

Practical Machine Learning Assignment

Huzaifah Munshi

18/08/2020

Synopsis

The aim of this study is to predict the manner (“classe”) in which some healthy subjects performed a weight lifting exercise.

The subjects carried out the exercise in different fashions (some correct and some wrong). Their movements were monitored using devices equipped with accelerometers and stored in datasets that are available in the “WayBack Machine” website: <http://web.archive.org/web/20161224072740/http://groupware.les.inf.puc-rio.br/har>.

Data download and required package loading

```
library(AppliedPredictiveModeling)
library(caret)

## Loading required package: lattice
## Loading required package: ggplot2
library(dplyr)

##
## Attaching package: 'dplyr'

## The following objects are masked from 'package:stats':
##
##   filter, lag

## The following objects are masked from 'package:base':
##
##   intersect, setdiff, setequal, union

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:dplyr':
##
##   combine

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

```
urlTrain <- "https://d396qusza40orc.cloudfront.net/predmachlearn/pml-training.csv"
urlTest <- "https://d396qusza40orc.cloudfront.net/predmachlearn/pml-testing.csv"

download.file(urlTrain, destfile = "./pml-training.csv")
download.file(urlTest, destfile = "./pml-testing.csv")

training <- read.csv("pml-training.csv", na.strings=c("", "NA"))
testing <- read.csv("pml-testing.csv", na.strings=c("", "NA"))
unique(training$classe)
```

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

Data Exploratory Analysis

The str() and table() functions are used to understand the basic structure of the dataset. Due to the high number of columns (160), the result is subsetting:

```
ncol(training)
```

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

```
str(training[,1:10]) # first 10 columns. The first variables are not actual predictors
```

```
## 'data.frame': 19622 obs. of 10 variables:
## $ X : int 1 2 3 4 5 6 7 8 9 10 ...
## $ user_name : chr "carlitos" "carlitos" "carlitos" "carlitos" ...
## $ raw_timestamp_part_1: int 1323084231 1323084231 1323084231 1323084232 1323084232 1323084232 1323084232 1323084232 1323084232 1323084232 ...
## $ raw_timestamp_part_2: int 788290 808298 820366 120339 196328 304277 368296 440390 484323 484434 ...
## $ cvtd_timestamp : chr "05/12/2011 11:23" "05/12/2011 11:23" "05/12/2011 11:23" "05/12/2011 11:23" "05/12/2011 11:23" ...
## $ new_window : chr "no" "no" "no" "no" ...
## $ 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 ...
```

