# Using Machine Learning to Diagnose Diabetes

Alan Feria

Sept, 17 2o23

## Contents

### 0. Abstract

Diabetes is more prevalent among the Pima Indian tribe of southern Arizona than any other population globally(1). This paper demonstrates the accurate diagnosis of diabetes in Pima women by optimizing parameters for three machine learning algorithms: Random Forest decision trees, a Support Vector Machine (SVM), and a Multilayer Perceptron (MLP) artificial neural network (ANN). These results underscore the potential of machine learning as a valuable tool for assisting doctors in disease diagnosis.

## 1. Introduction

This study explores the application of machine learning in the field of medicine. Can statistical learning algorithms aid doctors in making accurate diagnoses and informed decisions about patient health? To answer this question, we trained three distinct machine learning algorithms—Random Forest decision trees, a Support Vector Machine (SVM), and a Multilayer Perceptron (MLP) artificial neural network—using data on diabetes occurrences in Pima Indian Women. All analyses were conducted using the R statistical computing software package within the RStudio integrated development environment. Our findings indicate that all three algorithms achieved a high success rate in correctly diagnosing diabetes in Pima women, providing support for the hypothesis that machine learning has a role in a doctor's toolkit.

### 1.1 About the Data

The data on diabetes incidence in Pima women used in this analysis are freely available from the University of California, Irvine Machine Learning Repository at http://archive.ics.uci.edu/ml/. This dataset was created to predict the onset of diabetes in women from the Pima Indian Tribe, containing 768 observations of nine features. The features, as they appear in the dataset, are as follows:

| Feature Description | Feature Name in Dataset |
| --- | --- |
| Number of times pregnant | timesPregnant |
| Plasma glucose concentration at 2 hours in an oral glucose tolerance test | plasmaGlucose |
| Diastolic blood pressure (mm Hg) | diastolicPressure |
| Triceps skin fold thickness (mm) | tricepThickness |
| 2-Hour serum insulin (mu U/ml) | serumInsulin |
| Body mass index (weight in kg/(height in m)^2) | bmi |
| Diabetes pedigree function | pedigreeFunction |
| Age (years) | age |
| Class variable (diabetic or not diabetic) | diabetes |

A comprehensive description of this dataset is available at: https://www.ncbi.nlm.nih.gov/pmc/articles/PMC2245318/

### 1.2 Software Used in Analysis

All analyses were performed using the R statistical computing environment, with extended functionality provided by the following software libraries: `caret`, `mice`, `randomForest`, `kernlab`, `VIM`, and `RSNNS`. These libraries support data imputation and the three machine learning algorithms used in this analysis.

The analysis was conducted within the RStudio integrated development environment, and this report was generated using the `rmarkdown` library.

### 1.3 Outline of Analysis

The analysis workflow included data loading, exploratory analysis, missing data imputation, data standardization, splitting data into training and test sets, model construction, parameter optimization, and model evaluation.

## 2. Data Processing

### 2.1 Data Acquisition

The analysis started with importing the data into R and loading the required R packages.

```r
library(caret)
library(mice)
library(randomForest)
library(kernlab)
library(RSNNS)
library(VIM)

#Create a vector of the feature names
headers <- c("timesPregnant", "plasmaGlucose", "diastolicPressure", "tricepThickness",
             "serumInsulin", "bmi", "pedigreeFunction", "age", "diabetes")

#Import data
```

```
#www <- paste0("http://archive.ics.uci.edu/ml/machine-learning-databases/",
              #"pima-indians-diabetes/pima-indians-diabetes.data")
library(readr)
data <-  read_csv("diabetes.csv")

colnames(data) <- c("timesPregnant", "plasmaGlucose", "diastolicPressure", "tricepThickness",
                    "serumInsulin", "bmi", "pedigreeFunction", "age", "diabetes")
```

The data structure was examined to ensure correct import.

```
str(data)
```

```
## spc_tbl_ [768 x 9] (S3: spec_tbl_df/tbl_df/tbl/data.frame)
##  $ timesPregnant    : num [1:768] 6 1 8 1 0 5 3 10 2 8 ...
##  $ plasmaGlucose    : num [1:768] 148 85 183 89 137 116 78 115 197 125 ...
##  $ diastolicPressure: num [1:768] 72 66 64 66 40 74 50 0 70 96 ...
##  $ tricepThickness  : num [1:768] 35 29 0 23 35 0 32 0 45 0 ...
##  $ serumInsulin     : num [1:768] 0 0 0 94 168 0 88 0 543 0 ...
##  $ bmi              : num [1:768] 33.6 26.6 23.3 28.1 43.1 25.6 31 35.3 30.5 0 ...
##  $ pedigreeFunction : num [1:768] 0.627 0.351 0.672 0.167 2.288 ...
##  $ age              : num [1:768] 50 31 32 21 33 30 26 29 53 54 ...
##  $ diabetes         : num [1:768] 1 0 1 0 1 0 1 0 1 1 ...
##  - attr(*, "spec")=
##   .. cols(
##   ..   Pregnancies = col_double(),
##   ..   Glucose = col_double(),
##   ..   BloodPressure = col_double(),
##   ..   SkinThickness = col_double(),
##   ..   Insulin = col_double(),
##   ..   BMI = col_double(),
##   ..   DiabetesPedigreeFunction = col_double(),
##   ..   Age = col_double(),
##   ..   Outcome = col_double()
##   .. )
##  - attr(*, "problems")=<externalptr>
```

