

Practical Machine Learning Course Project

A. Sachin Vianney

October 7, 2018

Executive Summary

Three different Machine learning algorithm are applied to the training set and is used to predict using the test set with varied accuracies. We can infer from the analysis that Random Forest is the best model and is applied to the test set.

We first need to include all the essential libraries/packages

```
library(ggplot2)
library(rattle)

## Rattle: A free graphical interface for data science with R.
## Version 5.2.0 Copyright (c) 2006-2018 Togaware Pty Ltd.
## Type 'rattle()' to shake, rattle, and roll your data.

library(caret)

## Loading required package: lattice

library(rpart)
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:rattle':
##
##     importance

## The following object is masked from 'package:ggplot2':
##
##     margin
```

Initailly, we need to download the test and training datasets from the URLs given to us. We set our working directory and proceed with the download followed by reading of the datasets.

```
#INPUTTING THE DATA
training <-
```

```
read.csv(url("https://d396qusza40orc.cloudfront.net/predmachlearn/pml-
training.csv"),header=TRUE)
test <-
read.csv(url("https://d396qusza40orc.cloudfront.net/predmachlearn/pml-
testing.csv"),header=TRUE)
```

Displaying the raw uncleaned data

#DISPLAYING THE RAW DATA

```
str(training)

## 'data.frame':    19622 obs. of  160 variables:
## $ X                      : int  1 2 3 4 5 6 7 8 9 10 ...
## $ user_name              : Factor w/ 6 levels "adelmo","carlitos",...: 2
2 2 2 2 2 2 2 2 2 ...
## $ raw_timestamp_part_1   : int  1323084231 1323084231 1323084231
1323084232 1323084232 1323084232 1323084232 1323084232 1323084232
...
## $ raw_timestamp_part_2   : int  788290 808298 820366 120339 196328
304277 368296 440390 484323 484434 ...
## $ cvtd_timestamp        : Factor w/ 20 levels "02/12/2011 13:32",...: 9
9 9 9 9 9 9 9 9 9 ...
## $ new_window            : Factor w/ 2 levels "no","yes": 1 1 1 1 1 1 1
1 1 1 ...
## $ num_window            : int  11 11 11 12 12 12 12 12 12 12 ...
## $ roll_belt             : num  1.41 1.41 1.42 1.48 1.48 1.45 1.42 1.42
1.43 1.45 ...
## $ pitch_belt            : num  8.07 8.07 8.07 8.05 8.07 8.06 8.09 8.13
8.16 8.17 ...
## $ yaw_belt              : num  -94.4 -94.4 -94.4 -94.4 -94.4 -94.4 -
94.4 -94.4 -94.4 -94.4 ...
## $ total_accel_belt      : int  3 3 3 3 3 3 3 3 3 3 ...
## $ kurtosis_roll_belt    : Factor w/ 397 levels "", "-0.016850",...: 1 1 1
1 1 1 1 1 1 ...
## $ kurtosis_pitch_belt   : Factor w/ 317 levels "", "-0.021887",...: 1 1 1
1 1 1 1 1 1 ...
## $ kurtosis_yaw_belt     : Factor w/ 2 levels "", "#DIV/0!": 1 1 1 1 1 1
1 1 1 1 ...
## $ skewness_roll_belt    : Factor w/ 395 levels "", "-0.003095",...: 1 1 1
1 1 1 1 1 1 ...
## $ skewness_roll_belt.1  : Factor w/ 338 levels "", "-0.005928",...: 1 1 1
1 1 1 1 1 1 ...
## $ skewness_yaw_belt     : Factor w/ 2 levels "", "#DIV/0!": 1 1 1 1 1 1
1 1 1 1 ...
## $ max_roll_belt         : num  NA NA NA NA NA NA NA NA NA NA ...
## $ max_pitch_belt        : int  NA NA NA NA NA NA NA NA NA NA ...
## $ max_yaw_belt          : Factor w/ 68 levels "", "-0.1", "-0.2",...: 1 1
1 1 1 1 1 1 ...
