

PML_Exercise

Kristen Dardia

July 15, 2017

Peer-graded Assignment: Practical Machine Learning Assignment

Executive Summary

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://groupware.les.inf.puc-rio.br/har> (see the section on the Weight Lifting Exercise Dataset).

The goal of your project is to predict the manner in which they did the exercise. This is the "classe" variable in the training set. You may use any of the other variables to predict with. You should create a report describing how you built your model, how you used cross validation, what you think the expected out of sample error is, and why you made the choices you did. You will also use your prediction model to predict 20 different test cases.

data loading

```
setwd("H:/Desktop/R Programming/PML")
training <- read.csv("pml-training.csv")
testing  <- read.csv("pml-testing.csv")
```

look at data

```
dim(training)

## [1] 19622  160

dim(testing)

## [1]  20 160
```

```
names(training)
```

```
## [1] "X" "user_name"
## [3] "raw_timestamp_part_1" "raw_timestamp_part_2"
## [5] "cvtdd_timestamp" "new_window"
## [7] "num_window" "roll_belt"
## [9] "pitch_belt" "yaw_belt"
## [11] "total_accel_belt" "kurtosis_roll_belt"
## [13] "kurtosis_pitch_belt" "kurtosis_yaw_belt"
## [15] "skewness_roll_belt" "skewness_roll_belt.1"
## [17] "skewness_yaw_belt" "max_roll_belt"
## [19] "max_pitch_belt" "max_yaw_belt"
## [21] "min_roll_belt" "min_pitch_belt"
## [23] "min_yaw_belt" "amplitude_roll_belt"
## [25] "amplitude_pitch_belt" "amplitude_yaw_belt"
## [27] "var_total_accel_belt" "avg_roll_belt"
## [29] "stddev_roll_belt" "var_roll_belt"
## [31] "avg_pitch_belt" "stddev_pitch_belt"
## [33] "var_pitch_belt" "avg_yaw_belt"
## [35] "stddev_yaw_belt" "var_yaw_belt"
## [37] "gyros_belt_x" "gyros_belt_y"
## [39] "gyros_belt_z" "accel_belt_x"
## [41] "accel_belt_y" "accel_belt_z"
## [43] "magnet_belt_x" "magnet_belt_y"
## [45] "magnet_belt_z" "roll_arm"
## [47] "pitch_arm" "yaw_arm"
## [49] "total_accel_arm" "var_accel_arm"
## [51] "avg_roll_arm" "stddev_roll_arm"
## [53] "var_roll_arm" "avg_pitch_arm"
## [55] "stddev_pitch_arm" "var_pitch_arm"
## [57] "avg_yaw_arm" "stddev_yaw_arm"
## [59] "var_yaw_arm" "gyros_arm_x"
## [61] "gyros_arm_y" "gyros_arm_z"
## [63] "accel_arm_x" "accel_arm_y"
## [65] "accel_arm_z" "magnet_arm_x"
## [67] "magnet_arm_y" "magnet_arm_z"
## [69] "kurtosis_roll_arm" "kurtosis_pitch_arm"
