

Prediction Motivation

INTRODUCTION

Fitness band devices track exercise activities of the wearers. A weight lifting project conducted by Velloso, Bulling, Gellersen, Ugulino and Fuks in 2013 evaluated the "how well" an activity was performed by the device wearer. Their work focused on "the problem of specifying correct execution, the automatic and robust detection of execution mistakes, and how to provide feedback on the quality of execution to the user."

The data set used in this prediction project were collected from 6 young health male (20-28 years old) participants wearing accelerometers on their belt, forearm, arm and dumbbell. Because all the participants had little weight lifting experience, they presumably could easily simulate mistakes made in weight lifting exercises.

DATA SOURCE

The data from this project came from: <http://groupware.les.inf.puc-rio.br/har>. Both training and test data are provided by the course and downloaded from:

<https://d396qusza40orc.cloudfront.net/predmachlearn/pml-training.csv>

<https://d396qusza40orc.cloudfront.net/predmachlearn/pml-testing.csv>

The training data has 19,622 rows and 160 columns. The test data has 20 rows and 160 columns. Based on the structure and summaries of the data, there were a number of column that are predominantly populated with "NA" and blanks. Some also have "#DIV/0!" entries. The following code was used to pre-process the data:

```
# Read in data sets and define NA strings
# train = read.csv(file = "pml-training.csv", header=TRUE, na.strings=c("NA", "", "#DIV/0!"))
# test = read.csv(file = "pml-testing.csv", header=TRUE, na.strings=c("NA", "", "#DIV/0!"))

# Use str, describe and summary functions to check data for unique values
# Combine the two sets in preparation for data cleaning
# merge1 = train; merge2 = test
# merge1$type = "train"; merge2$type = "test"
# merge1$classe = NULL; merge2$problem_id = NULL
# merge = rbind(merge1, merge2)

# Delete columns that are predominantly "NA"
# mergeAll = merge[,complete.cases(t(merge))]

# Split the data back into training and testing sets
# training = subset(mergeAll, type == "train")
# testing = subset(mergeAll, type == "test")
# training$classe = train$classe
# training$type = NULL; testing$type = NULL
# testing$problem_id = test$problem_id
```

DATA SPLITTING

In preparation to the model development, the training data set was split 60:40 for the train and test sets, respectively. There were 11,776 rows in the train set and 7,846 rows on the test set. Both have the number of predictors reduced to 60.

```
# Create training, cross-validation and testing sets for model building
# set.seed(12345)
# fitIndex = createDataPartition(training$classe, p = 0.60, list=FALSE)
# fitTrain = training[fitIndex,]
# fitTest = training[-fitIndex,]
```

MODEL BUILDING

The initial baseline accuracy was calculated in the test set and is shown below:

```
# Baseline Accuracy
# table(fitTest$classe)
#   A   B   C   D   E
# 2232 1518 1368 1286 1442
```

Random forest was chosen as the predictive model. Because random forest work by building a large collections of trees, (where each tree is split on a random subset of the available independent variabel and built from a "bagged or bootstrapped" sample of the data), interpretability is hard but accuracy is increased over the predictiver models of classification and regression trees. The first model built was using all the predictors and resulted into an accuracy of 100%, which implied overfitting.

A second model was built with all the predictors but excluding the variables X and user_name. These 2 variables are deemed unnecessary as they are probably identifiers. The accuracy of the model decreased to 0.998598. Again, the model predictions indicate overfitting. A third model was built with 3 additional predictors removed. Shown below are the code and the results of this third mode with an accuracy of 0.9961764.

```
# set.seed(1)
# model2 = randomForest(classe~. -X -user_name -raw_timestamp_part_1 -raw_timestamp_part_2
#                         -cvtd_timestamp-new_window, data = fitTrain)
# model2P = predict(model2, newdata=fitTest)

# table(fitTest$classe,model2P)
#   model2P
#       A   B   C   D   E
# A 2232   0   0   0   0
# B   4 1511   3   0   0
# C   0   3 1364   1   0
# D   0   0  14 1270   2
# E   0   0   0   3 1439

# confusionMatrix(model2P, fitTest$classe)$overall["Accuracy"]
# Accuracy
# 0.9961764
```

PREDICTION METRICS

Some of the metrics that can be used to help identify the important predictors are: 1) the number of times, aggregated over all the trees in the random forest model, that a certain variable is selected for a split (Figure 1) and 2) average reduction in impurity, taken over all the times that variable is selected for splitting in all trees in the forest (Figure 2). Figure 1 shows that the following variables are the most important in terms of the number of splits: num_window, yaw_belt, pitch_belt, roll_belt, magnet_dumbell in x,y,z directions. Figure 2 shows that num_window was also the most important variable in terms of mean reduction in impurity and the other 5 in Figure 1 plus pitch_forearm.

SUMMARY

The predictive model is able to accurately predict the type of exercise. This model uses 54 predictors and could probably use less if the predictive model is combined with another predictive model to identify the best combination of predictors.

