# Practical Machine Learning Project

*nazila*

*January 23, 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 dumbell 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

## Reading and Cleaning Data

Start with reading both training and testing instances.

```
library(data.table)
library(caret)
```

```
## Warning: package 'caret' was built under R version 3.2.3
```

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

```
library(randomForest)
```

```
## randomForest 4.6-12
## Type rfNews() to see new features/changes/bug fixes.
```

```
library(foreach)
library(rpart)
library(rpart.plot)
library(corrplot)
```

```
#read training data
training_data <- read.csv("pml-training.csv", na.strings=c("#DIV/0!"," ", "", "NA", "NAs", "NULL"))
#read test data
testing_data <- read.csv("pml-testing.csv", na.strings=c("#DIV/0!"," ", "", "NA", "NAs", "NULL"))
```

Columns with NAs will be removed. Besides that highly correlated variables and variables with 0 (or approx to 0) variance will be removed.

## Cleaning data

```
clean_training <- training_data[, -which(names(training_data) %in% c("X", "user_name", "raw_timestamp_pa

#remove columns with NAs
clean_training = clean_training[, colSums(is.na(clean_training)) == 0]

#remove variables with 0 or near to 0 variance
zero_variance =nearZeroVar(clean_training[sapply(clean_training, is.numeric)], saveMetrics=TRUE)
clean_training = clean_training[, zero_variance[, 'nzv'] == 0]

correlation_matrix <- cor(na.omit(clean_training[sapply(clean_training, is.numeric)]))

dim(correlation_matrix)
```

```
## [1] 52 52
```

```
correlationmatrixdegreesoffreedom <- expand.grid(row = 1:52, col = 1:52)
 #this returns the correlation matrix in matrix format
correlationmatrixdegreesoffreedom$correlation <- as.vector(correlation_matrix)
removehighcorrelation <- findCorrelation(correlation_matrix, cutoff = .7, verbose = TRUE)
```

```
## Compare row 10  and column  1 with corr  0.992
##   Means:  0.27 vs 0.168 so flagging column 10
## Compare row 1  and column  9 with corr  0.925
##   Means:  0.25 vs 0.164 so flagging column 1
## Compare row 9  and column  22 with corr  0.722
##   Means:  0.233 vs 0.161 so flagging column 9
## Compare row 22  and column  4 with corr  0.759
##   Means:  0.224 vs 0.158 so flagging column 22
## Compare row 4  and column  3 with corr  0.762
##   Means:  0.2 vs 0.155 so flagging column 4
## Compare row 3  and column  8 with corr  0.708
##   Means:  0.2 vs 0.153 so flagging column 3
## Compare row 36  and column  29 with corr  0.849
##   Means:  0.257 vs 0.151 so flagging column 36
## Compare row 8  and column  2 with corr  0.966
##   Means:  0.229 vs 0.146 so flagging column 8
## Compare row 2  and column  11 with corr  0.884
##   Means:  0.212 vs 0.143 so flagging column 2
## Compare row 37  and column  38 with corr  0.769
```

```
##   Means:  0.198 vs 0.139 so flagging column 37
## Compare row 35  and column  30 with corr  0.773
##   Means:  0.195 vs 0.137 so flagging column 35
## Compare row 38  and column  5 with corr  0.781
##   Means:  0.177 vs 0.134 so flagging column 38
## Compare row 21  and column  24 with corr  0.814
##   Means:  0.176 vs 0.133 so flagging column 21
## Compare row 34  and column  28 with corr  0.808
##   Means:  0.176 vs 0.13 so flagging column 34
## Compare row 23  and column  26 with corr  0.779
##   Means:  0.137 vs 0.129 so flagging column 23
## Compare row 25  and column  24 with corr  0.792
##   Means:  0.145 vs 0.128 so flagging column 25
## Compare row 12  and column  13 with corr  0.779
##   Means:  0.122 vs 0.127 so flagging column 13
## Compare row 48  and column  51 with corr  0.772
##   Means:  0.145 vs 0.127 so flagging column 48
## Compare row 19  and column  18 with corr  0.918
##   Means:  0.095 vs 0.127 so flagging column 18
## Compare row 46  and column  45 with corr  0.846
##   Means:  0.131 vs 0.129 so flagging column 46
## Compare row 45  and column  31 with corr  0.71
##   Means:  0.098 vs 0.129 so flagging column 31
## Compare row 45  and column  33 with corr  0.716
##   Means:  0.078 vs 0.132 so flagging column 33
## All correlations <= 0.7
```

```r
#this removes highly correlated variables (in psychometric theory .7+ correlation is a high correlation)
clean_training <- clean_training[, -removehighcorrelation]

for(i in c(8:ncol(clean_training)-1)) {clean_training[,i] = as.numeric(as.character(clean_training[,i]))

for(i in c(8:ncol(testing_data)-1)) {testing_data[,i] = as.numeric(as.character(testing_data[,i]))}
```

