

Predictions using the given Dataset

Raghavandaar

Executive Summary

We will try to perform partial machine learning on the dataset

We will be doing the following:

Process the data Explore the data Data Modeling and predicting

Processing

We are going to process the data from the dataset

```
trainRaw <- read.csv("pml-training.csv")
testRaw <- read.csv("pml-testing.csv")
```

Exploratory data analyses

We will try to analyse and explore the data

```
dim(trainRaw)
```

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

```
dim(testRaw)
```

```
## [1] 20 160
```

A lot of NA values present. Therefore we remove that.

```
sum(complete.cases(trainRaw))
```

```
## [1] 406
```

```
trainRaw <- trainRaw[, colSums(is.na(trainRaw)) == 0]
testRaw <- testRaw[, colSums(is.na(testRaw)) == 0]
```

```

classe <- trainRaw$classe
trainRemove <- grepl("^X|timestamp|window", names(trainRaw))
trainRaw <- trainRaw[, !trainRemove]
trainCleaned <- trainRaw[, sapply(trainRaw, is.numeric)]
trainCleaned$classe <- classe
testRemove <- grepl("^X|timestamp|window", names(testRaw))
testRaw <- testRaw[, !testRemove]
testCleaned <- testRaw[, sapply(testRaw, is.numeric)]

```

```

set.seed(22519)
inTrain <- createDataPartition(trainCleaned$classe, p=0.70, list=F)
trainData <- trainCleaned[inTrain, ]
testData <- trainCleaned[-inTrain, ]

```

The data has been explored and cleansed. We now move on to model selection.

Data Modeling and predicting

We will try to modify the dataset for model selection using random forest

```

controlRf <- trainControl(method="cv", 5)
modelRf <- train(classe ~ ., data=trainData, method="rf", trControl=controlRf, ntree=250)
modelRf

```

```

## Random Forest
##
## 13737 samples
##    52 predictor
##    5 classes: 'A', 'B', 'C', 'D', 'E'
##
## No pre-processing
## Resampling: Cross-Validated (5 fold)
## Summary of sample sizes: 10988, 10989, 10989, 10991, 10991
## Resampling results across tuning parameters:
##
##  mtry  Accuracy   Kappa
##    2    0.9912654 0.9889499
##   27    0.9916291 0.9894104
##   52    0.9842766 0.9801110
##
## Accuracy was used to select the optimal model using the largest value.
## The final value used for the model was mtry = 27.

```

Performance and final analysis

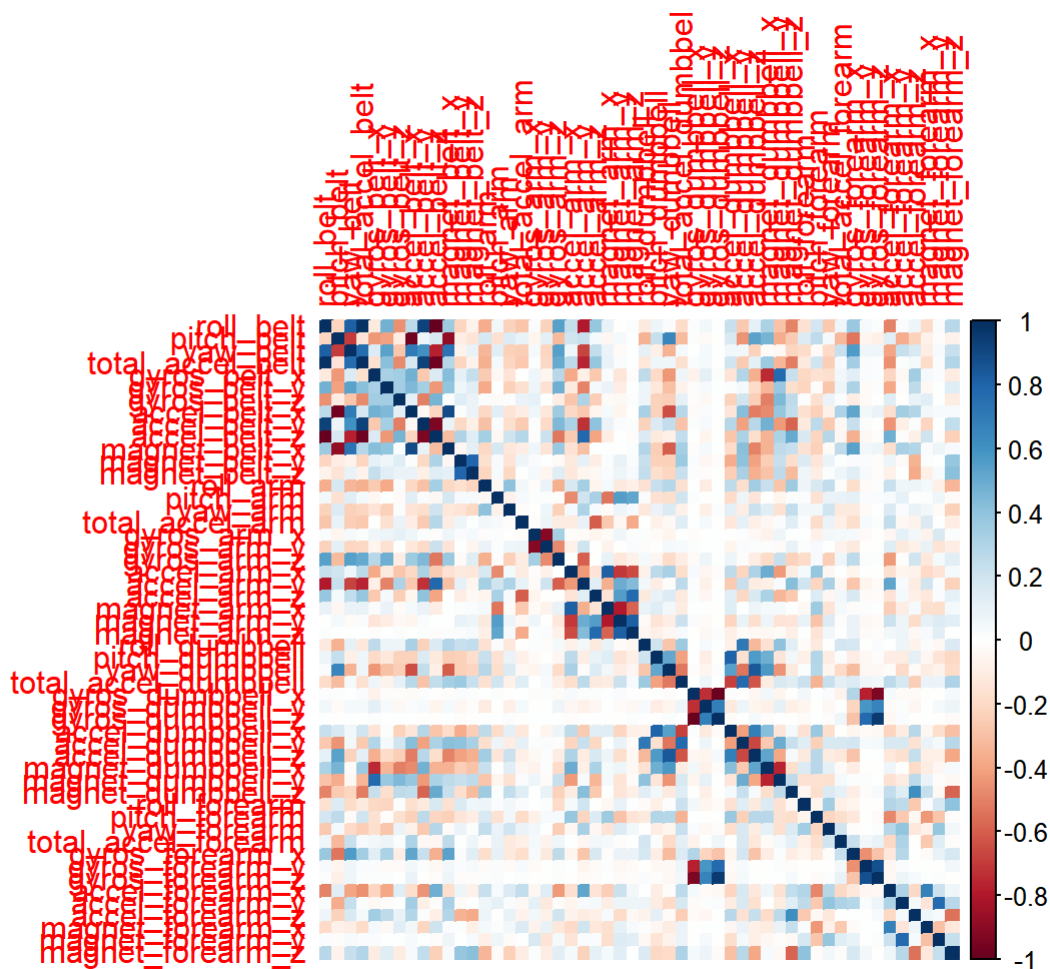
```

result <- predict(modelRf, testCleaned[, -length(names(testCleaned))])
result

```

```
## [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
```

```
corrPlot <- cor(trainData[, -length(names(trainData))])
corrplot(corrPlot, method="color")
```



```
treeModel <- rpart(classe ~ ., data=trainData, method="class")
prp(treeModel) # fast plot
```

