# Classification of Barbell Lifting

*Erik White*

*November 19, 2017*

## 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 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://web.archive.org/web/20161224072740/http:/groupware.les.inf.puc-rio.br/har (see the section on the Weight Lifting Exercise Dataset)." - https://www.coursera.org/learn/practical-machine-learning/peer/R43St/prediction-assignment-writeup

## Objective

We will attempt to implement a machine learning model that will utilize various measurements captured in the above study to predict the classification of how a particular barbell lift was performed. Each observation was classified into on one of five categories, stored within the 'classe' variable of the datasets that we will be utilizing.

- Class A - "exactly according to the specification" aka the correct implementation of the exercise
- Class B - "throwing the elbows to the front"
- Class C - "lifting the dumbbell only halfway"
- Class D - "lowering the dumbbell only halfway"
- Class E - "throwing the hips to the front"

## Pre-processing

We begin by loading the datasets and exploring some of the features:

```
wltrain <- read.csv("https://d396qusza40orc.cloudfront.net/predmachlearn/pml-training.csv", na.strings =
wlvalid <- read.csv("https://d396qusza40orc.cloudfront.net/predmachlearn/pml-testing.csv", na.strings =
str(wltrain)
```
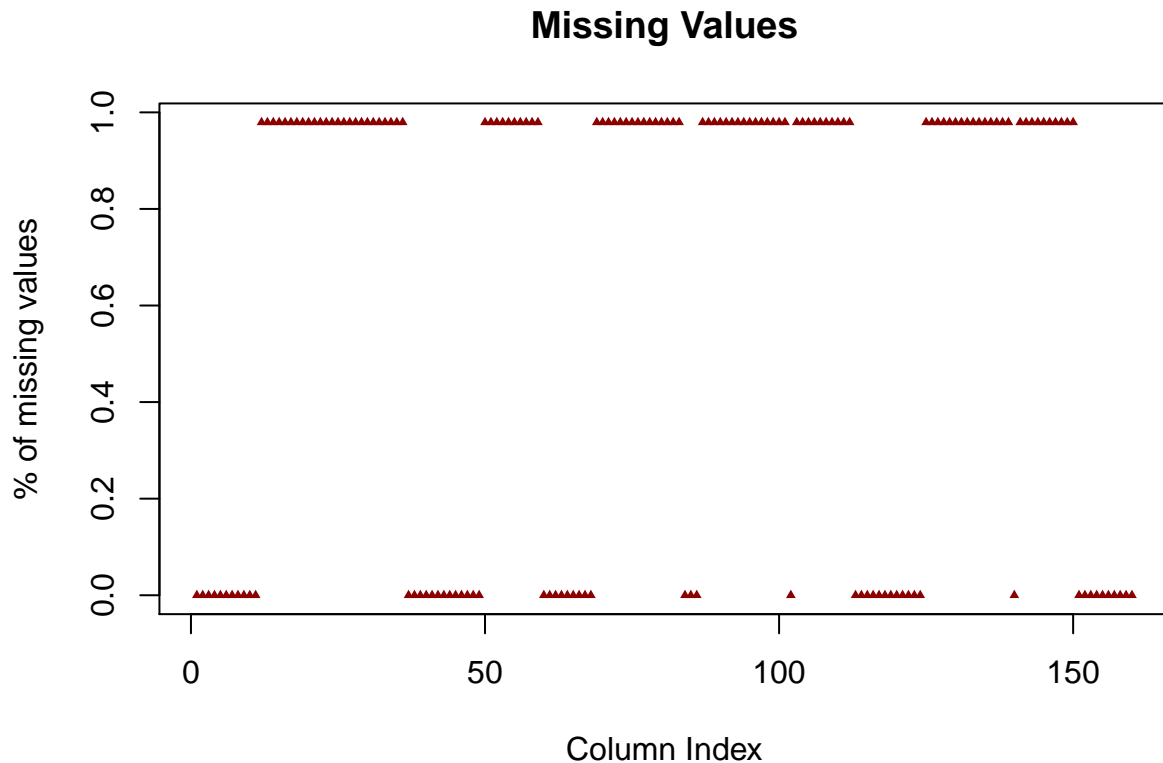
```
## 'data.frame':    19622 obs. of  160 variables:
##  $ X                    : int  1 2 3 4 5 6 7 8 9 10 ...
##  $ user_name            : Factor w/ 6 levels "adelmo","carlitos",..: 2 2 2 2 2 2 2 2 2 2 ...
##  $ raw_timestamp_part_1 : int  1323084231 1323084231 1323084231 1323084232 1323084232 1323084232
##  $ raw_timestamp_part_2 : int  788290 808298 820366 120339 196328 304277 368296 440390 484323 484
##  $ cvtd_timestamp       : Factor w/ 20 levels "02/12/2011 13:32",..: 9 9 9 9 9 9 9 9 9 9 ...
##  $ new_window           : Factor w/ 2 levels "no","yes": 1 1 1 1 1 1 1 1 1 1 ...
##  $ num_window           : int  11 11 11 12 12 12 12 12 12 12 ...
##  $ roll_belt            : num  1.41 1.41 1.42 1.48 1.48 1.45 1.42 1.42 1.43 1.45 ...
##  $ pitch_belt           : num  8.07 8.07 8.07 8.05 8.07 8.06 8.09 8.13 8.16 8.17 ...
##  $ yaw_belt             : num  -94.4 -94.4 -94.4 -94.4 -94.4 -94.4 -94.4 -94.4 -94.4 -94.4 ...
##  $ total_accel_belt     : int  3 3 3 3 3 3 3 3 3 3 ...
```

```
