

# Practical Machine Learning - Course Project

*gustavohrp*

*24 de dezembro de 2015*

GitHub repo: <https://github.com/gustavohrp/Practical-Machine-Learning—Course-Project.git>

RPubs: <http://rpubs.com/gustavohrp/138287>

## Background Introduction

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 Source

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.

## Project Goal

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.

1. 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 :-).
2. 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.

## Loading Library and Getting Data

The following Libraries were used for this project, which you should install - if not done yet - and load on your working environment.

```
# Libraries used for the data processing  
library(caret)
```

```
## Warning: package 'caret' was built under R version 3.2.2
```

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

```
library(rpart)  
library(rpart.plot)  
library(RColorBrewer)  
library(rattle)
```

```
## Warning: package 'rattle' was built under R version 3.2.2
```

```
## Rattle: A free graphical interface for data mining with R.  
## Version 4.0.5 Copyright (c) 2006-2015 Togaware Pty Ltd.  
## Type 'rattle()' to shake, rattle, and roll your data.
```

```
library(randomForest)
```

```
## randomForest 4.6-12  
## Type rfNews() to see new features/changes/bug fixes.
```

```
library(knitr)
```

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

```
rm(list = ls())  
if (!file.exists("pml-training.csv")) {  
  download.file("http://d396qusza40orc.cloudfront.net/predmachlearn/pml-training.csv", destfile = "pml-t  
}  
if (!file.exists("pml-testing.csv")) {  
  download.file("http://d396qusza40orc.cloudfront.net/predmachlearn/pml-testing.csv", destfile = "pml-t  
}  
training <- read.csv("pml-training.csv", na.strings=c("NA", "#DIV/0!", ""))  
testing <- read.csv("pml-testing.csv", na.strings=c("NA", "#DIV/0!", ""))
```

In order to reproduce the same results, you need to load the same seed.

```
set.seed(12345)
```

## Cleaning up the data and Bootstrap

Here, I remove columns full of NAs and remove features that are not in the testing set. The features containing NAs are the variance, mean and stddev within each window for each feature. Since the testing dataset has no time-dependence, these values are useless and can be disregarded. I also remove the first 7 features since they are related to the time-series or are not numeric.

Partitioning Training data set into two data sets, 60% for myTraining, 40% for myTesting.

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

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

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

I am now going to reduce the number of features by removing variables with nearly zero variance, variables that are almost always NA, and variables that don't make intuitive sense for prediction.

```
#removing variables nearly zero variance
nzv <- nearZeroVar(myTraining, saveMetrics=TRUE)
myTraining <- myTraining[,nzv$nzv==FALSE]

nzv<- nearZeroVar(myTesting,saveMetrics=TRUE)
myTesting <- myTesting[,nzv$nzv==FALSE]

#removing first ID variable so that it does not interfere with ML Algorithms
myTraining <- myTraining[,c(-1)]

#removing variables that have more than a 70% of NA's.
trainingV3 <- myTraining
for(i in 1:length(myTraining)) {
  if( sum( is.na( myTraining[, i] ) ) /nrow(myTraining) >= .7) {
    for(j in 1:length(trainingV3)) {
      if( length( grep(names(myTraining[i]), names(trainingV3)[j]) ) == 1) {
        trainingV3 <- trainingV3[ , -j]
      }
    }
  }
}

# Set back to the original variable name
myTraining <- trainingV3
rm(trainingV3)
```

Doing the same cleaning to myTesting and testing data sets

```
clean1 <- colnames(myTraining)
clean2 <- colnames(myTraining[, -58])
myTesting <- myTesting[clean1]
testing <- testing[clean2]
```

```

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])
    }
  }
}

# To get the same class between testing and myTraining
testing <- rbind(myTraining[2, -58] , testing)
testing <- testing[-1,]