## [71] "kurtosis_yaw_arm" "skewness_roll_arm"
## [73] "skewness_pitch_arm" "skewness_yaw_arm"
## [75] "max_roll_arm" "max_pitch_arm"
## [77] "max_yaw_arm" "min_roll_arm"
## [79] "min_pitch_arm" "min_yaw_arm"
## [81] "amplitude_roll_arm" "amplitude_pitch_arm"
## [83] "amplitude_yaw_arm" "roll_dumbbell"
## [85] "pitch_dumbbell" "yaw_dumbbell"
## [87] "kurtosis_roll_dumbbell" "kurtosis_pitch_dumbbell"
## [89] "kurtosis_yaw_dumbbell" "skewness_roll_dumbbell"
## [91] "skewness_pitch_dumbbell" "skewness_yaw_dumbbell"
## [93] "max_roll_dumbbell" "max_pitch_dumbbell"
## [95] "max_yaw_dumbbell" "min_roll_dumbbell"
```

```

## [97] "min_pitch_dumbbell"      "min_yaw_dumbbell"
## [99] "amplitude_roll_dumbbell" "amplitude_pitch_dumbbell"
## [101] "amplitude_yaw_dumbbell"  "total_accel_dumbbell"
## [103] "var_accel_dumbbell"      "avg_roll_dumbbell"
## [105] "stddev_roll_dumbbell"    "var_roll_dumbbell"
## [107] "avg_pitch_dumbbell"      "stddev_pitch_dumbbell"
## [109] "var_pitch_dumbbell"      "avg_yaw_dumbbell"
## [111] "stddev_yaw_dumbbell"     "var_yaw_dumbbell"
## [113] "gyros_dumbbell_x"        "gyros_dumbbell_y"
## [115] "gyros_dumbbell_z"        "accel_dumbbell_x"
## [117] "accel_dumbbell_y"        "accel_dumbbell_z"
## [119] "magnet_dumbbell_x"       "magnet_dumbbell_y"
## [121] "magnet_dumbbell_z"       "roll_forearm"
## [123] "pitch_forearm"           "yaw_forearm"
## [125] "kurtosis_roll_forearm"   "kurtosis_pitch_forearm"
## [127] "kurtosis_yaw_forearm"    "skewness_roll_forearm"
## [129] "skewness_pitch_forearm"  "skewness_yaw_forearm"
## [131] "max_roll_forearm"        "max_pitch_forearm"
## [133] "max_yaw_forearm"         "min_roll_forearm"
## [135] "min_pitch_forearm"       "min_yaw_forearm"
## [137] "amplitude_roll_forearm"  "amplitude_pitch_forearm"
## [139] "amplitude_yaw_forearm"   "total_accel_forearm"
## [141] "var_accel_forearm"       "avg_roll_forearm"
## [143] "stddev_roll_forearm"     "var_roll_forearm"
## [145] "avg_pitch_forearm"       "stddev_pitch_forearm"
## [147] "var_pitch_forearm"       "avg_yaw_forearm"
## [149] "stddev_yaw_forearm"      "var_yaw_forearm"
## [151] "gyros_forearm_x"         "gyros_forearm_y"
## [153] "gyros_forearm_z"         "accel_forearm_x"
## [155] "accel_forearm_y"         "accel_forearm_z"
## [157] "magnet_forearm_x"        "magnet_forearm_y"
## [159] "magnet_forearm_z"        "classe"

