

# “kNN, Linear regression, and multilinear regression”V

Alejandro Vergara Rincon 81190

2023-10-08

## Introduction

In this paper we will analyze the results after using the KNN, linear regression and Multilinear regression models, using the “diabetes\_012\_health\_indicators\_BRFSS2015.csv” dataset, in order to choose the best model according to the indications or requirements of the exercises.

In this regard, we begin by briefly explaining the data set variables:

Summary of the variables in the “diabetes\_012\_health\_indicators\_BRFSS2015.csv” dataset

Diabetes\_012: 0 = no diabetes 1 = prediabetes 2 = diabetes

HighBP: 0 = no high BP 1 = high BP

HighChol: 0 = no high cholesterol 1 = high cholesterol

CholCheck: 0 = no cholesterol check in 5 years 1 = yes cholesterol check in 5 years

BMI: Body Mass Index

Smoker: Have you smoked at least 100 cigarettes in your entire life? [Note: 5 packs = 100 cigarettes] 0 = no 1 = yes

Stroke: (Ever told) you had a stroke. 0 = no 1 = yes

HeartDiseaseorAttack: coronary heart disease (CHD) or myocardial infarction (MI) 0 = no 1 = yes

PhysActivity: physical activity in past 30 days - not including job 0 = no 1 = yes

Fruits: Consume Fruit 1 or more times per day 0 = no 1 = yes

Veggies: Consume Vegetables 1 or more times per day 0 = no 1 = yes

HvyAlcoholConsump: (adult men  $\geq 14$  drinks per week and adult women  $\geq 7$  drinks per week) 0 = no 1 = yes

AnyHealthcare: Have any kind of health care coverage, including health insurance, prepaid plans such as HMO, etc. 0 = no 1 = yes

NoDocbcCost: Was there a time in the past 12 months when you needed to see a doctor but could not because of cost? 0 = no 1 = yes

GenHlth: Would you say that in general your health is: scale 1-5 1 = excellent 2 = very good 3 = good 4 = fair 5 = poor

MentHlth: days of poor mental health scale 1-30 days

PhysHlth: physical illness or injury days in past 30 days scale 1-30

DiffWalk: Do you have serious difficulty walking or climbing stairs? 0 = no 1 = yes

Sex: 0 = female 1 = male

Age: 13-level age category (`_AGEG5YR` see codebook) 1 = 18-24 9 = 60-64 13 = 80 or older

Education: Education level (`EDUCA` see codebook) scale 1-6 1 = Never attended school or only kindergarten 2 =elementary etc.

Income: Income scale (`INCOME2` see codebook) scale 1-8 1 = less than \$10,000 5 = less than \$35,000 8 = \$75,000 or more

For more information about the dataset, please refer to: <https://www.kaggle.com/datasets/alexteboul/diabetes-health-indicators-dataset>

## First part Data exploration and data wrangling

First, we load all the necessary libraries for model creation and data analysis:

```
library(tidyverse)
library(caret)
library(class)
library(gmodels)
library(psych)
```

The Second thing that was done was to load the dataset that had previously been downloaded to the computer.

```
folder<-dirname(rstudioapi::getSourceEditorContext()$path)
parentFolder <-dirname(folder)
data <-
  read.csv(paste0(parentFolder,"/Dataset/diabetes_012.csv"))
```

so:

1. `folder <- dirname(rstudioapi::getSourceEditorContext()$path)`: Retrieves the current script's directory location.
2. `parentFolder <- dirname(folder)`: Obtains the parent directory of the current directory.
3. `data <- read.csv(paste0(parentFolder,"/Dataset/diabetes_012.csv"))`: Reads a CSV file named "diabetes\_012.csv" located in the "Dataset" directory, which is the parent directory, and loads the data into the `data` variable.

```
data$Diabetes_012 <- ifelse(data$Diabetes_012 == 0, 0, 1)
set.seed(1)
data_ <- data[sample(nrow(data), 3000), ]

table(data_.$Sex)
```

```
##
##      0      1
## 1700 1300
```

```
table(data_.$Smoker)
```

