

# Machine Learning Project

*Kobra*

*08/11/2019*

Outline:

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, our goal is to use data from accelerometers on the belt, forearm, arm, and dumbbell of 6 participants and predict the quality of the activities. In this dataset, quality is described with 5 different levels: A,B,C,D and E under variable 'Classe'.

## Loading data and required libraries

Loading required libraries:

Downloading training data:

```
URL1 <- "https://d396qusza40orc.cloudfront.net/predmachlearn/pml-training.csv"
download.file(URL1, destfile = "./training.csv", mode="wb")
```

```
data <- read.csv("./training.csv")
```

Downloading data for final evaluation:

```
URL2 <- "https://d396qusza40orc.cloudfront.net/predmachlearn/pml-testing.csv"
download.file (URL2, destfile = "./testing.csv", mode= "wb")
```

```
Quiz <- read.csv("./testing.csv")
```

```
dim (data)
```

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

```
str (data)
```

```
## '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_pitch_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_pitch_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_pitch_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 NA ...
## $ max_pitch_arm      : num   NA NA NA NA NA NA NA NA NA NA NA ...
## $ max_yaw_arm        : int   NA NA NA NA NA NA NA NA NA NA NA ...
## $ min_roll_arm       : num   NA NA NA NA NA NA NA NA NA NA NA ...
## $ min_pitch_arm      : num   NA NA NA NA NA NA NA NA NA NA NA ...
## $ min_yaw_arm        : int   NA NA NA NA NA NA NA NA NA NA NA ...
## $ amplitude_roll_arm : num   NA NA NA NA NA NA NA NA NA NA NA ...
## $ amplitude_pitch_arm : num   NA NA NA NA NA NA NA NA NA NA NA ...
## $ amplitude_yaw_arm   : int   NA 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_pitch_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 NA ...
## $ max_pitch_dumbbell    : num   NA 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 NA ...
## $ min_pitch_dumbbell    : num   NA 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 NA ...
## [list output truncated]
```

## Data Cleaning

The first 7 columns do not have any information that can be related to the quality of activities. So those columns will be removed.

```
data <- data[, -c(1:7)]
dim(data)
```

```
## [1] 19622 153
```

Also, there are NAs in the dataset which need to be handled. I have decided to take out the variables that have more than %90 NAs. If there is any remaining, `na.roughfix` will be used to replace NAs with either median (numeric variables) or mode (categorical variables).

```
ColRemove <- which(colSums(is.na(data)|data=="")>0.9*dim(data)[1])
dataClean <- data[,-ColRemove]
dim (dataClean)
```

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

Checking if there is any other variable with NAs:

```
sum(is.na (dataClean))
```

```
## [1] 0
```

The outcome shows that there is no more NA in dataset, so we are safe to move to the next step without any more preprocessing.

## Modelling- Deviding dataset to test and training

1. Dividing the data to training and test set:

```
set.seed(1234)
InTrain <- createDataPartition(dataClean$classe, p=0.7, list = FALSE)
train <- dataClean [InTrain, ]
test <- dataClean [-InTrain, ]
dim (train)
```

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

```
dim (test)
```

```
## [1] 5885    53
```

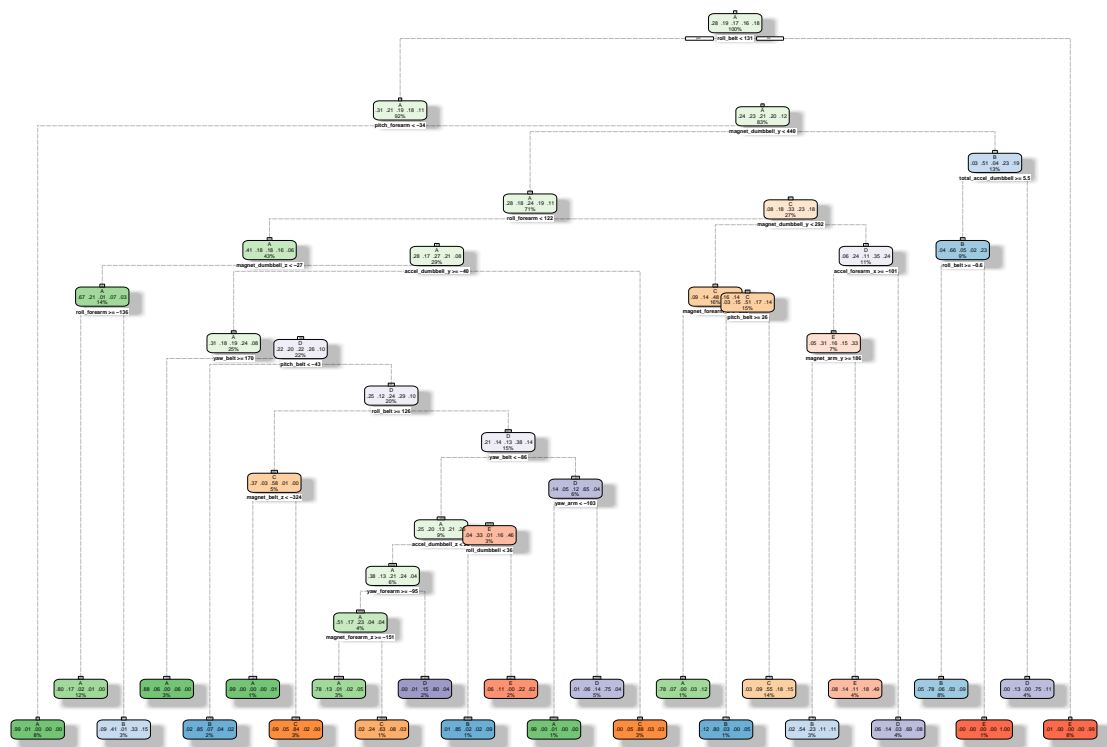
## Modelling using Decision Tree

Building the model:

```
modDecisionTree <- rpart(classe~., data= train, control = rpart.control(xval = 5))
```

```
fancyRpartPlot(modDecisionTree)
```

```
## Warning: labs do not fit even at cex 0.15, there may be some overplotting
```



