

Practical Machine Learning Assignment - week 3

Francois Ragnet

Sunday, September 27, 2015

Summary

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, we used 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. We built a Machine Learning model to try and predict the classe of outcome, then tested on a few validation sets. More information is available from the website here: <http://groupware.les.inf.puc-rio.br/har> (<http://groupware.les.inf.puc-rio.br/har>) (see the section on the Weight Lifting Exercise Dataset).

Data Processing and Analysis

Loading and preprocessing the data

We will download the main file from <https://d396qusza40orc.cloudfront.net/predmachlearn/pml-training.csv> (<https://d396qusza40orc.cloudfront.net/predmachlearn/pml-training.csv>) and the much smaller validation file from <https://d396qusza40orc.cloudfront.net/predmachlearn/pml-testing.csv> (<https://d396qusza40orc.cloudfront.net/predmachlearn/pml-testing.csv>). We will be using the first dataset for our model building (split into training and testing), then the second dataset as evaluation.

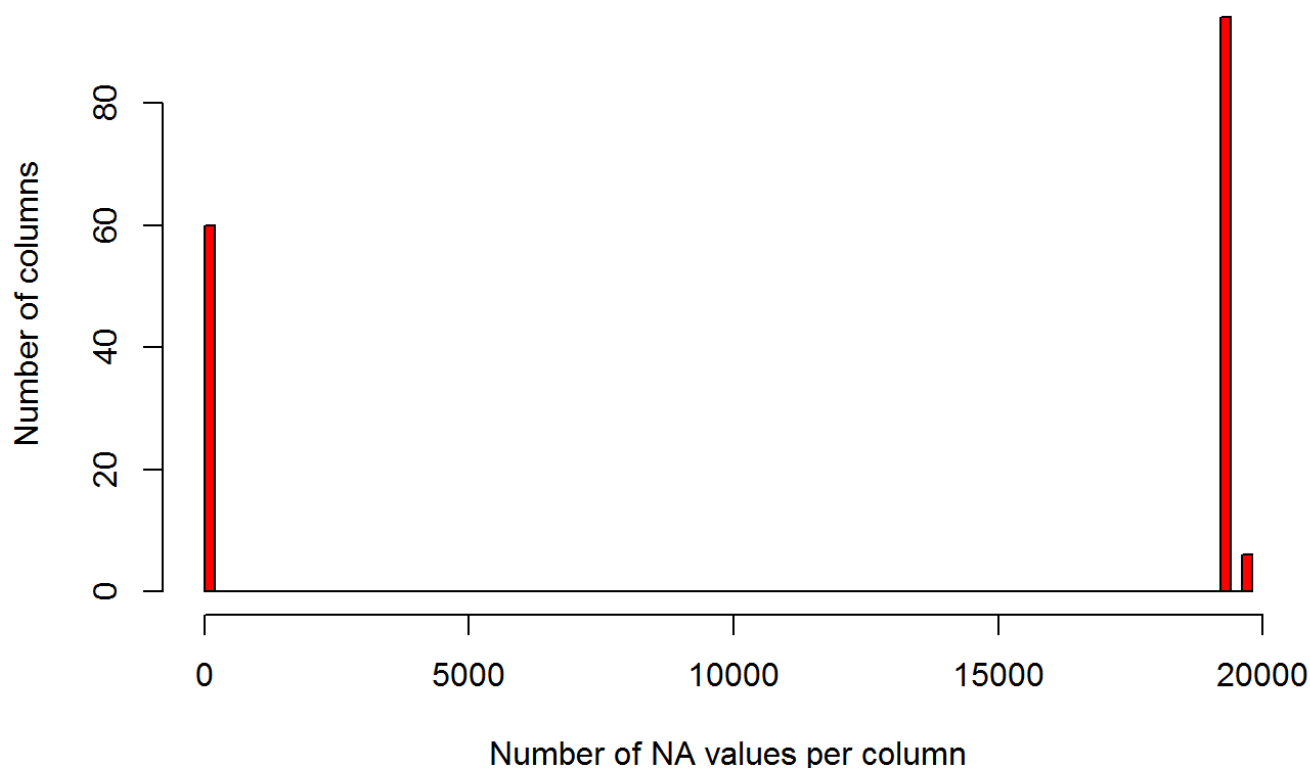
```
# setwd("C:\\Local\\My Local Documents\\Training\\Data Analytics\\Practical machine Learning\\Project")
# download.file("https://d396qusza40orc.cloudfront.net/predmachlearn/pml-training.csv",
# dest="pml-training.csv")
# download.file("https://d396qusza40orc.cloudfront.net/predmachlearn/pml-testing.csv",
# dest="pml-testing.csv")
#pmlMainFULL<-read.csv("pml-training.csv", na.strings=c("NA", "#DIV/0!", ""), stringsAsFactors=FALSE)
pmlMainFULL<-read.csv("pml-training.csv", na.strings=c("NA", "#DIV/0!", ""), stringsAsFactors=TRUE)
pmlValidation<-read.csv("pml-testing.csv", na.strings=c("NA", "#DIV/0!", ""), , stringsAsFactors=FALSE)
```