Subsequently, we encoded the class label of the `diabetes` feature as 'notDiabetic' for 0 and 'Diabetic' for 1, converting it into a factor.

```
data$diabetes <- as.factor(ifelse(data$diabetes == 0, "notDiabetic", "Diabetic"))
```

**2.2 Exploratory Analysis**

Following data import, an exploratory analysis was conducted, starting with the creation of a scatterplot matrix.

```
pairs(data)
```

The scatterplot matrix revealed several features with 0-valued observations.

### 2.3 Missing Data

The exploratory analysis highlighted features with 0-valued observations. For certain features, a 0 value is biologically implausible, particularly in `plasmaGlucose`, `diastolicPressure`, `tricepThickness`, `serumInsulin`, and `bmi`. Although the dataset did not explicitly contain missing values, these implicit missing values encoded as 0s were explicitly encoded as `NA`.

```
for (i in 2:6) {
  for (n in 1:nrow(data)) {
    if (data[n, i] == 0) {
      data[n, i] <- NA
    }
  }
}
```

An aggregation plot was then constructed to count the number of missing values.

```
aggr(data[, 2:6], cex.lab = 1, cex.axis = .5, numbers = TRUE, gap = 0)
```

The left plot displayed the proportion of missing values to total observations for each feature, showing that over half of all `serumInsulin` and nearly a third of `tricepThickness` observations were missing.

The right plot indicated that only slightly over half of the observations were complete. Due to the significant number of observations with missing values, it was decided not to remove them from the analysis. Instead, missing values were imputed using Imputation by Predicted Mean Matching, a method that imputes missing values by finding the nearest-neighbor donor based on the expected values of the missing variables conditional on the observed covariates.

Imputation by Predicted Mean Matching was chosen for its capability to provide valid inference when data are missing at random, an assumption tested by examining relationships between variables and missing values through a scatterplot matrix.

```
scattmatrixMiss(data)
```

No discernible relationship between variables and missing values was observed, validating the assumption of data missing at random. Imputation by Predicted Mean Matching was then applied.

```
tempdata <- mice(data, m = 3, method = 'pmm', seed = 100)
```

```
##
##  iter imp variable
##   1   1  plasmaGlucose  diastolicPressure  tricepThickness  serumInsulin  bmi
##   1   2  plasmaGlucose  diastolicPressure  tricepThickness  serumInsulin  bmi
##   1   3  plasmaGlucose  diastolicPressure  tricepThickness  serumInsulin  bmi
##   2   1  plasmaGlucose  diastolicPressure  tricepThickness  serumInsulin  bmi
##   2   2  plasmaGlucose  diastolicPressure  tricepThickness  serumInsulin  bmi
##   2   3  plasmaGlucose  diastolicPressure  tricepThickness  serumInsulin  bmi
##   3   1  plasmaGlucose  diastolicPressure  tricepThickness  serumInsulin  bmi
##   3   2  plasmaGlucose  diastolicPressure  tricepThickness  serumInsulin  bmi
##   3   3  plasmaGlucose  diastolicPressure  tricepThickness  serumInsulin  bmi
##   4   1  plasmaGlucose  diastolicPressure  tricepThickness  serumInsulin  bmi
##   4   2  plasmaGlucose  diastolicPressure  tricepThickness  serumInsulin  bmi
##   4   3  plasmaGlucose  diastolicPressure  tricepThickness  serumInsulin  bmi
##   5   1  plasmaGlucose  diastolicPressure  tricepThickness  serumInsulin  bmi
##   5   2  plasmaGlucose  diastolicPressure  tricepThickness  serumInsulin  bmi
##   5   3  plasmaGlucose  diastolicPressure  tricepThickness  serumInsulin  bmi
```

```
data <- complete(tempdata)
```

Distributions of the imputed data were compared to the original data.

```
densityplot(tempdata)
```

The distributions were found to be approximately equal, allowing for the continuation of feature selection.

### 2.4 Feature Selection

To ensure

that no two features were highly correlated, a correlation matrix was constructed. Features with a correlation coefficient greater than 0.7 were considered for removal.

```r
correlationMatrix <- cor(data[, 1:8])
findCorrelation(correlationMatrix, cutoff = 0.7)
```

```
## integer(0)
```

As no two features exhibited such high correlation, it was determined that no features needed to be eliminated from the analysis.

