

Practical Machine Learning

Mpho Godfrey Nkadimeng

17/10/2019

Practical Machine Learning

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, 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. More information is available from the website here: <http://groupware.les.inf.puc-rio.br/har> (<http://groupware.les.inf.puc-rio.br/har>) (see the section on the Weight Lifting Exercise Dataset).

Data

- The training data for this project are available here:
[<https://d396qusza40orc.cloudfront.net/predmachlearn/pml-training.csv>
(<https://d396qusza40orc.cloudfront.net/predmachlearn/pml-training.csv>)]
- The test data are available here:
[<https://d396qusza40orc.cloudfront.net/predmachlearn/pml-testing.csv>
(<https://d396qusza40orc.cloudfront.net/predmachlearn/pml-testing.csv>)]
- The data for this project come from this source: [<http://groupware.les.inf.puc-rio.br/har> (<http://groupware.les.inf.puc-rio.br/har>)]. If you use the document you create for this class for any purpose please cite them as they have been very generous in allowing their data to be used for this kind of assignment.

What you should submit

The goal of your project is to predict the manner in which they did the exercise. This is the “classe” variable in the training set. You may use any of the other variables to predict with. You should create a report describing how you built your model, how you used cross validation, what you think the expected out of sample error is, and why you made the choices you did. You will also use your prediction model to predict 20 different test cases.

Your submission should consist of a link to a Github repo with your R markdown and compiled HTML file describing your analysis. Please constrain the text of the writeup to < 2000 words and the number of figures to be less than 5. It will make it easier for the graders if you submit a repo with a gh-pages branch so the HTML page can be viewed online (and you always want to make it easy on graders :-). You should also apply your machine learning algorithm to the 20 test cases available in the test data above. Please submit your predictions in appropriate format to the programming assignment for automated grading. See the programming assignment for additional details.

Preliminary Work

Reproduceability

An overall pseudo-random number generator seed was set at 1234 for all code. In order to reproduce the results below, the same seed should be used. Different packages were downloaded and installed, such as caret and randomForest. These should also be installed in order to reproduce the results below (please see code below for ways and syntax to do so).

How the model was built

Our outcome variable is classe, a factor variable with 5 levels. For this data set, “participants were asked to perform one set of 10 repetitions of the Unilateral Dumbbell Biceps Curl in 5 different fashions:

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)
- throwing the hips to the front (Class E)?

Class A corresponds to the specified execution of the exercise, while the other 4 classes correspond to common mistakes." [1] Prediction evaluations will be based on maximizing the accuracy and minimizing the out-of-sample error. All other available variables after cleaning will be used for prediction. Two models will be tested using decision tree and random forest algorithms. The model with the highest accuracy will be chosen as our final model.

Cross-validation

Cross-validation will be performed by subsampling our training data set randomly without replacement into 2 subsamples: subTraining data (75% of the original Training data set) and subTesting data (25%). Our models will be fitted on the subTraining data set, and tested on the subTesting data. Once the most accurate model is chosen, it will be tested on the original Testing data set.

Expected out-of-sample error

The expected out-of-sample error will correspond to the quantity: 1-accuracy in the cross-validation data. Accuracy is the proportion of correct classified observation over the total sample in the subTesting data set. Expected accuracy is the expected accuracy in the out-of-sample data set (i.e. original testing data set). Thus, the expected value of the out-of-sample error will correspond to the expected number of missclassified observations/total observations in the Test data set, which is the quantity: 1-accuracy found from the cross-validation data set.

Our outcome variable “classe” is an unordered factor variable. Thus, we can choose our error type as 1-accuracy. We have a large sample size with $N = 19622$ in the Training data set. This allows us to divide our Training sample into subTraining and subTesting to allow cross-validation. Features with all missing values will be discarded as well as features that are irrelevant. All other features will be kept as relevant variables. Decision tree and random forest algorithms are known for their ability of detecting the features that are important for classification [2].

Packages, Libraries and Seed

Installing packages, loading libraries, and setting the seed for reproducibility:

```
library(caret)
```

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

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

```
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:ggplot2':  
##  
##      margin
```

```
library(rpart)
```

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

```
library(RColorBrewer)
```

```
library(rattle)
```

```
## Rattle: A free graphical interface for data science with R.  
## Version 5.2.0 Copyright (c) 2006-2018 Togaware Pty Ltd.  
## Type 'rattle()' to shake, rattle, and roll your data.
```

```
##  
## Attaching package: 'rattle'
```

```
## The following object is masked from 'package:randomForest':  
##  
##      importance
```

```
set.seed(1000)
```

Getting and cleaning data

The training data set can be found on the following URL:

```
trainUrl <- "http://d396qusza40orc.cloudfront.net/predmachlearn/pm
ml-training.csv"
```

```
testUrl <- "http://d396qusza40orc.cloudfront.net/predmachlearn/pm
ml-testing.csv"
```

Load data to memory.

