# Machine Learning Assignment

SofieB

12/4/2019

## Background

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, the goal will be to use data from accelerometers on the belt, forearm, arm, and dumbell of 6 participants and predict the manner in which they did the exercise. This is the "classe" variable in the training set. 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 (http://web.archive.org/web/20161224072740/http:/groupware.les.inf.puc-rio.br/har) (see the section on the Weight Lifting Exercise Dataset).

The data for this project come from this source: http://web.archive.org/web/20161224072740/http:/groupware.les.inf.puc-rio.br/har (http://web.archive.org/web/20161224072740/http:/groupware.les.inf.puc-rio.br/har). Thanks for sharing!

## Data processing

Download and read the data, replace empty observations and #DIV/0! errors by NA, to be able to delete them all together afterwards.

```
library(dplyr)
```

```
##
## Attaching package: 'dplyr'
```

```
## The following objects are masked from 'package:stats':
##
##     filter, lag
```

```
## The following objects are masked from 'package:base':
##
##     intersect, setdiff, setequal, union
```

```
library(caret)
```

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

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

```
library(corrplot)
```

```
## corrplot 0.84 loaded
```

```
library(rattle)
```

```
## Rattle: A free graphical interface for data science with R.
## Version 5.2.0 Copyright (c) 2006-2018 Togaware Pty Ltd.
## Type 'rattle()' to shake, rattle, and roll your data.
```

```
trainingUrl <- "https://d396qusza40orc.cloudfront.net/predmachlearn/pml-training.csv"
testingUrl <- "https://d396qusza40orc.cloudfront.net/predmachlearn/pml-testing.csv"
download.file(url = trainingUrl, destfile = "training.csv")
download.file(url = testingUrl, destfile = "testing.csv")
traindata <- read.csv(file ="training.csv", na.strings = c("", "#DIV/0!","NA"))
testdata <- read.csv(file ="testing.csv", na.strings = c("", "#DIV/0!","NA"))
dim(traindata)
```

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

```
dim(testdata)
```

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

```
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 1323084232
## 1323084232 1323084232 1323084232 ...
##  $ raw_timestamp_part_2     : int  788290 808298 820366 120339 196328 304277 368296 440390 484323 484434 ...
##  $ 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       : num  NA NA NA NA NA NA NA NA NA NA ...
##  $ kurtosis_picth_belt      : num  NA NA NA NA NA NA NA NA NA NA ...
##  $ kurtosis_yaw_belt        : logi  NA NA NA NA NA NA ...
##  $ skewness_roll_belt       : num  NA NA NA NA NA NA NA NA NA NA ...
##  $ skewness_roll_belt.1     : num  NA NA NA NA NA NA NA NA NA NA ...
##  $ skewness_yaw_belt        : logi  NA NA NA NA NA NA ...
##  $ 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             : num  NA NA NA NA NA NA NA NA NA NA ...
##  $ 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             : num  NA NA NA NA NA NA NA NA NA NA ...
##  $ 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       : num  NA NA NA NA NA NA NA NA NA NA ...
##  $ 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        : num  NA NA NA NA NA NA NA NA NA NA ...
##  $ kurtosis_picth_arm       : num  NA NA NA NA NA NA NA NA NA NA ...
##  $ kurtosis_yaw_arm         : num  NA NA NA NA NA NA NA NA NA NA ...
##  $ skewness_roll_arm        : num  NA NA NA NA NA NA NA NA NA NA ...
##  $ skewness_pitch_arm       : num  NA NA NA NA NA NA NA NA NA NA ...
##  $ skewness_yaw_arm         : num  NA NA NA NA NA NA NA NA NA NA ...
##  $ 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   : num  NA NA NA NA NA NA NA NA NA NA ...
##  $ kurtosis_picth_dumbbell  : num  NA NA NA NA NA NA NA NA NA NA ...
##  $ kurtosis_yaw_dumbbell    : logi  NA NA NA NA NA NA ...
##  $ skewness_roll_dumbbell   : num  NA NA NA NA NA NA NA NA NA NA ...
##  $ skewness_pitch_dumbbell  : num  NA NA NA NA NA NA NA NA NA NA ...
```

```
## $ skewness_yaw_dumbbell  : logi  NA NA NA NA NA NA ...
## $ 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       : num  NA NA NA NA NA NA NA NA NA NA ...
## $ 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       : num  NA NA NA NA NA NA NA NA NA NA ...
## $ amplitude_roll_dumbbell : num  NA NA NA NA NA NA NA NA NA NA ...
##  [list output truncated]
```