```

`names(testing)`

```

## [1] "X"                        "user_name"
## [3] "raw_timestamp_part_1"    "raw_timestamp_part_2"
## [5] "cvtd_timestamp"         "new_window"
## [7] "num_window"             "roll_belt"
## [9] "pitch_belt"             "yaw_belt"
## [11] "total_accel_belt"        "kurtosis_roll_belt"
## [13] "kurtosis_pitch_belt"     "kurtosis_yaw_belt"
## [15] "skewness_roll_belt"      "skewness_roll_belt.1"
## [17] "skewness_yaw_belt"       "max_roll_belt"
## [19] "max_pitch_belt"          "max_yaw_belt"
## [21] "min_roll_belt"           "min_pitch_belt"
## [23] "min_yaw_belt"            "amplitude_roll_belt"
## [25] "amplitude_pitch_belt"    "amplitude_yaw_belt"
## [27] "var_total_accel_belt"    "avg_roll_belt"
## [29] "stddev_roll_belt"        "var_roll_belt"

```

## [31]	"avg_pitch_belt"	"stddev_pitch_belt"
## [33]	"var_pitch_belt"	"avg_yaw_belt"
## [35]	"stddev_yaw_belt"	"var_yaw_belt"
## [37]	"gyros_belt_x"	"gyros_belt_y"
## [39]	"gyros_belt_z"	"accel_belt_x"
## [41]	"accel_belt_y"	"accel_belt_z"
## [43]	"magnet_belt_x"	"magnet_belt_y"
## [45]	"magnet_belt_z"	"roll_arm"
## [47]	"pitch_arm"	"yaw_arm"
## [49]	"total_accel_arm"	"var_accel_arm"
## [51]	"avg_roll_arm"	"stddev_roll_arm"
## [53]	"var_roll_arm"	"avg_pitch_arm"
## [55]	"stddev_pitch_arm"	"var_pitch_arm"
## [57]	"avg_yaw_arm"	"stddev_yaw_arm"
## [59]	"var_yaw_arm"	"gyros_arm_x"
## [61]	"gyros_arm_y"	"gyros_arm_z"
## [63]	"accel_arm_x"	"accel_arm_y"
## [65]	"accel_arm_z"	"magnet_arm_x"
## [67]	"magnet_arm_y"	"magnet_arm_z"
## [69]	"kurtosis_roll_arm"	"kurtosis_pitch_arm"
## [71]	"kurtosis_yaw_arm"	"skewness_roll_arm"
## [73]	"skewness_pitch_arm"	"skewness_yaw_arm"
## [75]	"max_roll_arm"	"max_pitch_arm"
## [77]	"max_yaw_arm"	"min_roll_arm"
## [79]	"min_pitch_arm"	"min_yaw_arm"
## [81]	"amplitude_roll_arm"	"amplitude_pitch_arm"
## [83]	"amplitude_yaw_arm"	"roll_dumbbell"
## [85]	"pitch_dumbbell"	"yaw_dumbbell"
## [87]	"kurtosis_roll_dumbbell"	"kurtosis_pitch_dumbbell"
## [89]	"kurtosis_yaw_dumbbell"	"skewness_roll_dumbbell"
## [91]	"skewness_pitch_dumbbell"	"skewness_yaw_dumbbell"
## [93]	"max_roll_dumbbell"	"max_pitch_dumbbell"
## [95]	"max_yaw_dumbbell"	"min_roll_dumbbell"
## [97]	"min_pitch_dumbbell"	"min_yaw_dumbbell"
## [99]	"amplitude_roll_dumbbell"	"amplitude_pitch_dumbbell"
## [101]	"amplitude_yaw_dumbbell"	"total_accel_dumbbell"
## [103]	"var_accel_dumbbell"	"avg_roll_dumbbell"
## [105]	"stddev_roll_dumbbell"	"var_roll_dumbbell"
## [107]	"avg_pitch_dumbbell"	"stddev_pitch_dumbbell"
## [109]	"var_pitch_dumbbell"	"avg_yaw_dumbbell"
## [111]	"stddev_yaw_dumbbell"	"var_yaw_dumbbell"
## [113]	"gyros_dumbbell_x"	"gyros_dumbbell_y"
## [115]	"gyros_dumbbell_z"	"accel_dumbbell_x"
## [117]	"accel_dumbbell_y"	"accel_dumbbell_z"
## [119]	"magnet_dumbbell_x"	"magnet_dumbbell_y"
## [121]	"magnet_dumbbell_z"	"roll_forearm"
## [123]	"pitch_forearm"	"yaw_forearm"
## [125]	"kurtosis_roll_forearm"	"kurtosis_pitch_forearm"
## [127]	"kurtosis_yaw_forearm"	"skewness_roll_forearm"
## [129]	"skewness_pitch_forearm"	"skewness_yaw_forearm"

```
## [131] "max_roll_forearm"          "max_pitch_forearm"
## [133] "max_yaw_forearm"          "min_roll_forearm"
## [135] "min_pitch_forearm"        "min_yaw_forearm"
## [137] "amplitude_roll_forearm"    "amplitude_pitch_forearm"
## [139] "amplitude_yaw_forearm"     "total_accel_forearm"
## [141] "var_accel_forearm"        "avg_roll_forearm"
## [143] "stddev_roll_forearm"      "var_roll_forearm"
## [145] "avg_pitch_forearm"        "stddev_pitch_forearm"
## [147] "var_pitch_forearm"        "avg_yaw_forearm"
## [149] "stddev_yaw_forearm"       "var_yaw_forearm"
## [151] "gyros_forearm_x"          "gyros_forearm_y"
## [153] "gyros_forearm_z"          "accel_forearm_x"
## [155] "accel_forearm_y"          "accel_forearm_z"
## [157] "magnet_forearm_x"         "magnet_forearm_y"
## [159] "magnet_forearm_z"         "problem_id"
```

