

Assignment: Machine learning

Jeroen van der A

November 13th, 2016

Executive summary

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, the goal is 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://groupware.les.inf.puc-rio.br/har>.

Six young health participants were asked to perform one set of 10 repetitions of the Unilateral Dumbbell Biceps Curl in five different fashions: exactly according to the specification (Class A), throwing the elbows to the front (Class B), lifting the dumbbell only halfway (Class C), lowering the dumbbell only halfway (Class D) and throwing the hips to the front (Class E).

Read more: <http://groupware.les.inf.puc-rio.br/har#ixzz4PvAXy4zU>

In our analysis we want to use the data from the from the devices to accurately predict which repetitions were performed correctly and if not What mistake was made. will use data from accelerometers on the belt, forearm, arm, and dumbbell of 6 participants to predict the manner in which they did the exercise.

Data Preprocessing

```
library(caret)
```

```
## Warning: package 'caret' was built under R version 3.2.5
```

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

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

```
## Warning: package 'ggplot2' was built under R version 3.2.4
```

```
library(randomForest)
```

```
## Warning: package 'randomForest' was built under R version 3.2.5
```

```
## 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
```

```
...
```

Read the Data After downloading the data from the data source, we can read the two csv files into two data frames.

```
ruwedata <- read.csv("http://d396qusza40orc.cloudfront.net/predmachlearn/pml-  
training.csv", sep=",", na.strings=c("", "NA", "#DIV/0!", "<NA>"))
```

```
testruw <- read.csv("http://d396qusza40orc.cloudfront.net/predmachlearn/pml-  
testing.csv", sep=",", na.strings=c("", "NA", "#DIV/0!", "<NA>"))
```

```
dim(ruwedata)
```

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

```
dim(testruw)
```

```
## [1] 20 160
```

The training data(trainingdata) has 19622 observations with 160 features. The test data (testingdata) has 20 observations with 160 features.

Clean the data

In this step, we will clean the data and get rid of observations with missing values as well as some meaningless variables.

```
sum(complete.cases(ruwedata))
```

```
## [1] 0
```

```
head(ruwedata)
```

```
##   X user_name raw_timestamp_part_1 raw_timestamp_part_2 cvtd_timestamp  
## 1 1 carlitos      1323084231      788290 05/12/2011 11:23  
## 2 2 carlitos      1323084231      808298 05/12/2011 11:23  
## 3 3 carlitos      1323084231      820366 05/12/2011 11:23  
## 4 4 carlitos      1323084232      120339 05/12/2011 11:23  
## 5 5 carlitos      1323084232      196328 05/12/2011 11:23  
## 6 6 carlitos      1323084232      304277 05/12/2011 11:23  
##   new_window num_window roll_belt pitch_belt yaw_belt total_accel_belt  
## 1         no         11      1.41      8.07    -94.4              3  
## 2         no         11      1.41      8.07    -94.4              3  
## 3         no         11      1.42      8.07    -94.4              3  
## 4         no         12      1.48      8.05    -94.4              3  
## 5         no         12      1.48      8.07    -94.4              3
```

## 6	no	12	1.45	8.06	-94.4	3
##	kurtosis_roll_belt	kurtosis_picth_belt	kurtosis_yaw_belt			
## 1	NA		NA	NA		
## 2	NA		NA	NA		
## 3	NA		NA	NA		
## 4	NA		NA	NA		
## 5	NA		NA	NA		
## 6	NA		NA	NA		
##	skewness_roll_belt	skewness_roll_belt.1	skewness_yaw_belt	max_roll_belt		
## 1	NA		NA	NA	NA	
## 2	NA		NA	NA	NA	
## 3	NA		NA	NA	NA	
## 4	NA		NA	NA	NA	
## 5	NA		NA	NA	NA	
## 6	NA		NA	NA	NA	
##	max_picth_belt	max_yaw_belt	min_roll_belt	min_pitch_belt	min_yaw_belt	
## 1	NA		NA	NA	NA	
## 2	NA		NA	NA	NA	
## 3	NA		NA	NA	NA	
## 4	NA		NA	NA	NA	
## 5	NA		NA	NA	NA	
## 6	NA		NA	NA	NA	
##	amplitude_roll_belt	amplitude_pitch_belt	amplitude_yaw_belt			
## 1	NA		NA	NA		
## 2	NA		NA	NA		
## 3	NA		NA	NA		
## 4	NA		NA	NA		
## 5	NA		NA	NA		
## 6	NA		NA	NA		
##	var_total_accel_belt	avg_roll_belt	stddev_roll_belt	var_roll_belt		
## 1	NA		NA	NA	NA	
## 2	NA		NA	NA	NA	
## 3	NA		NA	NA	NA	
## 4	NA		NA	NA	NA	
## 5	NA		NA	NA	NA	
## 6	NA		NA	NA	NA	
##	avg_pitch_belt	stddev_pitch_belt	var_pitch_belt	avg_yaw_belt		
## 1	NA		NA	NA	NA	
## 2	NA		NA	NA	NA	
## 3	NA		NA	NA	NA	
## 4	NA		NA	NA	NA	
## 5	NA		NA	NA	NA	
## 6	NA		NA	NA	NA	
##	stddev_yaw_belt	var_yaw_belt	gyros_belt_x	gyros_belt_y	gyros_belt_z	
## 1	NA		NA	0.00	0.00	-0.02
## 2	NA		NA	0.02	0.00	-0.02
## 3	NA		NA	0.00	0.00	-0.02
## 4	NA		NA	0.02	0.00	-0.03
## 5	NA		NA	0.02	0.02	-0.02
## 6	NA		NA	0.02	0.00	-0.02

