

Practical Machine Learning Wk 4

ekonomix

25 January 2017

Executive Summary

We have data from 4 sensors placed on participants bodies and object (belt, forearm, arm and the dumbbell); they measure how the different body parts and the dumbbell itself are moving as the participant is attempting to lift it.

Participants were asked to lift the dumbbell in 5 different ways, 1 correct way and 4 'wrong' ways. Our aim is to distinguish "how well" the exercise is taking place, hence using the sensor data distinguish between these different types of lift.

load some libraries we are likely to be using

```
library(caret)

## Warning: package 'caret' was built under R version 3.3.2

## Loading required package: lattice

## Loading required package: ggplot2

## Warning: package 'ggplot2' was built under R version 3.3.2

library(ggplot2)
```

Get data - download as csv and load

```
Urla <- "https://d396qusza40orc.cloudfront.net/predmachlearn/pml-
training.csv"
Urlb <- "https://d396qusza40orc.cloudfront.net/predmachlearn/pml-testing.csv"
training <- read.csv(url(Urla), na.strings=c("NA", "#DIV/0!", ""))
testing <- read.csv(url(Urlb), na.strings=c("NA", "#DIV/0!", ""))
```

take a look at the data

```
dim(training)

## [1] 19622 160

#str(training) # not shown in output
dim(testing)

## [1] 20 160

#str(testing) # not shown in output
```

it looks like the first 7 variables have no predictive value for this exercise

get rid of variables with many NAs and variables expected to have no predictive value in this case

```
NA_Count = sapply(1:dim(training)[2],function(x)sum(is.na(training[,x])))
NA_Count

## [1] 0 0 0 0 0 0 0 0 0 0 0 0
## [12] 19226 19248 19622 19225 19248 19622 19216 19216 19226 19216 19216
## [23] 19226 19216 19216 19226 19216 19216 19216 19216 19216 19216 19216
## [34] 19216 19216 19216 0 0 0 0 0 0 0 0
## [45] 0 0 0 0 0 19216 19216 19216 19216 19216 19216
## [56] 19216 19216 19216 19216 0 0 0 0 0 0 0
## [67] 0 0 19294 19296 19227 19293 19296 19227 19216 19216 19216
## [78] 19216 19216 19216 19216 19216 19216 0 0 0 19221 19218
## [89] 19622 19220 19217 19622 19216 19216 19221 19216 19216 19221 19216
## [100] 19216 19221 0 19216 19216 19216 19216 19216 19216 19216 19216
## [111] 19216 19216 0 0 0 0 0 0 0 0 0
## [122] 0 0 0 19300 19301 19622 19299 19301 19622 19216 19216
## [133] 19300 19216 19216 19300 19216 19216 19300 0 19216 19216 19216
## [144] 19216 19216 19216 19216 19216 19216 19216 0 0 0 0
## [155] 0 0 0 0 0 0 0

NA_list = which(NA_Count>0)
```

modify the training and test data sets to remove unnecessary columns and transforming the class into a factor

```
training_cleaning <- training[,-NA_list]
training_cleaning <- training_cleaning[,-c(1:7)]
training_cleaning$classe = factor(training_cleaning$classe)

inTrain <- createDataPartition(training_cleaning$classe, p=0.60, list=FALSE)
training_clean = training_cleaning[inTrain,]
validation_clean = training_cleaning[-inTrain,]

testing_clean <- testing[,-NA_list]
testing_clean <- testing_clean[,-c(1:7)]

# head(testing_clean) # not shown in output
```

build models and deciding which works best

this is a classification problem, and we'll try random forest and classification tree

These are methods used for supervised learning which means the class we are trying to predict is known.

A decision tree (both of these methods are some form of decision trees) basically takes the whole dataset, then runs through the variables and picks the one that causes the best split, which means the 2 groups are the most distinct. How this is measured can be specified in the algorithm.

The process of splitting each subset is repeated until the tree has reached a maximum depth, or the benefit of splitting in terms of increased distinctiveness cannot be reached.

The main difference between random forest and a std classification tree is that random forest, as the name suggests, build many trees and combines them, usually by voting. The advantage of random forest is that it is less likely to be influenced by quirks in the data (overfitting issue), however, on large datasets, it can be computationally expensive, and it is harder to explain.

In this case I will crossvalidate the random forest 3 times. Since the dataset is small and here I don't care about explainability, I expect random forest to perform better.

```
set.seed(2501)
```

Random Forest

```
rfFit <- train(classe ~ ., method = "rf", data = training_clean, importance =
T, trControl = trainControl(method = "cv", number = 3))

## Loading required package: randomForest
## Warning: package 'randomForest' was built under R version 3.3.2
## 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
#validation performance
validation_rfpred <- predict(rfFit, newdata=validation_clean)
rf_confusion <- confusionMatrix(validation_rfpred, validation_clean$classe)
rf_confusion
```

```
## Confusion Matrix and Statistics
##
##           Reference
## Prediction    A    B    C    D    E
##           A 2219    6    0    0    0
##           B    5 1508    6    1    0
##           C    7    4 1354   24    3
##           D    1    0    8 1260    7
##           E    0    0    0    1 1432
##
## Overall Statistics
##
##           Accuracy : 0.9907
##           95% CI : (0.9883, 0.9927)
##           No Information Rate : 0.2845
##           P-Value [Acc > NIR] : < 2.2e-16
##
##           Kappa : 0.9882
##           McNemar's Test P-Value : NA
##
## Statistics by Class:
##
##           Class: A Class: B Class: C Class: D Class: E
## Sensitivity      0.9942  0.9934  0.9898  0.9798  0.9931
## Specificity      0.9989  0.9981  0.9941  0.9976  0.9998
## Pos Pred Value   0.9973  0.9921  0.9727  0.9875  0.9993
## Neg Pred Value   0.9977  0.9984  0.9978  0.9960  0.9984
## Prevalence       0.2845  0.1935  0.1744  0.1639  0.1838
## Detection Rate   0.2828  0.1922  0.1726  0.1606  0.1825
## Detection Prevalence 0.2836  0.1937  0.1774  0.1626  0.1826
## Balanced Accuracy 0.9966  0.9958  0.9920  0.9887  0.9965
##
#Looks good
```

Random Forest Results look good!

