

# Machine\_Learning\_Project2

*Dennis Chakey*

*Thursday, June 18, 2015*

## Practical Machine Learning Course Project

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

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 this project is to predict the manner of performing unilateral dumbbell biceps curls based on data from accelerometers on the belt, forearm, arm, and dumbbell of 6 participants. The 5 possible methods include - .A: exactly according to the specification .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

```
## Loading required package: lattice
## Loading required package: ggplot2
## Loading required package: foreach
## Loading required package: iterators
## Loading required package: parallel
library(Hmisc)

## Warning: package 'Hmisc' was built under R version 3.1.3

## Loading required package: grid
## Loading required package: lattice
## Loading required package: survival
## Loading required package: splines
## Loading required package: Formula

## Warning: package 'Formula' was built under R version 3.1.3

## Loading required package: ggplot2
##
## Attaching package: 'Hmisc'
##
## The following object is masked from 'package:R.utils':
##
##     capitalize
##
## The following objects are masked from 'package:base':
##
##     format.pval, round.POSIXt, trunc.POSIXt, units

library(caret)

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

##
## Attaching package: 'caret'
##
```

```
## The following object is masked by_ '.GlobalEnv':
##
##      best
##
## The following object is masked from 'package:survival':
##
##      cluster

library(randomForest)

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

## randomForest 4.6-10
## Type rfNews() to see new features/changes/bug fixes.
##
## Attaching package: 'randomForest'
##
## The following object is masked from 'package:Hmisc':
##
##      combine

library(foreach)

## Warning: package 'foreach' was built under R version 3.1.3

library(doParallel)

## Warning: package 'doParallel' was built under R version 3.1.3

## Loading required package: iterators
## Loading required package: parallel
```

**The pml-training.csv data is used to create training and testing sets for fitting the model.**

**The pml-test.csv data is used to submit 20 test cases based on the fitted model.**

Read and Cleanse the data files All blank(''), '#DIV/0' and 'NA' values are converted to 'NA'.

```
trainingRaw <- read.csv(file="pml-training.csv", header=TRUE, as.is = TRUE, stringsAsFactors =
FALSE, sep=',', na.strings=c('NA','','#DIV/0!'))
testingRaw <- read.csv(file="pml-testing.csv", header=TRUE, as.is = TRUE, stringsAsFactors =
FALSE, sep=',', na.strings=c('NA','','#DIV/0!'))

trainingRaw$classe <- as.factor(trainingRaw$classe)
```

**Display the data and verify data integrity**

```
str(trainingRaw)

## 'data.frame':    19622 obs. of  160 variables:
## $ X                : int  1 2 3 4 5 6 7 8 9 10 ...
## $ user_name        : chr  "carlitos" "carlitos" "carlitos" "carlitos" ...
## $ 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     : chr  "05/12/2011 11:23" "05/12/2011 11:23" "05/12/2011 11:23"
"05/12/2011 11:23" ...
## $ new_window         : chr  "no" "no" "no" "no" ...
## $ 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 NA ...
```

##	\$	skewness_roll_belt	:	num	NA	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	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	NA	...
##	\$	max_picth_belt	:	int	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	...
##	\$	max_yaw_belt	:	num	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	...
##	\$	min_roll_belt	:	num	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	...
##	\$	min_pitch_belt	:	int	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	...
##	\$	min_yaw_belt	:	num	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	...
##	\$	amplitude_roll_belt	:	num	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	...
##	\$	amplitude_pitch_belt	:	int	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	...
##	\$	amplitude_yaw_belt	:	num	NA	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	NA	...
##	\$	avg_roll_belt	:	num	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	...
##	\$	stddev_roll_belt	:	num	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	...
##	\$	var_roll_belt	:	num	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	...
##	\$	avg_pitch_belt	:	num	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	...
##	\$	stddev_pitch_belt	:	num	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	...
##	\$	var_pitch_belt	:	num	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	...
##	\$	avg_yaw_belt	:	num	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	...
##	\$	stddev_yaw_belt	:	num	NA	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	0.02	0.03	...
##	\$	gyros_belt_y	:	num	0	0	0	0.02	0	0	0	0	...			
##	\$	gyros_belt_z	:	num	-0.02	-0.02	-0.02	-0.02	-0.03	-0.02	-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	-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	-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	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	0.02	...
##	\$	gyros_arm_y	:	num	0	-0.02	-0.02	-0.03	-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	0.02	0	0	0	-0.02	-0.02	...	
##	\$	accel_arm_x	:	int	-288	-290	-289	-289	-289	-289	-289	-289	-289	-288	-288	...
##	\$	accel_arm_y	:	int	109	110	110	111	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	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_picth_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_picth_dumbbell	:	num	NA	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	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	...					
##	\$	max_roll_dumbbell	:	num	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	NA	...
##	\$	max_picth_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]
```

**Discard the irrelevant variables (non-accelerometer measures) and additional problem data fields like the invalid dates**

**Compute the prediction only on the accelerometer data values of belt, forearm, arm and dumbbell.**

