

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

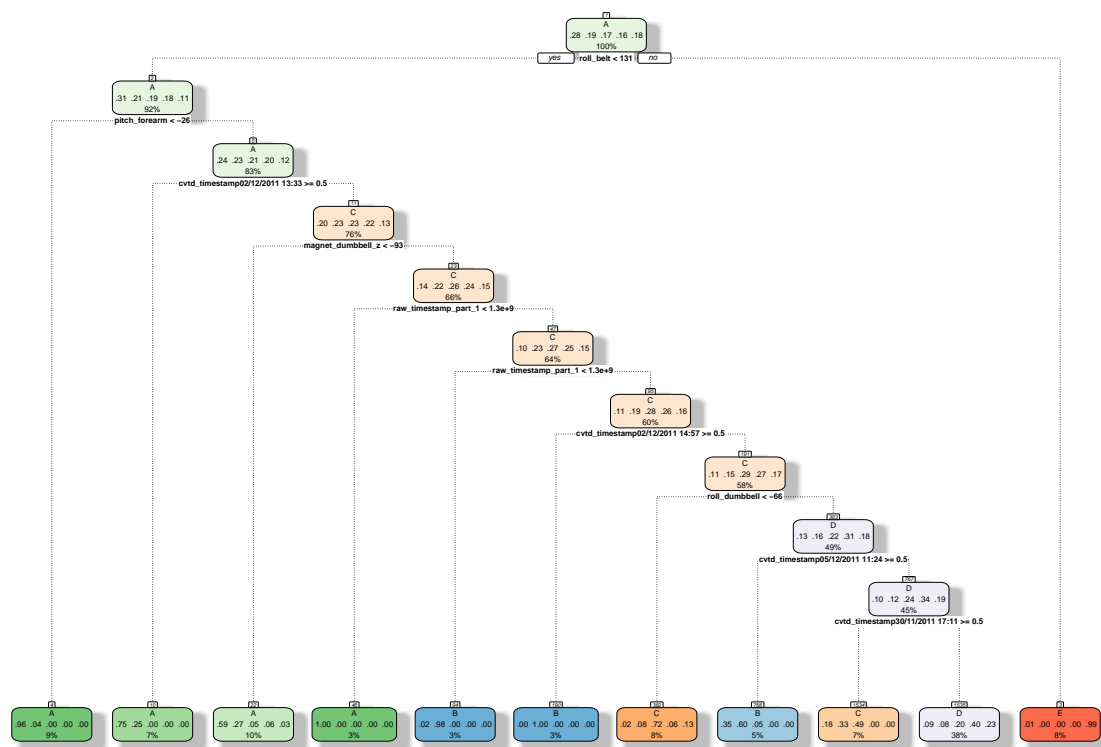
```
intrain<-createDataPartition(y=train1$class, p=0.7, list=F)
train1_train<-train1[intrain,]
train1_test<-train1[-intrain,]
NZV <- nearZeroVar(train1_train)
train1_train <- train1_train[, -NZV]
train1_test <- train1_test[, -NZV]
dim(train1_train)
```

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

Build a classification trees

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

fancyRpartPlot(mod$finalModel)
```



Rattle 2020-Jul-21 13:09:47 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 1271 297 27 26 9

B 91 508 16 0 0

C 86 153 547 22 54

D 221 181 436 916 541

E 5 0 0 0 478

##

Overall Statistics

##

Accuracy : 0.6321

95% CI : (0.6196, 0.6445)

No Information Rate : 0.2845

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

##

Kappa : 0.5381

##

```
## McNemar's Test P-Value : NA
##
## Statistics by Class:
##
##           Class: A Class: B Class: C Class: D Class: E
## Sensitivity      0.7593 0.44601 0.53314 0.9502 0.44177
## Specificity      0.9147 0.97745 0.93517 0.7198 0.99896
## Pos Pred Value   0.7798 0.82602 0.63457 0.3991 0.98965
## Neg Pred Value    0.9053 0.88027 0.90464 0.9866 0.88819
## Prevalence       0.2845 0.19354 0.17434 0.1638 0.18386
## Detection Rate    0.2160 0.08632 0.09295 0.1556 0.08122
## Detection Prevalence 0.2770 0.10450 0.14647 0.3900 0.08207
## Balanced Accuracy 0.8370 0.71173 0.73416 0.8350 0.72037
```

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