```
str(training[,149:160]) # last 12 columns. The outcome Classe appears at the end. Some columns appear to be NA
```

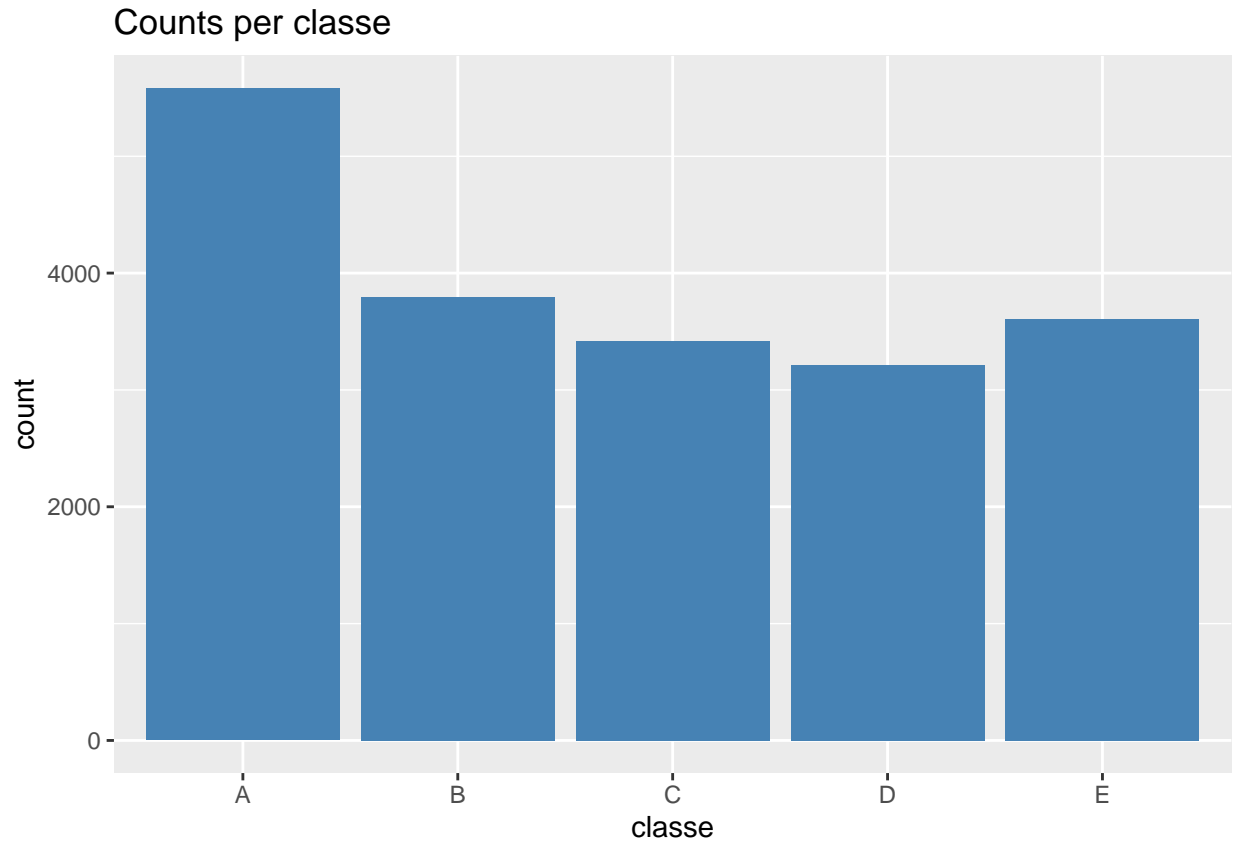
```
## 'data.frame': 19622 obs. of 12 variables:
## $ stddev_yaw_forearm: num NA NA NA NA NA NA NA NA NA NA NA ...
## $ var_yaw_forearm : num NA NA NA NA NA NA NA NA NA NA NA ...
## $ gyros_forearm_x : num 0.03 0.02 0.03 0.02 0.02 0.02 0.02 0.02 0.02 0.03 0.02 ...
## $ gyros_forearm_y : num 0 0 -0.02 -0.02 0 -0.02 0 -0.02 0 0 ...
## $ gyros_forearm_z : num -0.02 -0.02 0 0 -0.02 -0.03 -0.02 0 -0.02 -0.02 ...
## $ accel_forearm_x : int 192 192 196 189 189 193 195 193 193 190 ...
## $ accel_forearm_y : int 203 203 204 206 206 203 205 205 204 205 ...
## $ accel_forearm_z : int -215 -216 -213 -214 -214 -215 -215 -213 -214 -215 ...
## $ magnet_forearm_x : int -17 -18 -18 -16 -17 -9 -18 -9 -16 -22 ...
## $ magnet_forearm_y : num 654 661 658 658 655 660 659 660 653 656 ...
## $ magnet_forearm_z : num 476 473 469 469 473 478 470 474 476 473 ...
## $ classe : chr "A" "A" "A" "A" ...
```

```
table(training$classe,training$user_name) # number of observations per "user_name" and per "classe"
```

```
##
## adelmo carlitos charles eurico jeremy pedro
## A 1165 834 899 865 1177 640
## B 776 690 745 592 489 505
```

```
## C 750 493 539 489 652 499
## D 515 486 642 582 522 469
## E 686 609 711 542 562 497
```

```
ggplot(training, aes(classe)) + geom_bar(fill = "steelblue") + ggtitle("Counts per classe")
```



Data pre-processing

The outcome “classe” must be converted into a factor variable. Additionally, there are many columns which do not provide any relevant information, because they either have plenty of NAs or because they are not actual predictors obtained from accelerator measurements. Those columns will be removed:

```
training$classe <- as.factor(training$classe) # classe is converted into a factor variable.
```

```
trainingPrep <- training %>% select(8:160) # Non-predictors are removed.
```

```
trainingPrep <- trainingPrep %>% select_if(colSums(is.na(trainingPrep)) < 19000) # Only the columns with
```

