

# Practical Machine Learning - Course Project

José Manuel Mirás-Avalos

24 May 2020

## 1. Overview

This report documents the solution for the Peer Assessment Project from the Coursera's course on *Practical Machine Learning*, as part of the **Specialization in Data Science**. The report has been built as a *markdown* file using **R Studio**, meant to be published online in *html* format.

The main objective of this project was to predict the manner in which 6 participants carried out a given physical exercise, further described below. This is the “classe” variable in the training data set. Three machine learning algorithms were tested in the current analysis and the most accurate one was selected for being applied to the testing data set. The predictions obtained will be submitted in appropriate format to the Course Project Prediction Quiz for automated grading.

## 2. Background

Nowadays, it is possible to collect a large amount of data about personal activity relatively inexpensively using devices such as *Jawbone Up*, *Nike FuelBand*, and *Fitbit*. These type of devices are part of the quantified self movement – a group of enthusiasts who take measurements about themselves regularly for several reasons such as to improve their health, to find patterns in their behavior, or because they are tech geeks. For instance, one thing that people do is quantify how much of a particular activity they do, but they rarely quantify how well they do it.

In this context, the goal of this project is to use data from accelerometers on the belt, forearm, arm, and dumbbell of 6 participants. They were asked to perform barbell lifts correctly and incorrectly in 5 different ways. Further information can be found at: <http://groupware.les.inf.puc-rio.br/har> (<http://groupware.les.inf.puc-rio.br/har>) (see the section on the Weight Lifting Exercise Dataset).

Read more about this subject: <http://groupware.les.inf.puc-rio.br/har#ixzz3xsbS5bVX>  
(<http://groupware.les.inf.puc-rio.br/har#ixzz3xsbS5bVX>)

## 3. Loading data and exploratory analysis

### 3.1. Dataset Overview

The training data for this project are available here: <https://d396qusza40orc.cloudfront.net/predmachlearn/pml-training.csv> (<https://d396qusza40orc.cloudfront.net/predmachlearn/pml-training.csv>)

The test data are available here: <https://d396qusza40orc.cloudfront.net/predmachlearn/pml-testing.csv> (<https://d396qusza40orc.cloudfront.net/predmachlearn/pml-testing.csv>)

The data for this project come from this source: <http://groupware.les.inf.puc-rio.br/har> (<http://groupware.les.inf.puc-rio.br/har>).

This data has been gathered by Velloso et al. (2013), to which I am grateful for allowing the free use of their dataset in this assignment.

In the authors' website, a short description of the data collected can be found:

"Six young health participants were asked to perform one set of 10 repetitions of the Unilateral Dumbbell Biceps Curl in five different fashions: exactly according to the specification (Class A), throwing the elbows to the front (Class B), lifting the dumbbell only halfway (Class C), lowering the dumbbell only halfway (Class D) and throwing the hips to the front (Class E).

Class A corresponds to the specified execution of the exercise, while the other 4 classes correspond to common mistakes. Participants were supervised by an experienced weight lifter to make sure the execution complied to the manner they were supposed to simulate. The exercises were performed by six male participants aged between 20-28 years, with little weight lifting experience. We made sure that all participants could easily simulate the mistakes in a safe and controlled manner by using a relatively light dumbbell (1.25kg)."

## 3.2. Preparation of the R environment

First, the libraries required for completing the analyses in this assignment are loaded into R:

```
rm(list=ls())                # free up memory for the download of the data sets
library(knitr)
library(caret)
```

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

```
## Loading required package: ggplot2
```

```
library(rpart)
library(rpart.plot)
library(rattle)
```

```
## Rattle: A free graphical interface for data science with R.
## Versión 5.3.0 Copyright (c) 2006-2018 Togaware Pty Ltd.
## Escriba 'rattle()' para agitar, sacudir y rotar sus datos.
```

```
library(randomForest)
```

```
## randomForest 4.6-14
```

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

```
##
## Attaching package: 'randomForest'
```

```
## The following object is masked from 'package:rattle':
##
##     importance
```

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

```
library(corrplot)
```

```
## corrplot 0.84 loaded
```

```
set.seed(12345)
```

