



DRIVE/Cast

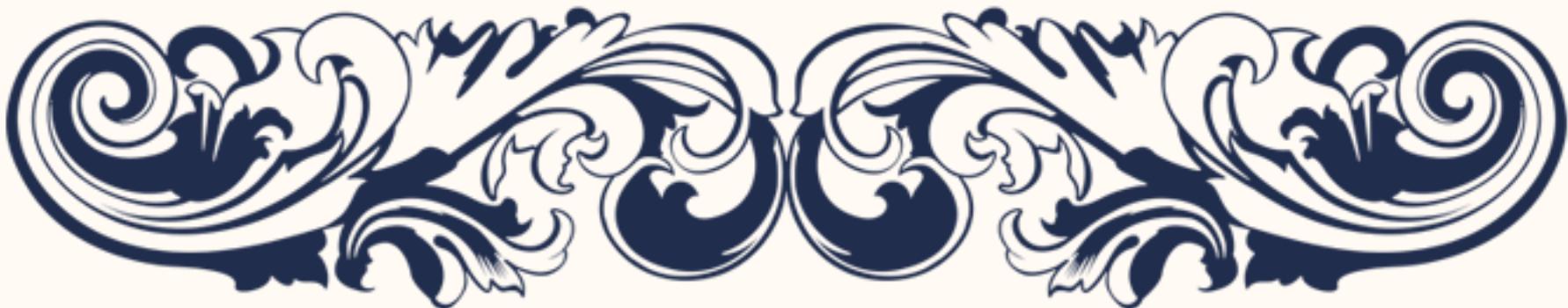
April 15, 21, and 27, 2021

#CASEdrive

Accelerate Smarter

READY, SET, BAKE

Recipes for Reproducible Reporting with RMarkdown



James Rogol | University of Virginia | 21 April 2021
DRIVE/Cast 2021

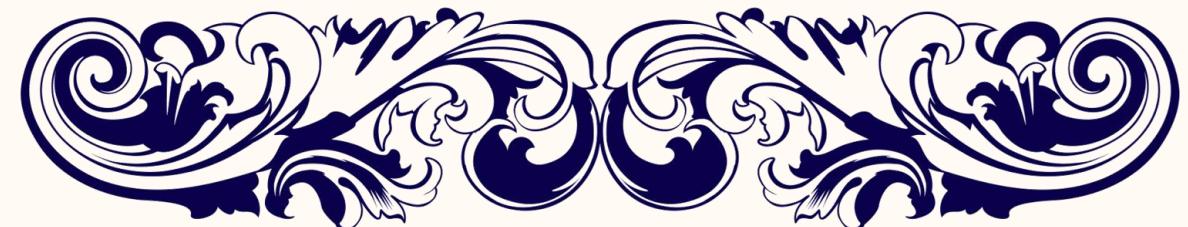




re·pro·duc·i·bil·i·ty, noun

- Repeatable
- Identical
- Shareable

THE GREAT BRITISH BAKING SHOW



#GBBO

THE GREAT BRITISH
BAKE OFF®



THE GREAT BRITISH BAKE OFF
STARTS TUESDAY 22ND SEPTEMBER. 8PM. CHANNEL 4.

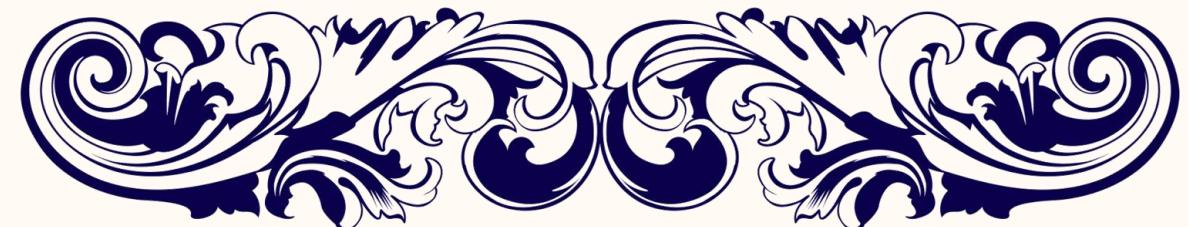


DRIVE/Cast 2021

re·pro·duc·i·bil·i·ty, noun

- Repeatable
- Identical
- Shareable

THE GREAT BRITISH BAKING SHOW



FOLLOW ALONG!



