

Muhammad Atif

Practical machine learning assignment

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. The aim of this project is to predict the manner in which participants perform a barbell lift. The data comes from <http://groupware.les.inf.puc-rio.br/har> wherein 6 participants were asked to perform the same set of exercises correctly and incorrectly with accelerometers placed on the belt, forearm, arm, and dumbbell.

For the purpose of this project, the following steps would be followed:

1. Data Preprocessing
2. Exploratory Analysis
3. Prediction Model Selection
4. Predicting Test Set Output

Data Preprocessing

First, we load the training and testing set from the online sources and then split the training set further into training and test sets.

```
library(caret)
setwd("~/Projects/R/Coursera-Practical-Machine-Learning-Assignment-1/")
trainURL <- "http://d396qusza40orc.cloudfront.net/predmachlearn/pml-training.csv"
testURL <- "http://d396qusza40orc.cloudfront.net/predmachlearn/pml-testing.csv"

training <- read.csv(url(trainURL))
testing <- read.csv(url(testURL))

label <- createDataPartition(training$classe, p = 0.7, list = FALSE)
train <- training[label, ]
test <- training[-label, ]
```

From among 160 variables present in the dataset, some variables have nearly zero variance whereas some contain a lot of NA terms which need to be excluded from the dataset. Moreover, other 5 variables used for identification can also be removed.

```
NZV <- nearZeroVar(train)
train <- train[ , -NZV]
test <- test[ , -NZV]

label <- apply(train, 2, function(x) mean(is.na(x))) > 0.95
train <- train[, -which(label, label == FALSE)]
test <- test[, -which(label, label == FALSE)]

train <- train[ , -(1:5)]
test <- test[ , -(1:5)]
```

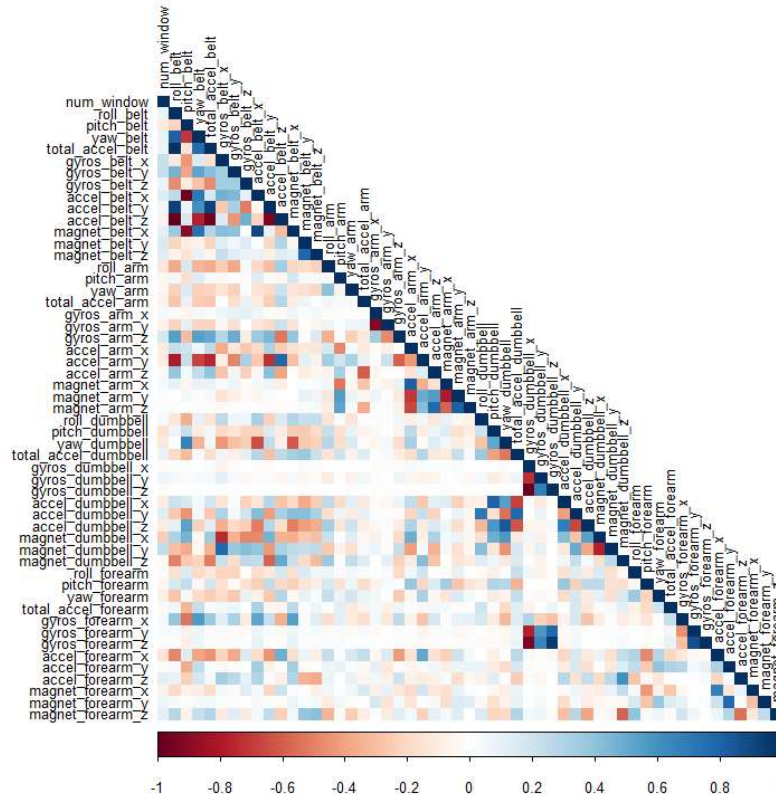
As a result of the preprocessing steps, we were able to reduce 160 variables to 54.

Exploratory Analysis

Now that we have cleaned the dataset off absolutely useless variables, we shall look at the dependence of these variables on each other through a correlation plot.

```
library(corrplot)
corrMat <- cor(train[, -54])
```

```
corrplot(corrMat, method = "color", type = "lower", tl.cex = 0.8, tl.col = rgb(0,0,0
```



In the plot above, darker gradient correspond to having high correlation. A Principal Component Analysis can be run to further reduce the correlated variables but we aren't doing that due to the number of correlations being quite few.

