

Practical machine learning project

Fuyu

July 21, 2020

Overview

In this project, we first download the data and split the train dataset to two separate datasets: training and testing. Because the dimensions of the data is so huge, and there is a lot of missing values. We need to drop these missin values and reduce the dimension of predictors. Then we fit a classification model and test it on the test dateset. Finally we apply the model to the new data and make predictions.

```
path<-getwd()
url <- "https://d396qusza40orc.cloudfront.net/predmachlearn/pml-training.csv"
download.file(url, file.path(path, "train.csv"))
url <- "https://d396qusza40orc.cloudfront.net/predmachlearn/pml-testing.csv"
download.file(url, file.path(path, "test.csv"))
train<-read.csv("train.csv")
test<-read.csv("test.csv")
library(caret)
```

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

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

```
library(AppliedPredictiveModeling)
library(dplyr)
```

```
##
```

```
## Attaching package: 'dplyr'
```

```
## The following objects are masked from 'package:stats':
```

```
##
```

```
##      filter, lag
```

```
## The following objects are masked from 'package:base':
```

```
##
```

```
##      intersect, setdiff, setequal, union
```

```
library(rpart.plot)
```

```
## Loading required package: rpart
```

```
library(rpart)
library(rattle)
```

```
## Loading required package: tibble
```

```
## Loading required package: bitops
```

```
## Rattle: A free graphical interface for data science with R.
```

```
## XXXX 5.4.0 Copyright (c) 2006-2020 Togaware Pty Ltd.
```

```
## Type 'rattle()' to shake, rattle, and roll your data.
names(train)[names(train)=="X"]="class"
names(test)[names(test)=="X"]="class"
train$class<-as.factor(train$class)
test$class<-as.factor(test$class)
```

Exclude missing value

```
drop_NA<-function(df){
  index<-NULL
  for( i in 1:dim(df)[2] ){
    if(is.na(df[1,i])==T){index<-c(index,i)}
  }
  df<-df[, -c(1,2,index)]
}
train1<-drop_NA(train)
test1<-drop_NA(test)
```

Reduce the dimension of predictors

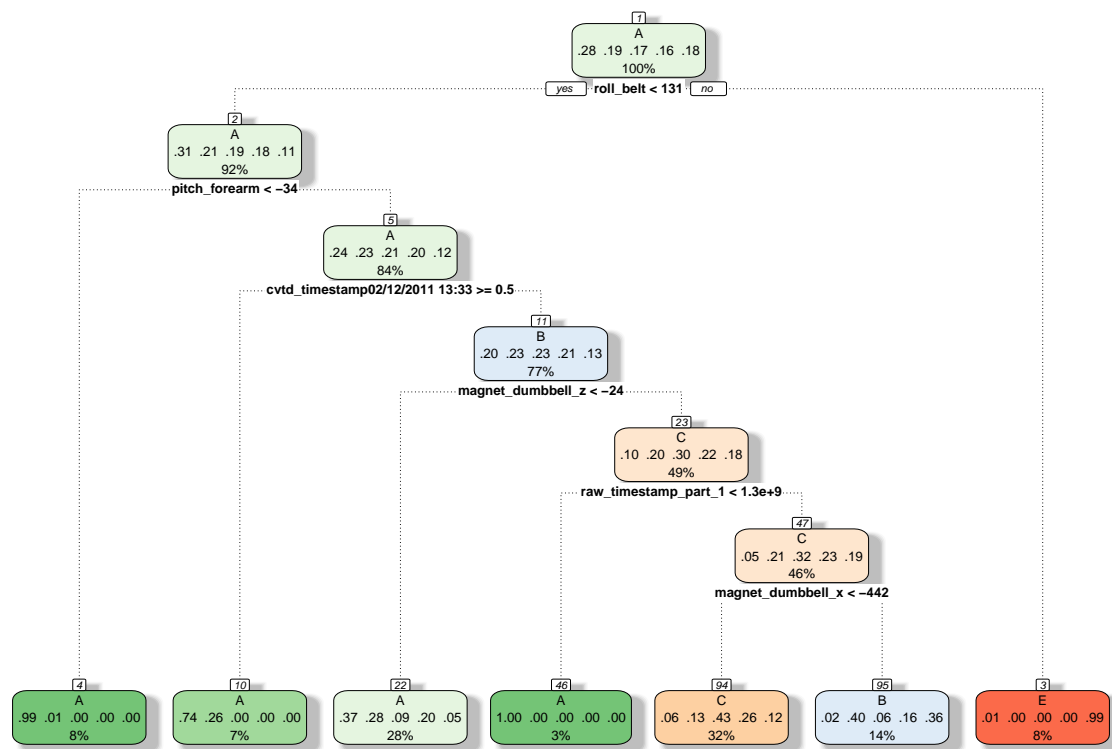
```
intrain<-createDataPartition(y=train1$class, p=0.7, list=F)
train1_train<-train1[intrain,]
train1_test<-train1[-intrain,]
NZV <- nearZeroVar(train1_train)
train1_train <- train1_train[, -NZV]
train1_test <- train1_test[, -NZV]
dim(train1_train)
```

```
## [1] 13737    57
```

Build a classification trees

```
set.seed(1)
mod<-train(class~., data=train1_train, method="rpart")

fancyRpartPlot(mod$finalModel)
```



Rattle 2020-Jul-21 12:33:38 feather

Crossvalidate the model

```
pred<-predict(mod,train1_test)
tb<-table(pred,train1_test$classe)
confusionMatrix(tb)
```

Confusion Matrix and Statistics

##

##

pred A B C D E

A 1538 593 161 337 95

B 21 304 58 138 250

C 110 242 807 489 237

D 0 0 0 0 0

E 5 0 0 0 500

##

Overall Statistics

##

Accuracy : 0.5351

95% CI : (0.5222, 0.5479)

No Information Rate : 0.2845

P-Value [Acc > NIR] : < 2.2e-16

##

Kappa : 0.3973

##

```
## McNemar's Test P-Value : NA
##
## Statistics by Class:
##
##           Class: A Class: B Class: C Class: D Class: E
## Sensitivity      0.9188 0.26690 0.7865 0.0000 0.46211
## Specificity      0.7184 0.90160 0.7781 1.0000 0.99896
## Pos Pred Value   0.5646 0.39429 0.4281      NaN 0.99010
## Neg Pred Value   0.9570 0.83672 0.9452 0.8362 0.89182
## Prevalence       0.2845 0.19354 0.1743 0.1638 0.18386
## Detection Rate   0.2613 0.05166 0.1371 0.0000 0.08496
## Detection Prevalence 0.4629 0.13101 0.3203 0.0000 0.08581
## Balanced Accuracy 0.8186 0.58425 0.7823 0.5000 0.73053
```

Therefore, here the **accuary** is **0.6503** and the **out of sample error** is **0.35**

Applying the model to the test data

```
predict(mod,test1)
```

```
## [1] A A C A A C C C A A C C B A C B B A A B
## Levels: A B C D E
```