

course project practical machine learning

Alex

4/30/2020

This document is the final project for the Coursera “Practical Machine Learning” course. It was produced using RStudio’s Markdown and Knitr.

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. 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 will 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.

The data consists of a Training data and a Test data (to be used to validate the selected model).

The goal of your project is to predict the manner in which they did the exercise. This is the “classe” variable in the training set. You may use any of the other variables to predict with.

Note: The dataset used in this project is a courtesy of “Ugulino, W.; Cardador, D.; Vega, K.; Velloso, E.; Milidui, R.; Fuks, H. Wearable Computing: Accelerometers’ Data Classification of Body Postures and Movements”

Data Loading and Processing

```
library(caret)
```

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

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

```
library(rpart)
```

```
library(RColorBrewer)
```

```
library(gbm)
```

```
## Loaded gbm 2.1.5
```

```
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
library(rpart.plot)
library(rattle)

## Rattle: A free graphical interface for data science with R.
## Version 5.3.0 Copyright (c) 2006-2018 Togaware Pty Ltd.
## Type 'rattle()' to shake, rattle, and roll your data.
##
## Attaching package: 'rattle'
## The following object is masked from 'package:randomForest':
##
##      importance
```

Getting, Cleaning and Exploring the data

```
train_in <- read.csv('./pml-training.csv', header=T)
valid_in <- read.csv('./pml-testing.csv', header=T)
dim(train_in)
```

```
## [1] 19622  160
```

```
dim(valid_in)
```

```
## [1]  20 160
```

As shown below there are 19622 observations and 160 variables in the Training dataset

Cleaning the input data

We will remove the variables that contains missing values. Note along the cleaning process we will display the dimension of the reduced dataset

```
trainData<- train_in[, colSums(is.na(train_in)) == 0]
validData <- valid_in[, colSums(is.na(valid_in)) == 0]
dim(trainData)
```

```
## [1] 19622   93
```

```
dim(validData)
```

```
## [1]  20  60
```

We now remove the first seven variables as they have little impact on the outcome classe

```
trainData <- trainData[, -c(1:7)]
validData <- validData[, -c(1:7)]
dim(trainData)
```

```
## [1] 19622   86
```

```
dim(validData)
```

```
## [1] 20 53
```

Preparing the datasets for prediction

Preparing the data for prediction by splitting the training data into 70% as train data and 30% as test data. This splitting will server also to compute the out-of-sample errors.

The test data renamed: valid_in (validate data) will stay as is and will be used later to test the prodction algorithm on the 20 cases.

```
set.seed(1234)
inTrain <- createDataPartition(trainData$classe, p = 0.7, list = FALSE)
trainData <- trainData[inTrain, ]
testData <- trainData[-inTrain, ]
dim(trainData)
```

```
## [1] 13737    86
```

```
dim(testData)
```

```
## [1] 4123    86
```

Cleaning even further by removing the variables that are near-zero-variance

```
NZV <- nearZeroVar(trainData)
trainData <- trainData[, -NZV]
testData <- testData[, -NZV]
dim(trainData)
```

```
## [1] 13737    53
```