data cleaning

Remove useless variables

```
badvar <- c("X", "user_name", "kurtosis_yaw_belt", "skewness_yaw_belt",
"amplitude_yaw_belt", "cvtd_timestamp", "kurtosis_yaw_dumbbell",
"skewness_yaw_dumbbell", "kurtosis_yaw_forearm", "skewness_yaw_forearm")
```

```
training <- training[ , -which(names(training) %in% badvar)]
```

convert to numeric if it should be numeric

```
numvar <- c(
  "kurtosis_roll_belt", "kurtosis_pitch_belt", "skewness_roll_belt",
  "skewness_roll_belt.1", "max_yaw_belt", "min_yaw_belt",
  "kurtosis_roll_arm", "kurtosis_pitch_arm", "kurtosis_yaw_arm",
  "skewness_roll_arm", "skewness_pitch_arm", "skewness_yaw_arm",
  "kurtosis_roll_dumbbell",
  "kurtosis_pitch_dumbbell", "skewness_roll_dumbbell",
  "skewness_pitch_dumbbell", "max_yaw_dumbbell", "min_yaw_dumbbell",
  "kurtosis_roll_forearm", "kurtosis_pitch_forearm", "skewness_roll_forearm",
  "skewness_pitch_forearm", "max_yaw_forearm", "min_yaw_forearm"
)
for (variable in numvar) {
  training[[variable]] <- as.numeric(as.character(training[[variable]]))
}
```

finally, remove NA variables and try to work with just the tidy data

```
training.naCounts <- colSums(sapply(training, is.na))
training.a <- training[, training.naCounts == 0]
```

now do it again for testing dataset

```
testing <- testing[ , -which(names(testing) %in% badvar)]
```

```

for (variable in numvar) {
  testing[[variable]] <- as.numeric(as.character(testing[[variable]]))
}

testing.naCounts <- colSums(sapply(testing, is.na))
testing.a <- testing[,testing.naCounts == 0]

```

Cross Validation

partition the original data into 70% training and 30% testing.

```

set.seed(2017)
library(caret)

## Warning: package 'caret' was built under R version 3.3.3
## Loading required package: lattice
## Loading required package: ggplot2
## Warning: package 'ggplot2' was built under R version 3.3.3

inTrain <- createDataPartition(training.a$classe, p=0.7, list=FALSE )
train    <- training.a[ inTrain, ]
test     <- training.a[ -inTrain, ]
table(train$class)

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

table(test$class)

##
##      A      B      C      D      E
## 1674 1139 1026  964 1082

```

now let's fit a random forest model which is good for classification. make sure to load the caret package

Random Forest model

```

modelFit <- train( classe~., data=train, method="rf", importance=TRUE)

## Loading required package: randomForest
## Warning: package 'randomForest' was built under R version 3.3.3
## 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
```

now let's see how good the model is on the holdout of the training data and also on the test data; fit first and then look at confusion matrix

```
predTrain <- predict( modelFit, train )
predTest  <- predict( modelFit, test )
confusionMatrix( predTrain, train$classe )
```

```
## Confusion Matrix and Statistics
```

```
##
##           Reference
## Prediction    A    B    C    D    E
##           A 3906    0    0    0    0
##           B    0 2658    0    0    0
##           C    0    0 2396    0    0
##           D    0    0    0 2252    0
##           E    0    0    0    0 2525
```

```
## Overall Statistics
```

```
##
##           Accuracy : 1
##           95% CI : (0.9997, 1)
##           No Information Rate : 0.2843
##           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       0.2843  0.1935  0.1744  0.1639  0.1838
## Detection Rate   0.2843  0.1935  0.1744  0.1639  0.1838
## Detection Prevalence 0.2843  0.1935  0.1744  0.1639  0.1838
## Balanced Accuracy 1.0000  1.0000  1.0000  1.0000  1.0000
```

```
confusionMatrix( predTest, test$classe )
```

```
## Confusion Matrix and Statistics
```

