

Machine Learning Course Project Writeup

Siyuan Cheng

June 11, 2017

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> (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>]

The test data are available here: [<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>]. 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.

Preliminary Work

Reproduceability

Please use the pseudo-random number generator seed 1234 for reproduceability.

Different packages were downloaded and installed and these should also be installed in order to reproduce the results below.

Dependent variables

Our outcome variable is classe, a factor variable with 5 levels.

- 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)

Cross-validation

Cross-validation will be performed by subsampling our training data set randomly without replacement into 2 subsamples: TrainTrainingSet data (75% of the original Training data set) and TestTrainingSet data (25%). Our models will be fitted on the TrainTrainingSet data set, and tested on the TestTrainingSet 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 TestTrainingSet data set. Expected accuracy is the expected accuracy in the out-of-sample 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 a factor variable. We split the Training dataset into TrainTrainingSet and TestTrainingSet datasets.

Approach and Result

Set up the environment

```
list.of.packages <- c("caret", "lattice", "ggplot2", "randomForest", "rpart", "rpart.plot")
new.packages <- list.of.packages[!(list.of.packages %in% installed.packages()[,"Package"])]
if(length(new.packages)) install.packages(new.packages)

# install.packages("caret"); install.packages("randomForest"); install.packages("rpart");
library(lattice); library(ggplot2); library(caret); library(randomForest); library(rpart); library(rpart.plot)

## 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
```

Set the seed for reproducibility

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

Clean the dataset

```
# data load and clean up
trainingset <- read.csv("pml-training.csv", na.strings=c("NA", "#DIV/0!", ""))
testingset <- read.csv("pml-testing.csv", na.strings=c("NA", "#DIV/0!", ""))

# Delete columns with all missing values
trainingset<-trainingset[,colSums(is.na(trainingset)) == 0]
testingset <-testingset[,colSums(is.na(testingset)) == 0]

# Delete variables are irrelevant to our current project: user_name, raw_timestamp_part_1, raw_timestamp_part_2
```

```

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

# partition the data so that 75% of the training dataset into training and the remaining 25% to testing
traintrainset <- createDataPartition(y=trainingset$classe, p=0.75, list=FALSE)
TrainTrainingSet <- trainingset[traintrainset, ]
TestTrainingSet <- trainingset[-traintrainset, ]

```

Plot response variables in the TrainTrainingSet data set

```
plot(TrainTrainingSet$classe, col="green", main="Plot of classe in TrainTrainingSet", xlab="classe", ylab="Frequency")
```



Model I: Decision Tree

```

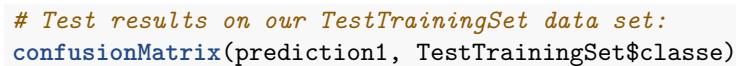
modell1 <- rpart(classe ~ ., data=TrainTrainingSet, method="class")

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

# Plot the Decision Tree
rpart.plot(modell1, main="Classification Tree", extra=102, under=TRUE, faclen=0)

```

	A
	B
	C
	D
	E



##

##

##

##

##

##

##

## 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

Model II: Random Forest

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

# Predicting:
prediction2 <- predict(model2, TestTrainingSet, type = "class")

# Test results on TestTrainingSet data set:
confusionMatrix(prediction2, TestTrainingSet$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
```

Compare two models

Random Forest algorithm performed better than Decision Trees. Accuracy for Random Forest model was 0.995 (95% CI: (0.993, 0.997)) compared to Decision Tree model with 0.739 (95% CI: (0.727, 0.752)). The Random Forests model is chosen. The expected out-of-sample error is estimated at 0.005, or 0.5%.

Submission

Here is the final outcome based on the Prediction Model 2 (Random Forest) applied against the Testing dataset

```
# 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
```