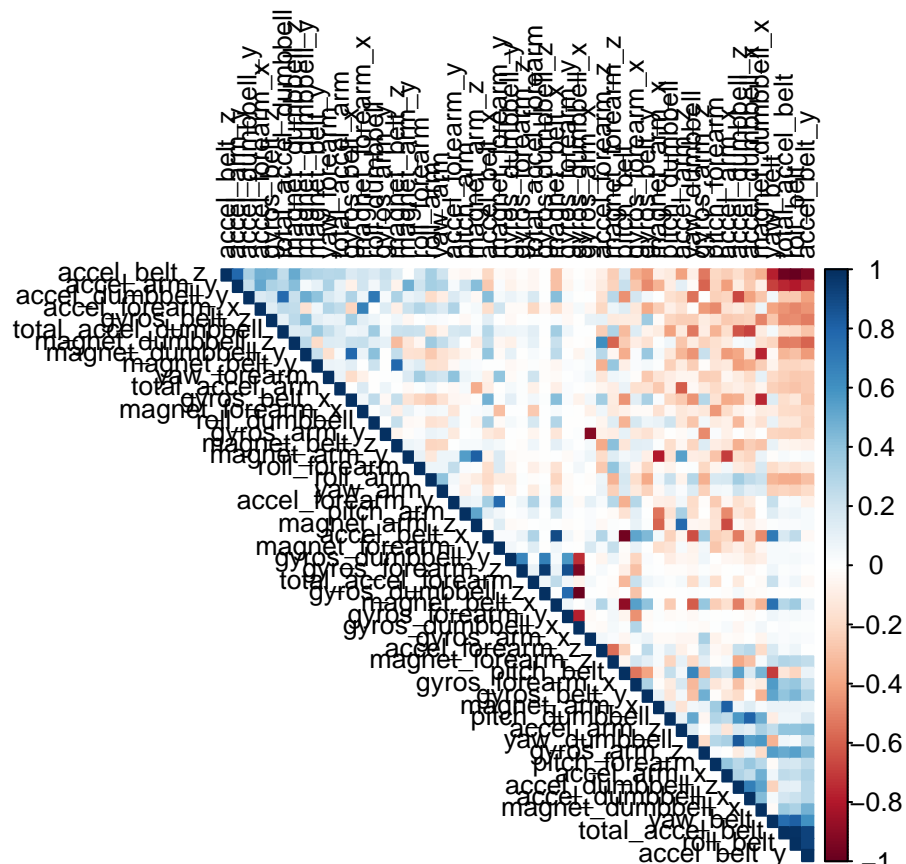
```
dim(testData)
```

```
## [1] 4123    53
```

After this cleaning we are down now to 53 variables

The following correlation plot uses the following parameters (source:CRAN Package ‘corrplot’) “FPC”: the first principal component order. “AOE”: the angular order tl.cex Numeric, for the size of text label (variable names) tl.col The color of text label.

```
cor_mat <- cor(trainData[, -53])
corrplot(cor_mat, order = "FPC", method = "color", type = "upper",
          tl.cex = 0.8, tl.col = rgb(0, 0, 0))
```



related predictors (variables) are those with a dark color intersection.

To obtain the names of the variables we do the following

we use the findCorrelation function to search for highly correlated attributes with a cut off equal to 0.75

```
highlyCorrelated = findCorrelation(cor_mat, cutoff=0.75)
```

We then obtain the names of highly correlated attributes

```
names(trainData)[highlyCorrelated]
```

```
## [1] "accel_belt_z"      "roll_belt"         "accel_belt_y"
## [4] "total_accel_belt"  "accel_dumbbell_z"  "accel_belt_x"
## [7] "pitch_belt"        "magnet_dumbbell_x" "accel_dumbbell_y"
## [10] "magnet_dumbbell_y" "accel_dumbbell_x"  "accel_arm_x"
## [13] "accel_arm_z"        "magnet_arm_y"      "magnet_belt_z"
## [16] "accel_forearm_y"    "gyros_forearm_y"   "gyros_dumbbell_x"
## [19] "gyros_dumbbell_z"   "gyros_arm_x"
```

Model building

For this project we will use two different algorithms, classification trees and random forests, to predict the outcome.