The cleaned dataset will only consist of complete columns. For a lighter dataset, user name, timestamps and windows will be removed.

```r
#drop blank column
featureset <- colnames(clean_training[colSums(is.na(clean_training)) == 0])[-(1:7)]
modeldata <- clean_training[featureset]

#cleansed data to build model
featureset
```

```
##  [1] "yaw_arm"               "total_accel_arm"      "gyros_arm_y"
##  [4] "gyros_arm_z"           "magnet_arm_x"         "magnet_arm_z"
##  [7] "roll_dumbbell"         "pitch_dumbbell"       "yaw_dumbbell"
## [10] "total_accel_dumbbell"  "gyros_dumbbell_y"     "magnet_dumbbell_z"
## [13] "roll_forearm"          "pitch_forearm"        "yaw_forearm"
## [16] "total_accel_forearm"   "gyros_forearm_x"      "gyros_forearm_y"
## [19] "accel_forearm_x"       "accel_forearm_z"      "magnet_forearm_x"
## [22] "magnet_forearm_y"      "magnet_forearm_z"     "classe"
```

Cross-Validation

## Split the sample in two samples. This is to divide training and testing for cross-validation.

```
idx <- createDataPartition(modeldata$classe, p=0.6, list=FALSE )
training <- modeldata[idx,]
testing <- modeldata[-idx,]

control <- trainControl(method="cv", 5)
model <- train(classe ~ ., data=training, method="rf", trControl=control, ntree=250)
model
```

```
## Random Forest
##
## 11776 samples
##    23 predictor
##     5 classes: 'A', 'B', 'C', 'D', 'E'
##
## No pre-processing
## Resampling: Cross-Validated (5 fold)
## Summary of sample sizes: 9420, 9420, 9423, 9421, 9420
## Resampling results across tuning parameters:
##
##   mtry  Accuracy   Kappa      Accuracy SD  Kappa SD
##    2    0.9691746  0.9609868  0.003456090  0.004382939
##   12    0.9651833  0.9559421  0.002458732  0.003116054
##   23    0.9537183  0.9414267  0.007106288  0.009010892
##
## Accuracy was used to select the optimal model using  the largest value.
## The final value used for the model was mtry = 2.
```

```
predict <- predict(model, testing)
confusionMatrix(testing$classe, predict)
```

```
## Confusion Matrix and Statistics
##
##           Reference
## Prediction    A    B    C    D    E
##          A 2225    2    1    3    1
##          B   37 1461   11    4    5
##          C    3   27 1330    7    1
##          D    2    0   46 1235    3
##          E    4    6    1   16 1415
##
## Overall Statistics
##
##                Accuracy : 0.9771
##                  95% CI : (0.9735, 0.9803)
##     No Information Rate : 0.2894
##     P-Value [Acc > NIR] : < 2.2e-16
```

```
## 
##                 Kappa : 0.971
##   Mcnemar's Test P-Value : 1.396e-13
## 
## Statistics by Class:
## 
##                     Class: A Class: B Class: C Class: D Class: E
## Sensitivity           0.9797   0.9766   0.9575   0.9763   0.9930
## Specificity           0.9987   0.9910   0.9941   0.9923   0.9958
## Pos Pred Value        0.9969   0.9625   0.9722   0.9603   0.9813
## Neg Pred Value        0.9918   0.9945   0.9909   0.9954   0.9984
## Prevalence            0.2894   0.1907   0.1770   0.1612   0.1816
## Detection Rate        0.2836   0.1862   0.1695   0.1574   0.1803
## Detection Prevalence  0.2845   0.1935   0.1744   0.1639   0.1838
## Balanced Accuracy     0.9892   0.9838   0.9758   0.9843   0.9944
```

```r
accuracy <- postResample(predict, testing$classe)
accuracy
```

```
##  Accuracy     Kappa
## 0.9770584 0.9709635
```

```r
result <- predict(model, training[, -length(names(training))])
result
```

```
##    [1] A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A
##   [35] A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A
##   [69] A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A
##  [103] A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A
##  [137] A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A
##  [171] A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A
##  [205] A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A
##  [239] A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A
##  [273] A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A
##  [307] A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A
##  [341] A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A
##  [375] A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A
##  [409] A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A
##  [443] A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A
##  [477] A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A
##  [511] A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A
##  [545] A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A
##  [579] A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A
##  [613] A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A
##  [647] A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A
##  [681] A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A
##  [715] A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A
##  [749] A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A
##  [783] A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A
##  [817] A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A
##  [851] A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A
##  [885] A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A
##  [919] A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A
```

```
##  [953] A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A
##  [987] A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A
## [1021] A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A
## [1055] A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A
## [1089] A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A
## [1123] A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A
## [1157] A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A
## [1191] A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A
## [1225] A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A
## [1259] A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A
## [1293] A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A
## [1327] A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A
## [1361] A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A
## [1395] A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A
## [1429] A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A
## [1463] A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A
## [1497] A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A
## [1531] A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A
## [1565] A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A
## [1599] A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A
## [1633] A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A
## [1667] A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A
## [1701] A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A
## [1735] A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A
## [1769] A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A
## [1803] A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A
## [1837] A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A
## [1871] A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A
## [1905] A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A
## [1939] A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A
## [1973] A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A
## [2007] A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A
## [2041] A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A
## [2075] A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A
## [2109] A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A
## [2143] A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A
## [2177] A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A
## [2211] A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A
## [2245] A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A
## [2279] A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A
## [2313] A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A
## [2347] A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A
## [2381] A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A
## [2415] A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A
## [2449] A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A
## [2483] A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A
## [2517] A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A
## [2551] A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A
## [2585] A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A
## [2619] A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A
## [2653] A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A
## [2687] A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A
## [2721] A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A
## [2755] A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A
```

```
