

fitbitdata_analysis

Aparna

December 11, 2016

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. The goal is to predict the manner in which they did the exercise. This is the “classe” variable in the training set

Reading the dataset

```
urltrain <- "https://d396qusza40orc.cloudfront.net/predmachlearn/pml-training.csv"
urltest <- "https://d396qusza40orc.cloudfront.net/predmachlearn/pml-testing.csv"
download.file(urltrain, destfile = "D:/OLD DATA/d/R learning/machineL/training.csv")
download.file(urltest, destfile = "D:/OLD DATA/d/R learning/machineL/testing.csv")
# Using the commands:
training = read.csv("D:/OLD DATA/d/R learning/machineL/training.csv")
testing <- read.csv("D:/OLD DATA/d/R learning/machineL/testing.csv")
```

Understanding the dataset

```
dim(training)
```

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

```
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 1323
084232 1323084232 1323084232 1323084232 1323084232 ...
## $ raw_timestamp_part_2 : int 788290 808298 820366 120339 196328 304277 368296 440390 484
323 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 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 ...
## $ 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 -0.02 -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 ...
## $ max_pitch_arm    : num   NA NA NA NA NA NA NA NA NA NA ...
## $ max_yaw_arm      : int   NA NA NA NA NA NA NA NA NA NA ...
## $ min_roll_arm     : num   NA NA NA NA NA NA NA NA NA NA ...
## $ min_pitch_arm    : num   NA NA NA NA NA NA NA NA NA NA ...
## $ min_yaw_arm      : int   NA NA NA NA NA NA NA NA NA NA ...
## $ amplitude_roll_arm : num   NA NA NA NA NA NA NA NA NA NA ...
## $ amplitude_pitch_arm : num   NA NA NA NA NA NA NA NA NA NA ...
## $ amplitude_yaw_arm : int   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 1 ...
## $ kurtosis_pitch_dumbbell : Factor w/ 401 levels "", "-0.0163", "-0.0233",...: 1 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 1 ...
## $ skewness_pitch_dumbbell : Factor w/ 402 levels "", "-0.0053", "-0.0084",...: 1 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 ...
## $ max_pitch_dumbbell : num   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 ...

```

```
## $ min_pitch_dumbbell      : num  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 ...  
## [list output truncated]
```

```
levels(training$classe)
```

```
## [1] "A" "B" "C" "D" "E"
```

Load all libraries

```
library(caret)
```

```
## Warning: package 'caret' was built under R version 3.3.2
```

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

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

```
## Warning: package 'ggplot2' was built under R version 3.3.2
```

```
library(kernlab)
```

```
## Warning: package 'kernlab' was built under R version 3.3.2
```

```
##  
## Attaching package: 'kernlab'
```

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

```
library(ISLR)
```

```
## Warning: package 'ISLR' was built under R version 3.3.2
```

```
library(Hmisc)
```

```
## Warning: package 'Hmisc' was built under R version 3.3.2
```

```
## Loading required package: survival
```

```
## Warning: package 'survival' was built under R version 3.3.2
```

```
##  
## Attaching package: 'survival'
```

```
## The following object is masked from 'package:caret':  
##  
##   cluster
```

```
## Loading required package: Formula
```

```
## Warning: package 'Formula' was built under R version 3.3.2
```

```
##  
## Attaching package: 'Hmisc'
```

```
## The following objects are masked from 'package:base':  
##  
##   format.pval, round.POSIXt, trunc.POSIXt, units
```

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

```
## Warning: package 'rpart.plot' was built under R version 3.3.2
```

```
library(AppliedPredictiveModeling)
```

```
## Warning: package 'AppliedPredictiveModeling' was built under R version  
## 3.3.2
```

```
library(e1071)
```

```
## Warning: package 'e1071' was built under R version 3.3.2
```

```
##  
## Attaching package: 'e1071'
```

```
## The following object is masked from 'package:Hmisc':  
##  
##   impute
```

```
library(randomForest)
```

```
## Warning: package 'randomForest' was built under R version 3.3.2
```

```
## randomForest 4.6-12
```

```
## Type rfNews() to see new features/changes/bug fixes.
```

```
##  
## Attaching package: 'randomForest'
```

```
## The following object is masked from 'package:Hmisc':  
##  
##      combine
```

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

```
library(gbm)
```

```
## Warning: package 'gbm' was built under R version 3.3.2
```

```
## Loading required package: splines
```

```
## Loading required package: parallel
```

```
## Loaded gbm 2.1.1
```