github.com/jrogol/ReadySetBake



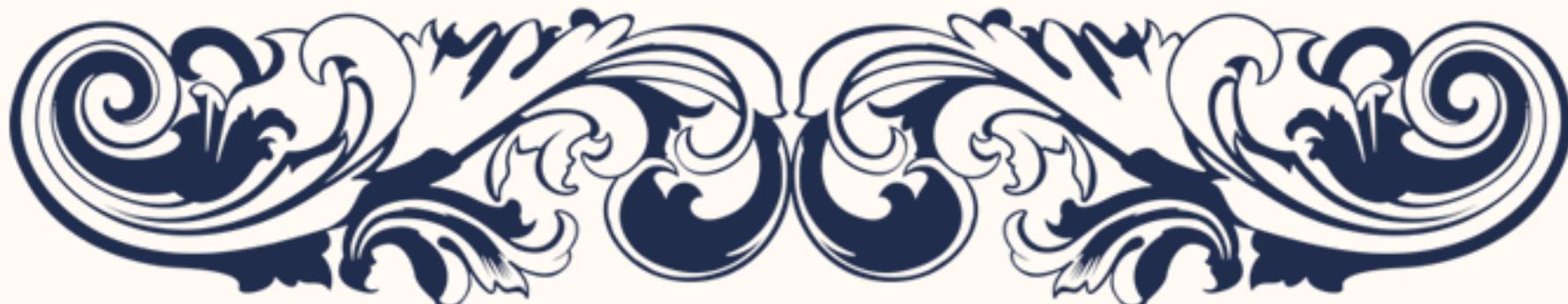
Reports (HTML & PDF)



Markdown (.Rmd)

CHALLENGE N° 1

THE SIGNATURE



DRIVE/Cast 2021

CHALLENGE N° 1

THE SIGNATURE

1-TheSignature.pdf

Contact Report Summary

Month of August

Major Gifts

05 March 2020

Reports filed in August, by Type.

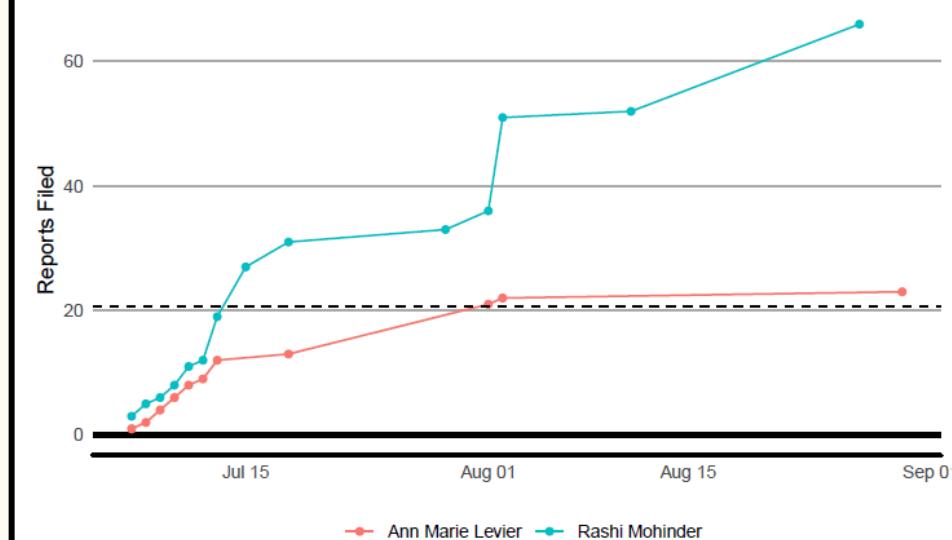
	Email	Phone	Visit	Total
Ann Marie Levier	3	3	4	10
Rashi Mohinder	10	5	18	33
Total	13	8	22	43

Reach and Outcomes

	Reports	Unique Donors	Positive Outcomes	Negative Outcomes
Ann Marie Levier	23	21	0	23
Rashi Mohinder	66	59	16	50

Year to Date Summary

Year to Date Activity (July to August)



The Dotted Line represents the Average number of reports filed across the organization.

DRIVE/Cast 2021

CHALLENGE N° 1

THE SIGNATURE

Contact Report Summary

Month of August

Major Gifts

05 March 2020

```
1 ---  
2 title: "Contact Report Summary"  
3 subtitle: "`r sprintf('Month of %s',params$endMonth)`"  
4 author: "`r params$dept`"  
5 date: "`r format(Sys.Date(),'%d %B %Y')`"  
6 output:  
7   pdf_document:  
8 params:  
9   endMonth: "August"  
10  dept: "Major Gifts"  
11 ---
```

Markdown/1-TheSignature.rmd
DRIVE/Cast 2021

DATA AVAILABILITY

“If the first line of your #rstats script is

```
setwd("C:\Users\jenny\path\that\only\I\have")
```

I will come into your lab and

SET YOUR COMPUTER ON FIRE.”

