual value are tabled, then the confusion matrix is created, and the recall value and precision value are calculated against the outcome column.