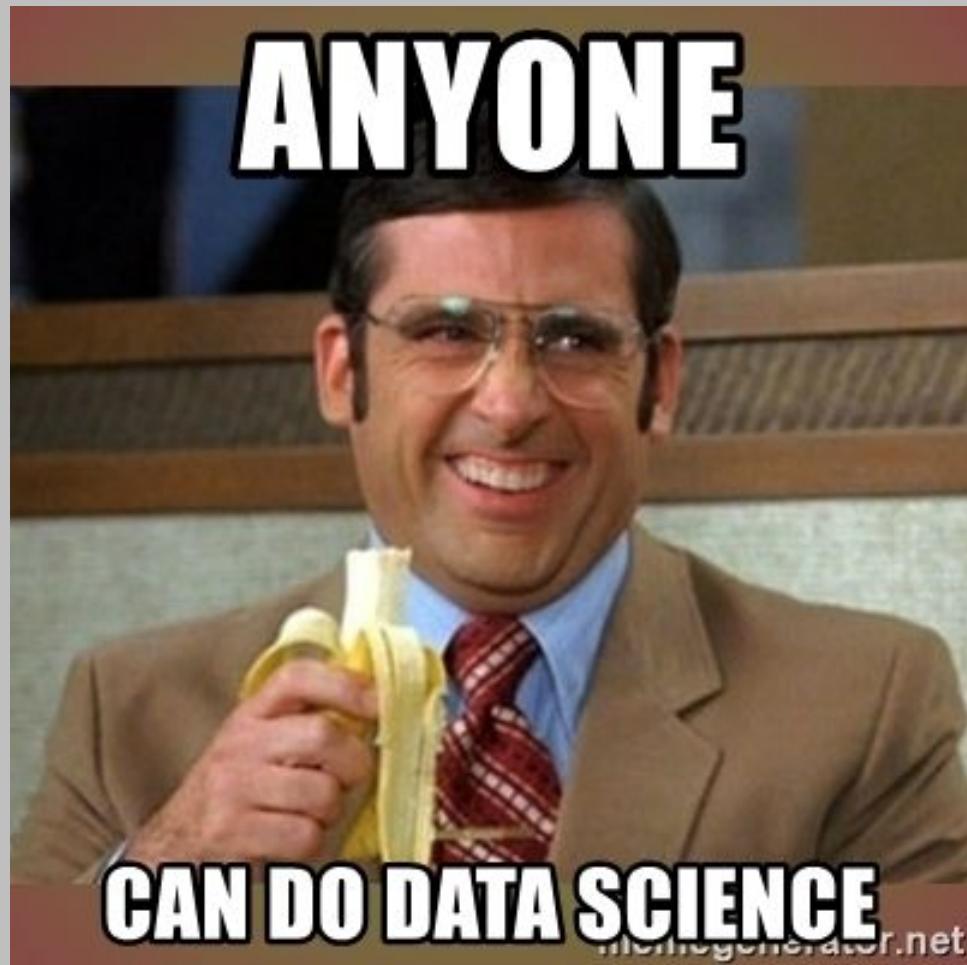


Big Data Analyst: Cohort 1 Graduation

The Struggle Was Real

Welcome To
Our Journey



The Meta Data

480
class hours

9
sprints

911
GitHub commits

1
machine learning
competition

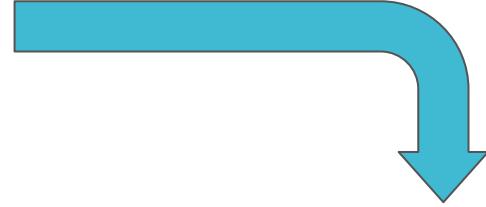
182
pounds of candy
consumed (we estimate)

n/a
days the class
temp was above
67 degrees

Once upon a time before data class...

Runjini manually copied and pasted data into Excel.

All Posts Published						
Published	Post	Type	Targeting	Reach	Engagement	Promote
01/10/2018 6:00 pm	Bright, bold, & unforgettable. For more than five decades			238	0 0 0	<button>Boost Post</button>
01/10/2018 10:05 am	Oh what a night! We had an amazing time at the National			1.3K	29 3 1	<button>Boost Post</button>
01/09/2018 12:00 pm	Names have power. Learn more about the origins of the			13.7K	210 1 14	<button>View Promotion</button>
01/08/2018 6:00 pm	By scouring the seafloor in search of materials for her			616	4 0 0	<button>Boost Post</button>
01/05/2018 12:00 pm	Looking to craft a brighter 2018 for yourself but not sure			803	7 0 0	<button>Boost Post</button>
01/02/2018 5:00 pm	Shaken, not stirred: Bar Leather Apron provides an			28.3K	489 5 25	<button>View Promotion</button>



1/10/18				867	20
	10:05 AM			Organic	20
				867	3
				Paid	3
				0	1
				1	Boost Post
1/9/18				9.9K	161
	12:00 PM			Organic	161
				1,019	1
				Paid	1
				8,855	9
				9	View Promot
1/8/18				608	4
	6:00 PM			Organic	4
				608	0
				Paid	0
				0	0
				0	Boost Post

