

Human_Activity_Recognition_1_0

Moh A

11/04/2017

Contents

Executive Summary	1
Background	1
Exploratory Data Analysis	2
Data PreProcessing	2
Data Modeling	3
Cross Validation & Out of Sample Error Estimation	3
Models Analysis and Evaluation	3
Conclusion	4
Reference	4
Appendix	4

Executive Summary

Background

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, **the goal will be to use data from accelerometers on the belt, forearm, arm, and dumbbell of 6 participants to predict the manner in which they did the exercise.** They were asked to perform barbell lifts correctly and incorrectly in 5 different ways. The participants were asked to perform one set of 10 repetitions of the Unilateral Dumbbell Biceps Curl in five different fashions: exactly according to the specification (Class A), while the other 4 classes (B, C, D and E) correspond to common mistakes.

Read more: <http://groupware.les.inf.puc-rio.br/har#ixzz4d1GCDX2m>

Exploratory Data Analysis

In the dataset given, there is 19622 observations and 160 variable.

The dependent variable, 'Classe', is a categorical variable corresponding to the five different ways of doing the weight lifting exercise: exactly according to the specification (Class A), throwing the elbows to the front (Class B), lifting the dumbbell only halfway (Class C), lowering the dumbbell only halfway (Class D) and throwing the hips to the front (Class E).

The rest of the variable corresponds to measurements from date, time, belt, arm, forearm and dumbbell sensors.

```
library(caret)
library(ggplot2)
library(knitr)
library(VIM)
library(grid)
library(gridExtra)
```

	A	B	C	D	E
classe_table	5580.00000	3797.00000	3422.00000	3216.00000	3607.00000
classe_percent	28.43747	19.35073	17.43961	16.38977	18.38243

As seen in the table above, class 'A' (corresponding to correctly performing the exercise) is most frequent class and counts to roughly 28% among the five classes available.

It is also worth noticing that frequency of class 'A' and hence the quality of performing the exercise varies significantly among the six participant, as evident in the barplot below.

The aim of the classification model is to classify if exercise was completed correctly (class A) as well as classifying the type of error made (class: B, C, D and E). Since the task given was to predict the class of 20 unknown observations. A model with an accuracy higher than 95% is needed. For example, the probability of a model with 95% accuracy to predict all cases (20/20) correctly is $p^{20} = 0.95^{20} = 0.36$. In contrast, a model with 99% accuracy will achieve $0.99^{20} = 0.8$ probability in predicting the 20 cases correctly.

Data PreProcessing

The raw data file has 19622 observations and 160 variable.

Firstly, looking at the plot below for the proportion of missing data among the dataset variables, it can be noticed that it has many variables with mostly missing and meaningless values (100 out of the 160 variables). Thus, their contribution to the prediction of the 'classe' variable is insignificant and we have to omit these variables.

Secondly, it is wise to remove variables that do not contribute to the prediction of the 'classe' variable. Those include non-sensory measurements like: 'X'(row number), time and date related variables.

The rest of the variables have no missing values and are ready for further processing.

Data Modeling

Cross Validation & Out of Sample Error Estimation

K-fold cross validation is used to estimate the out of sample error, which is a reasonable balance between prediction accuracy and the computation time. The advantage of this method is that it matters less how the data gets divided. Every data point gets to be in a test set exactly once, and gets to be in a training set $k-1$ times. The variance of the resulting estimate is reduced as k is increased. The disadvantage of this method is that the training algorithm has to be rerun from scratch k times, which means it takes k times as much computation to make an evaluation. The results are then aggregated to create an overall estimate of the out of sample error.

Model 1: Linear Discriminant Analysis

This model was selected since it is a relatively simple classification algorithm to find a linear combination of features that characterizes or separates two or more classes of objects or events. The resulting combination may be used as a linear classifier.

Model 2: Random Forest

The random forest model, being an ensemble learning method for classification, uses multiple learning algorithms to obtain better predictive performance than could be obtained from any of the constituent learning algorithms alone. Random forests have a high predictive accuracy, automatically learns feature interactions and handles lots of irrelevant features well. However, results have a less interpretability than linear discriminant analysis and it has a slow training speed. Same as LDA model, k -fold cross validation with five folds are used.

Model 3: Neural Networks

Neural Networks is a powerful computational data model that is able to capture and represent complex input/output relationships. A neural network can be thought of as a network of “neurons” organised in layers. The predictors (or inputs) form the bottom layer, and the forecasts (or outputs) form the top layer. There may be intermediate layers containing “hidden neurons”. NNet have a high predictive accuracy, automatically learns feature interactions and handles lots of irrelevant features well. However, results are difficult to interpret and has a slow training speed. Same as LDA model, k -fold cross validation with five folds are used.

