# Practical Machine Learning Course Project

*June 11, 2016*

## Overview

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.

The goal of this project is to build a machine learning algorithm to predict activity quality (classe) from activity monitors.

## Libraries

```r
library(caret)
```

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

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

```r
library(randomForest)
```

```
## randomForest 4.6-12
```

```
## Type rfNews() to see new features/changes/bug fixes.
```

```
##
## Attaching package: 'randomForest'
```

```
## The following object is masked from 'package:ggplot2':
##
##     margin
```

```r
library(corrplot)
library(rpart)
library(rpart.plot)
```

## Loading the Data

```r
trainUrl <- "https://d396qusza40orc.cloudfront.net/predmachlearn/pml-training.csv"
testUrl <- "https://d396qusza40orc.cloudfront.net/predmachlearn/pml-testing.csv"

training <- read.csv(trainUrl, na.strings=c("NA","#DIV/0!",""), header=TRUE)
testing <- read.csv(testUrl, na.strings=c("NA","#DIV/0!",""), header=TRUE)

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 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       : 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]
```

```
summary(training$classe)
```

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

## Cleaning the Data

I will remove columns with NA's, factor variables, and columns not useful as predictors

```
classe <- training$classe

#remove columns with NA's
training <- training[,colSums(is.na(training))==0]
testing <- testing[,colSums(is.na(testing))==0]

#remove factor variables
training <- training[, sapply(training, is.numeric)]
testing <- testing[, sapply(testing, is.numeric)]

#remove 1st 4 columns (X, timestamps, window)
training <- training[5:length(training)]
testing <- testing[5:length(testing)]

#add classe column back to training data set
training$classe <- classe
```

## Split Training Data for Cross Validation

Split the training data into a training data set (60%) and a validation data set (40%) that will be used for cross validation

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

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

```
dim(myTest)
```

```
## [1] 7846    53
```

## Data Modeling - Random Forest

I decided to use the random forest model to build my machine learning algorithm as it is appropriate for a classification problem as we have and based on information provided in class lectures this model tends to be more accurate than some other classification models.

We will use 5-fold cross validation when applying the algorithm.

Below I fit my model on my training data and then use my model to predict classe on my subset of data used for cross validation.

```
set.seed(1204)
controlRf <- trainControl(method="cv", 5)
modelRf <- train(classe ~ ., data=myTrain, method="rf", trControl=controlRf, ntree=250)
modelRf
```

```
## Random Forest
##
## 11776 samples
##    52 predictor
##     5 classes: 'A', 'B', 'C', 'D', 'E'
##
## No pre-processing
## Resampling: Cross-Validated (5 fold)
## Summary of sample sizes: 9422, 9421, 9420, 9420, 9421
## Resampling results across tuning parameters:
##
##   mtry  Accuracy   Kappa
##    2    0.9874319  0.9840989
##   27    0.9899791  0.9873227
##   52    0.9831852  0.9787271
##
## Accuracy was used to select the optimal model using  the largest value.
## The final value used for the model was mtry = 27.
```

```
predictRf <- predict(modelRf, myTest)
confusionMatrix(myTest$classe, predictRf)
```

```
## Confusion Matrix and Statistics
##
##           Reference
## Prediction    A    B    C    D    E
##          A 2227    2    2    1    0
##          B   14 1493   11    0    0
##          C    0   12 1349    7    0
##          D    0    1   17 1265    3
##          E    0    1    2    4 1435
##
## Overall Statistics
##
##                Accuracy : 0.9902
##                  95% CI : (0.9877, 0.9922)
##     No Information Rate : 0.2856
##     P-Value [Acc > NIR] : < 2.2e-16
##
##                   Kappa : 0.9876
##  Mcnemar's Test P-Value : NA
##
## Statistics by Class:
```

```
##
##                   Class: A Class: B Class: C Class: D Class: E
## Sensitivity          0.9938   0.9894   0.9768   0.9906   0.9979
## Specificity          0.9991   0.9961   0.9971   0.9968   0.9989
## Pos Pred Value        0.9978   0.9835   0.9861   0.9837   0.9951
## Neg Pred Value        0.9975   0.9975   0.9951   0.9982   0.9995
## Prevalence           0.2856   0.1923   0.1760   0.1628   0.1833
## Detection Rate        0.2838   0.1903   0.1719   0.1612   0.1829
## Detection Prevalence  0.2845   0.1935   0.1744   0.1639   0.1838
## Balanced Accuracy     0.9964   0.9927   0.9869   0.9937   0.9984
```

```r
#model accuarcy
accuracy <- postResample(predictRf, myTest$classe)
accuracy
```

```
##   Accuracy      Kappa
## 0.9901861 0.9875850
```

```r
accuracy1 <- accuracy[1]

#out-of-sample error
oose <- 1 - as.numeric(confusionMatrix(myTest$classe, predictRf)$overall[1])
oose
```

```
## [1] 0.009813918
```

