# Practical Machine Learning Project

*Nikitha Mohan*

*11/19/2018*

## Overview

We have data from accelerometers on the belt, forearm, arm, and dumbell of 6 participants who lifted them correctly and incorrectly. We want to find a way to use this data to predict the manner in which they did the exercise and predict 20 different test cases.

## Load in Libraries to be used

```
library(caret)
```

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

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

```
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.
```

## Loading the data

Training Data:

```
TrainData <- read.csv(url("https://d396qusza40orc.cloudfront.net/predmachlearn/pml-training.csv"),header=
```

Test Data:

```
TestData <- read.csv(url("https://d396qusza40orc.cloudfront.net/predmachlearn/pml-testing.csv"),header=T
```

```
str(TestData)
```

```
## 'data.frame':    20 obs. of  160 variables:
##  $ X                    : int  1 2 3 4 5 6 7 8 9 10 ...
##  $ user_name            : Factor w/ 6 levels "adelmo","carlitos",..: 6 5 5 1 4 5 5 5 2 3 ...
##  $ raw_timestamp_part_1 : int  1323095002 1322673067 1322673075 1322832789 1322489635 1322673149 
##  $ raw_timestamp_part_2 : int  868349 778725 342967 560311 814776 510661 766645 54671 916313 38428
##  $ cvtd_timestamp       : Factor w/ 11 levels "02/12/2011 13:33",..: 5 10 10 1 6 11 11 10 3 2 ...
##  $ new_window           : Factor w/ 1 level "no": 1 1 1 1 1 1 1 1 1 1 ...
##  $ num_window           : int  74 431 439 194 235 504 485 440 323 664 ...
##  $ roll_belt            : num  123 1.02 0.87 125 1.35 -5.92 1.2 0.43 0.93 114 ...
##  $ pitch_belt           : num  27 4.87 1.82 -41.6 3.33 1.59 4.44 4.15 6.72 22.4 ...
##  $ yaw_belt             : num  -4.75 -88.9 -88.5 162 -88.6 -87.7 -87.3 -88.5 -93.7 -13.1 ...
##  $ total_accel_belt     : int  20 4 5 17 3 4 4 4 4 18 ...
##  $ kurtosis_roll_belt   : logi  NA NA NA NA NA NA ...
##  $ kurtosis_picth_belt  : logi  NA NA NA NA NA NA ...
```

```
##  $ kurtosis_yaw_belt      : logi  NA NA NA NA NA NA ...
##  $ skewness_roll_belt     : logi  NA NA NA NA NA NA ...
##  $ skewness_roll_belt.1   : logi  NA NA NA NA NA NA ...
##  $ skewness_yaw_belt      : logi  NA NA NA NA NA NA ...
##  $ max_roll_belt          : logi  NA NA NA NA NA NA ...
##  $ max_picth_belt         : logi  NA NA NA NA NA NA ...
##  $ max_yaw_belt           : logi  NA NA NA NA NA NA ...
##  $ min_roll_belt          : logi  NA NA NA NA NA NA ...
##  $ min_pitch_belt         : logi  NA NA NA NA NA NA ...
##  $ min_yaw_belt           : logi  NA NA NA NA NA NA ...
##  $ amplitude_roll_belt    : logi  NA NA NA NA NA NA ...
##  $ amplitude_pitch_belt   : logi  NA NA NA NA NA NA ...
##  $ amplitude_yaw_belt     : logi  NA NA NA NA NA NA ...
##  $ var_total_accel_belt   : logi  NA NA NA NA NA NA ...
##  $ avg_roll_belt          : logi  NA NA NA NA NA NA ...
##  $ stddev_roll_belt       : logi  NA NA NA NA NA NA ...
##  $ var_roll_belt          : logi  NA NA NA NA NA NA ...
##  $ avg_pitch_belt         : logi  NA NA NA NA NA NA ...
##  $ stddev_pitch_belt      : logi  NA NA NA NA NA NA ...
##  $ var_pitch_belt         : logi  NA NA NA NA NA NA ...
##  $ avg_yaw_belt           : logi  NA NA NA NA NA NA ...
##  $ stddev_yaw_belt        : logi  NA NA NA NA NA NA ...
##  $ var_yaw_belt           : logi  NA NA NA NA NA NA ...
##  $ gyros_belt_x           : num  -0.5 -0.06 0.05 0.11 0.03 0.1 -0.06 -0.18 0.1 0.14 ...
##  $ gyros_belt_y           : num  -0.02 -0.02 0.02 0.11 0.02 0.05 0 -0.02 0 0.11 ...
##  $ gyros_belt_z           : num  -0.46 -0.07 0.03 -0.16 0 -0.13 0 -0.03 -0.02 -0.16 ...
##  $ accel_belt_x           : int  -38 -13 1 46 -8 -11 -14 -10 -15 -25 ...
##  $ accel_belt_y           : int  69 11 -1 45 4 -16 2 -2 1 63 ...
##  $ accel_belt_z           : int  -179 39 49 -156 27 38 35 42 32 -158 ...
