

# Practical Machine Learning Final Assignment

*Yann K*

*June 12, 2016*

## 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 dumbbell 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://groupware.les.inf.puc-rio.br/har> (see the section on the Weight Lifting Exercise Dataset).

## Data loading and model fitting

```
set.seed(99)
library(caret)
library(plyr)
library(parallel)
library(doParallel)
```

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

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

```
teurl<-"https://d396qusza40orc.cloudfront.net/predmachlearn/pml-testing.csv"
trurl<-"https://d396qusza40orc.cloudfront.net/predmachlearn/pml-training.csv"
tedf<-read.csv(teurl)
trdf<-read.csv(trurl)
names(tedf)
```

```
##      [1] "X"                                "user_name"
##      [3] "raw_timestamp_part_1"          "raw_timestamp_part_2"
##      [5] "cvtd_timestamp"              "new_window"
##      [7] "num_window"                  "roll_belt"
##      [9] "pitch_belt"                  "yaw_belt"
##     [11] "total_accel_belt"            "kurtosis_roll_belt"
##     [13] "kurtosis_pitch_belt"         "kurtosis_yaw_belt"
##     [15] "skewness_roll_belt"          "skewness_roll_belt.1"
##     [17] "skewness_yaw_belt"           "max_roll_belt"
##     [19] "max_pitch_belt"              "max_yaw_belt"
##     [21] "min_roll_belt"               "min_pitch_belt"
##     [23] "min_yaw_belt"                "amplitude_roll_belt"
##     [25] "amplitude_pitch_belt"        "amplitude_yaw_belt"
```