### 3.3. Loading and cleaning data

Once the R environment has been set up, the datasets (both for training and testing) are downloaded from the URLs provided in the assignment instructions. Then, the training dataset is partitioned into a training set (70% of the data) for constructing the model and a testing set (30% of the data) for validating the model. The downloaded testing dataset is not modified and will only be used for the quiz results generation.

```
# Set URL for downloading datasets
UrlTrain <- "http://d396qusza40orc.cloudfront.net/predmachlearn/pml-training.csv"
UrlTest  <- "http://d396qusza40orc.cloudfront.net/predmachlearn/pml-testing.csv"

# Downloading datasets and assigning them to R data.frames
training <- read.csv(url(UrlTrain))
testing  <- read.csv(url(UrlTest))

# Create a partition with the training dataset
inTrain <- createDataPartition(training$classe, p=0.7, list=FALSE)
TrainSet <- training[inTrain, ]
TestSet  <- training[-inTrain, ]
dim(TrainSet)
```

```
## [1] 13737 160
```

```
dim(TestSet)
```

```
## [1] 5885 160
```

Both, the training and testing sets have 160 variables and some of them have a large number of NA, so they can be removed with the procedures in the code chunk below. The Near Zero Variance (NZV) variables are also removed, as well as the ID variables.

```
# Remove variables with Nearly Zero Variance (NZV)
NZV <- nearZeroVar(TrainSet)
TrainSet <- TrainSet[, -NZV]
TestSet  <- TestSet[, -NZV]
dim(TrainSet)
```

```
## [1] 13737 104
```

```
dim(TestSet)
```

```
## [1] 5885 104
```

```
# Remove variables that are mostly NA
AllNA <- sapply(TrainSet, function(x) mean(is.na(x))) > 0.95
TrainSet <- TrainSet[, AllNA==FALSE]
TestSet <- TestSet[, AllNA==FALSE]
dim(TrainSet)
```

```
## [1] 13737 59
```

```
dim(TestSet)
```

```
## [1] 5885 59
```

```
# Remove identification only variables (columns 1 to 5)
TrainSet <- TrainSet[, -(1:5)]
TestSet <- TestSet[, -(1:5)]
dim(TrainSet)
```

```
## [1] 13737 54
```

```
dim(TestSet)
```