Creating the training and test set; Preparing the data

A quick analysis of the data shows that many columns are near-empty, or even totally empty.

```
hist(colSums(is.na(pm1MainFULL)), main="Histogram of number of NA values per column", b  
reaks = 100, xlab = "Number of NA values per column", ylab = "Number of columns", co  
l="red")
```

Histogram of number of NA values per column



```
sum((colSums(is.na(pm1MainFULL))>19000)==TRUE)
```

```
## [1] 100
```

```
sum((colSums(is.na(pm1MainFULL))==0)==TRUE)
```

```
## [1] 60
```

60 columns have no NA values, while the other 100 rows are empty or almost (over 19000 NA values, out of 19622). With only a few NA values, we could have tried to impute, but that will most likely not work here since there is so much missing data. Let's remove these last columns from our main dataset, and apply the same pre-processing to our Validation set.

```
pmlMain<-pmlMainFULL[,colSums(is.na(pmlMainFULL))==0]
pmlValidation<-pmlValidation[,colSums(is.na(pmlValidation)) == 0]
```

Let's see a summary of our variables, then pair-plot the first variables, which are not sensor variables.

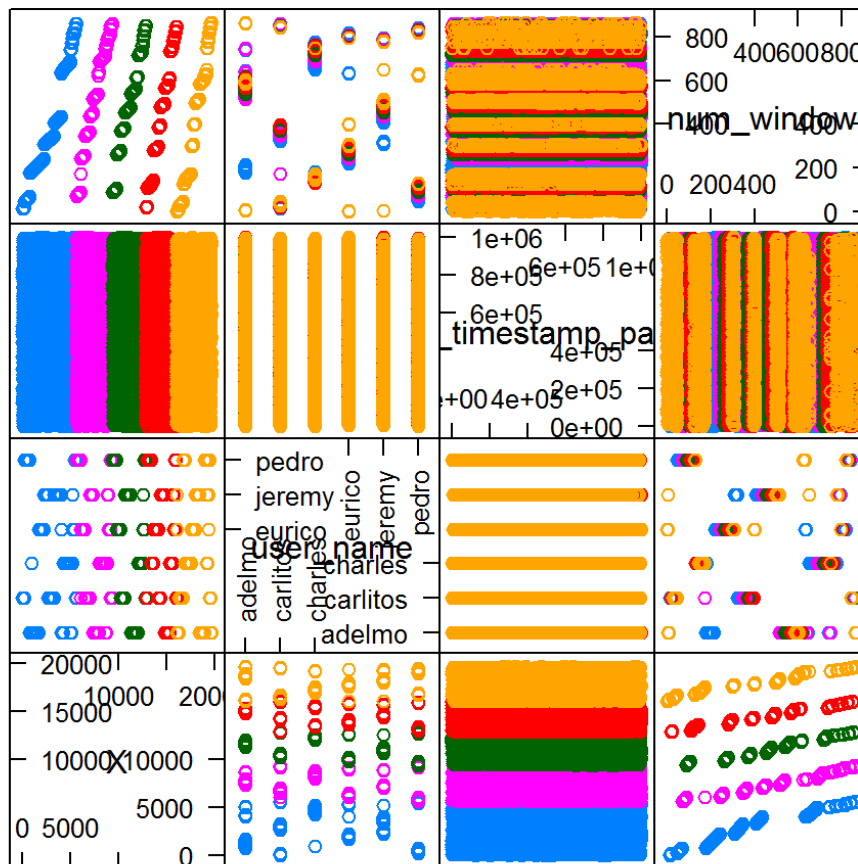
```
summary(pmlMain[,1:10])
```

```
##           X           user_name  raw_timestamp_part_1 raw_timestamp_part_2
## Min.      :    1   adelmo    :3892   Min.      :1.322e+09   Min.      :   294
## 1st Qu.: 4906   carlitos :3112   1st Qu.:1.323e+09   1st Qu.:252912
## Median : 9812   charles  :3536   Median :1.323e+09   Median :496380
## Mean    : 9812   eurico   :3070   Mean    :1.323e+09   Mean    :500656
## 3rd Qu.:14717   jeremy   :3402   3rd Qu.:1.323e+09   3rd Qu.:751891
## Max.    :19622   pedro    :2610   Max.    :1.323e+09   Max.    :998801
##
##           cvtd_timestamp  new_window  num_window  roll_belt
## 28/11/2011 14:14: 1498   no :19216   Min.      : 1.0   Min.      : -28.90
## 05/12/2011 11:24: 1497   yes:  406   1st Qu.:222.0   1st Qu.:   1.10
## 30/11/2011 17:11: 1440                                     Median :424.0   Median :113.00
## 05/12/2011 11:25: 1425                                     Mean    :430.6   Mean    :  64.41
## 02/12/2011 14:57: 1380                                     3rd Qu.:644.0   3rd Qu.:123.00
## 02/12/2011 13:34: 1375                                     Max.     :864.0   Max.     :162.00
## (Other)           :11007
##   pitch_belt      yaw_belt
## Min.     : -55.8000   Min.     : -180.00
## 1st Qu.:   1.7600   1st Qu.:  -88.30
## Median :   5.2800   Median :  -13.00
## Mean    :   0.3053   Mean    :  -11.21
## 3rd Qu.:  14.9000   3rd Qu.:   12.90
## Max.     :  60.3000   Max.     :  179.00
##
```

```
library(caret); library(kernlab)
```

