

Human Activity Recognition

Using Machine Learning to Predict the Correctness of Weight Lifting Exercises

Practical Machine Learning, JHU Data Science Specialization, Coursera

Introduction

The subject matter of this project is how correctly subjects are performing weight lifting exercises - specifically, weight lifting of a dumbbell. Sensors were placed on each of the subjects' belt, arm, forearm, and dumbbell. The sensors record measurements such as yaw, pitch, and roll, as well as 3-axis measurements for the gyroscope, accelerometer, and magnetometer. The subjects were asked to perform exercises in 5 distinct ways; one of the ways was the correct motion, while the other 4 were common mistakes made while weightlifting. The focus of this project is to predict which way the subjects performed the exercise based on the measurements from the sensors, using a specific set of data as the training set for the model and then applying it to a test set. More information and the original data can be found at <http://groupware.les.inf.puc-rio.br/har> (<http://groupware.les.inf.puc-rio.br/har>).

The resulting model developed was very successful, with a predicted accuracy of about 98% and able to predict the class of activity correctly in all 20 test cases.

Building a Model on the Data

Loading the data

```
trainURL <- "http://d396qusza40orc.cloudfront.net/predmachlearn/pml-training.csv"
testURL <- "http://d396qusza40orc.cloudfront.net/predmachlearn/pml-testing.csv"
download.file(trainURL, "pml-training.csv", mode = "wb")
download.file(testURL, "pml-testing.csv", mode = "wb")
traindata <- read.csv("pml-training.csv")
testdata <- read.csv("pml-testing.csv")
str(traindata)
```

```
## '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 ...
## $ raw_timestamp_part_1 : int  1323084231 1323084231 1323084231 1323084232 1323084232 132308423
2 1323084232 1323084232 1323084232 1323084232 ...
## $ raw_timestamp_part_2 : int  788290 808298 820366 120339 196328 304277 368296 440390 484323 4
84434 ...
## $ cvtd_timestamp      : Factor w/ 20 levels "02/12/2011 13:32",...: 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 ...
## $ num_window          : int  11 11 11 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 33 3 ...
## $ kurtosis_roll_belt    : Factor w/ 397 levels "", "-0.016850",...: 1 1 1 1 1 1 1 1 1 1 ...
## $ kurtosis_pitch_belt   : Factor w/ 317 levels "", "-0.021887",...: 1 1 1 1 1 1 1 1 1 1 ...
## $ kurtosis_yaw_belt     : Factor w/ 2 levels "", "#DIV/0!": 1 1 1 1 1 1 1 1 1 1 ...
## $ skewness_roll_belt    : Factor w/ 395 levels "", "-0.003095",...: 1 1 1 1 1 1 1 1 1 1 ...
## $ skewness_roll_belt.1  : Factor w/ 338 levels "", "-0.005928",...: 1 1 1 1 1 1 1 1 1 1 ...
## $ skewness_yaw_belt     : Factor w/ 2 levels "", "#DIV/0!": 1 1 1 1 1 1 1 1 1 1 ...
## $ max_roll_belt         : num  NA NA NA NA NA NA NA NA NA NA ...
## $ max_pitch_belt        : int  NA NA NA NA NA NA NA NA NA NA ...
## $ max_yaw_belt          : Factor w/ 68 levels "", "-0.1", "-0.2",...: 1 1 1 1 1 1 1 1 1 1 ...
## $ 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/ 68 levels "", "-0.1", "-0.2",...: 1 1 1 1 1 1 1 1 1 1 ...
## $ 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/ 4 levels "", "#DIV/0!", "0.00",...: 1 1 1 1 1 1 1 1 1 1 ...
## $ 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/ 330 levels "", "-0.02438",...: 1 1 1 1 1 1 1 1 1 1 ...
## $ kurtosis_picth_arm  : Factor w/ 328 levels "", "-0.00484",...: 1 1 1 1 1 1 1 1 1 1 ...
## $ kurtosis_yaw_arm    : Factor w/ 395 levels "", "-0.01548",...: 1 1 1 1 1 1 1 1 1 1 ...
## $ skewness_roll_arm   : Factor w/ 331 levels "", "-0.00051",...: 1 1 1 1 1 1 1 1 1 1 ...
## $ skewness_pitch_arm  : Factor w/ 328 levels "", "-0.00184",...: 1 1 1 1 1 1 1 1 1 1 ...
## $ skewness_yaw_arm    : Factor w/ 395 levels "", "-0.00311",...: 1 1 1 1 1 1 1 1 1 1 ...
## $ 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/ 398 levels "", "-0.0035", "-0.0073",...: 1 1 1 1 1 1 1 1 1 1 .
..
## $ kurtosis_picth_dumbbell : Factor w/ 401 levels "", "-0.0163", "-0.0233",...: 1 1 1 1 1 1 1 1 1 1 .
..
## $ kurtosis_yaw_dumbbell  : Factor w/ 2 levels "", "#DIV/0!": 1 1 1 1 1 1 1 1 1 1 ...
## $ skewness_roll_dumbbell : Factor w/ 401 levels "", "-0.0082", "-0.0096",...: 1 1 1 1 1 1 1 1 1 1 .
..
## $ skewness_pitch_dumbbell : Factor w/ 402 levels "", "-0.0053", "-0.0084",...: 1 1 1 1 1 1 1 1 1 1 .
..
## $ skewness_yaw_dumbbell  : Factor w/ 2 levels "", "#DIV/0!": 1 1 1 1 1 1 1 1 1 1 ...
## $ 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/ 73 levels "", "-0.1", "-0.2",...: 1 1 1 1 1 1 1 1 1 1 ...
## $ 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/ 73 levels "", "-0.1", "-0.2",...: 1 1 1 1 1 1 1 1 1 1 ...
## $ amplitude_roll_dumbbell : num  NA NA NA NA NA NA NA NA NA NA ...
## [list output truncated]

