# Course Project

*Kiran Joshi*

*December 10, 2017*

## Overview

We are going to analyze the personal activity monitors accelerometer information to predict the type of activity the person is doing. The data contains activities performed exactly as per specification of the exercise classified as A, and all the other errors into classed B-E as follows.

- exactly according to the specification (Class A)
- throwing the elbows to the front (Class B)
- lifting the dumbbell only halfway (Class C)
- lowering the dumbbell only halfway (Class D)
- throwing the hips to the front (Class E).

We will fit a model with the minimum features from the training data set to classify the activity in the test data set using Random Forest algorithm.

## Assumptions

- Our testing data will be a 70% split
- Model will be trained and validated on the training data, and prediction performed on the testing data
- Random forest with a 100 trees would be a good prediction model for the classification problem.

## Data Analysis

Lets obtain the data from the website, and do some introspection into the training data.

```r
trndatapth <- getURL("https://d396qusza40orc.cloudfront.net/predmachlearn/pml-training.csv")
tstdatapth <- getURL("https://d396qusza40orc.cloudfront.net/predmachlearn/pml-testing.csv")
training <- read.csv(textConnection(trndatapth),header = T)
testing <- read.csv(textConnection(tstdatapth),header = T)
dim(training)
```

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

```r
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     : Factor w/ 397 levels "","-0.016850",..: 1 1 1 1 1 1 1 1 1 1 ...
##  $ kurtosis_picth_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_picth_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_picth_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 ...
##  $ 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   : Factor w/ 398 levels "","-0.0035","-0.0073",..: 1 1 1 1 1 1 1 1 1 1 ...
##  $ kurtosis_picth_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 ...
##  $ max_picth_dumbbell       : num   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 ...
##  $ min_pitch_dumbbell       : num   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 ...
##   [list output truncated]
```

## Feature Selection

We first need to clean up the training data to remove the predictors that have NA values for all the records since they are bound to have no impact on the learning.

```
train <- training[, colSums(is.na(training)) == 0]
```

Remove all timestamp columns as the tests were random, and were independent of time factor from training data.Remove all the row Identities, as they should not make an impact on the learning, and avoid the algorithm from associating the predictors to specific identities.

```
trainRemove <- grepl("^X|timestamp|window|_id", names(train))
train <- train[, !trainRemove]
```

We will further strip all categorical variables as there are a lot of blank (not NA) values and cannot be imputed; except the dependent categorical information from the data set.

```r
trainData <- train[, sapply(train, is.numeric)]
#Adding the Classe back
trainData$classe <- train$classe
```