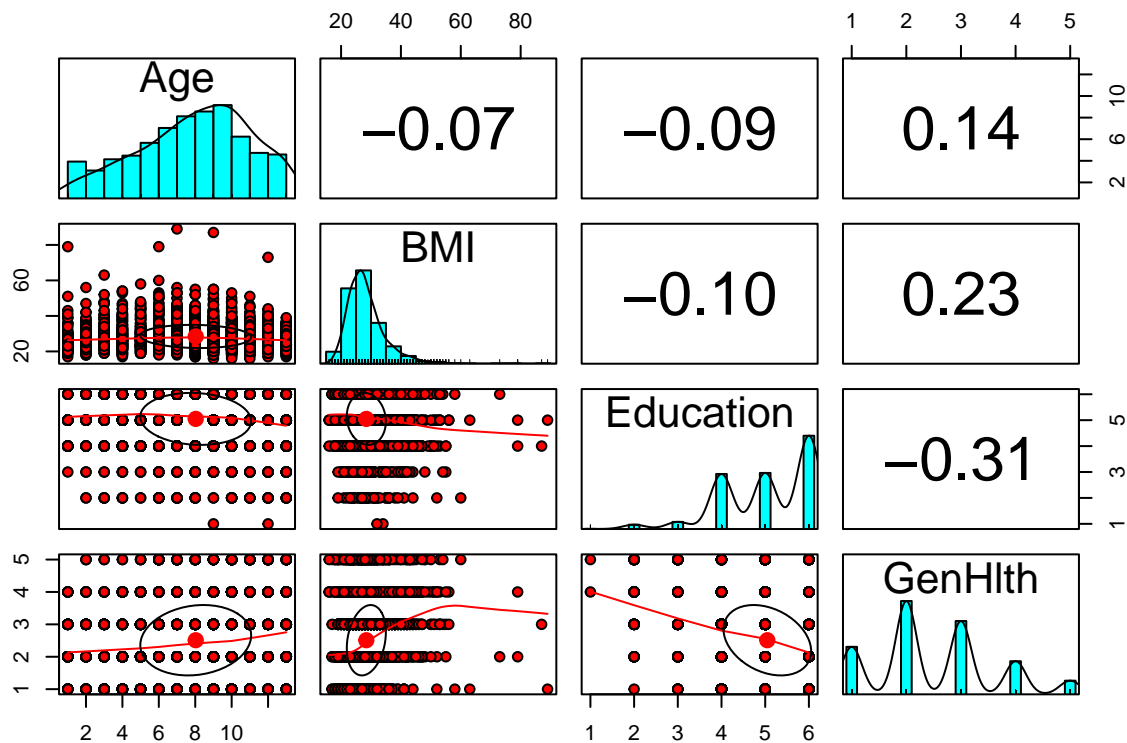
```
##
##      0      1
## 1646 1354
```

```
table(data_.$CholCheck)
```

```
##
##      0      1
## 119 2881
```

1. In this part, we first convert the variable “Diabetes\_012” into a binary variable where 0 represents the absence of diabetes, and 1 represents the presence of diabetes in the dataset. We set a random seed to ensure reproducibility of the results and alignment with the report. We create a subset of the data to explore it in a lighter and more efficient manner. We then calculate the frequency table for the variables “Sex,” “Smoker,” and “CholCheck” in the “data\_” dataset.
2. Based on this data, we can identify that most likely the majority of our sample will be female and non-smokers. However, the proportion of individuals who have had cholesterol checks in the last 5 years is expected to be higher. This will give us an understanding of our dataset to interpret future results.

```
pairs.panels(data_ [c("Age", "BMI", "Education", "GenHlth")],
  pch = 21,
  bg = c("red", "green3")[unclass(data_.$Diabetes_012)])
```



The `pairs.panels` code is used to create a matrix of scatterplots and histograms that display the relationships between selected variables and their distributions. Here’s a detailed explanation of the code:

- `data_[c("Age", "BMI", "Education", "GenHlth")]`: Selects the columns “Age,” “BMI” (Body Mass Index), “Education,” and “GenHlth” (General Health) from the dataset `data_`.
- `pch = 21`: Sets the point type in the scatterplots. In this case, it uses a circle-shaped point with a border.
- `bg = c("red", "green3", "blue", "orange", "yellow")[unclass(data_$Diabetes_012)]`: Defines the background color of the points based on the “Diabetes\_012” variable. Points will be colored based on whether the “Diabetes\_012” variable is 0 or 1.