##  $ kurtosis_roll_belt       : Factor w/ 396 levels "-0.016850","-0.021024",..: NA NA NA NA NA NA NA N
##  $ kurtosis_picth_belt      : Factor w/ 316 levels "-0.021887","-0.060755",..: NA NA NA NA NA NA NA N
##  $ kurtosis_yaw_belt        : Factor w/ 1 level "#DIV/0!": NA NA NA NA NA NA NA NA NA NA ...
##  $ skewness_roll_belt       : Factor w/ 394 levels "-0.003095","-0.010002",..: NA NA NA NA NA NA NA N
##  $ skewness_roll_belt.1     : Factor w/ 337 levels "-0.005928","-0.005960",..: NA NA NA NA NA NA NA N
##  $ skewness_yaw_belt        : Factor w/ 1 level "#DIV/0!": NA NA NA NA NA NA NA NA NA NA ...
##  $ max_roll_belt            : num  NA NA NA NA NA NA NA NA NA NA ...
##  $ max_picth_belt           : int  NA NA NA NA NA NA NA NA NA NA ...
##  $ max_yaw_belt             : Factor w/ 67 levels "-0.1","-0.2",..: NA NA NA NA NA NA NA NA NA NA ...
##  $ min_roll_belt            : num  NA NA NA NA NA NA NA NA NA NA ...
##  $ min_pitch_belt           : int  NA NA NA NA NA NA NA NA NA NA ...
##  $ min_yaw_belt             : Factor w/ 67 levels "-0.1","-0.2",..: NA NA NA NA NA NA NA NA NA NA ...
##  $ amplitude_roll_belt      : num  NA NA NA NA NA NA NA NA NA NA ...
##  $ amplitude_pitch_belt     : int  NA NA NA NA NA NA NA NA NA NA ...
##  $ amplitude_yaw_belt       : Factor w/ 3 levels "#DIV/0!","0.00",..: NA NA NA NA NA NA NA NA NA NA .
##  $ var_total_accel_belt     : num  NA NA NA NA NA NA NA NA NA NA ...
##  $ avg_roll_belt            : num  NA NA NA NA NA NA NA NA NA NA ...
##  $ stddev_roll_belt         : num  NA NA NA NA NA NA NA NA NA NA ...
##  $ var_roll_belt            : num  NA NA NA NA NA NA NA NA NA NA ...
##  $ avg_pitch_belt           : num  NA NA NA NA NA NA NA NA NA NA ...
##  $ stddev_pitch_belt        : num  NA NA NA NA NA NA NA NA NA NA ...
##  $ var_pitch_belt           : num  NA NA NA NA NA NA NA NA NA NA ...
##  $ avg_yaw_belt             : num  NA NA NA NA NA NA NA NA NA NA ...
##  $ stddev_yaw_belt          : num  NA NA NA NA NA NA NA NA NA NA ...
##  $ var_yaw_belt             : num  NA NA NA NA NA NA NA NA NA NA ...
##  $ gyros_belt_x             : num  0 0.02 0 0.02 0.02 0.02 0.02 0.02 0.02 0.03 ...
##  $ gyros_belt_y             : num  0 0 0 0 0.02 0 0 0 0 0 ...
##  $ gyros_belt_z             : num  -0.02 -0.02 -0.02 -0.03 -0.02 -0.02 -0.02 -0.02 -0.02 0 ...
##  $ accel_belt_x             : int  -21 -22 -20 -22 -21 -21 -22 -22 -20 -21 ...
##  $ accel_belt_y             : int  4 4 5 3 2 4 3 4 2 4 ...