We quickly analyse the data for variance, to see if there are any further variables that can be eliminated. Variables that show no variance are potentially constants that do not add value in the model.

```r
nsv <- nearZeroVar(trainData,saveMetrics = T)
nsv
```

```
##                      freqRatio percentUnique zeroVar   nzv
## roll_belt             1.101904     6.7781062   FALSE FALSE
## pitch_belt            1.036082     9.3772296   FALSE FALSE
## yaw_belt              1.058480     9.9734991   FALSE FALSE
## total_accel_belt      1.063160     0.1477933   FALSE FALSE
## gyros_belt_x          1.058651     0.7134849   FALSE FALSE
## gyros_belt_y          1.144000     0.3516461   FALSE FALSE
## gyros_belt_z          1.066214     0.8612782   FALSE FALSE
## accel_belt_x          1.055412     0.8357966   FALSE FALSE
## accel_belt_y          1.113725     0.7287738   FALSE FALSE
## accel_belt_z          1.078767     1.5237998   FALSE FALSE
## magnet_belt_x         1.090141     1.6664968   FALSE FALSE
## magnet_belt_y         1.099688     1.5187035   FALSE FALSE
## magnet_belt_z         1.006369     2.3290184   FALSE FALSE
## roll_arm             52.338462    13.5256345   FALSE FALSE
## pitch_arm            87.256410    15.7323412   FALSE FALSE
## yaw_arm              33.029126    14.6570176   FALSE FALSE
## total_accel_arm       1.024526     0.3363572   FALSE FALSE
## gyros_arm_x           1.015504     3.2769341   FALSE FALSE
## gyros_arm_y           1.454369     1.9162165   FALSE FALSE
## gyros_arm_z           1.110687     1.2638875   FALSE FALSE
## accel_arm_x           1.017341     3.9598410   FALSE FALSE
## accel_arm_y           1.140187     2.7367241   FALSE FALSE
## accel_arm_z           1.128000     4.0362858   FALSE FALSE
## magnet_arm_x          1.000000     6.8239731   FALSE FALSE
## magnet_arm_y          1.056818     4.4439914   FALSE FALSE
## magnet_arm_z          1.036364     6.4468454   FALSE FALSE
## roll_dumbbell         1.022388    84.2065029   FALSE FALSE
## pitch_dumbbell        2.277372    81.7449801   FALSE FALSE
## yaw_dumbbell          1.132231    83.4828254   FALSE FALSE
## total_accel_dumbbell  1.072634     0.2191418   FALSE FALSE
## gyros_dumbbell_x      1.003268     1.2282132   FALSE FALSE
## gyros_dumbbell_y      1.264957     1.4167771   FALSE FALSE
## gyros_dumbbell_z      1.060100     1.0498420   FALSE FALSE
## accel_dumbbell_x      1.018018     2.1659362   FALSE FALSE
## accel_dumbbell_y      1.053061     2.3748853   FALSE FALSE
## accel_dumbbell_z      1.133333     2.0894914   FALSE FALSE
## magnet_dumbbell_x     1.098266     5.7486495   FALSE FALSE
## magnet_dumbbell_y     1.197740     4.3012945   FALSE FALSE
## magnet_dumbbell_z     1.020833     3.4451126   FALSE FALSE
## roll_forearm         11.589286    11.0895933   FALSE FALSE
## pitch_forearm        65.983051    14.8557741   FALSE FALSE
## yaw_forearm          15.322835    10.1467740   FALSE FALSE
## total_accel_forearm   1.128928     0.3567424   FALSE FALSE
## gyros_forearm_x       1.059273     1.5187035   FALSE FALSE
```
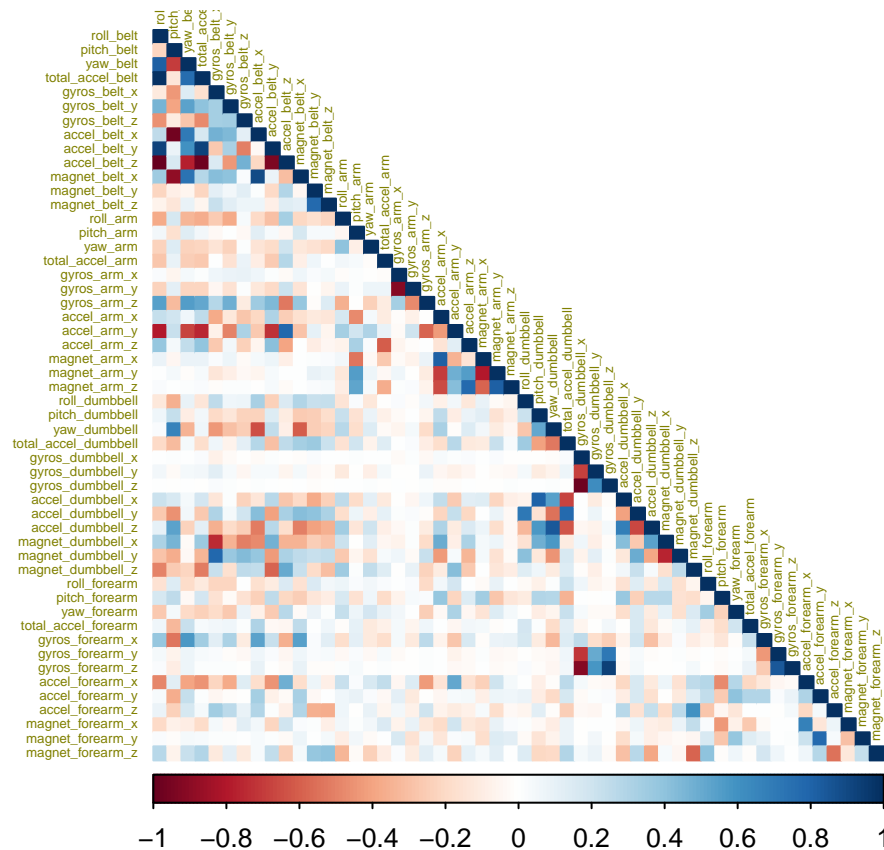
```
## gyros_forearm_y      1.036554     3.7763735    FALSE FALSE
## gyros_forearm_z      1.122917     1.5645704    FALSE FALSE
## accel_forearm_x      1.126437     4.0464784    FALSE FALSE
## accel_forearm_y      1.059406     5.1116094    FALSE FALSE
## accel_forearm_z      1.006250     2.9558659    FALSE FALSE
## magnet_forearm_x     1.012346     7.7667924    FALSE FALSE
## magnet_forearm_y     1.246914     9.5403119    FALSE FALSE
## magnet_forearm_z     1.000000     8.5771073    FALSE FALSE
## classe               1.469581     0.0254816    FALSE FALSE
```

Since none of the variables have true zero variance or near zero variance, it has passed the nsv test.