Classification Tree

```
rpartFit <- train(classe ~ ., method = "rpart", data = training_clean)
## Loading required package: rpart
#training performance
validation_rpartpred <- predict(rpartFit, newdata=validation_clean)
confusionMatrix(validation_rpartpred, validation_clean$classe)
## Confusion Matrix and Statistics
##
##           Reference
## Prediction    A    B    C    D    E
##           A 2018  656  637  597  211
```

```
##           B    30   369    15   193    72
##           C   176   493   716   496   475
##           D     0     0     0     0     0
##           E     8     0     0     0   684
##
## Overall Statistics
##
##               Accuracy : 0.4827
##               95% CI : (0.4716, 0.4938)
##       No Information Rate : 0.2845
##       P-Value [Acc > NIR] : < 2.2e-16
##
##               Kappa : 0.3241
##  Mcnemar's Test P-Value : NA
##
## Statistics by Class:
##
##               Class: A Class: B Class: C Class: D Class: E
## Sensitivity          0.9041  0.24308  0.52339  0.0000  0.47434
## Specificity          0.6258  0.95101  0.74684  1.0000  0.99875
## Pos Pred Value       0.4899  0.54345  0.30390      NaN  0.98844
## Neg Pred Value       0.9426  0.83968  0.88124  0.8361  0.89405
## Prevalence           0.2845  0.19347  0.17436  0.1639  0.18379
## Detection Rate       0.2572  0.04703  0.09126  0.0000  0.08718
## Detection Prevalence 0.5250  0.08654  0.30028  0.0000  0.08820
## Balanced Accuracy    0.7649  0.59705  0.63511  0.5000  0.73655

#not as good
```

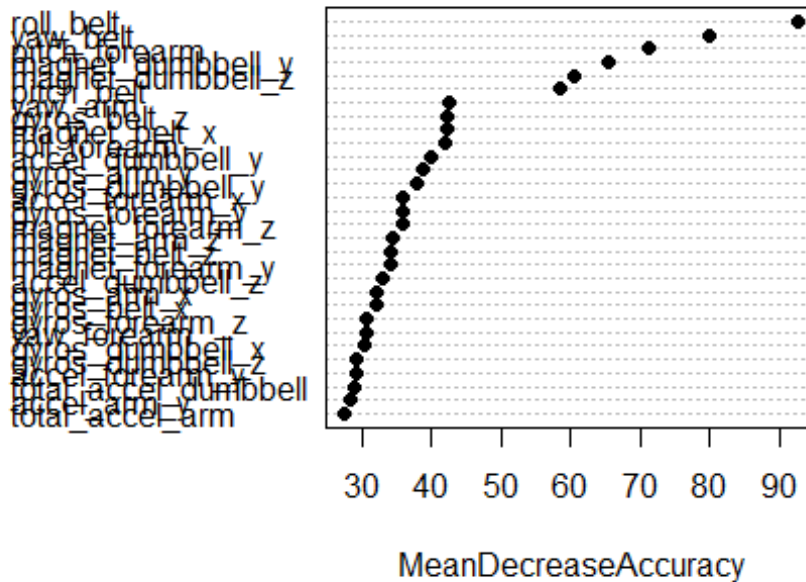
Regressions trees are not as good as random forest here.

Hence Random Forest is Chosen!

important variables, expected error (1-accuracy) and predictions for test data:

```
#Important Variables
imp_rf <- varImp(rfFit)$importance
varImpPlot(rfFit$finalModel, sort = TRUE, type = 1, pch = 19, col = 1, cex =
1, main = "Importance of the Predictors")
```

Importance of the Predictors



#accuracy and expected error

```
attributes(rf_confusion)
```

```
## $names
```

```
## [1] "positive" "table"    "overall"  "byClass"  "mode"     "dots"
```

##

```
## $class
```

```
## [1] "confusionMatrix"
```

```
rf_confusion$overall
```

##	Accuracy	Kappa	AccuracyLower	AccuracyUpper	AccuracyNull
----	----------	-------	---------------	---------------	--------------

##	0.9906959	0.9882329	0.9883155	0.9927002	0.2844762
----	-----------	-----------	-----------	-----------	-----------

```
## AccuracyPValue  McnemarPValue
```

```
##          0.0000000      NaN
```

```
rf_confusion$overall['Accuracy']
```

Accuracy

```
## 0.9906959
```

```
rf_confusion$overall['AccuracyUpper']
```

```
## AccuracyUpper
```

```
## 0.9927002
```

```
rf_confusion$overall['AccuracyLower']
```

```
## AccuracyLower
##      0.9883155

testing_rfpred <- predict(rfFit, newdata=testing_clean)
testing_rfpred

## [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
```

writing out the predictions

```
pml_write_files = function(x){
  n = length(x)
  for(i in 1:n){
    filename = paste0("./practical_ml_wk4_",i,".txt")

    write.table(x[i],file=filename,quote=FALSE,row.names=FALSE,col.names=FALSE)
  }
}

pml_write_files(testing_rfpred)

testing_rfpred

## [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
```