Jenny Bryan, Rstudio

Former Adjunct Professor

University of British Columbia

CHALLENGE N° 1

THE SIGNATURE

- Use `here`
 - RStudio Projects
 - Folder Structure
- R code Chunk

```
33  ````{r getData}
34  data <- read_csv(here("Data/contactReports.csv"))
35  ````
```

CHALLENGE N° 1

THE SIGNATURE

- Sections
- In-line code

```
69 • ## Reports filed in `r params$endMonth` , by Type.  
70  
71 •   ````{r byType}  
72 deptData %>%  
73   filter(month == params$endMonth) %>%  
74   count(staffName, method) %>%  
75   spread(method,n) %>%  
76   adorn_totals() %>%  
77   adorn_totals("col") %>%  
78   rename(" " = staffName) %>%  
79   kable(format = "latex")  
80   ...  
81  
82 • ## Reach and Outcomes  
83  
84 •   ````{r byOutcome}  
85 deptData %>%  
86   select(reportID,staffName, donor,outcome) %>%  
87   mutate(val = T) %>%  
88   spread(outcome, val, fill = F) %>%  
89   group_by(staffName) %>%  
90   summarize(reports = n(),  
91             uniqueDonors = n_distinct(donor),  
92             positive = sum(Positive),  
93             negative = sum(Negative)  
94           ) %>%  
95   kable(format = "latex",  
96         col.names = c("",  
97                         "Reports",  
98                         "Unique Donors",  
99                         "Positive Outcomes",  
100                        "Negative Outcomes"))  
101  
102 ...
```

CHALLENGE N° 1

THE SIGNATURE

- Rendering
- Parameters

```
170
171  ````{r render, eval=F}
172  rmarkdown::render(here("Markdown/1-TheSignature.Rmd"),
173                      output_dir = here("Reports"))
174
175
176  rmarkdown::render(here("Markdown/1-TheSignature.Rmd"),
177                      output_dir = here("Reports"),
178                      output_file = "1-TheSignature-AG_Sept",
179                      params = list(endMonth = "September",
180                        dept = "Planned Giving"))
181
```

```
|.....  
| 5%  
| inline R code fragments  
|.....  
| 10%  
label: setup (with options)  
List of 1  
$ include: logi FALSE  
|.....  
| 14%  
ordinary text without R code  
|.....  
| 19%  
label: getData  
|.....  
| 24%  
ordinary text without R code  
|.....  
| 29%  
label: cleanData1 (with options)  
List of 1
```

CHALLENGE N° 1

THE SIGNATURE

Contact Report Summary
Month of August
Major Gifts
05 March 2020

Reports filed in August, by Type.

	Email	Phone	Visit	Total
Ann Marie Levier	3	3	4	10
Rashi Mohinder	10	5	18	33
Total	13	8	22	43

Reach and Outcomes

	Reports	Unique Donors	Positive Outcomes	Negative Outcomes
Ann Marie Levier	23	21	0	23
Rashi Mohinder	66	59	16	50

Year to Date Summary

Year to Date Activity (July to August)

Date	Ann Marie Levier (Reports)	Rashi Mohinder (Reports)
Jul 15	2	2
Jul 22	4	4
Jul 29	6	6
Aug 05	8	8
Aug 12	12	12
Aug 19	14	14
Aug 26	16	16
Sep 02	21	21

The Dotted Line represents the Average number of reports filed across the organization.

1-TheSignature.pdf

Contact Report Summary
Month of September
Planned Giving
05 March 2020

Reports filed in September, by Type.

	Visit	Total
Deborah Mettier	2	2
Total	2	2

Reach and Outcomes

	Reports	Unique Donors	Positive Outcomes	Negative Outcomes
April Catson	32	29	12	20
Deborah Mettier	12	12	4	8

Year to Date Summary

Year to Date Activity (July to September)

Date	April Catson (Reports)	Deborah Mettier (Reports)
Jul 15	4	4
Jul 22	5	5
Jul 29	6	6
Aug 05	10	8
Aug 12	22	9
Aug 19	25	10
Aug 26	27	11
Sep 02	32	12

The Dotted Line represents the Average number of reports filed across the organization.

1-TheSignature-AG_sept.pdf

CHALLENGE N° 2

THE TECHNICAL



DRIVE/Cast 2021



KOUIGN AMANN





RASPBERRY BLANCMANGE

with Langues du Chat



CHALLENGE N° 2

THE TECHNICAL

```
1 ---  
2 title: "Thinking Through An Analysis"  
3 author: "James Rogol"  
4 date: "29 September 2018"  
5 output:  
6   html_document:  
7     toc: true  
8     toc_float: true  
9 ---
```

1-TheTechnical.html
1-TheTechnical.Rmd

Background
Approach
Data
Sampling
EDA
Analysis
Session Info

Thinking Through An Analysis

James Rogol

29 September 2018

Background

R Markdown can be a great tool when working through the approach to a new analysis, or troubleshooting (or reverse-engineering) old code. It offers the ability to combine code and prose in a single, self-contained document. As such, an analysts is able to capture the thought process behind the approach. This document will be presented as an example at the [2020 CASE Drive/ Conference](#).