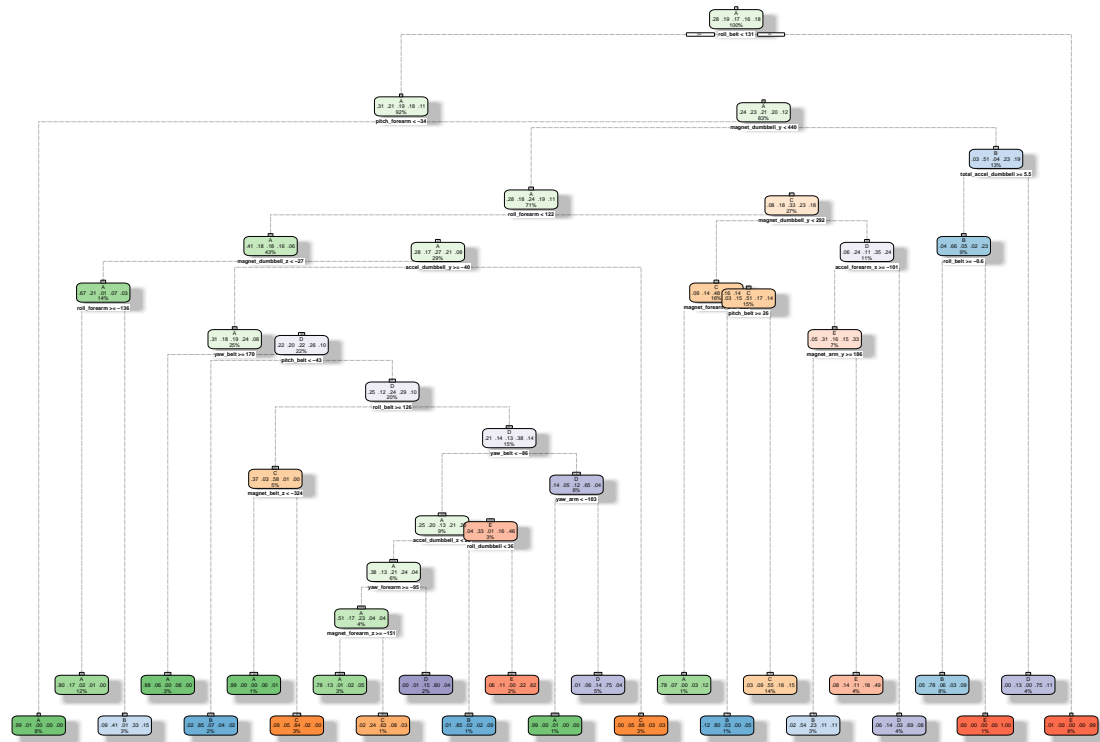
classification trees random forests Generalized Boosted Model

Prediction with classification trees

We first obtain the model, and then we use the `fancyRpartPlot()` function to plot the classification tree as a dendrogram.

```
set.seed(12345)
decisionTreeMod1 <- rpart(classe ~ ., data=trainData, method="class")
fancyRpartPlot(decisionTreeMod1)
```

Warning: labs do not fit even at cex 0.15, there may be some overplotting



Rattle 2020-Apr-30 17:58:11 alexandrosdymiotis

We then val-

idate the model "decisionTreeModel" on the testData to find out how well it performs by looking at the accuracy variable.

```
predictTreeMod1 <- predict(decisionTreeMod1, testData, type = "class")
cmtree <- confusionMatrix(predictTreeMod1, testData$classe)
cmtree
```

Confusion Matrix and Statistics

```
##
##           Reference
## Prediction   A    B    C    D    E
##           A 1067  105    9   24    9
##           B   40  502   59   63   77
##           C   28   90  611  116   86
##           D   11   49   41  423   41
##           E    19   41   18   46  548
```

Overall Statistics

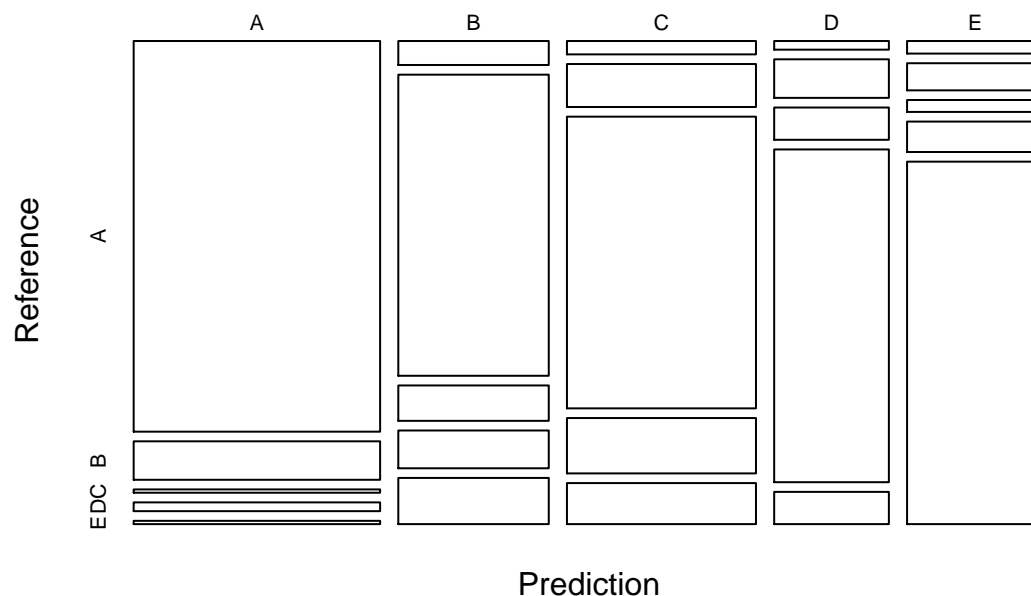
```
##
##           Accuracy : 0.7642
```

```
##                      95% CI : (0.751, 0.7771)
##      No Information Rate : 0.2826
##      P-Value [Acc > NIR] : < 2.2e-16
##
##                      Kappa : 0.7015
##
##      McNemar's Test P-Value : < 2.2e-16
##
## Statistics by Class:
##
##                      Class: A Class: B Class: C Class: D Class: E
## Sensitivity           0.9159   0.6379   0.8279   0.6295   0.7201
## Specificity           0.9503   0.9284   0.9055   0.9589   0.9631
## Pos Pred Value        0.8789   0.6775   0.6563   0.7487   0.8155
## Neg Pred Value        0.9663   0.9157   0.9602   0.9300   0.9383
## Prevalence            0.2826   0.1909   0.1790   0.1630   0.1846
## Detection Rate        0.2588   0.1218   0.1482   0.1026   0.1329
## Detection Prevalence  0.2944   0.1797   0.2258   0.1370   0.1630
## Balanced Accuracy      0.9331   0.7831   0.8667   0.7942   0.8416
```