## Cleaning

Removal of missing variables and variables that have no predictive information for outcome classe (column 1 to 7).

```
training <- select(traindata, 8:160)
testing <- select(testdata, 8:160)
training <- training[, colSums(is.na(training)) ==0]
testing <- testing[, colSums(is.na(testing)) ==0]
dim(training)
```
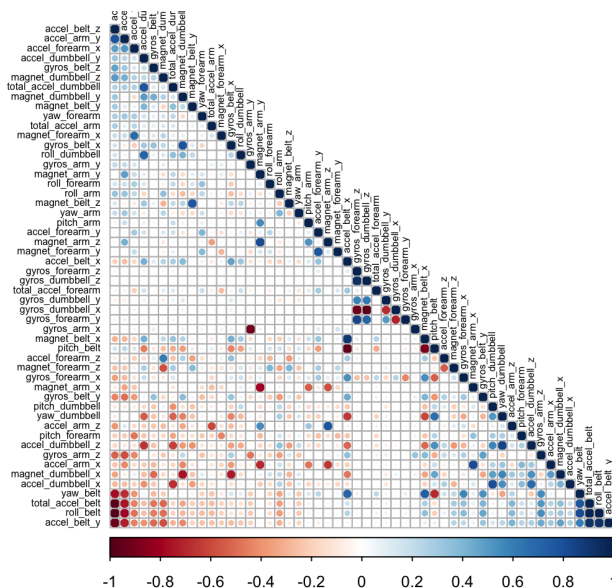
```
## [1] 19622    53
```

```
dim(testing)
```

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

## Exploring

Looking at correlations between the predictors, no further irregularities.

```
cor_matrix <- cor(training[,-53])
corrplot(cor_matrix, order = "FPC", method = "circle", type = "lower",  tl.col = rgb(0, 0, 0), tl.cex = 0.5)
```



# Prediction algorithms

## Splitting the data

Let's try to predict the outcome with classification trees and random forests to see which of both methods performs best. We create a validation set next to the training set in order to be able to compute out-of-sample errors.

```
set.seed(5487)
inTrain <- createDataPartition(training$classe, p = 0.7, list = FALSE)
training <- training[inTrain, ]
validating <- training[-inTrain, ]
dim(training)
```

```
## [1] 13737    53
```

```
dim(validating)
```

```
## [1] 4127   53
```
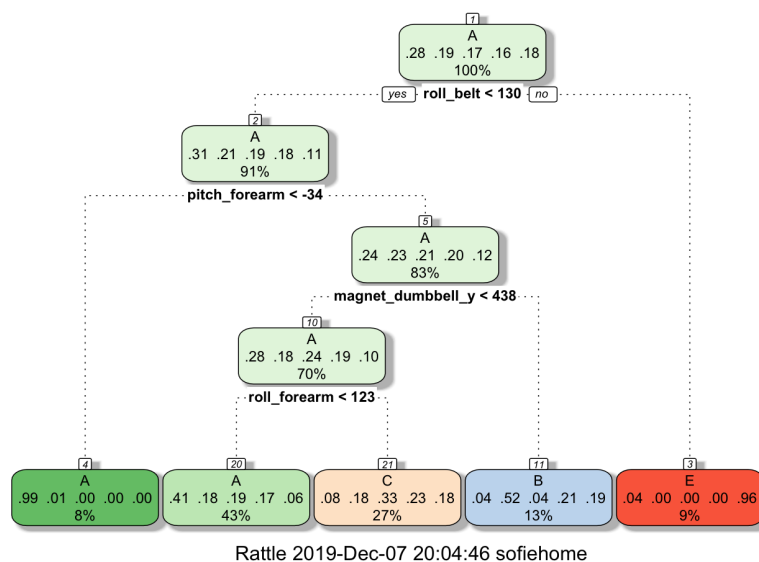
## Classification trees

Creating a model with classification trees, including 3-fold cross validation

```
(modfit_rpart <- train(classe ~ ., data=training, method = "rpart", trControl = trainControl(method = "cv", number=3)))
```

```
## CART
##
## 13737 samples
##    52 predictor
##     5 classes: 'A', 'B', 'C', 'D', 'E'
##
## No pre-processing
## Resampling: Cross-Validated (3 fold)
## Summary of sample sizes: 9158, 9157, 9159
## Resampling results across tuning parameters:
##
##   cp          Accuracy   Kappa
##   0.03438104  0.5105180  0.36144652
##   0.06018377  0.4087586  0.19694196
##   0.11595972  0.3104811  0.04015221
##
## Accuracy was used to select the optimal model using the largest value.
## The final value used for the model was cp = 0.03438104.
```

