

Practical Machine Learning Assignment

Darragh McLernon

29/01/2017

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.

The purpose of this assignment is to use data collected from the belt, forearm, arm and dumbbell of six participants while they are exercising, and classify the data into different “classe”. The final result should be an indication into how the participants are performing the exercise.

Data Analysis

The data for this assignment was kindly provided by <http://groupware.les.inf.puc-rio.br/har>. There are two data sets provided, a training data set and a test data set. The training data will be loaded and analysed. The test data will be used to verify the training model so it will be loaded but will remain untouched for the analysis. NA values will be removed from the training set.

```
test_data <- read.csv("./pml-testing.csv")
training_data <- read.csv("./pml-training.csv", na.strings = c("", "NA"))
dim(training_data)
```

```
## [1] 19622 160
```

The training data contains 19622 observations across 160 variables. At a glance there appears to be many variables containing mostly NA values, these will be removed. See appendix for a sample of the data. The next step is to remove variables that are not needed, starting with the ones mostly containing NA values.

```
#Find the colums with full sets of data (columns that contain data in every row)
is_data <- apply(!is.na(training_data), 2, sum) > 19621
```

```
#Create a new data set with only the colums with a full data set
training_set <- training_data[, is_data]

dim(training_set)
```

```
## [1] 19622 60
```

This leaves us with a training set with 19622 observations across 60 variables. There are further columns we can remove which will not be used to build the model.

```
columns_to_remove = c('X', 'user_name', 'raw_timestamp_part_1', 'raw_timestamp_part_2', 'cvtd_timestamp')

training_set <- training_set[, -which(names(training_set) %in% columns_to_remove)]

dim(training_set)
```

```
## [1] 19622    53
```

Our final training data set contains 19662 observations across 53 variables.

Training data subsampling and predictive analysis

In order to create a predictive model it is helpful to create a subsample of the training data. This will allow the model to be trained without using the testing set, which is best practice. The training set will be split into two smaller sets named, `sub_train` and `sub_test`. The `caret` package is required to accomplish this.

```
library(caret)
```

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

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

```
sub_sample <- createDataPartition(y=training_set$classe, p=0.6, list=FALSE)
sub_train <- training_set[sub_sample,]
sub_test <- training_set[-sub_sample, ]
dim(sub_train);dim(sub_test)
```

```
## [1] 11776    53
```

```
## [1] 7846    53
```

The sub-samples results in 11776 and 7846 observations respectively.

Random Forest

One method of analysis is to perform a random forest analysis. This method creates a number of decision trees to identify a classification. All of the trees classification weights are culumated to identify the most likely classification for each data set. The model is built below;

```
set.seed(12345)
```

```
model_control <- trainControl(method="cv", number=3, verboseIter=FALSE)
```

```
train_model_forest <- train(classe ~ ., data=sub_train, method="rf", trControl=model_control)
```

```
## Loading required package: randomForest
```

```
## 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
```

```
train_model_forest
```

```
## Random Forest
```

```
##
```

```
## 11776 samples
```

```
##      52 predictor
```

```
##      5 classes: 'A', 'B', 'C', 'D', 'E'
```

```
##
```

```
## No pre-processing
## Resampling: Cross-Validated (3 fold)
## Summary of sample sizes: 7852, 7850, 7850
## Resampling results across tuning parameters:
##
##      mtry  Accuracy   Kappa
##      2    0.9838652 0.9795834
##     27    0.9844599 0.9803404
##     52    0.9811479 0.9761482
##
## Accuracy was used to select the optimal model using the largest value.
## The final value used for the model was mtry = 27.
```

Applying the model to the test data set gives the predictions below;

```
forest_prediction <- predict(train_model_forest, newdata=sub_test)
pred_matrix <- confusionMatrix(forest_prediction, sub_test$classe)
pred_matrix
```

```
## Confusion Matrix and Statistics
##
##              Reference
## Prediction      A      B      C      D      E
##           A 2229    11      0      0      0
##           B      1 1504      7      0      1
##           C      1      3 1352     14      0
##           D      0      0      9 1270      0
##           E      1      0      0      2 1441
##
## Overall Statistics
##
##              Accuracy : 0.9936
##              95% CI : (0.9916, 0.9953)
##      No Information Rate : 0.2845
##      P-Value [Acc > NIR] : < 2.2e-16
##
##              Kappa : 0.9919
##      McNemar's Test P-Value : NA
##
## Statistics by Class:
##
##              Class: A Class: B Class: C Class: D Class: E
## Sensitivity          0.9987  0.9908  0.9883  0.9876  0.9993
## Specificity          0.9980  0.9986  0.9972  0.9986  0.9995
## Pos Pred Value       0.9951  0.9941  0.9869  0.9930  0.9979
## Neg Pred Value       0.9995  0.9978  0.9975  0.9976  0.9998
## Prevalence           0.2845  0.1935  0.1744  0.1639  0.1838
## Detection Rate       0.2841  0.1917  0.1723  0.1619  0.1837
## Detection Prevalence 0.2855  0.1928  0.1746  0.1630  0.1840
## Balanced Accuracy    0.9983  0.9947  0.9928  0.9931  0.9994
```

Judging from the confusion matrix it is apparent that this prediction model is very accurate. The highest number of miss-classifications is 14 for the 'D' classe and the overall accuracy is 0.9936. This model will be used to predict the classe of the original test data.

Predicting the original test data

The final step of this assignment is to predict the classe of the origin test data. The model created from the random forest will be used to accomplish this. Below is the predicted classe of each of the 20 test observations.

```
classification_prediction <- predict(train_model_forest,test_data,type='raw')
```

```
## Loading required package: randomForest
## 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
```

```
classification_prediction
```

```
## [1] B A B A A E D B A A B C B A E E A B B B
## Levels: A B C D E
```

References

Velloso, E.; Bulling, A.; Gellersen, H.; Ugulino, 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.

Appendix

Training data sample - untouched

```
test_data_sample <- read.csv("./pml-testing.csv")
```

```
str(test_data_sample)
```

```
## '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 1322673149 ...
## $ raw_timestamp_part_2 : int  868349 778725 342967 560311 814776 510661 766645 54671 916313 3842 ...
## $ 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 ...
## $ 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_pitch_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_pitch_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]

```