##	accel_belt_x	accel_belt_y	accel_belt_z	magnet_belt_x	magnet_belt_y	
## 1	-21	4	22	-3	599	
## 2	-22	4	22	-7	608	
## 3	-20	5	23	-2	600	
## 4	-22	3	21	-6	604	
## 5	-21	2	24	-6	600	
## 6	-21	4	21	0	603	
##	magnet_belt_z	roll_arm	pitch_arm	yaw_arm	total_accel_arm	var_accel_arm
## 1	-313	-128	22.5	-161	34	NA
## 2	-311	-128	22.5	-161	34	NA
## 3	-305	-128	22.5	-161	34	NA
## 4	-310	-128	22.1	-161	34	NA
## 5	-302	-128	22.1	-161	34	NA
## 6	-312	-128	22.0	-161	34	NA
##	avg_roll_arm	stddev_roll_arm	var_roll_arm	avg_pitch_arm	stddev_pitch_arm	
## 1	NA	NA	NA	NA	NA	NA
## 2	NA	NA	NA	NA	NA	NA
## 3	NA	NA	NA	NA	NA	NA
## 4	NA	NA	NA	NA	NA	NA
## 5	NA	NA	NA	NA	NA	NA
## 6	NA	NA	NA	NA	NA	NA
##	var_pitch_arm	avg_yaw_arm	stddev_yaw_arm	var_yaw_arm	gyros_arm_x	
## 1	NA	NA	NA	NA	0.00	
## 2	NA	NA	NA	NA	0.02	
## 3	NA	NA	NA	NA	0.02	
## 4	NA	NA	NA	NA	0.02	
## 5	NA	NA	NA	NA	0.00	
## 6	NA	NA	NA	NA	0.02	
##	gyros_arm_y	gyros_arm_z	accel_arm_x	accel_arm_y	accel_arm_z	magnet_arm_x
## 1	0.00	-0.02	-288	109	-123	-368
## 2	-0.02	-0.02	-290	110	-125	-369
## 3	-0.02	-0.02	-289	110	-126	-368
## 4	-0.03	0.02	-289	111	-123	-372
## 5	-0.03	0.00	-289	111	-123	-374
## 6	-0.03	0.00	-289	111	-122	-369
##	magnet_arm_y	magnet_arm_z	kurtosis_roll_arm	kurtosis_pitch_arm		
## 1	337	516	NA	NA		
## 2	337	513	NA	NA		
## 3	344	513	NA	NA		
## 4	344	512	NA	NA		
## 5	337	506	NA	NA		
## 6	342	513	NA	NA		
##	kurtosis_yaw_arm	skewness_roll_arm	skewness_pitch_arm	skewness_yaw_arm		
## 1	NA	NA	NA	NA		
## 2	NA	NA	NA	NA		
## 3	NA	NA	NA	NA		
## 4	NA	NA	NA	NA		
## 5	NA	NA	NA	NA		
## 6	NA	NA	NA	NA		
##	max_roll_arm	max_pitch_arm	max_yaw_arm	min_roll_arm	min_pitch_arm	