```

```
str(testdata)
```

```
## 'data.frame':    20 obs. of  160 variables:
## $ X                      : int  1 2 3 4 5 6 7 8 9 10 ...
## $ user_name              : Factor w/ 6 levels "adelmo","carlitos",...: 6 5 5 1 4 5 5 2 3 ...
## $ raw_timestamp_part_1   : int  1323095002 1322673067 1322673075 1322832789 1322489635 132267314
9 1322673128 1322673076 1323084240 1322837822 ...
## $ raw_timestamp_part_2   : int  868349 778725 342967 560311 814776 510661 766645 54671 916313 38
4285 ...
## $ cvtd_timestamp        : Factor w/ 11 levels "02/12/2011 13:33",...: 5 10 10 1 6 11 11 10 3 2 .
..
## $ new_window            : Factor w/ 1 level "no": 1 1 1 1 1 1 1 1 1 1 ...
## $ num_window            : int  74 431 439 194 235 504 485 440 323 664 ...
## $ roll_belt             : num  123 1.02 0.87 125 1.35 -5.92 1.2 0.43 0.93 114 ...
## $ pitch_belt            : num  27 4.87 1.82 -41.6 3.33 1.59 4.44 4.15 6.72 22.4 ...
## $ yaw_belt              : num  -4.75 -88.9 -88.5 162 -88.6 -87.7 -87.3 -88.5 -93.7 -13.1 ...
## $ total_accel_belt      : int  20 4 5 17 3 4 4 4 4 18 ...
## $ kurtosis_roll_belt    : logi  NA NA NA NA NA NA ...
## $ kurtosis_pitch_belt   : logi  NA NA NA NA NA NA ...
## $ kurtosis_yaw_belt     : logi  NA NA NA NA NA NA ...
## $ skewness_roll_belt    : logi  NA NA NA NA NA NA ...
## $ skewness_roll_belt.1  : logi  NA NA NA NA NA NA ...
## $ skewness_yaw_belt     : logi  NA NA NA NA NA NA ...
## $ max_roll_belt         : logi  NA NA NA NA NA NA ...
## $ max_pitch_belt        : logi  NA NA NA NA NA NA ...
## $ max_yaw_belt          : logi  NA NA NA NA NA NA ...
## $ min_roll_belt         : logi  NA NA NA NA NA NA ...
## $ min_pitch_belt        : logi  NA NA NA NA NA NA ...
## $ min_yaw_belt          : logi  NA NA NA NA NA NA ...
## $ amplitude_roll_belt   : logi  NA NA NA NA NA NA ...
## $ amplitude_pitch_belt  : logi  NA NA NA NA NA NA ...
## $ amplitude_yaw_belt    : logi  NA NA NA NA NA NA ...
## $ var_total_accel_belt  : logi  NA NA NA NA NA NA ...
## $ avg_roll_belt         : logi  NA NA NA NA NA NA ...
## $ stddev_roll_belt      : logi  NA NA NA NA NA NA ...
## $ var_roll_belt         : logi  NA NA NA NA NA NA ...
## $ avg_pitch_belt        : logi  NA NA NA NA NA NA ...
## $ stddev_pitch_belt     : logi  NA NA NA NA NA NA ...
## $ var_pitch_belt        : logi  NA NA NA NA NA NA ...
## $ avg_yaw_belt          : logi  NA NA NA NA NA NA ...
## $ stddev_yaw_belt       : logi  NA NA NA NA NA NA ...
## $ var_yaw_belt          : logi  NA NA NA NA NA NA ...
## $ gyros_belt_x          : num  -0.5 -0.06 0.05 0.11 0.03 0.1 -0.06 -0.18 0.1 0.14 ...
## $ gyros_belt_y          : num  -0.02 -0.02 0.02 0.11 0.02 0.05 0 -0.02 0 0.11 ...
## $ gyros_belt_z          : num  -0.46 -0.07 0.03 -0.16 0 -0.13 0 -0.03 -0.02 -0.16 ...
## $ accel_belt_x          : int  -38 -13 1 46 -8 -11 -14 -10 -15 -25 ...
```

```

