# PML-writeup

*Caroliem*

*Sunday, January 18, 2015*

## introduction & summary

This report predict the manner in which person performed their exercise, using data from accelerometers on the belt, forearm, arm, and dumbell of 6 participants. The data for this project come from this source: http://groupware.les.inf.puc-rio.br/har.

## Data loading and analysis

First loading required packages and data.

```
library(caret)
```

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

Then load data and do some exploratory data analysis. The files can be downloaded from: - https://d396qusza40orc.cloudfront.net/predmachlearn/pml-training.csv for training - https://d396qusza40orc.cloudfront.net/predmachlearn/pml-testing.csv for testing.

```
train <- read.csv('pml-training.csv' , na.strings=c("NA",""))
test <- read.csv('pml-testing.csv' , na.strings=c("NA",""))

## explorary data analysis
str(train)
```
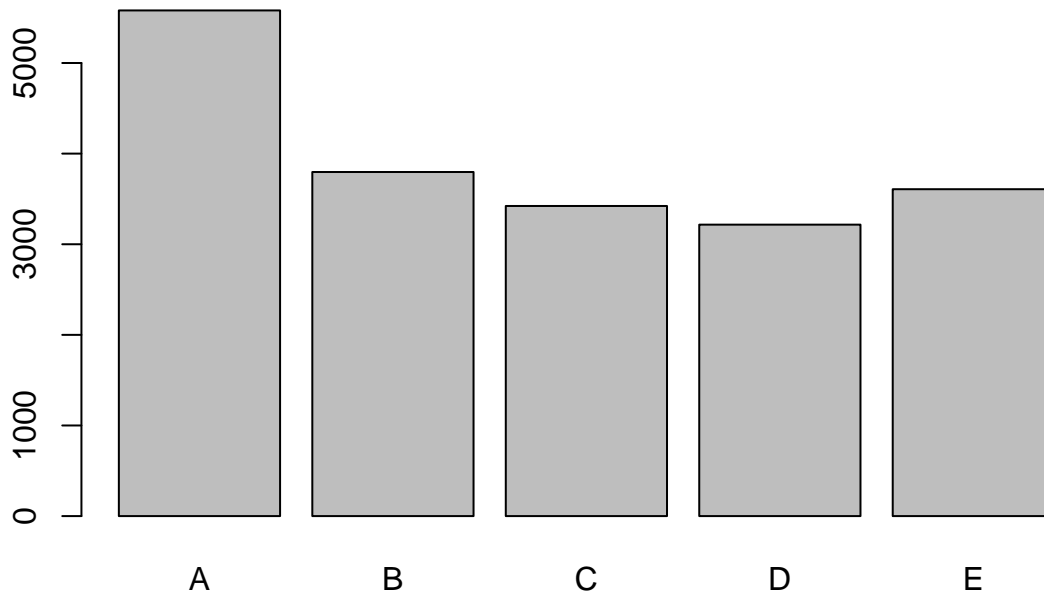
```
## 'data.frame':    19622 obs. of  160 variables:
##  $ X                   : int  1 2 3 4 5 6 7 8 9 10 ...
##  $ user_name           : Factor w/ 6 levels "adelmo","carlitos",..: 2 2 2 2 2 2 2 2 2 2 ...
##  $ raw_timestamp_part_1: int  1323084231 1323084231 1323084231 1323084232 1323084232 1323084232
##  $ raw_timestamp_part_2: int  788290 808298 820366 120339 196328 304277 368296 440390 484323 484
##  $ 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  : Factor w/ 396 levels "-0.016850","-0.021024",..: NA NA NA NA NA NA NA N
##  $ kurtosis_picth_belt : Factor w/ 316 levels "-0.021887","-0.060755",..: NA NA NA NA NA NA NA N
##  $ kurtosis_yaw_belt   : Factor w/ 1 level "#DIV/0!": NA NA NA NA NA NA NA NA NA NA ...
##  $ skewness_roll_belt  : Factor w/ 394 levels "-0.003095","-0.010002",..: NA NA NA NA NA NA NA N
##  $ skewness_roll_belt.1: Factor w/ 337 levels "-0.005928","-0.005960",..: NA NA NA NA NA NA NA N
##  $ skewness_yaw_belt   : Factor w/ 1 level "#DIV/0!": NA NA NA NA NA NA NA NA NA NA ...
##  $ max_roll_belt       : num  NA NA NA NA NA NA NA NA NA NA ...
##  $ max_picth_belt      : int  NA NA NA NA NA NA NA NA NA NA ...
```

```