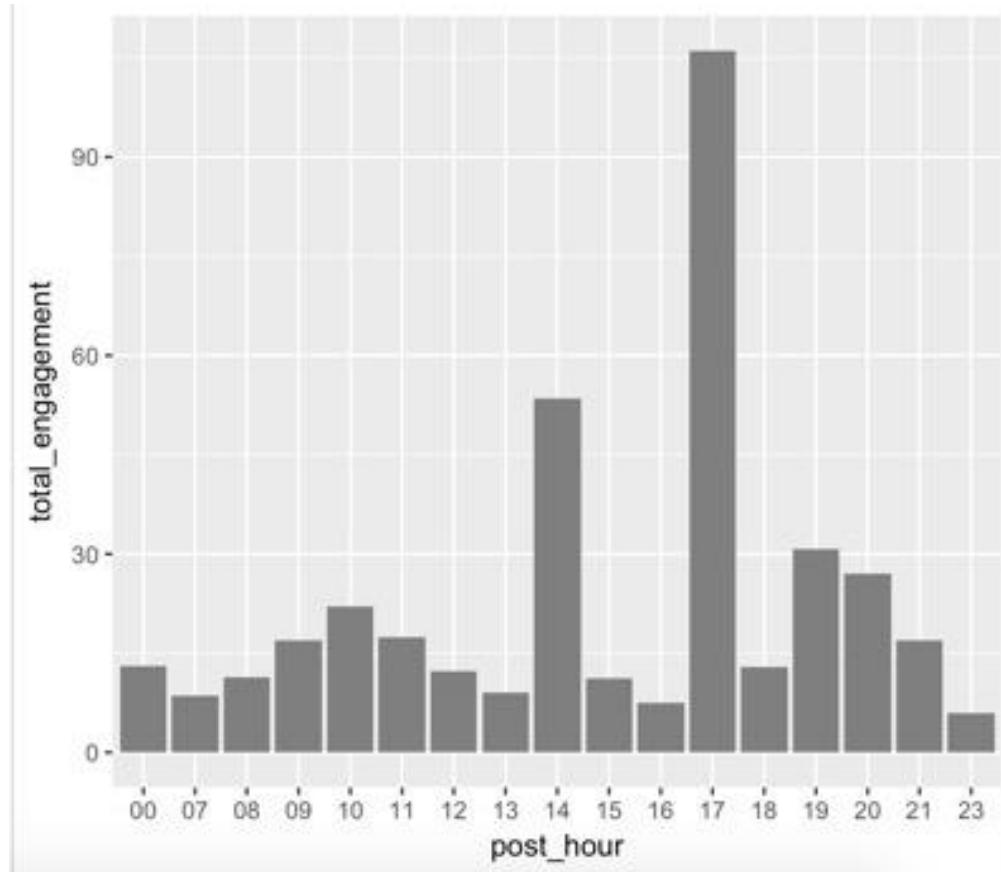
Then we learned about R packages!

19 lines of code replaced tedious copying, pasting, and reformatting.

```
1 install.packages('Rfacebook')
2 library("Rfacebook")
3 myoauth <- fbOAuth(app_id='app id goes here',app_secret='app secret goes here')
4 save(mymyauth, file="myoauth")
5 load("myoauth")
6 me <- getUsers("me", token = myoauth)
7 getpagedata <- getPage (page_number_from_website, token = myoauth, n = 250)
8 write.csv(getpagedata, "/Users/runjini.murthy/Desktop/GitHub/Runjini_Sprint_4/last-250-wv-fb-posts-Rfacebook.csv")
9 WVFBData <- read.csv(file="/Users/runjini.murthy/Desktop/GitHub/Runjini_Sprint_4/last-250-wv-fb-posts-Rfacebook.csv",
10 header=TRUE, sep=",")
11 WVFBData$total_engagement<-WVFBData$likes_count + WVFBData$comments_count + WVFBData$shares_count
12
13 install.packages("parsedate")
14 library(parsedate)
15 WVFBData$parsed_timestamp <- parse_iso_8601(WVFBData$created_time)
16 WVFBData$adjusted_timestamp <- WVFBData$parsed_timestamp - 10*60*60
17 WVFBData$parsed_date <- format(WVFBData$adjusted_timestamp, "%m/%d/%Y")
18 WVFBData$parsed_time <- format(WVFBData$adjusted_timestamp, "%H:%M")
19 WVFBData$post_hour <- format(WVFBData$adjusted_timestamp, "%H")
20 write.csv(WVFBData, "/Users/runjini.murthy/Desktop/GitHub/Runjini_Sprint_4/last-250-wv-fb-posts-Rfacebook.csv")
21 ggplot(WVFBData) + geom_bar(aes(x=post_hour,y=total_engagement),stat="summary", fun.y = "mean",fill=I("grey50"))
```