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

```
featurePlot(x=pmlMain[,c("X","user_name","raw_timestamp_part_2","num_window")],y=pmlMain$classe,plot="pairs")
```



Scatter Plot Matrix

These look like strange - X in has almost perfect correlation with classe, our outcome, and some of the time stamp variables appear to be the same. X is most likely an index after ordering by classe, or and/or experiments were ordered by time and by “classe”.

Let's not use those variables in dataset and focus our training solely on sensor data.

```
pmlMain<-pmlMain[,-c(1:7)]
pmlValidation<-pmlValidation[,-c(1:7)]
```

Preparing the dataset

First, let's create a training set and a testing set, using 75% of samples for training.

```
set.seed(7081971)
inTrain<-createDataPartition(y=pmlMain$classe, p=0.75, list = FALSE)
trainingPML<-pmlMain[inTrain,]
testingPML<-pmlMain[-inTrain,]
```

Building Machine Learning Models.

In this part we'll try a few categorization models that can produce a factor output, plot the importance of the variables, then build a few models and estimate the Out-of-Sample Error rate.

Let's first set our code to use multiple cores - if not, training certain models can be extremely long.

```
## Loading required package: foreach
## Loading required package: iterators
## Loading required package: parallel
```

Gradient Boosted Model training

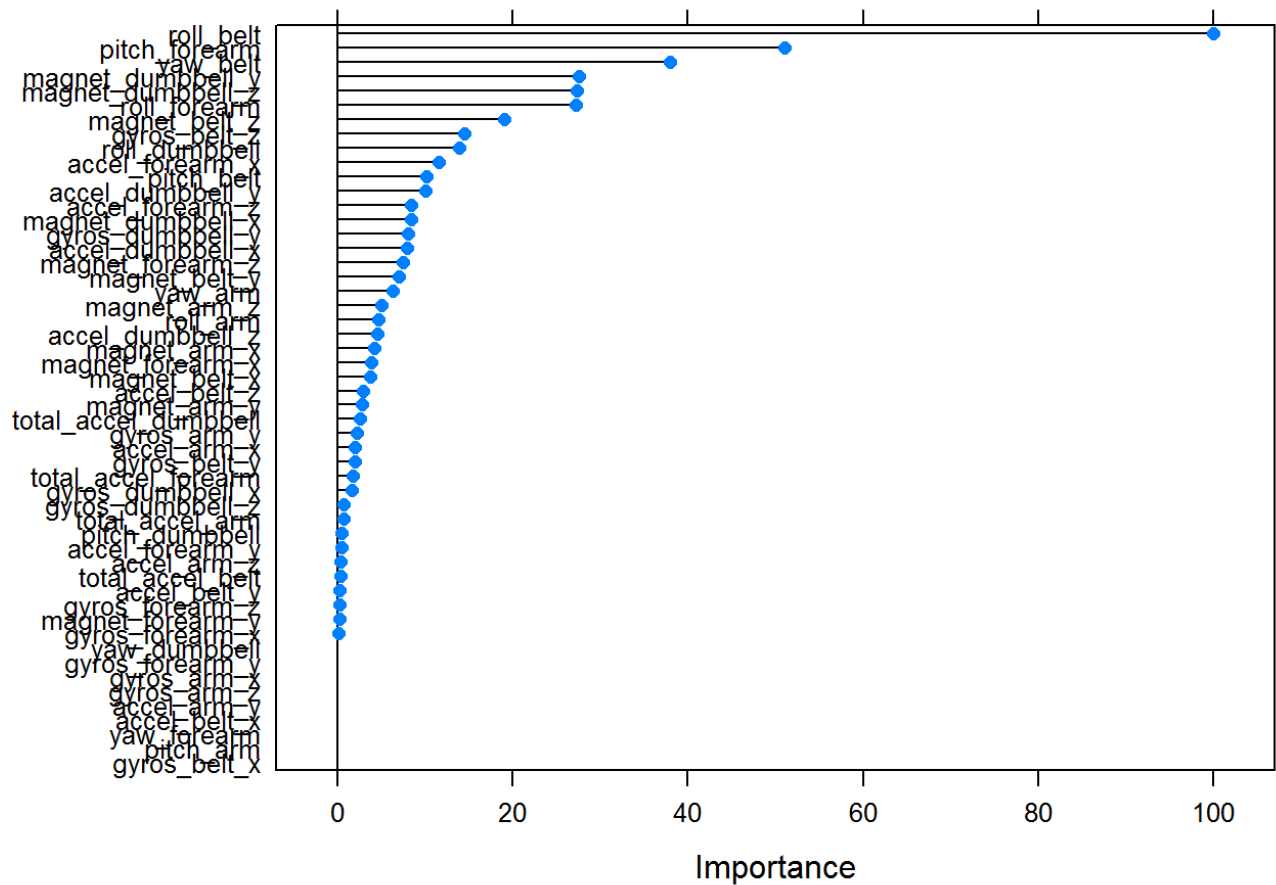
Let's start with a boosted tree model.

```
## Loading required package: gbm
## Loading required package: survival
##
## Attaching package: 'survival'
##
## The following object is masked from 'package:caret':
##
##   cluster
##
## Loading required package: splines
## Loaded gbm 2.1.1
## Loading required package: plyr
```

## Iter	TrainDeviance	ValidDeviance	StepSize	Improve
## 1	1.6094	nan	0.1000	0.2305
## 2	1.4607	nan	0.1000	0.1614
## 3	1.3597	nan	0.1000	0.1285
## 4	1.2795	nan	0.1000	0.0973
## 5	1.2170	nan	0.1000	0.0854
## 6	1.1626	nan	0.1000	0.0852
## 7	1.1099	nan	0.1000	0.0742
## 8	1.0639	nan	0.1000	0.0524
## 9	1.0300	nan	0.1000	0.0540
## 10	0.9946	nan	0.1000	0.0554
## 20	0.7626	nan	0.1000	0.0216
## 40	0.5351	nan	0.1000	0.0120
## 60	0.4111	nan	0.1000	0.0094
## 80	0.3274	nan	0.1000	0.0047
## 100	0.2696	nan	0.1000	0.0029
## 120	0.2275	nan	0.1000	0.0038
## 140	0.1929	nan	0.1000	0.0010
## 150	0.1791	nan	0.1000	0.0016

Let's look at a summary of our model, including variable importance.

```
#Plot variable importance
plot(varImp(modFitGBM))
```



