# Prediction Assignment Writeup

## BarthLdlg

## 27/04/2020

## Overview

Using devices such as Jawbone Up, Nike FuelBand, and Fitbit it is now possible to collect a large amount of data about personal activity relatively inexpensively. These type of devices are part of the quantified self movement – a group of enthusiasts who take measurements about themselves regularly to improve their health, to find patterns in their behavior, or because they are tech geeks. One thing that people regularly do is quantify how much of a particular activity they do, but they rarely quantify how well they do it. In this project, your goal will be to use data from accelerometers on the belt, forearm, arm, and dumbell of 6 participants. They were asked to perform barbell lifts correctly and incorrectly in 5 different ways. More information is available from the website here: http://web.archive.org/web/20161224072740/http://groupware.les.inf.puc-rio.br/har (see the section on the Weight Lifting Exercise Dataset).

The goal of the project is to predict the manner in which they did the exercise. This is the "classe" variable in the training set.

## Data and library loading

We load the libraries we need for the project:

```
library(caret)
```

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

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

```
library(gbm)
```

```
## Loaded gbm 2.1.5
```

```
library(ggplot2)
library(rpart)
library(lattice)
library(e1071)
library(AppliedPredictiveModeling)
library(rattle)
```

```
## Rattle: A free graphical interface for data science with R.
## Version 5.3.0 Copyright (c) 2006-2018 Togaware Pty Ltd.
## Entrez 'rattle()' pour secouer, faire vibrer, et faire défiler vos données.
```

We load the train data and the test data into R:

```
trainData <- read.csv("pml-training.csv")
testData <- read.csv("pml-testing.csv")
```

Let's see the dimension of the training and test data:

```
dim(trainData)
```

```
## [1] 19622   160
```

```
dim(testData)
```

```
## [1]  20 160
```

