

BARBELL LIFTS MODEL

Synopsis

This study is conducted to determine how the barbell lifts are done, based on data got from accelerometers on the back, forearm, arm, and dumbbell.

This "classification model" was tested using decision trees and random forest (too much time processing data on a "MAC Big ram and core i7").

The training data was split in Training and Testing data. Training 70%, Testing 30%.

The cross validation was conducted using three random samples, and the out of sample error is estimated by average of three models.

The model reports 66.17% accuracy average on testing data (extracted from training dataset): Model 1: 66.15% Model 2: 66.17% Model 3: 66.19%.

Out of sample error expected: 32%.

```
## Loading data and Exploratory Analysis

library(rattle)

## warning: package 'rattle' was built under R version 3.1.1

## rattle: A free graphical interface for data mining with R.
## Version 3.1.0. Copyright (c) 2000-2004 Thomas H. Dunning.
## escriba "rattle()" para ajudar, sacudir y rotar sus datos.

library(appliedPredictiveModeling)
library(caret)

## Loading required package: lattice
## Loading required package: applied

library(mlisc)

## Loading required package: grid
## Loading required package: plotrix
## Loading required package: upRmius
## Attaching package: 'survival'
## The following object is masked from 'package:caret':
##   cluster
## Loading required package: formula

## warning: package 'formula' was built under R version 3.1.1

## Attaching package: 'mlisc'
## The following objects are masked from 'package:base':
##   format.pval, round.POSIXt, trunc.POSIXt, units

library(s0071)

## Attaching package: 's0071'
## The following object is masked from 'package:mlisc':
##   lqpute

library(psych)

## Attaching package: 'psych'
## The following object is masked from 'package:mlisc':
##   describe
## The following object is masked from 'package:ggplot2':
##   xax

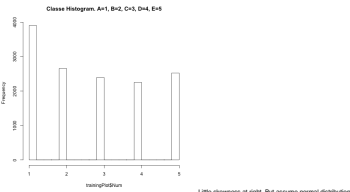
setwd("~/Coursera/Data Science/Practical Machine Learning/projects")
trainingset00 <- read.csv("pml-training.csv")
testingset00 <- read.csv("pml-testing.csv")

trainIndex <-
trainingset <- c(12,13,14,15,16,17,18,19,20,21,22,23,24,25,26,27,28,29,30,31,32,33,34,35,36,50,51,52,53,54,55,56,57,58,59,69,70,71,72,73,74,75,76,77,78,79,80,81,82,83,87,88,89,90,91,92,93,94,95,96,97,98,99,100,101,109,104,105,106,107,108,109,110,111,112,125,126,127,128,129,130,131,132,133,134,135,136,137,138,139,141,142,143,144,145,146,147,1
testingset <- c(12,13,14,15,16,17,18,19,20,21,22,23,24,25,26,27,28,29,30,31,32,33,34,35,36,50,51,52,53,54,55,56,57,58,59,69,70,71,72,73,74,75,76,77,78,79,80,81,82,83,87,88,89,90,91,92,93,94,95,96,97,98,99,100,101,109,104,105,106,107,108,109,110,111,112,125,126,127,128,129,130,131,132,133,134,135,136,137,138,139,141,142,143,144,145,146,147,14

#str(trainingset)
#str(testingset)
trainIndex <- createDataPartition(trainingset$class, p = 0.70, times=ns) #70% for training
training = trainingset[trainIndex,]
testing = trainingset[-trainIndex,]

#R<-train(class=..., data=training, method="rf", growTree)

# Assign values to Class ... for plotting purposes
# A=1, B=2, C=3, D=4, E=5
trainingset <- training
trainingset$train <- 0
columns1
trainingset[trainingset$class=="A",column] <- 1
trainingset[trainingset$class=="B",column] <- 2
trainingset[trainingset$class=="C",column] <- 3
trainingset[trainingset$class=="D",column] <- 4
trainingset[trainingset$class=="E",column] <- 5
hist(trainingset$train, main="Class histogram, A=1, B=2, C=3, D=4, E=5")
```



qplot(trainingset\$train, trainingset\$train, colour=trainingset\$user_name, data=trainingset, main="Class in Numbers: A=1, B=2, C=3, D=4, E=5")