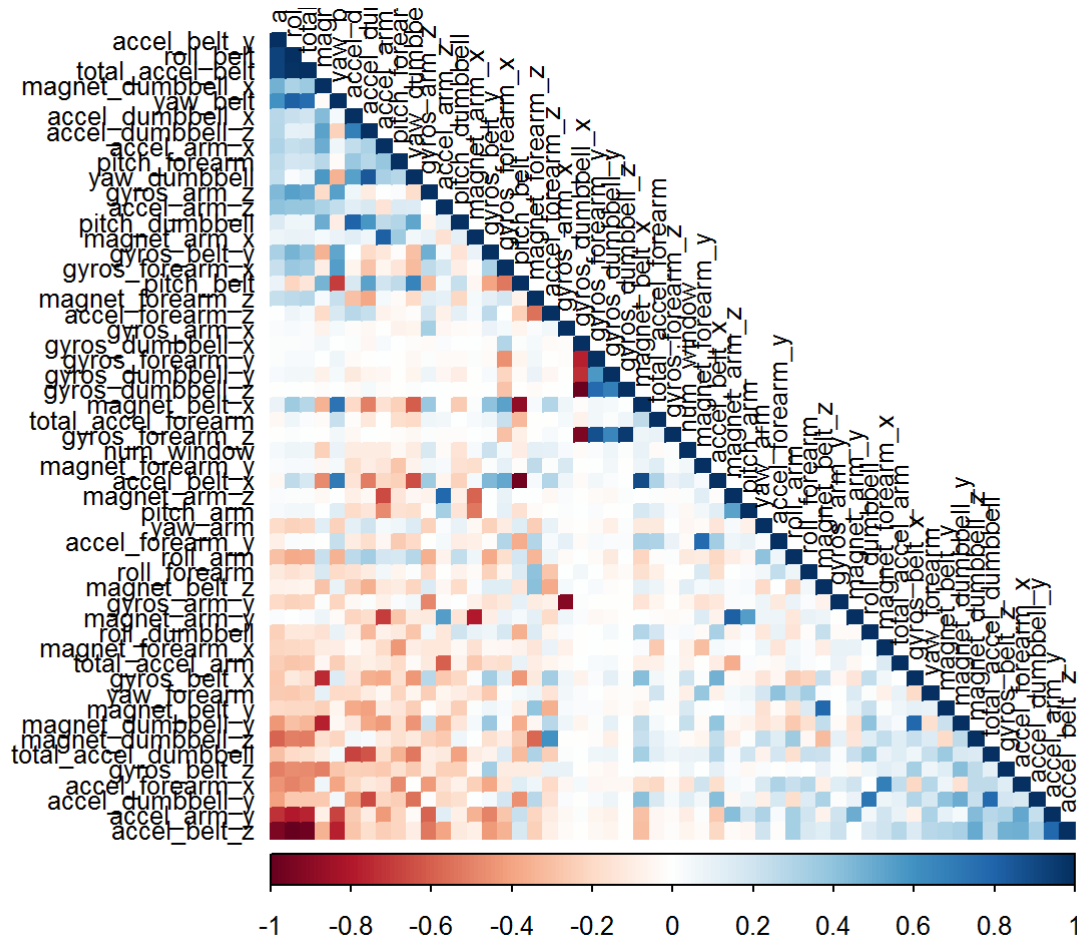
```
## [1] 5885 54
```

Using this cleaning procedure, the number of variables within the dataset has been reduced to 54.

## 3.4. Correlations among variables

In order to check if some variables within the dataset are significantly correlated among them, the following procedure was used:

```
corMatrix <- cor(TrainSet[, -54])
corrplot(corMatrix, order = "FPC", method = "color", type = "lower",
         tl.cex = 0.8, tl.col = rgb(0, 0, 0))
```



The figure shows the highly correlated variables in dark colours. These correlations are quite few and no further analysis was performed.

## 4. Building a prediction model

In order to construct the most accurate model for predicting the class of movement which they use for performing the physical exercise requested, three methods were tested in the training dataset: *Random Forest*, *Classification Trees* and *Generalized Boosted Model*.

Each method is described below and a confusion matrix is plotted at the end of each analysis for a visualization of the accuracy of each model.