## 1	NA	NA	NA	NA	NA
## 2	NA	NA	NA	NA	NA
## 3	NA	NA	NA	NA	NA
## 4	NA	NA	NA	NA	NA
## 5	NA	NA	NA	NA	NA
## 6	NA	NA	NA	NA	NA
##	min_yaw_arm	amplitude_roll_arm	amplitude_pitch_arm	amplitude_yaw_arm	
## 1	NA	NA	NA	NA	
## 2	NA	NA	NA	NA	
## 3	NA	NA	NA	NA	
## 4	NA	NA	NA	NA	
## 5	NA	NA	NA	NA	
## 6	NA	NA	NA	NA	
##	roll_dumbbell	pitch_dumbbell	yaw_dumbbell	kurtosis_roll_dumbbell	
## 1	13.05217	-70.49400	-84.87394	NA	
## 2	13.13074	-70.63751	-84.71065	NA	
## 3	12.85075	-70.27812	-85.14078	NA	
## 4	13.43120	-70.39379	-84.87363	NA	
## 5	13.37872	-70.42856	-84.85306	NA	
## 6	13.38246	-70.81759	-84.46500	NA	
##	kurtosis_pitch_dumbbell	kurtosis_yaw_dumbbell	skewness_roll_dumbbell		
## 1	NA	NA	NA		
## 2	NA	NA	NA		
## 3	NA	NA	NA		
## 4	NA	NA	NA		
## 5	NA	NA	NA		
## 6	NA	NA	NA		
##	skewness_pitch_dumbbell	skewness_yaw_dumbbell	max_roll_dumbbell		
## 1	NA	NA	NA		
## 2	NA	NA	NA		
## 3	NA	NA	NA		
## 4	NA	NA	NA		
## 5	NA	NA	NA		
## 6	NA	NA	NA		
##	max_pitch_dumbbell	max_yaw_dumbbell	min_roll_dumbbell	min_pitch_dumbbell	
## 1	NA	NA	NA	NA	
## 2	NA	NA	NA	NA	
## 3	NA	NA	NA	NA	
## 4	NA	NA	NA	NA	
## 5	NA	NA	NA	NA	
## 6	NA	NA	NA	NA	
##	min_yaw_dumbbell	amplitude_roll_dumbbell	amplitude_pitch_dumbbell		
## 1	NA	NA	NA		
## 2	NA	NA	NA		
## 3	NA	NA	NA		
## 4	NA	NA	NA		
## 5	NA	NA	NA		
## 6	NA	NA	NA		
##	amplitude_yaw_dumbbell	total_accel_dumbbell	var_accel_dumbbell		
## 1	NA	37	NA		

## 2	NA	37	NA
## 3	NA	37	NA
## 4	NA	37	NA
## 5	NA	37	NA
## 6	NA	37	NA
##	avg_roll_dumbbell	stddev_roll_dumbbell	var_roll_dumbbell
## 1	NA	NA	NA
## 2	NA	NA	NA
## 3	NA	NA	NA
## 4	NA	NA	NA
## 5	NA	NA	NA
## 6	NA	NA	NA
##	avg_pitch_dumbbell	stddev_pitch_dumbbell	var_pitch_dumbbell
## 1	NA	NA	NA
## 2	NA	NA	NA
## 3	NA	NA	NA
## 4	NA	NA	NA
## 5	NA	NA	NA
## 6	NA	NA	NA
##	avg_yaw_dumbbell	stddev_yaw_dumbbell	var_yaw_dumbbell
## 1	NA	NA	NA
## 2	NA	NA	NA
## 3	NA	NA	NA
## 4	NA	NA	NA
## 5	NA	NA	NA
## 6	NA	NA	NA
##	gyros_dumbbell_x		
## 1	0		
## 2	0		
## 3	0		
## 4	0		
## 5	0		
## 6	0		
##	gyros_dumbbell_y	gyros_dumbbell_z	accel_dumbbell_x
## 1	-0.02	0.00	-234
## 2	-0.02	0.00	-233
## 3	-0.02	0.00	-232
## 4	-0.02	-0.02	-232
## 5	-0.02	0.00	-233
## 6	-0.02	0.00	-234
##	accel_dumbbell_y		
## 1	47		
## 2	47		
## 3	46		
## 4	48		
## 5	48		
## 6	48		
##	accel_dumbbell_z	magnet_dumbbell_x	magnet_dumbbell_y
## 1	-271	-559	293
## 2	-269	-555	296
## 3	-270	-561	298
## 4	-269	-552	303
## 5	-270	-554	292
## 6	-269	-558	294
##	magnet_dumbbell_z		
## 1	-65		
## 2	-64		
## 3	-63		
## 4	-60		
## 5	-68		
## 6	-66		
##	roll_forearm	pitch_forearm	yaw_forearm
## 1	28.4	-63.9	-153
## 2	28.3	-63.9	-153
## 3	28.3	-63.9	-152
## 4	28.1	-63.9	-152
## 5	28.0	-63.9	-152
## 6	27.9	-63.9	-152
##	kurtosis_roll_forearm		
## 1	NA		
## 2	NA		
##	kurtosis_pitch_forearm	kurtosis_yaw_forearm	skewness_roll_forearm
## 1	NA	NA	NA
## 2	NA	NA	NA

