

Lab 8

Jaime Lin

11:59PM April 29, 2021

I want to make some use of my CART package. Everyone please try to run the following:

```
if (!pacman::p_isinstalled(YARF)){
  pacman::p_install_gh("kapelner/YARF/YARFJARs", ref = "dev")
  pacman::p_install_gh("kapelner/YARF/YARF", ref = "dev", force = TRUE)
}

## Installation failed: NULL : Failed to install 'unknown package' from
GitHub:
## HTTP error 403.
## API rate limit exceeded for 98.116.215.80. (But here's the good news:
Authenticated requests get a higher rate limit. Check out the documentation
for more details.)
##
## Rate limit remaining: 0/60
## Rate limit reset at: 2021-04-30 03:16:39 UTC
##
## To increase your GitHub API rate limit
## - Use `usethis::browse_github_pat()` to create a Personal Access Token.
## - Use `usethis::edit_r_environ()` and add the token as `GITHUB_PAT`.

##
## The following packages were not able to be installed:
## YARF

## Installation failed: NULL : Failed to install 'unknown package' from
GitHub:
## HTTP error 403.
## API rate limit exceeded for 98.116.215.80. (But here's the good news:
Authenticated requests get a higher rate limit. Check out the documentation
for more details.)
##
## Rate limit remaining: 0/60
## Rate limit reset at: 2021-04-30 03:16:39 UTC
##
## To increase your GitHub API rate limit
## - Use `usethis::browse_github_pat()` to create a Personal Access Token.
## - Use `usethis::edit_r_environ()` and add the token as `GITHUB_PAT`.

##
## The following packages were not able to be installed:
## YARF
```

```

options(java.parameters = "-Xmx4000m")
pacman::p_load(YARF)

## Warning: package 'YARF' is not available for this version of R
##
## A version of this package for your version of R might be available
## elsewhere,
## see the ideas at
## https://cran.r-project.org/doc/manuals/r-patched/R-admin.html#Installing-packages

## Warning: unable to access index for repository
http://www.stats.ox.ac.uk/pub/RWin/bin/windows/contrib/4.0:
## cannot open URL
'http://www.stats.ox.ac.uk/pub/RWin/bin/windows/contrib/4.0/PACKAGES'

## Warning in p_install(package, character.only = TRUE, ...):
## Warning in library(package, lib.loc = lib.loc, character.only = TRUE,
## logical.return = TRUE, : there is no package called 'YARF'

## Warning in pacman::p_load(YARF): Failed to install/load:
## YARF

```

For many of you it will not work. That's okay.

Throughout this part of this assignment you can use either the tidyverse package suite or data.table to answer but not base R. You can mix data.table with magrittr piping if you wish but don't go back and forth between tbl_df's and data.table objects.

```

pacman::p_load(tidyverse, magrittr, data.table)

## also installing the dependencies 'vctrs', 'dplyr', 'lubridate', 'tidyr'

## Warning: unable to access index for repository
http://www.stats.ox.ac.uk/pub/RWin/bin/windows/contrib/4.0:
## cannot open URL
'http://www.stats.ox.ac.uk/pub/RWin/bin/windows/contrib/4.0/PACKAGES'

##
## There is a binary version available but the source version is later:
## binary source needs_compilation
## vctrs 0.3.7 0.3.8 TRUE
##
## Binaries will be installed
## package 'vctrs' successfully unpacked and MD5 sums checked

## Warning: cannot remove prior installation of package 'vctrs'

## Warning in file.copy(savedcopy, lib, recursive = TRUE): problem copying
## D:\R language\R-4.0.3\library\00LOCK\vctrs\libs\x64\vctrs.dll to D:\R
## language\R-4.0.3\library\vctrs\libs\x64\vctrs.dll: Permission denied

```

```

## Warning: restored 'vctrs'

## package 'dplyr' successfully unpacked and MD5 sums checked

## Warning: cannot remove prior installation of package 'dplyr'

## Warning in file.copy(savedcopy, lib, recursive = TRUE): problem copying
## D:\R language\R-4.0.3\library\00LOCK\dplyr\libs\x64\dplyr.dll to D:\R
## language\R-4.0.3\library\dplyr\libs\x64\dplyr.dll: Permission denied

## Warning: restored 'dplyr'

## package 'lubridate' successfully unpacked and MD5 sums checked

## Warning: cannot remove prior installation of package 'lubridate'

## Warning in file.copy(savedcopy, lib, recursive = TRUE): problem copying D:
## \R language\R-4.0.3\library\00LOCK\lubridate\libs\x64\lubridate.dll to
D:\R
## language\R-4.0.3\library\lubridate\libs\x64\lubridate.dll: Permission
denied

## Warning: restored 'lubridate'

## package 'tidyr' successfully unpacked and MD5 sums checked

## Warning: cannot remove prior installation of package 'tidyr'

## Warning in file.copy(savedcopy, lib, recursive = TRUE): problem copying
## D:\R language\R-4.0.3\library\00LOCK\tidyr\libs\x64\tidyr.dll to D:\R
## language\R-4.0.3\library\tidyr\libs\x64\tidyr.dll: Permission denied

## Warning: restored 'tidyr'

## package 'tidyverse' successfully unpacked and MD5 sums checked
##
## The downloaded binary packages are in
## C:\Users\jaime\AppData\Local\Temp\RtmpYvYKFd\downloaded_packages
##
## tidyverse installed

## Warning: package 'tidyverse' was built under R version 4.0.5

## Warning in pacman::p_load(tidyverse, magrittr, data.table): Failed to
install/load:
## tidyverse