Approach

In this case, we'd like to identify potential donors in the 2018 fiscal year. We've already pulled the data in question, including giving totals from 2013 to 2017. Data from FY13 to FY16 will be used to predict 2017 donors. The best model will then use data from FY14 to FY17 to predict FY18 donors.

Data

First, we'll need to obtain the data for the analysis. This anonymized dataset was originally obtained [here](#), and modified.¹

```
basicReport <- read_csv(here("Data/basicReport.csv"))
```

There are 300 observations in the data, with 36 columns. 14 of the columns are `numeric`, 3 are `dates`.

Data Wrangling

Age

The `age` variable ranges from 0 to 85. Values of 0 do *not* make sense, and should be replaced with `NA`. Similarly, `age_bin` can be set to a value of "Unknown".

```
basicReport <- basicReport %>%  
  mutate(age = if_else(age == 0, NA_real_, age),  
        age_bin = if_else(is.na(age), "Unknown", age_bin))
```

Factor Variables

Of the 19 text columns above, `age_bin`, `gender`, and `address_type` are more like factors, and should be treated as such. `age_bin` can be an ordered factor, as well.

```
basicReport <- basicReport %>%  
  mutate_at(vars(gender, address_type), as.factor) %>%  
  mutate(age_bin = factor(age_bin, ordered = T))
```

CHALLENGE N° 2

THE TECHNICAL

- Prose
 - Text Formatting
- Code Chunks
 - Named
 - Simple, monotonic

```
36 ## Data
37
38 First, we'll need to obtain the data for the analysis.
This anonymized dataset was originally obtained
[here](https://www.kaggle.com/michaelpawlus/fundraising-data), and modified.[^modify]
39
40 {r getData}
41 basicReport <- read_csv(here("Data/basicReport.csv"))
42
43
44 There are `r nrow(basicReport)` observations in the data,
with `r ncol(basicReport)` columns. `r
ncol(select_if(basicReport, is.numeric))` of the columns
are **numeric**, `r
ncol(select_if(basicReport, !inherits(., "POSIXct")))` are
**dates**.
45
46
47 ### Data wrangling
48
49 #### Age
50
51 The `age` variable ranges from `r min(basicReport$age,
na.rm = T)` to `r max(basicReport$age, na.rm = T)` .
Values of 0 do _not_ make sense, and should be replaced
with `NA`. Similarly, `age_bin` can be set to a value of
"Unknown".
52
53 {r cleanAge}
54 basicReport <- basicReport %>%
55   mutate(age = if_else(age == 0, NA_real_, age),
56         age_bin = if_else(is.na(age), "Unknown", age_bin))
57
58
```

CHALLENGE N° 2

THE TECHNICAL

- Randomness
- Functions
 - Simple, monotonic

Sampling

```
trainingPer <- .8
```

In order to measure model performance, we'll need to split the data into *Training* and *Testing* sets. 80% of the data will be used for training, and the other 20% will be used for testing.

```
set.seed(2020)

trainingInd <- sample(nrow(basicReport), nrow(basicReport)*trainingPer)

training <- basicReport[trainingInd,]

testing <- basicReport[-trainingInd,]
```

In the future, we can use the above as a model for a function. It will take three arguments: `df` (the data frame in question), `seed` for the random seed, and `per` for the percentage of observations to use as training.

```
sampleData <- function(df, per, seed = 2020){

  # Basic Error Handling
  stopifnot(inherits(df, "data.frame"))

  set.seed(seed)

  ind <- sample(nrow(df), nrow(df)*per)

  training <- df[ind,]

  testing <- df[-ind,]

  out <- list()

  out$training <- training
  out$testing <- testing

  return(out)
}
```

CHALLENGE N° 2

THE TECHNICAL

- Randomness
- Functions
 - Simple, monotonic
- Dependencies

Session Info

Analysis and Report generated on Windows 10 x64 (build 17763) using R version 3.5.3 (2019-03-11).

```
# Results = 'asis' renders the output as if it were text in the markdown doc. In
# this case, the script renders as a bulleted list.

# Get a list of attached packages
packages <- sessioninfo::package_info()

# Combine Package name and version number
p <- sprintf("%s (%s)",
             packages$package[packages$attached == T],
             packages$loadedversion[packages$attached == T])
```