## $ accel_belt_y      : int  69 11 -1 45 4 -16 2 -2 1 63 ...
## $ accel_belt_z      : int -179 39 49 -156 27 38 35 42 32 -158 ...
## $ magnet_belt_x      : int -13 43 29 169 33 31 50 39 -6 10 ...
## $ magnet_belt_y      : int 581 636 631 608 566 638 622 635 600 601 ...
## $ magnet_belt_z      : int -382 -309 -312 -304 -418 -291 -315 -305 -302 -330 ...
## $ roll_arm           : num  40.7 0 0 -109 76.1 0 0 0 -137 -82.4 ...
## $ pitch_arm          : num -27.8 0 0 55 2.76 0 0 0 11.2 -63.8 ...
## $ yaw_arm            : num  178 0 0 -142 102 0 0 0 -167 -75.3 ...
## $ total_accel_arm     : int  10 38 44 25 29 14 15 22 34 32 ...
## $ var_accel_arm       : logi  NA NA NA NA NA NA ...
## $ avg_roll_arm        : logi  NA NA NA NA NA NA ...
## $ stddev_roll_arm     : logi  NA NA NA NA NA NA ...
## $ var_roll_arm        : logi  NA NA NA NA NA NA ...
## $ avg_pitch_arm       : logi  NA NA NA NA NA NA ...
## $ stddev_pitch_arm    : logi  NA NA NA NA NA NA ...
## $ var_pitch_arm       : logi  NA NA NA NA NA NA ...
## $ avg_yaw_arm         : logi  NA NA NA NA NA NA ...
## $ stddev_yaw_arm      : logi  NA NA NA NA NA NA ...
## $ var_yaw_arm         : logi  NA NA NA NA NA NA ...
## $ gyros_arm_x         : num -1.65 -1.17 2.1 0.22 -1.96 0.02 2.36 -3.71 0.03 0.26 ...
## $ gyros_arm_y         : num  0.48 0.85 -1.36 -0.51 0.79 0.05 -1.01 1.85 -0.02 -0.5 ...
## $ gyros_arm_z         : num -0.18 -0.43 1.13 0.92 -0.54 -0.07 0.89 -0.69 -0.02 0.79 ...
## $ accel_arm_x         : int  16 -290 -341 -238 -197 -26 99 -98 -287 -301 ...
## $ accel_arm_y         : int  38 215 245 -57 200 130 79 175 111 -42 ...
## $ accel_arm_z         : int  93 -90 -87 6 -30 -19 -67 -78 -122 -80 ...
## $ magnet_arm_x        : int -326 -325 -264 -173 -170 396 702 535 -367 -420 ...
## $ magnet_arm_y        : int  385 447 474 257 275 176 15 215 335 294 ...
## $ magnet_arm_z        : int  481 434 413 633 617 516 217 385 520 493 ...
## $ kurtosis_roll_arm   : logi  NA NA NA NA NA NA ...
## $ kurtosis_pitch_arm  : logi  NA NA NA NA NA NA ...
## $ kurtosis_yaw_arm    : logi  NA NA NA NA NA NA ...
## $ skewness_roll_arm   : logi  NA NA NA NA NA NA ...
## $ skewness_pitch_arm  : logi  NA NA NA NA NA NA ...
## $ skewness_yaw_arm    : logi  NA NA NA NA NA NA ...
## $ max_roll_arm        : logi  NA NA NA NA NA NA ...
## $ max_pitch_arm       : logi  NA NA NA NA NA NA ...
## $ max_yaw_arm         : logi  NA NA NA NA NA NA ...
## $ min_roll_arm        : logi  NA NA NA NA NA NA ...
## $ min_pitch_arm       : logi  NA NA NA NA NA NA ...
## $ min_yaw_arm         : logi  NA NA NA NA NA NA ...
## $ amplitude_roll_arm  : logi  NA NA NA NA NA NA ...
## $ amplitude_pitch_arm : logi  NA NA NA NA NA NA ...
## $ amplitude_yaw_arm   : logi  NA NA NA NA NA NA ...
## $ roll_dumbbell       : num -17.7 54.5 57.1 43.1 -101.4 ...
## $ pitch_dumbbell      : num  25 -53.7 -51.4 -30 -53.4 ...
## $ yaw_dumbbell        : num  126.2 -75.5 -75.2 -103.3 -14.2 ...
## $ kurtosis_roll_dumbbell : logi  NA NA NA NA NA NA ...