##  $ max_yaw_belt         : Factor w/ 67 levels "-0.1","-0.2",..: 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         : Factor w/ 67 levels "-0.1","-0.2",..: 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   : Factor w/ 3 levels "#DIV/0!","0.00",..: 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    : Factor w/ 329 levels "-0.02438","-0.04190",..: NA NA NA NA NA NA NA NA N
##  $ kurtosis_picth_arm   : Factor w/ 327 levels "-0.00484","-0.01311",..: NA NA NA NA NA NA NA NA N
##  $ kurtosis_yaw_arm     : Factor w/ 394 levels "-0.01548","-0.01749",..: NA NA NA NA NA NA NA NA N
##  $ skewness_roll_arm    : Factor w/ 330 levels "-0.00051","-0.00696",..: NA NA NA NA NA NA NA NA N
##  $ skewness_pitch_arm   : Factor w/ 327 levels "-0.00184","-0.01185",..: NA NA NA NA NA NA NA NA N
```

```
##  $ skewness_yaw_arm       : Factor w/ 394 levels "-0.00311","-0.00562",..: NA NA NA NA NA NA NA NA
##  $ max_roll_arm           : num  NA NA NA NA NA NA NA NA NA NA ...
##  $ max_picth_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 : Factor w/ 397 levels "-0.0035","-0.0073",..: NA NA NA NA NA NA NA NA NA
##  $ kurtosis_picth_dumbbell : Factor w/ 400 levels "-0.0163","-0.0233",..: NA NA NA NA NA NA NA NA NA
##  $ kurtosis_yaw_dumbbell  : Factor w/ 1 level "#DIV/0!": NA NA NA NA NA NA NA NA NA NA ...
##  $ skewness_roll_dumbbell : Factor w/ 400 levels "-0.0082","-0.0096",..: NA NA NA NA NA NA NA NA NA
##  $ skewness_pitch_dumbbell : Factor w/ 401 levels "-0.0053","-0.0084",..: NA NA NA NA NA NA NA NA NA
##  $ skewness_yaw_dumbbell  : Factor w/ 1 level "#DIV/0!": NA NA NA NA NA NA NA NA NA NA ...
##  $ max_roll_dumbbell      : num  NA NA NA NA NA NA NA NA NA NA ...
##  $ max_picth_dumbbell     : num  NA NA NA NA NA NA NA NA NA NA ...
##  $ max_yaw_dumbbell       : Factor w/ 72 levels "-0.1","-0.2",..: 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 ...
##  $ min_yaw_dumbbell       : Factor w/ 72 levels "-0.1","-0.2",..: NA NA NA NA NA NA NA NA NA NA ...
##  $ amplitude_roll_dumbbell : num  NA NA NA NA NA NA NA NA NA NA ...
##   [list output truncated]
```

```r
str(train$classe)
```

```
##  Factor w/ 5 levels "A","B","C","D",..: 1 1 1 1 1 1 1 1 1 1 ...
```

```r
dim(train)
```

```
## [1] 19622    160
```

```r
plot(train$classe, main = "histogram classe")
```