Prediction Model Selection

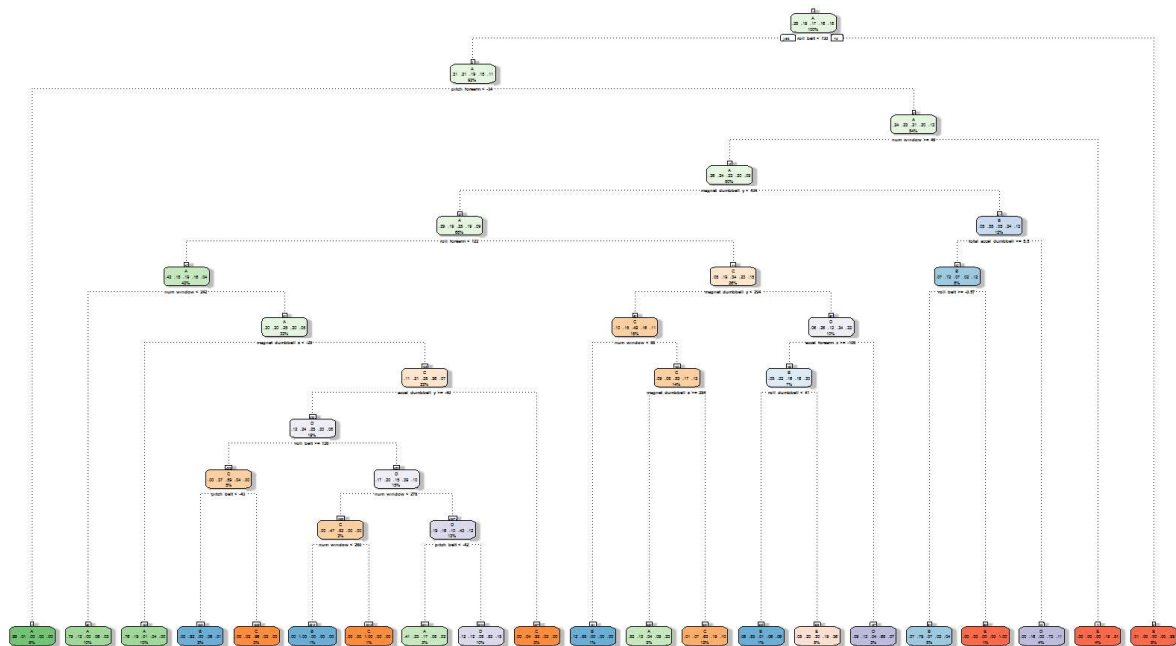
We will use 3 methods to model the training set and thereby choose the one having the best accuracy to predict the outcome variable in the testing set. The methods are Decision Tree, Random Forest and Generalized Boosted Model.

A confusion matrix plotted at the end of each model will help visualize the analysis better.

Decision Tree

```
library(rpart)
library(rpart.plot)
library(rattle)
set.seed(13908)
```

```
modelDT <- rpart(classe ~ ., data = train, method = "class")
fancyRpartPlot(modelDT)
```



```
predictDT <- predict(modelDT, test, type = "class")
confMatDT <- confusionMatrix(predictDT, test$classe)
confMatDT
```

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

```
##
```

```
##           Reference
```

```
## Prediction    A    B    C    D    E
##           A 1502  236   53   83   47
##           B   49  627   34   32   24
##           C    8   67  829  139   69
##           D   95  152   51  609  116
##           E   20   57   59  101  826
```

```
##
```

```
## Overall Statistics
```

```
##
```

```
##           Accuracy : 0.7465
##           95% CI : (0.7352, 0.7575)
```

```
##           No Information Rate : 0.2845
##           P-Value [Acc > NIR] : < 2.2e-16
```

```
##
```

```
##           Kappa : 0.678
```

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

```
##
## Statistics by Class:
##
##           Class: A Class: B Class: C Class: D Class: E
## Sensitivity      0.8973  0.5505  0.8080  0.6317  0.7634
## Specificity      0.9005  0.9707  0.9418  0.9159  0.9507
## Pos Pred Value   0.7819  0.8185  0.7455  0.5953  0.7770
## Neg Pred Value    0.9566  0.9000  0.9587  0.9270  0.9469
## Prevalence       0.2845  0.1935  0.1743  0.1638  0.1839
## Detection Rate    0.2552  0.1065  0.1409  0.1035  0.1404
## Detection Prevalence 0.3264  0.1302  0.1890  0.1738  0.1806
## Balanced Accuracy 0.8989  0.7606  0.8749  0.7738  0.8570
```

