# Weight_lifting_exercise_assignment

*kannanthegreat*

*December 6, 2016*

## Executive Summary

This document is the final report of the Peer Assessment project from Coursera's course Practical Machine Learning, as part of the Specialization in Data Science. It was built up in RStudio, using its knitr functions, meant to be published in html format. Quantified Self devices are becoming more and more common, and are able to collect a large amount of data about people and their personal health activities. The focus of this project is to utilize some sample data on the quality of certain exercises to predict the manner in which they did the exercise. This analysis is meant to be the basis for the course quiz and a prediction assignment writeup. The main goal of the project is to predict the manner in which 6 participants performed some exercise as described below. Three ML models, namely, Decision Tree, GBM and Randon Forest algorithms were considered among which Random forest was found to be most accurate, hence was applied to test the cases.

## Steps

- Process the data
- Explore the data
- Model selection
- Model cross validation
- Determining expected out of sample error
- Using prediction model to predict 20 different test cases.

## Data Source

The training data for this project are available at:

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

The test data are available at:

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

The data for this project come from : http://groupware.les.inf.puc-rio.br/har.

## Data Loading and Exploratory Analysis

The following Libraries were used for this project:

```
library(caret)
library(rpart)
library(rpart.plot)
library(RColorBrewer)
library(rattle)
library(randomForest)
library(ggplot2)
library(gbm)
library(survival)
```

```
library(splines)
library(parallel)
library(plyr)
```

## Loading and exploring Data:

The analysis starts by downloading the data into local files. There are 2 data sets, the training data set and the testing data set. The data exploration reveals many NAs in both data sets. When the data is loaded into dataframes, it is necessary to locate strings containing '#DIV/0!' in otherwise numeric data, a common sentinal error code for division by zero errors. These error codes are loaded into the data frame as NA fields.

```
set.seed(12345)
training <- read.csv(("C:/Users/kannanthegreat/Documents/Data_train/pml-training.csv"), na.strings=c("N
testing <- read.csv(("C:/Users/kannanthegreat/Documents/Data_test/pml-testing.csv"), na.strings=c("NA",
```

```
str(training)
```

```
## '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 4844
##  $ 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]
```

```
str(testing)
```

```
## '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]
```

## Partioning Training data set into two data sets:

To find an optimal model, with the best performance both in Accuracy as well as minimizing Out of Sample Error, the full testing data is split randomly with a set seed with 60% of the data into the training sample and 40% of the data used as cross-validation.

```
inTrain <- createDataPartition(y=training$classe, p=0.6, list=FALSE)
myTraining <- training[inTrain, ]; myTesting <- training[-inTrain, ]
dim(myTraining);
```

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

```
dim(myTesting)
```

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

## Cleaning the data: Transformation 1 - Cleaning Near Zero Variance Variables

```
myDataNZV <- nearZeroVar(myTraining, saveMetrics=TRUE)
```

## Create another subset without Near Zero Variance variables:

```
myNZVvars <- names(myTraining) %in% c("new_window", "kurtosis_roll_belt", "kurtosis_picth_belt",
                                      "kurtosis_yaw_belt", "skewness_roll_belt", "skewness_roll_belt.1"
                                      "max_yaw_belt", "min_yaw_belt", "amplitude_yaw_belt", "avg_roll_a
                                      "var_roll_arm", "avg_pitch_arm", "stddev_pitch_arm", "var_pitch_a
                                      "stddev_yaw_arm", "var_yaw_arm", "kurtosis_roll_arm", "kurtosis_p
                                      "kurtosis_yaw_arm", "skewness_roll_arm", "skewness_pitch_arm", "s
                                      "max_roll_arm", "min_roll_arm", "min_pitch_arm", "amplitude_roll_a
                                      "kurtosis_roll_dumbbell", "kurtosis_picth_dumbbell", "kurtosis_ya
                                      "skewness_pitch_dumbbell", "skewness_yaw_dumbbell", "max_yaw_dumbl
                                      "amplitude_yaw_dumbbell", "kurtosis_roll_forearm", "kurtosis_pictl
                                      "skewness_roll_forearm", "skewness_pitch_forearm", "skewness_yaw_
                                      "max_yaw_forearm", "min_roll_forearm", "min_yaw_forearm", "amplitu
                                      "amplitude_yaw_forearm", "avg_roll_forearm", "stddev_roll_forearm"
                                      "avg_pitch_forearm", "stddev_pitch_forearm", "var_pitch_forearm",
                                      "stddev_yaw_forearm", "var_yaw_forearm")
myTraining <- myTraining[!myNZVvars]
#To check the new N?? of observations
dim(myTraining)
```

```
## [1] 11776   100
```

