Practical Machine Learning Project by John Hopkins University and Coursera

Niño Lito Jake Briones

August 18, 2020

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 dumbell 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 SOURCE:

The training data for this project are available here:

[https://d396qusza40orc.cloudfront.net/predmachlearn/pml-training.csv]

The test data are available here:

[https://d396qusza40orc.cloudfront.net/predmachlearn/pml-testing.csv]

The data for this project come from this source: [http://groupware.les.inf.puc-rio.br/har]. If you use the document you create for this class for any purpose please cite them as they have been very generous in allowing their data to be used for this kind of assignment.

1. LOAD THE PACKAGES AND DOWNLOAD THE DATASET:

The dataset will be downloaded from the internet with two separate dataframes, - training and testing.

```
library(caret)

## Loading required package: lattice

## Loading required package: ggplot2

library(rpart)
library(rpart.plot)
library(randomForest)
```

- ## randomForest 4.6-14
- ## Type rfNews() to see new features/changes/bug fixes.

```
##
## Attaching package: 'randomForest'
## The following object is masked from 'package:ggplot2':
##
       margin
library(corrplot)
## corrplot 0.84 loaded
Download the Dataset
trainUrl <-"https://d396qusza40orc.cloudfront.net/predmachlearn/pml-training.csv"
testUrl <- "https://d396qusza40orc.cloudfront.net/predmachlearn/pml-testing.csv"
trainFile <- "./data/pml-training.csv"</pre>
testFile <- "./data/pml-testing.csv"</pre>
if (!file.exists("./data")) {
  dir.create("./data")
}
if (!file.exists(trainFile)) {
  download.file(trainUrl, destfile=trainFile, method="curl")
if (!file.exists(testFile)) {
  download.file(testUrl, destfile=testFile, method="curl")
}
```

2. READ THE DATA:

After downloading the data from the data source, we can read the two csv files into two data frames.

```
trainRaw <- read.csv("./data/pml-training.csv")
testRaw <- read.csv("./data/pml-testing.csv")
dim(trainRaw)
## [1] 19622 160
dim(testRaw)</pre>
```

[1] 20 160

The training data set contains 19622 observations and 160 variables, while the testing data set contains 20 observations and 160 variables. The "classe" variable in the training set is the outcome to predict.

3. DATA CLEANING AND PREPARATION:

We take off those data contains more than 95% of the observation to be NA. We filter out those records.

```
sum(complete.cases(trainRaw))
## [1] 406
a. rid of the columns that contain NA missing values.
trainRaw <- trainRaw[, colSums(is.na(trainRaw)) == 0]
testRaw <- testRaw[, colSums(is.na(testRaw)) == 0]</pre>
```

b. rid of some columns that do not contribute much to the accelerometer measurements.

```
classe <- trainRaw$classe
trainRemove <- grepl("^X|timestamp|window", names(trainRaw))
trainRaw <- trainRaw[, !trainRemove]
trainCleaned <- trainRaw[, sapply(trainRaw, is.numeric)]
trainCleaned$classe <- classe
testRemove <- grepl("^X|timestamp|window", names(testRaw))
testRaw <- testRaw[, !testRemove]
testCleaned <- testRaw[, sapply(testRaw, is.numeric)]</pre>
```

As a result, the cleaned training data set contains 19622 observations and 53 variables, while the testing data set contains 20 observations and 53 variables. The "classe" variable is still in the cleaned training set.

3. DATA PARTITION:

I will split the cleaned training set into a pure training data set (70%) and a validation data set (30%). I will use the validation data set to conduct cross validation in future steps.

```
set.seed(22519) # For reproducibile purpose
inTrain <- createDataPartition(trainCleaned$classe, p=0.70, list=F)
trainData <- trainCleaned[inTrain, ]
testData <- trainCleaned[-inTrain, ]</pre>
```

5. RANDOM FOREST; MACHINE LEARNING ALGORITHM:

I will use the random forest since it automatically selects important variables and is robust to correlated covariates & outliers in general. We will use **5-fold cross validation** when applying the algorithm.

```
controlRf <- trainControl(method="cv", 5)</pre>
modelRf <- train(classe ~ ., data=trainData, method="rf", trControl=controlRf, ntree=250)</pre>
modelRf
## Random Forest
##
## 13737 samples
##
      52 predictor
       5 classes: 'A', 'B', 'C', 'D', 'E'
##
##
## No pre-processing
## Resampling: Cross-Validated (5 fold)
## Summary of sample sizes: 10988, 10989, 10989, 10991, 10991
## Resampling results across tuning parameters:
##
##
     mtry Accuracy
                      Kappa
           0.9912654 0.9889499
##
     2
##
     27
           0.9916291 0.9894104
##
           0.9842766 0.9801110
##
## Accuracy was used to select the optimal model using the largest value.
## The final value used for the model was mtry = 27.
```

6. PREDICTION:

Now, we apply the model to the original testing data set downloaded from the data source. We remove the problem_id column first.

```
result <- predict(modelRf, testCleaned[, -length(names(testCleaned))])
result</pre>
```

```
## [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
```

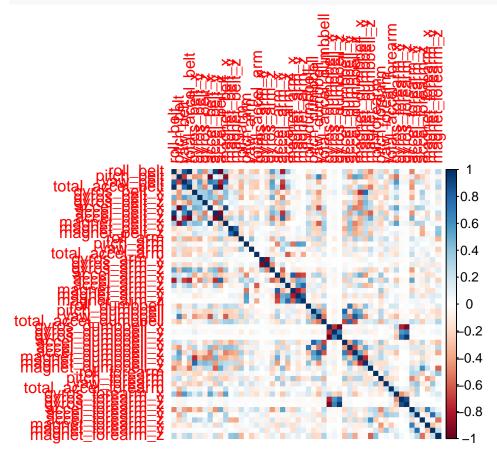
7. CONCLUSION:

Based on the result, in terms of accuracy random forest model has a better output than decision tree model. I get a 99.25% result from random forest and only 50% from decision tree model.

8. FIGURES:

1. Correlation Matrix Visualization

```
corrPlot <- cor(trainData[, -length(names(trainData))])
corrplot(corrPlot, method="color")</pre>
```



2. Decision Tree Visualization

```
treeModel <- rpart(classe ~ ., data=trainData, method="class")
prp(treeModel) # fast plot</pre>
```