Regarding observations about the distributions:

- the “Age” variable has a normal distribution, it means that age values are distributed more or less symmetrically around the mean. This can be an important feature in some statistical analyses, as some methods assume data normality.
- BMI has a distribution skewed to the left, it means that Body Mass Index values are skewed to the left. This suggests that most observations have lower BMI values, indicating a concentration of individuals with lower BMI in the sample.

## Second part KNN MODEL

### KNN Models and Experiments to Find Diabetes

```
set.seed(1)
data_estratificada <- data %>%
  group_by(Diabetes_012) %>%
  sample_n(1500, replace = TRUE) %>%
  ungroup()

sample.index <- sample(1:nrow(data_estratificada)
                      ,nrow(data_estratificada)*0.7
                      ,replace = F)

predictors <- c("HighBP", "HighChol", "CholCheck", "BMI", "Smoker", "Stroke", "HeartDiseaseorAttack", "I")

# Original data
train.data <- data_estratificada[sample.index, c(predictors, "Diabetes_012"), drop = FALSE]
test.data <- data_estratificada[-sample.index, c(predictors, "Diabetes_012"), drop = FALSE]

train.data$Diabetes_012 <- factor(train.data$Diabetes_012)
test.data$Diabetes_012 <- factor(test.data$Diabetes_012)
```

The first step would be the data preparation, in this case we are going to use all variables except Diabetes clearly

`set.seed(1)`: Set a fixed random seed for reproducibility.

`stratified_data <- ...`: Stratifies the data by the variable “Diabetes\_012” and creates balanced subsets.

`sample.index <- ...`: Randomly selects 70% of the data for training.

`predictors <- ...`: Lists predictor variables for analysis.

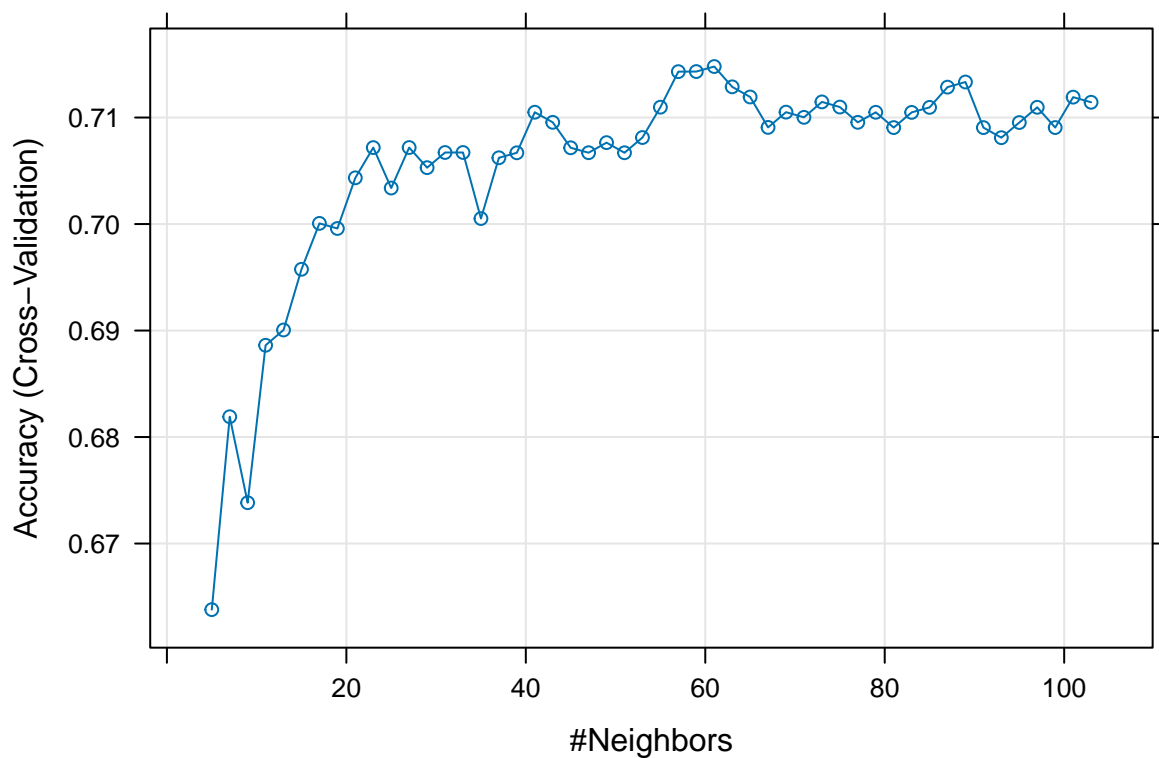
`train.data <- ...`: Creates the training dataset with the selected predictors and “Diabetes\_012” as factor.