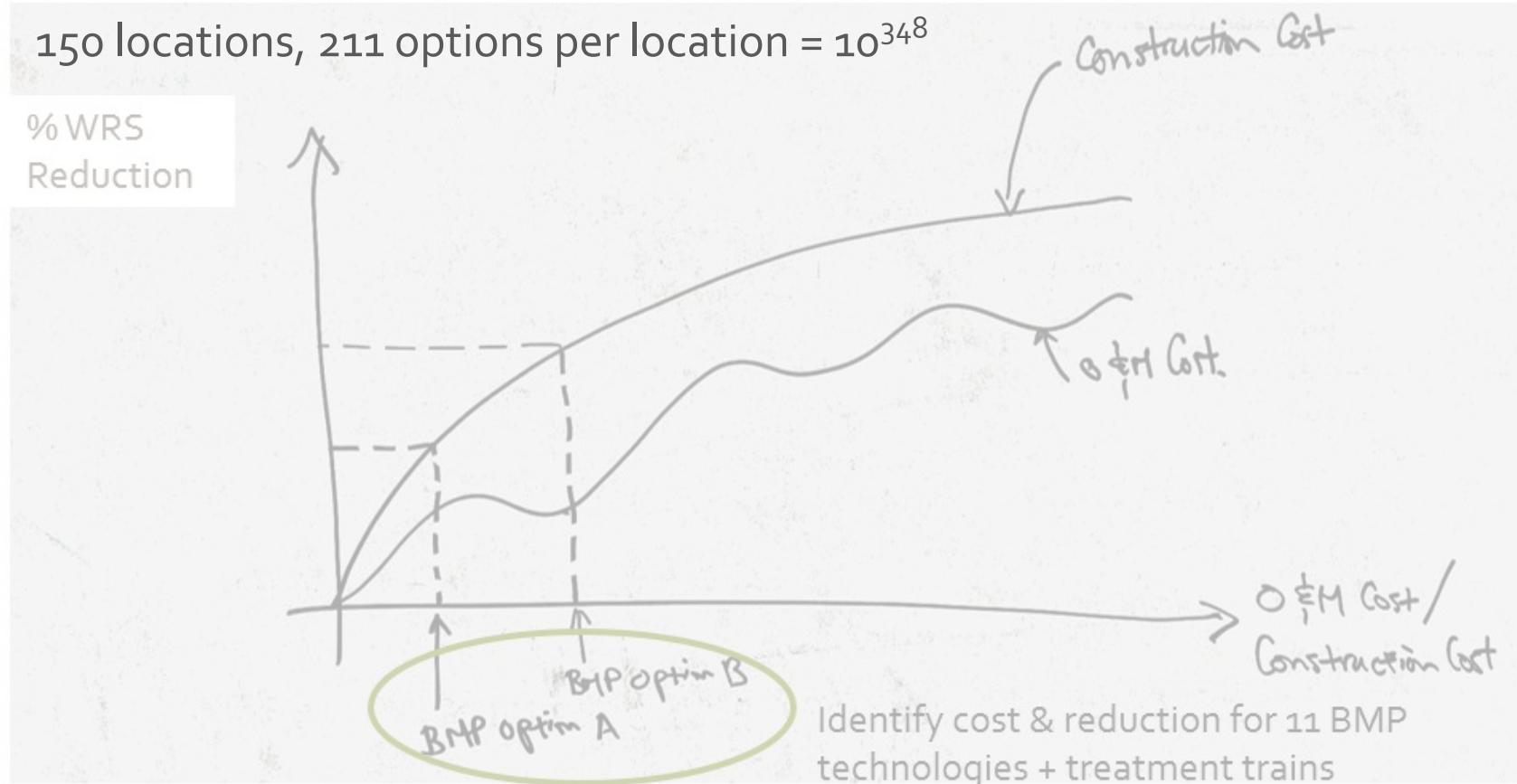
And lived
happily data
after.

Now it is easier to perform business analysis, like a plot of times to show optimal post hours to drive Facebook engagement.

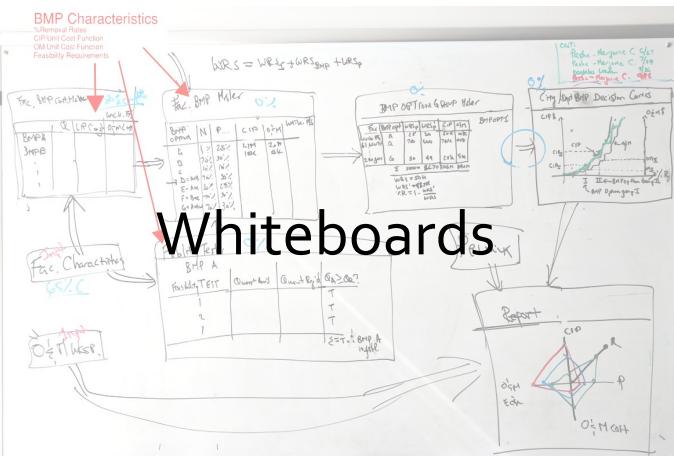


Once upon a time before data class...