```

```
## $ kurtosis_picth_dumbbell : logi NA NA NA NA NA NA ...
## $ kurtosis_yaw_dumbbell   : logi NA NA NA NA NA NA ...
## $ skewness_roll_dumbbell  : logi NA NA NA NA NA NA ...
## $ skewness_pitch_dumbbell : logi NA NA NA NA NA NA ...
## $ skewness_yaw_dumbbell   : logi NA NA NA NA NA NA ...
## $ max_roll_dumbbell       : logi NA NA NA NA NA NA ...
## $ max_picth_dumbbell      : logi NA NA NA NA NA NA ...
## $ max_yaw_dumbbell        : logi NA NA NA NA NA NA ...
## $ min_roll_dumbbell       : logi NA NA NA NA NA NA ...
## $ min_pitch_dumbbell      : logi NA NA NA NA NA NA ...
## $ min_yaw_dumbbell        : logi NA NA NA NA NA NA ...
## $ amplitude_roll_dumbbell : logi NA NA NA NA NA NA ...
## [list output truncated]
```

Cleaning the data (part 1)

This dataset has 160 variables - way too many to build a model on. Inspection of the test dataset reveals that several variables can be removed from the training set because the corresponding values to not exist in the test set. This reduces the training set to 60 variables, which helps, but is still not enough. At this point, I chose to split the training dataset.

```
traindata <- traindata[,as.vector(apply(apply(sapply(testdata, as.numeric),2,is.na),2,sum))==0]
```

Splitting the data

To set up cross-validation, as well as a final validation set, I split the initial training set into a validation set (to be used after building the model) and 3 other training/“test” sets (for k-fold cross-validation). Before splitting the non-validation data into 3 folds, I performed the following data cleaning on that data.

```
library(caret)
```

```
## Warning: package 'caret' was built under R version 3.1.2
```