GBM Prediction Accuracy and Error rate

We can now predict on the testing dataset and estimate accuracy of the model.

```
predictions<-predict(modFitGBM, newdata=testingPML)
#Estimate Prediction accuracy
confusionMatrix(predictions, testingPML$classe)
```

```
## Confusion Matrix and Statistics
##
##           Reference
## Prediction   A    B    C    D    E
##           A 1371   30    0    0    2
##           B   18  902   40    2   11
##           C    6   16  799   30    6
##           D    0    1   13  766   10
##           E    0    0    3    6  872
##
## Overall Statistics
##
##           Accuracy : 0.9604
##           95% CI : (0.9546, 0.9657)
##           No Information Rate : 0.2845
##           P-Value [Acc > NIR] : < 2.2e-16
##
##           Kappa : 0.9499
##           McNemar's Test P-Value : NA
##
## Statistics by Class:
##
##           Class: A Class: B Class: C Class: D Class: E
## Sensitivity           0.9828   0.9505   0.9345   0.9527   0.9678
## Specificity           0.9909   0.9820   0.9857   0.9941   0.9978
## Pos Pred Value        0.9772   0.9270   0.9323   0.9696   0.9898
## Neg Pred Value        0.9931   0.9880   0.9862   0.9908   0.9928
## Prevalence            0.2845   0.1935   0.1743   0.1639   0.1837
## Detection Rate        0.2796   0.1839   0.1629   0.1562   0.1778
## Detection Prevalence  0.2861   0.1984   0.1748   0.1611   0.1796
## Balanced Accuracy      0.9868   0.9663   0.9601   0.9734   0.9828
```

Results look good. The **out-of-sample error rate is around 4%**, with the confidence intervals above.

Evaluation on the test dataset

```
predictions<-predict(modFitGBM, newdata=pmlValidation)
```

```
## Loading required package: gbm
## Loading required package: survival
## Loading required package: lattice
## Loading required package: splines
## Loading required package: parallel
## Loaded gbm 2.1.1
## Loading required package: plyr
```

```
#Estimate Prediction accuracy
```

Let's create one file per answer using the following function:

```
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=FALSE,row.names=FALSE,col.names=FALSE)  
  }  
}
```

Let's write our predictions in the results subdirectory:

```
dir.create("GBMresults")  
setwd("results")  
pml_write_files(predictions)
```

Further investigations

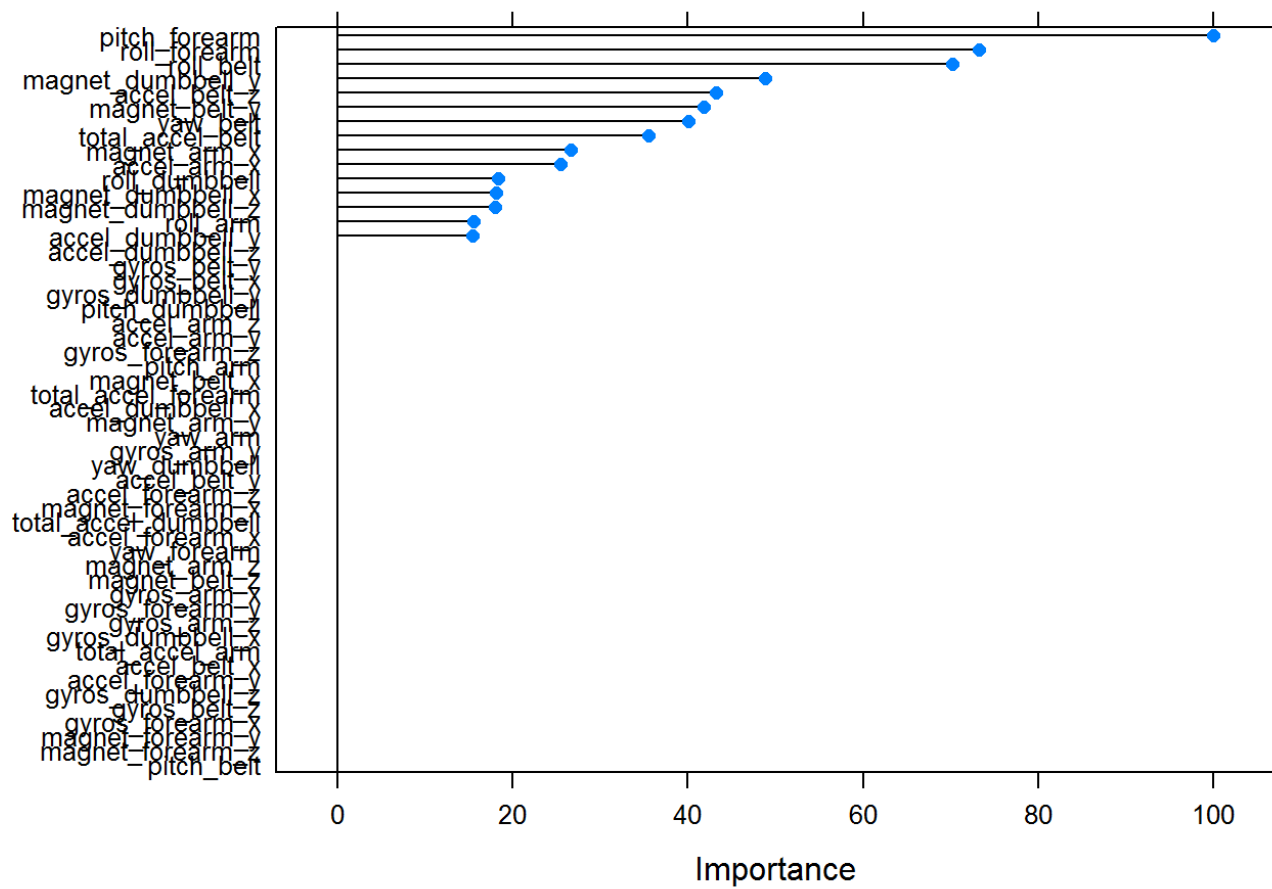
By lack of time, I was not able to re-train many other categorizers without the non-sensor data. I had trained most of the categorizers below on the full dataset (excluding X, but not the various timestamps), until realizing I should exclude them from the predictors as they seem to produce results too good to be true.