## Plotting Predictors

A correlation among variables is analysed before proceeding to the modeling procedures. We would just consider the predictors and thus the classe variable is removed.

```
trainCor <- cor(trainData[, -53])
corrplot(trainCor, order = "original", method = "color", type = "lower",tl.cex = 0.45, tl.col = rgb(.5,
```



The highly correlated variables are shown in dark colors in the graph above.As we can see, there are not too many highly correlated variables (ignoring the diagonal), and hence does not need more cleanup to avoid overfitting.

Although we have been given explicit testing and training data, we will split the data randomly to get 75% of the data found in the training data for fitting the model, so that we do not touch the test data provided for

tuning or cross-validation of the model.

# Prediction Model

We will try to predict the classe from the other variables in the dataset.

```
set.seed(54321)
inTrain <- createDataPartition(trainData$classe, p=0.75, list=FALSE)
train_data <- trainData[inTrain, ]
test_data <- trainData[-inTrain, ]
```

We will use the Random Forest method to fit the model as per our assumption

```
FitRandForest <- train(classe ~ ., data=train_data, method="rf", ntree=100)
```

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

Lets now check if we have a good fit, based on the accuracy of the model

```
FitRandForest$finalModel
```

```
##
## Call:
##  randomForest(x = x, y = y, ntree = 100, mtry = param$mtry)
##                Type of random forest: classification
##                      Number of trees: 100
## No. of variables tried at each split: 27
##
##         OOB estimate of  error rate: 0.71%
## Confusion matrix:
##      A    B    C    D    E  class.error
## A 4182    1    1    0    1 0.0007168459
## B   23 2813   12    0    0 0.0122893258
## C    0   16 2542    9    0 0.0097389949
## D    0    1   27 2381    3 0.0128524046
## E    0    2    3    6 2695 0.0040650407
```

As we can see, the error rate is 0.71%, which puts the model at 99.29% accuracy.

We can now validate the model against test_data, which is still a part of the training set and compare the predicted values against the true values using a confusion matrix.

```
PredRandForest<- predict(FitRandForest, newdata = test_data)
confusionMatrix(PredRandForest,test_data$classe)
```

```
## Confusion Matrix and Statistics
##
##           Reference
```

```
## Prediction    A    B    C    D    E
##         A 1395    9    0    0    0
##         B    0  939    5    0    1
##         C    0    1  847    5    3
##         D    0    0    3  799    5
##         E    0    0    0    0  892
##
## Overall Statistics
##
##                Accuracy : 0.9935
##                  95% CI : (0.9908, 0.9955)
##     No Information Rate : 0.2845
##     P-Value [Acc > NIR] : < 2.2e-16
##
##                   Kappa : 0.9917
##  Mcnemar's Test P-Value : NA
##
## Statistics by Class:
##
##                      Class: A Class: B Class: C Class: D Class: E
## Sensitivity            1.0000   0.9895   0.9906   0.9938   0.9900
## Specificity            0.9974   0.9985   0.9978   0.9980   1.0000
## Pos Pred Value         0.9936   0.9937   0.9895   0.9901   1.0000
## Neg Pred Value         1.0000   0.9975   0.9980   0.9988   0.9978
## Prevalence             0.2845   0.1935   0.1743   0.1639   0.1837
## Detection Rate         0.2845   0.1915   0.1727   0.1629   0.1819
## Detection Prevalence   0.2863   0.1927   0.1746   0.1646   0.1819
## Balanced Accuracy      0.9987   0.9940   0.9942   0.9959   0.9950
```

With 99.35% accuracy, we should be now confident to test it against the testing data.

```
features <- names(trainData)
features <- features[-53]
testData <- testing[,features]
PredtestData <- predict(FitRandForest, newdata = testData)
PredtestData
```

```
##  [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
```

# Conclusion

The Testing data Classes were predicted, and the output used to complete the Quiz section of the project.