Including Plot

Prediction Assignment Write up Mariana Martins

Overview

This analysis meant to be the basis for the course quiz and a prediction assignment write up. The main goal of the project is to predict the manner in which 6 participants performed some exercise as described below. This is the "class" variable in the training set. The machine learning algorithm described here is applied to the 20 test cases available in the test data and the predictions are submitted in appropriate format to the Course Project Prediction Quiz for automated grading.

Data Loading and Exploratory Analysis

Data Source The training data for this project are available here:

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

The test data are available here:

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

Environment Setup

library(knitr) library(caret) library(rpart) library(rpart.plot) library(rattle) library(randomForest) library(corrplot) set.seed(301)

```
Data Loading and Cleaning
The next step is loading the dataset from the URL provided above. The training dataset is then partitioned in 2 to create a Training set (70% of the
```

for the quiz results generation.

TrainUrl <- "https://d396qusza40orc.cloudfront.net/predmachlearn/pml-training.csv" TestUrl <- "https://d396qusza40orc.cloudfront.net/predmachlearn/pml-testing.csv"

data) for the modeling process and a Test set (with the remaining 30%) for the validations. The testing dataset is not changed and will only be used

```
TrainFile<-"pml-traininig.csv"
 TestFile<-"pml-testing.csv"
Download the datasets
 if(!file.exists(TrainFile)) {
   download.file(TrainUrl, destfile = TrainFile)
 training <- read.csv(TrainFile)</pre>
```

NZV <- nearZeroVar(TrainSet)

Remove variables that are mostly NA

TrainSet <- TrainSet[, -(1:5)]</pre> TestSet <- TestSet[, -(1:5)]</pre>

[1] 5885 59

dim(TrainSet)

dim(TrainSet)

```
if(!file.exists(TestFile)) {
   download.file(TestUrl,destfile = TestFile)
   }
 testing <- read.csv(TestFile)</pre>
Create a partition using caret with the training dataset on 70,30 ratio
 inTrain <- createDataPartition(training$classe, p=0.7, list=FALSE)</pre>
 TrainSet <- training[inTrain, ]</pre>
```

```
TestSet <- training[-inTrain, ]</pre>
dim(TrainSet)
```

```
## [1] 13737 160
```

dim(TestSet)

```
## [1] 5885 160
Both created datasets have 160 variables. Let's clean NA, The Near Zero variance (NZV) variables and the ID variables as well.
Remove variables with Nearly Zero Variance
```

TrainSet <- TrainSet[, -NZV]</pre> TestSet <- TestSet[, -NZV]</pre> dim(TestSet) ## [1] 5885 105

dim(TrainSet) ## [1] **13737 105**

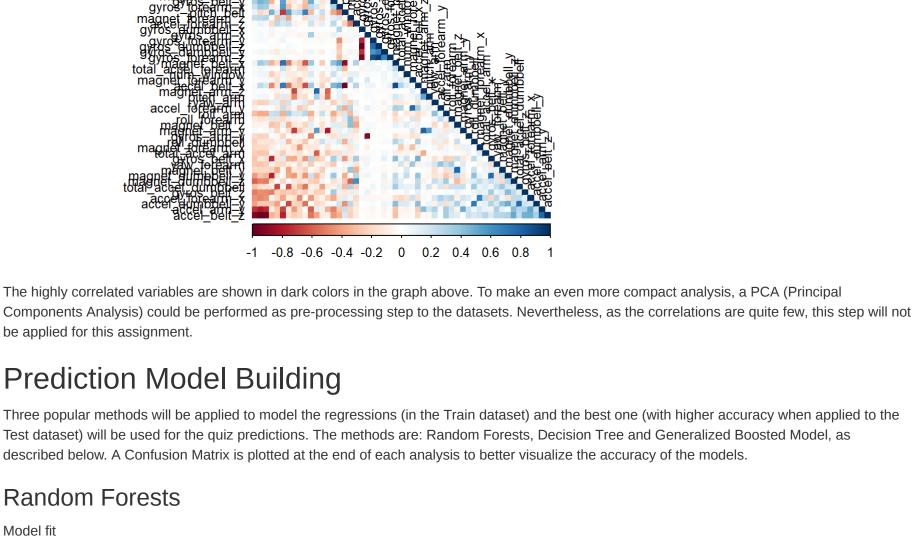
AllNA <- sapply(TrainSet, function(x) mean(is.na(x))) > 0.95 TrainSet <- TrainSet[, AllNA==FALSE]</pre> TestSet <- TestSet[, AllNA==FALSE]</pre> dim(TestSet)

[1] 13737 59 Remove identification only variables (columns 1 to 5)

[1] 13737 54 After cleaning, we can see that the number of variables for the analysis are now only 53. **Correlation Analysis**

corMatrix <- cor(TrainSet[, -54])</pre> corrplot(corMatrix, order = "FPC", method = "color", type = "lower", tl.cex = 0.8, tl.col = rgb(0, 0, 0))

