

Practical Machine Learning - Joe Pogson

Joe Pogson

4/21/2020

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, your goal will be to 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.

More information is available from the website here: <http://web.archive.org/web/20161224072740/http://groupware.les.inf.puc-rio.br/har>

Library and Data loading.

```
library(caret)
```

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

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

```
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.
```

```
library(e1071)
```

```
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
```

```

library(gbm)

## Loaded gbm 2.1.5
library(survival)

##
## Attaching package: 'survival'
## The following object is masked from 'package:caret':
##
##   cluster
library(splines)
library(parallel)
library(plyr)

TrainData <- read.csv(url("https://d396qusza40orc.cloudfront.net/predmachlearn/pml-training.csv"),header=T)
dim(TrainData)

## [1] 19622  160

TestData <-read.csv(url("https://d396qusza40orc.cloudfront.net/predmachlearn/pml-testing.csv"),header=T)
dim(TestData)

## [1]  20 160

str(TrainData)

## '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 ...
## $ raw_timestamp_part_2 : int  788290 808298 820366 120339 196328 304277 368296 440390 484323 484...
## $ 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 1 ...
## $ kurtosis_pitch_belt : Factor w/ 317 levels "","-0.021887",...: 1 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 1 ...
## $ skewness_roll_belt.1 : Factor w/ 338 levels "","-0.005928",...: 1 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 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 ...
## $ 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 ...
## $ amplitude_pitch_belt : int  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 1 ...
## $ var_total_accel_belt : num  NA NA NA NA NA NA NA NA NA NA ...
## $ avg_roll_belt : num  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 0.03 ...
## $ 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 ...
## $ kurtosis_pitch_dumbbell : Factor w/ 401 levels "", "-0.0163", "-0.0233", ...: 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 ...
## $ skewness_pitch_dumbbell : Factor w/ 402 levels "", "-0.0053", "-0.0084", ...: 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]
```

Here we can see that the data is made up of 160 parameters from 19622 observations. We can remove some of these parameters as they are mainly NAs or are purely personal information not relevant for this model.

Cleaning Data

```
removeCol <- which(colSums(is.na(TrainData)|TrainData=="")>0.9*dim(TrainData)[1])
CleanTrainData <- TrainData[, -removeCol]
CleanTrainData <- CleanTrainData[, -(1:5)]
dim(CleanTrainData)
```

```
## [1] 19622    55
```

Repeat for the Test Data.

```
removeCol <- which(colSums(is.na(TestData)|TestData=="")>0.9*dim(TestData)[1])
CleanTestData <- TestData[, -removeCol]
CleanTestData <- CleanTestData[, -(1:5)]
dim(CleanTestData)
```

```
## [1] 20 55
```

We have now reduced the number of parameters to 55.

Training

In order to limit overfitting cross-validation will be used on a partitioned data set using 5-folds.

```
set.seed(12345)
inTrain1 <- createDataPartition(CleanTrainData$classe, p=0.75, list=FALSE)
Train1 <- CleanTrainData[inTrain1,]
Test1 <- CleanTrainData[-inTrain1,]
dim(Train1)
```

```
## [1] 14718    55
```

```
dim(Test1)
```