Models Analysis and Evaluation

All three models implemented achieved a very high predictive accuracy, when applied to our test set (30% of the dataset). The highest accuracy was 0.9986 by Random Forest model, whereas Linear Discriminant Analysis and Neural Networks achieved 0.9984 and 0.9983 respectively. Such slight differences in the accuracy score (difference between highest and lowest accuracy is approx. 0.0003) is very subtle and could be negligible when making a decision of which algorithm to choose.

As seen in the figure below, all models achieve optimal results with only 30 predictors out 54 predictors resulted after cleaning data. As shown in the figure below for the variable importance plot, the top 8 most important predictors for all three models are: num_window, roll_belt, pitch_forearm, yaw_belt, magnet_drumbell_z, magnet_drumbell_y, pitch_belt and roll_forearm.

Since the three models achieved a very high prediction accuracy when applied to the held out test set. The Out of Sample Error Rate is expected to be less than 0.2%

Finally, testing all three models on the validation set, all three models successfully classified the 20 cases of the validation set.

Conclusion

As explained in the previous section, all three models achieved a near perfect accuracy and predicted all the 20 cases in the validation data set. Given that the results presented in the summary table above look to be very close, it appears that a choice between either of the three algorithms is about the same. However, due to the properties of the dataset used in this report. It is sensible to choose Random Forest model over Linear Discriminant Analysis and Neural Networks. This is because it has a higher tolerance to missing values, irrelevant attributes, dealing with discrete/binary/continuous attributes and model parameter handling than the other two algorithm.

Reference

- Ugulino, W.; Cardador, D.; Vega, K.; Velloso, E.; Milidui, R.; Fuks, H. Wearable Computing: Accelerometers' Data Classification of Body Postures and Movements. Proceedings of 21st Brazilian Symposium on Artificial Intelligence. Advances in Artificial Intelligence - SBIA 2012. In: Lecture Notes in Computer Science. , pp. 52-61. Curitiba, PR: Springer Berlin / Heidelberg, 2012. ISBN 978-3-642-34458-9. DOI: 10.1007/978-3-642-34459-6_6.

Appendix

```
library(caret)
library(ggplot2)
library(knitr)
library(VIM)
library(grid)
library(gridExtra)
```

```
classe_table <- table(raw_data$classe)
classe_percent <- (100*classe_table)/length(raw_data$classe)
kable(rbind(classe_table, classe_percent))
```