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

Partitioning the training set into two

```
inTrain <- createDataPartition(y=training$classe, p=0.6, list=FAL
SE)
myTraining <- training[inTrain, ]; myTesting <- training[-inTrain
, ]
dim(myTraining); dim(myTesting)
```

```
## [1] 11776    160
```

```
## [1] 7846     160
```

Cleaning the Data

The following transformations were used to clean the data:

Transformation 1: Cleaning NearZeroVariance Variables Run this code to view possible NZV Variables:

```
myDataNZV <- nearZeroVar(myTraining, saveMetrics=TRUE)
```

Another subset of NZV variables.

```

myNZVvars <- names(myTraining) %in% c("new_window", "kurtosis_roll_belt", "kurtosis_pitch_belt", "kurtosis_yaw_belt", "skewness_roll_belt", "skewness_roll_belt.1", "skewness_yaw_belt", "max_yaw_belt", "min_yaw_belt", "amplitude_yaw_belt", "avg_roll_arm", "stddev_roll_arm", "var_roll_arm", "avg_pitch_arm", "stddev_pitch_arm", "var_pitch_arm", "avg_yaw_arm", "stddev_yaw_arm", "var_yaw_arm", "kurtosis_roll_arm", "kurtosis_pitch_arm", "kurtosis_yaw_arm", "skewness_roll_arm", "skewness_pitch_arm", "skewness_yaw_arm", "max_roll_arm", "min_roll_arm", "min_pitch_arm", "amplitude_roll_arm", "amplitude_pitch_arm", "kurtosis_roll_dumbbell", "kurtosis_pitch_dumbbell", "kurtosis_yaw_dumbbell", "skewness_roll_dumbbell", "skewness_pitch_dumbbell", "skewness_yaw_dumbbell", "max_yaw_dumbbell", "min_yaw_dumbbell", "amplitude_yaw_dumbbell", "kurtosis_roll_forearm", "kurtosis_pitch_forearm", "kurtosis_yaw_forearm", "skewness_roll_forearm", "skewness_pitch_forearm", "skewness_yaw_forearm", "max_roll_forearm", "max_yaw_forearm", "min_roll_forearm", "min_yaw_forearm", "amplitude_roll_forearm", "amplitude_yaw_forearm", "avg_roll_forearm", "stddev_roll_forearm", "var_roll_forearm", "avg_pitch_forearm", "stddev_pitch_forearm", "var_pitch_forearm", "avg_yaw_forearm", "stddev_yaw_forearm", "var_yaw_forearm")
myTraining <- myTraining[!myNZVvars]

dim(myTraining)

```

```
## [1] 11776    100
```

Transformation 2: Killing first column of Dataset - ID Removing first ID variable so that it does not interfere with ML Algorithms:

```
myTraining <- myTraining[c(-1)]
```

Transformation 3: Cleaning Variables with too many NAs. For Variables that have more than a 60% threshold of NA's I'm going to leave them out:

```

trainingV3 <- myTraining #creating another subset to iterate in loop
for(i in 1:length(myTraining)) { #for every column in the training dataset
  if( sum( is.na( myTraining[, i] ) ) /nrow(myTraining) >=
.6 ) { #if n?? NAs > 60% of total observations
    for(j in 1:length(trainingV3)) {
      if( length( grep(names(myTraining[i]), names(trainingV3)[j])) ==1) { #if the columns are the same:
        trainingV3 <- trainingV3[ , -j] #Remove that column
      }
    }
  }
}
#To check the new N?? of observations
dim(trainingV3)

```

```
## [1] 11776    58
```

```

#Setting back to our set:
myTraining <- trainingV3
rm(trainingV3)

```

Now let us do the exact same 3 transformations for myTesting and testing data sets.

```

clean1 <- colnames(myTraining)
clean2 <- colnames(myTraining[, -58]) #already with classe column removed
myTesting <- myTesting[clean1]
testing <- testing[clean2]

#To check the new N?? of observations
dim(myTesting)

```

```
## [1] 7846    58
```

```

#To check the new N?? of observations
dim(testing)

```

```
## [1] 20 57
```

In order to ensure proper functioning of Decision Trees and especially RandomForest Algorithm with the Test data set (data set provided), we need to coerce the data into the same type.