##  $ accel_belt_z             : int  22 22 23 21 24 21 21 21 24 22 ...
##  $ magnet_belt_x            : int  -3 -7 -2 -6 -6 0 -4 -2 1 -3 ...
##  $ magnet_belt_y            : int  599 608 600 604 600 603 599 603 602 609 ...
##  $ magnet_belt_z            : int  -313 -311 -305 -310 -302 -312 -311 -313 -312 -308 ...
##  $ roll_arm                 : num  -128 -128 -128 -128 -128 -128 -128 -128 -128 -128 ...
##  $ pitch_arm                : num  22.5 22.5 22.5 22.1 22.1 22 21.9 21.8 21.7 21.6 ...
##  $ yaw_arm                  : num  -161 -161 -161 -161 -161 -161 -161 -161 -161 -161 ...
##  $ total_accel_arm          : int  34 34 34 34 34 34 34 34 34 34 ...
##  $ var_accel_arm            : num  NA NA NA NA NA NA NA NA NA NA ...
##  $ avg_roll_arm             : num  NA NA NA NA NA NA NA NA NA NA ...
##  $ stddev_roll_arm          : num  NA NA NA NA NA NA NA NA NA NA ...
##  $ var_roll_arm             : num  NA NA NA NA NA NA NA NA NA NA ...
##  $ avg_pitch_arm            : num  NA NA NA NA NA NA NA NA NA NA ...
##  $ stddev_pitch_arm         : num  NA NA NA NA NA NA NA NA NA NA ...
##  $ var_pitch_arm            : num  NA NA NA NA NA NA NA NA NA NA ...
##  $ avg_yaw_arm              : num  NA NA NA NA NA NA NA NA NA NA ...
##  $ stddev_yaw_arm           : num  NA NA NA NA NA NA NA NA NA NA ...
##  $ var_yaw_arm              : num  NA NA NA NA NA NA NA NA NA NA ...
##  $ gyros_arm_x              : num  0 0.02 0.02 0.02 0 0.02 0 0.02 0.02 0.02 ...
##  $ gyros_arm_y              : num  0 -0.02 -0.02 -0.03 -0.03 -0.03 -0.03 -0.02 -0.03 -0.03 ...
##  $ gyros_arm_z              : num  -0.02 -0.02 -0.02 0.02 0 0 0 0 -0.02 -0.02 ...
##  $ accel_arm_x              : int  -288 -290 -289 -289 -289 -289 -289 -289 -288 -288 ...
##  $ accel_arm_y              : int  109 110 110 111 111 111 111 111 109 110 ...
##  $ accel_arm_z              : int  -123 -125 -126 -123 -123 -122 -125 -124 -122 -124 ...
```

```
##  $ magnet_arm_x            : int   -368 -369 -368 -372 -374 -369 -373 -372 -369 -376 ...
##  $ magnet_arm_y            : int   337 337 344 344 337 342 336 338 341 334 ...
##  $ magnet_arm_z            : int   516 513 513 512 506 513 509 510 518 516 ...