Cleaning the data for variables without variance and removing variables with more than 60% NAs

```
nzv <- nearZeroVar(training)  
training <- training[, -nzv]  
dim(training)
```

```
## [1] 19622 100
```

```

training_1 <- training
for(i in 1:length(training))
  if( sum( is.na( training[, i] ) ) /nrow(training) >= .6 )
    for(j in 1:length(training_1))
      if( length( grep(names(training[i]), names(training_1)[j]) ) ==1)
        training_1 <- training_1[ , -j]

dim(training_1)

```

```
## [1] 19622    59
```

```
training <- training_1
```

Variables related with data acquisition (like: id, timestamps, individuals' names, etc.) are not suitable to be used in prediction and are removed

```

training <- training[, -(1:6)]
training$classe = factor(training$classe)

```

Within training, create a training and testing dataset for building the model and validating it

```

inTrain = createDataPartition(training$classe, p = 3/4)[[1]]
train = training[ inTrain,]
test = training[-inTrain,]

```

Train 2 different models and check the accuracy on out of sample dataset

```

tc <- trainControl(method = "cv", number = 7, verboseIter=FALSE , preProcOptions="pca", allowParallel=TRUE)
mod1 <- train(classe ~ ., data=train, method="rf", trControl= tc)
mod2<- train(classe ~ ., data=train, method ="svmRadial", trControl= tc)

pred1 <- predict(mod1, test)
pred2 <- predict(mod2, test)

confusionMatrix(pred1, test$classe)

```



```
## Confusion Matrix and Statistics
##
##           Reference
## Prediction   A    B    C    D    E
##           A 1393    2    0    0    0
##           B    0  945    2    0    0
##           C    1    2  849   11    0
##           D    0    0    4  793    1
##           E    1    0    0    0  900
##
## 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           0.9986   0.9958   0.9930   0.9863   0.9989
## Specificity           0.9994   0.9995   0.9965   0.9988   0.9998
## Pos Pred Value        0.9986   0.9979   0.9838   0.9937   0.9989
## Neg Pred Value        0.9994   0.9990   0.9985   0.9973   0.9998
## Prevalence            0.2845   0.1935   0.1743   0.1639   0.1837
## Detection Rate        0.2841   0.1927   0.1731   0.1617   0.1835
## Detection Prevalence  0.2845   0.1931   0.1760   0.1627   0.1837
## Balanced Accuracy      0.9990   0.9976   0.9948   0.9925   0.9993
```

```
confusionMatrix(pred2, test$classe)
```

```
## Confusion Matrix and Statistics
##
##           Reference
## Prediction   A    B    C    D    E
##           A 1387   82    2    5    0
##           B    2  835   41    1   12
##           C    4   32  789   85   24
##           D    0    0   14  711   24
##           E    2    0    9    2  841
##
## Overall Statistics
##
##           Accuracy : 0.9305
##           95% CI : (0.923, 0.9374)
##           No Information Rate : 0.2845
##           P-Value [Acc > NIR] : < 2.2e-16
##
##           Kappa : 0.9119
##           McNemar's Test P-Value : < 2.2e-16
##
## Statistics by Class:
##
##           Class: A Class: B Class: C Class: D Class: E
## Sensitivity       0.9943  0.8799  0.9228  0.8843  0.9334
## Specificity       0.9746  0.9858  0.9642  0.9907  0.9968
## Pos Pred Value    0.9397  0.9371  0.8448  0.9493  0.9848
## Neg Pred Value     0.9977  0.9716  0.9834  0.9776  0.9852
## Prevalence        0.2845  0.1935  0.1743  0.1639  0.1837
## Detection Rate    0.2828  0.1703  0.1609  0.1450  0.1715
## Detection Prevalence 0.3010  0.1817  0.1905  0.1527  0.1741
## Balanced Accuracy  0.9845  0.9329  0.9435  0.9375  0.9651
```

Random Forest gives a better model.