A quick summary:

```
str(trainData)
```

```
## 'data.frame':    19622 obs. of  160 variables:
##  $ X                       : int  1 2 3 4 5 6 7 8 9 10 ...
##  $ user_name               : Factor w/ 6 levels "adelmo","carlitos",..: 2 2 2 2 2 2 2 2 2 2 ...
##  $ raw_timestamp_part_1    : int  1323084231 1323084231 1323084231 1323084232 1323084232 1323084232
##  $ raw_timestamp_part_2    : int  788290 808298 820366 120339 196328 304277 368296 440390 484323 484
##  $ cvtd_timestamp          : Factor w/ 20 levels "02/12/2011 13:32",..: 9 9 9 9 9 9 9 9 9 9 ...
##  $ new_window              : Factor w/ 2 levels "no","yes": 1 1 1 1 1 1 1 1 1 1 ...
##  $ num_window              : int  11 11 11 12 12 12 12 12 12 12 ...
##  $ roll_belt               : num  1.41 1.41 1.42 1.48 1.48 1.45 1.42 1.42 1.43 1.45 ...
##  $ pitch_belt              : num  8.07 8.07 8.07 8.05 8.07 8.06 8.09 8.13 8.16 8.17 ...
##  $ yaw_belt                : num  -94.4 -94.4 -94.4 -94.4 -94.4 -94.4 -94.4 -94.4 -94.4 -94.4 ...
##  $ total_accel_belt        : int  3 3 3 3 3 3 3 3 3 3 ...
##  $ kurtosis_roll_belt      : Factor w/ 397 levels "","-0.016850",..: 1 1 1 1 1 1 1 1 1 1 ...
##  $ kurtosis_picth_belt     : Factor w/ 317 levels "","-0.021887",..: 1 1 1 1 1 1 1 1 1 1 ...
##  $ kurtosis_yaw_belt       : Factor w/ 2 levels "","#DIV/0!": 1 1 1 1 1 1 1 1 1 1 ...
##  $ skewness_roll_belt      : Factor w/ 395 levels "","-0.003095",..: 1 1 1 1 1 1 1 1 1 1 ...
##  $ skewness_roll_belt.1    : Factor w/ 338 levels "","-0.005928",..: 1 1 1 1 1 1 1 1 1 1 ...
##  $ skewness_yaw_belt       : Factor w/ 2 levels "","#DIV/0!": 1 1 1 1 1 1 1 1 1 1 ...
##  $ max_roll_belt           : num  NA NA NA NA NA NA NA NA NA NA ...
##  $ max_picth_belt          : int  NA NA NA NA NA NA NA NA NA NA ...
##  $ max_yaw_belt            : Factor w/ 68 levels "","-0.1","-0.2",..: 1 1 1 1 1 1 1 1 1 1 ...
##  $ min_roll_belt           : num  NA NA NA NA NA NA NA NA NA NA ...
##  $ min_pitch_belt          : int  NA NA NA NA NA NA NA NA NA NA ...
##  $ min_yaw_belt            : Factor w/ 68 levels "","-0.1","-0.2",..: 1 1 1 1 1 1 1 1 1 1 ...
##  $ amplitude_roll_belt     : num  NA NA NA NA NA NA NA NA NA NA ...
##  $ amplitude_pitch_belt    : int  NA NA NA NA NA NA NA NA NA NA ...
##  $ amplitude_yaw_belt      : Factor w/ 4 levels "","#DIV/0!","0.00",..: 1 1 1 1 1 1 1 1 1 1 ...
##  $ var_total_accel_belt    : num  NA NA NA NA NA NA NA NA NA NA ...
##  $ avg_roll_belt           : num  NA NA NA NA NA NA NA NA NA NA ...
##  $ stddev_roll_belt        : num  NA NA NA NA NA NA NA NA NA NA ...
##  $ var_roll_belt           : num  NA NA NA NA NA NA NA NA NA NA ...
##  $ avg_pitch_belt          : num  NA NA NA NA NA NA NA NA NA NA ...
```

```
