# ML_Assignment4

V Ha

March 21, 2021

# Practical Machine Learning Assignment

## Synopsis

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, your goal will be to use data from accelerometers on the belt, forearm, arm, and dumbell of 6 participants. Our goal for this project is to predict a "Test" group of 6 participants with 20 activities from a training set that was supplied

## 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, your goal will be to use data from accelerometers on the belt, forearm, arm, and dumbell 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 (http://web.archive.org/web/20161224072740/http:/groupware.les.inf.puc-rio.br/har) (see the section on the Weight Lifting Exercise Dataset).

## Data Download, Processing and Cleaning

```
## Warning: package 'ggplot2' was built under R version 4.0.3
```

```
## Warning: package 'caret' was built under R version 4.0.3
```

```
## Warning: package 'corrplot' was built under R version 4.0.4
```

```
## Warning: package 'randomForest' was built under R version 4.0.4
```

```
## 'data.frame':    19622 obs. of  160 variables:
##  $ X                        : int  1 2 3 4 5 6 7 8 9 10 ...
##  $ user_name                : chr  "carlitos" "carlitos" "carlitos" "carlitos" ...
##  $ raw_timestamp_part_1     : int  1323084231 1323084231 1323084231 1323084232 1323084232 1323
084232 1323084232 1323084232 1323084232 1323084232 ...
##  $ raw_timestamp_part_2     : int  788290 808298 820366 120339 196328 304277 368296 440390 484
323 484434 ...
##  $ cvtd_timestamp           : chr  "05/12/2011 11:23" "05/12/2011 11:23" "05/12/2011 11:23" "0
5/12/2011 11:23" ...
##  $ new_window               : chr  "no" "no" "no" "no" ...
##  $ 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       : chr  "" "" "" "" ...
##  $ kurtosis_picth_belt      : chr  "" "" "" "" ...
##  $ kurtosis_yaw_belt        : chr  "" "" "" "" ...
##  $ skewness_roll_belt       : chr  "" "" "" "" ...
##  $ skewness_roll_belt.1     : chr  "" "" "" "" ...
##  $ skewness_yaw_belt        : chr  "" "" "" "" ...
##  $ 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             : chr  "" "" "" "" ...
##  $ 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             : chr  "" "" "" "" ...
##  $ 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       : chr  "" "" "" "" ...
##  $ 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        : chr  "" "" "" "" ...
##  $ kurtosis_picth_arm       : chr  "" "" "" "" ...
##  $ kurtosis_yaw_arm         : chr  "" "" "" "" ...
##  $ skewness_roll_arm        : chr  "" "" "" "" ...
##  $ skewness_pitch_arm       : chr  "" "" "" "" ...
##  $ skewness_yaw_arm         : chr  "" "" "" "" ...
##  $ 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   : chr  "" "" "" "" ...
##  $ kurtosis_picth_dumbbell  : chr  "" "" "" "" ...
##  $ kurtosis_yaw_dumbbell    : chr  "" "" "" "" ...
##  $ skewness_roll_dumbbell   : chr  "" "" "" "" ...
##  $ skewness_pitch_dumbbell  : chr  "" "" "" "" ...
##  $ skewness_yaw_dumbbell    : chr  "" "" "" "" ...
##  $ 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         : chr  "" "" "" "" ...
##  $ 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         : chr  "" "" "" "" ...
##  $ amplitude_roll_dumbbell  : num  NA NA NA NA NA NA NA NA NA NA ...
##   [list output truncated]
```

Upon inspection of the training Data set, we can see that there are multiple issues with the raw data. We can see that:

1. the first few columns are not data sets, rather identifier and descriptor columns

2. There are many columns with data but the values contained show very little variance (ie, basically the same data throughout) **Initial runs had my personal PC take a significant amount of time to process without success so this part became born out of necessity

3. There are many 'NA' values which would serve us no useful purpose

Removing descriptors

```
trainingData<-trainingData %>% select(!c(1:7))
testingData<-testingData %>% select(!c(1:7))
```