## Transformation 2:

**Removing first column of Dataset (ID) it does not interfer with ML Algorithms**

```
myTraining <- myTraining[c(-1)]
```

## Transformation 3:

**Cleaning Variables with too many NAs. For Variables that have more than a 60% threshold of NA's I'm going to leave them out**

```
trainingV3 <- myTraining #creating another subset to iterate in loop
for(i in 1:length(myTraining)) { #for every column in the training dataset
        if( sum( is.na( myTraining[, i] ) ) /nrow(myTraining) >= .6 ) { #if n?? NAs > 60% of total obse
                for(j in 1:length(trainingV3)) {
                        if( length( grep(names(myTraining[i]), names(trainingV3)[j]) ) ==1)  { #if the
                                trainingV3 <- trainingV3[ , -j] #Remove that column
                        }
                }
        }
}
#To check the new N?? of observations
dim(trainingV3)
```

```
## [1] 11776    58
```

**Applying transformations to "myTesting" and "testing" data sets**

```
myTraining <- trainingV3
rm(trainingV3)
clean1 <- colnames(myTraining)
clean2 <- colnames(myTraining[, -58])
myTesting <- myTesting[clean1]
testing <- testing[clean2]
dim(myTesting)
```

```
## [1] 7846    58
```

```
dim(testing)
```

```
## [1] 20 57
```

**In order to ensure proper functioning of Decision Trees and especially RandomForest Algorithm with the Test data set (data set provided), we need to coerce the data into the same type.**

```
for (i in 1:length(testing) ) {
        for(j in 1:length(myTraining)) {
        if( length( grep(names(myTraining[i]), names(testing)[j]) ) ==1)  {
            class(testing[j]) <- class(myTraining[i])
        }
```

```
    }
}
#And to make sure Coertion really worked, simple smart ass technique:
testing <- rbind(myTraining[2, -58] , testing) #note removing row 2 as it does not mean anything
testing <- testing[-1,]
```
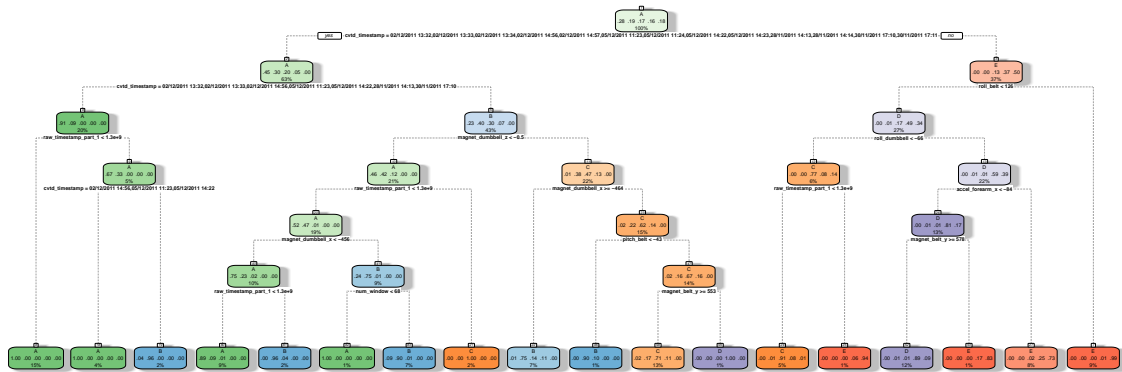
## Using ML algorithms for prediction: Decision Tree

```
modFitA1 <- rpart(classe ~ ., data=myTraining, method="class")
fancyRpartPlot(modFitA1)
```