## [27]	"var_total_accel_belt"	"avg_roll_belt"
## [29]	"stddev_roll_belt"	"var_roll_belt"
## [31]	"avg_pitch_belt"	"stddev_pitch_belt"
## [33]	"var_pitch_belt"	"avg_yaw_belt"
## [35]	"stddev_yaw_belt"	"var_yaw_belt"
## [37]	"gyros_belt_x"	"gyros_belt_y"
## [39]	"gyros_belt_z"	"accel_belt_x"
## [41]	"accel_belt_y"	"accel_belt_z"
## [43]	"magnet_belt_x"	"magnet_belt_y"
## [45]	"magnet_belt_z"	"roll_arm"
## [47]	"pitch_arm"	"yaw_arm"
## [49]	"total_accel_arm"	"var_accel_arm"
## [51]	"avg_roll_arm"	"stddev_roll_arm"
## [53]	"var_roll_arm"	"avg_pitch_arm"
## [55]	"stddev_pitch_arm"	"var_pitch_arm"
## [57]	"avg_yaw_arm"	"stddev_yaw_arm"
## [59]	"var_yaw_arm"	"gyros_arm_x"
## [61]	"gyros_arm_y"	"gyros_arm_z"
## [63]	"accel_arm_x"	"accel_arm_y"
## [65]	"accel_arm_z"	"magnet_arm_x"
## [67]	"magnet_arm_y"	"magnet_arm_z"
## [69]	"kurtosis_roll_arm"	"kurtosis_pitch_arm"
## [71]	"kurtosis_yaw_arm"	"skewness_roll_arm"
## [73]	"skewness_pitch_arm"	"skewness_yaw_arm"
## [75]	"max_roll_arm"	"max_pitch_arm"
## [77]	"max_yaw_arm"	"min_roll_arm"
## [79]	"min_pitch_arm"	"min_yaw_arm"
## [81]	"amplitude_roll_arm"	"amplitude_pitch_arm"
## [83]	"amplitude_yaw_arm"	"roll_dumbbell"
## [85]	"pitch_dumbbell"	"yaw_dumbbell"
## [87]	"kurtosis_roll_dumbbell"	"kurtosis_pitch_dumbbell"
## [89]	"kurtosis_yaw_dumbbell"	"skewness_roll_dumbbell"
## [91]	"skewness_pitch_dumbbell"	"skewness_yaw_dumbbell"
## [93]	"max_roll_dumbbell"	"max_pitch_dumbbell"
## [95]	"max_yaw_dumbbell"	"min_roll_dumbbell"
## [97]	"min_pitch_dumbbell"	"min_yaw_dumbbell"
## [99]	"amplitude_roll_dumbbell"	"amplitude_pitch_dumbbell"
## [101]	"amplitude_yaw_dumbbell"	"total_accel_dumbbell"
## [103]	"var_accel_dumbbell"	"avg_roll_dumbbell"
## [105]	"stddev_roll_dumbbell"	"var_roll_dumbbell"
## [107]	"avg_pitch_dumbbell"	"stddev_pitch_dumbbell"
## [109]	"var_pitch_dumbbell"	"avg_yaw_dumbbell"
## [111]	"stddev_yaw_dumbbell"	"var_yaw_dumbbell"
## [113]	"gyros_dumbbell_x"	"gyros_dumbbell_y"
## [115]	"gyros_dumbbell_z"	"accel_dumbbell_x"
## [117]	"accel_dumbbell_y"	"accel_dumbbell_z"
## [119]	"magnet_dumbbell_x"	"magnet_dumbbell_y"
## [121]	"magnet_dumbbell_z"	"roll_forearm"
## [123]	"pitch_forearm"	"yaw_forearm"
## [125]	"kurtosis_roll_forearm"	"kurtosis_pitch_forearm"
## [127]	"kurtosis_yaw_forearm"	"skewness_roll_forearm"
## [129]	"skewness_pitch_forearm"	"skewness_yaw_forearm"
## [131]	"max_roll_forearm"	"max_pitch_forearm"
## [133]	"max_yaw_forearm"	"min_roll_forearm"

```
## [135] "min_pitch_forearm"      "min_yaw_forearm"
## [137] "amplitude_roll_forearm" "amplitude_pitch_forearm"
## [139] "amplitude_yaw_forearm"  "total_accel_forearm"
## [141] "var_accel_forearm"      "avg_roll_forearm"
## [143] "stddev_roll_forearm"    "var_roll_forearm"
## [145] "avg_pitch_forearm"      "stddev_pitch_forearm"
## [147] "var_pitch_forearm"      "avg_yaw_forearm"
## [149] "stddev_yaw_forearm"     "var_yaw_forearm"
## [151] "gyros_forearm_x"        "gyros_forearm_y"
## [153] "gyros_forearm_z"        "accel_forearm_x"
## [155] "accel_forearm_y"        "accel_forearm_z"
## [157] "magnet_forearm_x"       "magnet_forearm_y"
## [159] "magnet_forearm_z"       "problem_id"
```

```
#Removing the first columns (X, username, timestamp, window) having non predictive data
tedf<-tedf[,8:160]
trdf<-trdf[,8:160]

dim(trdf)
```

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

```
#Checking for NAs in the training dataset
count(trdf[is.na(trdf)])
```

```
##      x      freq
## 1 <NA> 1287472
```

```
#Removing NA variables
trdfna<-trdf[,colSums(is.na(trdf))==0]
#Removing near zero variance columns
nz<-nearZeroVar(trdfna, saveMetrics=TRUE)
trdfna<-trdfna[, !as.logical(nz$nzv)]
dim(trdfna)
```

