

Practical Machine Learning

Thiago Almeida

18 de julho de 2018

```
knitr::opts_chunk$set(echo = TRUE, cache = TRUE, cache.lazy = FALSE)
library(caret)
## Loading required package: lattice
## Loading required package: ggplot2
library(rattle)
## Rattle: A free graphical interface for data science with R.
## Version 5.1.0 Copyright (c) 2006-2017 Togaware Pty Ltd.
## Type 'rattle()' to shake, rattle, and roll your data.
```

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, 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).

Human Activity Recognition - HAR - has emerged as a key research area in the last years and is gaining increasing attention by the pervasive computing research community (see picture below, that illustrates the increasing number of publications in HAR with wearable accelerometers), especially for the development of context-aware systems. There are many potential applications for HAR, like: elderly monitoring, life log systems for monitoring energy expenditure and for supporting weight-loss programs, and digital assistants for weight lifting exercises

Loading Data

Loading data from web.

```
#load train
urltrain <-
read.csv(url("https://d396qusza40orc.cloudfront.net/predmachlearn/pml-
training.csv"),header=TRUE,na.strings = c("", "NA"))
urltest <-
read.csv(url("https://d396qusza40orc.cloudfront.net/predmachlearn/pml-
testing.csv"),header=TRUE,na.strings = c("", "NA"))
```

Visualize Data

```
#visualize data
str(urltrain)
## '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 1323084232
1323084232 ...
##  $ raw_timestamp_part_2   : int  788290 808298 820366 120339
196328 304277 368296 440390 484323 484434 ...
##  $ 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/ 396 levels "-0.016850","-
0.021024",...: NA NA NA NA NA NA NA NA NA NA ...
##  $ kurtosis_pitch_belt   : Factor w/ 316 levels "-0.021887","-
0.060755",...: NA NA NA NA NA NA NA NA NA NA ...
##  $ kurtosis_yaw_belt     : Factor w/  1 level "#DIV/0!": NA NA NA
NA NA NA NA NA NA NA NA ...
##  $ skewness_roll_belt    : Factor w/ 394 levels "-0.003095","-
0.010002",...: NA NA NA NA NA NA NA NA NA NA ...
##  $ skewness_roll_belt.1  : Factor w/ 337 levels "-0.005928","-
0.005960",...: NA NA NA NA NA NA NA NA NA NA ...
##  $ skewness_yaw_belt     : Factor w/  1 level "#DIV/0!": NA NA NA
NA NA NA NA NA NA NA NA ...
##  $ 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/  67 levels "-0.1","-0.2",...:
NA NA NA NA NA NA NA NA NA NA ...
##  $ 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/  67 levels "-0.1","-0.2",...:
NA NA NA NA NA NA NA NA NA NA ...
##  $ 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/  3 levels
"#DIV/0!","0.00",...: NA NA NA NA NA NA NA NA NA ...
##  $ 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/ 329 levels "-0.02438","-
0.04190",...: NA NA NA NA NA NA NA NA NA NA ...
## $ kurtosis_pitch_arm : Factor w/ 327 levels "-0.00484","-
0.01311",...: NA NA NA NA NA NA NA NA NA NA ...
## $ kurtosis_yaw_arm   : Factor w/ 394 levels "-0.01548","-
0.01749",...: NA NA NA NA NA NA NA NA NA NA ...
## $ skewness_roll_arm  : Factor w/ 330 levels "-0.00051","-
0.00696",...: NA NA NA NA NA NA NA NA NA NA ...
## $ skewness_pitch_arm : Factor w/ 327 levels "-0.00184","-
0.01185",...: NA NA NA NA NA NA NA NA NA NA ...
## $ skewness_yaw_arm   : Factor w/ 394 levels "-0.00311","-
0.00562",...: NA NA NA NA NA NA NA NA NA NA ...

```

```
## $ max_roll_arm : num NA NA NA NA NA NA NA NA NA NA NA ...
## $ max_pitch_arm : num NA NA NA NA NA NA NA NA NA NA NA ...
## $ max_yaw_arm : int NA NA NA NA NA NA NA NA NA NA NA ...
## $ min_roll_arm : num NA NA NA NA NA NA NA NA NA NA NA ...
## $ min_pitch_arm : num NA NA NA NA NA NA NA NA NA NA NA ...
## $ min_yaw_arm : int NA NA NA NA NA NA NA NA NA NA NA ...
## $ amplitude_roll_arm : num NA NA NA NA NA NA NA NA NA NA NA ...
## $ amplitude_pitch_arm : num NA NA NA NA NA NA NA NA NA NA NA ...
## $ amplitude_yaw_arm : int NA 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/ 397 levels "-0.0035","-0.0073",...: NA NA NA NA NA NA NA NA NA NA NA ...
## $ kurtosis_pitch_dumbbell : Factor w/ 400 levels "-0.0163","-0.0233",...: NA NA NA NA NA NA NA NA NA NA NA ...
## $ kurtosis_yaw_dumbbell : Factor w/ 1 level "#DIV/0!": NA NA NA NA NA NA NA NA NA NA NA ...
## $ skewness_roll_dumbbell : Factor w/ 400 levels "-0.0082","-0.0096",...: NA NA NA NA NA NA NA NA NA NA NA ...
## $ skewness_pitch_dumbbell : Factor w/ 401 levels "-0.0053","-0.0084",...: NA NA NA NA NA NA NA NA NA NA NA ...
## $ skewness_yaw_dumbbell : Factor w/ 1 level "#DIV/0!": NA NA NA NA NA NA NA NA NA NA NA ...
## $ 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/ 72 levels "-0.1","-0.2",...: NA NA NA NA NA NA NA NA NA NA NA ...
## $ 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/ 72 levels "-0.1","-0.2",...: NA NA NA NA NA NA NA NA NA NA NA ...
## $ amplitude_roll_dumbbell : num NA NA NA NA NA NA NA NA NA NA NA ...
## [list output truncated]
```