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

```

set.seed(17)
inTrain <- createDataPartition(y=traindata$classe, p=0.7, list=FALSE)
sampletrain <- traindata[inTrain,]
testtrain <- traindata[-inTrain,]
folds <- createFolds(y=sampletrain$classe,k=3,
                     list=TRUE,returnTrain=FALSE)
sample1 <- sampletrain[-folds$Fold1,]
sample2 <- sampletrain[-folds$Fold2,]
sample3 <- sampletrain[-folds$Fold3,]
cv1 <- sampletrain[folds$Fold1,]
cv2 <- sampletrain[folds$Fold2,]
cv3 <- sampletrain[folds$Fold3,]

```

Cleaning the data (part 2)

```
names(sampletrain)
```

```

## [1] "X"                "user_name"        "raw_timestamp_part_1"
## [4] "raw_timestamp_part_2" "cvtd_timestamp"   "new_window"
## [7] "num_window"        "roll_belt"        "pitch_belt"
## [10] "yaw_belt"          "total_accel_belt" "gyros_belt_x"
## [13] "gyros_belt_y"      "gyros_belt_z"     "accel_belt_x"
## [16] "accel_belt_y"      "accel_belt_z"     "magnet_belt_x"
## [19] "magnet_belt_y"     "magnet_belt_z"    "roll_arm"
## [22] "pitch_arm"         "yaw_arm"          "total_accel_arm"
## [25] "gyros_arm_x"       "gyros_arm_y"      "gyros_arm_z"
## [28] "accel_arm_x"       "accel_arm_y"      "accel_arm_z"
## [31] "magnet_arm_x"      "magnet_arm_y"     "magnet_arm_z"
## [34] "roll_dumbbell"     "pitch_dumbbell"   "yaw_dumbbell"
## [37] "total_accel_dumbbell" "gyros_dumbbell_x" "gyros_dumbbell_y"
## [40] "gyros_dumbbell_z"  "accel_dumbbell_x" "accel_dumbbell_y"
## [43] "accel_dumbbell_z"  "magnet_dumbbell_x" "magnet_dumbbell_y"
## [46] "magnet_dumbbell_z" "roll_forearm"     "pitch_forearm"
## [49] "yaw_forearm"       "total_accel_forearm" "gyros_forearm_x"
## [52] "gyros_forearm_y"   "gyros_forearm_z"   "accel_forearm_x"
## [55] "accel_forearm_y"   "accel_forearm_z"   "magnet_forearm_x"
## [58] "magnet_forearm_y"  "magnet_forearm_z"  "classe"

```

The first remaining 7 variables are not related to the sensors like the following 52, so I chose to remove them. The 60th variable is “classe,” which is what we are trying to predict. (In hindsight, I could have left in the user for this exercise, but chose to remove it in favor of extending the prediction capabilities to users whom the model was not trained on. Also, adding this variable to my final model did not improve the accuracy.)

```
sampletrain_new <- sampletrain[,8:59]
```

Building the model

With the remaining 52 variables, I performed a near zero variance test to determine which ones could be removed. It turned out that none of them could, but I chose to build the model on the variables whose unique frequency was greater than 10. Having now narrowed down the predictive variables to 14, it was time to determine if these were the right ones.

```
s <- apply(sampletrain, 2, as.numeric)
s <- as.data.frame(s)
nzv <- nearZeroVar(s, saveMetrics = TRUE)
rownames(nzv[nzv$percentUnique > 10, ])
```

```
## [1] "pitch_belt"      "yaw_belt"        "roll_arm"
## [4] "pitch_arm"       "yaw_arm"         "roll_dumbbell"
## [7] "pitch_dumbbell"  "yaw_dumbbell"    "roll_forearm"
## [10] "pitch_forearm"   "yaw_forearm"     "magnet_forearm_x"
## [13] "magnet_forearm_y" "magnet_forearm_z"
```

The names of these variables actually makes sense - using the yaw, pitch, and roll from each of the four body sensors (except for the roll from the belt) as well as the 3-axis magnet measurements from the forearm. This way, full body motion and all 3 rotational axis dimensions are utilized, with an added finesse including the additional forearm measurement.