Modelling: 3 random models for cross validation

```
# Model 1
modf1 <- train(class=..., method="rpart", data=training) #Commented for publishing results

## Loading required package: rpart

confusionMatrix(testing$class, predict(modf1, newdata=testing))
```

```
## Confusion Matrix and Statistics
##
##      Reference
## Prediction  A      B      C      D      E
## A      1674    0      0      0      0
## B      0      1139    0      0      0
## C      0      0      0      1035    0
## D      0      0      0      364      0
## E      0      0      0      1082    0
##
## Overall Statistics
##
##      Accuracy : 0.662
##      95% CI   : (0.65, 0.674)
##      No Information Rate : 0.122
##      P-value [Acc > NA] : <2e-16
##      Kappa : 0.17
##      Mcnemar's Test P-value : NA
##
## Statistics by Class:
##
##      Class: A Class: B Class: C Class: D Class: E
## Sensitivity      1.000 1.000 NA      NA      0.332
## Specificity      1.000 1.000 0.826 0.836 1.000
## Pos Pred Value   1.000 1.000 NA      NA      1.000
## Neg Pred Value   1.000 1.000 NA      NA      0.186
## Prevalence       0.284 0.134 0.000 0.000 0.122
## Detection Rate   0.284 0.134 0.000 0.000 0.122
## Detection Prevalence 0.284 0.134 0.000 0.000 0.122
## Balanced Accuracy 1.000 1.000 NA      NA      0.676
```

"" Model 1 Accuracy: 66.17%. Out of sample Error = 34% ""

```
# Model 2
trainIndex2 <- createDataPartition(trainingset$class, p = 0.70, times=ns) #70% for training
training2 = trainingset[trainIndex2,]
testing2 = trainingset[-trainIndex2,]

modf2 <- train(class=..., method="rpart", data=training2) #Commented for publishing results
confusionMatrix(testing2$class, predict(modf2, newdata=testing2))
```

```
## Confusion Matrix and Statistics
##
##      Reference
## Prediction  A      B      C      D      E
## A      1674    0      0      0      0
## B      0      1139    0      0      0
## C      0      0      0      1035    0
## D      0      0      0      364      0
## E      0      0      0      1082    0
##
## Overall Statistics
##
##      Accuracy : 0.662
##      95% CI   : (0.65, 0.674)
##      No Information Rate : 0.122
##      P-value [Acc > NA] : <2e-16
##      Kappa : 0.17
##      Mcnemar's Test P-value : NA
##
## Statistics by Class:
##
##      Class: A Class: B Class: C Class: D Class: E
## Sensitivity      1.000 1.000 NA      NA      0.332
## Specificity      1.000 1.000 0.826 0.836 1.000
## Pos Pred Value   1.000 1.000 NA      NA      1.000
## Neg Pred Value   1.000 1.000 NA      NA      0.186
## Prevalence       0.284 0.134 0.000 0.000 0.122
## Detection Rate   0.284 0.134 0.000 0.000 0.122
## Detection Prevalence 0.284 0.134 0.000 0.000 0.122
## Balanced Accuracy 1.000 1.000 NA      NA      0.676
```

```
"Model 2 Accuracy: 66.19%. Out of sample Error = 34%"

# Model 3
trainIndex = createDataPartition(y=trainIndex$class, p = 0.70, list=FALSE) #70% for training
training = trainingSet[trainIndex,]
testing = trainingSet[-trainIndex,]

modFit1 <- train(class ~ ., method="part", data=training) #Commented for publishing results
confusionMatrix(testingSet[class, p=modFit1$class, newdata=testing])

## Confusion matrix and statistics
##
## Reference
##
## Prediction A B C D E
## A 1673 0 0 0
## B 0 1188 0 0
## C 0 0 0 1695
## D 0 0 0 1664
## E 0 0 0 1082
##
## Overall statistics
##
## Accuracy : 0.662
## Sens CI : (0.649, 0.674)
## No Information Rate : 0.332
## P-value [Acc > NRI] : <2e-16
## NRI : 0.369
## Mcnemar's test P-value : NA
##
## Statistics by Class:
##
## Class: A Class: B Class: C Class: D Class: E
## Sensitivity 1.000 1.000 0.828 0.838 1.000
## Pos Pred Value 0.999 0.999 NA NA 1.000
## Neg Pred Value 1.000 1.000 NA NA 0.999
## Prevalence 0.384 0.194 0.000 0.000 0.332
## Detection Rate 0.384 0.194 0.000 0.000 0.332
## Detection Prevalence 0.384 0.194 0.174 0.164 0.332
## Balanced Accuracy 1.000 0.999 NA NA 0.978

Model 3 Accuracy: 66.19%. Out of sample Error = 34%

Applying model to test dataset

results <- predict(modFit1, newdata=testingSet)
class(results)

## [1] "factor"

answers <- as.vector(results, mode = "character")
is.vector(answers)

## [1] TRUE

class(answers)

## [1] "character"

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

plot_write_files(answers)
```