## [2789] A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A
## [2823] A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A
## [2857] A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A
## [2891] A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A
## [2925] A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A
## [2959] A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A
## [2993] A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A
## [3027] A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A
## [3061] A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A
## [3095] A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A
## [3129] A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A
## [3163] A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A
## [3197] A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A
## [3231] A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A
## [3265] A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A
## [3299] A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A A
## [3333] A A A A A A A A A A A A A A A B B B B B B B B B B B B B B B B B B B
## [3367] B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B
## [3401] B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B
## [3435] B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B
## [3469] B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B
## [3503] B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B
## [3537] B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B
## [3571] B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B
## [3605] B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B
## [3639] B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B
## [3673] B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B
## [3707] B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B
## [3741] B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B
## [3775] B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B
## [3809] B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B
## [3843] B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B
## [3877] B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B
## [3911] B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B
## [3945] B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B
## [3979] B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B
## [4013] B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B
## [4047] B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B
## [4081] B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B
## [4115] B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B
## [4149] B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B
## [4183] B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B
## [4217] B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B
## [4251] B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B
## [4285] B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B
## [4319] B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B
## [4353] B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B
## [4387] B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B
## [4421] B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B
## [4455] B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B
## [4489] B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B
## [4523] B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B
## [4557] B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B
## [4591] B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B
```

```
##  [4625] B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B
##  [4659] B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B
##  [4693] B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B
##  [4727] B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B
##  [4761] B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B
##  [4795] B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B
##  [4829] B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B
##  [4863] B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B
##  [4897] B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B
##  [4931] B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B
##  [4965] B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B
##  [4999] B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B
##  [5033] B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B
##  [5067] B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B
##  [5101] B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B
##  [5135] B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B
##  [5169] B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B
##  [5203] B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B
##  [5237] B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B
##  [5271] B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B
##  [5305] B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B
##  [5339] B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B
##  [5373] B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B
##  [5407] B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B
##  [5441] B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B
##  [5475] B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B
##  [5509] B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B
##  [5543] B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B
##  [5577] B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B B
##  [5611] B B B B B B B B B B B B B B B B B C C C C C C C C C C C C C C C C C
##  [5645] C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C
##  [5679] C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C
##  [5713] C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C
##  [5747] C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C
##  [5781] C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C
##  [5815] C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C
##  [5849] C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C
##  [5883] C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C
##  [5917] C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C
##  [5951] C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C
##  [5985] C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C
##  [6019] C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C
##  [6053] C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C
##  [6087] C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C
##  [6121] C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C
##  [6155] C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C
##  [6189] C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C
##  [6223] C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C
##  [6257] C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C
##  [6291] C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C
##  [6325] C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C
##  [6359] C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C
##  [6393] C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C
##  [6427] C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C
```

```
## [6461] C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C
## [6495] C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C
## [6529] C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C
## [6563] C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C
## [6597] C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C
## [6631] C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C
## [6665] C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C
## [6699] C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C
## [6733] C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C
## [6767] C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C
## [6801] C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C
## [6835] C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C
## [6869] C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C
## [6903] C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C
## [6937] C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C
## [6971] C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C
## [7005] C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C
## [7039] C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C
## [7073] C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C
## [7107] C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C
## [7141] C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C
## [7175] C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C
## [7209] C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C
## [7243] C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C
## [7277] C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C
## [7311] C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C
## [7345] C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C
## [7379] C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C
## [7413] C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C
## [7447] C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C
## [7481] C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C
## [7515] C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C
## [7549] C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C
## [7583] C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C
## [7617] C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C
## [7651] C C C C C C C C C C C C C C C C C C C C C C C C C C C C C C D D D D
## [7685] D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D
## [7719] D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D
## [7753] D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D
## [7787] D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D
## [7821] D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D
## [7855] D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D
## [7889] D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D
## [7923] D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D
## [7957] D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D
## [7991] D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D
## [8025] D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D
## [8059] D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D
## [8093] D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D
## [8127] D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D
## [8161] D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D
## [8195] D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D
## [8229] D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D
## [8263] D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D
```