Rattle 2016–Dec–06 06:22:36 kannanthegreat

## Predicting and Using confusion Matrix to test results:

```
predictionsA1 <- predict(modFitA1, myTesting, type = "class")
cmtree <- confusionMatrix(predictionsA1, myTesting$classe)
cmtree
```

```
## Confusion Matrix and Statistics
##
##           Reference
## Prediction    A    B    C    D    E
##          A 2150   60    7    1    0
##          B   61 1260   69   64    0
##          C   21  188 1269  143    4
##          D    0   10   14  857   78
##          E    0    0    9  221 1360
##
## Overall Statistics
##
##                Accuracy : 0.8789
##                  95% CI : (0.8715, 0.8861)
##     No Information Rate : 0.2845
##     P-Value [Acc > NIR] : < 2.2e-16
```
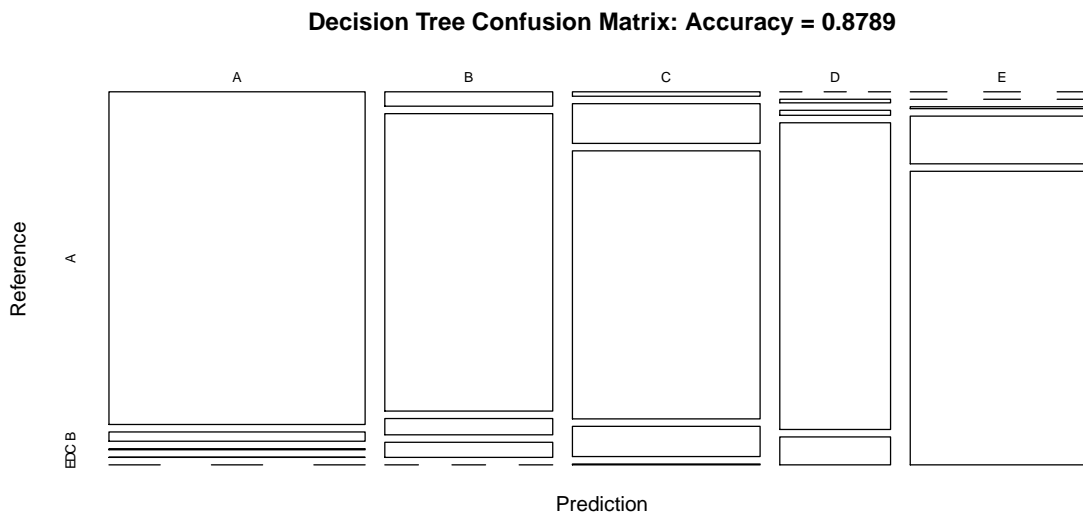
```
##
##                   Kappa : 0.8468
##   Mcnemar's Test P-Value : NA
##
## Statistics by Class:
##
##                    Class: A Class: B Class: C Class: D Class: E
## Sensitivity          0.9633   0.8300   0.9276   0.6664   0.9431
## Specificity          0.9879   0.9693   0.9450   0.9845   0.9641
## Pos Pred Value        0.9693   0.8666   0.7809   0.8936   0.8553
## Neg Pred Value        0.9854   0.9596   0.9841   0.9377   0.9869
## Prevalence           0.2845   0.1935   0.1744   0.1639   0.1838
## Detection Rate        0.2740   0.1606   0.1617   0.1092   0.1733
## Detection Prevalence  0.2827   0.1853   0.2071   0.1222   0.2027
## Balanced Accuracy     0.9756   0.8997   0.9363   0.8254   0.9536
```

```r
plot(cmtree$table, col = cmtree$byClass, main = paste("Decision Tree Confusion Matrix: Accuracy =", rou
```

