

PML_Project

R. Hosek

February 27, 2016

Final Project: Coursera Practical Machine Learning

Executive Summary:

Problem Statement (quoted from the source):

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.

The five possible ‘manners’ are:

- A. exactly according to the specifications
- B. throwing the elbows to the front
- C. lifting the dumbbell only halfway
- D. lowering the dumbbell only halfway
- E. throwing the hips to the front

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

The training data for this project are available here:

<https://d396qusza40orc.cloudfront.net/predmachlearn/pml-training.csv>

The test data are available here:

<https://d396qusza40orc.cloudfront.net/predmachlearn/pml-testing.csv>

The data for this project come from this source: <http://groupware.les.inf.puc-rio.br/har>.

Training and Testing data files previously downloaded to working directory from above-cited urls

Get the data

```
w <- setwd("~/GitHub/practical-machine-learning")

train_0 <- read.csv("pml-training.csv", header = TRUE, na.strings = c('NA',
'', '#DIV/0!'))
test_0 <- read.csv("pml-testing.csv", header = TRUE, na.strings = c('NA',
'', '#DIV/0!'))
```

Explore and Clean the Data

First, check to see if two data sets have same variables (except for the last one, which we know is different)

```
n_tr <- names(train_0)
n_te <- names(test_0)

all.equal(n_tr[length(n_tr)-1], n_te[length(n_te)-1])
## [1] TRUE
```

That looks fine; now preview the data

```
str(train_0)
## '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
##  $ raw_timestamp_part_1 : int  1323084231 1323084231 1323084231
##  $ raw_timestamp_part_2 : int  788290 808298 820366 120339 196328
##  $ raw_timestamp_part_3 : int  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 ...
## $ 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 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 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 ...
## $ 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    : num    NA NA NA NA NA NA NA NA NA NA ...
## $ kurtosis_pitch_arm   : num    NA NA NA NA NA NA NA NA NA NA ...
## $ kurtosis_yaw_arm     : num    NA NA NA NA NA NA NA NA NA NA ...
## $ skewness_roll_arm    : num    NA NA NA NA NA NA NA NA NA NA ...
## $ skewness_pitch_arm   : num    NA NA NA NA NA NA NA NA NA NA ...
## $ skewness_yaw_arm     : num    NA NA NA NA NA NA NA NA NA NA ...
## $ 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 : num    NA NA NA NA NA NA NA NA NA NA ...
## $ kurtosis_pitch_dumbbell : num    NA NA NA NA NA NA NA NA NA NA ...
## $ kurtosis_yaw_dumbbell : logi    NA NA NA NA NA NA ...
## $ skewness_roll_dumbbell : num    NA NA NA NA NA NA NA NA NA NA ...
## $ skewness_pitch_dumbbell : num    NA NA NA NA NA NA NA NA NA NA ...
## $ skewness_yaw_dumbbell : logi    NA NA NA NA NA NA ...
## $ max_roll_dumbbell    : num    NA NA NA NA NA NA NA NA NA NA ...
## $ max_pitch_dumbbell   : num    NA NA NA NA NA NA NA NA NA NA ...
## $ max_yaw_dumbbell     : num    NA NA NA NA NA NA NA NA NA NA ...
## $ min_roll_dumbbell    : num    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]
```

There appear to be many NA values; let's remove these

```
NAs <- apply(train_0, 2, function(x) {sum(is.na(x))})
train_1 <- train_0[which(NAs == 0)]
```

```
NAs <- apply(test_0, 2, function(x) {sum(is.na(x))})
test_1 <- test_0[which(NAs == 0)]
```

Next, remove the first 7 variables which will not be used

```
train_1 <- train_1[, 8:length(colnames(train_1))]
test_1 <- test_1[, 8:length(colnames(test_1))]
```

