

# Practical Machine Learning Course Project

*Oleg Rybkin*

*October 25, 2016*

## 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).

## 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>. If you use the document you create for this class for any purpose please cite them as they have been very generous in allowing their data to be used for this kind of assignment.

## Data Processing

Adding libraries

```
library(caret)
```

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

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

```
library(randomForest)
```

```
## 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
```

```
library(rpart)
```

Loading training and testing datasets

```
training <-read.csv("https://d396qusza40orc.cloudfront.net/predmachlearn/pml-training.csv", na.strings=
test <-read.csv("https://d396qusza40orc.cloudfront.net/predmachlearn/pml-testing.csv", na.strings=c("NA
```

Remove columns with all missing values

```
training<-training[,colSums(is.na(training)) == 0]
test <-test[,colSums(is.na(test)) == 0]
```

Remove columns with irrelevant values

```
training <-training[,-c(1:7)]
test <-test[,-c(1:7)]
```

Check the new datasets

```
str(training)
```

```
## 'data.frame': 19622 obs. of 53 variables:
## $ 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 ...
## $ 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 ...
## $ 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 ...
## $ 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 ...
```

```
## $ total_accel_dumbbell: int 37 37 37 37 37 37 37 37 37 37 ...
## $ gyros_dumbbell_x : num 0 0 0 0 0 0 0 0 0 0 ...
## $ gyros_dumbbell_y : num -0.02 -0.02 -0.02 -0.02 -0.02 -0.02 -0.02 -0.02 -0.02 -0.02 ...
## $ gyros_dumbbell_z : num 0 0 0 -0.02 0 0 0 0 0 0 ...
## $ accel_dumbbell_x : int -234 -233 -232 -232 -233 -234 -232 -234 -232 -235 ...
## $ accel_dumbbell_y : int 47 47 46 48 48 48 47 46 47 48 ...
## $ accel_dumbbell_z : int -271 -269 -270 -269 -270 -269 -270 -272 -269 -270 ...
## $ magnet_dumbbell_x : int -559 -555 -561 -552 -554 -558 -551 -555 -549 -558 ...
## $ magnet_dumbbell_y : int 293 296 298 303 292 294 295 300 292 291 ...
## $ magnet_dumbbell_z : num -65 -64 -63 -60 -68 -66 -70 -74 -65 -69 ...
## $ roll_forearm : num 28.4 28.3 28.3 28.1 28 27.9 27.9 27.8 27.7 27.7 ...
## $ pitch_forearm : num -63.9 -63.9 -63.9 -63.9 -63.9 -63.9 -63.9 -63.8 -63.8 -63.8 ...
## $ yaw_forearm : num -153 -153 -152 -152 -152 -152 -152 -152 -152 -152 ...
## $ total_accel_forearm : int 36 36 36 36 36 36 36 36 36 36 ...
## $ gyros_forearm_x : num 0.03 0.02 0.03 0.02 0.02 0.02 0.02 0.02 0.03 0.02 ...
## $ gyros_forearm_y : num 0 0 -0.02 -0.02 0 -0.02 0 -0.02 0 0 ...
## $ gyros_forearm_z : num -0.02 -0.02 0 0 -0.02 -0.03 -0.02 0 -0.02 -0.02 ...
## $ accel_forearm_x : int 192 192 196 189 189 193 195 193 193 190 ...
## $ accel_forearm_y : int 203 203 204 206 206 203 205 205 204 205 ...
## $ accel_forearm_z : int -215 -216 -213 -214 -214 -215 -215 -213 -214 -215 ...
## $ magnet_forearm_x : int -17 -18 -18 -16 -17 -9 -18 -9 -16 -22 ...
## $ magnet_forearm_y : num 654 661 658 658 655 660 659 660 653 656 ...
## $ magnet_forearm_z : num 476 473 469 469 473 478 470 474 476 473 ...
## $ classe : Factor w/ 5 levels "A","B","C","D",...: 1 1 1 1 1 1 1 1 1 1 ...
```

```
str(test)
```

```