	A	B	C	D	E
classe_table	5580.00000	3797.00000	3422.00000	3216.00000	3607.00000
classe_percent	28.43747	19.35073	17.43961	16.38977	18.38243

```
name_per_classe <-table(raw_data$classe, raw_data$user_name)

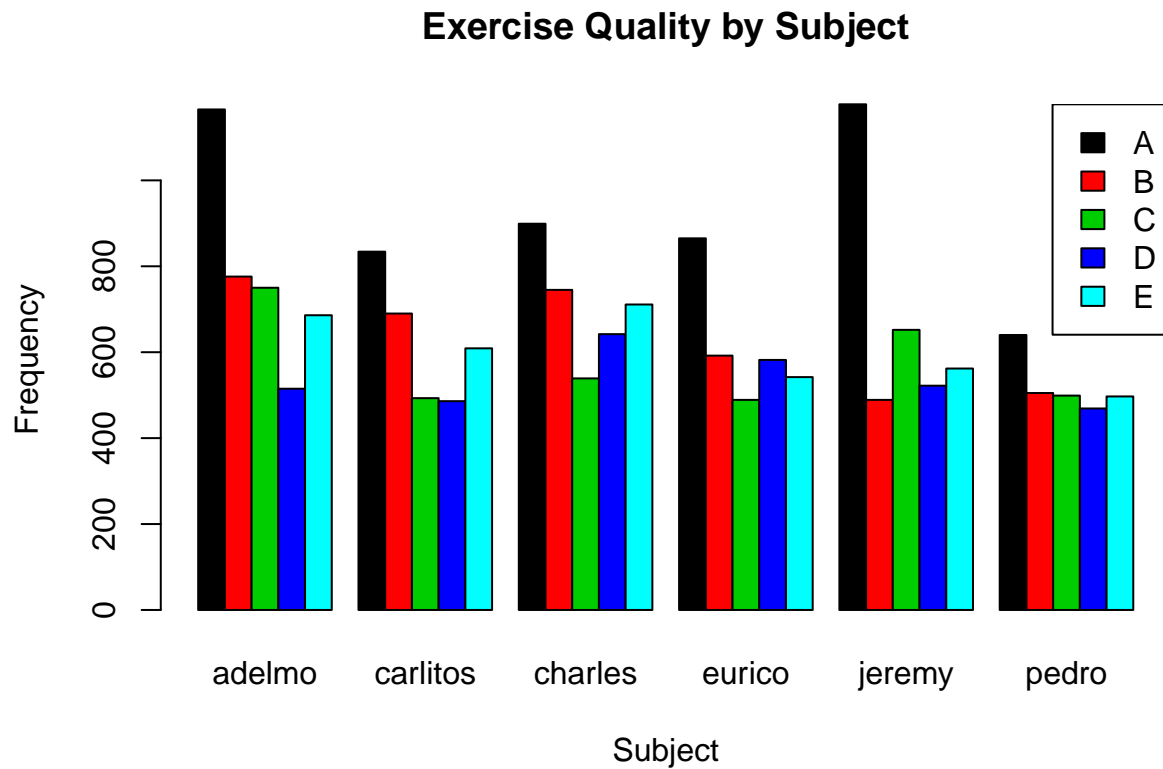
barplot(name_per_classe,
xlab = "Subject",
ylab = "Frequency",
main="Exercise Quality by Subject",
col = as.factor(rownames(name_per_classe)),
```

```

beside=TRUE
)

legend("topright", legend= rownames(name_per_classe), fill = as.factor(rownames(name_per_classe)))

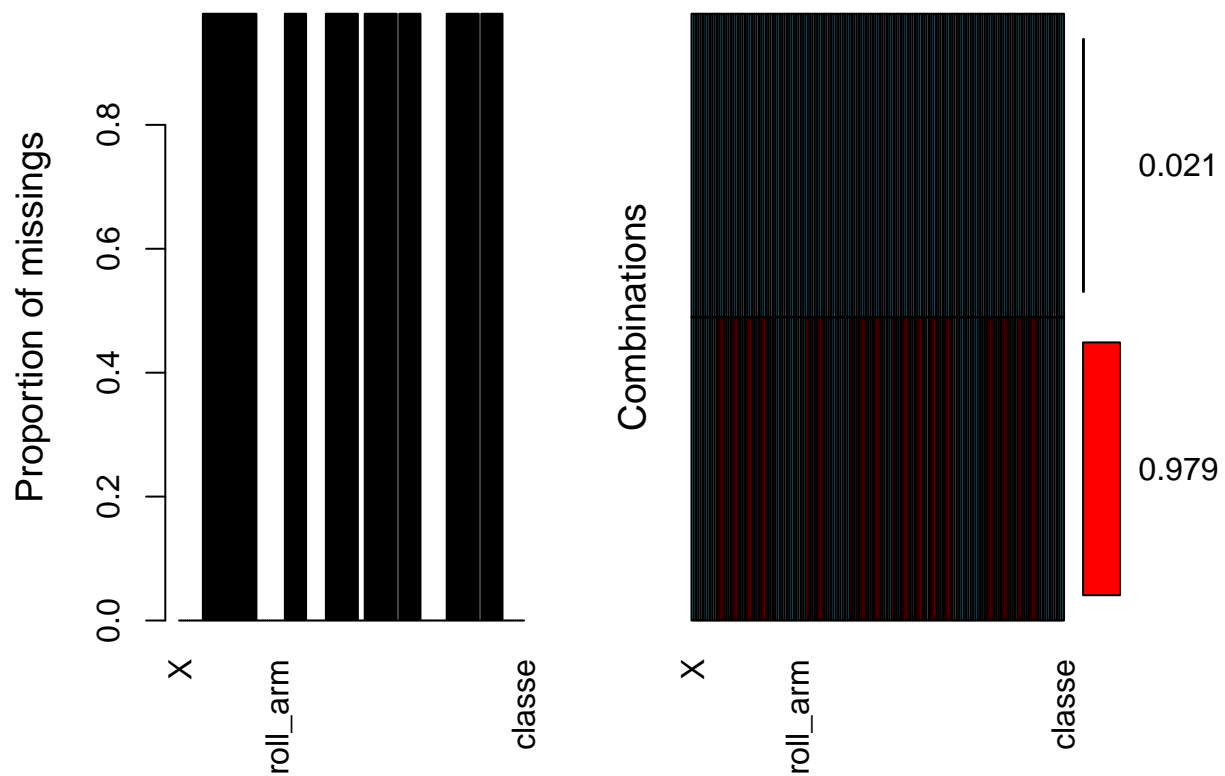
```



```

aggr(raw_data, prop=T, numbers=T)

```



```
majority_NA_columns <- which(sapply(raw_data, function(x) sum(is.na(x)) > (nrow(raw_data)*0.9)))
cleaned_data <- raw_data[, -majority_NA_columns]
cleaned_data <- cleaned_data[, -c(1,3,4,5,6)]
```

```
set.seed(55555)
inTrain <- createDataPartition(y=cleaned_data$classe,
p=0.7, list=FALSE)
training <- cleaned_data[inTrain,]
testing <- cleaned_data[-inTrain,]
```

```
library(parallel)
library(doParallel)
library(mlbench)
cluster <- makeCluster(detectCores())
registerDoParallel(cluster)
```

```
stopCluster(cluster)
registerDoSEQ()
```