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

```
names(trdfna)
```

```
## [1] "roll_belt"      "pitch_belt"      "yaw_belt"
## [4] "total_accel_belt" "gyros_belt_x"    "gyros_belt_y"
## [7] "gyros_belt_z"    "accel_belt_x"    "accel_belt_y"
## [10] "accel_belt_z"    "magnet_belt_x"   "magnet_belt_y"
## [13] "magnet_belt_z"   "roll_arm"        "pitch_arm"
## [16] "yaw_arm"         "total_accel_arm" "gyros_arm_x"
## [19] "gyros_arm_y"     "gyros_arm_z"     "accel_arm_x"
## [22] "accel_arm_y"     "accel_arm_z"     "magnet_arm_x"
## [25] "magnet_arm_y"    "magnet_arm_z"    "roll_dumbbell"
## [28] "pitch_dumbbell"  "yaw_dumbbell"    "total_accel_dumbbell"
## [31] "gyros_dumbbell_x" "gyros_dumbbell_y" "gyros_dumbbell_z"
## [34] "accel_dumbbell_x" "accel_dumbbell_y" "accel_dumbbell_z"
```

```

## [37] "magnet_dumbbell_x"      "magnet_dumbbell_y"      "magnet_dumbbell_z"
## [40] "roll_forearm"          "pitch_forearm"          "yaw_forearm"
## [43] "total_accel_forearm"    "gyros_forearm_x"        "gyros_forearm_y"
## [46] "gyros_forearm_z"        "accel_forearm_x"        "accel_forearm_y"
## [49] "accel_forearm_z"        "magnet_forearm_x"        "magnet_forearm_y"
## [52] "magnet_forearm_z"       "classe"

#Create training dataframe subset with only 20% of the data for performance on my laptop
it<-createDataPartition(trdfna$classe, p=.2, list=F)
training<-trdfna[it,]
testing<-trdfna[-it,]

#Parallel calculation clusters
cluster<-makeCluster(detectCores()-1)
registerDoParallel(cluster)

#Random Forest model fitting
tparam<-trainControl(allowParallel=TRUE, number=5)
trmodel<-train(classe~., data=training, method="rf", importance=TRUE, trainControl=tparam)

## 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

stopCluster(cluster)
trmodel

## Random Forest
##
## 3927 samples
## 52 predictor
## 5 classes: 'A', 'B', 'C', 'D', 'E'
##
## No pre-processing
## Resampling: Bootstrapped (25 reps)
## Summary of sample sizes: 3927, 3927, 3927, 3927, 3927, ...
## Resampling results across tuning parameters:
##
##  mtry  Accuracy  Kappa
##    2    0.9525031 0.9398955
##   27    0.9560089 0.9443492
##   52    0.9416067 0.9261450
##
## Accuracy was used to select the optimal model using the largest value.
## The final value used for the model was mtry = 27.

```

```
trmodel$finalModel
```

```
##
## Call:
## randomForest(x = x, y = y, mtry = param$mtry, importance = TRUE,      trainControl = ..2)
##           Type of random forest: classification
##           Number of trees: 500
## No. of variables tried at each split: 27
##
##           OOB estimate of  error rate: 2.85%
## Confusion matrix:
##      A   B   C   D   E class.error
## A 1106   7   1   2   0 0.008960573
## B   19 727  11   2   1 0.043421053
## C    1  20 658   5   1 0.039416058
## D    1   0  19 619   5 0.038819876
## E    0   3   8   6 705 0.023545706
```

```
#Evaluate training model on the testing dataset
```

```
pred<-predict(trmodel, testing)
confusionMatrix(pred, testing$classe)
```

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

```
##
##           Reference
## Prediction    A    B    C    D    E
##           A 4426   99    0    4    0
##           B   21 2898  100   15   15
##           C    2   30 2626   49   12
##           D   14    7   11 2494   10
##           E    1    3    0   10 2848
```

```
## Overall Statistics
```

```
##
##           Accuracy : 0.9743
##           95% CI : (0.9717, 0.9767)
## No Information Rate : 0.2844
## P-Value [Acc > NIR] : < 2.2e-16
```

```
##
##           Kappa : 0.9675
## McNemar's Test P-Value : < 2.2e-16
```