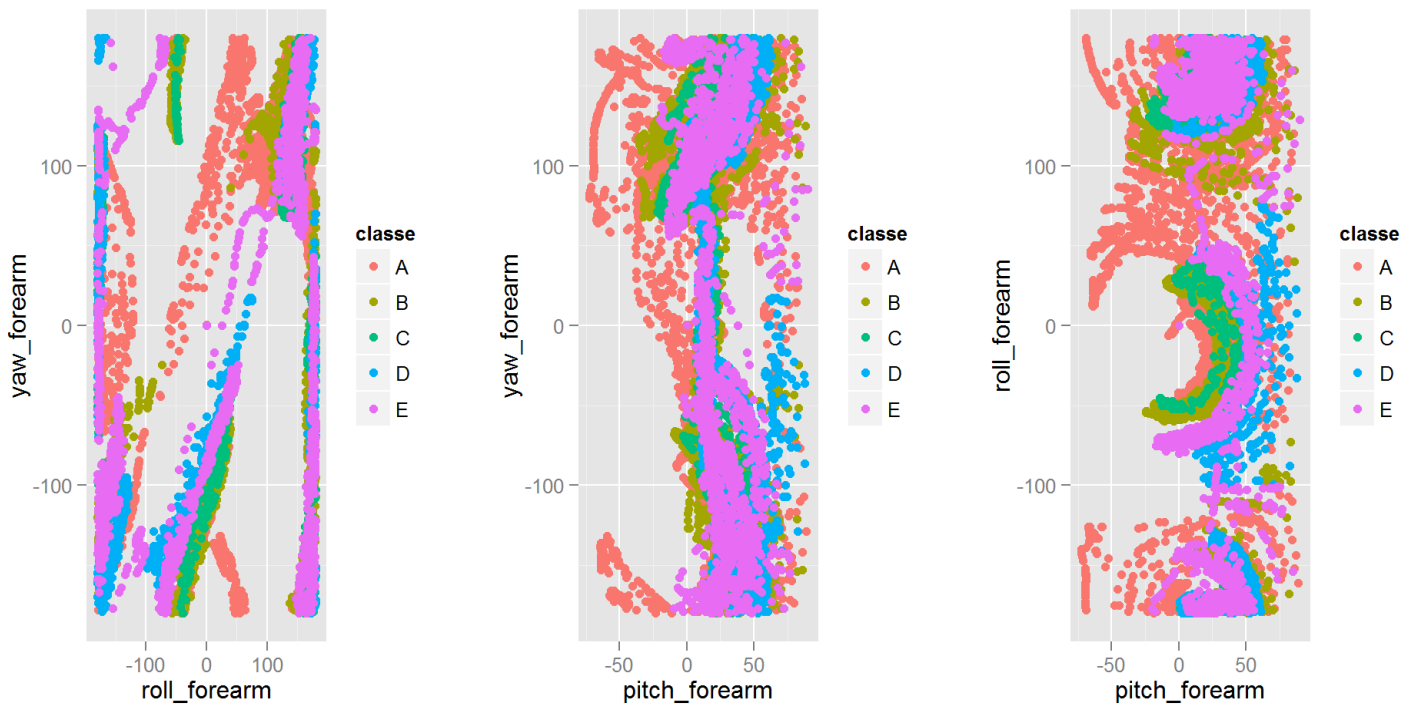
Of course, it is possible to perform a correlation analysis on the 52 variables, or to call the “pairs” function, but this would be quite a large, unreadable plot. Instead, I plotted many of the variables against each other, using color to show the “classe” variable, and the results were quite intriguing, showing that each variable contributes some unique information or pattern. Here, for example, I plotted yaw, roll, and pitch for the forearm sensor against each other. I invite anyone to plot the others as well, using a similar format.

```
plot1 <- qplot(roll_forearm, yaw_forearm, color = classe, data = sampletrain)
plot2 <- qplot(pitch_forearm, yaw_forearm, color = classe, data = sampletrain)
plot3 <- qplot(pitch_forearm, roll_forearm, color = classe, data = sampletrain)
library(gridExtra)
```

```
## Warning: package 'gridExtra' was built under R version 3.1.2
```

```
## Loading required package: grid
```

```
grid.arrange(plot1, plot2, plot3, ncol=3, nrow=1)
```