##  $ kurtosis_roll_arm       : Factor w/ 329 levels "-0.02438","-0.04190",..: NA NA NA NA NA NA NA NA N
##  $ kurtosis_picth_arm      : Factor w/ 327 levels "-0.00484","-0.01311",..: NA NA NA NA NA NA NA NA N
##  $ kurtosis_yaw_arm        : Factor w/ 394 levels "-0.01548","-0.01749",..: NA NA NA NA NA NA NA NA N
##  $ skewness_roll_arm       : Factor w/ 330 levels "-0.00051","-0.00696",..: NA NA NA NA NA NA NA NA N
##  $ skewness_pitch_arm      : Factor w/ 327 levels "-0.00184","-0.01185",..: NA NA NA NA NA NA NA NA N
##  $ skewness_yaw_arm        : Factor w/ 394 levels "-0.00311","-0.00562",..: NA NA NA NA NA NA NA NA N
##  $ max_roll_arm            : num   NA NA NA NA NA NA NA NA NA NA ...
##  $ max_picth_arm           : num   NA NA NA NA NA NA NA NA NA NA ...
##  $ max_yaw_arm             : int   NA NA NA NA NA NA NA NA NA NA ...
##  $ min_roll_arm            : num   NA NA NA NA NA NA NA NA NA NA ...
##  $ min_pitch_arm           : num   NA NA NA NA NA NA NA NA NA NA ...
##  $ min_yaw_arm             : int   NA NA NA NA NA NA NA NA NA NA ...
##  $ amplitude_roll_arm      : num   NA NA NA NA NA NA NA NA NA NA ...
##  $ amplitude_pitch_arm     : num   NA NA NA NA NA NA NA NA NA NA ...
##  $ amplitude_yaw_arm       : int   NA NA NA NA NA NA NA NA NA NA ...
##  $ roll_dumbbell           : num   13.1 13.1 12.9 13.4 13.4 ...
##  $ pitch_dumbbell          : num   -70.5 -70.6 -70.3 -70.4 -70.4 ...
##  $ yaw_dumbbell            : num   -84.9 -84.7 -85.1 -84.9 -84.9 ...
##  $ kurtosis_roll_dumbbell  : Factor w/ 397 levels "-0.0035","-0.0073",..: NA NA NA NA NA NA NA NA NA
##  $ kurtosis_picth_dumbbell : Factor w/ 400 levels "-0.0163","-0.0233",..: NA NA NA NA NA NA NA NA NA
##  $ kurtosis_yaw_dumbbell   : Factor w/ 1 level "#DIV/0!": NA NA NA NA NA NA NA NA NA NA ...
##  $ skewness_roll_dumbbell  : Factor w/ 400 levels "-0.0082","-0.0096",..: NA NA NA NA NA NA NA NA NA
##  $ skewness_pitch_dumbbell : Factor w/ 401 levels "-0.0053","-0.0084",..: NA NA NA NA NA NA NA NA NA
##  $ skewness_yaw_dumbbell   : Factor w/ 1 level "#DIV/0!": NA NA NA NA NA NA NA NA NA NA ...
##  $ max_roll_dumbbell       : num   NA NA NA NA NA NA NA NA NA NA ...
##  $ max_picth_dumbbell      : num   NA NA NA NA NA NA NA NA NA NA ...
##  $ max_yaw_dumbbell        : Factor w/ 72 levels "-0.1","-0.2",..: NA NA NA NA NA NA NA NA NA NA ...
##  $ min_roll_dumbbell       : num   NA NA NA NA NA NA NA NA NA NA ...
##  $ min_pitch_dumbbell      : num   NA NA NA NA NA NA NA NA NA NA ...
##  $ min_yaw_dumbbell        : Factor w/ 72 levels "-0.1","-0.2",..: NA NA NA NA NA NA NA NA NA NA ...
##  $ amplitude_roll_dumbbell : num   NA NA NA NA NA NA NA NA NA NA ...
##   [list output truncated]
```