Jon asked: How do you analyze 10^{348} options?

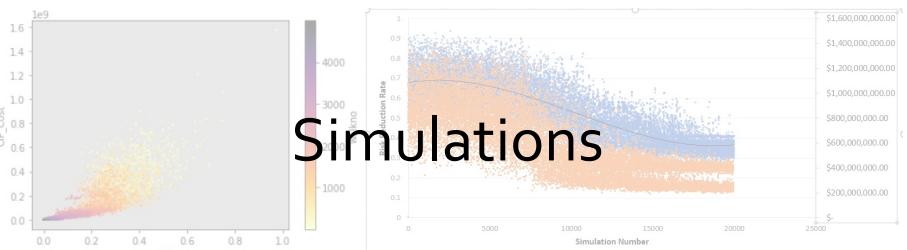


Along the way we met...



Facility_ID	id	January	February	March	April	Ma
1	81	3.39	3.85	4.05	2.87	1.8
2	171	4.16	3.02	2.96	1.25	1.6
3	65	4.96	4.13	5.50	2.98	2.5
4	91	5.71	5.16	6.07	4.72	3.4
5	97	5.28	4.98	5.43	3.63	2.8
6	103	3.66	2.01	1.93	1.26	1.0
7	112	3.11	3.53	3.67	2.19	1.3
8	121	5.25	4.90	5.96	4.11	3.3

Pandas



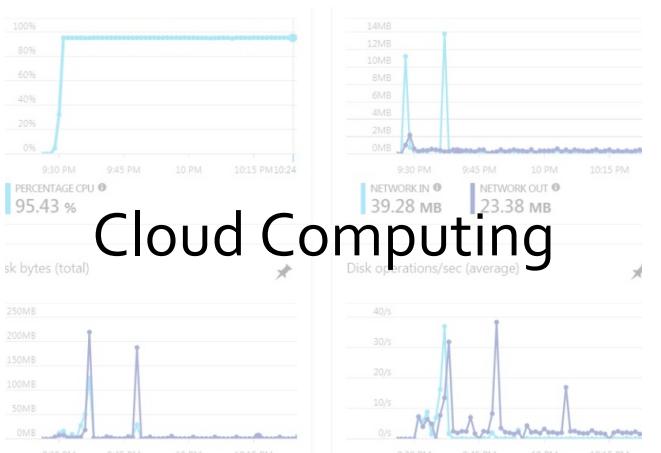
Link Samples By Date & Write to pd_exMaxConcs

is when we calculate Age Factor Weighted Avg in Step 4)

i sampled parameter by sample date. The latest sample date is assigned 0. Earlier s into new columns suffixed by 'SR'.

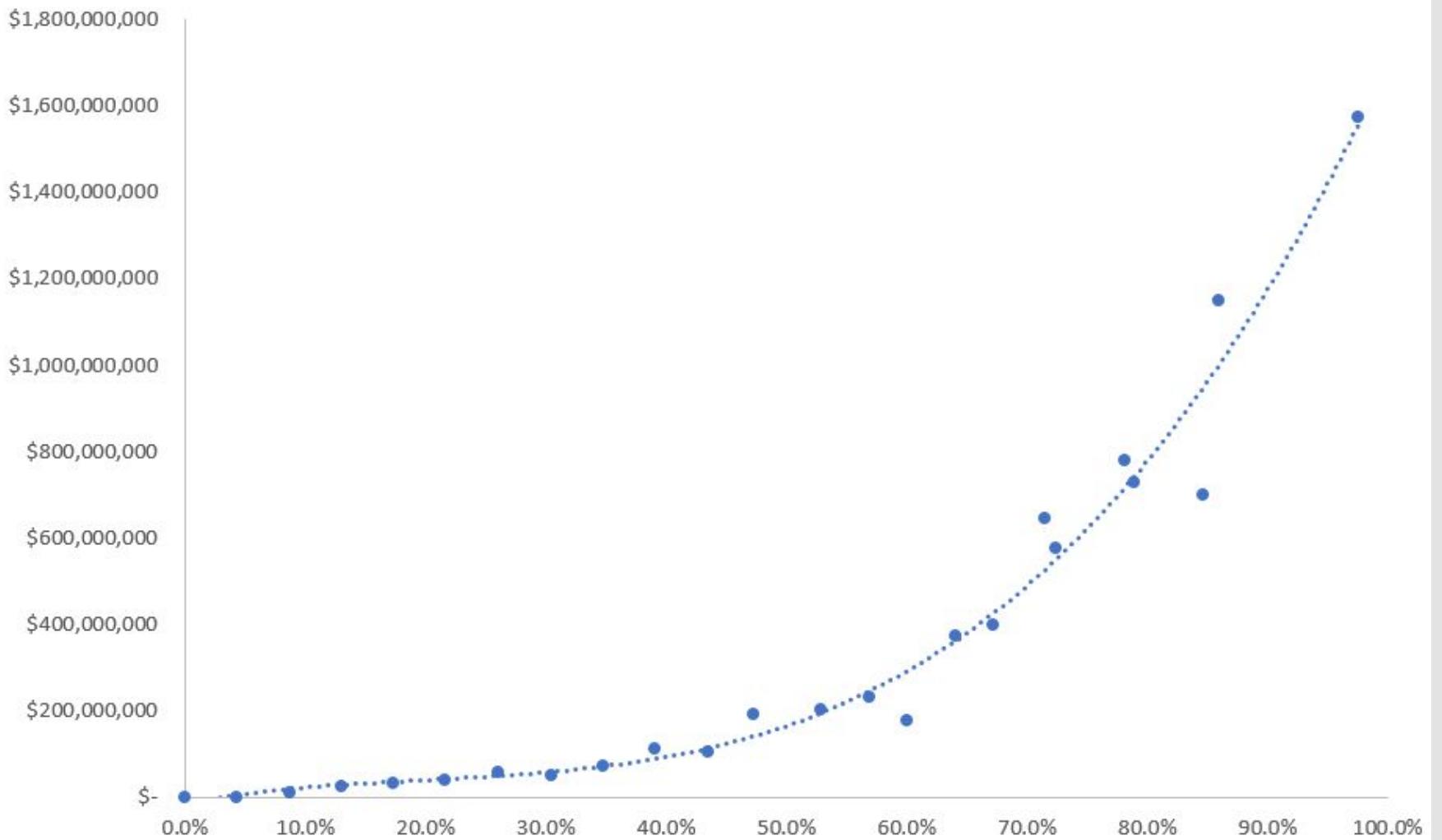
```
R_SampleRank(pd_Concs):
    type(datetime)
    int(str(datetime)[:10].replace('-', '')) #return a numeric va
    sampleRank(pd_Concs, pollIS):
        Constituent in pollIS:
            like helper column that expresses date as numeric:
            pd_Concs['c_' + Constituent + '_HelpSR'] = pd_Concs.apply(
                lambda row: _HELPER_SampleRank(row['sample_date']) if not (
                    rank sample dates for each constituent of each facility
```