## $ min_roll_belt         : num  NA NA NA NA NA NA NA NA NA NA ...
```

```

## $ min_pitch_belt      : int  NA NA NA NA NA NA NA NA NA NA NA ...
## $ min_yaw_belt        : Factor w/ 68 levels "", "-0.1", "-0.2", ...: 1 1
1 1 1 1 1 1 1 1 ...
## $ amplitude_roll_belt : num  NA NA NA NA NA NA NA NA NA NA NA ...
## $ amplitude_pitch_belt : int   NA NA NA NA NA NA NA NA NA NA NA ...
## $ amplitude_yaw_belt   : Factor w/ 4 levels "", "#DIV/0!", "0.00", ...: 1
1 1 1 1 1 1 1 1 ...
## $ var_total_accel_belt : num  NA NA NA NA NA NA NA NA NA NA NA ...
## $ avg_roll_belt        : num  NA NA NA NA NA NA NA NA NA NA NA ...
## $ stddev_roll_belt     : num  NA NA NA NA NA NA NA NA NA NA NA ...
## $ var_roll_belt        : num  NA NA NA NA NA NA NA NA NA NA NA ...
## $ avg_pitch_belt       : num  NA NA NA NA NA NA NA NA NA NA NA ...
## $ stddev_pitch_belt    : num  NA NA NA NA NA NA NA NA NA NA NA ...
## $ var_pitch_belt       : num  NA NA NA NA NA NA NA NA NA NA NA ...
## $ avg_yaw_belt         : num  NA NA NA NA NA NA NA NA NA NA NA ...
## $ stddev_yaw_belt      : num  NA NA NA NA NA NA NA NA NA NA NA ...
## $ var_yaw_belt         : num  NA NA NA NA NA NA NA NA NA NA NA ...
## $ gyros_belt_x         : num  0 0.02 0 0.02 0.02 0.02 0.02 0.02 0.02 0.02
0.03 ...
## $ gyros_belt_y         : num  0 0 0 0 0.02 0 0 0 0 0 ...
## $ gyros_belt_z         : num  -0.02 -0.02 -0.02 -0.03 -0.02 -0.02 -
0.02 -0.02 -0.02 0 ...
## $ accel_belt_x         : int    -21 -22 -20 -22 -21 -21 -22 -22 -20 -21
...
## $ accel_belt_y         : int     4 4 5 3 2 4 3 4 2 4 ...
## $ accel_belt_z         : int    22 22 23 21 24 21 21 21 24 22 ...
## $ magnet_belt_x        : int     -3 -7 -2 -6 -6 0 -4 -2 1 -3 ...
## $ magnet_belt_y        : int    599 608 600 604 600 603 599 603 602 609
...
## $ magnet_belt_z        : int   -313 -311 -305 -310 -302 -312 -311 -313
-312 -308 ...
## $ roll_arm             : num   -128 -128 -128 -128 -128 -128 -128 -128
-128 -128 ...
## $ pitch_arm            : num    22.5 22.5 22.5 22.1 22.1 22 21.9 21.8
21.7 21.6 ...
## $ yaw_arm              : num   -161 -161 -161 -161 -161 -161 -161 -161
-161 -161 ...
## $ total_accel_arm      : int    34 34 34 34 34 34 34 34 34 34 ...
## $ var_accel_arm        : num  NA NA NA NA NA NA NA NA NA NA NA ...
## $ avg_roll_arm         : num  NA NA NA NA NA NA NA NA NA NA NA ...
## $ stddev_roll_arm      : num  NA NA NA NA NA NA NA NA NA NA NA ...
## $ var_roll_arm         : num  NA NA NA NA NA NA NA NA NA NA NA ...
## $ avg_pitch_arm        : num  NA NA NA NA NA NA NA NA NA NA NA ...
## $ stddev_pitch_arm     : num  NA NA NA NA NA NA NA NA NA NA NA ...
## $ var_pitch_arm        : num  NA NA NA NA NA NA NA NA NA NA NA ...
## $ avg_yaw_arm          : num  NA NA NA NA NA NA NA NA NA NA NA ...
## $ stddev_yaw_arm       : num  NA NA NA NA NA NA NA NA NA NA NA ...
## $ var_yaw_arm          : num  NA NA NA NA NA NA NA NA NA NA NA ...
## $ gyros_arm_x          : num  0 0.02 0.02 0.02 0 0.02 0 0.02 0.02 0.02
...