```

## Prediction with Decision Trees

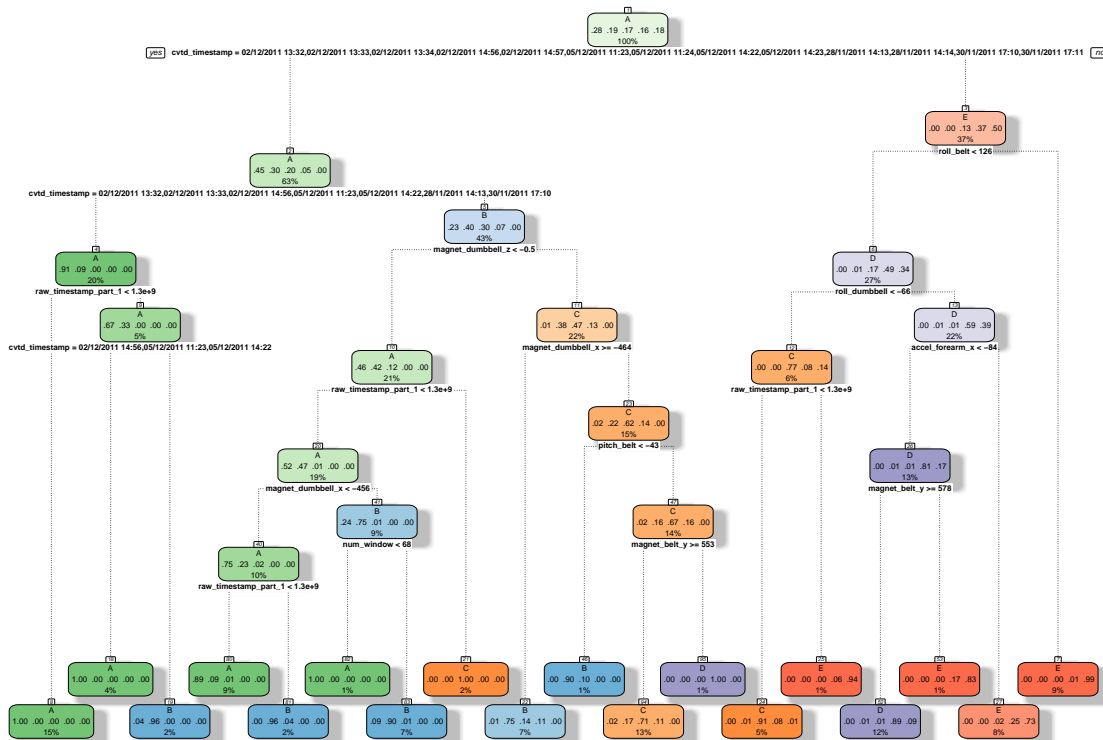
I decided to start with a Random Forest model, to see if it would have acceptable performance. I fit the model on ptrain1, and instruct the “train” function to use 3-fold cross-validation to select optimal tuning parameters for the model.

I see that it decided to use 500 trees and try 27 variables at each split. Now, I use the fitted model to predict the label (“classe”) in ptrain2, and show the confusion matrix to compare the predicted versus the actual labels:

```

set.seed(12345)
modFitA1 <- rpart(classe ~ ., data=myTraining, method="class")
fancyRpartPlot(modFitA1)

```



```

predictionsA1 <- predict(modFitA1, myTesting, type = "class")
cmtree <- confusionMatrix(predictionsA1, myTesting$classe)
cmtree

```

```
## Confusion Matrix and Statistics
```

```
##
##           Reference
## Prediction    A    B    C    D    E
##           A 2150   60    7    1    0
##           B   61 1260   69   64    0
##           C   21  188 1269  143    4
##           D    0   10   14  857   78
##           E    0    0    9  221 1360
##
```

```
## Overall Statistics
```

```
##
##           Accuracy : 0.8789
##           95% CI   : (0.8715, 0.8861)
##           No Information Rate : 0.2845
##           P-Value [Acc > NIR] : < 2.2e-16
##
##           Kappa : 0.8468
##           McNemar's Test P-Value : NA
##
```

```
## Statistics by Class:
```

```
##
##           Class: A Class: B Class: C Class: D Class: E
## Sensitivity      0.9633   0.8300   0.9276   0.6664   0.9431
## Specificity      0.9879   0.9693   0.9450   0.9845   0.9641
## Pos Pred Value   0.9693   0.8666   0.7809   0.8936   0.8553
## Neg Pred Value   0.9854   0.9596   0.9841   0.9377   0.9869
## Prevalence       0.2845   0.1935   0.1744   0.1639   0.1838
## Detection Rate   0.2740   0.1606   0.1617   0.1092   0.1733
## Detection Prevalence 0.2827   0.1853   0.2071   0.1222   0.2027
## Balanced Accuracy 0.9756   0.8997   0.9363   0.8254   0.9536

```

```
plot(cmtree$table, col = cmtree$byClass, main = paste("Decision Tree Confusion Matrix: Accuracy =", round(0.8789, 2)))

```

## Decision Tree Confusion Matrix: Accuracy = 0.8789

		A	B	C	D	E
Reference	A					
	B					
Reference	C					
	D					
Reference	E					
	ECB					
		Prediction				