##  $ stddev_pitch_belt      : num   NA NA NA NA NA NA NA NA NA NA ...
##  $ var_pitch_belt         : num   NA NA NA NA NA NA NA NA NA NA ...
##  $ avg_yaw_belt           : num   NA NA NA NA NA NA NA NA NA NA ...
##  $ stddev_yaw_belt        : num   NA NA NA NA NA NA NA NA NA NA ...
##  $ var_yaw_belt           : num   NA NA NA NA NA NA NA NA NA NA ...
##  $ gyros_belt_x           : num   0 0.02 0 0.02 0.02 0.02 0.02 0.02 0.02 0.03 ...
##  $ gyros_belt_y           : num   0 0 0 0 0.02 0 0 0 0 0 ...
##  $ gyros_belt_z           : num   -0.02 -0.02 -0.02 -0.03 -0.02 -0.02 -0.02 -0.02 -0.02 0 ...
##  $ accel_belt_x           : int   -21 -22 -20 -22 -21 -21 -22 -22 -20 -21 ...
##  $ accel_belt_y           : int   4 4 5 3 2 4 3 4 2 4 ...
##  $ accel_belt_z           : int   22 22 23 21 24 21 21 21 24 22 ...
##  $ magnet_belt_x          : int   -3 -7 -2 -6 -6 0 -4 -2 1 -3 ...
##  $ magnet_belt_y          : int   599 608 600 604 600 603 599 603 602 609 ...
##  $ magnet_belt_z          : int   -313 -311 -305 -310 -302 -312 -311 -313 -312 -308 ...
##  $ roll_arm               : num   -128 -128 -128 -128 -128 -128 -128 -128 -128 -128 ...
##  $ pitch_arm              : num   22.5 22.5 22.5 22.1 22.1 22 21.9 21.8 21.7 21.6 ...
##  $ yaw_arm                : num   -161 -161 -161 -161 -161 -161 -161 -161 -161 -161 ...
##  $ total_accel_arm        : int   34 34 34 34 34 34 34 34 34 34 ...
##  $ var_accel_arm          : num   NA NA NA NA NA NA NA NA NA NA ...
##  $ avg_roll_arm           : num   NA NA NA NA NA NA NA NA NA NA ...
##  $ stddev_roll_arm        : num   NA NA NA NA NA NA NA NA NA NA ...
##  $ var_roll_arm           : num   NA NA NA NA NA NA NA NA NA NA ...
##  $ avg_pitch_arm          : num   NA NA NA NA NA NA NA NA NA NA ...
##  $ stddev_pitch_arm       : num   NA NA NA NA NA NA NA NA NA NA ...
##  $ var_pitch_arm          : num   NA NA NA NA NA NA NA NA NA NA ...
##  $ avg_yaw_arm            : num   NA NA NA NA NA NA NA NA NA NA ...
##  $ stddev_yaw_arm         : num   NA NA NA NA NA NA NA NA NA NA ...
##  $ var_yaw_arm            : num   NA NA NA NA NA NA NA NA NA NA ...
##  $ gyros_arm_x            : num   0 0.02 0.02 0.02 0 0.02 0 0.02 0.02 0.02 ...
##  $ gyros_arm_y            : num   0 -0.02 -0.02 -0.03 -0.03 -0.03 -0.03 -0.02 -0.03 -0.03 ...
##  $ gyros_arm_z            : num   -0.02 -0.02 -0.02 0.02 0 0 0 0 -0.02 -0.02 ...
##  $ accel_arm_x            : int   -288 -290 -289 -289 -289 -289 -289 -289 -288 -288 ...
##  $ accel_arm_y            : int   109 110 110 111 111 111 111 111 109 110 ...
##  $ accel_arm_z            : int   -123 -125 -126 -123 -123 -122 -125 -124 -122 -124 ...
##  $ magnet_arm_x           : int   -368 -369 -368 -372 -374 -369 -373 -372 -369 -376 ...