```

```

## $ gyros_arm_y      : num  0 -0.02 -0.02 -0.03 -0.03 -0.03 -0.03 -
0.02 -0.03 -0.03 ...
## $ gyros_arm_z      : num  -0.02 -0.02 -0.02 0.02 0 0 0 0 -0.02 -
0.02 ...
## $ accel_arm_x      : int   -288 -290 -289 -289 -289 -289 -289 -289
-288 -288 ...
## $ accel_arm_y      : int   109 110 110 111 111 111 111 111 109 110
...
## $ accel_arm_z      : int   -123 -125 -126 -123 -123 -122 -125 -124
-122 -124 ...
## $ magnet_arm_x     : int   -368 -369 -368 -372 -374 -369 -373 -372
-369 -376 ...
## $ magnet_arm_y     : int   337 337 344 344 337 342 336 338 341 334
...
## $ magnet_arm_z     : int   516 513 513 512 506 513 509 510 518 516
...
## $ kurtosis_roll_arm : Factor w/ 330 levels "", "-0.02438",...: 1 1 1
1 1 1 1 1 1 1 ...
## $ kurtosis_pitch_arm : Factor w/ 328 levels "", "-0.00484",...: 1 1 1
1 1 1 1 1 1 1 ...
## $ kurtosis_yaw_arm   : Factor w/ 395 levels "", "-0.01548",...: 1 1 1
1 1 1 1 1 1 1 ...
## $ skewness_roll_arm  : Factor w/ 331 levels "", "-0.00051",...: 1 1 1
1 1 1 1 1 1 1 ...
## $ skewness_pitch_arm : Factor w/ 328 levels "", "-0.00184",...: 1 1 1
1 1 1 1 1 1 1 ...
## $ skewness_yaw_arm   : Factor w/ 395 levels "", "-0.00311",...: 1 1 1
1 1 1 1 1 1 1 ...
## $ max_roll_arm       : num   NA NA NA NA NA NA NA NA NA NA NA ...
## $ max_pitch_arm      : num   NA NA NA NA NA NA NA NA NA NA NA ...
## $ max_yaw_arm        : int   NA NA NA NA NA NA NA NA NA NA NA ...
## $ min_roll_arm       : num   NA NA NA NA NA NA NA NA NA NA NA ...
## $ min_pitch_arm      : num   NA NA NA NA NA NA NA NA NA NA NA ...
## $ min_yaw_arm        : int   NA NA NA NA NA NA NA NA NA NA NA ...
## $ amplitude_roll_arm : num   NA NA NA NA NA NA NA NA NA NA NA ...
## $ amplitude_pitch_arm : num   NA NA NA NA NA NA NA NA NA NA NA ...
## $ amplitude_yaw_arm   : int   NA NA NA NA NA NA NA NA NA NA NA ...
## $ roll_dumbbell      : num   13.1 13.1 12.9 13.4 13.4 ...
## $ pitch_dumbbell     : num   -70.5 -70.6 -70.3 -70.4 -70.4 ...
## $ yaw_dumbbell       : num   -84.9 -84.7 -85.1 -84.9 -84.9 ...
## $ kurtosis_roll_dumbbell : Factor w/ 398 levels "", "-0.0035", "-
0.0073",...: 1 1 1 1 1 1 1 1 1 1 1 ...
## $ kurtosis_pitch_dumbbell : Factor w/ 401 levels "", "-0.0163", "-
0.0233",...: 1 1 1 1 1 1 1 1 1 1 1 ...
## $ kurtosis_yaw_dumbbell  : Factor w/ 2 levels "", "#DIV/0!": 1 1 1 1 1 1
1 1 1 1 ...
## $ skewness_roll_dumbbell : Factor w/ 401 levels "", "-0.0082", "-
0.0096",...: 1 1 1 1 1 1 1 1 1 1 1 ...
## $ skewness_pitch_dumbbell : Factor w/ 402 levels "", "-0.0053", "-
0.0084",...: 1 1 1 1 1 1 1 1 1 1 1 ...

```

```
## $ skewness_yaw_dumbbell : Factor w/ 2 levels "", "#DIV/0!": 1 1 1 1 1 1
1 1 1 1 ...
## $ max_roll_dumbbell : num NA NA NA NA NA NA NA NA NA NA NA ...
## $ max_pitch_dumbbell : num NA NA NA NA NA NA NA NA NA NA NA ...
## $ max_yaw_dumbbell : Factor w/ 73 levels "", "-0.1", "-0.2", ...: 1 1
1 1 1 1 1 1 1 1 ...
## $ min_roll_dumbbell : num NA NA NA NA NA NA NA NA NA NA NA ...
## $ min_pitch_dumbbell : num NA NA NA NA NA NA NA NA NA NA NA ...
## $ min_yaw_dumbbell : Factor w/ 73 levels "", "-0.1", "-0.2", ...: 1 1
1 1 1 1 1 1 1 1 ...
## $ amplitude_roll_dumbbell : num NA NA NA NA NA NA NA NA NA NA NA ...
## [list output truncated]
```