### 4.1. Random Forest

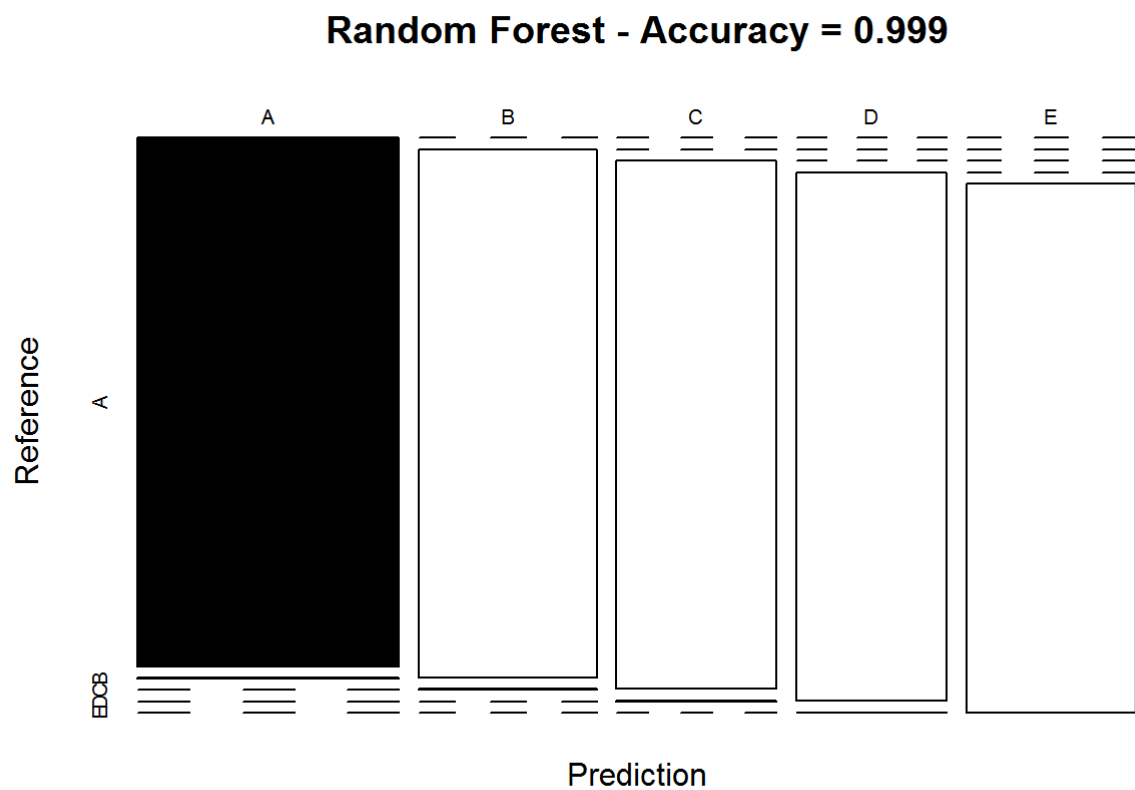
```
# model fit
set.seed(12345)
controlRF <- trainControl(method="cv", number=3, verboseIter=FALSE)
modFitRandForest <- train(classe ~ ., data=TrainSet, method="rf",
                           trControl=controlRF)
modFitRandForest$finalModel
```

```
##
## Call:
##  randomForest(x = x, y = y, mtry = param$mtry)
##              Type of random forest: classification
##              Number of trees: 500
## No. of variables tried at each split: 27
##
##              OOB estimate of  error rate: 0.23%
## Confusion matrix:
##      A    B    C    D    E  class.error
## A 3904     2     0     0     0 0.0005120328
## B   6 2647     4     1     0 0.0041384500
## C   0   5 2391     0     0 0.0020868114
## D   0   0   9 2243     0 0.0039964476
## E   0   0   0   5 2520 0.0019801980
```

```
# prediction on Test dataset
predictRandForest <- predict(modFitRandForest, newdata=TestSet)
confMatRandForest <- confusionMatrix(predictRandForest, TestSet$classe)
confMatRandForest
```