An initial review of the training set suggests that there are a material number of missing values throughout our dataset. Further analysis confirms this:

```
natrain <- colSums(is.na(wltrain))/nrow(wltrain)
natest <- colSums(is.na(wlvalid))/nrow(wlvalid)
plot(natrain, main = "Missing Values", xlab = "Column Index", ylab = "% of missing values", type = "n")
points(natrain, pch = 17, cex = .5, col = "dark red")
```

## Missing Values



There's a significant population of columns that have over 97% of their values missing. We choose to remove these columns from our data sets. All remaining columns do not have any missing values.

```r
wltrain <- wltrain[,names(which(natrain<.97))]
wlvalid <- wlvalid[,names(which(natest<.97))]
```

We can also note that there are a number of columns that are not directly related to the execution of the exercises. We choose to remove these from our model as they could potentially introduce unwanted bias into our model.

```r
drop <- c("X", "user_name","raw_timestamp_part_1", "raw_timestamp_part_2","cvtd_timestamp","new_window"
wltrain <- wltrain[,!names(wltrain) %in% drop]
wlvalid <- wlvalid[,!names(wlvalid) %in% drop]
```

## Cross Validation

In order to evaluate the performance of our model, we choose to cross validate by splitting the 'wltrain' data into a training and a testing dataset. We will keep the 'wlvalid' dataset off to the side to use as a validation dataset once we have finalized our model.

```r
library(caret)
```

```
## Warning: package 'caret' was built under R version 3.4.2
```

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

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

```
set.seed(123)
inTrain <- createDataPartition(wltrain$classe, p=.75, list = FALSE)
training <-  wltrain[inTrain,]
testing <- wltrain[-inTrain,]
```

## Model Selection

We choose to implement a random forest algorithm for our prediction model. We then create a confusion matrix between the

```
library(randomForest)
```

```
## Warning: package 'randomForest' was built under R version 3.4.2
```

```
## 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.seed(456)
RFmodel <- randomForest(classe ~ ., data = training)
RFprediction <- predict(RFmodel, testing)
confusionMatrix(RFprediction, testing$classe)
```

```
## Confusion Matrix and Statistics
##
##           Reference
## Prediction    A    B    C    D    E
##          A 1394    1    0    0    0
##          B    1  947    8    0    0
##          C    0    1  847   10    0
##          D    0    0    0  792    1
##          E    0    0    0    2  900
##
## 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   0.9979   0.9906   0.9851   0.9989
## Specificity           0.9997   0.9977   0.9973   0.9998   0.9995
## Pos Pred Value        0.9993   0.9906   0.9872   0.9987   0.9978
```

```
## Neg Pred Value       0.9997   0.9995   0.9980   0.9971   0.9998
## Prevalence           0.2845   0.1935   0.1743   0.1639   0.1837
## Detection Rate       0.2843   0.1931   0.1727   0.1615   0.1835
## Detection Prevalence 0.2845   0.1949   0.1750   0.1617   0.1839
## Balanced Accuracy    0.9995   0.9978   0.9940   0.9924   0.9992
```

When applied to our testing dataset, our model is 99.51% accurate. This suggests an expected out-of-sample error rate of .49%.