

Practical Machine Learning - Course project

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.

In this project, the goal will be to use data from accelerometers on the belt, forearm, arm, and dumbbell of 6 participants to predict the manner in which they did the exercise. This is the “classe” variable in the training set. More information is available from the website here: [<http://web.archive.org/web/20161224072740/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)
- throwing the hips to the front (Class E)

Body-sensors were placed on four different parts: the arm, the belt, the forearm and on the dumbbell.

Part 1: Loading and exploring the data

Load the necessary packages.

```
## Loading required package: lattice

## Loading required package: ggplot2

## Rattle: A free graphical interface for data science with R.
## Version 5.3.0 Copyright (c) 2006-2018 Togaware Pty Ltd.
## Type 'rattle()' to shake, rattle, and roll your data.
```

Load the data and explore it to get a sense of what is in the dataset.

```
data_training = read.csv("pml-training.csv")
data_testing = read.csv("pml-testing.csv")
```

```
print(dim(data_testing))
```

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

```
print(dim(data_training))
```

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

```
summary(data_training$classe)
```

```
## A B C D E  
## 5580 3797 3422 3216 3607
```

```
str(data_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 1323084232 1323084232 ...  
## $ raw_timestamp_part_2 : int 788290 808298 820366 120339 196328 304277 368296 440390 484323 484323 ...  
## $ 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 1.43 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 NA ...
## $ max_pitch_dumbbell : num NA 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 NA ...
## $ min_pitch_dumbbell : num NA 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 NA ...
## [list output truncated]
```

Part 2: Preparing the data

Before starting with the actual analysis the data has to be prepared. The current dataset contains 160 variables, which is a lot and hence we want to reduce this.

First we remove columns 1:7, since they contain information that is not relevant to the investigation. These include features like index, name and timestamp.

```
## [1] 19622 153
```

Next we will delete the columns where over 90% of the observations is NA, and thus only keep those with less than 90% of NA values.

```
na_cols <- colSums(is.na(data_training))/nrow(data_training) < 0.9
data_training <- data_training[,c(na_cols)]
data_testing <- data_testing[,c(na_cols)]
print(dim(data_training))
```

```
## [1] 19622 86
```

This reduces the number of variables from 153 to 86.

To further reduce the number of variables for analysis we will look at the variance of the variables. The variables with a variance near zero will be excluded from the analysis.

```
nz_cols <- nearZeroVar(data_training)
data_training <- data_training[, -nz_cols]
data_testing <- data_testing[, -nz_cols]
print(dim(data_training))
```

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

As can be seen the set of predictors is now reduced to 53 variables.

Now we will split the training data into a training set and a test set. Important is to set the seed so that the results are reproducible.

```
set.seed(12345)
inTrain <- createDataPartition(y=data_training$classe, p=0.7, list=FALSE)
training <- data_training[inTrain,]
testing <- data_training[-inTrain,]
```