```
ncol(trainingPrep) # The resulting amount of columns in the dataset is 53.
```

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

Create Data Partition

This dataset is further divided into train (75%) and test (25%) parts for cross-validation:

```
inTrain = createDataPartition(trainingPrep$classe, p = 3/4)[[1]]
trainPart = trainingPrep[ inTrain,]
```

```
testPart = trainingPrep[-inTrain,]
```

Model training

A couple of models will be trained and tested with cross validation to find out which of them has the highest accuracy level. More precisely, a random forest model and an LDA model will be tested:

```
set.seed(1234)
modfitrf <- randomForest(classe ~ ., method = "class", data = trainPart)
predrf <- predict(modfitrf, newdata = testPart, type = "class")
confusionMatrix(predrf, testPart$classe)
```

```
## Confusion Matrix and Statistics
##
##              Reference
## Prediction    A    B    C    D    E
##      A 1394     9     0     0     0
##      B     1   940     1     0     0
##      C     0     0   854     7     2
##      D     0     0     0   797     1
##      E     0     0     0     0   898
##
## Overall Statistics
##
##              Accuracy : 0.9957
##              95% CI : (0.9935, 0.9973)
##      No Information Rate : 0.2845
##      P-Value [Acc > NIR] : < 2.2e-16
##
##              Kappa : 0.9946
##
##  Mcnemar's Test P-Value : NA
##
## Statistics by Class:
##
##              Class: A Class: B Class: C Class: D Class: E
## Sensitivity          0.9993  0.9905  0.9988  0.9913  0.9967
## Specificity          0.9974  0.9995  0.9978  0.9998  1.0000
## Pos Pred Value       0.9936  0.9979  0.9896  0.9987  1.0000
## Neg Pred Value       0.9997  0.9977  0.9998  0.9983  0.9993
## Prevalence           0.2845  0.1935  0.1743  0.1639  0.1837
## Detection Rate       0.2843  0.1917  0.1741  0.1625  0.1831
## Detection Prevalence 0.2861  0.1921  0.1760  0.1627  0.1831
## Balanced Accuracy    0.9984  0.9950  0.9983  0.9955  0.9983
```

```
set.seed(1234)
modfitlda <- train(classe ~ ., method = "lda", data = trainPart)
predlda <- predict(modfitlda, newdata = testPart)
confusionMatrix(predlda, testPart$classe)
```

```
## Confusion Matrix and Statistics
##
##              Reference
## Prediction    A    B    C    D    E
##      A 1117   148    94    60   36
```

```
##           B    28  616   92   28  156
##           C   125  105  555   95   85
##           D   116   37   96  587   81
##           E    9   43   18   34  543
##
## Overall Statistics
##
##           Accuracy : 0.697
##           95% CI : (0.6839, 0.7098)
##           No Information Rate : 0.2845
##           P-Value [Acc > NIR] : < 2.2e-16
##
##           Kappa : 0.6165
##
## McNemar's Test P-Value : < 2.2e-16
##
## Statistics by Class:
##
##           Class: A Class: B Class: C Class: D Class: E
## Sensitivity      0.8007  0.6491  0.6491  0.7301  0.6027
## Specificity      0.9037  0.9231  0.8987  0.9195  0.9740
## Pos Pred Value   0.7677  0.6696  0.5751  0.6401  0.8393
## Neg Pred Value   0.9194  0.9164  0.9238  0.9456  0.9159
## Prevalence       0.2845  0.1935  0.1743  0.1639  0.1837
## Detection Rate   0.2278  0.1256  0.1132  0.1197  0.1107
## Detection Prevalence 0.2967  0.1876  0.1968  0.1870  0.1319
## Balanced Accuracy 0.8522  0.7861  0.7739  0.8248  0.7883
```

Model selection

The accuracy level of the random forest model (higher than 99%) is clearly higher than that of the LDA model (close to 70%). Therefore, the random forest model is selected.

Cross validation and expected out of sample error

The out of sample error (calculated as 1 - Accuracy Level) is below 1%, therefore very low.

Prediction on 20 test cases

```
predrf20 <- predict(modfitrf, newdata = testing, type = "class")
print(predrf20)
```

```
##  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
```

This is the prediction achieved with the selected model (random forest) for the 20 test cases.