

But How Well Are You Doing It? Predicting from Fitbit Data

Travis Fell

November 15, 2015

I. Introduction

This analysis predicts how well test subjects would perform a dumbbell bicep curl using data from personal movement tracking devices like Nike FuelBand and Fitbits.

Using a training set of measurements from these personal exercise movement devices like Nike FuelBands and FitBits as well as assessments from professional trainers on the correctness of exercises performed, this analysis will predict how well 20 test subjects performed the same exercise. Exercise performance is categorized by the classe column in the training set using the following levels:

- A: Exactly according to the specification
- B: Throwing the elbows to the front
- C: Lifting the dumbbell only halfway
- D: Lowering the dumbbell only halfway
- E: Throwing the hips to the front

Read more here: <http://groupware.les.inf.puc-rio.br/har> (<http://groupware.les.inf.puc-rio.br/har>)

II. Set Up Work

First, we'll start out by loading libraries, loading data and doing a bit of exploratory analysis.

```
setwd("C:/Users/fellt/Desktop/Data Science/Coursera Data Science Specialization/08 - Practical Machine Learning/8-PracticalMachineLearning")
set.seed(333)
library(caret)
library(randomForest)
library(ggplot2)
library(plyr) #need this before dplyr for GBM model
library(dplyr)
library(parallel, quietly=T)
library(doParallel, quietly=T)
setInternet2(use = TRUE)
download.file("https://d396qusza40orc.cloudfront.net/predmachlearn/pml-training.csv", destfile = "pml-training.csv")
download.file("https://d396qusza40orc.cloudfront.net/predmachlearn/pml-testing.csv", destfile = "pml-testing.csv")
train <- read.csv("pml-training.csv", stringsAsFactors = FALSE, na.strings=c("", " ", "N A"))
test <- read.csv("pml-testing.csv", stringsAsFactors = FALSE, na.strings=c("", " ", "N A"))
summary(train)
str(train)
View(train)
```

This data set has a lot of extraneous data. The first 7 columns have no predictive value with respect to performing an exercise. Furthermore, there are a LOT of columns with little to no data.

III. Preparing the Data Set

We need tidy up this data by squeezing out columns that are not useful for predicting and removing rows with missing data.

```

#remove first 7 columns from both data sets as they are only informational about the subjects and test
train <- train[,8:160]
test <- test[,8:160]

# filter out columns with majority NA values
train <- train[sapply(train, function(x) !any(is.na(x)))]
test <- test[sapply(test, function(x) !any(is.na(x)))]

# find columns with more than 15% of rows = 0 and remove
zerotest <- NULL
for(i in 1:length(train))
{
  zerotest <- c(zerotest, nrow(train[train[,i] == 0,])/nrow(train))
}
removecols <- colnames(train[,which(zerotest > .15)])
train[,removecols] <- list(NULL)
test[,removecols] <- list(NULL)

```

In addition to removing columns with missing data, we are also going to find pairs of columns that correlate with each other and remove one from the data set. Then, we'll drop records missing data.

```

#find correlations in remaining features to ID add'l variables to exclude
traincor <- findCorrelation(cor(train[,1:length(train)-1]), cutoff = .8, verbose = FALSE)
#Leave off dependent variable for this command
train <- train[,-traincor] #exclude 1 feature from each pair of highly correlated features from training set
test <- test[, -traincor] #exclude same features from test set

# find records with any 0 values and remove
train[train == 0] <- NA
train <- train[complete.cases(train),]
train$classe <- as.factor(train$classe)

```

The processing above yields 33 columns with predictive value in both the training and test set.

III. Analyze the Data

Now that we have a nice and tight data set, we'll perform some cross-validation on the training set which will enable us to estimate out of sample error on our models.

```

# split the training set into 70/30 training-validation sets
inTrain <- createDataPartition(train$classe, p = .7, list = FALSE)
trainMod <- train[inTrain,]
validMod <- train[-inTrain,]

```

Now, for the fun stuff. We will use the updated training set, `trainMod`, to build some models and the validation set, `validMod`, to estimate the out of sample error. We'll start with a generalized boosting model.

```
# turn on parallel processing to help improve performance
cluster <- makeCluster(detectCores() - 1)
registerDoParallel(cluster)

#create boosting model
load("boostingMod.rda") #to enable knitr to run quickly, will simply load up saved GBM model
#trainGBM <- train(classe ~ ., method = "gbm", data = trainMod, verbose = FALSE) #note: commented out this line to enable knitr to run in a timely manner
#save(trainGBM, file = "boostingMod.rda") #since this takes a while to run, let's save the model to quickly load and reuse later
```

Now that the model is created, let's run against the validation set and see the results. To estimate the OOS accuracy rate, we'll sum the number of prediction values that match the validation set then divide by the total number of prediction values. To find the estimated OOS error rate, we'll subtract the estimated accuracy rate from 1.

```
#perform cross validation and estimate OOS error on boosting model
validationPredGBM <- predict(trainGBM, validMod)
gbmOOSAccuracy <- sum(validationPredGBM == validMod$classe)/length(validationPredGBM)
gbmOOSError <- 1- gbmOOSAccuracy
gbmOOSAccuracy <- paste(round((gbmOOSAccuracy) * 100, digits = 1), "%", sep="")
gbmOOSError <- paste(round((gbmOOSError) * 100, digits = 1), "%", sep="")
```

The estimated OOS accuracy is 95.3% and the estimated OOS error is 4.7%. On its face, these are so-so results. Can we do better with a random forest model?

```
#create random forest model
trainRF <- randomForest(classe~., trainMod)
save(trainRF, file = "rfMod.rda")

#perform cross validation and estimate OOS error on random forest model
validationPredRF <- predict(trainRF, validMod)
rfOOSAccuracy <- sum(validationPredRF == validMod$classe)/length(validationPredRF)
rfOOSError <- 1- rfOOSAccuracy
rfOOSAccuracy <- paste(round((rfOOSAccuracy) * 100, digits = 1), "%", sep="")
rfOOSError <- paste(round((rfOOSError) * 100, digits = 1), "%", sep="")
```

Bam! Using the same approach as above, the Random Forest model gives us an accuracy rate of 98.5% and an OOS error rate of 1.5%. Let's use this model to run against the test set and produce files for the final assignment submission.

IV. Predict on the Test Set and Submit

Here we'll put the test set through the RF prediction model above and create the files for the assignment submission.

```
#make predictions on test data, save to file for submission
testPredRF <- predict(trainRF, test)

# create answers vector with values from test set prediction
testPredRF <- as.character(testPredRF)

# write to file
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(testPredRF)
```

The results returned 19 of 20 correct, which seems reasonable given the estimated accuracy rate above.