

# Peer graded assignment Practical machine learning

*MvL*

*16 oktober 2017*

In this report data from accelerometers on the belt, forearm, arm, and dumbbell of 6 participants will be used to predict in a test set of the (in)correctness of the movements. The participants were asked to perform barbell lifts correctly and incorrectly in 5 different ways.

There are five classes indicating the performance of the exercise, Class A corresponds to the specified execution of the exercise, the others indicate common mistakes. The exercises were performed by six male participants aged between 20-28 years, with little weight lifting experience. The data for this project is kindly provided from this source: <http://web.archive.org/web/20161224072740/http://groupware.les.inf.puc-rio.br/har>

## Packages needed

```
#Load the packages needed to run this report.
library(caret)

## Warning: package 'caret' was built under R version 3.4.2
## Loading required package: lattice
## Loading required package: ggplot2
library(rpart)

## Warning: package 'rpart' was built under R version 3.4.2
library(rpart)
library(rpart.plot)

## Warning: package 'rpart.plot' was built under R version 3.4.2
```

## Data Exploration

The raw data is downloaded and processed for making a prediction on performance of weight lifting exercise on 6 participants.

```
# Download the dataset
trainUrl <- "https://d396qusza40orc.cloudfront.net/predmachlearn/pml-training.csv"
testUrl <- "https://d396qusza40orc.cloudfront.net/predmachlearn/pml-testing.csv"

training <- read.csv(url(trainUrl), na.strings = c("NA", "#DIV/0!", ""))
testing <- read.csv(url(testUrl), na.strings = c("NA", "#DIV/0!", ""))
```

There is a lot of missing data in the datasets, in the next steps the columns with missing data and the first columns with data that can not be included to predict on, is filtered out of the datasets. A validation set and a training set are created from the training data to train and validate the models.

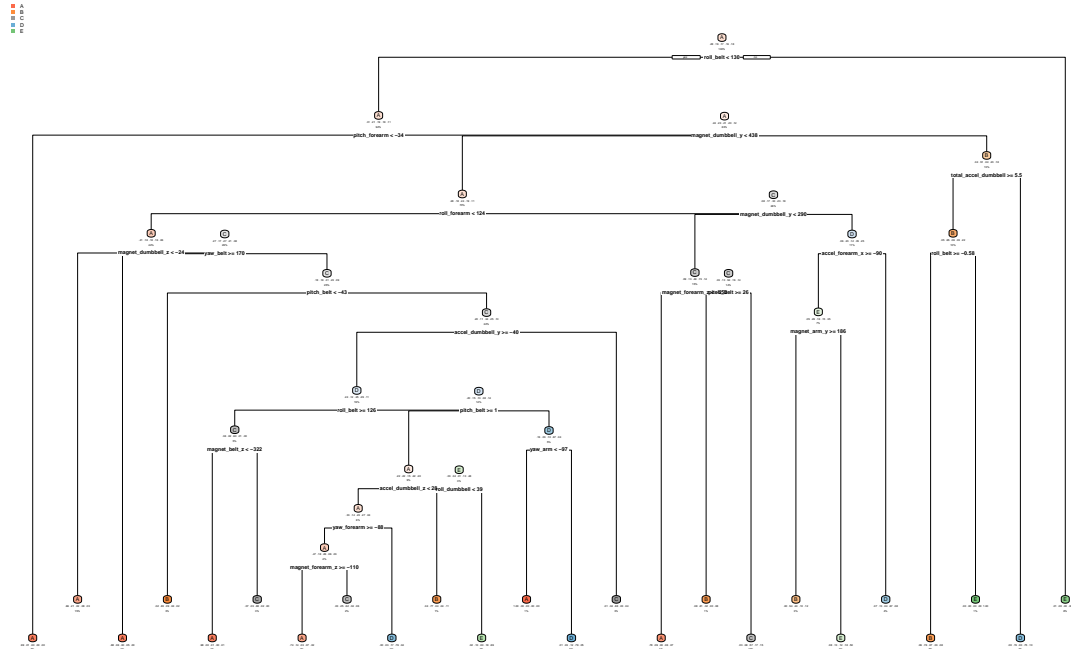
```
## [1] 13737    53
## [1] 5885     53
```

## Prediction models validation

Two different models will be tested to predict class in the weight lifting experiment. First a decision tree model is trained and validated, on the other hand a random forest model will be executed.