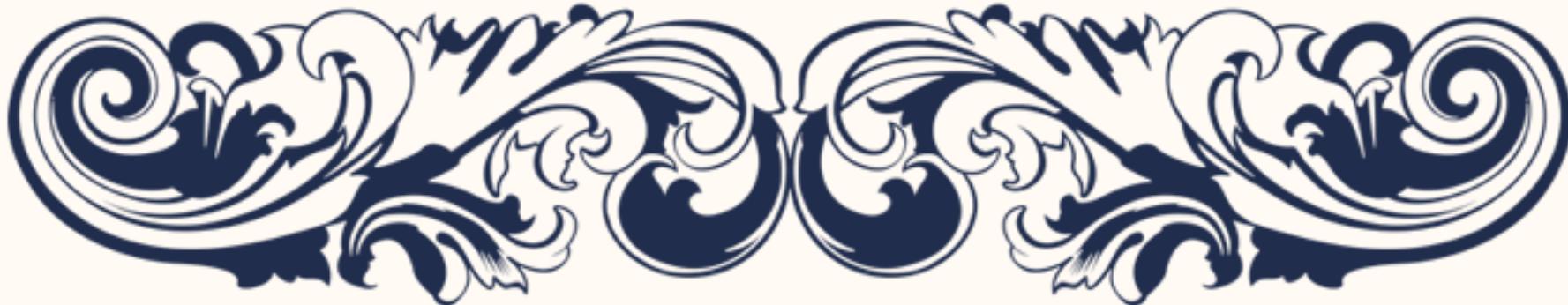
```
# Iterate over the packages, adding a hyphen and space at the start of the string
# (to create a bulleted list) and a new line character at the end of the string
# to list each package as a new bullet point.
```

```
for (x in p) {
  cat("- ", x, "\n")}
```

- dplyr (0.8.3)
- forcats (0.4.0)
- ggplot2 (3.2.1)
- here (0.1)
- janitor (1.1.1)
- kableExtra (1.1.0)
- knitr (1.28)
- readr (1.3.1)
- tidyverse (0.8.3)

CHALLENGE N° 3

THE SHOWSTOPPER





FY 2020 Giving Profile

School of Engineering Alumni

10 March 2020

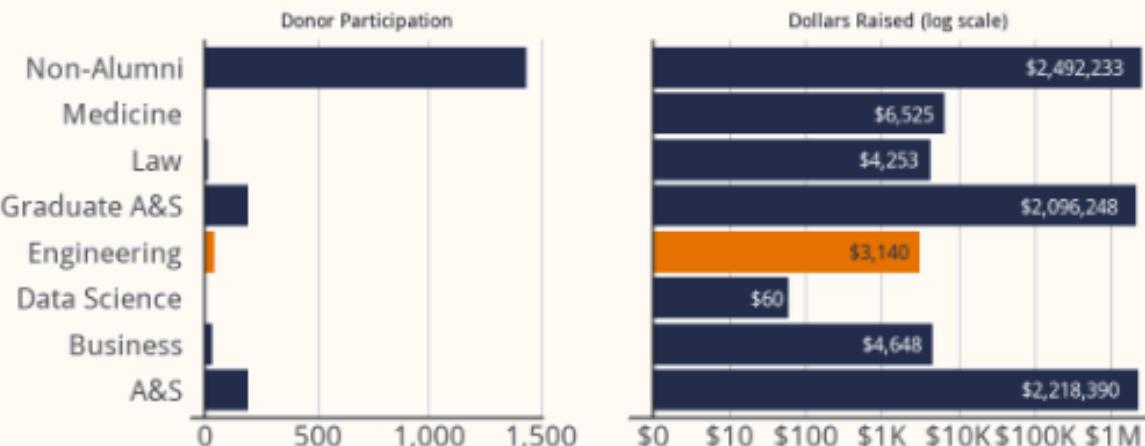
Overview

By Region

Session Info

Overview

43 Alumni from the School of Engineering gave in the 2020 Fiscal Year. This accounts for 2% of donor participation.



By Region

Of the 1,909 donors in FY 2020, 1,905 (100%) have a valid Zip Code.

418 entities reside outside of a CBSA. The table below shows the top 10 regions by count.

CHALLENGE N° 3

THE SHOWSTOPPER

- Custom CSS
 - JQuery Logo



Overview

By Region

Session Info

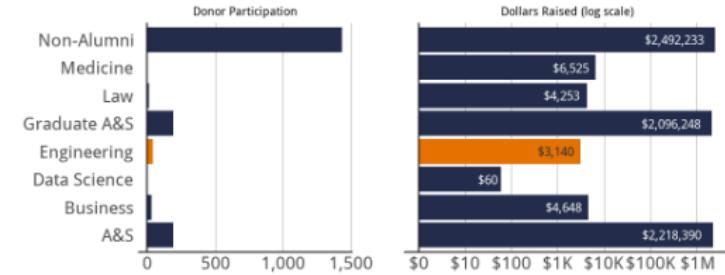
FY 2020 Giving Profile

School of Engineering Alumni

10 March 2020

Overview

43 Alumni from the School of Engineering gave in the 2020 Fiscal Year. This accounts for 2% of donor participation.



By Region

Of the 1,909 donors in FY 2020, 1,905 (100%) have a valid Zip Code.