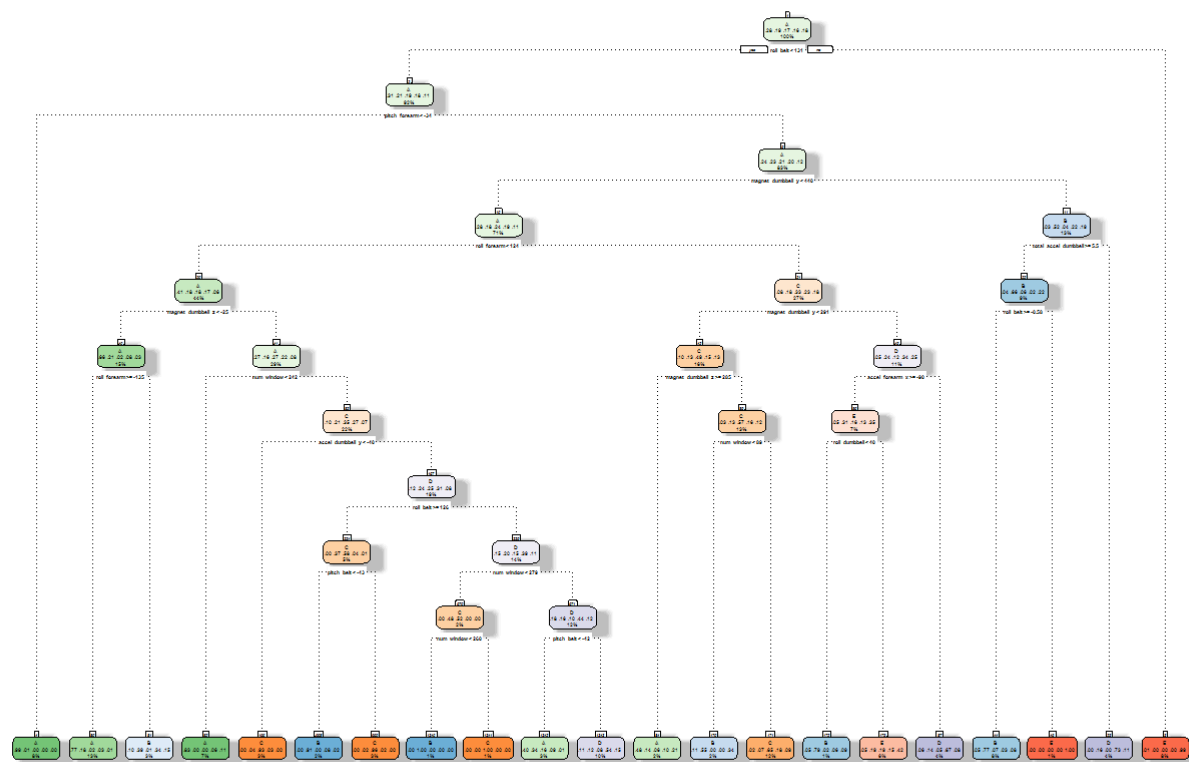
```
## Confusion Matrix and Statistics
##
##              Reference
## Prediction    A    B    C    D    E
##      A 1674     1     0     0     0
##      B   0 1138     2     0     0
##      C   0   0 1024     2     0
##      D   0   0   0 962     1
##      E   0   0   0   0 1081
##
## Overall Statistics
##
##              Accuracy : 0.999
##              95% CI : (0.9978, 0.9996)
##      No Information Rate : 0.2845
##      P-Value [Acc > NIR] : < 2.2e-16
##
##              Kappa : 0.9987
##
##      McNemar's Test P-Value : NA
##
## Statistics by Class:
##
##              Class: A Class: B Class: C Class: D Class: E
## Sensitivity          1.0000   0.9991   0.9981   0.9979   0.9991
## Specificity          0.9998   0.9996   0.9996   0.9998   1.0000
## Pos Pred Value       0.9994   0.9982   0.9981   0.9990   1.0000
## Neg Pred Value       1.0000   0.9998   0.9996   0.9996   0.9998
## Prevalence           0.2845   0.1935   0.1743   0.1638   0.1839
## Detection Rate       0.2845   0.1934   0.1740   0.1635   0.1837
## Detection Prevalence 0.2846   0.1937   0.1743   0.1636   0.1837
## Balanced Accuracy     0.9999   0.9994   0.9988   0.9989   0.9995
```

```
# plot matrix results
plot(confMatRandForest$table, col = confMatRandForest$byClass,
     main = paste("Random Forest - Accuracy =",
                   round(confMatRandForest$overall['Accuracy'], 4)))
```