```

We will be using the storms dataset from the dplyr package. Filter this dataset on all storms that have no missing measurements for the two diameter variables, “ts_diameter” and “hu_diameter”.

```

'
data("storms")

```

```

storms2 = storms %>% filter(!is.na(ts_diameter) & !is.na(hu_diameter) &
ts_diameter> 0 & hu_diameter> 0 )
storms2
'

## [1] "\ndata(\"storms\")\nstorms2 = storms %>% filter(!is.na(ts_diameter) &
!is.na(hu_diameter) & ts_diameter> 0 & hu_diameter> 0 )\nstorms2\n"

```

From this subset, create a data frame that only has storm, observation period number for each storm (i.e., 1, 2, ..., T) and the “ts_diameter” and “hu_diameter” metrics.

```

'
storms2 = storms2 %>%
  select(name, ts_diameter, hu_diameter) %>%
  group_by(name) %>%
  mutate(period = row_number())
'

## [1] "\nstorms2 = storms2 %>%\n  select(name, ts_diameter, hu_diameter)
%>%\n  group_by(name) %>%\n  mutate(period = row_number())\n"

```

Create a data frame in long format with columns “diameter” for the measurement and “diameter_type” which will be categorical taking on the values “hu” or “ts”.

```

'
storms_long = pivot_longer(storms2, cols = matches("diameter"), names_to =
"diameter")
storms_long
'

## [1] "\nstorms_long = pivot_longer(storms2, cols = matches(\"diameter\"),
names_to = \"diameter\")\nstorms_long\n"

```

Using this long-formatted data frame, use a line plot to illustrate both “ts_diameter” and “hu_diameter” metrics by observation period for four random storms using a 2x2 faceting. The two diameters should appear in two different colors and there should be an appropriate legend.

```

'
storms_sample = sample(unique(storms2$name), 4)
ggplot(storms_long %>% filter(name %in% storms_sample)) +
  geom_line(aes(x = period, y = value, col = diameter)) +
  facet_wrap(name ~., nrow = 2)
'

## [1] "\nstorms_sample = sample(unique(storms2$name), 4)\nggplot(storms_long
%>% filter(name %in% storms_sample)) + \n  geom_line(aes(x = period, y =
value, col = diameter)) +\n  facet_wrap(name ~., nrow = 2)\n"

```

In this next first part of this lab, we will be joining three datasets in an effort to make a design matrix that predicts if a bill will be paid on time. Clean up and load up the three files. Then I'll rename a few features and then we can examine the data frames:

```
,
rm(list = ls())
pacman::p_load(tidyverse, magrittr, data.table, R.utils)
bills =
fread("https://github.com/kapelner/QC_MATH_342W_Spring_2021/raw/master/labs/b
ills_dataset/bills.csv.bz2")
payments =
fread("https://github.com/kapelner/QC_MATH_342W_Spring_2021/raw/master/labs/b
ills_dataset/payments.csv.bz2")
discounts =
fread("https://github.com/kapelner/QC_MATH_342W_Spring_2021/raw/master/labs/b
ills_dataset/discounts.csv.bz2")
setnames(bills, "amount", "tot_amount")
setnames(payments, "amount", "paid_amount")
head(bills)
head(payments)
head(discounts)
bills = as_tibble(bills)
payments = as_tibble(payments)
discounts = as_tibble(discounts)
,

## [1] "\nrm(list = ls())\npacman::p_load(tidyverse, magrittr, data.table,
R.utils)\nbills =
fread(\"https://github.com/kapelner/QC_MATH_342W_Spring_2021/raw/master/labs/
bills_dataset/bills.csv.bz2\")\npayments =
fread(\"https://github.com/kapelner/QC_MATH_342W_Spring_2021/raw/master/labs/
bills_dataset/payments.csv.bz2\")\ndiscounts =
fread(\"https://github.com/kapelner/QC_MATH_342W_Spring_2021/raw/master/labs/
bills_dataset/discounts.csv.bz2\")\nsetnames(bills, \"amount\",
\"tot_amount\")\nsetnames(payments, \"amount\",
\"paid_amount\")\nhead(bills)\nhead(payments)\nhead(discounts)\nbills =
as_tibble(bills)\npayments = as_tibble(payments)\ndiscounts =
as_tibble(discounts)\n"
```

The unit we care about is the bill. The y metric we care about will be “paid in full” which is 1 if the company paid their total amount (we will generate this y metric later).

