

# Practical\_Machine\_Learning

## Project Description

### 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. Peer Review Portion

Your submission for the Peer Review portion 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 :-). Course Project Prediction Quiz Portion

Apply your machine learning algorithm to the 20 test cases available in the test data above and submit your predictions in appropriate format to the Course Project Prediction Quiz for automated grading.

### Reproducibility

Due to security concerns with the exchange of R code, your code will not be run during the evaluation by your classmates. Please be sure that if they download the repo, they will be able to view the compiled HTML version of your analysis.

Load packages:

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

```
## 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
```

```
## Rattle: A free graphical interface for data mining with R.  
## Version 4.1.0 Copyright (c) 2006-2015 Togaware Pty Ltd.  
## Type 'rattle()' to shake, rattle, and roll your data.
```

Read data:

Partition the training set in two:

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

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

Clean the data:

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

Verify that all column names in both data sets are the same:

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

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

There are more options, data could be split in multiple training sets and after removing zeros and NAs the model can be created or use the data as a single set and compare the test set and training set after cleaning the data. Using the whole data is a better approach.

Removing the columns with more than 60% NA. Partition the data.

```
## [1] 5.1 4.1 4.1 2.1
```

```
## Confusion Matrix and Statistics
##
##           Reference
## Prediction    A     B     C     D     E
##           A 2153    58     3     1     0
##           B   51 1258    87    57     0
##           C   28  194 1254   205     2
##           D    0    8   16  800    71
##           E    0    0    8  223 1369
##
## Overall Statistics
##
##           Accuracy : 0.871
##           95% CI : (0.8634, 0.8784)
##           No Information Rate : 0.2845
##           P-Value [Acc > NIR] : < 2.2e-16
##
##           Kappa : 0.8368
##           McNemar's Test P-Value : NA
##
## Statistics by Class:
##
##           Class: A Class: B Class: C Class: D Class: E
## Sensitivity           0.9646    0.8287    0.9167    0.6221    0.9494
## Specificity           0.9890    0.9692    0.9338    0.9855    0.9639
## Pos Pred Value        0.9720    0.8658    0.7451    0.8939    0.8556
## Neg Pred Value        0.9860    0.9593    0.9815    0.9301    0.9883
## Prevalence            0.2845    0.1935    0.1744    0.1639    0.1838
## Detection Rate        0.2744    0.1603    0.1598    0.1020    0.1745
## Detection Prevalence  0.2823    0.1852    0.2145    0.1141    0.2039
## Balanced Accuracy      0.9768    0.8990    0.9252    0.8038    0.9567
```

Random Forest Prediction Model

## Decision Tree

```

## Confusion Matrix and Statistics
##
##           Reference
## Prediction    A    B    C    D    E
##           A 2231    0    0    0    0
##           B   1 1518    3    0    0
##           C    0    0 1363    7    0
##           D    0    0    2 1279    1
##           E    0    0    0    0 1441
##
## Overall Statistics
##
##           Accuracy : 0.9982
##           95% CI : (0.997, 0.999)
##           No Information Rate : 0.2845
##           P-Value [Acc > NIR] : < 2.2e-16
##
##           Kappa : 0.9977
##           McNemar's Test P-Value : NA
##
## Statistics by Class:
##
##           Class: A Class: B Class: C Class: D Class: E
## Sensitivity          0.9996   1.0000   0.9963   0.9946   0.9993
## Specificity          1.0000   0.9994   0.9989   0.9995   1.0000
## Pos Pred Value       1.0000   0.9974   0.9949   0.9977   1.0000
## Neg Pred Value       0.9998   1.0000   0.9992   0.9989   0.9998
## Prevalence           0.2845   0.1935   0.1744   0.1639   0.1838
## Detection Rate       0.2843   0.1935   0.1737   0.1630   0.1837
## Detection Prevalence 0.2843   0.1940   0.1746   0.1634   0.1837
## Balanced Accuracy     0.9998   0.9997   0.9976   0.9970   0.9997

```