## Warning: labs do not fit even at cex 0.15, there may be some overplotting

### Decision Tree



### ## Confusion Matrix and Statistics

		Reference				
##	Prediction	A	B	C	D	E
##	A	1544	223	36	118	45
##	B	38	646	92	34	84
##	C	47	129	814	148	146
##	D	15	81	66	608	59
##	E	30	60	18	56	748

### ## Overall Statistics

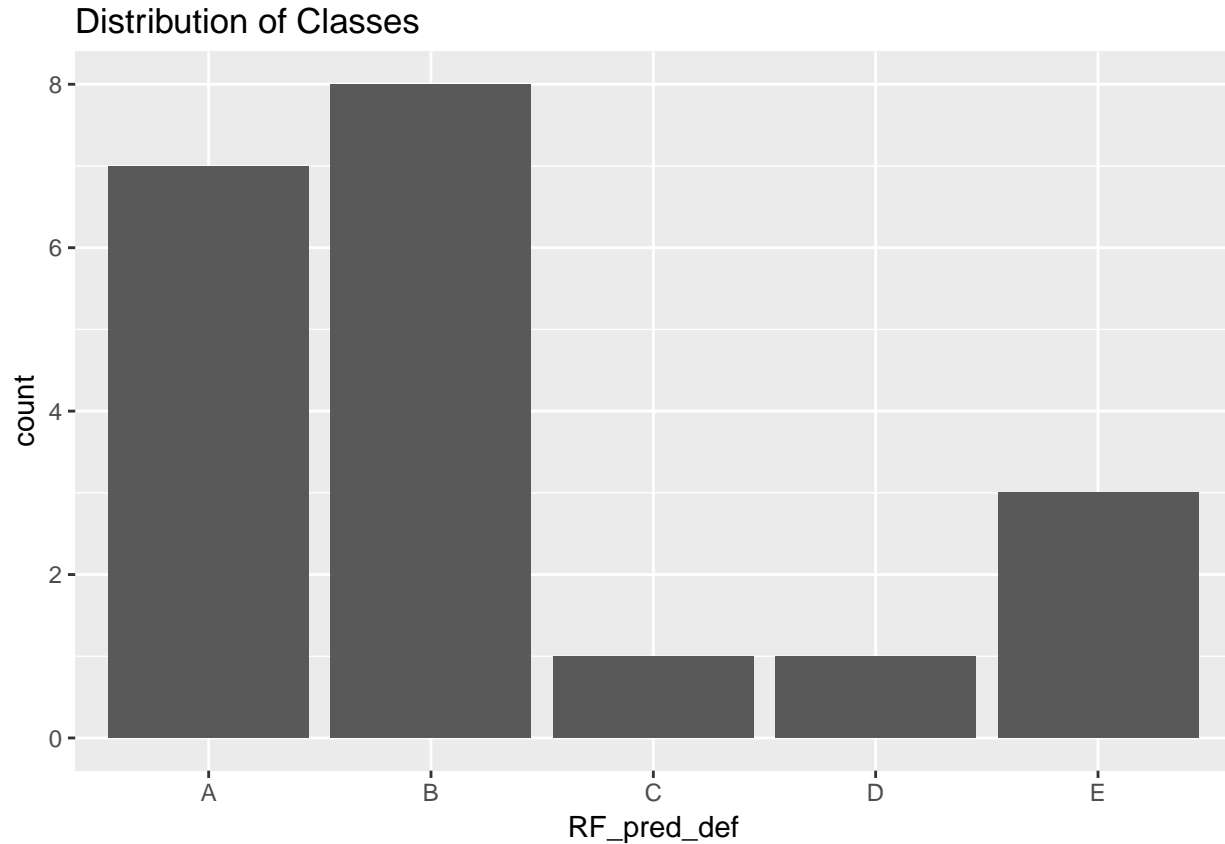
## Accuracy : 0.7409  
 ## 95% CI : (0.7295, 0.752)  
 ## No Information Rate : 0.2845  
 ## P-Value [Acc > NIR] : < 2.2e-16  
 ## Kappa : 0.6703  
 ## McNemar's Test P-Value : < 2.2e-16

```
## Statistics by Class:
##
##
##          Class: A Class: B Class: C Class: D Class: E
## Sensitivity      0.9223  0.5672  0.7934  0.6307  0.6913
## Specificity      0.8998  0.9477  0.9033  0.9551  0.9659
## Pos Pred Value   0.7854  0.7226  0.6340  0.7334  0.8202
## Neg Pred Value   0.9668  0.9012  0.9539  0.9296  0.9328
## Prevalence       0.2845  0.1935  0.1743  0.1638  0.1839
## Detection Rate   0.2624  0.1098  0.1383  0.1033  0.1271
## Detection Prevalence 0.3341  0.1519  0.2182  0.1409  0.1550
## Balanced Accuracy 0.9111  0.7575  0.8483  0.7929  0.8286

## Warning: package 'randomForest' was built under R version 3.4.2
## randomForest 4.6-12
## Type rfNews() to see new features/changes/bug fixes.
##
## Attaching package: 'randomForest'
## The following object is masked from 'package:ggplot2':
##
##     margin
## Confusion Matrix and Statistics
##
##          Reference
## Prediction    A    B    C    D    E
##          A 1668    6    1    0    0
##          B    4 1126   11    2    2
##          C    2    5 1011    8    1
##          D    0    2    3  954    6
##          