`test.data <- ...`: Creates the test data set in a similar fashion.

this code prepares the data for machine learning by splitting it into training and test sets and selecting the relevant predictor variables for analysis. It also ensures that the analysis results are reproducible thanks to the fixed random seed.

```
ctrl <- trainControl(method = "cv", p = 0.7)
knnFit <- train(Diabetes_012 ~ .
  , data = train.data
  , method = "knn", trControl = ctrl
  , preprocess = c("range") # c("center", "scale") for z-score
  , tuneLength = 50)

plot(knnFit)
```



1. The model is trained using the training data (**train.data**), and it specifies that the target variable is “Diabetes\_012,” and all other available predictor variables (“.” indicates all predictor variables).
2. **trainControl**: Here, control parameters for model training are set. It uses cross-validation (**method** = “cv”) with a 70% data partition for training (**p** = 0.7).
3. **train**: This function trains the k-NN model with the specified parameters. It uses the training data, the “knn” method, and the control parameters defined earlier. Additionally, it applies preprocessing

to scale the data within the range (“range”). The **tuneLength** parameter is set to 50, indicating that 50 iterations will be performed to find the optimal value of k.

4. **plot(knnFit)**: Finally, this line of code generates a plot that shows the performance of the trained k-NN model.

```
# Make predictions
knnPredict <- predict(knnFit, newdata = test.data)

# Creates the confusion matrix
confusionMatrix(data = knnPredict, reference = test.data$Diabetes_012)
```

In this code snippet, we are making predictions using the trained k-Nearest Neighbors (k-NN) model and evaluating its performance. In this model, we calculate K using the caret. The reason for doing this is because of the model’s accuracy.

1. **Make predictions**: The **predict** function is used to make predictions on new data. In this case, we’re applying the trained k-NN model (**knnFit**) to the test data (**test.data**) to predict the values of the “Diabetes\_012” variable for the test dataset. These predictions are stored in the **knnPredict** variable.
2. **Creates the confusion matrix**: The **confusionMatrix** function is used to create a confusion matrix to evaluate the performance of the k-NN model’s predictions. It takes two arguments:
  - data**: The vector of predicted values (**knnPredict**).
  - reference**: The true values of the target variable from the test data (**test.data\$Diabetes\_012**). It’s used as a reference to calculate metrics like accuracy, precision, recall, etc.

```
## Confusion Matrix and Statistics
##
##           Reference
## Prediction  0    1
##           0 329 123
##           1 117 331
##
##           Accuracy : 0.7333
##           95% CI : (0.7032, 0.762)
##           No Information Rate : 0.5044
##           P-Value [Acc > NIR] : <2e-16
##
##           Kappa : 0.4667
##
## Mcnemar's Test P-Value : 0.7469
##
##           Sensitivity : 0.7377
##           Specificity : 0.7291
##           Pos Pred Value : 0.7279
##           Neg Pred Value : 0.7388
##           Prevalence : 0.4956
##           Detection Rate : 0.3656
##           Detection Prevalence : 0.5022
##           Balanced Accuracy : 0.7334
##
##           'Positive' Class : 0
##
```