Cleaning Data

In this dataset there are too many NA values. In this step the NA columns will be excluded. The same columns will be excluded in test data.

```
#exclude NA columns
train <- urltrain[, colSums(is.na(urltrain)) == 0]
test <- urltest[, colSums(is.na(urltrain)) == 0]
```

There are correlated columns that can be excluded.

```
#exclude columns with high correlation
cor.matrix <- cor(train[apply(train, is.numeric)])
c <- findCorrelation(cor.matrix, cutoff = .90)
train <- train[,-c]
```

Splitting Validation Dataset

The train dataset will be splitted for validation.

```
#seed
```

```
set.seed(9876)
```

```
#split 75%
inTrain <- createDataPartition(train$classe, p=0.75, list=FALSE)
train <- train[inTrain,]
valid <- train[-inTrain,]
```

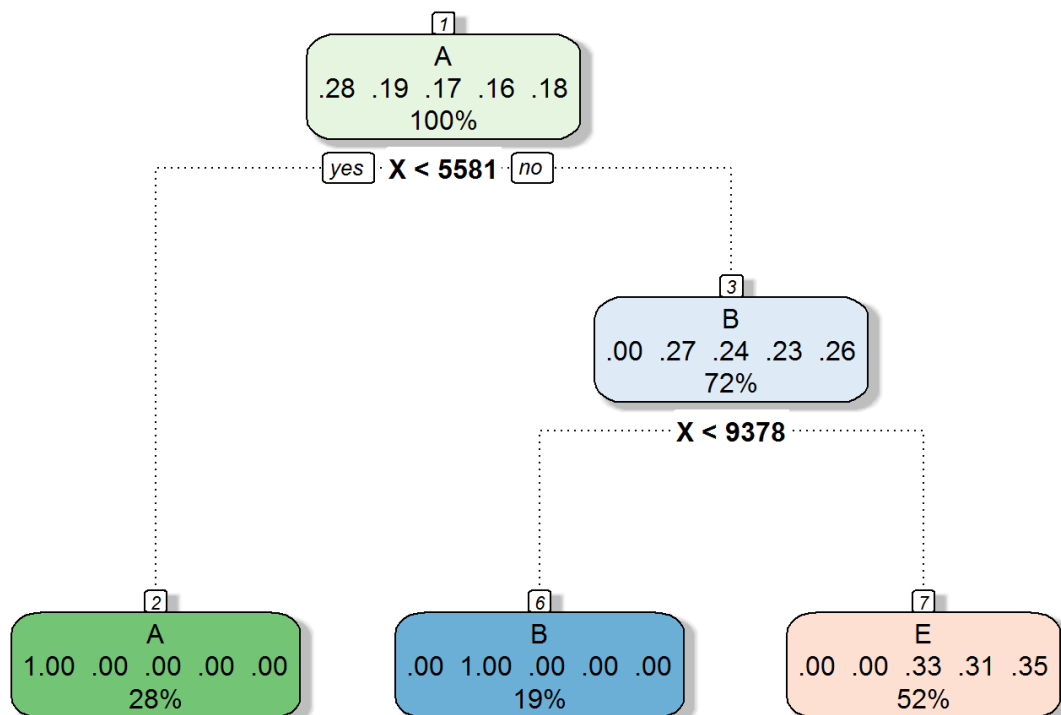
Decision Tree

The first method will be Decision Tree using the Caret library

```
#control
control <- trainControl(method="cv", number=3, classProbs=TRUE,
allowParallel=TRUE, verboseIter = FALSE)
```

```
#model
fit_rpart <- train(classe ~ ., data = train, method = "rpart",
trControl = control)
```

```
#plot
fancyRpartPlot(fit_rpart$finalModel)
## Warning: Bad 'data' field in model 'call' field.
##       To make this warning go away:
##       Call prp with roundint=FALSE,
##       or rebuild the rpart model with model=TRUE.
```



Rattle 2018-jul-19 14:01:15 talmeida

```
#predict
pred_rpart <- predict(fit_rpart, valid)

confusion_rpart <- confusionMatrix(valid$classe, pred_rpart)
```

The accuracy of this method is 0.659875

Random Forest

The second method will be Random Forest using the Caret library

```
#model
fit_rf <- train(classe ~ ., data = train, method = "rf", trControl =
control, ntree=500, keep.forest=TRUE, importance=TRUE)

#predict
pred_rf <- predict(fit_rf, valid)

confusion_rf <- confusionMatrix(valid$classe, pred_rf)
```

The accuracy of this method is 1

Gradient Boosting Method

The second method will be Gradient Boosting Method using the Caret library

```
#model
fit_gbm <- train(classe ~ ., data = train, method = "gbm", trControl =
control)
#predict
pred_gbm <- predict(fit_gbm, valid)

confusion_gbm <- confusionMatrix(valid$classe, pred_gbm)
```

The accuracy of this method is 1

Using Test Data

The best accuracy is using the Random Forest method, so the test will use this model.

```
#predict
pred_test <- predict(fit_rf, test)
```