Cloud Computing



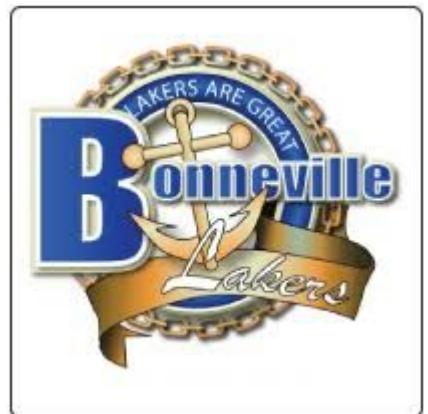
And lived
happily data
after.

Pollutant Removal Rate vs CIP Cost



Once upon a time before data class...

For Tori, it all started in the 9th grade with Ms. Carlston...



Fast forward to Tori's sophomore year of college where she takes her first HTML class.



Fast forward to January 2016. Girl meets boy. Girl meets Devleague.



The love-hate relationship continues to grow...

And lived
happily data
after.

```
Victorias-MacBook:~ victorialarson$ ls
Applications          Library
Desktop               Movies
Documents             Music
Downloads            Pictures
Public
US_diseases.R
USdisease_questions.R
anaconda
```



```
#getting page 7 from the file
get_file_from_date <- function(start_date) {
  filename <- paste(start_date, ".pdf", sep = "")
  file_text <- pdf_text(filename)
  # Want to focus on page 3 of the vector in the list.
  page7 <- file_text[[7]]
  #Using regexpr to find the location of RevPAR\n
  tindex <- regexpr("Running\nTop 25 Markets", page7)
  #deleting "Running\nTop 25 Markets" and everything before it using substr
  page7 <- substring(page7,tindex+23)
  # String split with "\n" - creates new lines
  newline <- strsplit(page7, "\n")
  #using substring turned my page into a list, I need it to be a character
  newchar <- unlist(newline)
  #taking out all of the rows I dont need.
  newchar <- newchar[-(1)]
  newchar <- newchar[-(26:28)]
  # Extracted all of the floats from the vectors
  extract <- (str_extract_all(newchar, "[+-]?([0-9]*[.])?[0-9]+"))
  # turning the list back into a character
  un_list <- as.numeric(unlist(extract))
  # telling it there need to be 18 lines
  rep_list <- rep(1:18, times=length(un_list)/18)
  # I dont know whats going on here
  split_list <- split(un_list, rep_list)
  # making the columns
  df <- cbind.data.frame(split_list, stringsAsFactors=F)
  # taking out the columns that I dont want
  df <- df[-(8:18)]
  #adding column names
  row.names(df) <- c("Anaheim/SantaAna_CA", "Atlanta_GA", "Boston_MA", "Chicago_IL", "Dallas_TX", "Denver_CO", "Houston_TX", "Las_Vegas_NV", "Los_Angeles_CA", "Miami_FL", "Minneapolis_STP", "Nashville_TN", "New_York_NY", "Philadelphia_PA", "Phoenix_AZ", "Portland_OR", "Seattle_WA", "St_Louis_MO", "Tampa_Bay_FL", "Washington_DC")
```

A Data Offering for DevLeague

DATA PREP	Intro To Data and Mathematical Foundations
MODEL 1	Data Analytics, Data Collection & The Analytical Process
MODEL 2	Business Statistics & Probability Theory with Microsoft Excel
MODEL 3	Exploratory Data Analysis & Business Intelligence Reporting with Tableau
MODEL 4	Introduction to SQL Databases & Relational Data Modeling
MODEL 5	Introduction to Programming in R & R Studio

Oh wait,
there's more!

