

Practical Machine Learning - Course Project

Paul Stables

21/09/2020

Background

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 dumbbell 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).

Project

The goal of this project is to build a model using a given training set of data that predicts what exercise is being performed.

Model

Load Environment and Data

I first load the caret package which will be used to build the model, read the data into training and testing variables, and produce a summary of the columns to investigate what the data contains.

```
library(caret)
```

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

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

```
training <- read.csv("pml-training.csv")
testing <- read.csv("pml-testing.csv")
str(training)
```

```
## '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/ 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_pitch_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_pitch_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_pitch_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_pitch_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]
```

Clean Data

The first stage in cleaning the data is to remove any variable whose variance is near 0. The summary of the data also shows that there are some variables that are mostly NA. These will not be useful for prediction purposes, so should also be removed.

```
nzv <- nearZeroVar(training)
training <- training[, -nzv]
testing <- testing[, -nzv]
```

```

mostlyNA <- sapply(training, function(x) mean(is.na(x))) > 0.95
training <- training[, mostlyNA==F]
testing <- testing[, mostlyNA==F]

```

Looking at the remaining variables, there are also a few that shouldn't be predictive (e.g. user name or time stamps). These are also removed.

```

training <- training[, -(1:5)]
testing <- testing[, -(1:5)]

```

Train the Model

I have chosen to use a Random Forest Model thanks to their high level of accuracy. I fit the model using the training data, including all variables as predictors.

```

controlRF <- trainControl(method="cv", number=3, verboseIter=FALSE)
forestMod <- train(classe ~ . , method = "rf", data=training, trControl=controlRF)

```

Accuracy of model

After completing the model, I then check the accuracy of the prediction based on the training set versus the actual classifications and find that the model has a good level of accuracy.

```

trainingPrediction <- predict(forestMod, training)
confusionMatrix(trainingPrediction, training$classe)

```

```

## Confusion Matrix and Statistics
##
##              Reference
## Prediction    A    B    C    D    E
##      A 5580     0     0     0     0
##      B     0 3797     0     0     0
##      C     0     0 3422     0     0
##      D     0     0     0 3216     0
##      E     0     0     0     0 3607
##
## Overall Statistics
##
##              Accuracy : 1
##              95% CI : (0.9998, 1)
##      No Information Rate : 0.2844
##      P-Value [Acc > NIR] : < 2.2e-16
##
##              Kappa : 1
##
##      McNemar's Test P-Value : NA
##
## Statistics by Class:
##
##              Class: A Class: B Class: C Class: D Class: E

```

## Sensitivity	1.0000	1.0000	1.0000	1.0000	1.0000
## Specificity	1.0000	1.0000	1.0000	1.0000	1.0000
## Pos Pred Value	1.0000	1.0000	1.0000	1.0000	1.0000
## Neg Pred Value	1.0000	1.0000	1.0000	1.0000	1.0000
## Prevalence	0.2844	0.1935	0.1744	0.1639	0.1838
## Detection Rate	0.2844	0.1935	0.1744	0.1639	0.1838
## Detection Prevalence	0.2844	0.1935	0.1744	0.1639	0.1838
## Balanced Accuracy	1.0000	1.0000	1.0000	1.0000	1.0000

Results

Applying the Random Forest model to the test data gives the required predictions.

```
testingPrediction <- predict(forestMod, newdata = testing)
testingPrediction
```

```
## [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
```