Random Forest

```
library(caret)
set.seed(13908)
control <- trainControl(method = "cv", number = 3, verboseIter=FALSE)
modelRF <- train(classe ~ ., data = train, method = "rf", trControl = control)
modelRF$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: 0.21%
## Confusion matrix:
##      A    B    C    D    E  class.error
## A 3904     1     0     0     1 0.0005120328
## B   42651     3     0     0 0.0026335591
## C     0     2 2394     0     0 0.0008347245
## D     0     0   13 2238     1 0.0062166963
## E     0     1     0     3 2521 0.0015841584
```

```
predictRF <- predict(modelRF, test)
confMatRF <- confusionMatrix(predictRF, test$classe)
confMatRF
```

```
## Confusion Matrix and Statistics
##
##           Reference
## Prediction   A    B    C    D    E
##           A 1674    2    0    0    0
##           B    0 1136    4    0    0
##           C    0    0 1022    7    0
##           D    0    1    0 957    1
##           E    0    0    0    0 1081
##
## Overall Statistics
##
##           Accuracy : 0.9975
##           95% CI : (0.9958, 0.9986)
##           No Information Rate : 0.2845
##           P-Value [Acc > NIR] : < 2.2e-16
##
##           Kappa : 0.9968
##           McNemar's Test P-Value : NA
##
## Statistics by Class:
##
##           Class: A Class: B Class: C Class: D Class: E
## Sensitivity          1.0000   0.9974   0.9961   0.9927   0.9991
## Specificity          0.9995   0.9992   0.9986   0.9996   1.0000
## Pos Pred Value       0.9988   0.9965   0.9932   0.9979   1.0000
## Neg Pred Value       1.0000   0.9994   0.9992   0.9986   0.9998
## Prevalence           0.2845   0.1935   0.1743   0.1638   0.1839
## Detection Rate       0.2845   0.1930   0.1737   0.1626   0.1837
## Detection Prevalence 0.2848   0.1937   0.1749   0.1630   0.1837
## Balanced Accuracy     0.9998   0.9983   0.9973   0.9962   0.9995
```

Generalized Boosted Model

```
library(caret)
set.seed(13908)
control <- trainControl(method = "repeatedcv", number = 5, repeats = 1, verboseIter
modelGBM <- train(classe ~ ., data = train, trControl = control, method = "gbm", ver
modelGBM$finalModel
```

```
## A gradient boosted model with multinomial loss function.
## 150 iterations were performed.
```

```
## There were 53 predictors of which 41 had non-zero influence.
```

```
predictGBM <- predict(modelGBM, test)
confMatGBM <- confusionMatrix(predictGBM, test$classe)
confMatGBM
```

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

```
##
```

```
##           Reference
```

```
## Prediction    A    B    C    D    E
##           A 1670    10    0    0    1
##           B   3 1117    11    5    2
##           C   0   11 1009   10    3
##           D   1    1    3  949    2
##           E   0    0    3    0 1074
```

```
##
```

```
## Overall Statistics
```

```
##
```

```
##           Accuracy : 0.9888
##           95% CI : (0.9858, 0.9913)
##           No Information Rate : 0.2845
##           P-Value [Acc > NIR] : < 2.2e-16
```

```
##
```

```
##           Kappa : 0.9858
```

```
## McNemar's Test P-Value : NA
```

```
##
```

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

```
##
```

```
##           Class: A Class: B Class: C Class: D Class: E
## Sensitivity      0.9976  0.9807  0.9834  0.9844  0.9926
## Specificity      0.9974  0.9956  0.9951  0.9986  0.9994
## Pos Pred Value   0.9935  0.9815  0.9768  0.9927  0.9972
## Neg Pred Value    0.9990  0.9954  0.9965  0.9970  0.9983
## Prevalence       0.2845  0.1935  0.1743  0.1638  0.1839
## Detection Rate   0.2838  0.1898  0.1715  0.1613  0.1825
## Detection Prevalence 0.2856  0.1934  0.1755  0.1624  0.1830
## Balanced Accuracy 0.9975  0.9881  0.9892  0.9915  0.9960
```

As Random Forest offers the maximum accuracy of 99.75%, we will go with Random Forest Model to predict our test data class variable.

Predicting Test Set Output

```
predictRF <- predict(modelRF, testing)
predictRF
```

```
## [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
```