##  $ magnet_belt_x          : int  -13 43 29 169 33 31 50 39 -6 10 ...
##  $ magnet_belt_y          : int  581 636 631 608 566 638 622 635 600 601 ...
##  $ magnet_belt_z          : int  -382 -309 -312 -304 -418 -291 -315 -305 -302 -330 ...
##  $ roll_arm               : num  40.7 0 0 -109 76.1 0 0 0 -137 -82.4 ...
##  $ pitch_arm              : num  -27.8 0 0 55 2.76 0 0 0 11.2 -63.8 ...
##  $ yaw_arm                : num  178 0 0 -142 102 0 0 0 -167 -75.3 ...
##  $ total_accel_arm        : int  10 38 44 25 29 14 15 22 34 32 ...
##  $ var_accel_arm          : logi  NA NA NA NA NA NA ...
##  $ avg_roll_arm           : logi  NA NA NA NA NA NA ...
##  $ stddev_roll_arm        : logi  NA NA NA NA NA NA ...
##  $ var_roll_arm           : logi  NA NA NA NA NA NA ...
##  $ avg_pitch_arm          : logi  NA NA NA NA NA NA ...
##  $ stddev_pitch_arm       : logi  NA NA NA NA NA NA ...
##  $ var_pitch_arm          : logi  NA NA NA NA NA NA ...
##  $ avg_yaw_arm            : logi  NA NA NA NA NA NA ...
##  $ stddev_yaw_arm         : logi  NA NA NA NA NA NA ...
##  $ var_yaw_arm            : logi  NA NA NA NA NA NA ...
##  $ gyros_arm_x            : num  -1.65 -1.17 2.1 0.22 -1.96 0.02 2.36 -3.71 0.03 0.26 ...
##  $ gyros_arm_y            : num  0.48 0.85 -1.36 -0.51 0.79 0.05 -1.01 1.85 -0.02 -0.5 ...
##  $ gyros_arm_z            : num  -0.18 -0.43 1.13 0.92 -0.54 -0.07 0.89 -0.69 -0.02 0.79 ...
##  $ accel_arm_x            : int  16 -290 -341 -238 -197 -26 99 -98 -287 -301 ...
##  $ accel_arm_y            : int  38 215 245 -57 200 130 79 175 111 -42 ...
##  $ accel_arm_z            : int  93 -90 -87 6 -30 -19 -67 -78 -122 -80 ...
##  $ magnet_arm_x           : int  -326 -325 -264 -173 -170 396 702 535 -367 -420 ...
##  $ magnet_arm_y           : int  385 447 474 257 275 176 15 215 335 294 ...
```

```
##  $ magnet_arm_z            : int   481 434 413 633 617 516 217 385 520 493 ...
##  $ kurtosis_roll_arm       : logi  NA NA NA NA NA NA ...
##  $ kurtosis_picth_arm      : logi  NA NA NA NA NA NA ...
##  $ kurtosis_yaw_arm        : logi  NA NA NA NA NA NA ...
##  $ skewness_roll_arm       : logi  NA NA NA NA NA NA ...
##  $ skewness_pitch_arm      : logi  NA NA NA NA NA NA ...
##  $ skewness_yaw_arm        : logi  NA NA NA NA NA NA ...
##  $ max_roll_arm            : logi  NA NA NA NA NA NA ...
##  $ max_picth_arm           : logi  NA NA NA NA NA NA ...
##  $ max_yaw_arm             : logi  NA NA NA NA NA NA ...
##  $ min_roll_arm            : logi  NA NA NA NA NA NA ...
##  $ min_pitch_arm           : logi  NA NA NA NA NA NA ...
##  $ min_yaw_arm             : logi  NA NA NA NA NA NA ...
##  $ amplitude_roll_arm      : logi  NA NA NA NA NA NA ...
##  $ amplitude_pitch_arm     : logi  NA NA NA NA NA NA ...
##  $ amplitude_yaw_arm       : logi  NA NA NA NA NA NA ...
##  $ roll_dumbbell           : num   -17.7 54.5 57.1 43.1 -101.4 ...
##  $ pitch_dumbbell          : num   25 -53.7 -51.4 -30 -53.4 ...
##  $ yaw_dumbbell            : num   126.2 -75.5 -75.2 -103.3 -14.2 ...
##  $ kurtosis_roll_dumbbell  : logi  NA NA NA NA NA NA ...
##  $ kurtosis_picth_dumbbell : logi  NA NA NA NA NA NA ...
##  $ kurtosis_yaw_dumbbell   : logi  NA NA NA NA NA NA ...
##  $ skewness_roll_dumbbell  : logi  NA NA NA NA NA NA ...
##  $ skewness_pitch_dumbbell : logi  NA NA NA NA NA NA ...
##  $ skewness_yaw_dumbbell   : logi  NA NA NA NA NA NA ...
##  $ max_roll_dumbbell       : logi  NA NA NA NA NA NA ...
##  $ max_picth_dumbbell      : logi  NA NA NA NA NA NA ...
##  $ max_yaw_dumbbell        : logi  NA NA NA NA NA NA ...
##  $ min_roll_dumbbell       : logi  NA NA NA NA NA NA ...
##  $ min_pitch_dumbbell      : logi  NA NA NA NA NA NA ...
##  $ min_yaw_dumbbell        : logi  NA NA NA NA NA NA ...
##  $ amplitude_roll_dumbbell : logi  NA NA NA NA NA NA ...
##   [list output truncated]
```