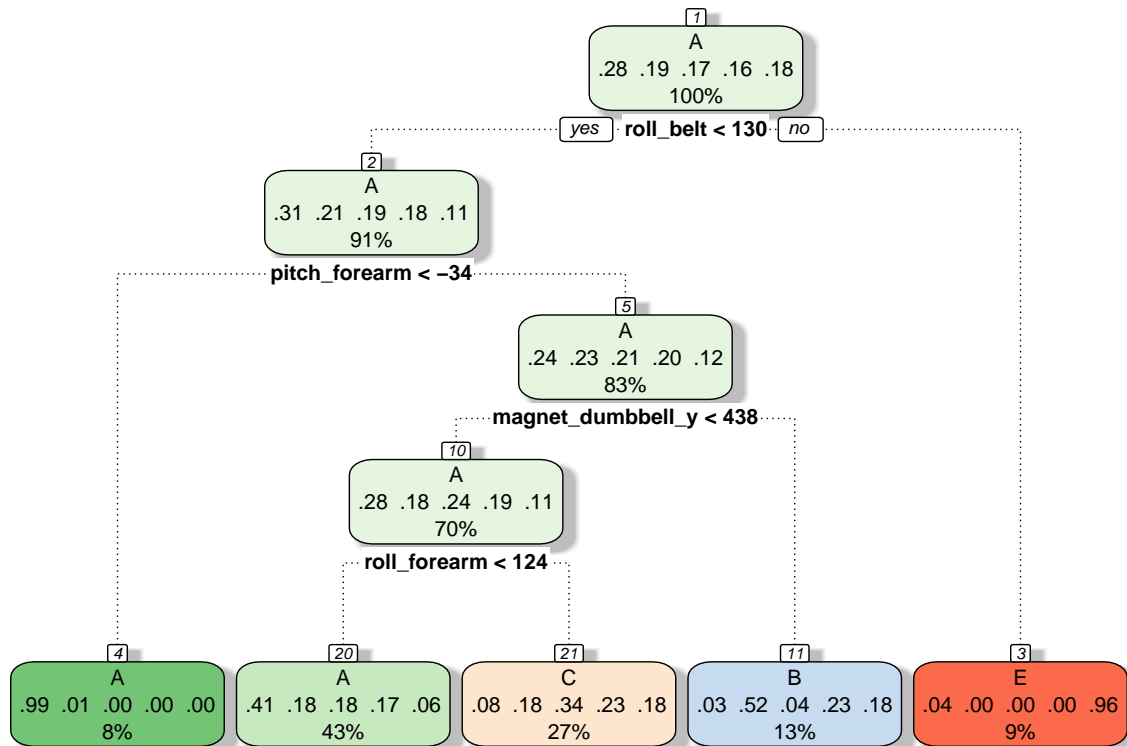
```
## [1] 4904    55
```

Classification tree

First we will test the data using a classification tree,

```
trControl <- trainControl(method="cv", number=5)
model_CT <- train(classe~., data=Train1, method="rpart", trControl=trControl)
```

```
fancyRpartPlot(model_CT$finalModel)
```



Rattle 2020-Apr-22 12:40:52 admin

Show the confusion matrix and accuracy of this model.

```
PredTrain <- predict(model_CT, newdata = Test1)
ConfMtxCT <- confusionMatrix(Test1$class, PredTrain)
ConfMtxCT$table
```

```
##           Reference
## Prediction   A    B    C    D    E
##           A 1252   30   90    0   23
##           B  396  317  236    0    0
##           C  434   24  397    0    0
##           D  343  151  310    0    0
##           E   114  132  229    0  426
```

```
ConfMtxCT$overall[1]
```

```
## Accuracy
## 0.4877651
```

Here we can see that the accuracy for this Model is 49% not a bad start.

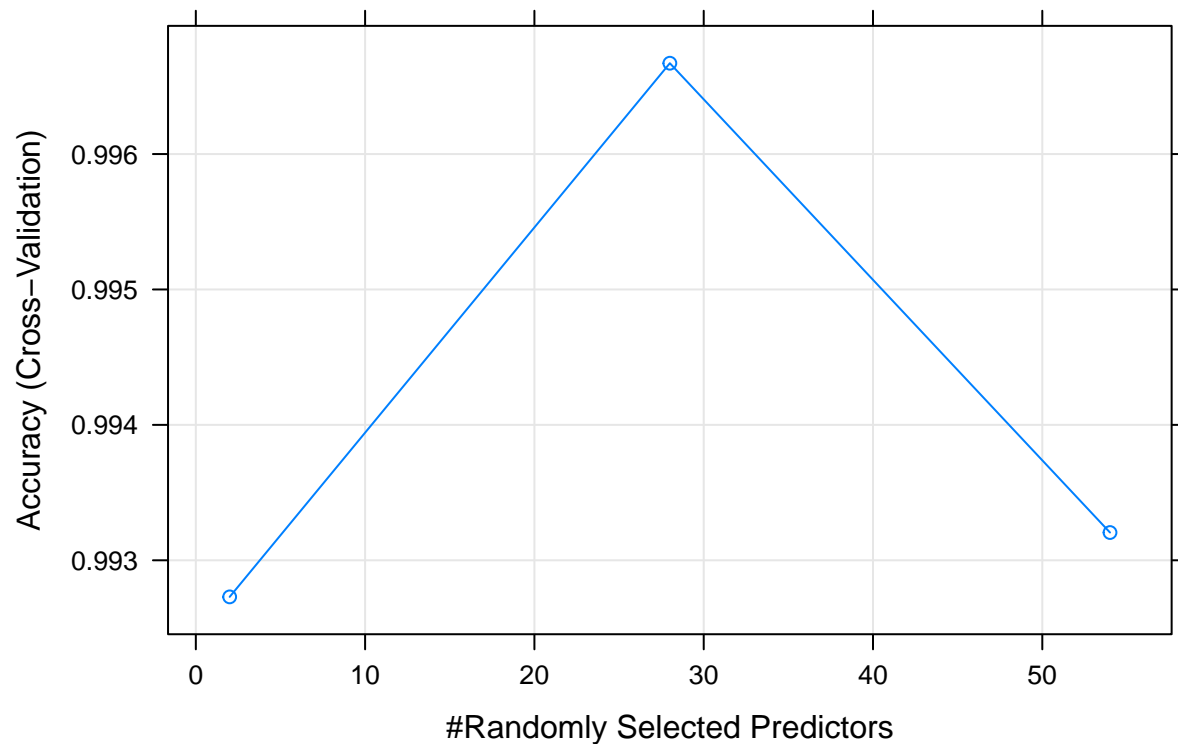
Random forests

Lets try again but with random forests this time.

```
controlRF <- trainControl(method="cv", number=3, verboseIter=FALSE)
RandFor <- train(classe ~ ., data=Train1, method="rf",
                 trControl=controlRF)
RandFor$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: 28
##
##           OOB estimate of  error rate: 0.2%
## Confusion matrix:
##      A    B    C    D    E class.error
## A 4183     1     0     0     1 0.0004778973
## B     5 2839     3     1     0 0.0031601124
## C     0     4 2562     1     0 0.0019477990
## D     0     0     9 2403     0 0.0037313433
## E     0     1     0     4 2701 0.0018477458
plot(RandFor,main="Accuracy vs number of predictors")
```