RPART

```
modFitRPART<-train(classe~.,method="rpart",data=trainingPML)
```

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

```
plot(varImp(modFitRPART))
```

```
predictions<-predict(modFitRPART, newdata=testingPML)
confusionMatrix(predictions, testingPML$classe)
```

```
## Confusion Matrix and Statistics
##
##           Reference
## Prediction   A    B    C    D    E
##           A 1268  384  389  351  137
##           B   18  331   35  156  120
##           C  104  234  431  297  238
##           D    0    0    0    0    0
##           E    5    0    0    0  406
##
## Overall Statistics
##
##           Accuracy : 0.4967
##           95% CI : (0.4826, 0.5108)
##           No Information Rate : 0.2845
##           P-Value [Acc > NIR] : < 2.2e-16
##
##           Kappa : 0.3426
##           McNemar's Test P-Value : NA
##
## Statistics by Class:
##
##           Class: A Class: B Class: C Class: D Class: E
## Sensitivity           0.9090   0.3488   0.50409   0.0000   0.45061
## Specificity           0.6406   0.9168   0.78439   1.0000   0.99875
## Pos Pred Value        0.5014   0.5015   0.33052      NaN   0.98783
## Neg Pred Value        0.9465   0.8544   0.88222   0.8361   0.88983
## Prevalence            0.2845   0.1935   0.17435   0.1639   0.18373
## Detection Rate        0.2586   0.0675   0.08789   0.0000   0.08279
## Detection Prevalence  0.5157   0.1346   0.26591   0.0000   0.08381
## Balanced Accuracy      0.7748   0.6328   0.64424   0.5000   0.72468
```

The important variables are not the same, which is suspicious. This model is not working well, as the Accuracy shows.

GBM PCA

```
#Estimate Error rates
# modFitGBMPCA<-train(classe~.,method="gbm",preProcess="pca",data=trainingPML)
# predictions<-predict(modFitGBMPCA, newdata=testingPML)
# confusionMatrix(predictions, testingPML$classe)
```

Random Forest

```
#Estimate Error rates
# modFitRF<-train(classe~.,method="rf",data=trainingPML)
# plot(varImp(modFitRF))
# predictions<-predict(modFitRF, newdata=testingPML)
# confusionMatrix(predictions, testingPML$classe)
```

Treebag

```
# modFitTreeBag<-train(classe~.,method="treebag",data=trainingPML)
# plot(varImp(modFitTreeBag))
# predictions<-predict(modFitTreeBag, newdata=testingPML)
# confusionMatrix(predictions, testingPML$classe)
```

BagFDA

```
#Estimate Error rates
# modFitBagFDA<-train(classe~.,method="bagFDA",data=trainingPML)
# plot(varImp(modFitBagFDA))
# predictions<-predict(modFitBagFDA, newdata=testingPML)
# confusionMatrix(predictions, testingPML$classe)
```