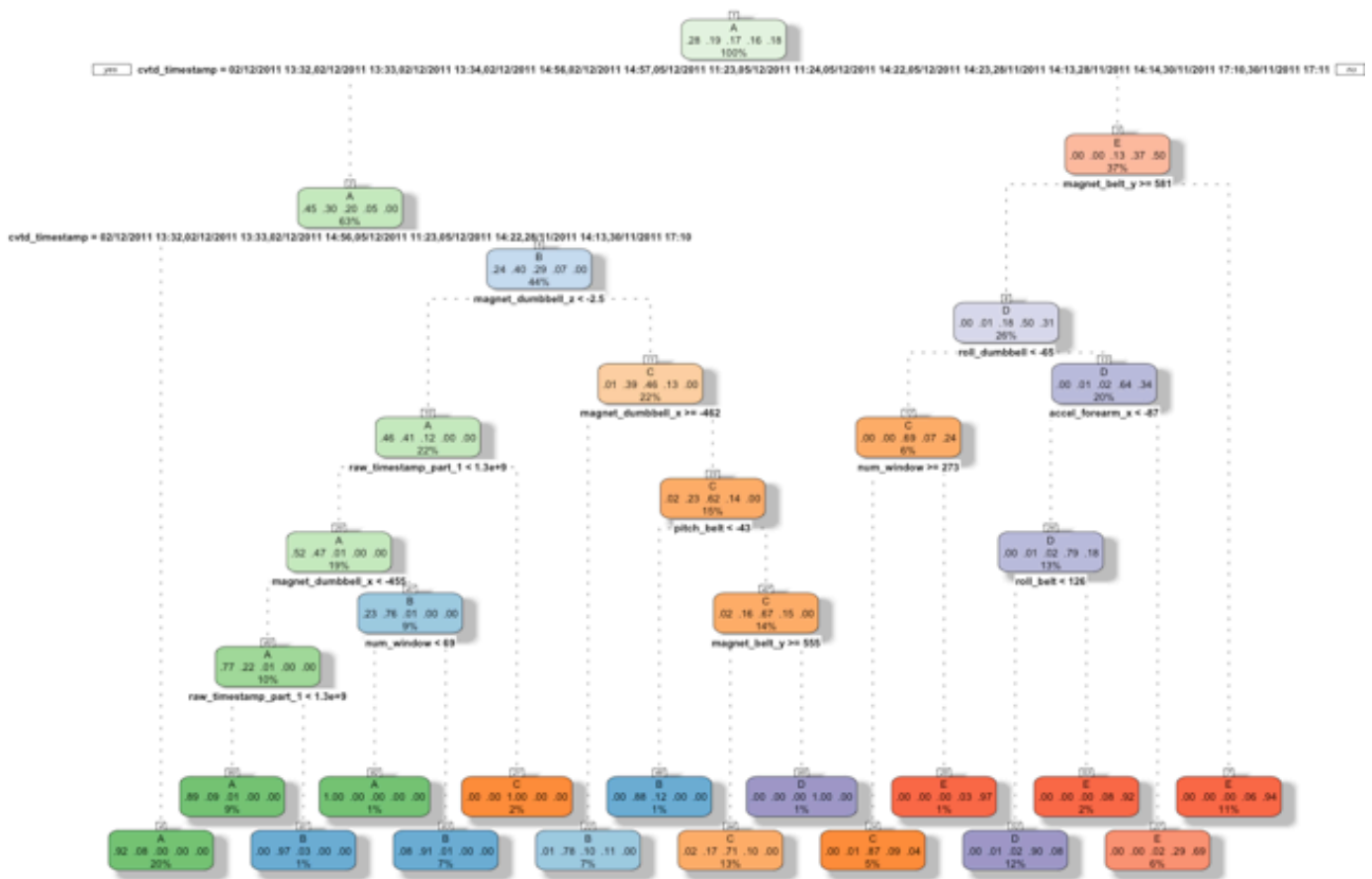
```
for (i in 1:length(testing) ) {  
  for(j in 1:length(myTraining)) {  
    if( length( grep(names(myTraining[i]), names(testing)[j])  
  ) ==1) {  
      class(testing[j]) <- class(myTraining[i])  
    }  
  }  
}  
#And to make sure Coertion really worked, simple smart ass techni  
que:  
testing <- rbind(myTraining[2, -58] , testing) #note row 2 does n  
ot mean anything, this will be removed right.. now:  
testing <- testing[-1,]
```

Using ML algorithms for prediction: Decision Tree

```
modFitA1 <- rpart(classe ~ ., data=myTraining, method="class")
```