### 2.5 Data Standardization

Since the features in the dataset represented different units, data standardization was performed by adjusting the mean of each column to 0 and the standard deviation to 1.

```r
data[, 1:8] <- scale(data[, 1:8], center = TRUE, scale = TRUE)
```

## 3. Training the Algorithms

With data processing complete, the analysis proceeded to model construction. Models for each of the three methods were built using a training set of observations with 10-fold cross-validation to prevent overfitting, facilitated by the caret package.

```r
tenFoldCV <- trainControl(method = "repeatedcv",
                          number = 10,
                          repeats = 10,
```

```
                              classProbs = TRUE,
                              summaryFunction = twoClassSummary)
```

The best model for each method was selected based on the area under the ROC curve (AUC), representing model accuracy. A high AUC score, close to 1, indicated high accuracy.

### 3.1 Predicting the Most Common Label

Of the 768 observations, approximately 65% were labeled 'notDiabetic,' and the remaining 35% were labeled 'Diabetic.' Predicting 'notDiabetic' for every instance would already achieve 65% accuracy. Therefore, any model needed to exceed this accuracy threshold to be considered viable.

### 3.2 Training and Test Sets

Seventy percent of the observations were used to train the models, while the remaining 30% were reserved for testing.

```
sampleSize <- floor(.7 * nrow(data))

set.seed(131)
trainIndices <- sample(seq_len(nrow(data)), size = sampleSize)

x.train <- data[trainIndices, 1:8]
y.train <- data[trainIndices, 9]

x.test <- data[-trainIndices, 1:8]
y.test <- data[-trainIndices, 9]
```

### 3.3 Random Forest

Random Forest, a tree-based classification method, was employed in the analysis. It randomly selected a subset of predictor variables as split candidates at each tree node. Models were trained for different random sample sizes (denoted as `mtry`).

```
rf.expand <- expand.grid(mtry = 2:8)

set.seed(100)
rf <- caret::train(x.train,
                   y.train,
                   method = "rf",
                   metric = "ROC",
                   trControl = tenFoldCV,
                   tuneGrid = rf.expand)

rf
```

```
## Random Forest
##
## 537 samples
##   8 predictor
##   2 classes: 'Diabetic', 'notDiabetic'
##
## No pre-processing
## Resampling: Cross-Validated (10 fold, repeated 10 times)
## Summary of sample sizes: 484, 483, 482, 484, 484, 483, ...
## Resampling results across tuning parameters:
```

```
##
##    mtry  ROC         Sens        Spec
##    2     0.8375074   0.5878070   0.8627857
##    3     0.8345722   0.5790643   0.8568730
##    4     0.8320079   0.5774269   0.8545952
##    5     0.8287727   0.5812281   0.8554365
##    6     0.8272897   0.5818421   0.8517778
##    7     0.8238832   0.5762865   0.8506825
##    8     0.8211825   0.5736842   0.8484206
##
## ROC was used to select the optimal model using the largest value.
## The final value used for the model was mtry = 2.
```

The best model was achieved when considering only 2 randomly chosen predictors. The importance of variables was visualized.

```
varImpPlot(rf$finalModel, type = 2, main = "Random Forest")
```



The plot listed the importance of each variable to the model, indicating that all variables appeared significant, and no feature elimination or retraining was required. Thus, the Random Forest model with `mtry = 2` was chosen as the final model.