418 entities reside outside of a CBSA. The table below shows the top 10 regions by count.

School of Engineering Alumni

Donor Participation in FY 2020

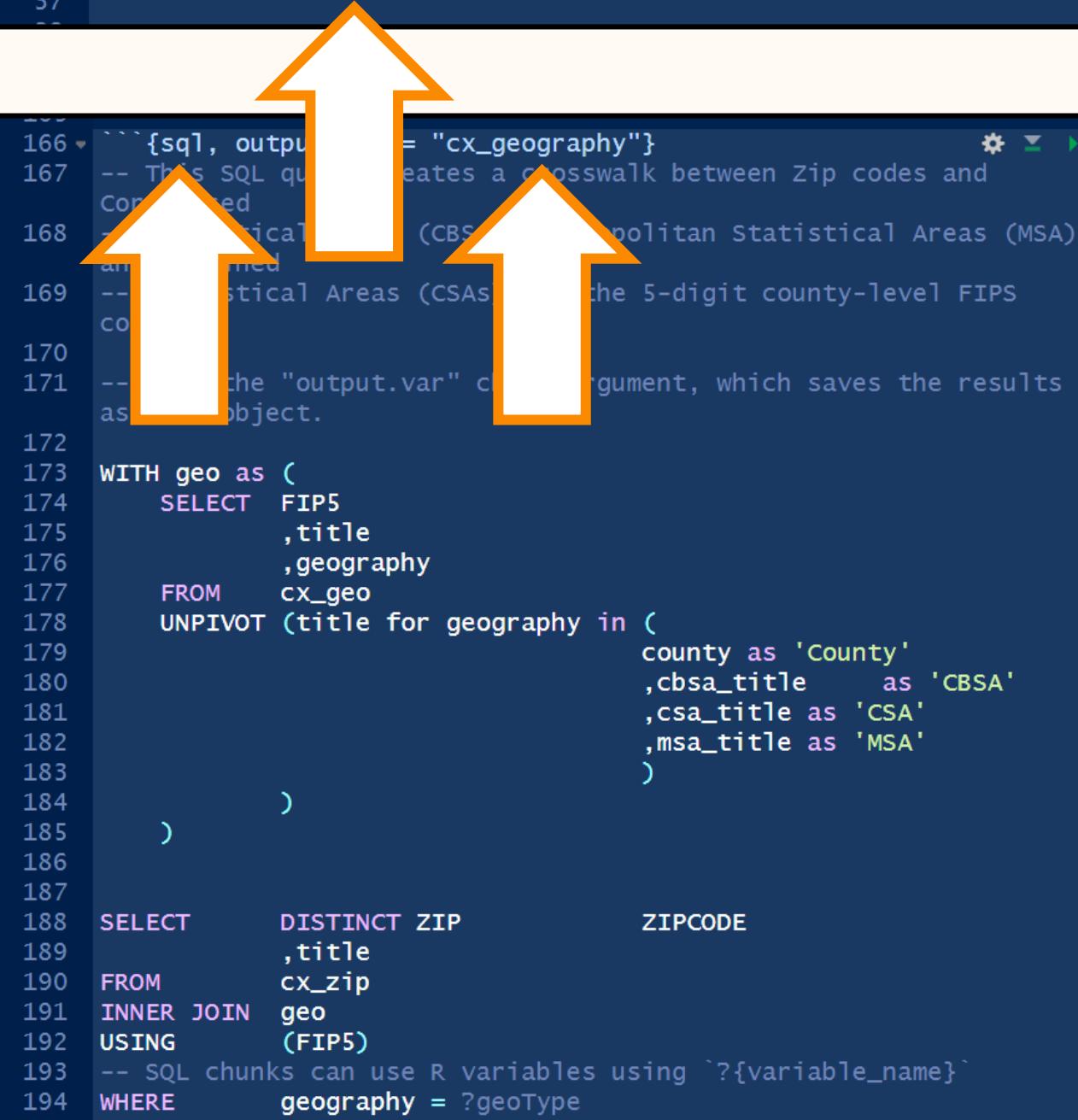
CBSA	Alumni	Total Donors	Alumni %
Portland-South Portland, ME	2	7	29%
Evanston, WY	1	1	100%
Heber, UT	1	1	100%
Ithaca, NY	1	1	100%
Monroe, MI	1	1	100%
Allentown-Bethlehem-Easton, PA-NJ	1	2	50%
Concord, NH	1	2	50%

CHALLENGE N° 3

THE SHOWSTOPPER

- Custom CSS
 - JQuery Logo
- SQL Chunks

```
35 con <- advancementtools::cx_Oracle(usr = "BI",
36                                         pw =
37                                         key_get("advance","BI","adv"))
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
```