```
##  [8297] D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D
##  [8331] D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D
##  [8365] D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D
##  [8399] D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D
##  [8433] D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D
##  [8467] D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D
##  [8501] D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D
##  [8535] D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D
##  [8569] D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D
##  [8603] D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D
##  [8637] D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D
##  [8671] D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D
##  [8705] D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D
##  [8739] D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D
##  [8773] D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D
##  [8807] D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D
##  [8841] D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D
##  [8875] D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D
##  [8909] D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D
##  [8943] D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D
##  [8977] D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D
##  [9011] D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D
##  [9045] D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D
##  [9079] D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D
##  [9113] D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D
##  [9147] D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D
##  [9181] D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D
##  [9215] D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D
##  [9249] D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D
##  [9283] D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D
##  [9317] D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D
##  [9351] D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D
##  [9385] D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D
##  [9419] D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D
##  [9453] D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D
##  [9487] D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D
##  [9521] D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D
##  [9555] D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D D
##  [9589] D D D D D D D D D D D D D D D D D D D D D D D D E E E E E E E E E E
##  [9623] E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E
##  [9657] E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E
##  [9691] E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E
##  [9725] E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E
##  [9759] E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E
##  [9793] E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E
##  [9827] E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E
##  [9861] E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E
##  [9895] E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E
##  [9929] E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E
##  [9963] E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E
##  [9997] E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E
## [10031] E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E
## [10065] E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E
## [10099] E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E
```

```
## [10133] E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E
## [10167] E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E
## [10201] E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E
## [10235] E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E
## [10269] E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E
## [10303] E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E
## [10337] E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E
## [10371] E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E
## [10405] E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E
## [10439] E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E
## [10473] E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E
## [10507] E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E
## [10541] E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E
## [10575] E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E
## [10609] E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E
## [10643] E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E
## [10677] E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E
## [10711] E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E
## [10745] E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E
## [10779] E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E
## [10813] E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E
## [10847] E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E
## [10881] E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E
## [10915] E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E
## [10949] E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E
## [10983] E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E
## [11017] E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E
## [11051] E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E
## [11085] E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E
## [11119] E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E
## [11153] E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E
## [11187] E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E
## [11221] E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E
## [11255] E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E
## [11289] E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E
## [11323] E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E
## [11357] E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E
## [11391] E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E
## [11425] E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E
## [11459] E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E
## [11493] E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E
## [11527] E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E
## [11561] E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E
## [11595] E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E
## [11629] E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E
## [11663] E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E
## [11697] E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E
## [11731] E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E E
## [11765] E E E E E E E E E E E E
## Levels: A B C D E
```

```r
treeModel <- rpart(classe ~ ., data=clean_training, method="class")
prp(treeModel)
```

yes  **pitch_forearm < –34**  no

**magnet_belt_y >= 556**

(A)

**accel_forearm_z >= 120**

**roll_dumbbell < 64**

**total_accel_dum >= 6.5**  (E)

(D)

**roll_forearm < 122**

**magnet_dumbbell >= 284**

(D)

**roll_dumbbell >= 79**

**magnet_dumbbell < –28**

**accel_forearm_x >= –108**

**magnet_belt_y >= 592**

**roll_dumbbell >= –88**  (A)

**total_accel_dum < 4.5**

**gyros_belt_z >= –0.26**

(A)  (C)

(B)  **magnet_arm_x < 126**  (E)

(E)

**magnet_dumbbell < 50**  (A)

(C)  (B)

**magnet_dumbbell < 22**

**yaw_dumbbell < –93**  (C)  (D)  (A)  (D)

(A)  (B)

(C)

## Conclusions and Test Data Submit

As can be seen from the confusion matrix the proposed model is very accurate.

Prepare the submission. (using COURSERA provided code)

```
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)
  }
}


x <- evaluation_data
x <- x[feature_set[feature_set!='classe']]
answers <- predict(rf, newdata=x)

answers

pml_write_files(answers)
```