

# Assignment: Prediction Assignment Writeup

Carlos Diaz

2 de septiembre de 2018

## OBJECTIVE:

The objective of this assignment is to predict the way in which the subjects of the experiment actually perform the exercise based on the information and variables found in the training set. The results will be tested in a test set and deviations will be quantified.

## 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 dumbbell 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).

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).

## DATA:

The training data for this project are available here:

<https://d396qusza40orc.cloudfront.net/predmachlearn/pml-training.csv>

The test data are available here:

<https://d396qusza40orc.cloudfront.net/predmachlearn/pml-testing.csv>

## REFERENCES:

Velloso, E.; Bulling, A.; Gellersen, H.; Ugulino, W.; Fuks, H. Qualitative Activity Recognition of Weight Lifting Exercises. Proceedings of 4th International Conference in Cooperation with SIGCHI (Augmented Human '13) . Stuttgart, Germany: ACM SIGCHI, 2013.

## LOADING DATA AND LIBRARIES:

```
## Loading libraries:

setwd("C:/Users/ADMIN/Desktop/Data Scientist Specialization/Course 8 -
Practical Machine Learning")
library(ggplot2)
library(lattice)
library(caret)
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:ggplot2':
##
##     margin
```

```

library(e1071)

## Downloading data:

trainURL <- "https://d396qusza40orc.cloudfront.net/predmachlearn/pml-
training.csv"
download.file(url=trainURL, destfile="training.csv")

testURL <- "https://d396qusza40orc.cloudfront.net/predmachlearn/pml-
testing.csv"
download.file(url=testURL, destfile="testing.csv")

## Reading data:

train <- read.csv("training.csv", na.strings=c("NA", "#DIV/0!", ""))
test <- read.csv("testing.csv", na.strings=c("NA", "#DIV/0!", ""))
str(train)

## '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_pitch_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_pitch_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 NA ...
## $ amplitude_roll_belt : num  NA NA NA NA NA NA NA NA NA NA NA ...
## $ amplitude_pitch_belt : int   NA NA NA NA NA NA NA NA NA NA NA ...
## $ amplitude_yaw_belt  : num  NA 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 NA ...
## $ avg_roll_belt      : num  NA NA NA NA NA NA NA NA NA NA NA ...
## $ stddev_roll_belt   : num  NA NA NA NA NA NA NA NA NA NA NA ...
## $ var_roll_belt      : num  NA NA NA NA NA NA NA NA NA NA NA ...
## $ avg_pitch_belt     : num  NA NA NA NA NA NA NA NA NA NA NA ...
## $ stddev_pitch_belt  : num  NA NA NA NA NA NA NA NA NA NA NA ...
## $ var_pitch_belt     : num  NA NA NA NA NA NA NA NA NA NA NA ...
## $ avg_yaw_belt       : num  NA NA NA NA NA NA NA NA NA NA NA ...
## $ stddev_yaw_belt    : num  NA NA NA NA NA NA NA NA NA NA NA ...
## $ var_yaw_belt       : num  NA 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
0.02 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 NA ...
## $ avg_roll_arm       : num  NA NA NA NA NA NA NA NA NA NA NA ...
## $ stddev_roll_arm    : num  NA NA NA NA NA NA NA NA NA NA NA ...
## $ var_roll_arm       : num  NA NA NA NA NA NA NA NA NA NA NA ...
## $ avg_pitch_arm      : num  NA NA NA NA NA NA NA NA NA NA NA ...
## $ stddev_pitch_arm   : num  NA NA NA NA NA NA NA NA NA NA NA ...
## $ var_pitch_arm      : num  NA NA NA NA NA NA NA NA NA NA NA ...
## $ avg_yaw_arm        : num  NA NA NA NA NA NA NA NA NA NA NA ...
## $ stddev_yaw_arm     : num  NA NA NA NA NA NA NA NA NA NA NA ...
## $ var_yaw_arm        : num  NA 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_pitch_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_pitch_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_pitch_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_pitch_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]

## Outcome to be predicted:

summary(train$classe)

```

```
##      A      B      C      D      E
## 5580 3797 3422 3216 3607
```

## TRAINING AND TESTING DATA:

```
## Subsetting train data for cross-validation (60% - 40%):
```

```
dpTrain <- createDataPartition(y=train$classe, p=0.6, list=FALSE)
Train_1 <- train[dpTrain, ]
Test_1 <- train[-dpTrain, ]
dim(Train_1)
```

```
## [1] 11776 160
```

```
dim(Test_1)
```

```
## [1] 7846 160
```