**3.4 Support Vector Machine**

The Support Vector Machine (SVM) aimed to linearly separate observations based on class labels. A linear kernel was initially employed, and models were trained with different cost (C) values.

```
linear.svm.expand <- expand.grid(C = c(.1, 1, 10))
set.seed(131)
linear.svm <- caret::train(x.train,
                           y.train,
                           method = "svmLinear",
```

```
                          metric = "ROC",
                          trControl = tenFoldCV,
                          tuneLength = 10,
                          tuneGrid = linear.svm.expand)
linear.svm
```

```
## Support Vector Machines with Linear Kernel
##
## 537 samples
##   8 predictor
##   2 classes: 'Diabetic', 'notDiabetic'
##
## No pre-processing
## Resampling: Cross-Validated (10 fold, repeated 10 times)
## Summary of sample sizes: 484, 484, 483, 484, 483, 483, ...
## Resampling results across tuning parameters:
##
##   C      ROC        Sens       Spec
##    0.1   0.8356491  0.5433333  0.8754127
##    1.0   0.8348555  0.5461404  0.8757063
##   10.0   0.8343674  0.5471930  0.8771032
##
## ROC was used to select the optimal model using the largest value.
## The final value used for the model was C = 0.1.
```
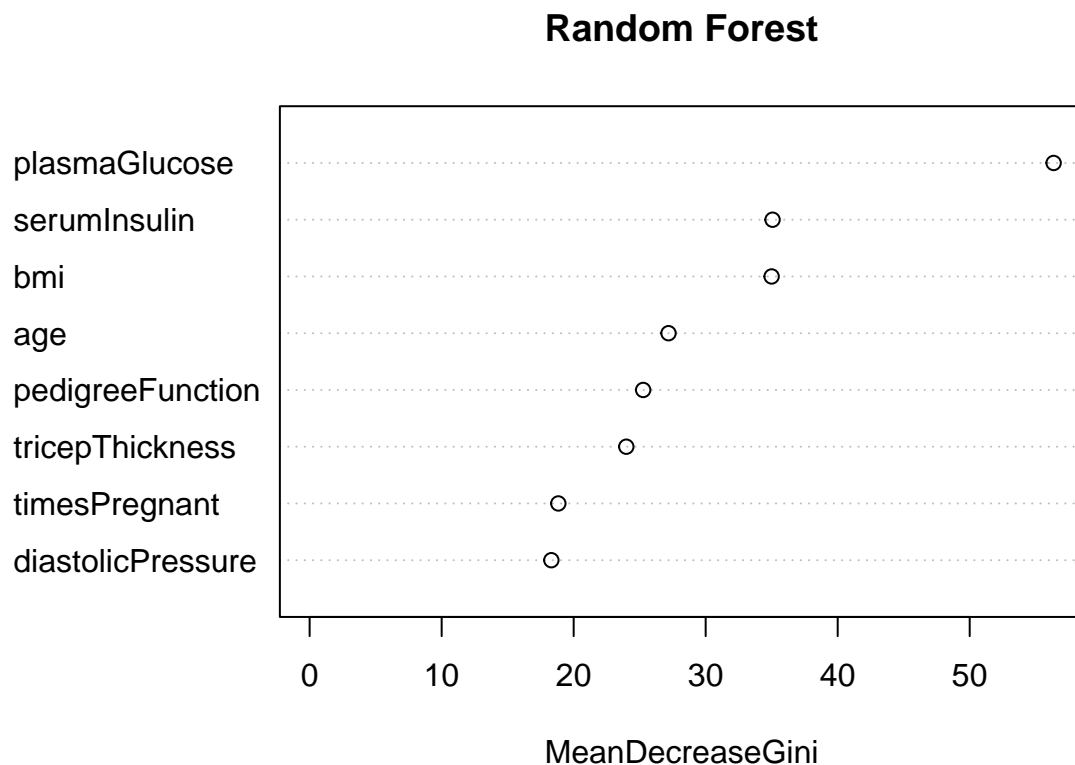
After finding that `C = 0.1` produced the best model, a narrower search was conducted.

```
linear.svm.expand2 <- expand.grid(C = c(.05, .1, .15))
set.seed(131)
linear.svm2 <- caret::train(x.train,
                            y.train,
                            method = "svmLinear",
                            metric = "ROC",
                            trControl = tenFoldCV,
                            tuneGrid = linear.svm.expand2)
linear.svm2
```

```
## Support Vector Machines with Linear Kernel
##
## 537 samples
##   8 predictor
##   2 classes: 'Diabetic', 'notDiabetic'
##
## No pre-processing
## Resampling: Cross-Validated (10 fold, repeated 10 times)
## Summary of sample sizes: 484, 484, 483, 484, 483, 483, ...
## Resampling results across tuning parameters:
##
##   C     ROC        Sens       Spec
##   0.05  0.8375866  0.5406140  0.8751349
##   0.10  0.8356491  0.5405556  0.8740159
##   0.15  0.8357497  0.5427193  0.8779603
##
## ROC was used to select the optimal model using the largest value.
## The final value used for the model was C = 0.05.
```

