HONOR the FUTURE

THE CAMPAIGN FOR THE UNIVERSITY OF VIRGINIA



Advancement

READY, SET, BAKE

Recipes for Reproducible Reporting with RMarkdown



James Rogol | CASE Drive/ | 25 March 2020





re·pro·duc·i·bil·i·ty, noun

Repeatable

Identical

Shareable



FOLLOW ALONG:









1-TheSignature.pdf

Contact Report Summary

Month of August

Major Gifts

05 March 2020

Reports filed in August, by Type.

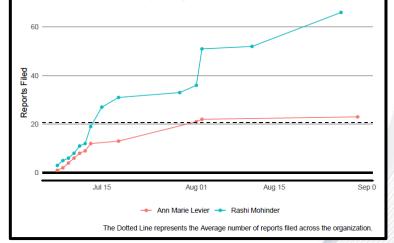
	Email	Phone	Visit	Total
Ann Marie Levier	3	3	4	10
Rashi Mohinder	10	5	18	33
Total	13	8	22	43

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)



Contact Report Summary

Month of August

Major Gifts

05 March 2020

Markdown/1-TheSignature.rmd

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

• Use `here`

R code Chunk

```
```{r chunkName}
```

Sections

In-line code

`r params\$endMonth`

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

170

173

174

```{r render, eval=F}

Rendering

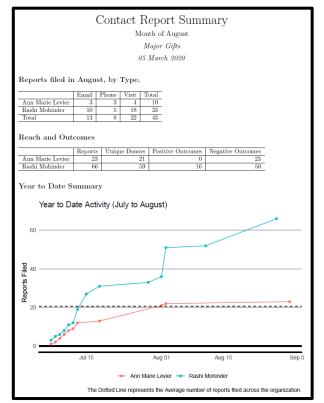
Parameters

```
175
     rmarkdown::render(here("Markdown/1-TheSignature.Rmd"),
                       output_dir = here("Reports"),
178
                       output_file = "1-TheSignature-AG_Sept",
                       params = list(endMonth = "September",
179
                                     dept = "Planned Giving"))
         inline R code fragments
       10%
      label: setup (with options)
      List of 1
       $ include: logi FALSE
       14%
        ordinary text without R code
      1 19%
      label: getData
       24%
        ordinary text without R code
       29%
      label: cleanData1 (with options)
     List of 1
```

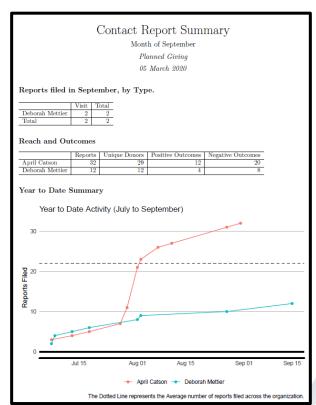
rmarkdown::render(here("Markdown/1-TheSignature.Rmd"),

output_dir = here("Reports"))

THE SIGNATURE



1-TheSignature.pdf



1-TheSignature-AG_sept.pdf







Background

Approach

.

Sampling

EDA

Analysis

Session Info

```
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.md

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 was originally obtained here, and modified.¹

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

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".

Factor Variables

Of the 19 text columns above, <code>age_bin</code>, <code>gender</code>, and <code>address_type</code> are more like factors, and should be treated as such. <code>age_bin</code> 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))
```

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

```
First, we'll need to obtain the data for the analysis.
    This anonymized dataset was was originally obtained
    [here](https://www.kaggle.com/michaelpawlus/fundraising-
    ata), and modified. [^modify]
    ```{r getData}
 basicReport <- read_csv(here("Data/basicReport.csv"))</pre>
 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.
47 ### Data Wrangling
49 - #### Age
 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".
 {r cleanAge}
 basicReport <- basicReport %>%
 mutate(age = if_else(age == 0, NA_real_,age),
 age_bin = if_else(is.na(age),"Unknown",age_bin);
```

36 + ## Data

- 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,]</pre>
```

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)
}</pre>
```

- 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 bulletted list.
Get a list of attached packages
packages <- sessioninfo::package info()</pre>
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 hypen and space at the start of th
e string
(to create a bulletted list) and a new line character at the end of th
e 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)
- tidyr (0.8.3)

# THE SHOWSTOPPER





# Overview By Region Session Info

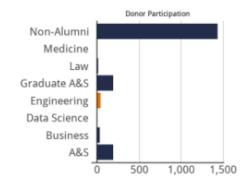
### FY 2020 Giving Profile

#### 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.





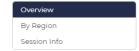
### By Region

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.







### Custom CSS

JQuery Logo

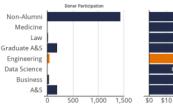
#### FY 2020 Giving Profile

#### 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

Concord NH

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

| Donor r arciolpaci                    | 011111111111 |                 |             |
|---------------------------------------|--------------|-----------------|-------------|
| CBSA                                  | Alumni       | Total<br>Donors | Alumni<br>% |
| 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-<br>NJ | 1            | 2               | 50%         |
|                                       |              |                 |             |