Accuracy vs number of predictors



```
trainpred <- predict(RandFor, newdata=Test1)
RFConfMTX <- confusionMatrix(Test1$classe, trainpred)
```

```
RFConfMTX$overall[1]
```

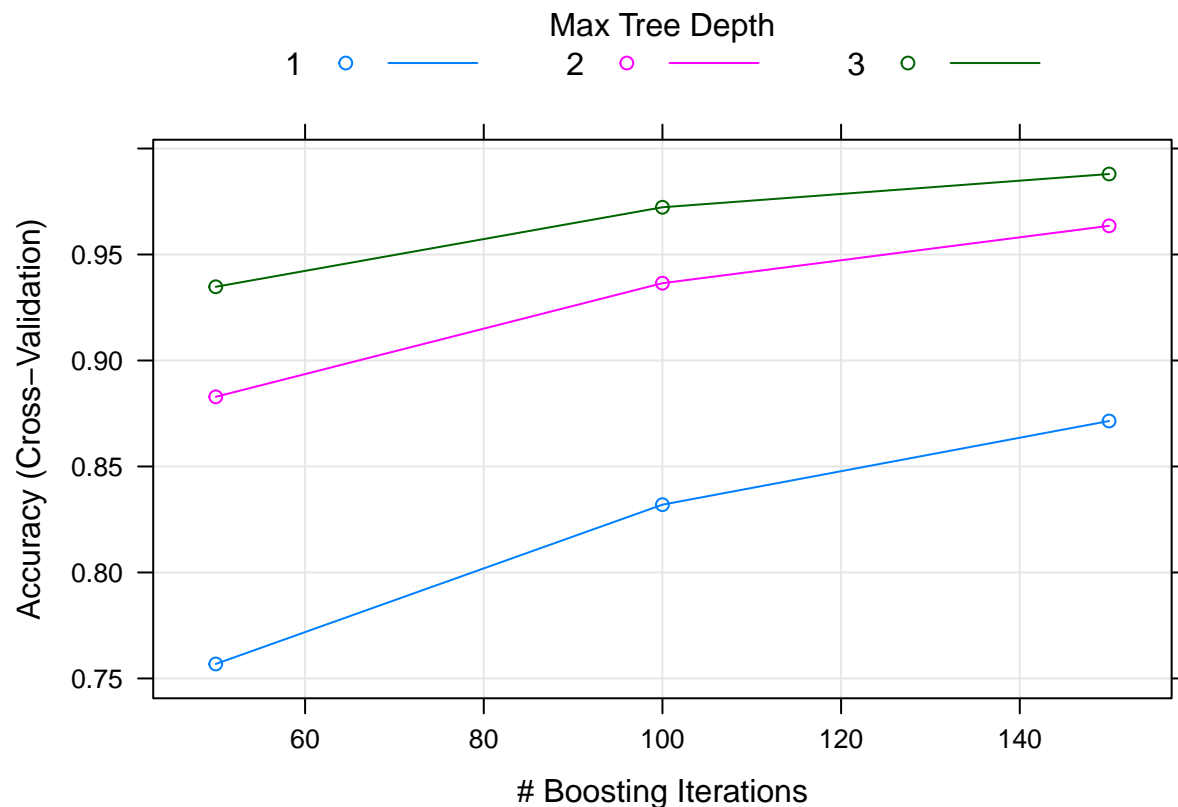
```
## Accuracy
## 0.9983687
```

With random forest we have increased our accuracy to almost 100% WOW.

Generalised Boosted Model

Just to see if we can possibly improve on the random forest method we will use the generalised boosted model.

```
model_GBM <- train(classe~., data=Train1, method="gbm", trControl=trControl, verbose=FALSE)
plot(model_GBM)
```



```
trainpred <- predict(model_GBM, newdata=Test1)
confMatGBM <- confusionMatrix(Test1$classe, trainpred)
confMatGBM$table
```

```
##           Reference
## Prediction   A   B   C   D   E
##           A 1391   4   0   0   0
##           B   10  928  10   1   0
```

```
##           C      0      4  850      1      0
##           D      0      1   18  779      6
##           E      0      5    2    4  890
```

```
confMatGBM$overall[1]
```

```
## Accuracy
## 0.9865416
```

With 5 folds the precision is 98.6%. Just a bit less than the Random Forest.

Conculsion

With an accuracy of 99.9% random forest is the best method to use.

```
FinalTestPred <- predict(RandFor,newdata=CleanTestData)
FinalTestPred
```

```
## [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
```