As we can see, there are a lot of NA values, we need to clean this data in order to make our analysis easier and more meaningful

#DATA CLEANING

```
na_count <- sapply(test, function(y) sum((is.na(y))))
na <- na_count[na_count == 20]
good <- names(na)
training <- training[,!(names(training) %in% good)]
test <- test[,!(names(test) %in% good)]
good2 <- c('user_name', 'raw_timestamp_part_1', 'raw_timestamp_part_2',
'cvtd_timestamp', 'new_window', 'num_window', 'X')
training <- training[,!(names(training) %in% good2)]
test <- test[,!(names(test) %in% good2)]
```

Displaying the clean data, we can see a reduction in the number of rows

#DISPLAYING THE CLEAN DATA

```
str(training)

## 'data.frame': 19622 obs. of 53 variables:
## $ roll_belt : num 1.41 1.41 1.42 1.48 1.48 1.45 1.42 1.42 1.43
1.45 ...
## $ pitch_belt : num 8.07 8.07 8.07 8.05 8.07 8.06 8.09 8.13 8.16
8.17 ...
## $ yaw_belt : num -94.4 -94.4 -94.4 -94.4 -94.4 -94.4 -94.4 -
94.4 -94.4 -94.4 ...
## $ total_accel_belt : int 3 3 3 3 3 3 3 3 3 3 ...
## $ gyros_belt_x : num 0 0.02 0 0.02 0.02 0.02 0.02 0.02 0.02 0.03
...
## $ gyros_belt_y : num 0 0 0 0.02 0 0 0 0 0 ...
## $ gyros_belt_z : num -0.02 -0.02 -0.02 -0.03 -0.02 -0.02 -0.02 -
0.02 -0.02 0 ...
## $ accel_belt_x : int -21 -22 -20 -22 -21 -21 -22 -22 -20 -21 ...
## $ accel_belt_y : int 4 4 5 3 2 4 3 4 2 4 ...
## $ accel_belt_z : int 22 22 23 21 24 21 21 21 24 22 ...
## $ magnet_belt_x : int -3 -7 -2 -6 -6 0 -4 -2 1 -3 ...
```

```

## $ magnet_belt_y      : int  599 608 600 604 600 603 599 603 602 609 ...
## $ magnet_belt_z      : int  -313 -311 -305 -310 -302 -312 -311 -313 -312
-308 ...
## $ roll_arm           : num  -128 -128 -128 -128 -128 -128 -128 -128 -128
-128 ...
## $ pitch_arm          : num   22.5 22.5 22.5 22.1 22.1 22 21.9 21.8 21.7
21.6 ...
## $ yaw_arm            : num  -161 -161 -161 -161 -161 -161 -161 -161 -161
-161 ...
## $ total_accel_arm    : int   34 34 34 34 34 34 34 34 34 34 ...
## $ gyros_arm_x        : num   0 0.02 0.02 0.02 0 0.02 0 0.02 0.02 0.02 ...
## $ gyros_arm_y        : num   0 -0.02 -0.02 -0.03 -0.03 -0.03 -0.03 -0.02
-0.03 -0.03 ...
## $ gyros_arm_z        : num  -0.02 -0.02 -0.02 0.02 0 0 0 0 -0.02 -0.02
...
## $ accel_arm_x        : int  -288 -290 -289 -289 -289 -289 -289 -289 -288
-288 ...
## $ accel_arm_y        : int   109 110 110 111 111 111 111 111 109 110 ...
## $ accel_arm_z        : int  -123 -125 -126 -123 -123 -122 -125 -124 -122
-124 ...
## $ magnet_arm_x       : int  -368 -369 -368 -372 -374 -369 -373 -372 -369
-376 ...
## $ magnet_arm_y       : int   337 337 344 344 337 342 336 338 341 334 ...
## $ magnet_arm_z       : int   516 513 513 512 506 513 509 510 518 516 ...
## $ roll_dumbbell      : num   13.1 13.1 12.9 13.4 13.4 ...
## $ pitch_dumbbell     : num  -70.5 -70.6 -70.3 -70.4 -70.4 ...
## $ yaw_dumbbell       : num  -84.9 -84.7 -85.1 -84.9 -84.9 ...
## $ total_accel_dumbbell : int   37 37 37 37 37 37 37 37 37 37 ...
## $ gyros_dumbbell_x    : num   0 0 0 0 0 0 0 0 0 0 ...
## $ gyros_dumbbell_y    : num  -0.02 -0.02 -0.02 -0.02 -0.02 -0.02 -0.02 -
0.02 -0.02 -0.02 ...
## $ gyros_dumbbell_z    : num   0 0 0 -0.02 0 0 0 0 0 0 ...
## $ accel_dumbbell_x    : int  -234 -233 -232 -232 -233 -234 -232 -234 -232
-235 ...
## $ accel_dumbbell_y    : int   47 47 46 48 48 48 47 46 47 48 ...
## $ accel_dumbbell_z    : int  -271 -269 -270 -269 -270 -269 -270 -272 -269
-270 ...
## $ magnet_dumbbell_x   : int  -559 -555 -561 -552 -554 -558 -551 -555 -549
-558 ...
## $ magnet_dumbbell_y   : int   293 296 298 303 292 294 295 300 292 291 ...
## $ magnet_dumbbell_z   : num  -65 -64 -63 -60 -68 -66 -70 -74 -65 -69 ...
## $ roll_forearm       : num   28.4 28.3 28.3 28.1 28 27.9 27.9 27.8 27.7
27.7 ...
## $ pitch_forearm       : num  -63.9 -63.9 -63.9 -63.9 -63.9 -63.9 -63.9 -
63.8 -63.8 -63.8 ...
## $ yaw_forearm        : num  -153 -153 -152 -152 -152 -152 -152 -152 -152
-152 ...
## $ total_accel_forearm : int   36 36 36 36 36 36 36 36 36 36 ...
## $ gyros_forearm_x     : num   0.03 0.02 0.03 0.02 0.02 0.02 0.02 0.02 0.02
0.03
0.02 ...