## histogram classe



As we see there are lots of variabrles with NA's lets remove those columns with more than 10% NA's. Also remove columns related to index, username or timestamp

```
train <- train[,7:160]
test <- test[,7:160]
NoNa <- apply(!is.na(train),2,sum)>(19622-0.1*19622)
train <- train[,NoNa]
test <- test[,NoNa]
dim(train)
```

```
## [1] 19622    54
```

As I was not able to run a mode wit 19622 observation, so for speed I needed to reduced the training set to only include 45%

```
InTrain <- createDataPartition(y=train$classe, p=0.45,list=FALSE)
trainset <- train[InTrain,]
testset <- train[-InTrain,]
```

For the training we use a 5-fold cross validation. As we already need to reduce the training set, we don't make k greater than 5, in order to avoid the sample set become to small.

```
control <- trainControl(method = "cv", number = 5)
```

As boosting and random forest are the most accurate models, we are using both and check accuracy.

```
#randomforest
set.seed(32333)
fitrf <- train(classe ~  ., data = trainset, method = "rf", proxy=TRUE, trControl=control)
```

```
## Loading required package: randomForest
## randomForest 4.6-10
## Type rfNews() to see new features/changes/bug fixes.
```

```
#gbm
set.seed(32333)
fitb <- train(classe ~  ., data = trainset, method = "gbm",  trControl=control,verbose=FALSE)
```

```
## Loading required package: gbm
## Loading required package: survival
## Loading required package: splines
##
## Attaching package: 'survival'
##
## The following object is masked from 'package:caret':
##
##     cluster
##
## Loading required package: parallel
## Loaded gbm 2.1
## Loading required package: plyr
```

```
fitb$finalModel
```

```
## A gradient boosted model with multinomial loss function.
## 150 iterations were performed.
## There were 53 predictors of which 44 had non-zero influence.
```

Let's compare accuracy of both models:

```
results <- data.frame(method = c(fitrf$method,fitb$method),
                     accuracy=c( fitrf$results[fitrf$results$mtry %in% fitrf$bestTune,2],
                                 tail(fitb$results[order(fitb$results$Accuracy),],1)$Accuracy)
                     )
results
```

```
##    method  accuracy
## 1      rf 0.9933193
## 2     gbm 0.9840356
```

Accuracy of rf model is better so lets use and view the Random forest model.

```
fitrf$finalModel
```

```
##
## Call:
```

```
##  randomForest(x = x, y = y, mtry = param$mtry, proxy = TRUE)
##                 Type of random forest: classification
##                       Number of trees: 500
## No. of variables tried at each split: 27
##
##          OOB estimate of  error rate: 0.51%
## Confusion matrix:
##       A    B    C    D    E  class.error
## A 2510    0    0    0    1 0.0003982477
## B    5 1699    4    1    0 0.0058513751
## C    0    9 1531    0    0 0.0058441558
## D    0    1   13 1433    1 0.0103591160
## E    0    1    0    9 1614 0.0061576355
```

```
predictionsrf <- predict(fitrf, trainset)
confusionMatrix(predictionsrf ,trainset$classe)
```

```
## Confusion Matrix and Statistics
##
##           Reference
## Prediction    A    B    C    D    E
##          A 2511    0    0    0    0
##          B    0 1709    0    0    0
##          C    0    0 1540    0    0
##          D    0    0    0 1448    0
##          E    0    0    0    0 1624
##
## Overall Statistics
##
##                Accuracy : 1
##                  95% CI : (0.9996, 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
```

As we only used part of the training data set for prediction we can use the other part for testing and estimating
the out of sample error.

```
predictionstest <- predict(fitrf, testset)
confusionMatrix(predictionstest ,testset$classe)
```

```
## Confusion Matrix and Statistics
##
##           Reference
## Prediction    A    B    C    D    E
##          A 3068   11    0    0    0
##          B    1 2073   12    0    1
##          C    0    3 1870   17    0
##          D    0    1    0 1751    7
##          E    0    0    0    0 1975
##
## Overall Statistics
##
##                Accuracy : 0.9951
##                  95% CI : (0.9936, 0.9963)
##     No Information Rate : 0.2844
##     P-Value [Acc > NIR] : < 2.2e-16
##
##                   Kappa : 0.9938
##  Mcnemar's Test P-Value : NA
##
## Statistics by Class:
##
##                      Class: A Class: B Class: C Class: D Class: E
## Sensitivity            0.9997   0.9928   0.9936   0.9904   0.9960
## Specificity            0.9986   0.9984   0.9978   0.9991   1.0000
## Pos Pred Value         0.9964   0.9933   0.9894   0.9955   1.0000
## Neg Pred Value         0.9999   0.9983   0.9987   0.9981   0.9991
## Prevalence             0.2844   0.1935   0.1744   0.1639   0.1838
## Detection Rate         0.2843   0.1921   0.1733   0.1623   0.1830
## Detection Prevalence   0.2854   0.1934   0.1752   0.1630   0.1830
## Balanced Accuracy      0.9991   0.9956   0.9957   0.9947   0.9980
```

```
outOfSampleError <- 1-sum(predictionstest == testset$classe)/length(predictionstest)
print(outOfSampleError)
```

```
## [1] 0.004911956
```