plot matrix results

```
# plot matrix results
plot(cmtree$table, col = cmtree$byClass,
     main = paste("Decision Tree - Accuracy =", round(cmtree$overall['Accuracy'], 4)))
```

Decision Tree – Accuracy = 0.7642



We see that the accuracy rate of the model is low: 0.6967 and therefore the out-of-sample-error is about 0.3 which is considerable.

Prediction with Random Forest

We first determine the model

```
controlRF <- trainControl(method="cv", number=3, verboseIter=FALSE)
modRF1 <- train(classe ~ ., data=trainData, method="rf", trControl=controlRF)

modRF1$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: 0.7%
## Confusion matrix:
##      A    B    C    D    E class.error
## A 3902     3     0     0     1 0.001024066
## B   19 2634     5     0     0 0.009029345
## C    0   17 2369    10     0 0.011268781
## D    0    1  26 2224     1 0.012433393
## E    0    2    5    6 2512 0.005148515
```

We then validate the model obtained model “modRF1” on the test data to find out how well it performs by looking at the Accuracy variable

```
predictRF1 <- predict(modRF1, newdata=testData)
cmrf <- confusionMatrix(predictRF1, testData$classe)
cmrf
```

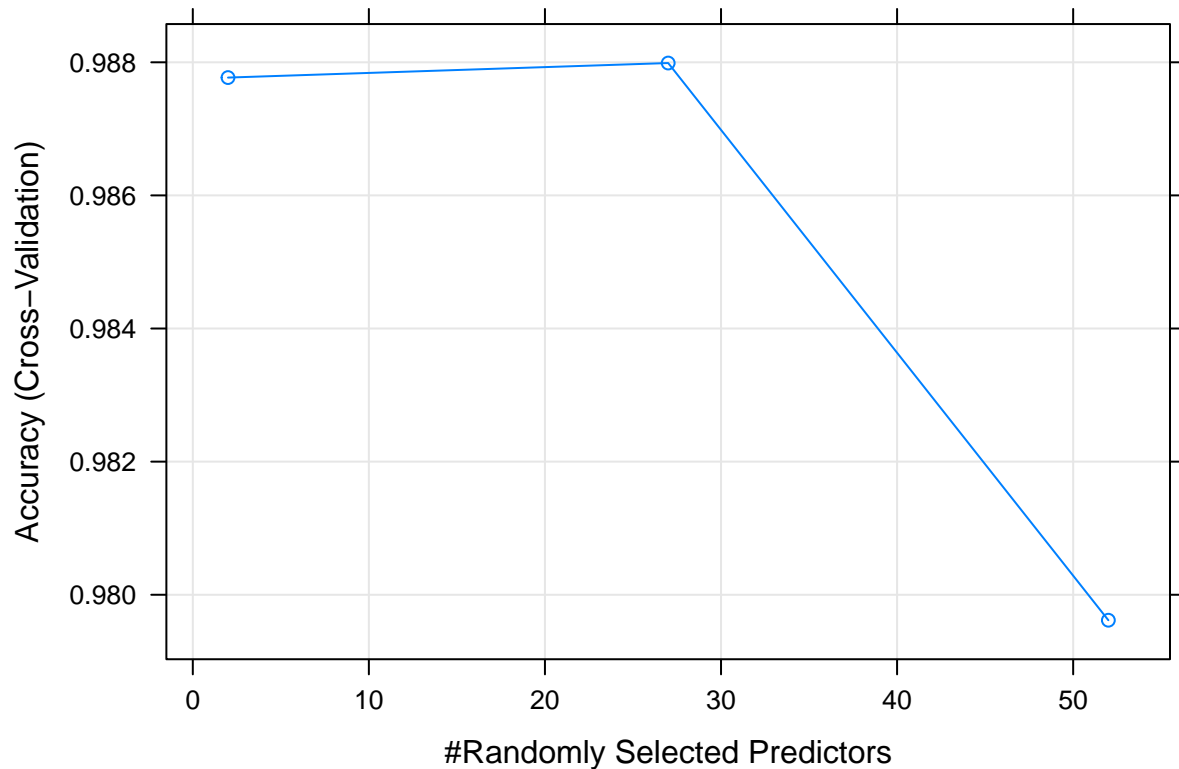
```