```

```
## $ gyros_forearm_y      : num  0 0 -0.02 -0.02 0 -0.02 0 -0.02 0 0 ...
## $ gyros_forearm_z      : num  -0.02 -0.02 0 0 -0.02 -0.03 -0.02 0 -0.02 -
0.02 ...
## $ accel_forearm_x      : int   192 192 196 189 189 193 195 193 193 190 ...
## $ accel_forearm_y      : int   203 203 204 206 206 203 205 205 204 205 ...
## $ accel_forearm_z      : int   -215 -216 -213 -214 -214 -215 -215 -213 -214
-215 ...
## $ magnet_forearm_x     : int   -17 -18 -18 -16 -17 -9 -18 -9 -16 -22 ...
## $ magnet_forearm_y     : num   654 661 658 658 655 660 659 660 653 656 ...
## $ magnet_forearm_z     : num   476 473 469 469 473 478 470 474 476 473 ...
## $ classe               : Factor w/ 5 levels "A","B","C","D",...: 1 1 1 1 1
1 1 1 1 1 ...
```

Machine Learning Algorithms

For this analysis, we will be using three algorithmic approaches The Three algorithms applied are: 1)Decision Tree 2)Random Forest 3)Generalized Boosted Regression

```
#MACHINE LEARNING ALGORITHMS IMPLEMENTATION
```

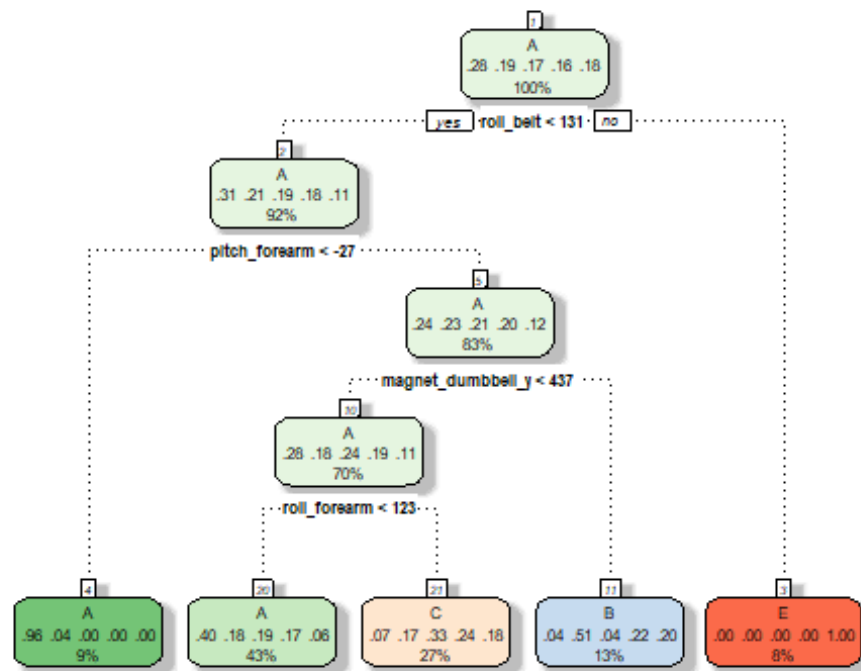
```
set.seed(49)
part <- createDataPartition(training$classe, p = 0.75, list = FALSE)
training1 <- training[part,]
test1 <- training[-part,]
```