Random forest was used because it tends to be one of the most accurate methods, and the interpretability is not as critical here because the individual measurements themselves are difficult to understand. The model was trained using each of the 3 folds, and the accuracy for each was measured against the respective cross-validation set and the training set.

```
library(randomForest)
modFit1 <- train(classe ~ yaw_belt + pitch_belt +
  roll_arm + yaw_arm + pitch_arm +
  roll_dumbbell + yaw_dumbbell + pitch_dumbbell +
  roll_forearm + yaw_forearm + pitch_forearm +
  magnet_forearm_x + magnet_forearm_y + magnet_forearm_z,
  method = "rf", data = sample1)
confusionMatrix(sample1$classe, predict(modFit1, newdata=sample1))$overall[1]
```

```
## Accuracy
##      1
```

```
confusionMatrix(cv1$classe, predict(modFit1, newdata=cv1))$overall[1]
```

```
## Accuracy
## 0.9814
```

```
modFit2 <- train(classe ~ yaw_belt + pitch_belt +
                roll_arm + yaw_arm + pitch_arm +
                roll_dumbbell + yaw_dumbbell + pitch_dumbbell +
                roll_forearm + yaw_forearm + pitch_forearm +
                magnet_forearm_x + magnet_forearm_y + magnet_forearm_z,
                method = "rf", data = sample2)
confusionMatrix(sample2$classe, predict(modFit2, newdata=sample2))$overall[1]
```

```
## Accuracy
##      1
```

```
confusionMatrix(cv2$classe, predict(modFit2, newdata=cv2))$overall[1]
```

```
## Accuracy
##    0.9827
```

```
modFit3 <- train(classe ~ yaw_belt + pitch_belt +
                roll_arm + yaw_arm + pitch_arm +
                roll_dumbbell + yaw_dumbbell + pitch_dumbbell +
                roll_forearm + yaw_forearm + pitch_forearm +
                magnet_forearm_x + magnet_forearm_y + magnet_forearm_z,
                method = "rf", data = sample3)
confusionMatrix(sample3$classe, predict(modFit3, newdata=sample3))$overall[1]
```

```
## Accuracy
##      1
```

```
confusionMatrix(cv3$classe, predict(modFit3, newdata=cv3))$overall[1]
```

```
## Accuracy
##    0.9793
```

Model selection and results

```
mean(confusionMatrix(cv1$classe, predict(modFit1, newdata=cv1))$overall[1],
     confusionMatrix(cv2$classe, predict(modFit2, newdata=cv2))$overall[1],
     confusionMatrix(cv3$classe, predict(modFit3, newdata=cv3))$overall[1])
```

```
## Loading required package: 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.
```

```
## [1] 0.9814
```

All 3 models performed well, with the 2nd model slightly outperforming the other two, and an average accuracy of 98.1%. Interestingly enough, all 3 models resulted in an accuracy of 100% on their own training sets, which may be a result of overfitting. Still, applying the best model to the validation set results in an accuracy of 97.8%, which is similar to the average accuracy, so we expect the out-of-sample accuracy (one of measures of out-of-sample error) to be somewhere in this range (98%).

```
confusionMatrix(testtrain$classe,predict(modFit2,newdata=testtrain))
```

Confusion Matrix and Statistics

##

Reference

Prediction A B C D E

A 1665 7 0 2 0

B 20 1099 16 2 2

C 0 15 995 13 3

D 1 2 11 948 2

E 0 9 7 17 1049

##

Overall Statistics

##

Accuracy : 0.978

95% CI : (0.974, 0.982)

No Information Rate : 0.286

P-Value [Acc > NIR] : <2e-16

##

Kappa : 0.972

McNemar's Test P-Value : NA

##

Statistics by Class:

##

Class: A Class: B Class: C Class: D Class: E

Sensitivity 0.988 0.971 0.967 0.965 0.993

Specificity 0.998 0.992 0.994 0.997 0.993

Pos Pred Value 0.995 0.965 0.970 0.983 0.970

Neg Pred Value 0.995 0.993 0.993 0.993 0.999

Prevalence 0.286 0.192 0.175 0.167 0.179

Detection Rate 0.283 0.187 0.169 0.161 0.178

Detection Prevalence 0.284 0.194 0.174 0.164 0.184

Balanced Accuracy 0.993 0.981 0.980 0.981 0.993