

ML_PROJECT

AYUSH BHARGAVA

November 28, 2017

Goal

We will analyze the exercise data collected on six people using wearable devices such as awbone Up, Nike FuelBand, and Fitbit. We are particularly interested in finding out how well they perform barbell lifts. We will build some predictive models to predict this characteristic on new datasets and try to find the best model.

Data

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>

Strategy

-This is a classification problem where we have to build predictive models for predicting outcome of a factor variable "classe" having 5 levels -As this is a classification problem we will use accuracy as a metric to check errors in our predictions -Random Forest and classification trees are best suited for classification problems so we will work with these algorithms and find which performs best

Loading libraries,data and doing preliminary analysis

```
suppressWarnings(suppressMessages(library(caret)))
suppressWarnings(suppressMessages(library(randomForest)))
suppressWarnings(suppressMessages(library(rpart)))
suppressWarnings(suppressMessages(library(rpart.plot)))
# setting the overall seed for reproducibility
set.seed(1234)
#loading training data set
training <- read.csv("C:/Users/Ayush Bhargava/Desktop/Practical ml project/pml-training.csv", na.strings=)
#loading testing data set
testing <- read.csv("C:/Users/Ayush Bhargava/Desktop/Practical ml project/pml-testing.csv", na.strings=)
#looking at the training 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 ...
## $ 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 ...
## $ kurtosis_roll_belt : num NA NA NA NA NA NA NA NA NA NA ...
## $ kurtosis_pitch_belt : num NA NA NA NA NA NA NA NA NA NA ...
## $ kurtosis_yaw_belt : logi NA NA NA NA NA NA ...
## $ skewness_roll_belt : num NA NA NA NA NA NA NA NA NA NA ...
## $ skewness_roll_belt.1 : num NA NA NA NA NA NA NA NA NA NA ...
## $ skewness_yaw_belt : logi NA NA NA NA NA NA ...
## $ 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 : num NA NA NA NA NA NA NA NA NA NA ...
## $ 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 : num NA NA NA NA NA NA NA NA NA NA ...
## $ 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 : num NA NA NA NA NA NA NA NA NA NA ...
## $ 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 ...
## $ var_roll_belt : num NA NA NA NA NA NA NA NA NA NA ...
## $ avg_pitch_belt : num NA NA NA NA NA NA NA NA NA NA ...
## $ stddev_pitch_belt : num NA NA NA NA NA NA NA NA NA NA ...
## $ var_pitch_belt : num NA NA NA NA NA NA NA NA NA NA ...
## $ avg_yaw_belt : num NA NA NA NA NA NA NA NA NA NA ...
## $ stddev_yaw_belt : num NA NA NA NA NA NA NA NA NA NA ...
## $ var_yaw_belt : num 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 ...
## $ avg_roll_arm : num NA NA NA NA NA NA NA NA NA NA ...
## $ stddev_roll_arm : num NA NA NA NA NA NA NA NA NA NA ...
## $ var_roll_arm : num NA NA NA NA NA NA NA NA NA NA ...
## $ avg_pitch_arm : num NA NA NA NA NA NA NA NA NA NA ...
## $ stddev_pitch_arm : num NA NA NA NA NA NA NA NA NA NA ...
## $ var_pitch_arm : num NA NA NA NA NA NA NA NA NA NA ...
## $ avg_yaw_arm : num NA NA NA NA NA NA NA NA NA NA ...
## $ stddev_yaw_arm : num NA NA NA NA NA NA NA NA NA NA ...
## $ var_yaw_arm : num 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 : num  NA NA NA NA NA NA NA NA NA NA NA ...
## $ kurtosis_pitch_arm : num  NA NA NA NA NA NA NA NA NA NA NA ...
## $ kurtosis_yaw_arm  : num  NA NA NA NA NA NA NA NA NA NA NA ...
## $ skewness_roll_arm : num  NA NA NA NA NA NA NA NA NA NA NA ...
## $ skewness_pitch_arm : num  NA NA NA NA NA NA NA NA NA NA NA ...
## $ skewness_yaw_arm  : num  NA NA NA NA NA NA NA NA NA NA NA ...
## $ 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 : num  NA NA NA NA NA NA NA NA NA NA NA ...
## $ kurtosis_pitch_dumbbell : num  NA NA NA NA NA NA NA NA NA NA NA ...
## $ kurtosis_yaw_dumbbell : logi  NA NA NA NA NA NA NA ...
## $ skewness_roll_dumbbell : num  NA NA NA NA NA NA NA NA NA NA NA ...
## $ skewness_pitch_dumbbell : num  NA NA NA NA NA NA NA NA NA NA NA ...
## $ skewness_yaw_dumbbell : logi  NA NA NA NA NA NA NA ...
## $ 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   : num  NA NA NA NA NA NA NA NA NA NA NA ...
## $ 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   : num  NA NA NA NA NA NA NA NA NA NA NA ...
## $ amplitude_roll_dumbbell : num  NA NA NA NA NA NA NA NA NA NA NA ...
## [list output truncated]
```

