

Machine learning course project

Yuva

September 21, 2016

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>).

Preliminary work

Reproducibility

In order to reproduce the result, I set seed at 1234. Different R packages such as caret and randomForest were installed.

Model Building

Our outcome variable is class, 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. 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.

Reasons for my choices

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. Feature selection is inherent, so it is not so necessary at the data preparation phase. Thus, there won't be any feature selection section in this report.

Results

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

```
library(caret)
```

```
## Warning: package 'caret' was built under R version 3.2.5
```

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

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

```
## Warning: package 'ggplot2' was built under R version 3.2.5
```

```
library(rpart)  
library(randomForest)
```

```
## Warning: package 'randomForest' was built under R version 3.2.5
```

```
## 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
```

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

```
## Warning: package 'rpart.plot' was built under R version 3.2.5
```

```
set.seed(1234)
```

loading the data and preliminary cleaning

```
#Loading the data and changing missing values to "NA"
```

```
trainingset <- read.csv("C:\\Users\\Subodha\\Rprogramming\\Machinelearning\\pml-training.csv", na.  
a.strings=c("NA", "#DIV/0!", ""))
```

```
testingset <- read.csv("C:\\Users\\Subodha\\Rprogramming\\Machinelearning\\pml-testing.csv", na.s  
trings = c("NA", "#DIV/0!", ""))
```

```
# checking dimensions of the data set
```

```
dim(trainingset)
```

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

```
dim(testingset)
```

```
## [1] 20 160
```

```
# Lets delete column with all missing value
trainingset<-trainingset[,colSums(is.na(trainingset)) == 0]
testingset <-testingset[,colSums(is.na(testingset)) == 0]

#lets delete some irrrelevent columns from the dataset such as user_name, raw_timestamp_part_1,
raw_timestamp_part_2 cvtd_timestamp, new_window, and num_window (columns 1 to 7)

trainingset <-trainingset[,-c(1:7)]
testingset <-testingset[,-c(1:7)]

#lets look at the new dimensions of the cleaned dataset
dim(trainingset)
```

```
## [1] 19622 53
```

```
dim(testingset)
```

```
## [1] 20 53
```

```
#lets look at the classes of the dataset( the details are hidden in the output file as it takes
lots of space)
head(trainingset)
head(testingset)
```

partitioning training data set for cross-validation

```
subsamples <- createDataPartition(y=trainingset$classe, p=0.75, list=FALSE)
subTraining <- trainingset[subsamples, ]
subTesting <- trainingset[-subsamples, ]
dim(subTraining)
```

```
## [1] 14718 53
```

```
dim(subTesting)
```

```
## [1] 4904 53
```

First prediction model: Using Decision Tree

```

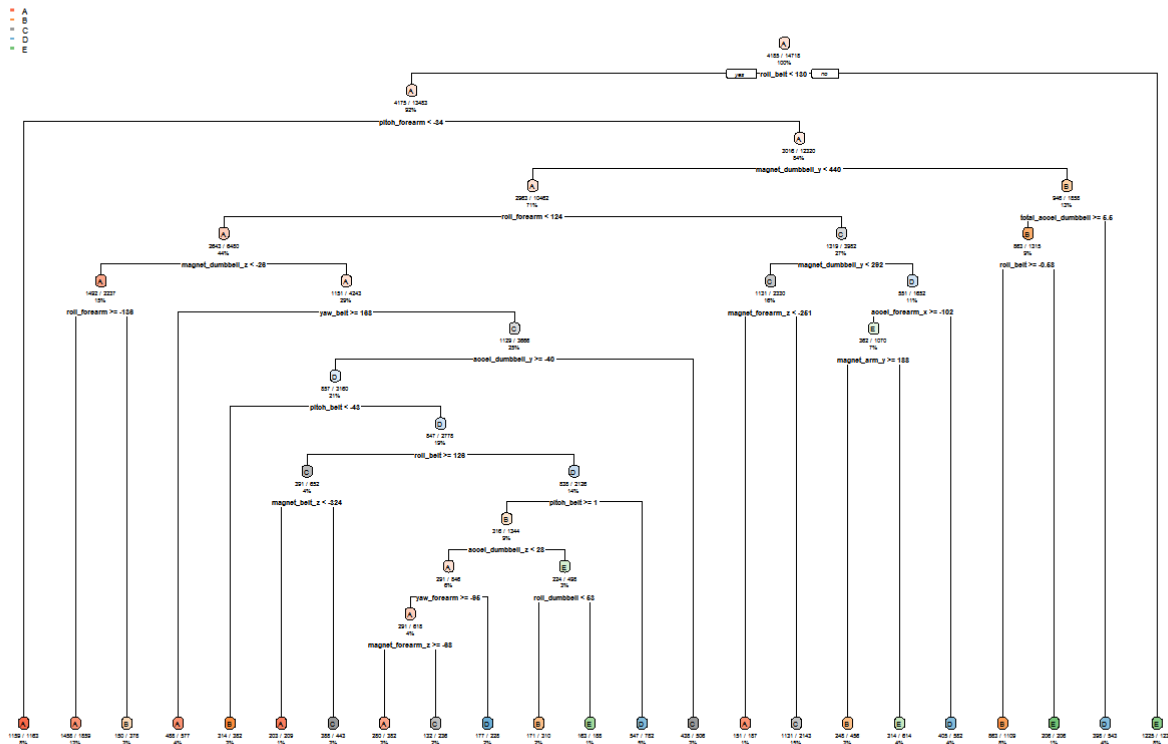
modell1 <- rpart(classe ~ ., data=subTraining, method="class")

prediction1 <- predict(modell1, subTesting, type = "class")

#plot of decesion tree
rpart.plot(modell1, main="Classification Tree", extra=102, under=TRUE, faclen=0)