## FILTERING USEFUL DATA:

```
## Removing variables with NAs >= 75%:
```

```
Train_1_NAs <- Train_1
for (i in 1:length(Train_1)) {
  if (sum(is.na(Train_1[, i])) / nrow(Train_1) >= .75) {
    for (j in 1:length(Train_1_NAs)) {
      if (length(grep(names(Train_1[i]), names(Train_1_NAs)[j]))==1) {
        Train_1_NAs <- Train_1_NAs[, -j]
      }}
    }
}
dim(Train_1_NAs)
```

```
## [1] 11776 60
```

```
## Removing columns that are not predictors:
```

```
Train_FIL1 <- Train_1_NAs[, 8:length(Train_1_NAs)]
dim(Train_FIL1)
```

```
## [1] 11776 53
```

## RANDOM FOREST ESTIMATION:

```
## A Random Forest Estimation to predict classe is made on the filtered
training dataset which afterwards is cross-validated:
```

```
set.seed(666)
```

```
RFest <- randomForest(classe~., data = Train_FIL1)
print(RFest)
```

```
##
```

```
## Call:
```

```
## randomForest(formula = classe ~ ., data = Train_FIL1)
```

```

##           Type of random forest: classification
##           Number of trees: 500
## No. of variables tried at each split: 7
##
##           OOB estimate of  error rate: 0.65%
## Confusion matrix:
##      A      B      C      D      E  class.error
## A 3345      3      0      0      0 0.0008960573
## B   16 2257      6      0      0 0.0096533567
## C      0   19 2032      3      0 0.0107108082
## D      0      0   19 1910      1 0.0103626943
## E      0      0      3      6 2156 0.0041570439

## Cross validation:

CrossV <- predict(RFest, Test_1, type = "class")
confusionMatrix(Test_1$classe, CrossV)

## Confusion Matrix and Statistics
##
##           Reference
## Prediction      A      B      C      D      E
##           A 2226      4      1      0      1
##           B      5 1507      6      0      0
##           C      0   17 1350      1      0
##           D      0      0   19 1267      0
##           E      0      0      1      3 1438
##
## Overall Statistics
##
##           Accuracy : 0.9926
##           95% CI : (0.9905, 0.9944)
##           No Information Rate : 0.2843
##           P-Value [Acc > NIR] : < 2.2e-16
##
##           Kappa : 0.9906
##           McNemar's Test P-Value : NA
##
## Statistics by Class:
##
##           Class: A Class: B Class: C Class: D Class: E
## Sensitivity      0.9978  0.9863  0.9804  0.9969  0.9993
## Specificity      0.9989  0.9983  0.9972  0.9971  0.9994
## Pos Pred Value   0.9973  0.9928  0.9868  0.9852  0.9972
## Neg Pred Value   0.9991  0.9967  0.9958  0.9994  0.9998
## Prevalence       0.2843  0.1947  0.1755  0.1620  0.1834
## Detection Rate   0.2837  0.1921  0.1721  0.1615  0.1833
## Detection Prevalence 0.2845  0.1935  0.1744  0.1639  0.1838
## Balanced Accuracy 0.9983  0.9923  0.9888  0.9970  0.9993

```

```

CrossV2 <- predict(RFest, Train_1, type = "class")
confusionMatrix(Train_1$classe, CrossV2)

## Confusion Matrix and Statistics
##
##              Reference
## Prediction      A      B      C      D      E
##      A 3348      0      0      0      0
##      B      0 2279      0      0      0
##      C      0      0 2054      0      0
##      D      0      0      0 1930      0
##      E      0      0      0      0 2165
##
## Overall Statistics
##
##              Accuracy : 1
##              95% CI : (0.9997, 1)
##      No Information Rate : 0.2843
##      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   1.0000   1.0000   1.0000
## Specificity          1.0000   1.0000   1.0000   1.0000   1.0000
## Pos Pred Value       1.0000   1.0000   1.0000   1.0000   1.0000
## Neg Pred Value       1.0000   1.0000   1.0000   1.0000   1.0000
## Prevalence           0.2843   0.1935   0.1744   0.1639   0.1838
## Detection Rate       0.2843   0.1935   0.1744   0.1639   0.1838
## Detection Prevalence 0.2843   0.1935   0.1744   0.1639   0.1838
## Balanced Accuracy    1.0000   1.0000   1.0000   1.0000   1.0000

```

## PARTIAL RESULTS:

As it may be seen from the estimations results, it is possible to evidenciate that running the model on test data for cross validation it is possible to find an accuracy of 99.3%. When the model is confronted to training data used to build the model, it shows a 100% accuracy.

## TESTING MODEL WITH TEST DATASET:

```

FINAL_forecasting <- predict(RFest, test, type = "class")
print(FINAL_forecasting)

##  1  2  3  4  5  6  7  8  9 10 11 12 13 14 15 16 17 18 19 20
##  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

```