**Decision Tree Confusion Matrix: Accuracy = 0.8789**



## Using ML algorithms for prediction: Random Forests ### Predicting in-sample error and Using confusion Matrix to test results:
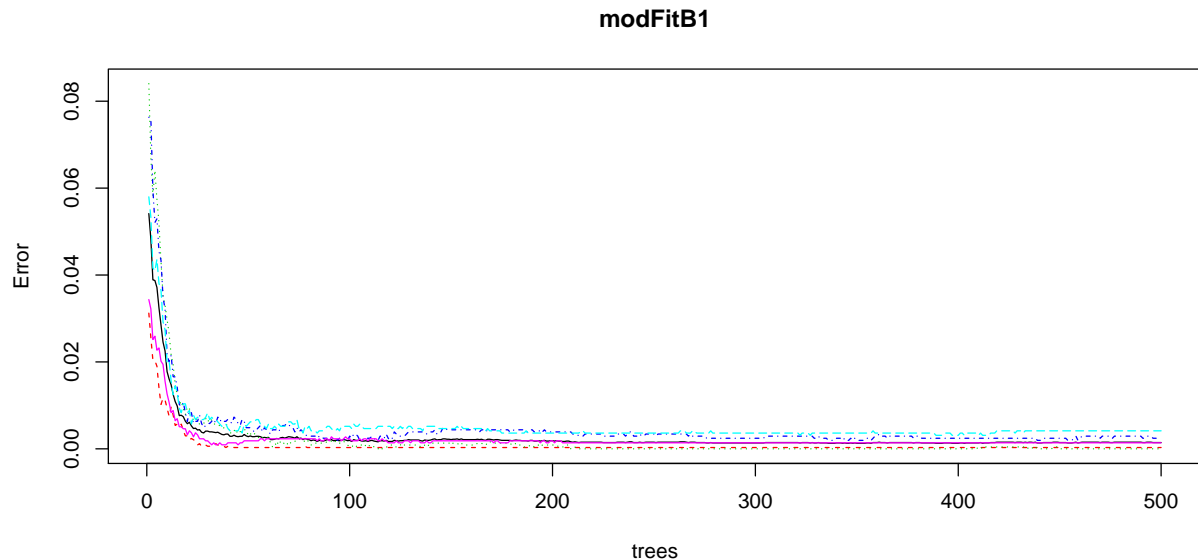
```r
set.seed(12345)
modFitB1 <- randomForest(classe ~ ., data=myTraining)
predictionB1 <- predict(modFitB1, myTesting, type = "class")
cmrf <- confusionMatrix(predictionB1, myTesting$classe)
cmrf
```

```
## Confusion Matrix and Statistics
##
##           Reference
## Prediction    A    B    C    D    E
##          A 2231    2    0    0    0
##          B    1 1516    0    0    0
##          C    0    0 1367    3    0
##          D    0    0    1 1282    1
##          E    0    0    0    1 1441
```
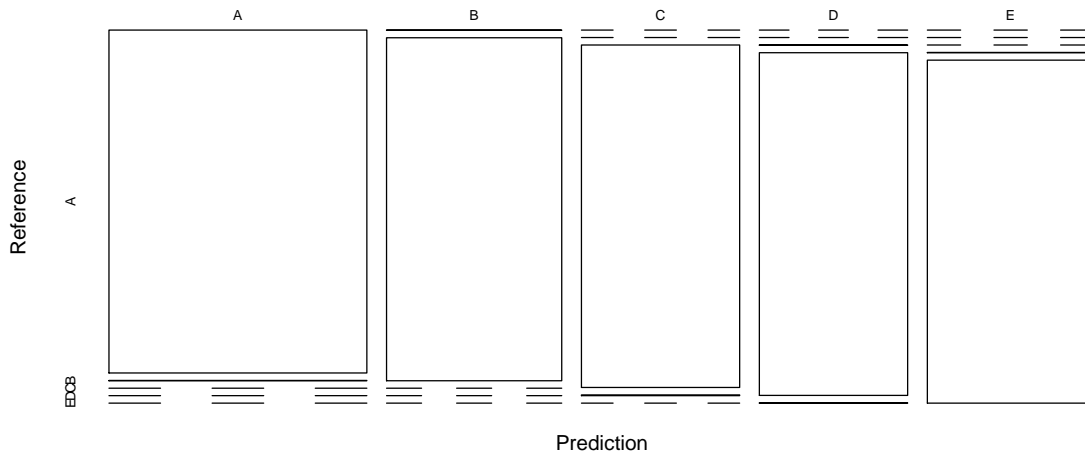
```
##
## Overall Statistics
##
##                Accuracy : 0.9989
##                  95% CI : (0.9978, 0.9995)
##     No Information Rate : 0.2845
##     P-Value [Acc > NIR] : < 2.2e-16
##
##                   Kappa : 0.9985
##  Mcnemar's Test P-Value : NA
##
## Statistics by Class:
##
##                      Class: A Class: B Class: C Class: D Class: E
## Sensitivity            0.9996   0.9987   0.9993   0.9969   0.9993
## Specificity            0.9996   0.9998   0.9995   0.9997   0.9998
## Pos Pred Value         0.9991   0.9993   0.9978   0.9984   0.9993
## Neg Pred Value         0.9998   0.9997   0.9998   0.9994   0.9998
## Prevalence             0.2845   0.1935   0.1744   0.1639   0.1838
## Detection Rate         0.2843   0.1932   0.1742   0.1634   0.1837
## Detection Prevalence   0.2846   0.1933   0.1746   0.1637   0.1838
## Balanced Accuracy      0.9996   0.9993   0.9994   0.9983   0.9996
```

```r
plot(modFitB1)
```

**modFitB1**



```r
plot(cmrf$table, col = cmtree$byClass, main = paste("Random Forest Confusion Matrix: Accuracy =", round
```