A correlation among variables is analysed before proceeding to the modeling procedures.



controlRF <- trainControl(method="cv", number=3, verboseIter=FALSE)</pre>

0 0.0005120328

0 0.0037622272

0 0.0016694491

1 0.0039964476

Kappa: 0.9985

Mcnemar's Test P-Value : NA

Statistics by Class:

modFitRandForest <- train(classe ~ ., data=TrainSet, method="rf", trContro =controlRF)</pre> modFitRandForest\$finalModel

randomForest(x = x, y = y, mtry = min(param\$mtry, ncol(x)), trContro = ...1)Type of random forest: classification Number of trees: 500

8 2243

Model fit

A 3904

C

D

No. of variables tried at each split: 27 OOB estimate of error rate: 0.23% ## Confusion matrix: C D E class.error

```
5 2519 0.0023762376
Prediction on Test dataset
 predictRandForest <- predict(modFitRandForest, newdata=TestSet)</pre>
 confMatRandForest <- confusionMatrix(predictRandForest, as.factor(TestSet$classe))</pre>
 confMatRandForest
 ## Confusion Matrix and Statistics
               Reference
 ## Prediction
                  Α
             A 1674
                        0
                   0 1138
                        1 1026
                                   2
 ##
                                   0 1078
```

Overall Statistics ##

##

##

##

Accuracy : 0.9988 ## 95% CI: (0.9976, 0.9995) No Information Rate: 0.2845 ## P-Value [Acc > NIR] : < 2.2e-16

```
##
                       Class: A Class: B Class: C Class: D Class: E
                      1.0000 0.9991 1.0000 0.9979
## Sensitivity
                                                          0.9963
 ## Specificity
                        1.0000 1.0000
                                                 0.9992
                                        0.9994
                                                          1.0000
 ## Pos Pred Value
                        1.0000 1.0000 0.9971 0.9959
                                                          1.0000
 ## Neg Pred Value
                        1.0000 0.9998 1.0000
                                                  0.9996
                                                          0.9992
 ## Prevalence
                         0.2845
                                 0.1935
                                         0.1743
                                                  0.1638
                                                          0.1839
                   0.2845
 ## Detection Rate
                                 0.1934
                                         0.1743
                                                 0.1635
                                                          0.1832
 ## Detection Prevalence 0.2845
                                                  0.1641 0.1832
                                 0.1934 0.1749
 ## Balanced Accuracy
                        1.0000
                                 0.9996 0.9997
                                                  0.9986
                                                          0.9982
Plot matrix results
plot(confMatRandForest$table, col = confMatRandForest$byClass, main = paste("Random Forest - Accuracy =", round(c
onfMatRandForest$overall['Accuracy'], 4)))
                       Random Forest - Accuracy = 0.9988
    Reference
```

Prediction

Decision Tree Model fit

set.seed(301)

Prediction on Test dataset

Confusion Matrix and Statistics

Reference

209

A 1497

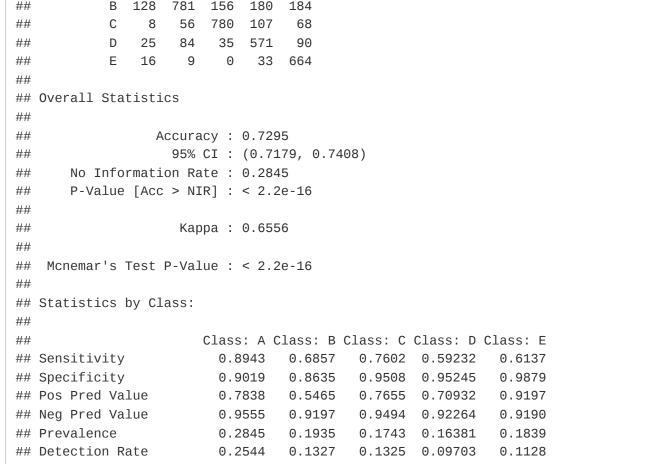
confMatDecTree

Prediction

Plot matrix results

DecTree\$overall['Accuracy'], 4)))

modFitDecTree <- rpart(classe ~ ., data=TrainSet, method="class")</pre> fancyRpartPlot(modFitDecTree)