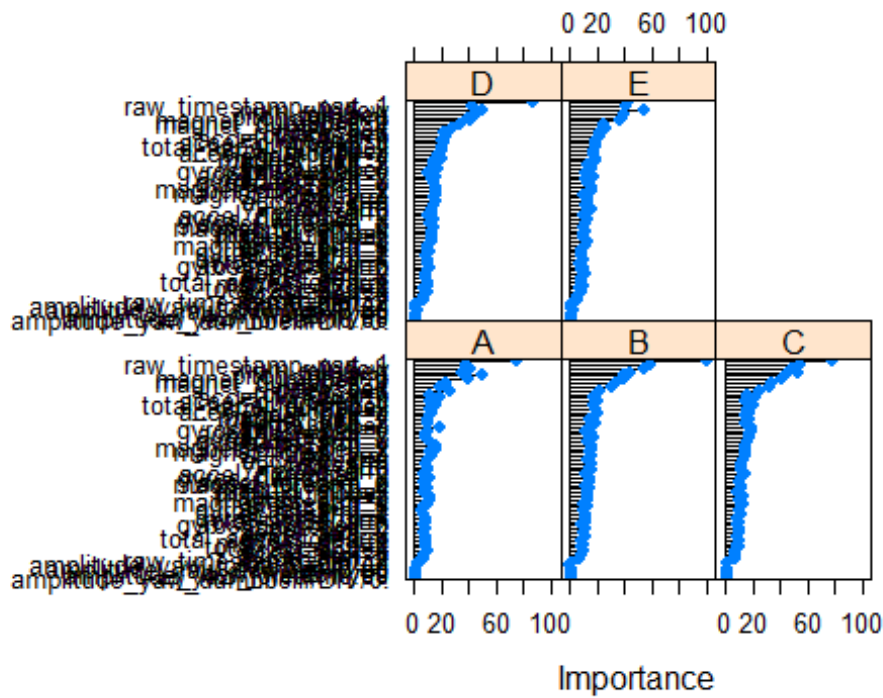
```
##
##           Reference
## Prediction    A    B    C    D    E
##           A 1674    0    0    0    0
##           B    0 1139    1    0    0
```

```
##           C      0      0 1025      1      0
##           D      0      0      0 963      1
##           E      0      0      0      0 1081
##
## Overall Statistics
##
##           Accuracy : 0.9995
##           95% CI : (0.9985, 0.9999)
##           No Information Rate : 0.2845
##           P-Value [Acc > NIR] : < 2.2e-16
##
##           Kappa : 0.9994
##           McNemar's Test P-Value : NA
##
## Statistics by Class:
##
##           Class: A Class: B Class: C Class: D Class: E
## Sensitivity      1.0000   1.0000   0.9990   0.9990   0.9991
## Specificity      1.0000   0.9998   0.9998   0.9998   1.0000
## Pos Pred Value    1.0000   0.9991   0.9990   0.9990   1.0000
## Neg Pred Value    1.0000   1.0000   0.9998   0.9998   0.9998
## Prevalence       0.2845   0.1935   0.1743   0.1638   0.1839
## Detection Rate    0.2845   0.1935   0.1742   0.1636   0.1837
## Detection Prevalence 0.2845   0.1937   0.1743   0.1638   0.1837
## Balanced Accuracy 1.0000   0.9999   0.9994   0.9994   0.9995
```

results: The model has a very good accuracy of 99.95% on the testing set.

finally, let's see which variables actually mattered:

```
variable.importances <- varImp(modelFit)
plot(variable.importances)
```

###looks like just a few mattered; next step here would be to try to limit the number of variables to cut down on processing time. That's for another day :)