## Confusion Matrix and Statistics
##
##           Reference
## Prediction    A    B    C    D    E
##           A 1165     0     0     0     0
##           B    0  787     0     0     0
##           C    0    0  738     0     0
##           D    0    0    0  672     0
##           E    0    0    0    0  761
##
## Overall Statistics
##
##           Accuracy : 1
##           95% CI : (0.9991, 1)
##           No Information Rate : 0.2826
##           P-Value [Acc > NIR] : < 2.2e-16
##
##           Kappa : 1
##
##           McNemar's Test P-Value : NA
##
## Statistics by Class:
##
##           Class: A Class: B Class: C Class: D Class: E
## Sensitivity          1.0000   1.0000   1.000   1.000   1.0000
## Specificity          1.0000   1.0000   1.000   1.000   1.0000
## Pos Pred Value       1.0000   1.0000   1.000   1.000   1.0000
```

## Neg Pred Value	1.0000	1.0000	1.000	1.000	1.0000
## Prevalence	0.2826	0.1909	0.179	0.163	0.1846
## Detection Rate	0.2826	0.1909	0.179	0.163	0.1846
## Detection Prevalence	0.2826	0.1909	0.179	0.163	0.1846
## Balanced Accuracy	1.0000	1.0000	1.000	1.000	1.0000

The accuracy rate using the random forest is very high: Accuracy : 1 and therefore the out-of-sample-error is equal to 0***. But it might be due to overfitting.

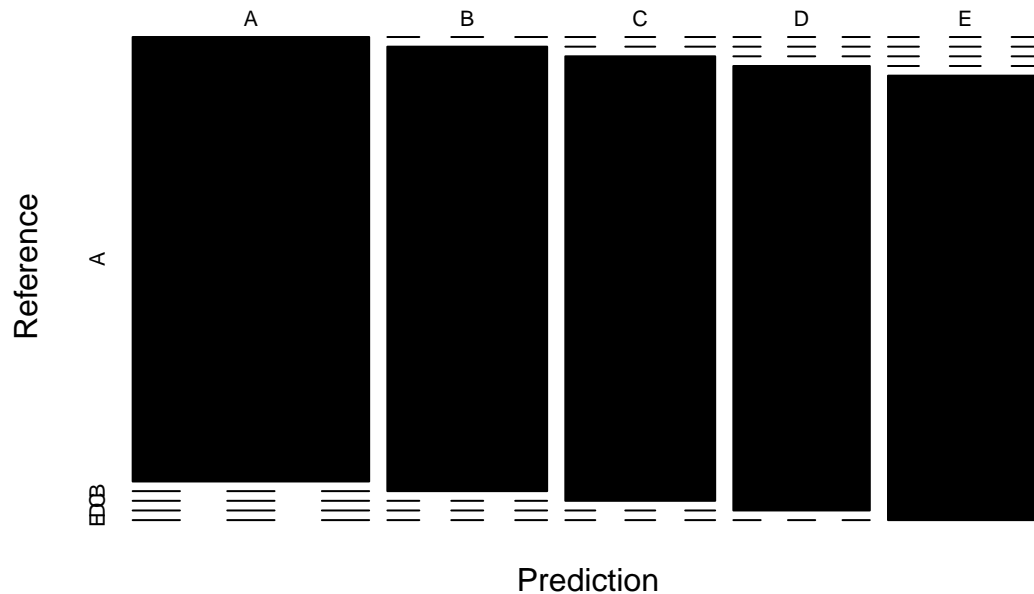
Let's plot the model

```
plot(modRF1)
```



```
plot(cmrf$table, col = cmrf$byClass, main = paste("Random Forest Confusion Matrix: Accuracy =", round(c
```


