

Practical Machine Learning: An Analysis of the Weight Lifting Exercises Dataset

Bingfang Xu

Friday, January 29, 2016

Introduction

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 predict the manner in which people did the exercise. "classe" variable will be outcome, any of other variables are predictors. Model building, cross evaluation and prediction were performed. Data are collected from accelerometers on the belt, forearm, arm, and dumbbell of 6 participants.

load necessary packages. If the packages is not pre-installed, use `install.packages()` function to install.

```
library(caret)
```

```
## Warning: package 'caret' was built under R version 3.1.3
```

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

```
## Warning: package 'ggplot2' was built under R version 3.1.3
```

```
library(randomForest)
```

```
## Warning: package 'randomForest' was built under R version 3.1.3
```

```
## randomForest 4.6-12  
## Type rfNews() to see new features/changes/bug fixes.
```

```
library(rpart)
```

```
## Warning: package 'rpart' was built under R version 3.1.3
```

Put the data in the work directory. Load the data. Define missing data. Look at the variables in the data

```
trainingset <- read.csv("./pml-training.csv", na.strings=c("NA", "#DIV/0!", ""))  
testingset <- read.csv('./pml-testing.csv', na.strings=c("NA", "#DIV/0!", ""))  
str(trainingset)
```

```

## '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 132
3084232 1323084232 1323084232 1323084232 1323084232 1323084232 ...
## $ raw_timestamp_part_2 : int  788290 808298 820366 120339 196328 304277 36829
6 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 : num  NA NA NA NA NA NA NA NA NA NA ...
## $ kurtosis_pitch_belt : num  NA NA NA NA NA NA NA NA NA NA ...
## $ kurtosis_yaw_belt : logi  NA NA NA NA NA NA ...
## $ skewness_roll_belt : num  NA NA NA NA NA NA NA NA NA NA ...
## $ skewness_roll_belt.1 : num  NA NA NA NA NA NA NA NA NA NA ...
## $ skewness_yaw_belt : logi  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 : num  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 : num  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 : num  NA 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 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 -0.0
2 -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 -3
08 ...
## $ roll_arm          : num  -128 -128 -128 -128 -128 -128 -128 -128 -128 -1
28 ...
## $ 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 -1
61 ...
## $ 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 NA ...
## $ avg_roll_arm      : num  NA NA NA NA NA NA NA NA NA NA NA ...
## $ stddev_roll_arm   : num  NA NA NA NA NA NA NA NA NA NA NA ...
## $ var_roll_arm      : num  NA NA NA NA NA NA NA NA NA NA NA ...
## $ avg_pitch_arm     : num  NA NA NA NA NA NA NA NA NA NA NA ...
## $ stddev_pitch_arm  : num  NA NA NA NA NA NA NA NA NA NA NA ...
## $ var_pitch_arm     : num  NA NA NA NA NA NA NA NA NA NA NA ...
## $ avg_yaw_arm       : num  NA NA NA NA NA NA NA NA NA NA NA ...
## $ stddev_yaw_arm    : num  NA NA NA NA NA NA NA NA NA NA NA ...
## $ var_yaw_arm       : num  NA 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 -2
88 ...
## $ 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 -1
24 ...
## $ magnet_arm_x      : int  -368 -369 -368 -372 -374 -369 -373 -372 -369 -3
76 ...
## $ 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 : num  NA NA NA NA NA NA NA NA NA NA NA ...
## $ kurtosis_picth_arm : num  NA NA NA NA NA NA NA NA NA NA NA ...

```

```
## $ kurtosis_yaw_arm      : num  NA NA NA NA NA NA NA NA NA NA NA ...
## $ skewness_roll_arm    : num  NA NA NA NA NA NA NA NA NA NA NA ...
## $ skewness_pitch_arm   : num  NA NA NA NA NA NA NA NA NA NA NA ...
## $ skewness_yaw_arm     : num  NA 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 : num  NA NA NA NA NA NA NA NA NA NA NA ...
## $ kurtosis_pitch_dumbbell : num  NA NA NA NA NA NA NA NA NA NA NA ...
## $ kurtosis_yaw_dumbbell  : logi  NA NA NA NA NA NA NA ...
## $ skewness_roll_dumbbell : num  NA NA NA NA NA NA NA NA NA NA NA ...
## $ skewness_pitch_dumbbell : num  NA NA NA NA NA NA NA NA NA NA NA ...
## $ skewness_yaw_dumbbell  : logi  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      : num  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      : num  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]
```

```
set.seed(1111)
trainingset<-trainingset[,colSums(is.na(trainingset)) == 0]
testingset <-testingset[,colSums(is.na(testingset)) == 0]
trainingset  <-trainingset[,-c(1:7)]
testingset <-testingset[,-c(1:7)]
```