```

Classification Tree



```

# Test results on our subTesting data set:
confusionMatrix(prediction1, subTesting$classe)

```

```
## Confusion Matrix and Statistics
##
##           Reference
## Prediction   A    B    C    D    E
##           A 1235  157   16   50   20
##           B   55  568   73   80  102
##           C   44  125  690  118  116
##           D   41   64   50  508   38
##           E   20   35   26   48  625
##
## Overall Statistics
##
##           Accuracy : 0.7394
##           95% CI : (0.7269, 0.7516)
##           No Information Rate : 0.2845
##           P-Value [Acc > NIR] : < 2.2e-16
##
##           Kappa : 0.6697
##           McNemar's Test P-Value : < 2.2e-16
##
## Statistics by Class:
##
##           Class: A Class: B Class: C Class: D Class: E
## Sensitivity           0.8853   0.5985   0.8070   0.6318   0.6937
## Specificity           0.9307   0.9216   0.9005   0.9529   0.9678
## Pos Pred Value        0.8356   0.6469   0.6313   0.7247   0.8289
## Neg Pred Value        0.9533   0.9054   0.9567   0.9296   0.9335
## Prevalence            0.2845   0.1935   0.1743   0.1639   0.1837
## Detection Rate        0.2518   0.1158   0.1407   0.1036   0.1274
## Detection Prevalence  0.3014   0.1790   0.2229   0.1429   0.1538
## Balanced Accuracy      0.9080   0.7601   0.8537   0.7924   0.8307
```

Second prediction model: Using Random Forest

```
model2 <- randomForest(classe ~. , data=subTraining, method="class")

# Predicting:

prediction2 <- predict(model2, subTesting, type = "class")

# Test results on subTesting data set:

confusionMatrix(prediction2, subTesting$classe)
```

```
## Confusion Matrix and Statistics
##
##           Reference
## Prediction   A    B    C    D    E
##           A 1394    3    0    0    0
##           B    1  944   10    0    0
##           C    0    2  843    6    0
##           D    0    0    2  798    0
##           E    0    0    0    0  901
##
## 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           0.9993   0.9947   0.9860   0.9925   1.0000
## Specificity           0.9991   0.9972   0.9980   0.9995   1.0000
## Pos Pred Value        0.9979   0.9885   0.9906   0.9975   1.0000
## Neg Pred Value        0.9997   0.9987   0.9970   0.9985   1.0000
## Prevalence            0.2845   0.1935   0.1743   0.1639   0.1837
## Detection Rate        0.2843   0.1925   0.1719   0.1627   0.1837
## Detection Prevalence  0.2849   0.1947   0.1735   0.1631   0.1837
## Balanced Accuracy      0.9992   0.9960   0.9920   0.9960   1.0000
```

Decision

Accuracy for Random Forest model was 0.995 (95% CI: (0.993, 0.997)) compared to 0.739 (95% CI: (0.727, 0.752)) for Decision Tree model. The expected out-of-sample error is estimated at 0.005, or 0.5%. The expected out-of-sample error is calculated as 1 - accuracy for predictions made against the cross-validation set. Our Test data set comprises 20 cases. With an accuracy above 99% on our cross-validation data, we can expect that very few, or none, of the test samples will be missclassified.

Generating Files to be submitted as answers for the Assignment

```
# predict outcome levels on the original Testing data set using Random Forest algorithm
predictfinal <- predict(model2, testingset, type="class")
predictfinal
```

```
## 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
```

```
# Write files for submission
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(predictfinal)
```