The final model with an AUC of 0.891 was achieved with `C = 0.05`. Subsequently, a Support Vector Machine with a radial basis function kernel was constructed for comparison.

```
radial.svm.expand <- expand.grid(sigma = c(.2, .4, .6, .8),
                                 C = c(.1, 1, 5, 10, 100))
set.seed(131)
radial.svm <- caret::train(x.train,
                           y.train,
                           method = "svmRadial",
                           metric = "ROC",
                           trControl = tenFoldCV,
                           tuneGrid = radial.svm.expand)
radial.svm
```

```
## Support Vector Machines with Radial Basis Function Kernel
##
## 537 samples
##    8 predictor
##    2 classes: 'Diabetic', 'notDiabetic'
##
## No pre-processing
## Resampling: Cross-Validated (10 fold, repeated 10 times)
## Summary of sample sizes: 484, 484, 483, 484, 483, 483, ...
## Resampling results across tuning parameters:
##
##    sigma  C      ROC        Sens       Spec
##    0.2      0.1  0.8112930  0.6571345  0.8087143
##    0.2      1.0  0.8040346  0.5141520  0.8658571
##    0.2      5.0  0.7757416  0.4461404  0.8728016
##    0.2     10.0  0.7535821  0.3968713  0.8733413
##    0.2    100.0  0.7252477  0.3424854  0.8521270
##    0.4      0.1  0.7988646  0.5971930  0.8233175
##    0.4      1.0  0.7823711  0.4652632  0.8653651
##    0.4      5.0  0.7400638  0.3856725  0.8608810
##    0.4     10.0  0.7320281  0.3455263  0.8583492
##    0.4    100.0  0.7328023  0.3517836  0.8524683
##    0.6      0.1  0.7831319  0.5356725  0.8347143
##    0.6      1.0  0.7757746  0.4537427  0.8685556
##    0.6      5.0  0.7387671  0.3604971  0.8572857
##    0.6     10.0  0.7367051  0.3539474  0.8516984
##    0.6    100.0  0.7380561  0.3508480  0.8593016
##    0.8      0.1  0.7726678  0.4722807  0.8536667
##    0.8      1.0  0.7697746  0.4319298  0.8693730
##    0.8      5.0  0.7391022  0.3488889  0.8623492
##    0.8     10.0  0.7395639  0.3461988  0.8601111
##    0.8    100.0  0.7395639  0.3473099  0.8603810
##
## ROC was used to select the optimal model using the largest value.
## The final values used for the model were sigma = 0.2 and C = 0.1.
```

A narrower search was conducted around `sigma = 0.2` and `C = 0.1`.

```
radial.svm.expand2 <- expand.grid(sigma = c(.15, .2, .25),
                                  C = c(.01, .05, .1, .15, .25))
set.seed(131)
radial.svm2 <- caret::train(x.train,
```

```
                            y.train,
                            method = "svmRadial",
                            metric = "ROC",
                            trControl = tenFoldCV,
                            tuneGrid = radial.svm.expand2)
radial.svm2
```

```
## Support Vector Machines with Radial Basis Function Kernel
##
## 537 samples
##    8 predictor
##    2 classes: 'Diabetic', 'notDiabetic'
##
## No pre-processing
## Resampling: Cross-Validated (10 fold, repeated 10 times)
## Summary of sample sizes: 484, 484, 483, 484, 483, 483, ...
## Resampling results across tuning parameters:
##
##   sigma  C     ROC        Sens       Spec
##   0.15   0.01  0.8159053  0.6638596  0.8101429
##   0.15   0.05  0.8158322  0.6621053  0.8092937
##   0.15   0.10  0.8157709  0.6681579  0.8098492
##   0.15   0.15  0.8151147  0.6111111  0.8360714
##   0.15   0.25  0.8107846  0.5627485  0.8554841
##   0.20   0.01  0.8110942  0.6577193  0.8100873
##   0.20   0.05  0.8112947  0.6511696  0.8115000
##   0.20   0.10  0.8112930  0.6560526  0.8115397
##   0.20   0.15  0.8108263  0.6385380  0.8208571
##   0.20   0.25  0.8084707  0.5593567  0.8492540
##   0.25   0.01  0.8086369  0.6467251  0.8132063
##   0.25   0.05  0.8086970  0.6472222  0.8104206
##   0.25   0.10  0.8086820  0.6416374  0.8135079
##   0.25   0.15  0.8086828  0.6411404  0.8140794
##   0.25   0.25  0.8070392  0.5758480  0.8410873
##
## ROC was used to select the optimal model using the largest value.
## The final values used for the model were sigma = 0.15 and C = 0.01.
```

The final SVM model with the radial basis function kernel achieved an AUC of 0.864 with `sigma = 0.15`
and `C = 0.01`. As this AUC was lower than that of the linear kernel, the SVM model with the linear kernel
and `C = 0.05` was chosen as the final Support Vector Machine model.