I have used the cross-validation technique in order to improve efficiency and limit overfitting.I have used 10 folds.

```
tr_cont <- trainControl(method="cv", number=10)
```

Descision Trees

```
modell1 <- train(classe~., data=training1, method="rpart", trControl=tr_cont)
fancyRpartPlot(modell1$finalModel)
```



Rattle 2018-Oct-07 20:42:08 mahe

```
prediction1 <- predict(model1,newdata=test1)
conf_mat1 <- confusionMatrix(test1$classe,prediction1)
conf_mat1$table
```

```
##           Reference
## Prediction   A    B    C    D    E
##           A 1267   22   98    0    8
##           B  398  332  219    0    0
##           C  396   35  424    0    0
##           D  352  168  284    0    0
##           E  140  123  235    0  403
```

```
conf_mat1$overall[1]
```

```
## Accuracy
## 0.4946982
```

Random Forest

```
model2 <- train(classe~., data=training1,
method="rf",trControl=tr_cont,verbose = FALSE)
print(model2)
```

```
## Random Forest
##
## 14718 samples
## 52 predictor
```



```
##      5 classes: 'A', 'B', 'C', 'D', 'E'
##
## No pre-processing
## Resampling: Cross-Validated (10 fold)
## Summary of sample sizes: 13246, 13248, 13246, 13246, 13245, 13246, ...
## Resampling results across tuning parameters:
##
##  mtry  Accuracy   Kappa
##    2    0.9925262 0.9905445
##   27    0.9927292 0.9908018
##   52    0.9880410 0.9848702
##
## Accuracy was used to select the optimal model using the largest value.
## The final value used for the model was mtry = 27.

prediction2 <- predict(model2,newdata=test1)
conf_mat2 <- confusionMatrix(test1$classe,prediction2)
conf_mat2$table

##              Reference
## Prediction      A      B      C      D      E
##           A 1394      1      0      0      0
##           B    9   939      1      0      0
##           C    0    7   846      2      0
##           D    0    1    4   798      1
##           E    0    1    2    1   897

conf_mat2$overall[1]

## Accuracy
## 0.9938825
```

Generalized Boosted Regression

```
model3 <- train(classe~., data=training1, method="gbm", trControl=tr_cont,
verbose=FALSE)
print(model3)

## Stochastic Gradient Boosting
##
## 14718 samples
##   52 predictor
##   5 classes: 'A', 'B', 'C', 'D', 'E'
##
## No pre-processing
## Resampling: Cross-Validated (10 fold)
## Summary of sample sizes: 13247, 13248, 13246, 13245, 13246, 13245, ...
## Resampling results across tuning parameters:
##
##  interaction.depth  n.trees  Accuracy   Kappa
##         1              50    0.7576432 0.6927629
```

```
##      1      100      0.8187920 0.7706153
##      1      150      0.8537841 0.8149065
##      2       50      0.8543951 0.8155150
##      2      100      0.9075262 0.8829763
##      2      150      0.9298819 0.9112773
##      3       50      0.8980840 0.8709724
##      3      100      0.9415677 0.9260695
##      3      150      0.9629701 0.9531596
##
## 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.

prediction3 <- predict(model3,newdata=test1)
conf_mat3 <- confusionMatrix(test1$classe,prediction3)
conf_mat3$table

##           Reference
## Prediction    A    B    C    D    E
##           A 1375   16    4    0    0
##           B   36  888   24    0    1
##           C    0   28  816    8    3
##           D    0    5   19   77    3
##           E    3   13   10   18  857

conf_mat3$overall[1]

## Accuracy
## 0.9610522
```

The Final Result/Prediction

```
final <- predict(model2,newdata= test)
final

## [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
```