Visualized in a plot.

```
fancyRpartPlot(modfit_rpart$finalModel)
```



Rattle 2019-Dec-07 20:04:46 sofiehome

Predicting on the validation set.

```
predict_rpart <- predict(modfit_rpart, validating)
(cm_rpart <- confusionMatrix(validating$classe, predict_rpart))
```

```
## Confusion Matrix and Statistics
##
##           Reference
## Prediction    A    B    C    D    E
##          A 1067   21   97    0   15
##          B  308  286  210    0    0
##          C  322   13  350    0    0
##          D  309  114  263    0    0
##          E   96  108  181    0  367
##
## Overall Statistics
##
##                Accuracy : 0.5016
##                  95% CI : (0.4862, 0.5169)
##     No Information Rate : 0.5093
##     P-Value [Acc > NIR] : 0.8442
##
##                   Kappa : 0.3486
##
##  Mcnemar's Test P-Value : NA
##
## Statistics by Class:
##
##                      Class: A Class: B Class: C Class: D Class: E
## Sensitivity            0.5076   0.5277  0.31789       NA  0.96073
## Specificity            0.9343   0.8555  0.88929   0.8338  0.89720
## Pos Pred Value         0.8892   0.3557  0.51095       NA  0.48803
## Neg Pred Value         0.6464   0.9230  0.78181       NA  0.99556
## Prevalence             0.5093   0.1313  0.26678   0.0000  0.09256
## Detection Rate         0.2585   0.0693  0.08481   0.0000  0.08893
## Detection Prevalence   0.2908   0.1948  0.16598   0.1662  0.18221
## Balanced Accuracy      0.7210   0.6916  0.60359       NA  0.92896
```

With an accuracy rate of 0.5 the classification tree model performance is rather low. We notice a big difference in the accuracy results between the 3 sample folds. The best performing option has no classe D as a predicted outcome which seems a bit strange given the normal presence of D within the trainingdata.

Let's see if we can do better.

## Random forests

Modeling with random forests, including 3-fold cross-validation.

```
modfit_rf <- train(classe ~ ., data=training, method = "rf", trControl = trainControl(method = "cv",number = 3))
modfit_rf
```

```
## Random Forest
##
## 13737 samples
##    52 predictor
##     5 classes: 'A', 'B', 'C', 'D', 'E'
##
## No pre-processing
## Resampling: Cross-Validated (3 fold)
## Summary of sample sizes: 9158, 9159, 9157
## Resampling results across tuning parameters:
##
##   mtry  Accuracy   Kappa
##    2    0.9884256  0.9853566
##   27    0.9868969  0.9834228
##   52    0.9758323  0.9694240
##
## Accuracy was used to select the optimal model using the largest value.
## The final value used for the model was mtry = 2.
```

Predicting on the validation set.

```
predict_rf <- predict(modfit_rf, validating)
(cm_rf <- confusionMatrix(validating$classe, predict_rf))
```

```
## Confusion Matrix and Statistics
##
##           Reference
## Prediction    A    B    C    D    E
##          A 1200    0    0    0    0
##          B    0  804    0    0    0
##          C    0    0  685    0    0
##          D    0    0    0  686    0
##          E    0    0    0    0  752
##
## Overall Statistics
##
##                Accuracy : 1
##                  95% CI : (0.9991, 1)
##     No Information Rate : 0.2908
##     P-Value [Acc > NIR] : < 2.2e-16
##
##                   Kappa : 1
##
##  Mcnemar's Test P-Value : NA
##
## Statistics by Class:
##
##                      Class: A Class: B Class: C Class: D Class: E
## Sensitivity            1.0000   1.0000    1.000   1.0000   1.0000
## Specificity            1.0000   1.0000    1.000   1.0000   1.0000
## Pos Pred Value         1.0000   1.0000    1.000   1.0000   1.0000
## Neg Pred Value         1.0000   1.0000    1.000   1.0000   1.0000
## Prevalence             0.2908   0.1948    0.166   0.1662   0.1822
## Detection Rate         0.2908   0.1948    0.166   0.1662   0.1822
## Detection Prevalence   0.2908   0.1948    0.166   0.1662   0.1822
## Balanced Accuracy      1.0000   1.0000    1.000   1.0000   1.0000
```

The performance of the random forests model is clearly much better. The accuracy level is over 98.8 in the training data and goes up to 100% on the validation set. Out-of sample error = 0, but this might be due to overfitting. Anyhow, the random forests model shows the best results.

## Prediction on the testing set

We will continu with the random forests model to predict on the testset.

```
(predict_final <- predict(modfit_rf, testing))
```

```
##  [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
```