## 3	NA	NA	NA	
## 4	NA	NA	NA	
## 5	NA	NA	NA	
## 6	NA	NA	NA	
##	skewness_pitch_forearm	skewness_yaw_forearm	max_roll_forearm	
## 1	NA	NA	NA	
## 2	NA	NA	NA	
## 3	NA	NA	NA	
## 4	NA	NA	NA	
## 5	NA	NA	NA	
## 6	NA	NA	NA	
##	max_picth_forearm	max_yaw_forearm	min_roll_forearm	min_pitch_forearm
## 1	NA	NA	NA	NA
## 2	NA	NA	NA	NA
## 3	NA	NA	NA	NA
## 4	NA	NA	NA	NA
## 5	NA	NA	NA	NA
## 6	NA	NA	NA	NA
##	min_yaw_forearm	amplitude_roll_forearm	amplitude_pitch_forearm	
## 1	NA	NA	NA	
## 2	NA	NA	NA	
## 3	NA	NA	NA	
## 4	NA	NA	NA	
## 5	NA	NA	NA	
## 6	NA	NA	NA	
##	amplitude_yaw_forearm	total_accel_forearm	var_accel_forearm	
## 1	NA	36	NA	
## 2	NA	36	NA	
## 3	NA	36	NA	
## 4	NA	36	NA	
## 5	NA	36	NA	
## 6	NA	36	NA	
##	avg_roll_forearm	stddev_roll_forearm	var_roll_forearm	avg_pitch_forearm
## 1	NA	NA	NA	NA
## 2	NA	NA	NA	NA
## 3	NA	NA	NA	NA
## 4	NA	NA	NA	NA
## 5	NA	NA	NA	NA
## 6	NA	NA	NA	NA
##	stddev_pitch_forearm	var_pitch_forearm	avg_yaw_forearm	
## 1	NA	NA	NA	
## 2	NA	NA	NA	
## 3	NA	NA	NA	
## 4	NA	NA	NA	
## 5	NA	NA	NA	
## 6	NA	NA	NA	
##	stddev_yaw_forearm	var_yaw_forearm	gyros_forearm_x	gyros_forearm_y
## 1	NA	NA	0.03	0.00
## 2	NA	NA	0.02	0.00
## 3	NA	NA	0.03	-0.02

```
## 4      NA      NA      0.02      -0.02
## 5      NA      NA      0.02      0.00
## 6      NA      NA      0.02      -0.02
## gyros_forearm_z accel_forearm_x accel_forearm_y accel_forearm_z
## 1      -0.02      192      203      -215
## 2      -0.02      192      203      -216
## 3       0.00      196      204      -213
## 4       0.00      189      206      -214
## 5      -0.02      189      206      -214
## 6      -0.03      193      203      -215
## magnet_forearm_x magnet_forearm_y magnet_forearm_z classe
## 1      -17      654      476      A
## 2      -18      661      473      A
## 3      -18      658      469      A
## 4      -16      658      469      A
## 5      -17      655      473      A
## 6       -9      660      478      A
```

Looking through the data we see a lot of columns with NA values. We remove these to reduce the number of columns. First, we remove columns that contain NA missing values.

```
ruwedata <- ruwedata[, colSums(is.na(ruwedata)) == 0]
testruw <- testruw[, colSums(is.na(testruw)) == 0]

classe <- ruwedata$classe
deletecolumns <- grepl("^X|timestamp|window", names(ruwedata))
ruwedata <- ruwedata[, !deletecolumns]
trainCleaned <- ruwedata[, sapply(ruwedata, is.numeric)]
trainCleaned$classe <- classe
testRemove <- grepl("^X|timestamp|window", names(testruw))
testruw <- testruw[, !testRemove]
testCleaned <- testruw[, sapply(testruw, is.numeric)]
```

Now, the cleaned training data set contains 19622 observations and 53 variables, while the testing data set contains 20 observations and 53 variables.