Rattle 2022-jul-29 14:25:17 mmartins

predictDecTree <- predict(modFitDecTree, newdata=TestSet, type="class")</pre>

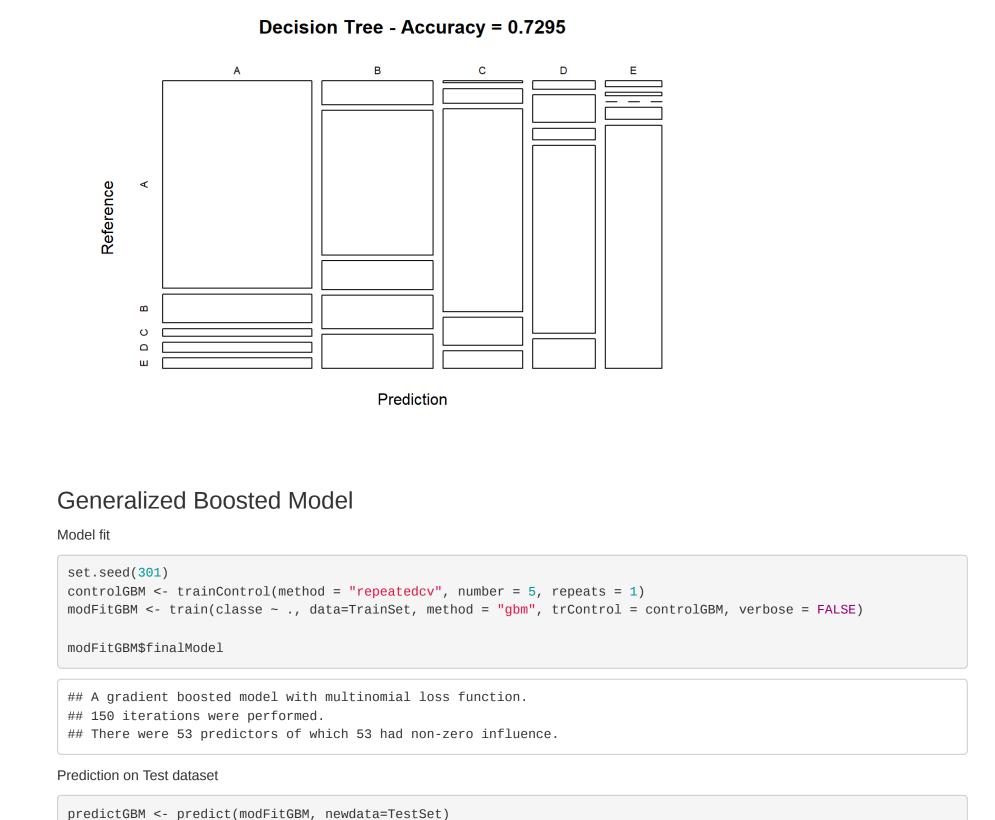
73

Detection Prevalence 0.3246 0.2428 0.1732 0.13679

Balanced Accuracy 0.8981 0.7746 0.8555 0.77239

76

confMatDecTree <- confusionMatrix(predictDecTree, as.factor(TestSet\$classe))</pre>



0.1227

plot(confMatDecTree\$table, col = confMatDecTree\$byClass, main = paste("Decision Tree - Accuracy =", round(confMat

0.8008

Accuracy: 0.9879 ## No Information Rate : 0.2845

Confusion Matrix and Statistics

Reference

A 1671 11

0

1

14 1020 12

Prediction A B

С

Overall Statistics

confMatGBM

confMatGBM <- confusionMatrix(predictGBM, as.factor(TestSet\$classe))</pre>

1

2 945 11

```
95% CI: (0.9848, 0.9906)
       P-Value [Acc > NIR] : < 2.2e-16
 ##
 ##
                   Kappa: 0.9847
 ##
   Mcnemar's Test P-Value : NA
## Statistics by Class:
##
                      Class: A Class: B Class: C Class: D Class: E
## Sensitivity
                     0.9982 0.9772 0.9942 0.9803 0.9843
 ## Specificity
                       0.9967 0.9971 0.9944 0.9972
                                                       0.9996
                     0.9917 0.9876 0.9742 0.9854 0.9981
## Pos Pred Value
 ## Neg Pred Value
                     0.9993 0.9945 0.9988 0.9961 0.9965
## Prevalence
                       0.2845 0.1935 0.1743 0.1638
                                                       0.1839
## Detection Rate 0.2839 0.1891 0.1733 0.1606 0.1810
## Detection Prevalence 0.2863 0.1915 0.1779 0.1630 0.1813
## Balanced Accuracy 0.9974 0.9871 0.9943 0.9887 0.9919
Plot matrix results
plot(confMatGBM$table, col = confMatGBM$byClass, main = paste("GBM - Accuracy =", round(confMatGBM$overall['Accur
acy'], 4)))
                           GBM - Accuracy = 0.9879
```

Applying the selected model The accuracy of the 3 regression modeling methods above are: Random Forest: 0.9986

In that case, the Random Forest model will be applied to predict the 20 quiz results (testing dataset) as shown below.

predictTEST ## [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

Reference

Prediction

• GBM: 0.989

predictTEST <- predict(modFitRandForest, newdata=testing)</pre>

• Decision Tree: 0.7295