The accuracy is 0.9901861 and the out-of-sample error is 0.0098139.
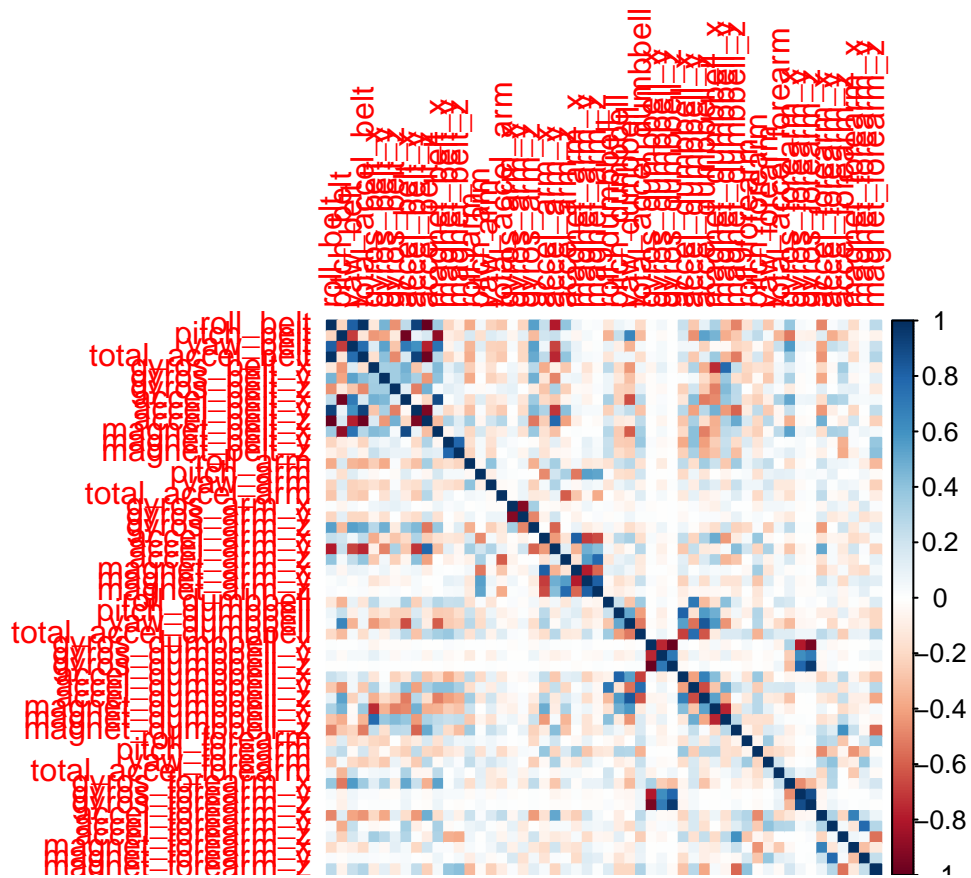
## Predicting for Test Data Set

```r
#remove problem_id from test data set
x <- length(testing)-1
testing <- testing[1:x]
final <- predict(modelRf, testing)
final
```

```
##  [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
```

## Appendix

**Correlation Matrix**

```r
x <- length(myTrain)-1
myTrain2 <- myTrain[1:x]
corrPlot <- cor(myTrain2)
corrplot(corrPlot, method="color")
```

**Decision Tree**

```
treeModel <- rpart(classe ~ ., data=myTrain, method="class")
prp(treeModel)
```

yes **roll_belt < 130** no

E

**pitch_forearm < −34**

A **magnet_dumbbell_y < 426**

**total_accel_dumbb >= 5.5**

**roll_forearm < 124**

**magnet_dumbbell_y < 290**

D

**roll_belt >= −0.6**

**magnet_dumbbell_z < −28**

**yaw_belt >= 170**

**magnet_forearm_z < −251**

**accel_forearm_x >= −100**

B

E

**roll_forearm >= −136**

B

**accel_dumbbell_y >= −40**

A

**pitch_belt >= 26**

D

A

**roll_dumbbell < 41**

A

A

C

**pitch_belt < −43**

B

C

B

E

B

**roll_belt >= 126**

B

**magnet_belt_z < −328**

D

A

C

8