Since this is the response, we would like to construct the very best design matrix in order to predict y.

I will create the basic steps for you guys. First, join the three datasets in an intelligent way. You will need to examine the datasets beforehand.

```
,
bills_with_payments = left_join(bills, payments, by = c("id"="bill_id"))
bills_with_payments
```

```

bills_with_payments_with_discounts = left_join(bills_with_payments,
discounts, by = c("discount_id" = "id"))
bills_with_payments_with_discounts
'

## [1] "\nbills_with_payments = left_join(bills, payments, by =
c(\"id\"=\"bill_id\"))\nbills_with_payments\nbills_with_payments_with_discounts = left_join(bills_with_payments, discounts, by = c(\"discount_id\" = \"id\"))\nbills_with_payments_with_discounts\n"

```

Now create the binary response metric `paid_in_full` as the last column and create the beginnings of a design matrix `bills_data`. Ensure the unit / observation is bill i.e. each row should be one bill!

```

'
bills_data = bills_with_payments_with_discounts %>%
  mutate(tot_amount = if_else(is.na(pct_off), tot_amount, tot_amount*(1-
pct_off/100))) %>%
  group_by(id) %>%
  mutate(sum_of_payments = sum(paid_amount)) %>%
  mutate(paid_in_full = if_else(sum_of_payments >= tot_amount, 1, 0, missing
= 0)) %>%
  slice(1)
table(bills_data$paid_in_full, useNA = "always")
'

## [1] "\nbills_data = bills_with_payments_with_discounts %>%\n
mutate(tot_amount = if_else(is.na(pct_off), tot_amount, tot_amount*(1-
pct_off/100))) %>%\n  group_by(id) %>%\n  mutate(sum_of_payments =
sum(paid_amount)) %>%\n  mutate(paid_in_full = if_else(sum_of_payments >=
tot_amount, 1, 0, missing = 0)) %>%\n  slice(1)\n
\ntable(bills_data$paid_in_full, useNA = \"always\") \n"

```

How should you add features from transformations (called “featurization”)? What data type(s) should they be? Make some features below if you think of any useful ones. Name the columns appropriately so another data scientist can easily understand what information is in your variables.

```

#pacman::p_load("lubricate")
#bills_data %>%
#  select(-id, -id.y, -num_days, -transaction_date, -pct_off, -
days_until_discount, -paid_amount) %>%
#  mutate(num_days_to_pay = as.integer(ymd(due_date) - ymd(invoice_date)))
%>%
#  select(-due_date, -invoice_date) %>%
#  mutate(discount_id = as.factor(discount_id)) %>%
#  group_by(customer_id) %>%
#  mutate(bill_num = row_number()) %>%
#  ungroup() %>%
#bills_data

```

Now let's do this exercise. Let's retain 25% of our data for test.

```
'  
K = 4  
test_indices = sample(1 : nrow(bills_data), round(nrow(bills_data) / K))  
train_indices = setdiff(1 : nrow(bills_data), test_indices)  
bills_data_test = bills_data[test_indices, ]  
bills_data_train = bills_data[train_indices, ]  
'  
  
## [1] "\nK = 4\ntest_indices = sample(1 : nrow(bills_data),  
round(nrow(bills_data) / K))\ntrain_indices = setdiff(1 : nrow(bills_data),  
test_indices)\nbills_data_test = bills_data[test_indices, ]\nbills_data_train  
= bills_data[train_indices, ]\n"
```

Now try to build a classification tree model for `paid_in_full` with the features (use the `Xy` parameter in `YARF`). If you cannot get `YARF` to install, use the package `rpart` (the standard R tree package) instead. You will need to install it and read through some documentation to find the correct syntax.

Warning: this data is highly anonymized and there is likely zero signal! So don't expect to get predictive accuracy. The value of the exercise is in the practice. I think this exercise (with the joining exercise above) may be one of the most useful exercises in the entire semester.

```
# paid_in_full = rpart(payments$paid_amount~ payments$transaction_date, data  
= bills_data)
```

For those of you who installed `YARF`, what are the number of nodes and depth of the tree?

```
#TO-DO
```

For those of you who installed `YARF`, print out an image of the tree.

```
#TO-DO
```