##  $ magnet_arm_y           : int   337 337 344 344 337 342 336 338 341 334 ...
##  $ magnet_arm_z           : int   516 513 513 512 506 513 509 510 518 516 ...
##  $ kurtosis_roll_arm      : Factor w/ 330 levels "","-0.02438",..: 1 1 1 1 1 1 1 1 1 1 ...
##  $ kurtosis_picth_arm     : Factor w/ 328 levels "","-0.00484",..: 1 1 1 1 1 1 1 1 1 1 ...
##  $ kurtosis_yaw_arm       : Factor w/ 395 levels "","-0.01548",..: 1 1 1 1 1 1 1 1 1 1 ...
##  $ skewness_roll_arm      : Factor w/ 331 levels "","-0.00051",..: 1 1 1 1 1 1 1 1 1 1 ...
##  $ skewness_pitch_arm     : Factor w/ 328 levels "","-0.00184",..: 1 1 1 1 1 1 1 1 1 1 ...
##  $ skewness_yaw_arm       : Factor w/ 395 levels "","-0.00311",..: 1 1 1 1 1 1 1 1 1 1 ...
##  $ max_roll_arm           : num   NA NA NA NA NA NA NA NA NA NA ...
##  $ max_picth_arm          : num   NA NA NA NA NA NA NA NA NA NA ...
##  $ max_yaw_arm            : int   NA NA NA NA NA NA NA NA NA NA ...
##  $ min_roll_arm           : num   NA NA NA NA NA NA NA NA NA NA ...
##  $ min_pitch_arm          : num   NA NA NA NA NA NA NA NA NA NA ...
##  $ min_yaw_arm            : int   NA NA NA NA NA NA NA NA NA NA ...
##  $ amplitude_roll_arm     : num   NA NA NA NA NA NA NA NA NA NA ...
##  $ amplitude_pitch_arm    : num   NA NA NA NA NA NA NA NA NA NA ...
##  $ amplitude_yaw_arm      : int   NA NA NA NA NA NA NA NA NA NA ...
##  $ roll_dumbbell          : num   13.1 13.1 12.9 13.4 13.4 ...
##  $ pitch_dumbbell         : num   -70.5 -70.6 -70.3 -70.4 -70.4 ...
```

```
##  $ yaw_dumbbell          : num  -84.9 -84.7 -85.1 -84.9 -84.9 ...
##  $ kurtosis_roll_dumbbell : Factor w/ 398 levels "","-0.0035","-0.0073",..: 1 1 1 1 1 1 1 1 1 1 ...
##  $ kurtosis_picth_dumbbell : Factor w/ 401 levels "","-0.0163","-0.0233",..: 1 1 1 1 1 1 1 1 1 1 ...
##  $ kurtosis_yaw_dumbbell  : Factor w/ 2 levels "","#DIV/0!": 1 1 1 1 1 1 1 1 1 1 ...
##  $ skewness_roll_dumbbell : Factor w/ 401 levels "","-0.0082","-0.0096",..: 1 1 1 1 1 1 1 1 1 1 ...
##  $ skewness_pitch_dumbbell : Factor w/ 402 levels "","-0.0053","-0.0084",..: 1 1 1 1 1 1 1 1 1 1 ...
##  $ skewness_yaw_dumbbell  : Factor w/ 2 levels "","#DIV/0!": 1 1 1 1 1 1 1 1 1 1 ...
##  $ max_roll_dumbbell      : num  NA NA NA NA NA NA NA NA NA NA ...
##  $ max_picth_dumbbell     : num  NA NA NA NA NA NA NA NA NA NA ...
##  $ max_yaw_dumbbell       : Factor w/ 73 levels "","-0.1","-0.2",..: 1 1 1 1 1 1 1 1 1 1 ...
##  $ min_roll_dumbbell      : num  NA NA NA NA NA NA NA NA NA NA ...
##  $ min_pitch_dumbbell     : num  NA NA NA NA NA NA NA NA NA NA ...
##  $ min_yaw_dumbbell       : Factor w/ 73 levels "","-0.1","-0.2",..: 1 1 1 1 1 1 1 1 1 1 ...
##  $ amplitude_roll_dumbbell : num  NA NA NA NA NA NA NA NA NA NA ...
##   [list output truncated]
```