Confusion matrix: The confusion matrix provides a summary of the model's performance in classifying instances into two classes (0 and 1), where:

The rows represent the predicted classes (0 and 1).

The columns represent the actual or reference classes (0 and 1).

The four values of the matrix are

True Negative (TN): 330 - Instances correctly predicted as 0 (no diabetes).

False Positives (FP): 122 - Instances incorrectly predicted as 1 (diabetes).

False Negatives (FN): 116 - Instances incorrectly predicted as 0 (no diabetes).

True Positives (TP): 332 - Instances correctly predicted as 1 (diabetes).

Accuracy: The accuracy of the model is 73.56%, indicating that 73.56% of the predictions made by the model are correct. This metric measures the overall correctness.

Kappa Index: The Kappa index (Kappa) is a measure of inter-rater agreement, in this case, the agreement between model predictions and actual classes. It adjusts the accuracy for the possibility of random agreement and takes values between -1 and 1. In this case, Kappa is 0.4711.

A Kappa of 1 indicates perfect agreement.

A Kappa of 0 indicates agreement equivalent to random agreement.

A Kappa of less than 0 indicates agreement worse than chance.

In this case, a Kappa of 0.4711 indicates moderate agreement between model predictions and actual classes. It suggests that the model performance is better than chance, but that there is room for improvement.

For clarification, please analyze the image regarding the Kappa index:

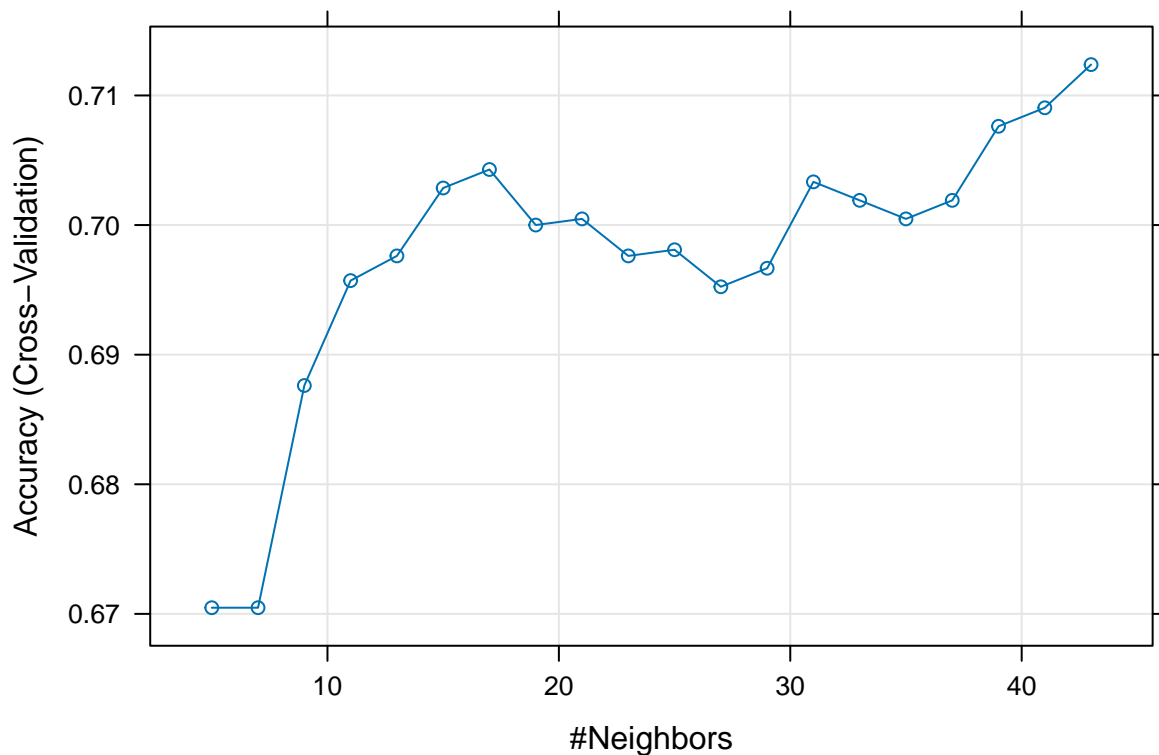
Values of kappa	Interpretation
< 0	No agreement
0-0.19	Poor agreement
0.20-0.39	Fair agreement
0.40-0.59	Moderate agreement
0.60-0.79	Substantial agreement
0.80-1.00	Almost perfect agreement