**3.5 Multilayer Perceptron Artificial Neural Network**

The Multilayer Perceptron Artificial Neural Network (MLP ANN) processes data through layers of artificial
neurons, each weighting inputs to reach a final decision. Models with varying numbers of hidden layers were
trained.

```
set.seed(131)
mlpnn <- caret::train(x.train,
                      y.train,
                      method = "mlpML",
                      metric = "ROC",
                      trControl = tenFoldCV)
```

```
## Warning: At least one layer had zero units and were removed. The new structure
```

```
## is 1

## Warning: At least one layer had zero units and were removed. The new structure
## is 3

## Warning: At least one layer had zero units and were removed. The new structure
## is 5

## Warning: At least one layer had zero units and were removed. The new structure
## is 1

## Warning: At least one layer had zero units and were removed. The new structure
## is 3

## Warning: At least one layer had zero units and were removed. The new structure
## is 5

## Warning: At least one layer had zero units and were removed. The new structure
## is 1

## Warning: At least one layer had zero units and were removed. The new structure
## is 3

## Warning: At least one layer had zero units and were removed. The new structure
## is 5

## Warning: At least one layer had zero units and were removed. The new structure
## is 1

## Warning: At least one layer had zero units and were removed. The new structure
## is 3

## Warning: At least one layer had zero units and were removed. The new structure
## is 5

## Warning: At least one layer had zero units and were removed. The new structure
## is 1

## Warning: At least one layer had zero units and were removed. The new structure
## is 3

## Warning: At least one layer had zero units and were removed. The new structure
## is 5

## Warning: At least one layer had zero units and were removed. The new structure
## is 1

## Warning: At least one layer had zero units and were removed. The new structure
## is 3

## Warning: At least one layer had zero units and were removed. The new structure
## is 5

## Warning: At least one layer had zero units and were removed. The new structure
## is 1

## Warning: At least one layer had zero units and were removed. The new structure
## is 3

## Warning: At least one layer had zero units and were removed. The new structure
## is 5

## Warning: At least one layer had zero units and were removed. The new structure
## is 1
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 3
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 5
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 1
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 3
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 5
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 1
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 3
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 5
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 1
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 3
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 5
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 1
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 3
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 5
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 1
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 3
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 5
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 1
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 3
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 5
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 1
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 3
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 5
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 1
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 3
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 5
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 1
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 3
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 5
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 1
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 3
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 5
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 1
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 3
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 5
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 1
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 3
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 5
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 1
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 3
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 5
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 1
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 3
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 5
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 1
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 3
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 5
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 1
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 3
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 5
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 1
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 3
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 5
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 1
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 3
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 5
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 1
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 3
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 5
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 1
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 3
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 5
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 1
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 3

## Warning: At least one layer had zero units and were removed. The new structure
## is 5

## Warning: At least one layer had zero units and were removed. The new structure
## is 1

## Warning: At least one layer had zero units and were removed. The new structure
## is 3

## Warning: At least one layer had zero units and were removed. The new structure
## is 5

## Warning: At least one layer had zero units and were removed. The new structure
## is 1

## Warning: At least one layer had zero units and were removed. The new structure
## is 3

## Warning: At least one layer had zero units and were removed. The new structure
## is 5

## Warning: At least one layer had zero units and were removed. The new structure
## is 1

## Warning: At least one layer had zero units and were removed. The new structure
## is 3

## Warning: At least one layer had zero units and were removed. The new structure
## is 5

## Warning: At least one layer had zero units and were removed. The new structure
## is 1

## Warning: At least one layer had zero units and were removed. The new structure
## is 3

## Warning: At least one layer had zero units and were removed. The new structure
## is 5

## Warning: At least one layer had zero units and were removed. The new structure
## is 1

## Warning: At least one layer had zero units and were removed. The new structure
## is 3

## Warning: At least one layer had zero units and were removed. The new structure
## is 5

## Warning: At least one layer had zero units and were removed. The new structure
## is 1

## Warning: At least one layer had zero units and were removed. The new structure
## is 3

## Warning: At least one layer had zero units and were removed. The new structure
## is 5

## Warning: At least one layer had zero units and were removed. The new structure
## is 1
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 3
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 5
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 1
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 3
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 5
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 1
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 3
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 5
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 1
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 3
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 5
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 1
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 3
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 5
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 1
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 3
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 5
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 1
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 3
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 5
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 1
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 3

## Warning: At least one layer had zero units and were removed. The new structure
## is 5

## Warning: At least one layer had zero units and were removed. The new structure
## is 1

## Warning: At least one layer had zero units and were removed. The new structure
## is 3

## Warning: At least one layer had zero units and were removed. The new structure
## is 5

## Warning: At least one layer had zero units and were removed. The new structure
## is 1

## Warning: At least one layer had zero units and were removed. The new structure
## is 3

## Warning: At least one layer had zero units and were removed. The new structure
## is 5

## Warning: At least one layer had zero units and were removed. The new structure
## is 1

## Warning: At least one layer had zero units and were removed. The new structure
## is 3

## Warning: At least one layer had zero units and were removed. The new structure
## is 5

## Warning: At least one layer had zero units and were removed. The new structure
## is 1

## Warning: At least one layer had zero units and were removed. The new structure
## is 3

## Warning: At least one layer had zero units and were removed. The new structure
## is 5

## Warning: At least one layer had zero units and were removed. The new structure
## is 1

## Warning: At least one layer had zero units and were removed. The new structure
## is 3

## Warning: At least one layer had zero units and were removed. The new structure
## is 5

## Warning: At least one layer had zero units and were removed. The new structure
## is 1

## Warning: At least one layer had zero units and were removed. The new structure
## is 3

## Warning: At least one layer had zero units and were removed. The new structure
## is 5

## Warning: At least one layer had zero units and were removed. The new structure
## is 1
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 3

