

# Practical Machine Learning Course Project

*Ong Kwee Hian*

*Saturday, August 22, 2015*

## Introduction

In this project, the goal is to use data from accelerometers on the belt, forearm, arm, and dumbbell of 6 participants who were asked to perform barbell lifts correctly and incorrectly in 5 different ways as training data set to develop a model, and finally use it to predict the manner in which they did the exercise in the 20 test data set observations.

## Load and prepare/clean Training data

Attempt to first investigate the characteristics of raw training data set.

```
library(lattice)
library(ggplot2)
library(caret)

rawTrainData <- read.csv("pml-training.csv") # Training data set
rawTestData <- read.csv("pml-testing.csv")   # Test data set

str(rawTrainData)
```

```
## '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 ...
## $ 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 4 ...
## $ cvtd_timestamp        : Factor w/ 20 levels "02/12/2011 13:32",...: 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 ...
## $ 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 ...
## $ kurtosis_pitch_belt   : Factor w/ 317 levels "", "-0.021887",...: 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 ...
## $ skewness_roll_belt    : Factor w/ 395 levels "", "-0.003095",...: 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 NA ...
## $ max_pitch_belt : int NA 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 NA ...
## $ min_pitch_belt : int NA 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 NA ...
## $ amplitude_pitch_belt : int NA 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 NA ...
## $ avg_roll_belt : num NA NA NA NA NA NA NA NA NA NA NA ...
## $ stddev_roll_belt : num NA NA NA NA NA NA NA NA NA NA NA ...
## $ var_roll_belt : num NA NA NA NA NA NA NA NA NA NA NA ...
## $ avg_pitch_belt : num NA NA NA NA NA NA NA NA NA NA NA ...
## $ stddev_pitch_belt : num NA NA NA NA NA NA NA NA NA NA NA ...
## $ var_pitch_belt : num NA NA NA NA NA NA NA NA NA NA NA ...
## $ avg_yaw_belt : num NA NA NA NA NA NA NA NA NA NA NA ...
## $ stddev_yaw_belt : num NA NA NA NA NA NA NA NA NA NA NA ...
## $ var_yaw_belt : num NA 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 NA ...
## $ avg_roll_arm : num NA NA NA NA NA NA NA NA NA NA NA ...
## $ stddev_roll_arm : num NA NA NA NA NA NA NA NA NA NA NA ...
## $ var_roll_arm : num NA NA NA NA NA NA NA NA NA NA NA ...
## $ avg_pitch_arm : num NA NA NA NA NA NA NA NA NA NA NA ...
## $ stddev_pitch_arm : num NA NA NA NA NA NA NA NA NA NA NA ...
## $ var_pitch_arm : num NA NA NA NA NA NA NA NA NA NA NA ...
## $ avg_yaw_arm : num NA NA NA NA NA NA NA NA NA NA NA ...
## $ stddev_yaw_arm : num NA NA NA NA NA NA NA NA NA NA NA ...
## $ var_yaw_arm : num NA 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 -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 ...
## $ kurtosis_pitch_arm : Factor w/ 328 levels "", "-0.00484", ...: 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 ...
## $ skewness_roll_arm  : Factor w/ 331 levels "", "-0.00051", ...: 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 ...
## $ skewness_yaw_arm   : Factor w/ 395 levels "", "-0.00311", ...: 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 ...
## $ 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 ...
## $ 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 ...
## $ 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 ...
## $ amplitude_roll_dumbbell : num  NA NA NA NA NA NA NA NA NA NA NA ...
## [list output truncated]

```

There a quite a large number of NA and also variable not useful for analysis. To clean away unnecessary variables, i will remove variables with nearly zero variance, variables that has NA, and variables that are irrelevant for prediction e.g. user\_name, time stanp, etc.

```

# remove near zero variance variables
nzv <- nearZeroVar(rawTrainData)
rawTrainData <- rawTrainData[, -nzv]

# remove NA variables
NAs <- apply(rawTrainData, 2, function(x) { sum(is.na(x)) })
rawTrainData <- rawTrainData[, which(NAs == 0)]

# remove first five irrelevant variables (X, user_name, raw_timestamp_part_1, raw_timestamp_part_2, cv
td_timestamp)
rawTrainData <- rawTrainData[, -(1:5)]

str(rawTrainData)

```

```

## 'data.frame': 19622 obs. of 54 variables:
## $ 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 ...
## $ 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 ...
## $ 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 ...
## $ 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 ...
## $ total_accel_dumbbell: int 37 37 37 37 37 37 37 37 37 37 ...