### THE SHOWSTOPPER

- Custom CSS
  - JQuery Logo
- SQL Chunks
  - Connection

```
key_get("advance", "BI", "adv"))
166 \cdot\```{sql, output.var = "cx_geography"}
167 -- This SQL query creates a crosswalk between Zip codes and
 Core-Based
168 -- Statistical Areas (CBSA), Metropolitan Statistical Areas (MSA)
169 -- Statistical Areas (CSAs) via the 5-digit county-level FIPS
170
 as an R object.
173
 WITH geo as (
174
 SELECT FIP5
175
 ,title
 , geography
 FROM
 cx_geo
 UNPIVOT (title for geography in (
179
 county as 'County'
 ,cbsa_title
 as 'CBSA'
 .csa_title as 'CSA'
 ,msa_title as 'MSA'
184
 SELECT
 DISTINCT ZIP
 ZIPCODE
 ,title
 FROM
 cx_zip
 INNER JOIN
 geo
 USING
 (FIP5)
 -- SQL chunks can use R variables using `?{variable_name}
 geography = ?geoType
 WHERE
```

con <- advancementtools::cx\_Oracle(usr = "BI",</pre>

pw =

### THE SHOWSTOPPER

- Custom CSS
  - JQuery Logo
- SQL Chunks
  - Connection

connection = "con")

pw =

con <- advancementtools::cx\_Oracle(usr = "BI",</pre>

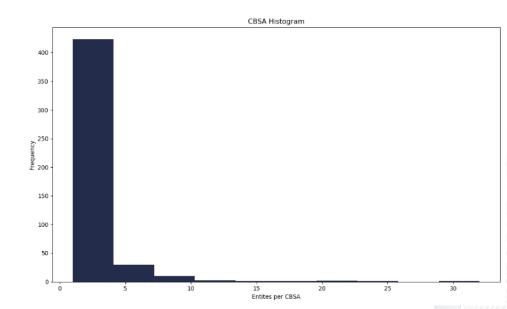
```
WITH geo as (
174
 SELECT FIP5
 ,title
 , geography
 FROM
 cx_geo
 UNPIVOT (title for geography in (
179
 county as 'County'
 ,cbsa_title
 as 'CBSA'
 ,csa_title as 'CSA'
 ,msa_title as 'MSA'
184
 SELECT
 DISTINCT ZIP
 ZIPCODE
 .title
 FROM
 cx_zip
 INNER JOIN
 geo
 USING
 (FIP5)
```

geography = ?geoType

WHERE

# CHALLENGE N° 3 THE SHOWSTOPPER

- Custom CSS
  - JQuery Logo
- SQL Chunks
  - Connection



### THE SHOWSTOPPER

- Custom CSS
  - JQuery Logo
- SQL Chunks
  - Connection
- Python Chunks 268 268 269
  - Reticulate

```
{python geoHist}
 import pandas as pd
 import numpy as np
249
 np.random.seed(2020)
 titles = r.combo['TITLE'].value_counts()
253
 t2 = titles.drop('NA')
256
 pl = t2.plot.hist(figsize = (14,8),
258
 title = "%s Histogram" % (r.geoType),
259
 color = '#232D4B')
 pl.set_xlabel("Entites per %s" % (r.geoType))
262
 # pl.show()
264
265
 `r scales::comma(py$titles["NA"])` individuals do __not__
 valid `r geoType`.
 {r geoTable}
 geoTable <- table(py$titles)</pre>
270
271
 maxRegions <- names(geoTable[geoTable == max(geoTable)])</pre>
```

# CHALLENGE N° 3 THE SHOWSTOPPER

Schools

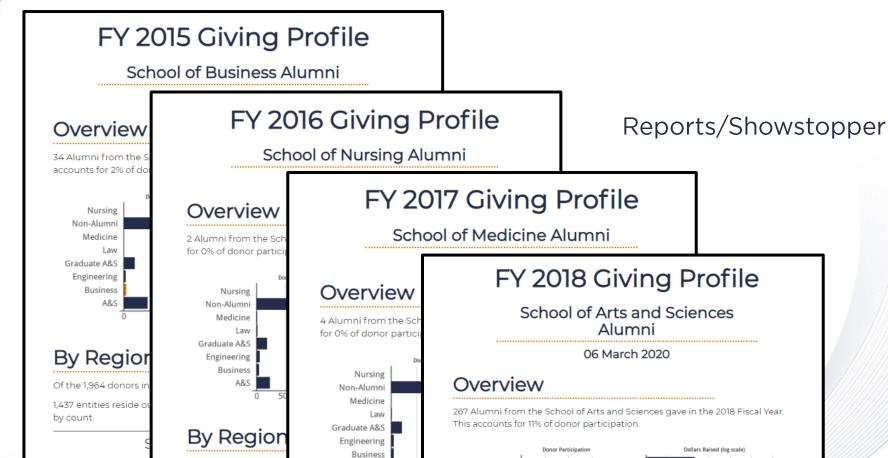
**G**years

54 reports

### THE SHOWSTOPPER

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

### THE SHOWSTOPPER



# CHALLENGE N° 3 THE SHOWSTOPPER



### CONTACT

- rogol@virginia.edu
  in /in/jrogol
  /jrogol
- 7 /jrogol

### RESOURCES

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

### REFERENCES

Great British Baking Show Logo from Twin Cities PBS.

"The Great British Bake Off." Wikipedia.

Jenny Bryan. "Ode to the here package."

Kouign Amann recipe and picture by <u>David Lebovitz</u>.

Raspberry Blancmange recipe and picture from Prue Leith.

Icons made by <u>Smashicons</u>, <u>Gregor Cresnar</u> and <u>Pixel Perfect</u> from <u>www.flaticon.com</u>, licensed <u>CC 3.0 BY</u>.

Data adapted from Kaggle's <u>"Fundraising Data"</u>, curated by Michael Pawlus.

Ashutosh Nandeshwar and Rodger Devine. Data Science for Fundraising.

Package Hex stickers from <u>rmarkdown</u>, <u>bookdown</u>, <u>blogdown</u>, <u>pkgdown</u>, <u>pagedown</u> and <u>xaringan</u>.