## Data cleaning and exploring

The goal is to clean the data and to select the best predictors for the outcome **classe**.

In the summary we see that many columns have NA or blank values. We want to remove the columns that have more than 90% of NA or blank values.

```
naSum <- colSums(is.na(trainData) |trainData=="")/dim(trainData)[1]
colToRemove <- which(naSum>0.9)
trainData <- trainData[-colToRemove]

naSum <- colSums(is.na(testData) |testData=="")/dim(testData)[1]
colToRemove <- which(naSum>0.9)
testData <- testData[-colToRemove]
```

We also see that the first 7 columns are user information and so not relevant for our analysis. We decide to remove them.

```
trainData <- trainData[,-c(1:7)]
testData <- testData[,-c(1:7)]
```

Now the dimension of our data sets are:

```
dim(trainData)
```

```
## [1] 19622    53
```

```
dim(testData)
```

```
## [1] 20 53
```

Then, we remove the columns that are near-zero-variance:

```
NZV <- nearZeroVar(trainData)
NZV
```

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

There no column like that. So the data sets are not changed.

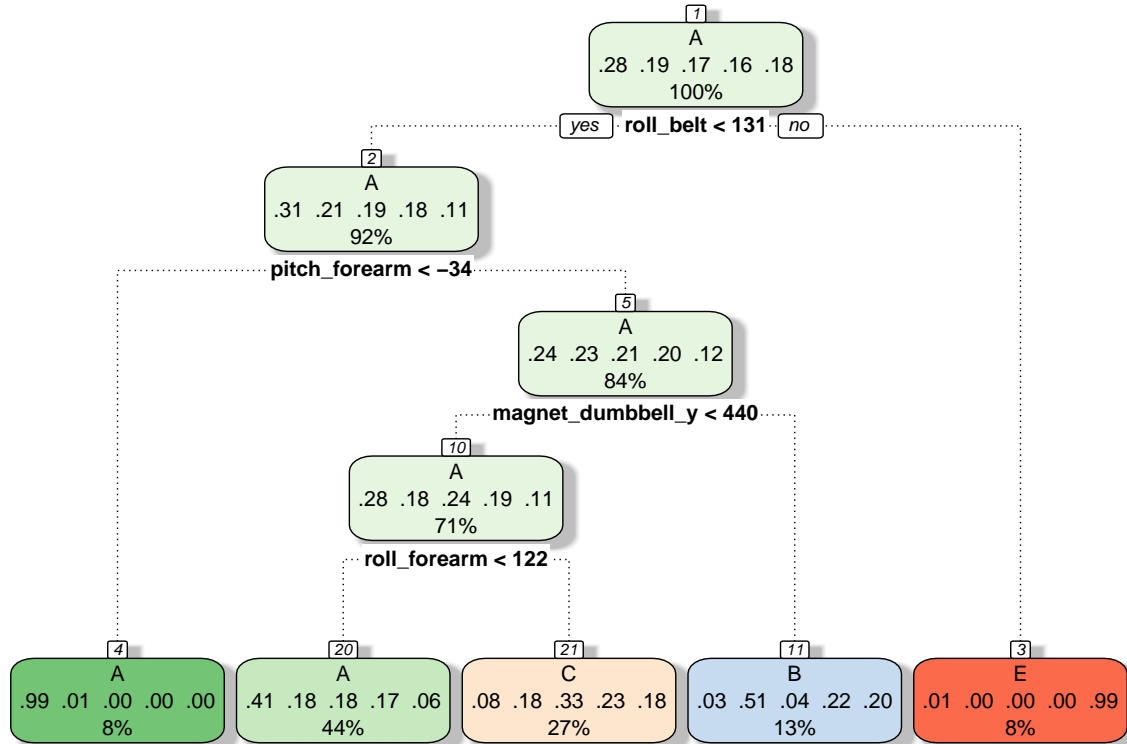Finally, we create a partition 75/25 with the training data set:

```
dataPartition <- createDataPartition(trainData$classe, p=0.75, list=FALSE)
train1 <- trainData[dataPartition,]
test1 <- trainData[-dataPartition,]
```

## Models

We will test three different models in our analysis: classification tree, generalized boosted regression and random forest. We will use the cross-validation technique to improve the efficiency of the models. We will use 5 folds.

**Classification tree**

```
trControl <- trainControl(method="cv", number=5)
model_CT <- train(classe~., data=train1, method="rpart", trControl=trControl)
fancyRpartPlot(model_CT$finalModel)
```



Rattle 2020–avr.–28 09:36:27 Barth

```
trainPred <- predict(model_CT,newdata=test1)
confMatCT <- confusionMatrix(test1$classe,trainPred)
confMatCT
```

```
## Confusion Matrix and Statistics
##
##          Reference
## Prediction    A    B    C    D    E
##          A 1267   19  106    0    3
##          B  397  315  237    0    0
##          C  395   32  428    0    0
##          D  361  147  296    0    0
##          E  127  111  243    0  420
##
## Overall Statistics
##
##                Accuracy : 0.4955
##                  95% CI : (0.4814, 0.5096)
##     No Information Rate : 0.5194
##     P-Value [Acc > NIR] : 0.9996
##
##                   Kappa : 0.3407
##
##  Mcnemar's Test P-Value : NA
##
## Statistics by Class:
##
##                     Class: A Class: B Class: C Class: D Class: E
## Sensitivity           0.4974  0.50481  0.32672       NA  0.99291
## Specificity           0.9457  0.85187  0.88119   0.8361  0.89266
## Pos Pred Value        0.9082  0.33193  0.50058       NA  0.46615
## Neg Pred Value        0.6352  0.92187  0.78217       NA  0.99925
## Prevalence            0.5194  0.12724  0.26713   0.0000  0.08626
## Detection Rate        0.2584  0.06423  0.08728   0.0000  0.08564
## Detection Prevalence  0.2845  0.19352  0.17435   0.1639  0.18373
## Balanced Accuracy     0.7216  0.67834  0.60395       NA  0.94278
```

We can notice that the accuracy of the model is very low (49%). We propably won't keep this model.