Predict on the test set and compute a confusion matrix.

```
'  
bills_data = na.omit(bills_data)  
set.seed(1)  
train_size = 5000  
test_indices = sample(1 : nrow(bills_data), round(nrow(bills_data) / K))  
train_indices = setdiff(1 : nrow(bills_data), test_indices)  
bills_data_test = bills_data[test_indices, ]  
bills_data_train = bills_data[train_indices, ]  
  
logistic_mod = glm(paid_amount ~ ., bills_data_train, family = "binomial")  
p_hats_train = predict(logistic_mod, bills_data_train, type = "response")  
p_hats_test = predict(logistic_mod, bills_data_test, type = "response")  
'
```

```

y_hats_train = factor(ifelse(p_hats_train >= 0.5, ">50K", "<=50K"))
mean(y_hats_train != y_train)
table(y_train, y_hats_train)
,

## [1] "\nbills_data = na.omit(bills_data) \nset.seed(1)\ntrain_size =
5000\ntest_indices = sample(1 : nrow(bills_data), round(nrow(bills_data) /
K))\ntrain_indices = setdiff(1 : nrow(bills_data),
test_indices)\nbills_data_test = bills_data[test_indices, ]\nbills_data_train
= bills_data[train_indices, ]\n\n\n logistic_mod = glm(paid_amount ~ .,
bills_data_train, family = \"binomial\")\n p_hats_train =
predict(logistic_mod, bills_data_train, type = \"response\")\n p_hats_test =
predict(logistic_mod, bills_data_test, type = \"response\")\n y_hats_train =
factor(ifelse(p_hats_train >= 0.5, \">50K\", \"<=50K\"))\n mean(y_hats_train
!= y_train)\n table(y_train, y_hats_train)\n"

```

Report the following error metrics: misclassification error, precision, recall, F1, FDR, FOR.

```

# n = sum(oos_conf_table)
# fp = oos_conf_table[1, 2]
# fn = oos_conf_table[2, 1]
# tp = oos_conf_table[2, 2]
# tn = oos_conf_table[1, 1]
# num_pred_pos = sum(oos_conf_table[, 2])
# num_pred_neg = sum(oos_conf_table[, 1])
# num_pos = sum(oos_conf_table[2, ])
# num_neg = sum(oos_conf_table[1, ])
# Misscalculation_error = (fn + fp) / n
# precision = tp / num_pred_pos
# cat("precision", round(precision * 100, 2), "%\n")
# recall = tp / num_pos
# cat("recall", round(recall * 100, 2), "%\n")
# f_1 = 2 / ((1/recall) + (1/precision))
# false_discovery_rate = 1 - precision
# cat("false_discovery_rate", round(false_discovery_rate * 100, 2), "%\n")
# false_omission_rate = fn / num_pred_neg
# cat("false_omission_rate", round(false_omission_rate * 100, 2), "%\n")

```

Is this a good model? (yes/no and explain).

#TO-DO

There are probability asymmetric costs to the two types of errors. Assign the costs below and calculate oos total cost.

```

# y_hats_test = factor(ifelse(p_hats_test >= 0.9, ">50K", "<=50K"))
# mean(y_hats_test != y_test)
# oos_conf_table = table(y_test, y_hats_test)
# oos_conf_table

```


We now wish to do asymmetric cost classification. Fit a logistic regression model to this data.

#to-do

Use the function from class to calculate all the error metrics for the values of the probability threshold being 0.001, 0.002, ..., 0.999 in a data frame.

```
treshold_data = cbind(1:999 / 1000)
treshold_t = factor(ifelse(treshold_data >= 0.5, ">0.5", "<=0.5"))
t = table(treshold_t, treshold_data)
p_th = nrow(treshold_data) / 1000
N = sum(t[, 1])
P = sum(t[, 2])
PN = sum(t[, 1])
PP = sum(t[, 2])
n = N + P
fn = p_th * N
fp = (1 - p_th) * N
tp = (1 - p_th) * P
tn = p_th * P
Misscalculation_error = (fn + fp) / n
precision = tp / PP
cat("precision", round(precision * 100, 2), "%\n")

## precision 0.1 %

recall = tp / P
cat("recall", round(recall * 100, 2), "%\n")

## recall 0.1 %

f_1 = 2 / ((1/recall) + (1/precision))
false_discovery_rate = 1 - precision
cat("false_discovery_rate", round(false_discovery_rate * 100, 2), "%\n")

## false_discovery_rate 99.9 %

false_omission_rate = fn / PN
cat("false_omission_rate", round(false_omission_rate * 100, 2), "%\n")

## false_omission_rate 99.9 %
```

Calculate the column total_cost and append it to this data frame.

colSums(total_cost)

Which is the winning probability threshold value and the total cost at that threshold?

#TO-DO

Plot an ROC curve and interpret.

#TO-DO

#TO-DO interpretation

Calculate AUC and interpret.

#TO-DO

#TO-DO interpretation

Plot a DET curve and interpret.

#TO-DO

#TO-DO interpretation