Dimensions of Training Data

```r
dim(TrainData)
```

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

Dimentions of Test Data

```r
dim(TestData)
```

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

We can see that the training data has 19622 observations on 160 columns and some of those column have NAs and blank values. We need to remove them because they wont produce any information. In addition the first seven columns have information about people who did the test and also had the timestamps, we do not need these.

# Removing Unnecessary Data

In this step, we will clean the data and get rid of observations with missing values as well as some meaningless variables.

```
sum(complete.cases(TrainData))
```

## [1] 406

First, we remove columns that contain NA missing values.

```
TrainData <- TrainData[, colSums(is.na(TrainData)) == 0]
TestData <- TestData[, colSums(is.na(TestData)) == 0]
```

Next, we get rid of some columns that do not contribute much to the accelerometer measurements.

```
classe <- TrainData$classe
trainRemove <- grepl("^X|timestamp|window", names(TrainData))
TrainData <- TrainData[, !trainRemove]
TrainClean <- TrainData[, sapply(TrainData, is.numeric)]
TrainClean$classe <- classe
testRemove <- grepl("^X|timestamp|window", names(TestData))
TestData <- TestData[, !testRemove]
TestClean <- TestData[, sapply(TestData, is.numeric)]
```

Dimensions of cleaned Trained Data

```
dim(TrainClean)
```

## [1] 19622    53

Dimensions of cleaned Test Data

```
dim(TestClean)
```

## [1] 20 53

Now, the cleaned training data set contains 19622 observations and 53 variables, while the testing data set contains 20 observations and 53 variables. The "classe" variable is still in the cleaned training set. After cleaning the data there is only 53 columns that have data present

# Splitting up data for Machine learning

We want to split the data into 80% training data to 20% test data

```
set.seed(15697)
inTrain <- createDataPartition(TrainClean$classe, p=0.80, list=F)
TrainData <- TrainClean[inTrain, ]
TestData <- TrainClean[-inTrain, ]
```

# Machine Learning Method: Random Forest

```
set.seed(15697)
controlRF <- trainControl(method="cv", number=3, verboseIter=FALSE)
modFitRandForest <- train(classe ~ ., data=TrainData, 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: 0.61%
## Confusion matrix:
##      A    B    C    D    E class.error
## A 4457    5    1    0    1 0.001568100
## B   20 3012    4    2    0 0.008558262
## C    0   10 2717   11    0 0.007669832
## D    0    1   28 2543    1 0.011659541
## E    0    1    3    7 2875 0.003811504
```

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

```
## Confusion Matrix and Statistics
##
##           Reference
## Prediction    A    B    C    D    E
##          A 1116    2    0    0    0
##          B    0  755    6    0    0
##          C    0    2  674    6    1
##          D    0    0    4  635    1
##          E    0    0    0    2  719
##
## Overall Statistics
##
##                Accuracy : 0.9939
##                  95% CI : (0.9909, 0.9961)
##     No Information Rate : 0.2845
##     P-Value [Acc > NIR] : < 2.2e-16
##
##                   Kappa : 0.9923
##  Mcnemar's Test P-Value : NA
##
## Statistics by Class:
##
##                      Class: A Class: B Class: C Class: D Class: E
## Sensitivity            1.0000   0.9947   0.9854   0.9876   0.9972
## Specificity            0.9993   0.9981   0.9972   0.9985   0.9994
## Pos Pred Value         0.9982   0.9921   0.9868   0.9922   0.9972
## Neg Pred Value         1.0000   0.9987   0.9969   0.9976   0.9994
## Prevalence             0.2845   0.1935   0.1744   0.1639   0.1838
## Detection Rate         0.2845   0.1925   0.1718   0.1619   0.1833
## Detection Prevalence   0.2850   0.1940   0.1741   0.1631   0.1838
## Balanced Accuracy      0.9996   0.9964   0.9913   0.9930   0.9983
```

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

# Random Forest – Accuracy = 0.9939