MODEL 6

Python Programming & Introduction to Big Data Analysis

MODEL 7

Introduction to Machine Learning & Predictive Modeling

MODEL 8

Advanced Data Visualization & Storytelling

MODEL 9

Data Governance & Security in The Enterprise

CAPSTONE
PROJECT

The Big Reveal

Where were we for the last 40 weeks?



Haiku: An Ode to Data Analysis

Google is your friend.
Wrangling takes most of the time.
One typo ruins all.

What's the deal with R vs. Python?

Excel

The screenshot shows a Microsoft Excel spreadsheet with four columns: Year, State, WeeksReporting, and Disease. The data table has two rows: Row 1 (header) and Row 6 (data). Row 6 contains values 1953, Utah, 52, and Smallpox respectively. To the right of the table is a filter dialog. The 'Equals' dropdown is set to 'Smallpox'. The 'And' radio button is selected. Below the dropdown is a 'Choose One' dropdown and a search bar. At the bottom of the dialog, there are checkboxes: '(Select All)', 'Measles' (unchecked), and 'Smallpox' (checked).

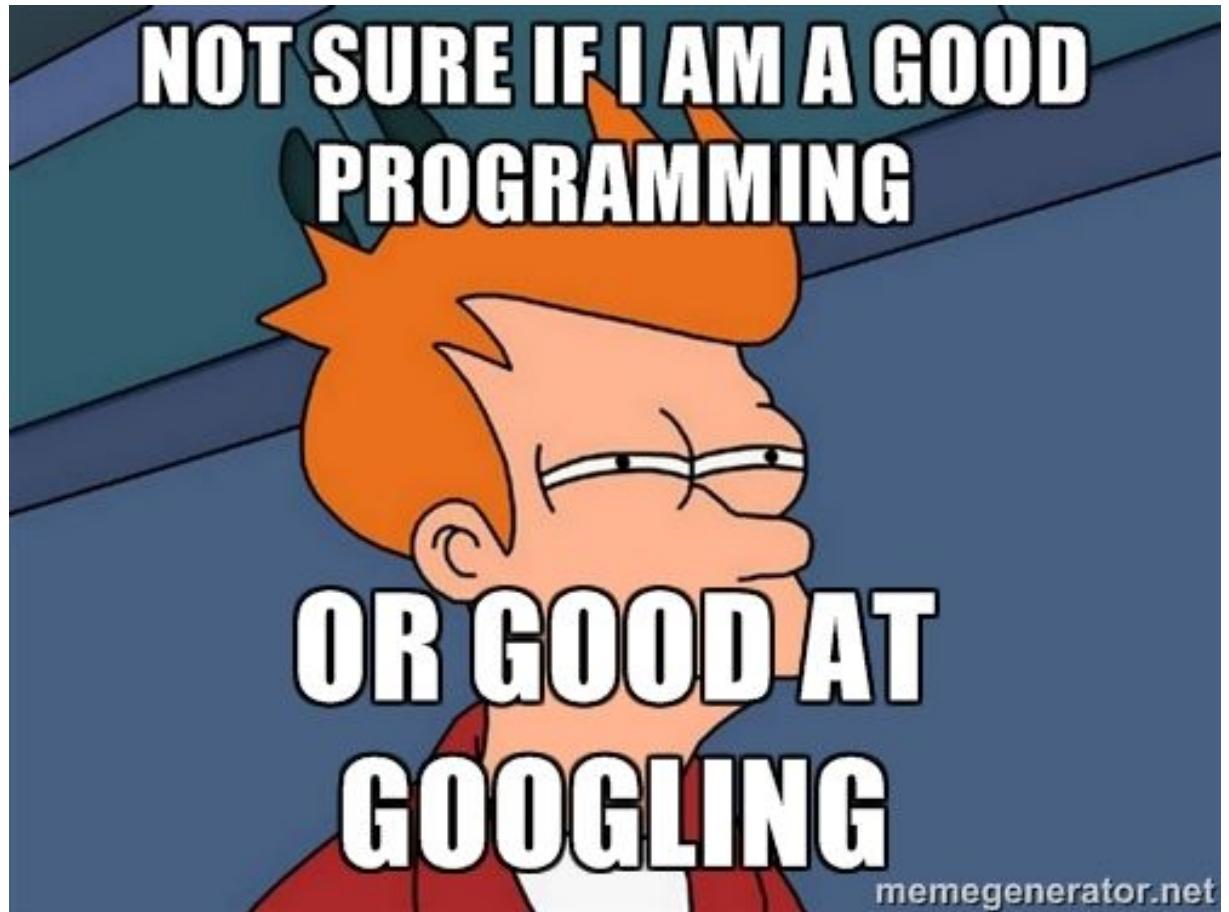
R

```
#filtering out just Small pox from the weeks reporting 52  
Smlpx52wks <- subset(Week52_reporting, disease=="Smallpox")
```