```
NAindex <- apply(trainingRaw,2,function(x) {sum(is.na(x))})
trainingRaw <- trainingRaw[,which(NAindex == 0)]
NAindex <- apply(testingRaw,2,function(x) {sum(is.na(x))})
testingRaw <- testingRaw[,which(NAindex == 0)]
```

## Preprocessing variables

```
v <- which(lapply(trainingRaw, class) %in% "numeric")

preObj <-preProcess(trainingRaw[,v],method=c('knnImpute', 'center', 'scale'))
trainLess1 <- predict(preObj, trainingRaw[,v])
trainLess1$classe <- trainingRaw$classe

testLess1 <-predict(preObj,testingRaw[,v])
```

**Eliminate the non zero variables to improve accuracy**

```
nzv <- nearZeroVar(trainLess1,saveMetrics=TRUE)
trainLess1 <- trainLess1[,nzv$nzv==FALSE]

nzv <- nearZeroVar(testLess1,saveMetrics=TRUE)
testLess1 <- testLess1[,nzv$nzv==FALSE]
```

## Create cross validation set

The training set is divided in two parts, one for training and the other for cross validation

```
set.seed(20150618)
inTrain = createDataPartition(trainLess1$classe, p = 3/4, list=FALSE)
training = trainLess1[inTrain,]
crossValidation = trainLess1[-inTrain,]
```

## Training Model

Train model with Random Forest technique to achieve a greater accuracy level. The model is build on a training set of 28 variables from the initial 160. Cross validation is implemented as the train control method of choice.

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

```
## 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: 13245, 13246, 13246, 13246, 13247, 13247, ...
##
## Resampling results across tuning parameters:
##
##      mtry  Accuracy   Kappa      Accuracy SD   Kappa SD
##      2    0.9937485  0.9920924  0.002122178   0.002684251
##     14    0.9924580  0.9904603  0.002321199   0.002936228
##     27    0.9916428  0.9894292  0.002482747   0.003140459
##
## Accuracy was used to select the optimal model using  the largest value.
## The final value used for the model was mtry = 2.
```

## Training Prediction Set

```
trainingPred <- predict(modFit, training)
confusionMatrix(trainingPred, training$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
```

## Cross Validation Set

```
cvPred <- predict(modFit, crossValidation)
confusionMatrix(cvPred, crossValidation$classe)

## Confusion Matrix and Statistics
##
##              Reference
## Prediction      A      B      C      D      E
##      A 1391      11      0      0      0
##      B      3  930      7      0      0
##      C      0      8  844     11      0
##      D      0      0      4    791      1
##      E      1      0      0      2    900
##
## Overall Statistics
##
##              Accuracy : 0.9902
##              95% CI : (0.987, 0.9928)
##      No Information Rate : 0.2845
##      P-Value [Acc > NIR] : < 2.2e-16
##
##              Kappa : 0.9876
##      McNemar's Test P-Value : NA
##
## Statistics by Class:
```

```
##
##
##      Class: A Class: B Class: C Class: D Class: E
## Sensitivity      0.9971  0.9800  0.9871  0.9838  0.9989
## Specificity      0.9969  0.9975  0.9953  0.9988  0.9993
## Pos Pred Value   0.9922  0.9894  0.9780  0.9937  0.9967
## Neg Pred Value   0.9989  0.9952  0.9973  0.9968  0.9998
## Prevalence       0.2845  0.1935  0.1743  0.1639  0.1837
## Detection Rate   0.2836  0.1896  0.1721  0.1613  0.1835
## Detection Prevalence 0.2859  0.1917  0.1760  0.1623  0.1841
## Balanced Accuracy 0.9970  0.9887  0.9912  0.9913  0.9991
```

## Predictions on the Real Testing Set

```
testingPrediction <- predict(modFit, testLess1)
testingPrediction

##      [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
```

## Create the Answer Submittal Files Function

```
pml_write_files = function(x){
  n = length(x)
  for(i in 1:n){
    filename = paste0("problem_id_",i,".txt")
    write.table(x[i],file=filename,quote=FALSE,row.names=FALSE,col.names=FALSE)
  }
}
```

## Create the Answer Submittal Files

```
answers <- predict(modFit, testLess1)
pml_write_files(answers)
```