Analysis summary and Preprocessing the data

-There are many columns which have mostly “na” values thus we will get rid of them

```
training<-training[,colSums(is.na(training)) == 0]
testing <-testing[,colSums(is.na(testing)) == 0]
```

-There are many columns such as participant name etc(column no. 1 to 7) which are not usefull for this project thus we will get rid of these columns as well

```
training<-training[,~c(1:7)]
testing <-testing[,~c(1:7)]
dim(training)
```

```
## [1] 19622    53
```

```
dim(testing)
```

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

-We have enough observations in training data set to further divide training data in to subtraining(75%) and subtesting(25%) datasets for cross validation purposes

```
subsamples <- createDataPartition(y=training$classe, p=0.75, list=FALSE)
subTraining <- training[subsamples, ]
subTesting <- training[-subsamples, ]
dim(subTraining)
```

```
## [1] 14718 53
```

```
dim(subTesting)
```

```
## [1] 4904 53
```

Prediction Models

Classification tree

```
model1 <- rpart(classe ~ ., data=subTraining, method="class")
```

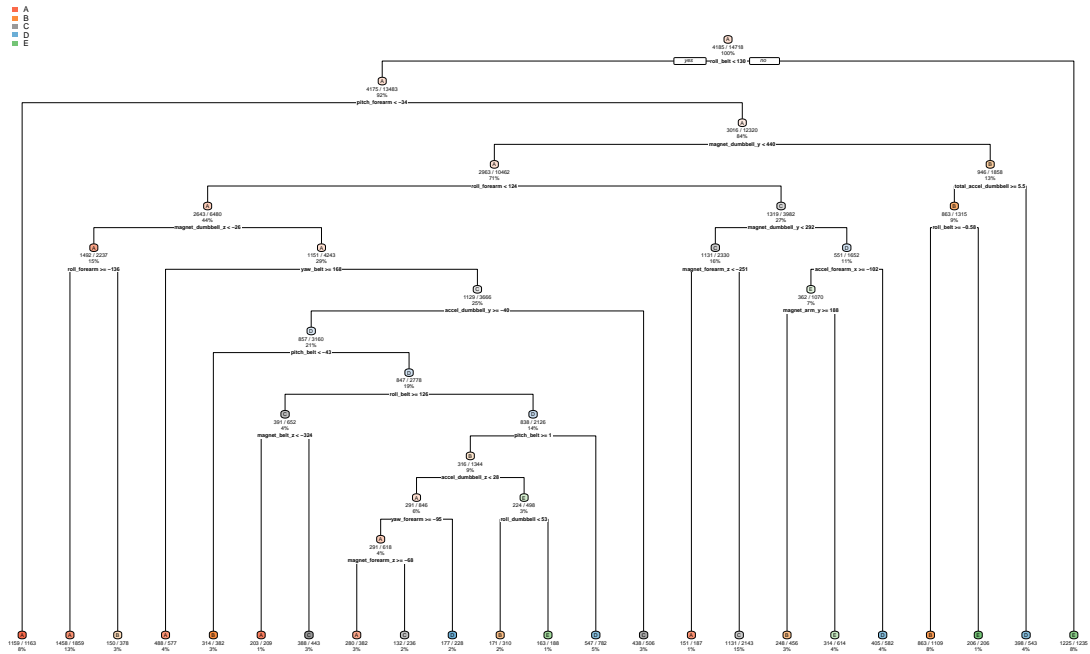
```
# Predictions:
```

```
prediction1 <- predict(model1, subTesting, type = "class")
```

```
# The Decision Tree
```

```
rpart.plot(model1, main="Classification Tree", extra=102, under=TRUE, faclen=0)
```