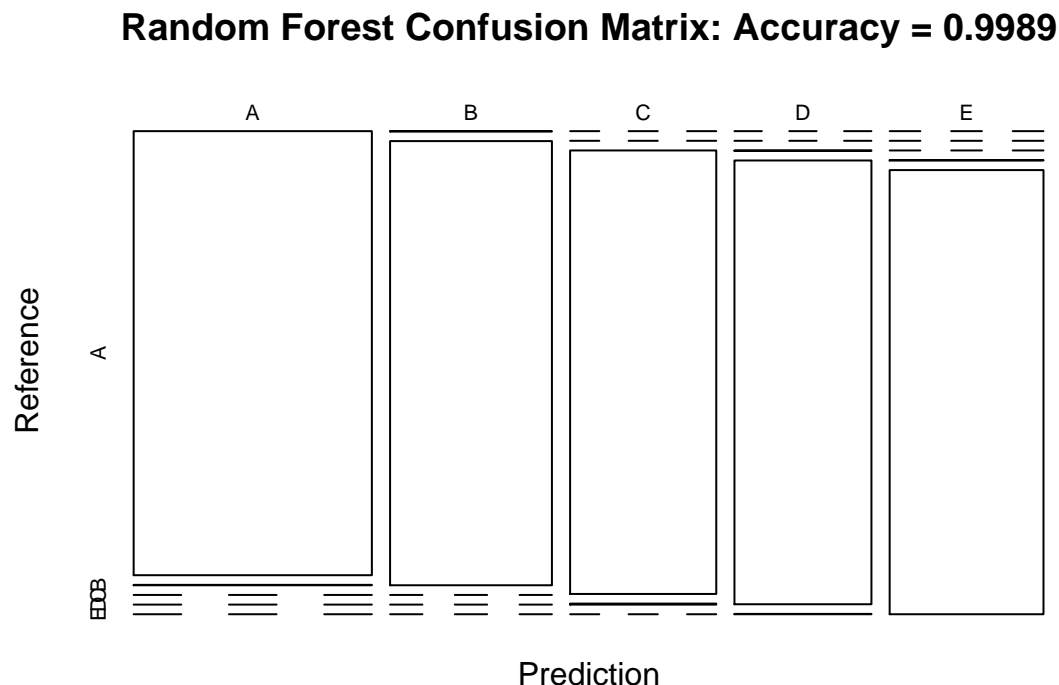
## Prediction with Random Forests

```
set.seed(12345)
modFitB1 <- randomForest(classe ~ ., data=myTraining)
predictionB1 <- predict(modFitB1, myTesting, type = "class")
cmrf <- confusionMatrix(predictionB1, myTesting$classe)
cmrf
```

```
## Confusion Matrix and Statistics
##
##           Reference
## Prediction   A    B    C    D    E
##           A 2231     2     0     0     0
##           B   1 1516     0     0     0
##           C    0     0 1367     3     0
##           D    0     0   1 1282     1
##           E    0     0    0     1 1441
##
## Overall Statistics
##
##           Accuracy : 0.9989
##           95% CI : (0.9978, 0.9995)
##           No Information Rate : 0.2845
##           P-Value [Acc > NIR] : < 2.2e-16
##
##           Kappa : 0.9985
##           McNemar's Test P-Value : NA
##
## Statistics by Class:
```

```
##
##               Class: A Class: B Class: C Class: D Class: E
## Sensitivity      0.9996   0.9987   0.9993   0.9969   0.9993
## Specificity      0.9996   0.9998   0.9995   0.9997   0.9998
## Pos Pred Value   0.9991   0.9993   0.9978   0.9984   0.9993
## Neg Pred Value   0.9998   0.9997   0.9998   0.9994   0.9998
## Prevalence       0.2845   0.1935   0.1744   0.1639   0.1838
## Detection Rate   0.2843   0.1932   0.1742   0.1634   0.1837
## Detection Prevalence 0.2846 0.1933 0.1746 0.1637 0.1838
## Balanced Accuracy 0.9996   0.9993   0.9994   0.9983   0.9996
```

```
plot(cmrnf$table, col = cmtree$byClass, main = paste("Random Forest Confusion Matrix: Accuracy =", round
```



The accuracy is 99.8%, thus my predicted accuracy for the out-of-sample error is 0.2%.

This is an excellent result, so rather than trying additional algorithms, I will use Random Forests to predict on the test set.

## Predicting Results on the Test Data

Random Forests gave an Accuracy in the myTesting dataset of 99.89%, which was more accurate than what I got from the Decision Trees. The expected out-of-sample error is  $100 - 99.89 = 0.11\%$ .

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

# Write the results to a text file 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(predictionsB2)
```