

R Notebook: random forest on dataset with classes - full modeal

```
library(caret)

## Loading required package: ggplot2
## Loading required package: lattice
library(ggplot2)
library(randomForest)

## randomForest 4.6-14
## Type rfNews() to see new features/changes/bug fixes.
##
## Attaching package: 'randomForest'
## The following object is masked from 'package:ggplot2':
##
##   margin
library(tidyverse)

## -- Attaching packages ----- tidyverse 1.3.1 --
## v tibble  3.1.5    v dplyr   1.0.7
## v tidyr   1.1.4    v stringr 1.4.0
## v readr   2.0.2    v forcats 0.5.1
## v purrr   0.3.4

## -- Conflicts ----- tidyverse_conflicts() --
## x dplyr::combine()      masks randomForest::combine()
## x dplyr::filter()       masks stats::filter()
## x dplyr::lag()          masks stats::lag()
## x purrr::lift()         masks caret::lift()
## x randomForest::margin() masks ggplot2::margin()

setwd("~/CSP571ProjectGroup")
df <- read_csv("df_with_class.csv")

## New names:
## * `` -> ...1

## Rows: 114660 Columns: 17

## -- Column specification -----
## Delimiter: ","
## chr (10): number, incident_state, sys_updated_by, contact_type, category, su...
## dbl (6): ...1, reassignment_count, reopen_count, sys_mod_count, problem_id,...
## lgl (1): made_sla

##
## i Use `spec()` to retrieve the full column specification for this data.
## i Specify the column types or set `show_col_types = FALSE` to quiet this message.
```

```

df %>% group_by(class) %>% summarise(n = n())

## # A tibble: 19 x 2
##   class      n
##   <chr>   <int>
## 1 >14days  5807
## 2 0mins    27618
## 3 10days   3499
## 4 10mins   2326
## 5 14days   5356
## 6 1day     9859
## 7 1hr      2934
## 8 28days   6425
## 9 2days   8289
## 10 30mins  4331
## 11 3days   5244
## 12 3hr     5316
## 13 4days   4079
## 14 5days   3609
## 15 5mins   5890
## 16 6days   3551
## 17 6hr     3928
## 18 7days   3617
## 19 8days   2982

df = subset(df, select=-c(sys_updated_by, number, subcategory))
df$class = as.factor(df$class)

# try on the full data set (df_small from previous version of the code)
df_small = df
train_idx = createDataPartition(df_small$resolved_updated_diff, p = 0.8, list=F)
df_train = df_small[train_idx,]
df_test = df_small[-train_idx,]
# try a randomforest model
rf_mod = randomForest(class~., data=df_train, ntree=1000)
df_test_X = subset(df_test, select=-c(class))
df_test_y = df_test$class
pred = predict(rf_mod, df_test)

columns = c("pred", "actual")
acc_df = data.frame(matrix(nrow=nrow(df_test), ncol = length(columns)))
colnames(acc_df) = columns
acc_df$pred = c(pred)
acc_df$actual = c(df_test_y)
#acc_df
acc = nrow(acc_df[acc_df$pred == acc_df$actual,])
acc/nrow(df_test)

## [1] 0.8789795

table(pred, df_test$class)

##
## pred      >14days 0mins 10days 10mins 14days 1day  1hr 28days 2days 30mins

```

```

## >14days 1193 0 0 0 0 0 0 0 0 0
## 0mins 0 5560 0 0 0 0 0 0 0 0
## 10days 0 0 520 0 0 0 0 0 0 0
## 10mins 0 0 0 294 0 0 0 0 0 0
## 14days 0 0 9 0 963 0 0 0 0 0
## 1day 0 0 103 95 43 1936 234 0 0 63
## 1hr 0 0 0 0 0 0 305 0 0 6
## 28days 2 0 0 0 2 0 0 1261 0 0
## 2days 0 0 36 3 0 7 3 0 1628 0
## 30mins 0 0 5 23 6 10 9 6 14 717
## 3days 0 0 5 0 0 0 0 0 0 0
## 3hr 0 0 0 51 0 0 76 0 0 57
## 4days 0 0 3 0 0 0 0 0 0 0
## 5days 0 0 5 0 0 0 0 0 0 0
## 5mins 0 0 0 5 0 0 1 0 0 2
## 6days 0 0 10 0 0 0 0 0 0 0
## 6hr 0 0 0 0 0 0 2 0 0 0
## 7days 0 0 3 0 0 0 0 0 0 0
## 8days 0 0 12 0 0 0 0 0 0 0
##
## pred 3days 3hr 4days 5days 5mins 6days 6hr 7days 8days
## >14days 0 0 0 0 0 0 0 0 0 0
## 0mins 0 0 0 0 0 0 0 0 0 0
## 10days 0 0 0 0 0 0 0 0 6 11
## 10mins 0 2 0 0 0 0 0 0 0 0
## 14days 0 0 1 0 0 0 0 0 10 6
## 1day 92 71 146 137 2 145 231 121 117
## 1hr 0 1 0 0 0 0 1 0 0
## 28days 0 0 0 1 0 0 0 1 1
## 2days 53 14 67 52 0 72 6 67 77
## 30mins 11 19 13 13 1 3 14 7 5
## 3days 883 0 18 1 0 5 0 4 11
## 3hr 0 877 0 0 5 0 9 0 0
## 4days 2 0 600 7 0 5 0 6 5
## 5days 0 0 8 467 0 12 0 7 10
## 5mins 0 1 0 1 1136 0 0 0 0
## 6days 0 0 0 12 0 467 0 23 26
## 6hr 0 7 0 0 0 0 552 0 0
## 7days 0 0 3 8 0 18 0 466 40
## 8days 0 0 0 3 0 4 0 6 330

```

```
summary(rf_mod)
```

```

##          Length Class Mode
## call          4 -none- call
## type           1 -none- character
## predicted     91730 factor numeric
## err.rate      20000 -none- numeric
## confusion      380 -none- numeric
## votes        1742870 matrix numeric
## oob.times      91730 -none- numeric
## classes        19 -none- character
## importance      13 -none- numeric
## importanceSD     0 -none- NULL
## localImportance  0 -none- NULL

```

## proximity	0 -none- NULL
## ntree	1 -none- numeric
## mtry	1 -none- numeric
## forest	14 -none- list
## y	91730 factor numeric
## test	0 -none- NULL
## inbag	0 -none- NULL
## terms	3 terms call