## 'data.frame': 20 obs. of 53 variables:
## $ roll_belt : num 123 1.02 0.87 125 1.35 -5.92 1.2 0.43 0.93 114 ...
## $ pitch_belt : num 27 4.87 1.82 -41.6 3.33 1.59 4.44 4.15 6.72 22.4 ...
## $ yaw_belt : num -4.75 -88.9 -88.5 162 -88.6 -87.7 -87.3 -88.5 -93.7 -13.1 ...
## $ total_accel_belt : int 20 4 5 17 3 4 4 4 4 18 ...
## $ gyros_belt_x : num -0.5 -0.06 0.05 0.11 0.03 0.1 -0.06 -0.18 0.1 0.14 ...
## $ gyros_belt_y : num -0.02 -0.02 0.02 0.11 0.02 0.05 0 -0.02 0 0.11 ...
## $ gyros_belt_z : num -0.46 -0.07 0.03 -0.16 0 -0.13 0 -0.03 -0.02 -0.16 ...
## $ accel_belt_x : int -38 -13 1 46 -8 -11 -14 -10 -15 -25 ...
## $ accel_belt_y : int 69 11 -1 45 4 -16 2 -2 1 63 ...
## $ accel_belt_z : int -179 39 49 -156 27 38 35 42 32 -158 ...
## $ magnet_belt_x : int -13 43 29 169 33 31 50 39 -6 10 ...
## $ magnet_belt_y : int 581 636 631 608 566 638 622 635 600 601 ...
## $ magnet_belt_z : int -382 -309 -312 -304 -418 -291 -315 -305 -302 -330 ...
## $ roll_arm : num 40.7 0 0 -109 76.1 0 0 0 -137 -82.4 ...
## $ pitch_arm : num -27.8 0 0 55 2.76 0 0 0 11.2 -63.8 ...
## $ yaw_arm : num 178 0 0 -142 102 0 0 0 -167 -75.3 ...
## $ total_accel_arm : int 10 38 44 25 29 14 15 22 34 32 ...
## $ gyros_arm_x : num -1.65 -1.17 2.1 0.22 -1.96 0.02 2.36 -3.71 0.03 0.26 ...
## $ gyros_arm_y : num 0.48 0.85 -1.36 -0.51 0.79 0.05 -1.01 1.85 -0.02 -0.5 ...
## $ gyros_arm_z : num -0.18 -0.43 1.13 0.92 -0.54 -0.07 0.89 -0.69 -0.02 0.79 ...
## $ accel_arm_x : int 16 -290 -341 -238 -197 -26 99 -98 -287 -301 ...
## $ accel_arm_y : int 38 215 245 -57 200 130 79 175 111 -42 ...
## $ accel_arm_z : int 93 -90 -87 6 -30 -19 -67 -78 -122 -80 ...
## $ magnet_arm_x : int -326 -325 -264 -173 -170 396 702 535 -367 -420 ...
## $ magnet_arm_y : int 385 447 474 257 275 176 15 215 335 294 ...
## $ magnet_arm_z : int 481 434 413 633 617 516 217 385 520 493 ...
```

```
## $ roll_dumbbell      : num -17.7 54.5 57.1 43.1 -101.4 ...
## $ pitch_dumbbell     : num 25 -53.7 -51.4 -30 -53.4 ...
## $ yaw_dumbbell       : num 126.2 -75.5 -75.2 -103.3 -14.2 ...
## $ total_accel_dumbbell: int 9 31 29 18 4 29 29 29 3 2 ...
## $ gyros_dumbbell_x   : num 0.64 0.34 0.39 0.1 0.29 -0.59 0.34 0.37 0.03 0.42 ...
## $ gyros_dumbbell_y   : num 0.06 0.05 0.14 -0.02 -0.47 0.8 0.16 0.14 -0.21 0.51 ...
## $ gyros_dumbbell_z   : num -0.61 -0.71 -0.34 0.05 -0.46 1.1 -0.23 -0.39 -0.21 -0.03 ...
## $ accel_dumbbell_x   : int 21 -153 -141 -51 -18 -138 -145 -140 0 -7 ...
## $ accel_dumbbell_y   : int -15 155 155 72 -30 166 150 159 25 -20 ...
## $ accel_dumbbell_z   : int 81 -205 -196 -148 -5 -186 -190 -191 9 7 ...
## $ magnet_dumbbell_x  : int 523 -502 -506 -576 -424 -543 -484 -515 -519 -531 ...
## $ magnet_dumbbell_y  : int -528 388 349 238 252 262 354 350 348 321 ...
## $ magnet_dumbbell_z  : int -56 -36 41 53 312 96 97 53 -32 -164 ...
## $ roll_forearm       : num 141 109 131 0 -176 150 155 -161 15.5 13.2 ...
## $ pitch_forearm      : num 49.3 -17.6 -32.6 0 -2.16 1.46 34.5 43.6 -63.5 19.4 ...
## $ yaw_forearm        : num 156 106 93 0 -47.9 89.7 152 -89.5 -139 -105 ...
## $ total_accel_forearm: int 33 39 34 43 24 43 32 47 36 24 ...
## $ gyros_forearm_x    : num 0.74 1.12 0.18 1.38 -0.75 -0.88 -0.53 0.63 0.03 0.02 ...
## $ gyros_forearm_y    : num -3.34 -2.78 -0.79 0.69 3.1 4.26 1.8 -0.74 0.02 0.13 ...
## $ gyros_forearm_z    : num -0.59 -0.18 0.28 1.8 0.8 1.35 0.75 0.49 -0.02 -0.07 ...
## $ accel_forearm_x    : int -110 212 154 -92 131 230 -192 -151 195 -212 ...
## $ accel_forearm_y    : int 267 297 271 406 -93 322 170 -331 204 98 ...
## $ accel_forearm_z    : int -149 -118 -129 -39 172 -144 -175 -282 -217 -7 ...
## $ magnet_forearm_x   : int -714 -237 -51 -233 375 -300 -678 -109 0 -403 ...
## $ magnet_forearm_y   : int 419 791 698 783 -787 800 284 -619 652 723 ...
## $ magnet_forearm_z   : int 617 873 783 521 91 884 585 -32 469 512 ...
## $ problem_id         : int 1 2 3 4 5 6 7 8 9 10 ...
```