Classification Tree



```
# Test results on subTesting data set:
confusionMatrix(prediction1, subTesting$classe)
```

```
## Confusion Matrix and Statistics
```

```
##
##           Reference
## Prediction   A    B    C    D    E
##           A 1235  157   16   50   20
##           B   55  568   73   80  102
##           C   44  125  690  118  116
##           D   41   64   50  508   38
##           E   20   35   26   48  625
```

```
## Overall Statistics
```

```
##
##           Accuracy : 0.7394
##           95% CI : (0.7269, 0.7516)
##           No Information Rate : 0.2845
##           P-Value [Acc > NIR] : < 2.2e-16
```

```
##
##           Kappa : 0.6697
##           McNemar's Test P-Value : < 2.2e-16
```

```
## Statistics by Class:
```

```
##
##           Class: A Class: B Class: C Class: D Class: E
```

## Sensitivity	0.8853	0.5985	0.8070	0.6318	0.6937
## Specificity	0.9307	0.9216	0.9005	0.9529	0.9678
## Pos Pred Value	0.8356	0.6469	0.6313	0.7247	0.8289
## Neg Pred Value	0.9533	0.9054	0.9567	0.9296	0.9335
## Prevalence	0.2845	0.1935	0.1743	0.1639	0.1837
## Detection Rate	0.2518	0.1158	0.1407	0.1036	0.1274
## Detection Prevalence	0.3014	0.1790	0.2229	0.1429	0.1538
## Balanced Accuracy	0.9080	0.7601	0.8537	0.7924	0.8307

From the above confusion matrix result it is clear that classification tree did not perform well and was able to predict on our subTesting data with 74% accuracy. Lets build our second model

Random Forest

```
model2 <- randomForest(classe ~. , data=subTraining, method="class")

# Predicting:
prediction2 <- predict(model2, subTesting, type = "class")

# Test results on subTesting data set:
confusionMatrix(prediction2, subTesting$classe)
```

```
## Confusion Matrix and Statistics
##
##           Reference
## Prediction   A    B    C    D    E
##           A 1394     3     0     0     0
##           B     1   944    10     0     0
##           C     0     2   843     6     0
##           D     0     0     2   798     0
##           E     0     0     0     0   901
##
## Overall Statistics
##
##               Accuracy : 0.9951
##               95% CI : (0.9927, 0.9969)
##       No Information Rate : 0.2845
##       P-Value [Acc > NIR] : < 2.2e-16
##
##               Kappa : 0.9938
##   Mcnemar's Test P-Value : NA
##
## Statistics by Class:
##
##               Class: A Class: B Class: C Class: D Class: E
## Sensitivity      0.9993  0.9947  0.9860  0.9925  1.0000
## Specificity      0.9991  0.9972  0.9980  0.9995  1.0000
## Pos Pred Value   0.9979  0.9885  0.9906  0.9975  1.0000
## Neg Pred Value   0.9997  0.9987  0.9970  0.9985  1.0000
## Prevalence       0.2845  0.1935  0.1743  0.1639  0.1837
## Detection Rate   0.2843  0.1925  0.1719  0.1627  0.1837
## Detection Prevalence 0.2849  0.1947  0.1735  0.1631  0.1837
## Balanced Accuracy 0.9992  0.9960  0.9920  0.9960  1.0000
```

Random forest is able to predict with 99.5% accuracy and thus outperformed classification tree and as it is able to predict quite accurately on the subTesting dataset we will move ahead and predict on our actual testing data with model2

```
# predict outcome levels on the actual Testing data set using Random Forest algorithm
predictfinal <- predict(model2, testing, type="class")
predictfinal
```

```
##  1  2  3  4  5  6  7  8  9 10 11 12 13 14 15 16 17 18 19 20
##  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
```