## 4.2. Classification Trees

```
# model fit
set.seed(12345)
modFitDecTree <- rpart(classe ~ ., data=TrainSet, method="class")
fancyRpartPlot(modFitDecTree)
```



Rattle 2020-may-24 13:37:22 Josema

*# prediction on Test dataset*

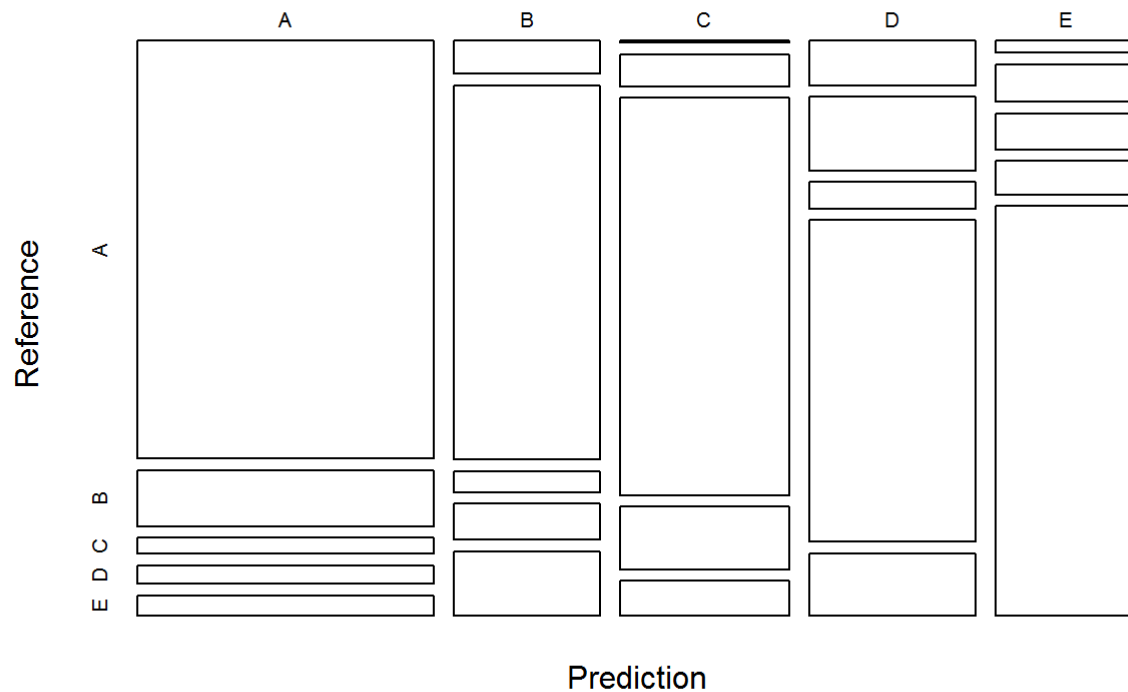
```
predictDecTree <- predict(modFitDecTree, newdata=TestSet, type="class")
confMatDecTree <- confusionMatrix(predictDecTree, TestSet$classe)
confMatDecTree
```



```
## Confusion Matrix and Statistics
##
##           Reference
## Prediction   A    B    C    D    E
##           A 1502  201   59   66   74
##           B   58 660   37   64  114
##           C    4  66 815  129   72
##           D   90 148  54 648  126
##           E   20  64  61  57 696
##
## Overall Statistics
##
##           Accuracy : 0.7342
##           95% CI : (0.7228, 0.7455)
##           No Information Rate : 0.2845
##           P-Value [Acc > NIR] : < 2.2e-16
##
##           Kappa : 0.6625
##
## Mcnemar's Test P-Value : < 2.2e-16
##
## Statistics by Class:
##
##           Class: A Class: B Class: C Class: D Class: E
## Sensitivity      0.8973  0.5795  0.7943  0.6722  0.6433
## Specificity      0.9050  0.9425  0.9442  0.9151  0.9579
## Pos Pred Value   0.7897  0.7074  0.7505  0.6079  0.7751
## Neg Pred Value   0.9568  0.9033  0.9560  0.9344  0.9226
## Prevalence       0.2845  0.1935  0.1743  0.1638  0.1839
## Detection Rate   0.2552  0.1121  0.1385  0.1101  0.1183
## Detection Prevalence 0.3232  0.1585  0.1845  0.1811  0.1526
## Balanced Accuracy 0.9011  0.7610  0.8693  0.7936  0.8006
```

```
# plot matrix results
plot(confMatDecTree$table, col = confMatDecTree$byClass,
     main = paste("Decision Tree - Accuracy =",
                  round(confMatDecTree$overall['Accuracy'], 4)))
```

## Decision Tree - Accuracy = 0.7342



### 4.3. Generalized Boosted Model

```
# model fit
set.seed(12345)
controlGBM <- trainControl(method = "repeatedcv", number = 5, repeats = 1)
modFitGBM <- train(classe ~ ., data=TrainSet, method = "gbm",
                   trControl = controlGBM, verbose = FALSE)
modFitGBM$finalModel
```

