# Practical Maching Learning Course Project

Dan Greiber

Tuesday, March 17, 2015

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 dumbell 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.

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. 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 choosen, 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.

Choice reasons The 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 allow 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.

## Code and Results

Packages and Libraries Installing packages, loading libraries and setting the seed for reproduceability:

```
library(caret)

## Warning: package 'caret' was built under R version 3.1.3

## Loading required package: lattice
## Loading required package: ggplot2

## Warning: package 'ggplot2' was built under R version 3.1.2

library(randomForest)
```

## Warning: package 'randomForest' was built under R version 3.1.2

```
## randomForest 4.6-10
## Type rfNews() to see new features/changes/bug fixes.
library(rpart)
set.seed(1234)
# Loading the training data set into my R session replacing all missing with "NA"
trainingset <- read.csv(file="c:/Coursera/Practical Machine Learning/pml-training.csv", sep=",", na.str
# Loading the testing data set
testingset <- read.csv(file="c:/Coursera/Practical Machine Learning/pml-testing.csv", sep=",", na.string.csv", sep=", na.string.csv", sep=",", na.
# Check dimensions for number of variables and number of observations
dim(trainingset)
Loading data sets
## [1] 19622
                                     160
dim(testingset)
## [1] 20 160
# Delete columns with all missing values
trainingset<-trainingset[,colSums(is.na(trainingset)) == 0]</pre>
testingset <-testingset[,colSums(is.na(testingset)) == 0]</pre>
# Remove some variables that are irrelevant to our project
trainingset <-trainingset[,-c(1:7)]
testingset <-testingset[,-c(1:7)]</pre>
dim(trainingset)
## [1] 19622
                                       53
dim(testingset)
## [1] 20 53
Partitioning the training data set to allow cross-validation The training data set contains 53
variables and 19622 obs. The testing data set contains 53 variables and 20 obs. In order to perform
cross-validation, the training data set is partionned into 2 sets: subTraining (75%) and subTest (25%).
subsamples <- createDataPartition(y=trainingset$classe, p=0.75, list=FALSE)
subTraining <- trainingset[subsamples, ]</pre>
subTesting <- trainingset[-subsamples, ]</pre>
```

dim(subTraining)

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

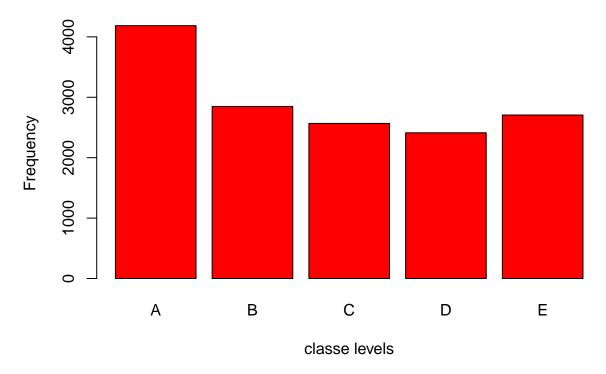
```
dim(subTesting)
```

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

A graphical look at our test data The variable "classe" contains 5 levels: A, B, C, D and E. A plot of the outcome variable will allow us to see the frequency of each levels in the subTraining data set and compare one another.

plot(subTraining\$classe, col="red", main="Bar Plot of levels of the variable classe within the subTrain

# Bar Plot of levels of the variable classe within the subTraining data s



From the graph, we can see that each level frequency is within the same order of magnitude of each other. Level A is the most frequent with more than 4000 occurrences while level D is the least frequent with about 2500 occurrences.

```
model1 <- rpart(classe ~ ., data=subTraining, method="class")
# Predicting:
prediction1 <- predict(model1, subTesting, type = "class")
confusionMatrix(prediction1, subTesting$classe)</pre>
```

First prediction model:

```
## Confusion Matrix and Statistics
##
##
             Reference
                Α
                           С
                                D
                                     Ε
## Prediction
                      В
##
            A 1235
                    157
                          16
                               50
                                    20
##
            В
                55
                    568
                          73
                               80 102
##
            С
                44
                    125
                         690
                              118
                                  116
                41
                     64
                              508
                                    38
##
            D
                          50
##
            Ε
                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
                                                      0.7924
## Balanced Accuracy
                          0.9080 0.7601
                                            0.8537
                                                               0.8307
model2 <- randomForest(classe ~. , data=subTraining, method="class")</pre>
# Predicting:
prediction2 <- predict(model2, subTesting, type = "class")</pre>
# Test results on subTesting data set:
confusionMatrix(prediction2, subTesting$classe)
Second prediction model:
```

```
## Confusion Matrix and Statistics
##
##
              Reference
## Prediction
                             С
                                   D
                                        Ε
                        В
                  Α
##
             A 1395
                        2
                             0
                                   0
##
            В
                     945
                             9
                                   0
                                        0
                  0
##
             С
                  0
                        2
                           844
                                   6
                                        0
##
            D
                        0
                                        0
                  0
                             2
                                798
##
            Ε
                                   0
                                     901
##
```

```
## 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
                                    0.9958
                                             0.9871
                                                       0.9925
                                                                1.0000
                           1.0000
## Specificity
                           0.9994
                                    0.9977
                                             0.9980
                                                       0.9995
                                                                1.0000
## Pos Pred Value
                           0.9986
                                    0.9906
                                             0.9906
                                                       0.9975
                                                                1.0000
## Neg Pred Value
                                                                1.0000
                           1.0000
                                    0.9990
                                             0.9973
                                                       0.9985
## Prevalence
                           0.2845
                                    0.1935
                                              0.1743
                                                       0.1639
                                                                0.1837
## Detection Rate
                                                                0.1837
                           0.2845
                                    0.1927
                                             0.1721
                                                       0.1627
## Detection Prevalence
                           0.2849
                                    0.1945
                                             0.1737
                                                       0.1631
                                                                0.1837
## Balanced Accuracy
                           0.9997
                                    0.9968
                                             0.9926
                                                       0.9960
                                                                1.0000
```

Summary of results As expected, Random Forest algorithm performed better than Decision Trees. 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 random Forest model is choosen. The accuracy of the model is 0.995. 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.

#### Submission

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

```
## 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(predictall)
```

## References

- [1] 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.
- [2] Krzysztof Gra??bczewski and Norbert Jankowski. Feature Selection with Decision Tree Criterion.