Partition dataset in 2 parts, 75% and 25%

```
subsamples <- createDataPartition(y=training$classe, p=0.75, list=FALSE)
subTraining <- training[subsamples, ]
subTesting <- training[-subsamples, ]
```

Use Decision Tree Model for prediction

```
modell1 <- rpart(classe ~ ., data=subTraining, method="class")

prediction1 <- predict(modell1, subTesting, type = "class")
table(prediction1, subTesting$classe)
```

```
##
## prediction1      A      B      C      D      E
##           A 1269  162    17    53    16
##           B   49  530    59    66    71
##           C   38  145   705   120   111
##           D   13   74    67   510    61
##           E    26   38     7    55   642
```

```
confusionMatrix(prediction1, subTesting$classe)
```

## Confusion Matrix and Statistics

```
##
##           Reference
## Prediction    A    B    C    D    E
##           A 1269  162   17   53   16
##           B   49  530   59   66   71
##           C   38  145  705  120  111
##           D   13   74   67  510   61
##           E   26   38    7   55  642
##
## Overall Statistics
##
##           Accuracy : 0.7455
##           95% CI : (0.7331, 0.7577)
##           No Information Rate : 0.2845
##           P-Value [Acc > NIR] : < 2.2e-16
##
##           Kappa : 0.6773
##           McNemar's Test P-Value : < 2.2e-16
##
## Statistics by Class:
##
##           Class: A Class: B Class: C Class: D Class: E
## Sensitivity      0.9097  0.5585  0.8246  0.6343  0.7125
## Specificity      0.9293  0.9381  0.8978  0.9476  0.9685
## Pos Pred Value   0.8365  0.6839  0.6300  0.7034  0.8359
## Neg Pred Value   0.9628  0.8985  0.9604  0.9296  0.9374
## Prevalence       0.2845  0.1935  0.1743  0.1639  0.1837
## Detection Rate   0.2588  0.1081  0.1438  0.1040  0.1309
## Detection Prevalence 0.3093  0.1580  0.2282  0.1478  0.1566
## Balanced Accuracy 0.9195  0.7483  0.8612  0.7909  0.8405
```

Use Random Forest Model for prediction

```
model2 <- randomForest(classe ~. , data=subTraining, method="class")
prediction2 <- predict(model2, subTesting, type = "class")
table(prediction2, subTesting$classe)
```

```
##
## prediction2    A    B    C    D    E
##           A 1395    7    0    0    0
##           B    0  938    7    0    0
##           C    0    4  846   12    0
##           D    0    0    2  791    3
##           E    0    0    0    1  898
```

```
confusionMatrix(prediction2, subTesting$classe)
```

```
## Confusion Matrix and Statistics
##
##           Reference
## Prediction    A    B    C    D    E
##           A 1395    7    0    0    0
```

```

##           B      0  938      7      0      0
##           C      0   4  846     12      0
##           D      0   0   2  791      3
##           E      0   0   0   1  898
##
## 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  0.9884  0.9895  0.9838  0.9967
## Specificity      0.9980  0.9982  0.9960  0.9988  0.9998
## Pos Pred Value   0.9950  0.9926  0.9814  0.9937  0.9989
## Neg Pred Value    1.0000  0.9972  0.9978  0.9968  0.9993
## Prevalence       0.2845  0.1935  0.1743  0.1639  0.1837
## Detection Rate   0.2845  0.1913  0.1725  0.1613  0.1831
## Detection Prevalence 0.2859  0.1927  0.1758  0.1623  0.1833
## Balanced Accuracy 0.9990  0.9933  0.9928  0.9913  0.9982

```

## Conclusion

Random Forest algorithm performed better than Decision Trees.