Rattle 2019–Nov–18 21:32:44 Kobi

Predicting using Decision tree model:

```
PredDT <- predict(modDecisionTree,newdata=test, type = "class")
```

Estimating the accuracy of the model:

```
confMatCT <- confusionMatrix(test$classes,PredDT)
```

```
confMatCT$table
```

```
##           Reference
## Prediction   A    B    C    D    E
##           A 1522   58   47   25   22
##           B  167  706  109   94   63
##           C   12  100  819   67   28
##           D   49   79  148  609   79
##           E   13   96  139   52  782
```

```
confMatCT$overall[1]
```

```
## Accuracy
## 0.7541206
```

## Modelling using Radnom Forest

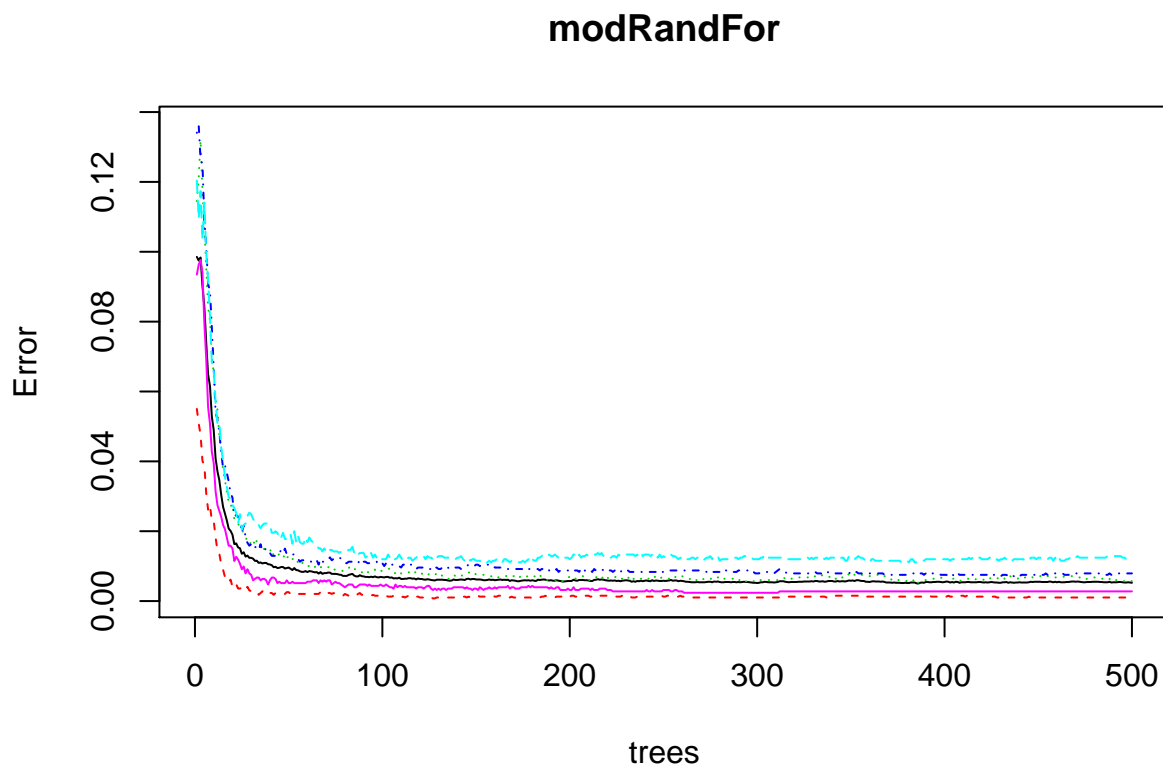
Bulding the model:

```
modRandFor <- randomForest (classe~., data= train, rfcv=rfcv(cv.fold = 5))
```

Predicting using Random forest model:

```
PredRF <- predict(modRandFor,newdata=test, type = "class")
```

```
plot (modRandFor)
```



Estimating the accuracy of the model:

```
confMatRF <- confusionMatrix(test$classe,PredRF)  
confMatRF$table
```

```
##           Reference  
## Prediction    A    B    C    D    E  
##           A 1673     1    0    0    0  
##           B   2 1133     4    0    0  
##           C   0   11 1014     1    0  
##           D   0    0    6  957     1  
##           E   0    0    0    0 1082
```

```
confMatRF$overall[1]
```

```
## Accuracy  
## 0.995582
```

## Conclusion:

Based on the results above, modelling using Random Forest algorithm gives us a very high accuracy (99.56%) which gives us enough confidence to use it for our final prediction.

## Predicting the final test data using Random Forest Model

```
PredRF_final <- predict(modRandFor,newdata=Quiz, type = "class")
```

```
final result: PredRF_final 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 > print(modRandFor)
```