## Warning: At least one layer had zero units and were removed. The new structure
## is 5

## Warning: At least one layer had zero units and were removed. The new structure
## is 1

## Warning: At least one layer had zero units and were removed. The new structure
## is 3

## Warning: At least one layer had zero units and were removed. The new structure
## is 5

## Warning: At least one layer had zero units and were removed. The new structure
## is 1

## Warning: At least one layer had zero units and were removed. The new structure
## is 3

## Warning: At least one layer had zero units and were removed. The new structure
## is 5

## Warning: At least one layer had zero units and were removed. The new structure
## is 1

## Warning: At least one layer had zero units and were removed. The new structure
## is 3

## Warning: At least one layer had zero units and were removed. The new structure
## is 5

## Warning: At least one layer had zero units and were removed. The new structure
## is 1

## Warning: At least one layer had zero units and were removed. The new structure
## is 3

## Warning: At least one layer had zero units and were removed. The new structure
## is 5

## Warning: At least one layer had zero units and were removed. The new structure
## is 1

## Warning: At least one layer had zero units and were removed. The new structure
## is 3

## Warning: At least one layer had zero units and were removed. The new structure
## is 5

## Warning: At least one layer had zero units and were removed. The new structure
## is 1

## Warning: At least one layer had zero units and were removed. The new structure
## is 3

## Warning: At least one layer had zero units and were removed. The new structure
## is 5

## Warning: At least one layer had zero units and were removed. The new structure
## is 1
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 3

## Warning: At least one layer had zero units and were removed. The new structure
## is 5

## Warning: At least one layer had zero units and were removed. The new structure
## is 1

## Warning: At least one layer had zero units and were removed. The new structure
## is 3

## Warning: At least one layer had zero units and were removed. The new structure
## is 5

## Warning: At least one layer had zero units and were removed. The new structure
## is 1

## Warning: At least one layer had zero units and were removed. The new structure
## is 3

## Warning: At least one layer had zero units and were removed. The new structure
## is 5

## Warning: At least one layer had zero units and were removed. The new structure
## is 1

## Warning: At least one layer had zero units and were removed. The new structure
## is 3

## Warning: At least one layer had zero units and were removed. The new structure
## is 5

## Warning: At least one layer had zero units and were removed. The new structure
## is 1

## Warning: At least one layer had zero units and were removed. The new structure
## is 3

## Warning: At least one layer had zero units and were removed. The new structure
## is 5

## Warning: At least one layer had zero units and were removed. The new structure
## is 1

## Warning: At least one layer had zero units and were removed. The new structure
## is 3

## Warning: At least one layer had zero units and were removed. The new structure
## is 5

## Warning: At least one layer had zero units and were removed. The new structure
## is 1

## Warning: At least one layer had zero units and were removed. The new structure
## is 3

## Warning: At least one layer had zero units and were removed. The new structure
## is 5

## Warning: At least one layer had zero units and were removed. The new structure
## is 1
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 3

## Warning: At least one layer had zero units and were removed. The new structure
## is 5

## Warning: At least one layer had zero units and were removed. The new structure
## is 1

## Warning: At least one layer had zero units and were removed. The new structure
## is 3

## Warning: At least one layer had zero units and were removed. The new structure
## is 5

## Warning: At least one layer had zero units and were removed. The new structure
## is 1

## Warning: At least one layer had zero units and were removed. The new structure
## is 3

## Warning: At least one layer had zero units and were removed. The new structure
## is 5

## Warning: At least one layer had zero units and were removed. The new structure
## is 1

## Warning: At least one layer had zero units and were removed. The new structure
## is 3

## Warning: At least one layer had zero units and were removed. The new structure
## is 5

## Warning: At least one layer had zero units and were removed. The new structure
## is 1

## Warning: At least one layer had zero units and were removed. The new structure
## is 3

## Warning: At least one layer had zero units and were removed. The new structure
## is 5

## Warning: At least one layer had zero units and were removed. The new structure
## is 1

## Warning: At least one layer had zero units and were removed. The new structure
## is 3

## Warning: At least one layer had zero units and were removed. The new structure
## is 5

## Warning: At least one layer had zero units and were removed. The new structure
## is 1

## Warning: At least one layer had zero units and were removed. The new structure
## is 3

## Warning: At least one layer had zero units and were removed. The new structure
## is 5

## Warning: At least one layer had zero units and were removed. The new structure
## is 1
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 3
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 5
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 1
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 3
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 5
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 1
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 3
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 5
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 1
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 3
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 5
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 1
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 3
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 5
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 1
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 3
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 5
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 1
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 3
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 5
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 1
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 3
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 5
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 1
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 3
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 5
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 1
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 3
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 5
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 1
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 3
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 5
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 1
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 3
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 5
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 1
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 3
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 5
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 1
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 3
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 5
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 1
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 3
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 5
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 1
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 3
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 5
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 1
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 3
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 5
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 1
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 3
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 5
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 1
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 3
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 5
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 1
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 3
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 5
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 1
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 3
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 5
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 1
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 3

## Warning: At least one layer had zero units and were removed. The new structure
## is 5

## Warning: At least one layer had zero units and were removed. The new structure
## is 1

## Warning: At least one layer had zero units and were removed. The new structure
## is 3

## Warning: At least one layer had zero units and were removed. The new structure
## is 5

## Warning: At least one layer had zero units and were removed. The new structure
## is 1

## Warning: At least one layer had zero units and were removed. The new structure
## is 3

## Warning: At least one layer had zero units and were removed. The new structure
## is 5

## Warning: At least one layer had zero units and were removed. The new structure
## is 1

## Warning: At least one layer had zero units and were removed. The new structure
## is 3

## Warning: At least one layer had zero units and were removed. The new structure
## is 5

## Warning: At least one layer had zero units and were removed. The new structure
## is 1

## Warning: At least one layer had zero units and were removed. The new structure
## is 3

## Warning: At least one layer had zero units and were removed. The new structure
## is 5

## Warning: At least one layer had zero units and were removed. The new structure
## is 1

## Warning: At least one layer had zero units and were removed. The new structure
## is 3

## Warning: At least one layer had zero units and were removed. The new structure
## is 5

## Warning: At least one layer had zero units and were removed. The new structure
## is 1

## Warning: At least one layer had zero units and were removed. The new structure
## is 3

## Warning: At least one layer had zero units and were removed. The new structure
## is 5

## Warning: At least one layer had zero units and were removed. The new structure
## is 1
```

```
## Warning: At least one layer had zero units and were removed. The new structure
## is 3

## Warning: At least one layer had zero units and were removed. The new structure
## is 5

## Warning: At least one layer had zero units and were removed. The new structure
## is 1

## Warning: At least one layer had zero units and were removed. The new structure
## is 3

## Warning: At least one layer had zero units and were removed. The new structure
## is 5

## Warning: At least one layer had zero units and were removed. The new structure
## is 1
```

```
mlpnn
```

```
## Multi-Layer Perceptron, with multiple layers
##
## 537 samples
##    8 predictor
##    2 classes: 'Diabetic', 'notDiabetic'
##
## No pre-processing
## Resampling: Cross-Validated (10 fold, repeated 10 times)
## Summary of sample sizes: 484, 484, 483, 484, 483, 483, ...
## Resampling results across tuning parameters:
##
##    layer1  ROC         Sens        Spec
##    1       0.8321527   0.6528655   0.8196905
##    3       0.8229064   0.5819883   0.8494921
##    5       0.8130643   0.5851170   0.8250079
##
## Tuning parameter 'layer2' was held constant at a value of 0
## Tuning
##  parameter 'layer3' was held constant at a value of 0
## ROC was used to select the optimal model using the largest value.
## The final values used for the model were layer1 = 1, layer2 = 0 and layer3 = 0.
```