CHALLENGE N° 3

THE SHOWSTOPPER

- Custom CSS
 - JQuery Logo
- SQL Chunks
 - Connection

```
35 con <- advancementtools::cx_Oracle(usr = "BI",
36                                         pw =
37                                         key_get("advance","BI","adv"))
```

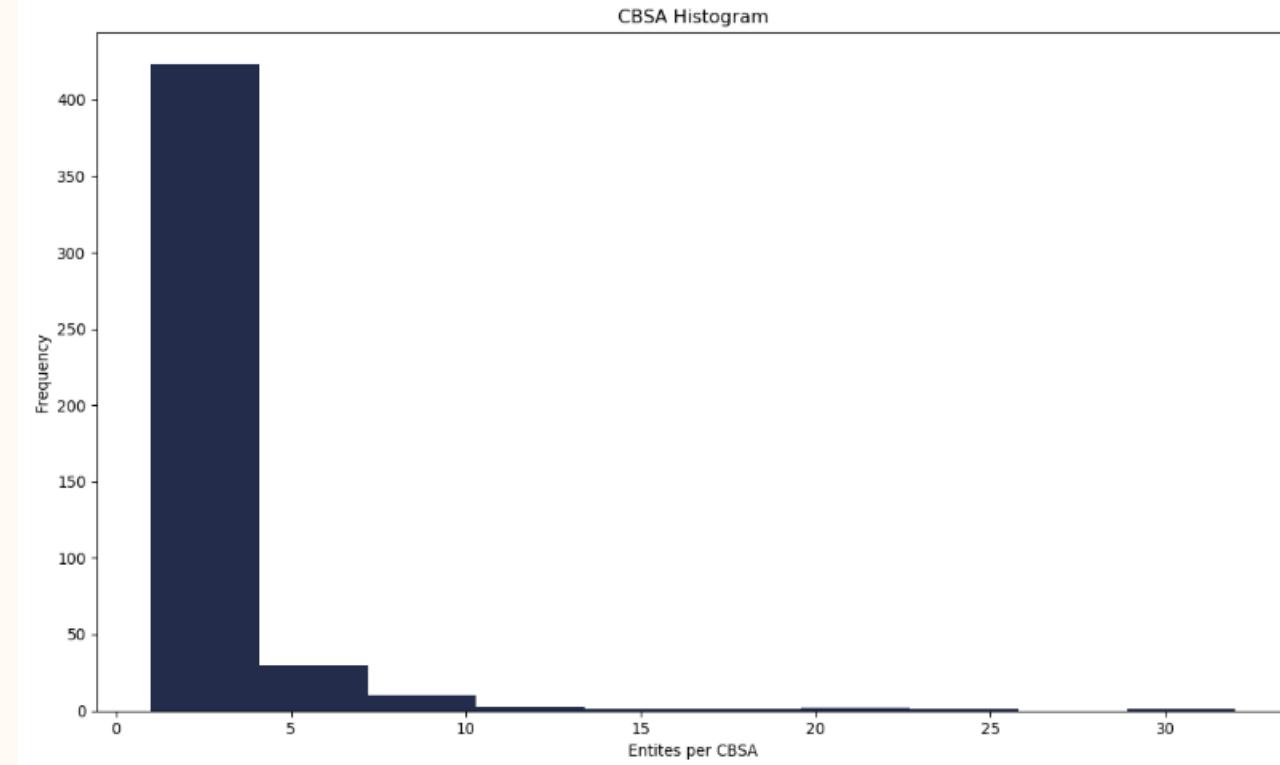
```
40 knitr::opts_chunk$set(echo = FALSE,
41                         message = F,
42                         warning = F,
43                         fig.showtext = TRUE,
44                         # dpi = 72,
45                         # The SQL connection needs to be passed to
46                         any SQL chunks,
47                         # and can be set here.
48                         connection = "con")  
```
```

```
173 WITH geo as (
174 SELECT FIP5
175 ,title
176 ,geography
177 FROM cx_geo
178 UNPIVOT (title for geography in (
179 county as 'County'
180 ,cbsa_title as 'CBSA'
181 ,csa_title as 'CSA'
182 ,msa_title as 'MSA'
183)
184)
185)
186
187
188 SELECT DISTINCT ZIP ZIPCODE
189 ,title
190 FROM cx_zip
191 INNER JOIN geo
192 USING (FIP5)
193 -- SQL chunks can use R variables using `?{variable_name}`
194 WHERE geography = ?geoType
```

CHALLENGE N° 3

# THE SHOWSTOPPER

- Custom CSS
  - JQuery Logo
- SQL Chunks
  - Connection



# CHALLENGE N° 3

# THE SHOWSTOPPER

- Custom CSS

- JQuery Logo

- SQL Chunks

- Connection

- Python Chunks

- `reticulate`