Random Forest Confusion Matrix: Accuracy = 1



Prediction with Generalized Boosted Regression Models

```
set.seed(12345)
controlGBM <- trainControl(method = "repeatedcv", number = 5, repeats = 1)
modGBM <- train(classe ~ ., data=trainData, method = "gbm", trControl = controlGBM, verbose = FALSE)
modGBM$finalModel
```

```
## A gradient boosted model with multinomial loss function.
## 150 iterations were performed.
## There were 52 predictors of which 52 had non-zero influence.
```

```
# print model summary
print(modGBM)
```

```
## Stochastic Gradient Boosting
##
## 13737 samples
## 52 predictor
## 5 classes: 'A', 'B', 'C', 'D', 'E'
##
## No pre-processing
## Resampling: Cross-Validated (5 fold, repeated 1 times)
## Summary of sample sizes: 10990, 10990, 10989, 10991, 10988
## Resampling results across tuning parameters:
##
## interaction.depth  n.trees  Accuracy  Kappa
## 1                  50       0.7521285  0.6858434
## 1                  100       0.8227397  0.7756753
## 1                  150       0.8521496  0.8129547
## 2                   50       0.8563724  0.8180344
## 2                  100       0.9059465  0.8809760
## 2                  150       0.9302623  0.9117412
```

```
##      3              50      0.8969931 0.8695557
##      3              100     0.9398712 0.9238994
##      3              150     0.9593802 0.9486037
##
## Tuning parameter 'shrinkage' was held constant at a value of 0.1
##
## Tuning parameter 'n.minobsinnode' was held constant at a value of 10
## Accuracy was used to select the optimal model using the largest value.
## The final values used for the model were n.trees = 150, interaction.depth =
## 3, shrinkage = 0.1 and n.minobsinnode = 10.
```

Validate the GBM model and

```
predictGBM <- predict(modGBM, newdata=testData)
cmGBM <- confusionMatrix(predictGBM, testData$classe)
cmGBM
```

```
## Confusion Matrix and Statistics
##
##              Reference
## Prediction    A    B    C    D    E
##      A 1155    20    0    0    1
##      B   9   754   17    5    6
##      C   1   12  713   16    3
##      D   0    1    6  647    8
##      E   0    0    2    4  743
##
## Overall Statistics
##
##              Accuracy : 0.9731
##              95% CI : (0.9677, 0.9778)
##      No Information Rate : 0.2826
##      P-Value [Acc > NIR] : < 2.2e-16
##
##              Kappa : 0.966
##
## Mcnemar's Test P-Value : NA
##
## Statistics by Class:
##
##              Class: A Class: B Class: C Class: D Class: E
## Sensitivity          0.9914   0.9581   0.9661   0.9628   0.9763
## Specificity          0.9929   0.9889   0.9905   0.9957   0.9982
## Pos Pred Value       0.9821   0.9532   0.9570   0.9773   0.9920
## Neg Pred Value       0.9966   0.9901   0.9926   0.9928   0.9947
## Prevalence           0.2826   0.1909   0.1790   0.1630   0.1846
## Detection Rate       0.2801   0.1829   0.1729   0.1569   0.1802
## Detection Prevalence 0.2852   0.1919   0.1807   0.1606   0.1817
## Balanced Accuracy    0.9922   0.9735   0.9783   0.9792   0.9873
```

The accuracy rate using the random forest is very high: Accuracy : 0.9736 and therefore the *out-of-sample-error is equal to 0.0264**.

Applying the best model to the validation data

By comparing the accuracy rate values of the three models, it is clear the the 'Random Forest' model is the winner. So will use it on the validation data

```
Results <- predict(modRF1, newdata=validData)
Results
```

```
## [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
```

These Results of the output will be used to answer the “Course Project Prediction Quiz”