Remove new zero variance time. Online research pointed me to the below for the removal of Near zero variance time. Please see link for additional info https://www.rdocumentation.org/packages/caret/versions/6.0-86/topics/nearZeroVar (https://www.rdocumentation.org/packages/caret/versions/6.0-86/topics/nearZeroVar)

```
NZV<-nearZeroVar(trainingData)
trainingData<-trainingData[,-NZV]
testingData<-testingData[,-NZV]
```
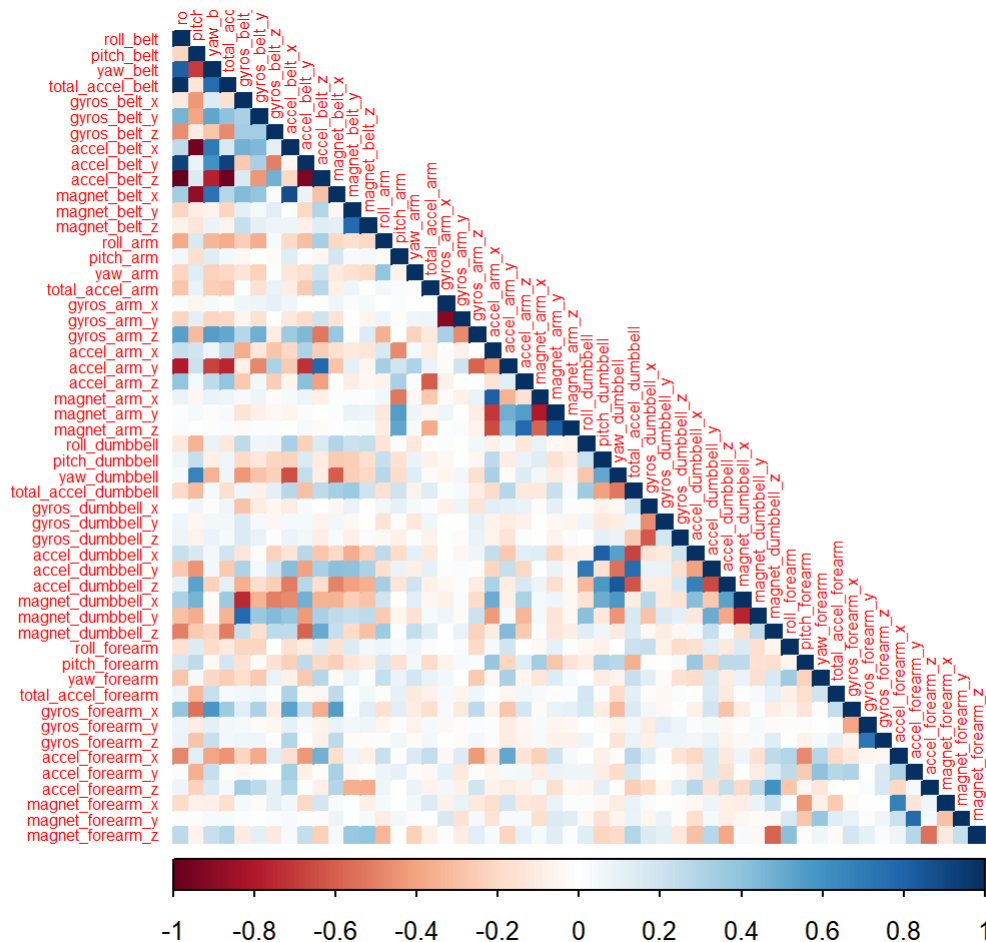
Remove all columns where all values are NA

```
trainingData<-trainingData[, colSums(is.na(trainingData)) == 0]
testingData<-testingData[, colSums(is.na(testingData)) == 0]
trainingData$classe<-factor(trainingData$classe)
testingData$problem_id<-factor(testingData$problem_id)
```

# Data Seperation for Training and Testing

A standard 70/30 split for data partitioning was created on the training data set along with a correlation plot which graphically shows us the strength of correlated variables

```
inTrain<-createDataPartition(y=trainingData$classe, p=0.7, list=F)
training<-trainingData[inTrain,]
testing<-trainingData[-inTrain,]
corData <- cor(training[,-53])
corrplot(corData, method = "color", type = "lower", tl.cex = 0.5)
```

This analysis is a part of the principle component analysis. The closer the color towards dark blue, the greater the correlation. There are minimal correlations so we move right to the predictive modelling

# Predictive modelling

Three (3) models have been selected for trialing. 1) Random Forest 2) Decision Tree 3) Gradient Boosting

# Model Generation

```
mod1<-randomForest(classe~., data=training, method='class')#random forest model
mod2<-train(classe ~ .,  data=training, method="rpart",  trControl = trainControl(method = "cv"
))#decision tree model
#generalized boosting
fitcontrol <- trainControl(method="cv",number=5,allowParallel=TRUE)
#Gradient boosting
mod3 <- train(classe ~ ., data = training, method = "gbm", trControl = fitcontrol,
              verbose = FALSE, na.action = na.omit)
```

# Prediction models on testing with accuracy

```
#predictions
pred1<-predict(mod1, testing)
pred2<-predict(mod2, testing)
pred3<-predict(mod3, testing)
#confusion matrix
accuracy1<-confusionMatrix(pred1, testing$classe)
accuracy2<-confusionMatrix(pred2, testing$classe)
accuracy3<-confusionMatrix(pred3, testing$classe)
accuracy1$overall['Accuracy']
```

```
##  Accuracy
## 0.9954121
```

```
accuracy2$overall['Accuracy']
```

```
##  Accuracy
## 0.4939677
```

```
accuracy3$overall['Accuracy']
```

```
##  Accuracy
## 0.9619371
```

With the accuracy of the random forest and GBM models being above 95%, both models would be a good predictor. The decision tree modelling has much poorer prediction accuracy below 50%

I decided to run the 2 final prediction models for both random forest and GBM on the testing data set

# Final prediction on testing set

```
predict(mod1, testingData)
```

```
##  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
```

```
predict(mod3,testingData)
```

```
##  [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
```

Outcomes are as per above with no differences between either the RF or GBM model As a reminder, outcomes are as follows: A: exactly according to the specification B: throwing the elbows to the front C: lifting the dumbbell only halfway D: lowering the dumbbell only halfway E: throwing the hips to the front