Check for congruity

```
names(train_1)
## [1] "roll_belt"      "pitch_belt"      "yaw_belt"
## [4] "total_accel_belt" "gyros_belt_x"    "gyros_belt_y"
## [7] "gyros_belt_z"    "accel_belt_x"    "accel_belt_y"
## [10] "accel_belt_z"    "magnet_belt_x"   "magnet_belt_y"
## [13] "magnet_belt_z"   "roll_arm"        "pitch_arm"
## [16] "yaw_arm"         "total_accel_arm" "gyros_arm_x"
## [19] "gyros_arm_y"     "gyros_arm_z"     "accel_arm_x"
## [22] "accel_arm_y"     "accel_arm_z"     "magnet_arm_x"
## [25] "magnet_arm_y"    "magnet_arm_z"    "roll_dumbbell"
## [28] "pitch_dumbbell"  "yaw_dumbbell"    "total_accel_dumbbell"
## [31] "gyros_dumbbell_x" "gyros_dumbbell_y" "gyros_dumbbell_z"
## [34] "accel_dumbbell_x" "accel_dumbbell_y" "accel_dumbbell_z"
## [37] "magnet_dumbbell_x" "magnet_dumbbell_y" "magnet_dumbbell_z"
## [40] "roll_forearm"    "pitch_forearm"   "yaw_forearm"
## [43] "total_accel_forearm" "gyros_forearm_x" "gyros_forearm_y"
## [46] "gyros_forearm_z" "accel_forearm_x" "accel_forearm_y"
## [49] "accel_forearm_z" "magnet_forearm_x" "magnet_forearm_y"
## [52] "magnet_forearm_z" "classe"

names(test_1)
## [1] "roll_belt"      "pitch_belt"      "yaw_belt"
## [4] "total_accel_belt" "gyros_belt_x"    "gyros_belt_y"
## [7] "gyros_belt_z"    "accel_belt_x"    "accel_belt_y"
## [10] "accel_belt_z"    "magnet_belt_x"   "magnet_belt_y"
## [13] "magnet_belt_z"   "roll_arm"        "pitch_arm"
## [16] "yaw_arm"         "total_accel_arm" "gyros_arm_x"
## [19] "gyros_arm_y"     "gyros_arm_z"     "accel_arm_x"
## [22] "accel_arm_y"     "accel_arm_z"     "magnet_arm_x"
## [25] "magnet_arm_y"    "magnet_arm_z"    "roll_dumbbell"
## [28] "pitch_dumbbell"  "yaw_dumbbell"    "total_accel_dumbbell"
## [31] "gyros_dumbbell_x" "gyros_dumbbell_y" "gyros_dumbbell_z"
## [34] "accel_dumbbell_x" "accel_dumbbell_y" "accel_dumbbell_z"
## [37] "magnet_dumbbell_x" "magnet_dumbbell_y" "magnet_dumbbell_z"
## [40] "roll_forearm"    "pitch_forearm"   "yaw_forearm"
## [43] "total_accel_forearm" "gyros_forearm_x" "gyros_forearm_y"
```

```
## [46] "gyros_forearm_z"      "accel_forearm_x"      "accel_forearm_y"
## [49] "accel_forearm_z"      "magnet_forearm_x"      "magnet_forearm_y"
## [52] "magnet_forearm_z"      "problem_id"
```

This checks out OK.

Load Caret

Preprocessing: Imputation of missing data and checking for near-zero variance

Applies only to numeric variables

```
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
var_num <- which(lapply(train_1, class) %in% "numeric")
var_preproc <- preProcess(train_1[,var_num], method = c("knnImpute",
"center", "scale"))
train_2 <- predict(var_preproc, train_1[,var_num])
train_2$classe <- train_1$classe
test_2 <- predict(var_preproc, test_1[, var_num])

nzv_train <- nearZeroVar(train_2, saveMetrics = TRUE)
train_3 <- train_2[, nzv_train$nzv == FALSE]

nzv_test <- nearZeroVar(test_2, saveMetrics = TRUE)
test_3 <- test_2[, nzv_test$nzv == FALSE]

dim(train_3)
## [1] 19622      28
dim(test_3)
## [1] 20 27
```