## Second experiment

```
predictors_to_remove <- c("AnyHealthcare", "NoDocbcCost", "DiffWalk", "Education", "Income")
train.data2 <- train.data[, !(names(train.data) %in% predictors_to_remove)]
test.data2 <- test.data[, !(names(test.data) %in% predictors_to_remove)]
```

```
ctrl <- trainControl(method = "cv", number = 5)
knnFit2 <- train(Diabetes_012 ~ .
  , data = train.data2
  , method = "knn", trControl = ctrl
  , preProcess = c("range") # c("center", "scale") for z-score
  , tuneLength = 20)

plot(knnFit2)
```



In this code:

- **predictors\_to\_remove** is a vector containing the names of predictors (features) that are to be removed from the dataset.
- **train.data2** and **test.data2** are created by subsetting the original **train.data** and **test.data** datasets, respectively, to remove the predictors listed in **predictors\_to\_remove**.

The code essentially removes five predictors (“AnyHealthcare,” “NoDocbcCost,” “DiffWalk,” “Education,” and “Income”) from the training and testing datasets.

this code removes specific predictors from the dataset and trains a k-NN model on the modified data with 5-fold cross-validation while tuning the hyperparameter (k) using a range of values to optimize the model’s performance.



```
# Make predictions
knnPredict2 <- predict(knnFit2, newdata = test.data2)

# Creates the confusion matrix
confusionMatrix(data = knnPredict2, reference = test.data2$Diabetes_012)
```

```
## Confusion Matrix and Statistics
##
##           Reference
## Prediction  0    1
##           0 315 108
##           1 131 346
##
##           Accuracy : 0.7344
##           95% CI : (0.7043, 0.763)
##       No Information Rate : 0.5044
##       P-Value [Acc > NIR] : <2e-16
##
##           Kappa : 0.4686
##
##  McNemar's Test P-Value : 0.1547
##
##           Sensitivity : 0.7063
##           Specificity : 0.7621
##           Pos Pred Value : 0.7447
##           Neg Pred Value : 0.7254
##           Prevalence : 0.4956
##           Detection Rate : 0.3500
##       Detection Prevalence : 0.4700
##           Balanced Accuracy : 0.7342
##
##           'Positive' Class : 0
##
```

The confusion matrix provides a summary of the model's performance in classifying instances into two classes, 0 and 1. Here's what each part of the matrix represents:

- True Negative (TN): 315 - Instances correctly predicted as 0 (no diabetes).
- False Positives (FP): 109 - Instances incorrectly predicted as 1 (diabetes).
- False Negatives (FN): 131 - Instances incorrectly predicted as 0 (no diabetes).
- True Positives (TP): 345 - Instances correctly predicted as 1 (diabetes).

Here are the key performance metrics:

**Accuracy:** The accuracy of the model is 73.33%, indicating that 73.33% of the predictions made by the model are correct. This metric measures overall correctness.

**Kappa:** The Kappa value is 0.4664, indicating moderate agreement between the model's predictions and the actual classes. It suggests that the model's performance is better than random chance, but there is still room for improvement. A higher Kappa value would indicate stronger agreement.