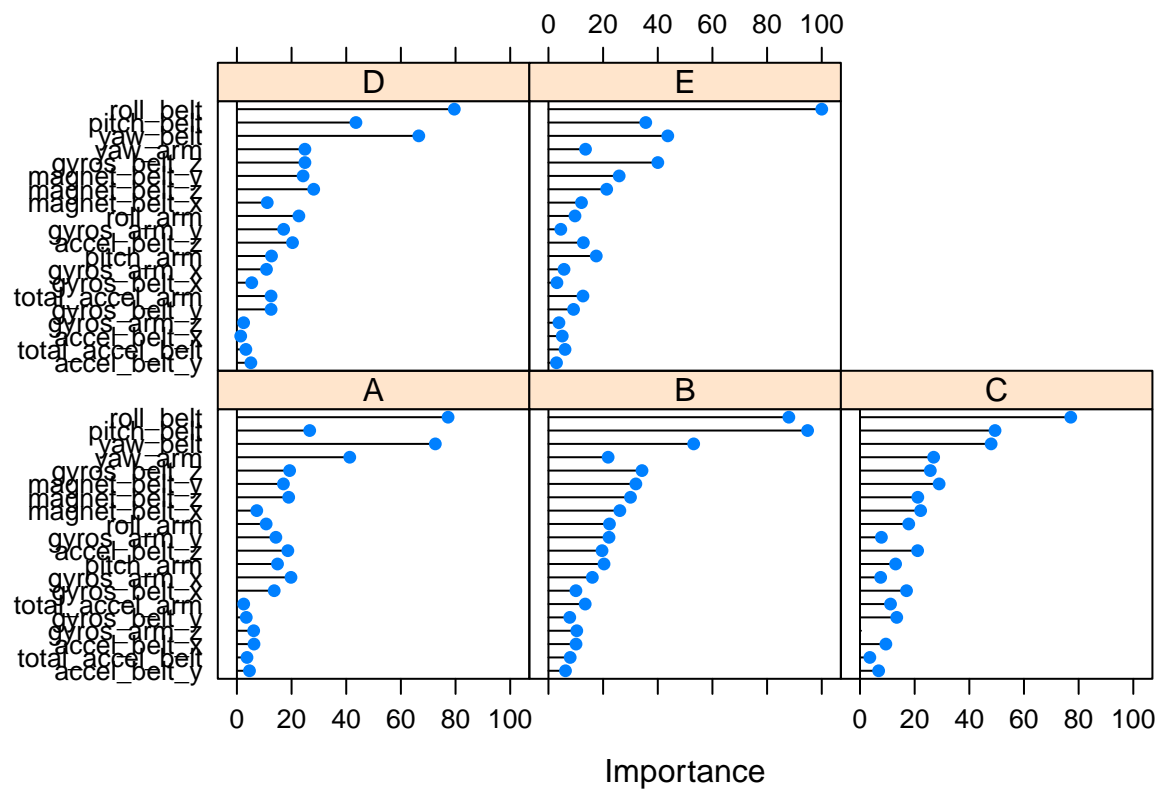
```
##
```

```
## Statistics by Class:
```

```
##
```

```
##           Class: A Class: B Class: C Class: D Class: E
## Sensitivity      0.9915  0.9542  0.9594  0.9697  0.9872
## Specificity      0.9908  0.9881  0.9928  0.9968  0.9989
## Pos Pred Value   0.9773  0.9505  0.9658  0.9834  0.9951
## Neg Pred Value   0.9966  0.9890  0.9914  0.9941  0.9971
## Prevalence       0.2844  0.1935  0.1744  0.1639  0.1838
## Detection Rate   0.2820  0.1846  0.1673  0.1589  0.1815
## Detection Prevalence 0.2886  0.1943  0.1732  0.1616  0.1824
## Balanced Accuracy 0.9912  0.9712  0.9761  0.9832  0.9930
```

```
#Plotting the top 20 most important predictors
#The most important predictors are roll_belt, pitch_belt, yaw_belt
vi<-varImp(trmodel)
vi[[1]]<-vi[[1]][1:20,]
plot(vi)
```



## Applying model to the testing dataset

```
pred<-predict(trmodel, tedf)
pred
```

```
## [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
```