Subset trainingset into newtraining and validation sets.

```
inTrain <- createDataPartition(y=trainingset$classe,p = 3/4,list=FALSE)
newtrainingset <- trainingset[inTrain, ]
validationset <- trainingset[-inTrain, ]
```

Build a model using newtrainingset data and decision tree, and use the model to predict validationset data. Preform cross-validation.

```
#Building a model using decsion tree:
modell <- rpart(classe ~ ., data=newtrainingset, method="class")

# Predicting:
prediction1 <- predict(modell, validationset, type = "class")

#cross validation
confusionMatrix(prediction1, validationset$classe)
```

```

## Confusion Matrix and Statistics
##
##           Reference
## Prediction   A    B    C    D    E
##           A 1255  139   14   44   17
##           B   45  542   70   68   83
##           C   51  142  694  124  131
##           D   19   71   44  521   54
##           E   25   55   33   47  616
##
## Overall Statistics
##
##           Accuracy : 0.7398
##           95% CI : (0.7273, 0.752)
##           No Information Rate : 0.2845
##           P-Value [Acc > NIR] : < 2.2e-16
##
##           Kappa : 0.6704
##           McNemar's Test P-Value : < 2.2e-16
##
## Statistics by Class:
##
##           Class: A Class: B Class: C Class: D Class: E
## Sensitivity      0.8996   0.5711   0.8117   0.6480   0.6837
## Specificity      0.9390   0.9327   0.8894   0.9541   0.9600
## Pos Pred Value   0.8543   0.6708   0.6077   0.7348   0.7938
## Neg Pred Value   0.9592   0.9006   0.9572   0.9325   0.9310
## Prevalence       0.2845   0.1935   0.1743   0.1639   0.1837
## Detection Rate   0.2559   0.1105   0.1415   0.1062   0.1256
## Detection Prevalence 0.2996   0.1648   0.2329   0.1446   0.1582
## Balanced Accuracy 0.9193   0.7519   0.8505   0.8011   0.8219

```

Build a model using newtrainingset data and random forest, and use the model to predict validationset data. Preform cross-validation.

```
model2 <- randomForest(classe ~. , data=newtrainingset, method="class")

# Predicting:
prediction2 <- predict(model2, validationset, type = "class")

# Test results on subTesting data set:
confusionMatrix(prediction2, validationset$classe)
```



```
## Confusion Matrix and Statistics
##
##           Reference
## Prediction   A     B     C     D     E
##           A 1394     2     0     0     0
##           B    1  944     3     0     0
##           C    0    3  852     8     1
##           D    0    0    0  795     1
##           E    0    0    0    1  899
##
## Overall Statistics
##
##           Accuracy : 0.9959
##           95% CI : (0.9937, 0.9975)
##           No Information Rate : 0.2845
##           P-Value [Acc > NIR] : < 2.2e-16
##
##           Kappa : 0.9948
##           McNemar's Test P-Value : NA
##
## Statistics by Class:
##
##           Class: A Class: B Class: C Class: D Class: E
## Sensitivity      0.9993  0.9947  0.9965  0.9888  0.9978
## Specificity      0.9994  0.9990  0.9970  0.9998  0.9998
## Pos Pred Value   0.9986  0.9958  0.9861  0.9987  0.9989
## Neg Pred Value   0.9997  0.9987  0.9993  0.9978  0.9995
## Prevalence       0.2845  0.1935  0.1743  0.1639  0.1837
## Detection Rate   0.2843  0.1925  0.1737  0.1621  0.1833
## Detection Prevalence 0.2847  0.1933  0.1762  0.1623  0.1835
## Balanced Accuracy 0.9994  0.9969  0.9968  0.9943  0.9988
```

Conclusion:

model made from random forest gave us 99.6% accuracy while model made from decision tree gave us 74.0% accuracy. So we chose model2 to predict testing data set. Since testing data set only have 20 samples, with 99.6% accuracy we would accurately predict the data.

Use the random forest model to predict testing data.

```
prediction3<- predict(model2, testingset, type="class")
prediction3
```

```
## 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20
## 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
```