**Sensitivity:** Sensitivity, also known as the True Positive Rate or Recall, is 70.63%. It represents the proportion of actual positive cases (diabetes) correctly predicted by the model.

Specificity: Specificity is 75.99%, indicating the proportion of actual negative cases (no diabetes) correctly predicted by the model.

Positive Predictive Value (Pos Pred Value): Pos Pred Value is 74.29%, which represents the probability that a predicted positive case is truly positive.

Negative Predictive Value (Neg Pred Value): Neg Pred Value is 72.48%, indicating the probability that a predicted negative case is truly negative.

Prevalence: Prevalence is 49.56%, representing the proportion of actual positive cases in the dataset.

Detection Rate: Detection Rate is 35%, indicating the proportion of true positive cases detected by the model.

Detection Prevalence: Detection Prevalence is 47.11%, representing the proportion of the dataset predicted as positive by the model.

Balanced Accuracy: Balanced Accuracy is 73.31%, which takes into account both sensitivity and specificity and provides a balanced measure of model performance.

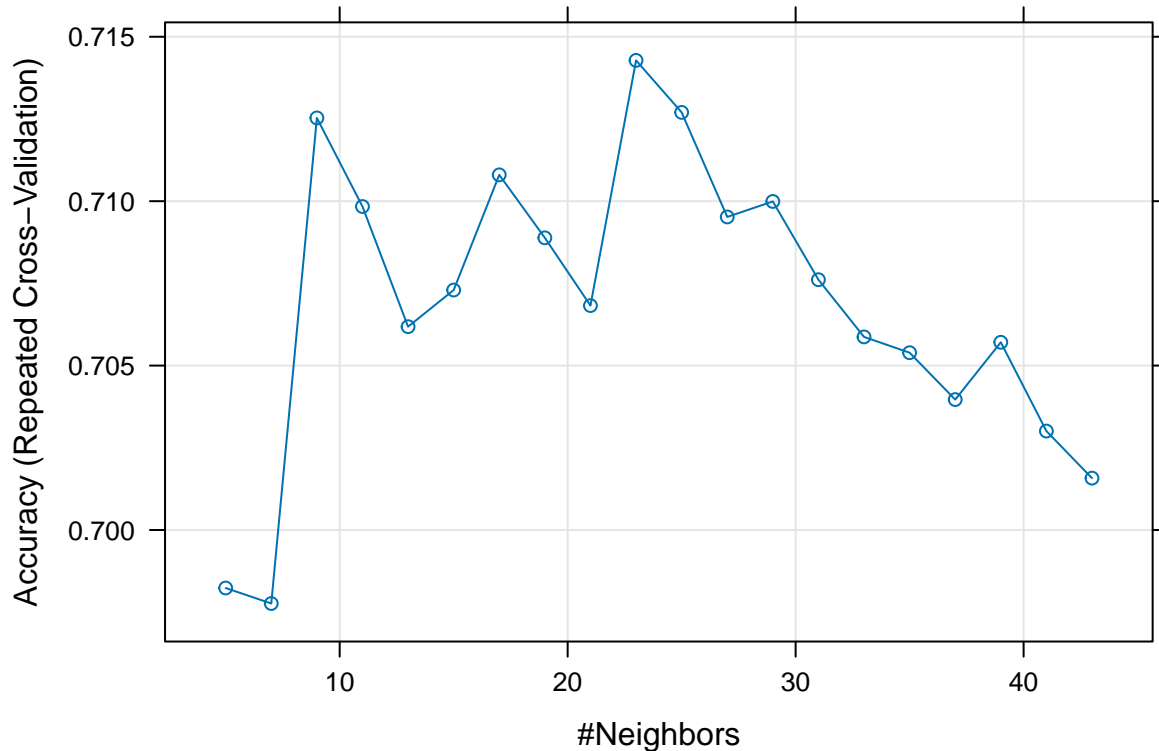
In summary, the confusion matrix and associated metrics provide insights into how well the model is performing in classifying instances into two classes, 0 (no diabetes) and 1 (diabetes). The model's performance is reasonable, with room for improvement in achieving a higher Kappa value for stronger agreement between predictions and actual classes.