To view the decision tree with fancy :

```
fancyRpartPlot(modFitA1)
```

Rattle 2019-Oct-17 19:40:58 sibongiledlamini

Predictiong

```
predictionsA1 <- predict(modFitA1, myTesting, type = "class")
```

Using confusion Matrix to test results:

```
confusionMatrix(predictionsA1, myTesting$classe)
```

Confusion Matrix and Statistics

##

| | | Reference | | | | |
|------------|------|-----------|------|-----|------|---|
| Prediction | | A | B | C | D | E |
| A | 2152 | 212 | 5 | 4 | 0 | |
| B | 56 | 1118 | 97 | 67 | 0 | |
| C | 24 | 180 | 1251 | 137 | 22 | |
| D | 0 | 8 | 6 | 834 | 81 | |
| E | 0 | 0 | 9 | 244 | 1339 | |

##

Overall Statistics

##

Accuracy : 0.8532
95% CI : (0.8451, 0.8609)
No Information Rate : 0.2845
P-Value [Acc > NIR] : < 2.2e-16

##

Kappa : 0.8138

##

McNemar's Test P-Value : NA

##

Statistics by Class:

##

| | | Class: A | Class: B | Class: C | Class: D | Class |
|----------------------|-----|----------|----------|----------|----------|-------|
| : E | | | | | | |
| Sensitivity | 286 | 0.9642 | 0.7365 | 0.9145 | 0.6485 | 0.9 |
| Specificity | 605 | 0.9606 | 0.9652 | 0.9440 | 0.9855 | 0.9 |
| Pos Pred Value | 411 | 0.9069 | 0.8356 | 0.7751 | 0.8977 | 0.8 |
| Neg Pred Value | 835 | 0.9854 | 0.9385 | 0.9812 | 0.9347 | 0.9 |
| Prevalence | 838 | 0.2845 | 0.1935 | 0.1744 | 0.1639 | 0.1 |
| Detection Rate | 707 | 0.2743 | 0.1425 | 0.1594 | 0.1063 | 0.1 |
| Detection Prevalence | 029 | 0.3024 | 0.1705 | 0.2057 | 0.1184 | 0.2 |
| Balanced Accuracy | 445 | 0.9624 | 0.8509 | 0.9292 | 0.8170 | 0.9 |

Using ML algorithms for prediction: Random Forests

```
modFitB1 <- randomForest(classe ~. , data=myTraining)
```

Predicting in-sample error:

```
predictionsB1 <- predict(modFitB1, myTesting, type = "class")
```

Using confusion Matrix to test results:

```
confusionMatrix(predictionsB1, myTesting$classe)
```

```
## Confusion Matrix and Statistics
##
##           Reference
## Prediction      A      B      C      D      E
##           A 2232      1      0      0      0
##           B   0 1517      2      0      0
##           C   0      0 1366      2      0
##           D   0      0      0 1284      3
##           E   0      0      0      0 1439
##
## Overall Statistics
##
##           Accuracy : 0.999
##           95% CI : (0.998, 0.9996)
##           No Information Rate : 0.2845
##           P-Value [Acc > NIR] : < 2.2e-16
##
##           Kappa : 0.9987
##
##           McNemar's Test P-Value : NA
##
## Statistics by Class:
##
##           Class: A Class: B Class: C Class: D Class
: E
## Sensitivity           1.0000    0.9993    0.9985    0.9984    0.9
979
## Specificity           0.9998    0.9997    0.9997    0.9995    1.0
000
## Pos Pred Value           0.9996    0.9987    0.9985    0.9977    1.0
000
## Neg Pred Value           1.0000    0.9998    0.9997    0.9997    0.9
995
## Prevalence           0.2845    0.1935    0.1744    0.1639    0.1
838
## Detection Rate           0.2845    0.1933    0.1741    0.1637    0.1
834
## Detection Prevalence    0.2846    0.1936    0.1744    0.1640    0.1
834
## Balanced Accuracy           0.9999    0.9995    0.9991    0.9990    0.9
990
```

Random Forests yielded better Results.

Generating Files to submit as answers for the Assignment:

Finally, using the provided Test Set out-of-sample error.

For Random Forests we use the following formula, which yielded a much better prediction in in-sample:

```
predictionsB2 <- predict(modFitB1, testing, type = "class")
```

Function to generate files with predictions to submit for assignment

```
pml_write_files = function(x){  
  n = length(x)  
  for(i in 1:n){  
    filename = paste0("problem_id_",i,".txt")  
    write.table(x[i],file=filename,quote=FALSE,row.names=FALSE,col.names=FALSE)  
  }  
}  
  
pml_write_files(predictionsB2)
```