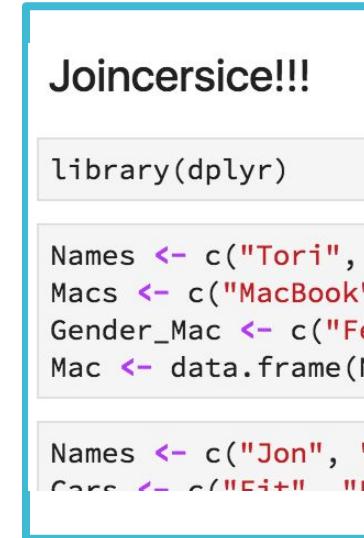
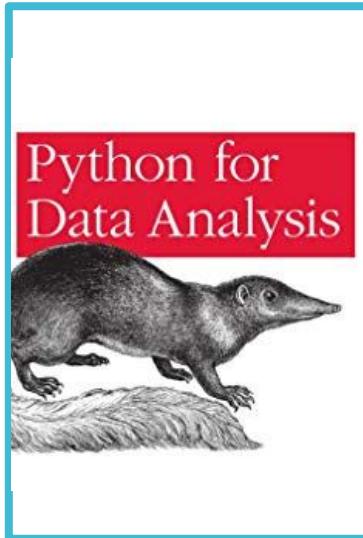
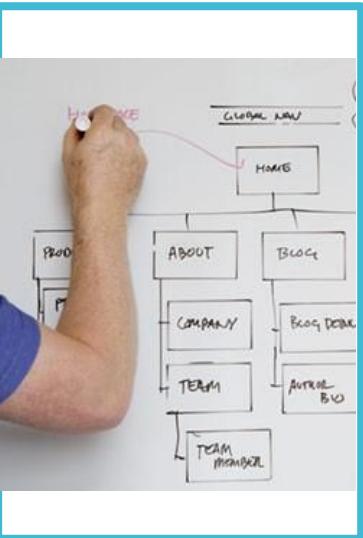
Python

```
Smlpx52wks = Week52_reporting.loc[Week52_reporting['disease'] == 'Smallpox']
```

How We Learned



Helpful Resources



Helpful Resources

Googling Like Ninjas

All Videos Images News Shopping More Settings Tools

About 223,000,000 results (0.71 seconds)

A Tutorial on Loops in R - Usage and Alternatives - DataCamp

Ad www.datacamp.com/ ▾

Learn data science with our free online and interactive tutorials. Courses: Python for Data Science, R Programming, Applied Finance, Data Manipulation, Data Visualization, Dplyr, Ggplot2, Machine Learning, Pytho...

A Tutorial on Loops in R - Usage and Alternatives (article) - DataCamp

<https://www.datacamp.com/community/tutorials/tutorial-on-loops-in-r> ▾

Sep 28, 2016 - A tutorial on **loops** in R that looks at the constructs available in R for ... and how to make use of alternatives, such as R's vectorization feature, ...

You've visited this page 2 times. Last visit: 8/29/18

How to write the first for loop in R | R-bloggers

<https://www.r-bloggers.com/how-to-write-the-first-for-loop-in-r> ▾

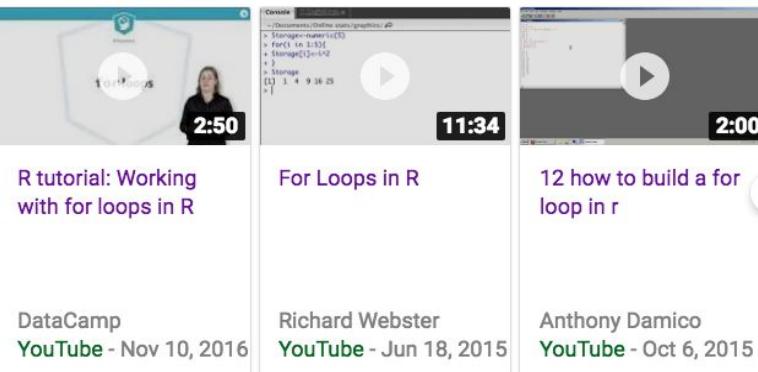
Dec 2, 2015 - Before you dive into writing **loops** in R, there is one important thing you ... For example, solutions that make use of loops are less efficient than ...

Programming with R: Loops in R - Our Lessons

<https://swcarpentry.github.io/r-novice-inflammation/15-supp-loops-in-depth/> ▾

How can I do the same thing multiple times more efficiently in R? ... Instead of using i in a to make our loop variable, we use the function seq_along to generate ...