Slice the data

Then, we can split the cleaned training set into a pure training data set (70%) and a validation data set (30%). We will use the validation data set to conduct cross validation in future steps.

```
set.seed(22519) # For reproducible purpose
inTrain <- createDataPartition(trainCleaned$classe, p=0.70, list=F)
trainData <- trainCleaned[inTrain, ]
testData <- trainCleaned[-inTrain, ]
```

Analysis

```
classes <- table(trainData$classe)
classes
```



```
##
##      A      B      C      D      E
## 3906 2658 2396 2252 2525
```

Data Modeling

WE use the random forest model. Tree models are good in picking the most important variables.

```
controlRf <- trainControl(method="cv", 5)
modelRf <- train(classe ~ ., data=trainData, method="rf",
trControl=controlRf, ntree=250)
modelRf

## Random Forest
##
## 13737 samples
##      52 predictor
##      5 classes: 'A', 'B', 'C', 'D', 'E'
##
## No pre-processing
## Resampling: Cross-Validated (5 fold)
## Summary of sample sizes: 10989, 10989, 10991, 10990, 10989
## Resampling results across tuning parameters:
##
##      mtry  Accuracy   Kappa
##      2     0.9901727 0.9875673
##      27     0.9917015 0.9895017
##      52     0.9840572 0.9798282
##
## Accuracy was used to select the optimal model using the largest value.
## The final value used for the model was mtry = 27.
```

The next step is to check the precision of the model.

```
prediction <- predict(modelRf, testData)
confusionMatrix(testData$classe, prediction)

## $positive
## NULL
##
## $table
##           Reference
## Prediction      A      B      C      D      E
##           A 1673      0      0      0      1
##           B      5 1131      3      0      0
##           C      0      0 1021      5      0
##           D      0      0  13  949      2
##           E      0      0   1   6 1075
##
## $overall
```

```
##      Accuracy      Kappa AccuracyLower AccuracyUpper AccuracyNull
##      0.9938828      0.9922620      0.9915411      0.9957120      0.2851317
## AccuracyPValue McNemarPValue
##      0.0000000      NaN
##
## $byClass
##      Sensitivity Specificity Pos Pred Value Neg Pred Value Precision
## Class: A      0.9970203      0.9997623      0.9994026      0.9988126 0.9994026
## Class: B      1.0000000      0.9983172      0.9929763      1.0000000 0.9929763
## Class: C      0.9836224      0.9989684      0.9951267      0.9965013 0.9951267
## Class: D      0.9885417      0.9969543      0.9844398      0.9977647 0.9844398
## Class: E      0.9972171      0.9985438      0.9935305      0.9993754 0.9935305
##      Recall      F1 Prevalence Detection Rate
## Class: A 0.9970203 0.9982100 0.2851317      0.2842821
## Class: B 1.0000000 0.9964758 0.1921835      0.1921835
## Class: C 0.9836224 0.9893411 0.1763806      0.1734919
## Class: D 0.9885417 0.9864865 0.1631266      0.1612574
## Class: E 0.9972171 0.9953704 0.1831776      0.1826678
##      Detection Prevalence Balanced Accuracy
## Class: A      0.2844520      0.9983913
## Class: B      0.1935429      0.9991586
## Class: C      0.1743415      0.9912954
## Class: D      0.1638063      0.9927480
## Class: E      0.1838573      0.9978804
##
## $mode
## [1] "sens_spec"
##
## $dots
## list()
##
## attr(,"class")
## [1] "confusionMatrix"

precisie <- postResample(prediction, testData$classe)
precisie

## Accuracy      Kappa
## 0.9938828 0.9922620

outofsample <- 1 - as.numeric(confusionMatrix(testData$classe,
prediction)$overall[1])
outofsample

## [1] 0.006117247
```

Estimated precision of the model 99.21% estimated out-of-sample error 0.78%.

Using the model to predict

Now we apply the model to the test set

```
result <- predict(modelRf, testCleaned[, -length(names(testCleaned))])
result

## [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
```