E    0    0    0    0 1073
##
## Overall Statistics
##
##          Accuracy : 0.991
##          95% CI : (0.9882, 0.9932)
##          No Information Rate : 0.2845
##          P-Value [Acc > NIR] : < 2.2e-16
##
##          Kappa : 0.9886
##          McNemar's Test P-Value : NA
##
## Statistics by Class:
##
##          Class: A Class: B Class: C Class: D Class: E
## Sensitivity      0.9964  0.9886  0.9854  0.9896  0.9917
## Specificity      0.9983  0.9960  0.9967  0.9978  1.0000
## Pos Pred Value   0.9958  0.9834  0.9844  0.9886  1.0000
## Neg Pred Value   0.9986  0.9973  0.9969  0.9980  0.9981
## Prevalence       0.2845  0.1935  0.1743  0.1638  0.1839
## Detection Rate   0.2834  0.1913  0.1718  0.1621  0.1823
## Detection Prevalence 0.2846  0.1946  0.1745  0.1640  0.1823
## Balanced Accuracy 0.9974  0.9923  0.9910  0.9937  0.9958
```

## Model selection

The random forest model performance has a higher accuracy than the decision tree model, i.e. 0.99 vs 0.74, therefore the random forest model is selected as final model to predict the correct execution of a weight lifting exercise. A large part of the dataset was not usefull to predict performance and not included in the prediction. However even excluding this data gives a near perfect prediction on the validation set.



## Conclusion

The prediction shows that the predicted classes of the 20 test cases.

## Reference

Velloso, E.; Bulling, A.; Gellersen, H.; Ugulino, W.; Fuks, H. Qualitative Activity Recognition of Weight Lifting Exercises. Proceedings of 4th International Conference in Cooperation with SIGCHI (Augmented Human '13) . Stuttgart, Germany: ACM SIGCHI, 2013.