**Generalized Boosted Regression**

```
GBMcontrol <- trainControl(method = "repeatedcv", number = 5, repeats = 1)
modGBM  <- train(classe ~ ., data=train1, method = "gbm", trControl = GBMcontrol, verbose = FALSE)
```

```
predictGBM <- predict(modGBM, newdata=test1)
cmGBM <- confusionMatrix(test1$classe, predictGBM)
cmGBM
```

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

```
##           Reference
## Prediction    A    B    C    D    E
##         A 1373   11    5    4    2
##         B   20  909   20    0    0
##         C    0   31  813    7    4
##         D    0    3   28  766    7
##         E    5   14   13   12  857
##
## Overall Statistics
##
##                Accuracy : 0.9621
##                  95% CI : (0.9563, 0.9672)
##     No Information Rate : 0.2851
##     P-Value [Acc > NIR] : < 2.2e-16
##
##                   Kappa : 0.952
##
##  Mcnemar's Test P-Value : 1.783e-07
##
## Statistics by Class:
##
##                    Class: A Class: B Class: C Class: D Class: E
## Sensitivity          0.9821   0.9390   0.9249   0.9708   0.9851
## Specificity          0.9937   0.9898   0.9896   0.9908   0.9891
## Pos Pred Value       0.9842   0.9579   0.9509   0.9527   0.9512
## Neg Pred Value       0.9929   0.9851   0.9837   0.9944   0.9968
## Prevalence           0.2851   0.1974   0.1792   0.1609   0.1774
## Detection Rate       0.2800   0.1854   0.1658   0.1562   0.1748
## Detection Prevalence 0.2845   0.1935   0.1743   0.1639   0.1837
## Balanced Accuracy    0.9879   0.9644   0.9572   0.9808   0.9871
```

The accuracy of this model is good (96%).

**Random Forest**

```
model_RF <- train(classe~., data=train1, method="rf", trControl=trControl, verbose=FALSE)
```

```
trainPred <- predict(model_RF,newdata=test1)
confMatRF <- confusionMatrix(test1$classe,trainPred)
confMatRF
```

```
## Confusion Matrix and Statistics
##
##           Reference
## Prediction    A    B    C    D    E
##         A 1394    0    0    0    1
##         B    6  943    0    0    0
##         C    0    4  849    2    0
##         D    0    0    5  799    0
##         E    0    0    0    0  901
##
```

```
## Overall Statistics
##
##                Accuracy : 0.9963
##                  95% CI : (0.9942, 0.9978)
##     No Information Rate : 0.2855
##     P-Value [Acc > NIR] : < 2.2e-16
##
##                   Kappa : 0.9954
##
##  Mcnemar's Test P-Value : NA
##
## Statistics by Class:
##
##                      Class: A Class: B Class: C Class: D Class: E
## Sensitivity            0.9957   0.9958   0.9941   0.9975   0.9989
## Specificity            0.9997   0.9985   0.9985   0.9988   1.0000
## Pos Pred Value         0.9993   0.9937   0.9930   0.9938   1.0000
## Neg Pred Value         0.9983   0.9990   0.9988   0.9995   0.9998
## Prevalence             0.2855   0.1931   0.1741   0.1633   0.1839
## Detection Rate         0.2843   0.1923   0.1731   0.1629   0.1837
## Detection Prevalence   0.2845   0.1935   0.1743   0.1639   0.1837
## Balanced Accuracy      0.9977   0.9971   0.9963   0.9981   0.9994
```

With the random forest model we get 99% accuracy. This is the highest we got among the three models selected.

## Conclusion

We will use the random forest model to predict the value of the test set.

```
results <- predict(model_RF, newdata=testData)
results
```

```
## [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
```