Near-zero variance not an issue; all variables retained

Next, partition training set to use cross-validation

```
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(caret)
library(e1071)
## Warning: package 'e1071' was built under R version 3.1.3
set.seed(92542)
train_tot <- createDataPartition(train_3$classe, p = .75, list=FALSE)
train_tr <- train_3[train_tot,]
test_tr <- train_3[-train_tot,]
```

Based on the lectures, the forum and reading, RANDOM FORESTS seems to be the best initial approach

```
model_1<-train(classe~.,data=train_tr, method="rf",
trControl=trainControl(method="cv"), number=5, allowParallel=TRUE)

print(model_1)
## Random Forest
##
## 14718 samples
##    27 predictor
##    5 classes: 'A', 'B', 'C', 'D', 'E'
##
## No pre-processing
## Resampling: Cross-Validated (10 fold)
## Summary of sample sizes: 13246, 13245, 13247, 13247, 13246, 13246, ...
## Resampling results across tuning parameters:
##
##   mtry Accuracy      Kappa      Accuracy SD   Kappa SD
##    2   0.9922543  0.9902011  0.002525836  0.003196485
##   14   0.9922543  0.9902018  0.002973869  0.003761870
##   27   0.9896721  0.9869359  0.004062373  0.005139260
##
## Accuracy was used to select the optimal model using the largest value.
## The final value used for the model was mtry = 14.
```

What is the prediction accuracy for the training and cross-validation test sets?

First, for the training set

```
pred_tr <- predict(model_1, train_tr)
confusionMatrix(pred_tr, train_tr$classe)
## Confusion Matrix and Statistics
##
##              Reference
## Prediction      A      B      C      D      E
##           A 4185      0      0      0      0
##           B      0 2848      0      0      0
##           C      0      0 2567      0      0
##           D      0      0      0 2412      0
##           E      0      0      0      0 2706
##
## 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    0.1839
## Detection Rate       0.2843    0.1935    0.1744    0.1639    0.1839
## Detection Prevalence 0.2843    0.1935    0.1744    0.1639    0.1839
## Balanced Accuracy     1.0000    1.0000    1.0000    1.0000    1.0000
```

Next, for the cross-validation test set

```
pred_cv_test <- predict(model_1, test_tr)
confusionMatrix(pred_cv_test, test_tr$classe)
## Confusion Matrix and Statistics
##
##              Reference
## Prediction      A      B      C      D      E
##      A 1395      5      0      0      0
##      B      0  944      9      1      3
##      C      0      0  837      5      1
##      D      0      0      9  796      1
##      E      0      0      0      2  896
##
## Overall Statistics
##
##              Accuracy : 0.9927
##              95% CI : (0.9899, 0.9949)
##      No Information Rate : 0.2845
##      P-Value [Acc > NIR] : < 2.2e-16
##
##              Kappa : 0.9907
##      McNemar's Test P-Value : NA
##
## Statistics by Class:
##
##              Class: A Class: B Class: C Class: D Class: E
## Sensitivity          1.0000    0.9947    0.9789    0.9900    0.9945
## Specificity          0.9986    0.9967    0.9985    0.9976    0.9995
## Pos Pred Value       0.9964    0.9864    0.9929    0.9876    0.9978
## Neg Pred Value       1.0000    0.9987    0.9956    0.9980    0.9988
## Prevalence           0.2845    0.1935    0.1743    0.1639    0.1837
## Detection Rate       0.2845    0.1925    0.1707    0.1623    0.1827
## Detection Prevalence 0.2855    0.1951    0.1719    0.1644    0.1831
## Balanced Accuracy     0.9993    0.9957    0.9887    0.9938    0.9970
```

Make prediction of the real test set

```
pred_real_test <- predict(model_1, test_3)

pred_real_test
## [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
```