```

```
## $ gyros_dumbbell_x      : num  0 0 0 0 0 0 0 0 0 0 ...
## $ gyros_dumbbell_y      : num  -0.02 -0.02 -0.02 -0.02 -0.02 -0.02 -0.02 -0.02 -0.02 -0.02 ...
## $ gyros_dumbbell_z      : num  0 0 0 -0.02 0 0 0 0 0 0 ...
## $ accel_dumbbell_x      : int   -234 -233 -232 -232 -233 -234 -232 -234 -232 -235 ...
## $ accel_dumbbell_y      : int    47 47 46 48 48 48 47 46 47 48 ...
## $ accel_dumbbell_z      : int   -271 -269 -270 -269 -270 -269 -270 -272 -269 -270 ...
## $ magnet_dumbbell_x     : int   -559 -555 -561 -552 -554 -558 -551 -555 -549 -558 ...
## $ magnet_dumbbell_y     : int    293 296 298 303 292 294 295 300 292 291 ...
## $ magnet_dumbbell_z     : num   -65 -64 -63 -60 -68 -66 -70 -74 -65 -69 ...
## $ roll_forearm          : num   28.4 28.3 28.3 28.1 28 27.9 27.9 27.8 27.7 27.7 ...
## $ pitch_forearm         : num   -63.9 -63.9 -63.9 -63.9 -63.9 -63.9 -63.9 -63.9 -63.8 -63.8 ...
## $ yaw_forearm           : num   -153 -153 -152 -152 -152 -152 -152 -152 -152 -152 ...
## $ total_accel_forearm   : int    36 36 36 36 36 36 36 36 36 36 ...
## $ gyros_forearm_x       : num    0.03 0.02 0.03 0.02 0.02 0.02 0.02 0.02 0.03 0.02 ...
## $ gyros_forearm_y       : num    0 0 -0.02 -0.02 0 -0.02 0 -0.02 0 0 ...
## $ gyros_forearm_z       : num   -0.02 -0.02 0 0 -0.02 -0.03 -0.02 0 -0.02 -0.02 ...
## $ accel_forearm_x       : int   192 192 196 189 189 193 195 193 193 190 ...
## $ accel_forearm_y       : int   203 203 204 206 206 203 205 205 204 205 ...
## $ accel_forearm_z       : int   -215 -216 -213 -214 -214 -215 -215 -213 -214 -215 ...
## $ magnet_forearm_x      : int   -17 -18 -18 -16 -17 -9 -18 -9 -16 -22 ...
## $ magnet_forearm_y      : num   654 661 658 658 655 660 659 660 653 656 ...
## $ magnet_forearm_z      : num   476 473 469 469 473 478 470 474 476 473 ...
## $ classe                : Factor w/ 5 levels "A","B","C","D",...: 1 1 1 1 1 1 1 1 1 1 ...
```

To be able to estimate the out-of-sample error, I will make use of this training set by splitting the Full training data set (rawTrainData) into 2 smaller sets for use, one for training purpose (rawTrainData1) and the other for validation purpose (rawTrainData2):

```
set.seed(082015)
inTrain <- createDataPartition(y=rawTrainData$classe, p=0.7, list=F)
rawTrainData1 <- rawTrainData[inTrain, ]      # For training
rawTrainData2 <- rawTrainData[-inTrain, ]     # For testing
```

## Build Model

Applying Random Forest model on rawTrainData1, and “train” using 3-fold cross-validation to select optimal tuning parameters for acceptable performance for the model.

```
# Use 3-fold cross-validation to select optimal tuning parameters
modControl <- trainControl(method="cv", number=3, verboseIter=FALSE)

# fit model on rawTrainData1
modFit <- train(classe ~ ., data=rawTrainData1, method="rf", trControl=modControl)
```