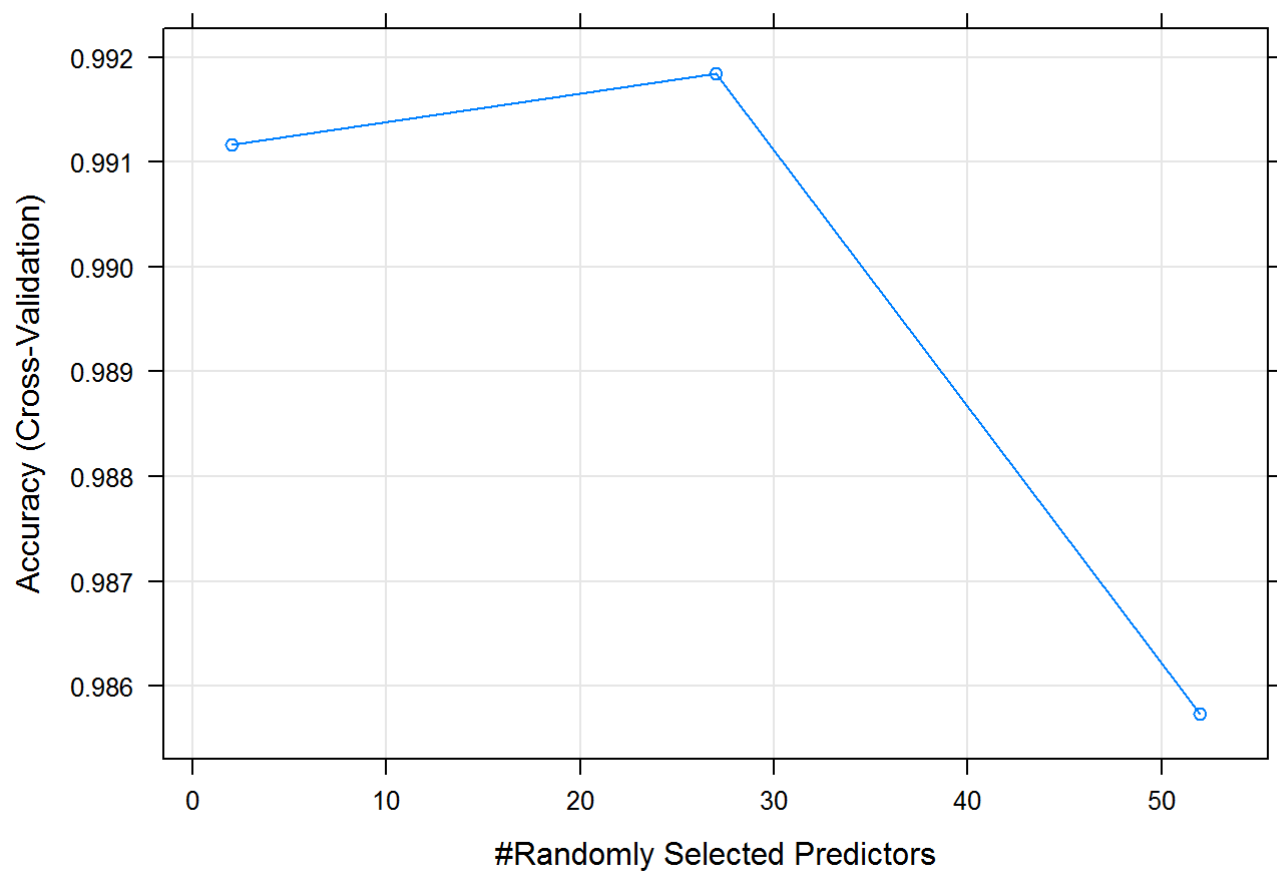
```
mod1
```

```
## Random Forest
##
## 14718 samples
##    52 predictor
##    5 classes: 'A', 'B', 'C', 'D', 'E'
##
## No pre-processing
## Resampling: Cross-Validated (7 fold)
## Summary of sample sizes: 12617, 12614, 12616, 12615, 12616, 12614, ...
## Resampling results across tuning parameters:
##
##  mtry  Accuracy   Kappa
##    2    0.9911670 0.9888258
##   27    0.9918461 0.9896846
##   52    0.9857302 0.9819469
##
## Accuracy was used to select the optimal model using  the largest value.
## The final value used for the model was mtry = 27.
```

```
varImp(mod1)
```

```
## rf variable importance
##
##    only 20 most important variables shown (out of 52)
##
##              Overall
## roll_belt      100.000
## pitch_forearm   59.037
## yaw_belt        53.013
## pitch_belt      44.696
## magnet_dumbbell_z 44.475
## magnet_dumbbell_y 44.469
## roll_forearm    40.914
## accel_dumbbell_y 21.899
## accel_forearm_x 17.873
## roll_dumbbell    17.708
## magnet_dumbbell_x 16.392
## accel_dumbbell_z 14.747
## accel_belt_z     14.183
## magnet_forearm_z 13.874
## magnet_belt_z    13.865
## total_accel_dumbbell 13.607
## magnet_belt_y    11.566
## gyros_belt_z     11.434
## yaw_arm          10.974
## magnet_belt_x     9.822
```

Including Plots



#Predicting Results on the Test Data using Random Forest since that had higher accuracy

```
rfPredictions <- predict(mod1, newdata = testing)
rfPredictions
```

```
## [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
```