**Random Forest Confusion Matrix: Accuracy = 0.9989**



## Prediction with Generalized Boosted Regression

```r
set.seed(12345)

fitControl <- trainControl(method = "repeatedcv",
                           number = 5,
                           repeats = 1)

gbmFit1 <- train(classe ~ ., data=myTraining, method = "gbm",
                 trControl = fitControl,
                 verbose = FALSE)



gbmFinMod1 <- gbmFit1$finalModel

gbmPredTest <- predict(gbmFit1, newdata=myTesting)
gbmAccuracyTest <- confusionMatrix(gbmPredTest, myTesting$classe)
gbmAccuracyTest
```
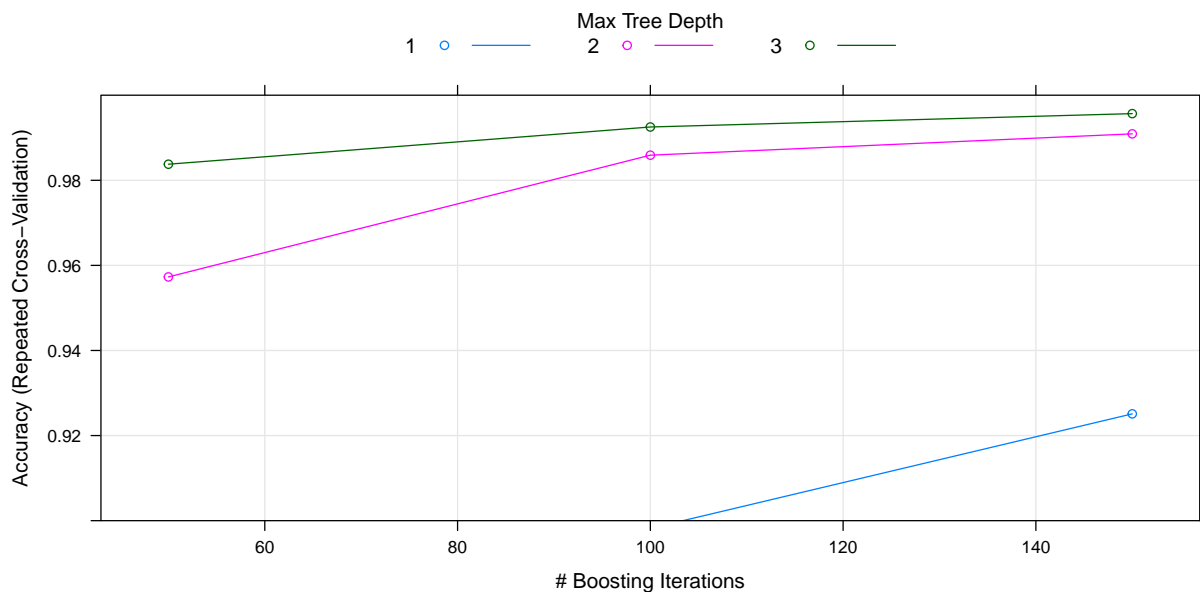
```
## Confusion Matrix and Statistics
##
##           Reference
## Prediction    A    B    C    D    E
##          A 2231    4    0    0    0
##          B    1 1512    1    0    0
##          C    0    2 1361    3    0
##          D    0    0    6 1274    1
##          E    0    0    0    9 1441
##
## Overall Statistics
##
##                Accuracy : 0.9966
##                  95% CI : (0.995, 0.9977)
##     No Information Rate : 0.2845
##     P-Value [Acc > NIR] : < 2.2e-16
```

```
##
##                   Kappa : 0.9956
##   Mcnemar's Test P-Value : NA
##
## Statistics by Class:
##
##                     Class: A Class: B Class: C Class: D Class: E
## Sensitivity           0.9996   0.9960   0.9949   0.9907   0.9993
## Specificity           0.9993   0.9997   0.9992   0.9989   0.9986
## Pos Pred Value        0.9982   0.9987   0.9963   0.9945   0.9938
## Neg Pred Value        0.9998   0.9991   0.9989   0.9982   0.9998
## Prevalence            0.2845   0.1935   0.1744   0.1639   0.1838
## Detection Rate        0.2843   0.1927   0.1735   0.1624   0.1837
## Detection Prevalence  0.2849   0.1930   0.1741   0.1633   0.1848
## Balanced Accuracy     0.9994   0.9979   0.9971   0.9948   0.9990
```

```
plot(gbmFit1, ylim=c(0.9, 1))
```



## Predicting Results on the Test Data Random Forests gave an Accuracy in the myTesting dataset of 99.89%, which was more accurate that what I got from the Decision Trees or GBM. The expected out-of-sample error is 100-99.89 = 0.11%.

## Generating Files to submit as answers for the Assignment

```
predictionB2 <- predict(modFitB1, testing, type = "class")
predictionB2
```

```
##  1  2 31  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
```

**Write the results to a text file for submission**

```
predictionsB2 <- predict(modFitB1, testing, type = "class")
pml_write_files = function(x){
  n = length(x)
  for(i in 1:n){
    filename = paste0("problem_id_",i,".txt")
    write.table(x[i],file=filename,quote=FALSE,row.names=FALSE,col.names=FALSE)
  }
}

pml_write_files(predictionsB2)
```

```
predictionsB2
```

```
##   1  2 31  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
```