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

```
## Warning: package 'randomForest' was built under R version 3.1.3
```

```
## randomForest 4.6-10  
## Type rfNews() to see new features/changes/bug fixes.
```

```
# print tuning parameters choosen in final model  
modFit$finalModel
```

```
##  
## Call:  
## randomForest(x = x, y = y, mtry = param$mtry)  
##           Type of random forest: classification  
##           Number of trees: 500  
## No. of variables tried at each split: 27  
##  
##           OOB estimate of  error rate: 0.2%  
## Confusion matrix:  
##      A    B    C    D    E  class.error  
## A 3904     1     0     0     1 0.0005120328  
## B   6 2650     2     0     0 0.0030097818  
## C    0   6 2389     1     0 0.0029215359  
## D    0    0   6 2246     0 0.0026642984  
## E    0    1    0    4 2520 0.0019801980
```

The result return is 500 trees and try 27 variables at each split.

## Model Evaluation and Selection

This model is used to predict the label (“classe”) in rawTrainData2, and confusion matrix is used to compare the predicted versus the actual labels:

```
# Predict classe label in validation set (rawTrainData2) using model  
predictions <- predict(modFit, newdata=rawTrainData2)  
  
# Print confusion matrix to understand out-of-sample error  
confusionMatrix(rawTrainData2$classe, predictions)
```

```

## Confusion Matrix and Statistics
##
##           Reference
## Prediction    A    B    C    D    E
##           A 1674    0    0    0    0
##           B    1 1135    3    0    0
##           C    0    0 1026    0    0
##           D    0    0    3  961    0
##           E    0    0    0    1 1081
##
## Overall Statistics
##
##           Accuracy : 0.9986
##           95% CI : (0.9973, 0.9994)
##           No Information Rate : 0.2846
##           P-Value [Acc > NIR] : < 2.2e-16
##
##           Kappa : 0.9983
##           McNemar's Test P-Value : NA
##
## Statistics by Class:
##
##           Class: A Class: B Class: C Class: D Class: E
## Sensitivity          0.9994   1.0000   0.9942   0.9990   1.0000
## Specificity          1.0000   0.9992   1.0000   0.9994   0.9998
## Pos Pred Value       1.0000   0.9965   1.0000   0.9969   0.9991
## Neg Pred Value       0.9998   1.0000   0.9988   0.9998   1.0000
## Prevalence           0.2846   0.1929   0.1754   0.1635   0.1837
## Detection Rate       0.2845   0.1929   0.1743   0.1633   0.1837
## Detection Prevalence 0.2845   0.1935   0.1743   0.1638   0.1839
## Balanced Accuracy     0.9997   0.9996   0.9971   0.9992   0.9999

```

The accuracy is 99.8%, thus my predicted accuracy for the out-of-sample error with cross validation is 0.2%, which is very good.

## Re-training the Model using Full training data

As the result is very good, Random Forests will be used to predict labels for the test set (20 test cases) provided. I will now retrain the model using the Full training set (rawTrainData) before predicting on the test set given.

```
# re-model using full training set (rawTrainData)
modControl <- trainControl(method="cv", number=3, verboseIter=F)
modFit <- train(classe ~ ., data=rawTrainData, method="rf", trControl=modControl)

# print tuning parameters choosen in final model
modFit$finalModel
```

```
##
## Call:
## randomForest(x = x, y = y, mtry = param$mtry)
##           Type of random forest: classification
##           Number of trees: 500
## No. of variables tried at each split: 27
##
##           OOB estimate of  error rate: 0.13%
## Confusion matrix:
##      A    B    C    D    E  class.error
## A 5578     1     0     0     1 0.0003584229
## B   6378     2     1     0 0.0023702923
## C    43418     0     0 0.0011689071
## D    83207     1 0.0027985075
## E   23605 0.0005544774
```

## Using the model for Test Set Predictions

The new model fit is applied on activity monitors in rawTestData to predict the activity quality label, and output to files:

```
# predict on test set
predictions <- predict(modFit, newdata=rawTestData)

# convert predictions to list of vector
predictions <- as.character(predictions)

# create function to write predictions to files
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=F, row.names=F, col.names=F)
  }
}

# create prediction files to submit
pml_write_files(predictions)
```



The created result file input are used for project submission. As the model take a long time to generate, the model output are only show in html file.