The best performance was achieved with a single hidden layer, making it the final model for this method.

## 4. Selecting the Best Overall Model

Three final models were obtained, one for each method. The model with the highest classification accuracy on the test set was chosen as the best model for this problem.

```
# Random Forest Test Accuracy44
rf.predict <- predict(rf$finalModel, x.test)
rf.test.accuracy <- mean(rf.predict == y.test)
rf.test.accuracy
```

```
## [1] 0.7619048
```

```
# SVM Test Accuracy
svm.predict <- predict(linear.svm2$finalModel, x.test)
svm.test.accuracy <- mean(svm.predict == y.test)
```

```
svm.test.accuracy
```

```
## [1] 0.7575758
```

```
# MLPNN Test Accuracy
mlpnn.predict <- predict(mlpnn$finalModel, x.test)
mlpnn.predict <- as.data.frame(mlpnn.predict)
mlpnn.predict$prediction <- ifelse(mlpnn.predict$V1 >= .5, "Diabetic", "notDiabetic")
mlpnn.test.accuracy <- mean(mlpnn.predict$prediction == y.test)
mlpnn.test.accuracy
```

```
## [1] 0.7575758
```

The Random Forest achieved the highest classification accuracy on the test set, making it the best model for this problem.

## 5. Conclusion

Despite a small and incomplete training set, data mining demonstrated its viability for diabetes diagnosis in Pima Indian women. Of the three algorithms tested, all surpassed the minimum viable accuracy rate. The Random Forest algorithm outperformed the others and was selected as the best model. Further improvements are expected with more data free of missing values. Despite these challenges, machine learning proved to be a valuable addition to the medical industry.

## 6. References

(1) Diabetes Incidence in Pima Indians

(2) Semicontinuous Longitudinal Data

(3) Missing Not at Random Data

(4) James et al., "Introduction to Statistical Learning"

(5) Random Forest Variable Importance

(6) Multilayer Perceptron

**Software:** - R Project - RStudio - caret Package - mice Package - randomForest Package - kernlab Package - RSNNS Package - VIM Package