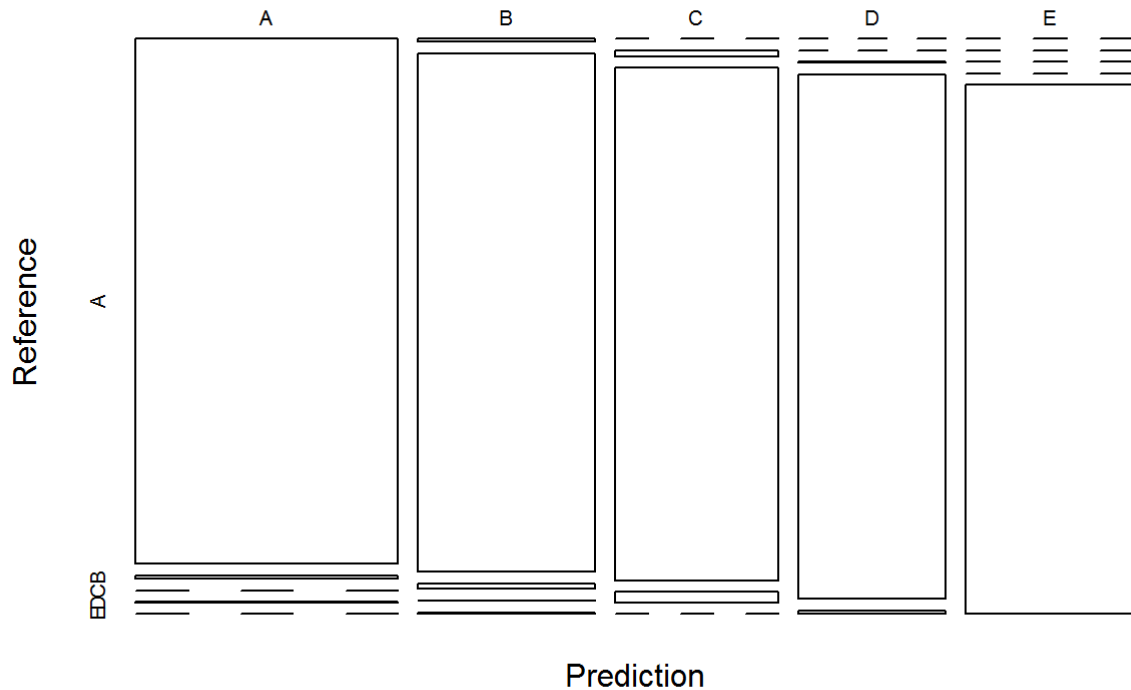
```
## A gradient boosted model with multinomial loss function.
## 150 iterations were performed.
## There were 53 predictors of which 53 had non-zero influence.
```

```
# prediction on Test dataset
predictGBM <- predict(modFitGBM, newdata=TestSet)
confMatGBM <- confusionMatrix(predictGBM, TestSet$classe)
confMatGBM
```

```
## Confusion Matrix and Statistics
##
##           Reference
## Prediction   A    B    C    D    E
##           A 1668   12    0    1    0
##           B    6 1115   12    1    3
##           C    0   12 1012   21    0
##           D    0    0    2  941    6
##           E    0    0    0    0 1073
##
## Overall Statistics
##
##           Accuracy : 0.9871
##           95% CI : (0.9839, 0.9898)
##           No Information Rate : 0.2845
##           P-Value [Acc > NIR] : < 2.2e-16
##
##           Kappa : 0.9837
##
## Mcnemar's Test P-Value : NA
##
## Statistics by Class:
##
##           Class: A Class: B Class: C Class: D Class: E
## Sensitivity      0.9964  0.9789  0.9864  0.9761  0.9917
## Specificity      0.9969  0.9954  0.9932  0.9984  1.0000
## Pos Pred Value   0.9923  0.9807  0.9684  0.9916  1.0000
## Neg Pred Value   0.9986  0.9949  0.9971  0.9953  0.9981
## Prevalence       0.2845  0.1935  0.1743  0.1638  0.1839
## Detection Rate   0.2834  0.1895  0.1720  0.1599  0.1823
## Detection Prevalence 0.2856  0.1932  0.1776  0.1613  0.1823
## Balanced Accuracy 0.9967  0.9871  0.9898  0.9873  0.9958
```

```
# plot matrix results
plot(confMatGBM$table, col = confMatGBM$byClass,
     main = paste("GBM - Accuracy =", round(confMatGBM$overall['Accuracy'], 4)))
```

## GBM - Accuracy = 0.9871



## 5. Application of the selected model to the test data

The previously fitted models had the following accuracies:

- Random Forest: 0.999
- Classification Trees: 0.7342
- Generalized Boosted Model: 0.9871

Therefore, the Random Forest model was selected for applying to the Test set and predict the 20 quiz results as shown in this code:

```
predictTEST <- predict(modFitRandForest, newdata=testing)
predictTEST
```

```
## [1] B A B A A E D B A A B C B A E E A B B B
## Levels: A B C D E
```

## 6. Reference

Velloso, E., Bulling, A., Gellersen, H., Ugulino, W., Fuks, H. (2013). Qualitative activity recognition of weight lifting exercises. *Proceedings of 4th International Conference in Cooperation with SIGCHI (Augmented Human '13)*. Stuttgart, Germany: ACM SIGCHI.