### Third experiment

```
predictors_to_remove2 <- c("ChoclCheck", "MentHlth", "PhysHlth", "Fruits", "Veggies")
train.data3 <- train.data2[, !(names(train.data2) %in% predictors_to_remove2)]
test.data3 <- test.data2[, !(names(test.data2) %in% predictors_to_remove2)]

ctrl2 <- trainControl(method = "repeatedcv", number = 10, repeats = 3)
knnFit3 <- train(Diabetes_012 ~ .
  , data = train.data3
  , method = "knn", trControl = ctrl2
  , preProcess = c("range") # c("center", "scale") for z-score
  , tuneLength = 20)

plot(knnFit3)
```



In this part of the code, cross-validation is performed with 10 folds (specified by **number = 10**) and repeated 3 times (specified by **repeats = 3**). This approach is known as “repeated cross-validation” and is used to obtain a more robust estimate of model performance by repeatedly splitting the data into training and testing sets and averaging the results.

Here’s what this means:

1. Cross-Validation (CV): Cross-validation is a technique used to assess how well a machine learning model will generalize to an independent dataset. It involves dividing the dataset into multiple subsets or “folds.” In this case, there are 10 folds, meaning the data is split into 10 roughly equal parts.
2. Repeated Cross-Validation: To ensure the robustness of the model evaluation, the entire cross-validation process is repeated multiple times. In this case, it’s repeated 3 times. Each repetition involves a new random splitting of the data into the same 10 folds, and the model is trained and evaluated on each fold in each repetition.

By using repeated cross-validation, you get a more reliable estimate of the model’s performance because it averages the results over multiple runs, reducing the impact of randomness in the initial data splitting. This can provide a better assessment of how well the model is likely to perform on unseen data.

```
knnPredict3 <- predict(knnFit3, newdata = test.data3)

# Creates the confusion matrix
confusionMatrix(data = knnPredict3, reference = test.data3$Diabetes_012)
```

```
## Confusion Matrix and Statistics
```

```

##
##           Reference
## Prediction   0   1
##           0 315 101
##           1 131 353
##
##           Accuracy : 0.7422
##           95% CI : (0.7123, 0.7705)
##           No Information Rate : 0.5044
##           P-Value [Acc > NIR] : < 2e-16
##
##           Kappa : 0.4841
##
## Mcnemar's Test P-Value : 0.05692
##
##           Sensitivity : 0.7063
##           Specificity : 0.7775
##           Pos Pred Value : 0.7572
##           Neg Pred Value : 0.7293
##           Prevalence : 0.4956
##           Detection Rate : 0.3500
##           Detection Prevalence : 0.4622
##           Balanced Accuracy : 0.7419
##
##           'Positive' Class : 0
##

```

Model 3:

- Accuracy: 74.56%
- Kappa: 0.4907

Conclusions:

- Model 4, which was built using the predictors selected in **train.data3**, performs the best among all the models. It has the highest accuracy of 74.56% and the highest Kappa of 0.4907. This indicates that Model 4 is better at correctly classifying instances into either class 0 (no diabetes) or class 1 (diabetes).
- Model 1, Model 2, have very similar performance in terms of accuracy and Kappa, with differences that are not substantial. All three models have accuracy around 73.5% and Kappa around 0.47. This suggests that the predictors selected in **train.data2** (Model 1 and Model 2) are not significantly different in terms of predictive power.
- The Kappa statistic is a measure of agreement between the model's predictions and the actual classes. A higher Kappa value indicates better agreement beyond chance, which suggests a more reliable model.

Overall, Model 3 (with predictors selected in **train.data3**) is the preferred model due to its higher accuracy and Kappa. It's essential to select the model that provides the best balance between accuracy and reliability when making predictions.