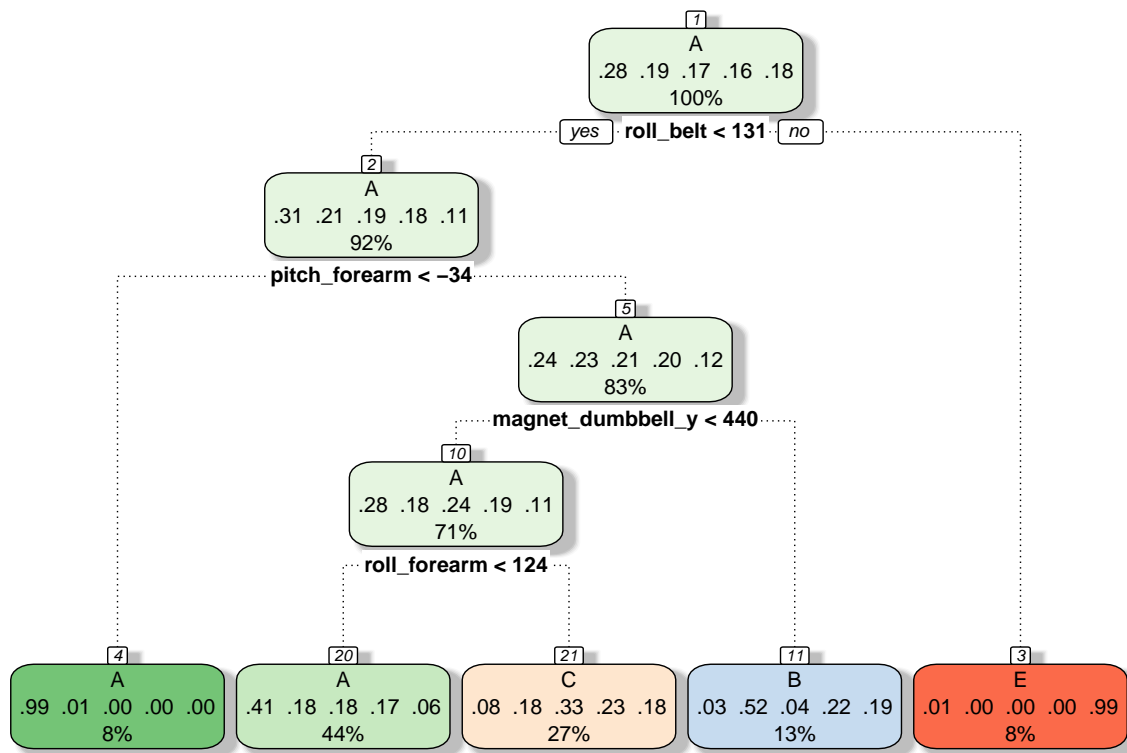
Part 3: Fitting the model

Decision tree

For this classification problem we first try to fit a decision tree model using the training function of the caret package with rpart as specified method.

```
model_tree <- train(classe ~., data=training, method="rpart")
```

```
fancyRpartPlot(model_tree$finalModel, sub = "Decision tree")
```



Decision tree

Prediction using the testing data

```
prediction_tree <- predict(model_tree, newdata=testing)
confusionMatrix(prediction_tree, testing$classe)
```

```
## Confusion Matrix and Statistics
##
##           Reference
## Prediction    A    B    C    D    E
##           A 1525  484  499  423  153
##           B   29  385   37  187  159
##           C  116  270  490  354  289
##           D    0    0    0    0    0
##           E    4    0    0    0  481
```

```
##
## Overall Statistics
##
##           Accuracy : 0.4895
##           95% CI : (0.4767, 0.5024)
##       No Information Rate : 0.2845
##       P-Value [Acc > NIR] : < 2.2e-16
##
##           Kappa : 0.3324
##
## McNemar's Test P-Value : NA
##
## Statistics by Class:
##
##           Class: A Class: B Class: C Class: D Class: E
## Sensitivity      0.9110 0.33802 0.47758 0.0000 0.44455
## Specificity      0.6298 0.91319 0.78823 1.0000 0.99917
## Pos Pred Value   0.4945 0.48306 0.32258      NaN 0.99175
## Neg Pred Value   0.9468 0.85181 0.87723 0.8362 0.88870
## Prevalence       0.2845 0.19354 0.17434 0.1638 0.18386
## Detection Rate   0.2591 0.06542 0.08326 0.0000 0.08173
## Detection Prevalence 0.5240 0.13543 0.25811 0.0000 0.08241
## Balanced Accuracy 0.7704 0.62560 0.63291 0.5000 0.72186
```

The prediction accuracy of this model is 0.489, which is quite low. We ideally want a higher accuracy.

Random forest

The second model we will fit is the random forest model. We do this again by using the train function from the caret package but now with “rf” as the specified method.

```
train_control <- trainControl(method="cv", number=5)
model_rf <- train(classe ~., data=training, method="rf", trControl=train_control)
```

Prediction using the testing data

```
prediction_rf <- predict(model_rf, newdata=testing)
confusionMatrix(prediction_rf, testing$classe)
```

```
## Confusion Matrix and Statistics
##
##           Reference
## Prediction    A    B    C    D    E
##           A 1672    6    0    0    0
##           B    1 1131    5    0    0
##           C    1    2 1018    8    1
##           D    0    0    3  956    1
##           E    0    0    0    0 1080
##
## Overall Statistics
##
##           Accuracy : 0.9952
```

```
##          95% CI : (0.9931, 0.9968)
##    No Information Rate : 0.2845
##    P-Value [Acc > NIR] : < 2.2e-16
##
##          Kappa : 0.994
##
##    McNemar's Test P-Value : NA
##
## Statistics by Class:
##
##          Class: A Class: B Class: C Class: D Class: E
## Sensitivity      0.9988   0.9930   0.9922   0.9917   0.9982
## Specificity      0.9986   0.9987   0.9975   0.9992   1.0000
## Pos Pred Value   0.9964   0.9947   0.9883   0.9958   1.0000
## Neg Pred Value   0.9995   0.9983   0.9984   0.9984   0.9996
## Prevalence       0.2845   0.1935   0.1743   0.1638   0.1839
## Detection Rate   0.2841   0.1922   0.1730   0.1624   0.1835
## Detection Prevalence 0.2851   0.1932   0.1750   0.1631   0.1835
## Balanced Accuracy 0.9987   0.9959   0.9949   0.9954   0.9991
```

The accuracy of the random forest model is 0.995 when using cross-validation with k=5 folds. This accuracy is very high and clearly higher than the accuracy of the decision tree model. Thus we choose the random forest model as the best prediction model for this dataset.

```
print(model_rf)
```

```
## 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: 10990, 10989, 10991, 10990, 10988
## Resampling results across tuning parameters:
##
##    mtry  Accuracy  Kappa
##    2     0.9911188  0.9887641
##    27    0.9911917  0.9888565
##    52    0.9834027  0.9790024
##
## Accuracy was used to select the optimal model using the largest value.
## The final value used for the model was mtry = 27.
```

Conclusion

The random forest model has shown to be a better model than the decision tree and hence this model will be used to predict the outcomes, the classes, of the test dataset.

```
pred <- predict(model_rf, newdata = data_testing)
pred
```

```
## [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
```