```
246 ````{python geoHist}
247 import pandas as pd
248 import numpy as np
249
250 np.random.seed(2020)
251
252 titles = r.combo['TITLE'].value_counts()
253
254 # Remove NA values
255 t2 = titles.drop('NA')
256
257 p1 = t2.plot.hist(figsize = (14,8),
258 title = "%s Histogram" % (r.geoType),
259 color = '#232D4B')
260
261 p1.set_xlabel("Entites per %s" % (r.geoType))
262
263 ````{r
264
265 `r scales::comma(py$titles["NA"])` individuals do not
266 `r` valid `r geoType`.
267
268 ````{r geoTable}
269 geoTable <- table(py$titles)
270
271 maxRegions <- names(geoTable[geoTable == max(geoTable)])
272
273 ````
```

CHALLENGE N° 3  
THE SHOWSTOPPER

9  
*schools*

6  
*years*

54  
*reports*

CHALLENGE N° 3

# THE SHOWSTOPPER

```
1 library(dplyr)
2
3 schools <- readr::read_csv(here::here("Data/donorBio.csv")) %>%
4 # Only Alumni have a SCHOOL, replace missing values
5 distinct(SCHOOL = tidyr::replace_na(SCHOOL, "Non-Alumni")) %>%
6 pull()
7
8 years <- 2015:2020
9
10
11 grid <- expand.grid(schools, years)
12
13 start <- Sys.time()
14 purrr::pwalk(grid,
15 ~rmarkdown::render(here::here("Markdown/3-TheShowstopper.Rmd"),
16 output_file = paste0(stringr::str_remove_all(..1, "[[:space:]]|[[:punct:]]"),
17 "FY", ..2, ".html"),
18 output_dir = here::here("Reports>Showstopper"),
19 params = list(school = ..1,
20 year = ..2,
21 geoType = "MSA")))
```

CHALLENGE N° 3

# THE SHOWSTOPPER

## FY 2015 Giving Profile

School of Business Alumni

### Overview

34 Alumni from the S  
accounts for 2% of do



### By Region

Of the 1,964 donors in  
1,437 entities reside ou  
by count.

## FY 2016 Giving Profile

School of Nursing Alumni

### Overview

2 Alumni from the Sch  
for 0% of donor particip



### By Region

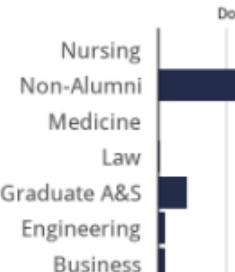
Reports>Showstopper

## FY 2017 Giving Profile

School of Medicine Alumni

### Overview

4 Alumni from the Sch  
for 0% of donor particip



## FY 2018 Giving Profile

School of Arts and Sciences  
Alumni

06 March 2020

### Overview

267 Alumni from the School of Arts and Sciences gave in the 2018 Fiscal Year.  
This accounts for 11% of donor participation.



CHALLENGE N° 3

# THE SHOWSTOPPER



# CONTACT



[rogol@virginia.edu](mailto:rogol@virginia.edu)



[/in/jrogol](https://www.linkedin.com/in/jrogol)



[@jrogol](https://twitter.com/jrogol)



[@jrogol](https://github.com/jrogol)

# **RESOURCES**

- Yihui Xie, J. J. Allaire and Garrett Grolemund. *R Markdown: The Definitive Guide.*
- Yihui Xie. “R Markdown Cookbook.” 27 February 2020.
- R Markdown Cheat Sheet. RStudio.
- R Markdown. RStudio.
- Emily Riederer. “R Markdown Driven Development: the Technical Appendix.” 1 February 2020.
- Hadley Wickham. *R for Data Science*. “Chapter 27: R Markdown.”
- Paul Hively “Reproducible Data Science.” DRIVE/ 2019.

# REFERENCES

Great British Baking Show Logo from [Twin Cities PBS](#).

Great British Bake Off tent image from the GBBO [twitter](#).

Great British Bake Off Series 11 Cast from the [greatbritishbakeoff.co.uk](#)

[“The Great British Bake Off.” Wikipedia](#).

Jenny Bryan. [“Ode to the here package.”](#)

Kouign Amann recipe and picture by [David Lebovitz](#).

Raspberry Blancmange recipe and picture from [Prue Leith](#).

Icons made by [Smashicons](#), [Gregor Cresnar](#) and [Pixel Perfect](#) from [www.flaticon.com](#), licensed [CC 3.0 BY](#).

Data adapted from Kaggle’s [“Fundraising Data”](#), curated by Michael Pawlus.

Ashutosh Nandeshwar and Rodger Devine. *Data Science for Fundraising*.

Emily Robinson and Jacqueline Nolis. [Build a Career in Data Science](#).

Package Hex stickers from [rmarkdown](#), [bookdown](#), [blogdown](#), [pkgdown](#), [pagedown](#) and [xaringan](#).

# DRIVE/Cast

*Accelerate Smarter*