REFERENCE

Velloso, E; Bulling, A.; Gellersen, H.; Ugulin W., W; Fuks, H. Qualitative Activity Recognition of Weight Lifting Exercises. Proceedings of 4th International Conference in Cooperation with SIGCHI (Augmented Human '13), Stuttgart, Germany: ACM SIGCHI, 2013

Figure 1. Most important variables in terms of the number of splits.

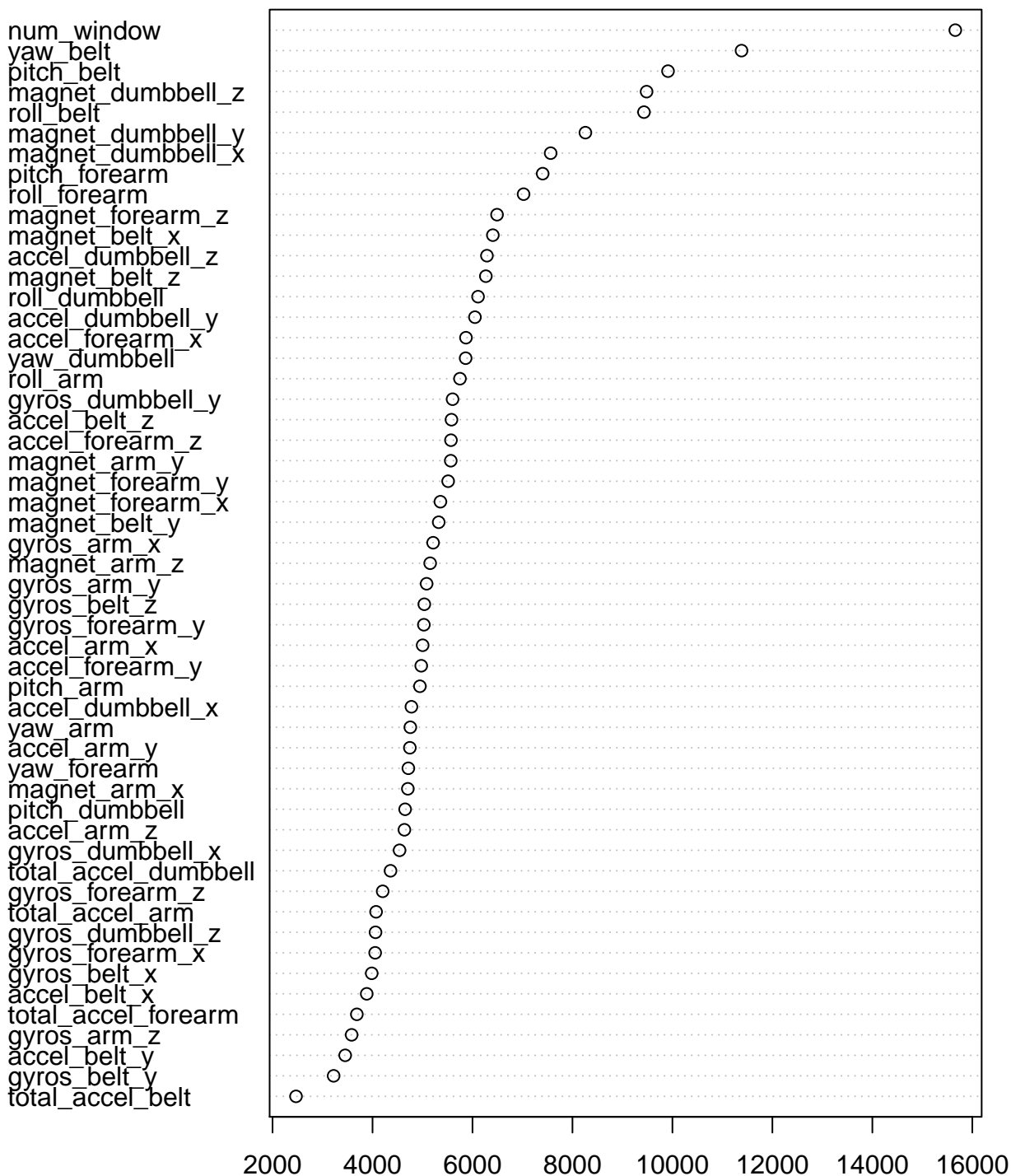


Figure 2. Most important variables in terms of mean reduction in impurity.

model2

num_window
roll_belt
yaw_belt
magnet_dumbbell_z
pitch_forearm
pitch_belt
magnet_dumbbell_y
roll_forearm
magnet_dumbbell_x
accel_belt_z
accel_dumbbell_y
magnet_belt_z
roll_dumbbell
magnet_belt_y
accel_dumbbell_z
accel_forearm_x
roll_arm
gyros_belt_z
accel_dumbbell_x
yaw_dumbbell
total_accel_dumbbell
magnet_forearm_z
magnet_belt_x
magnet_arm_x
accel_forearm_z
accel_arm_x
gyros_dumbbell_y
magnet_arm_y
total_accel_belt
yaw_arm