Videos



R for Loop (With Examples)

<https://www.datamentor.io/r-programming/for-loop/> ▾

Loops are used in programming to repeat a specific block of code. In this article, you will learn to create a for loop in R programming.

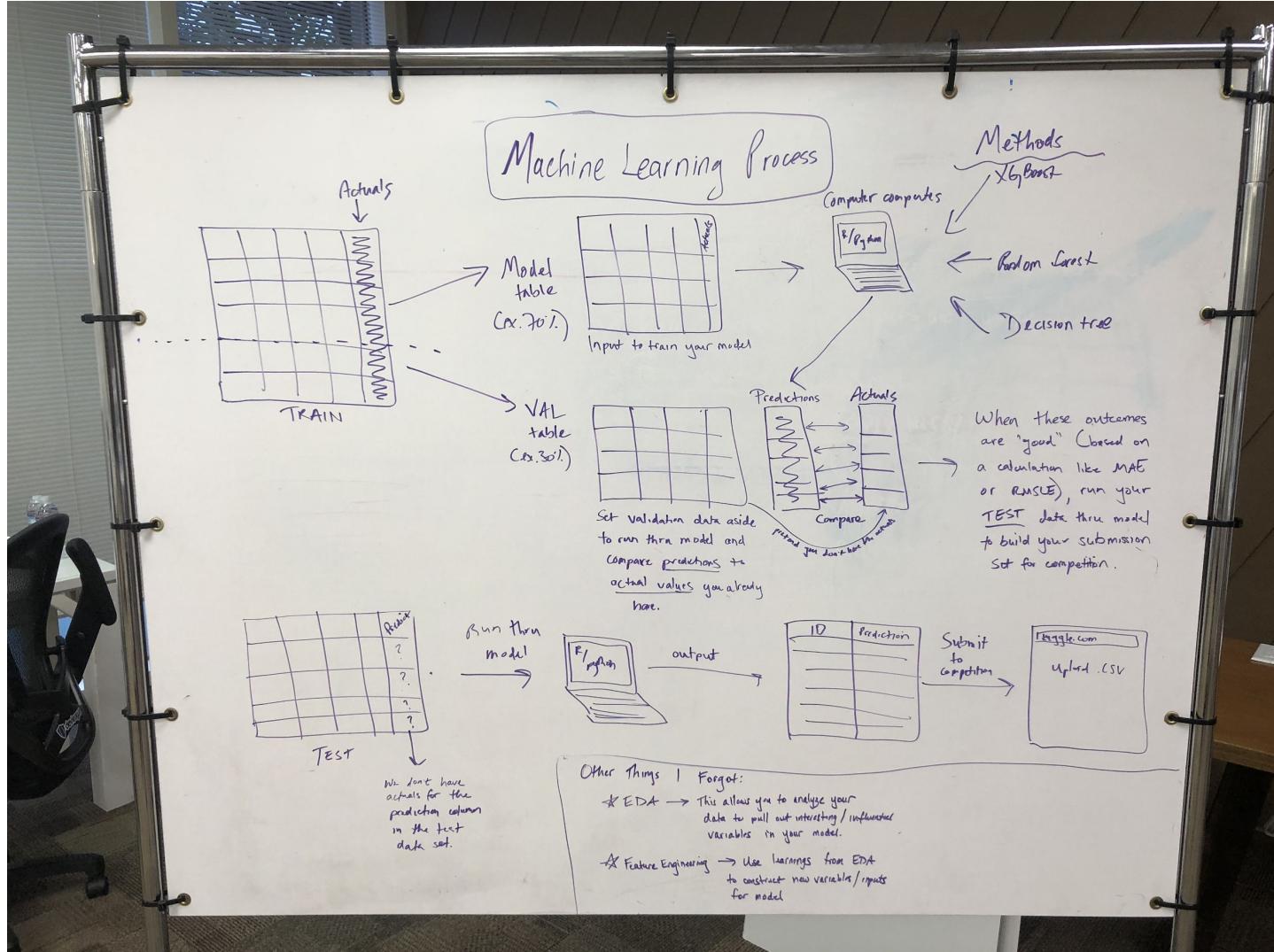
18.05 R Tutorial: 'For Loops'

<https://ocw.mit.edu/ans7870/18/18.05/s14/html/r-tut-forloop.html> ▾

R makes this easy with the replicate function rep() # rep(0, 10) makes a ...

Helpful Resources

Whiteboarding Like Champs



Helpful Resources

Reading Like It's Our Job

Marketing Dashboard ≡

f Facebook
G AdWords
Website Traffic

Demographics and Spend

Select Age and/or Gender

Age

18-24
 25-34
 35-44
 45-54

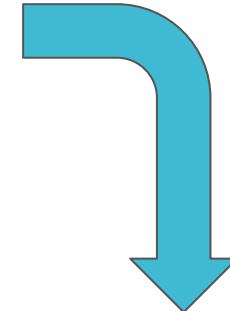
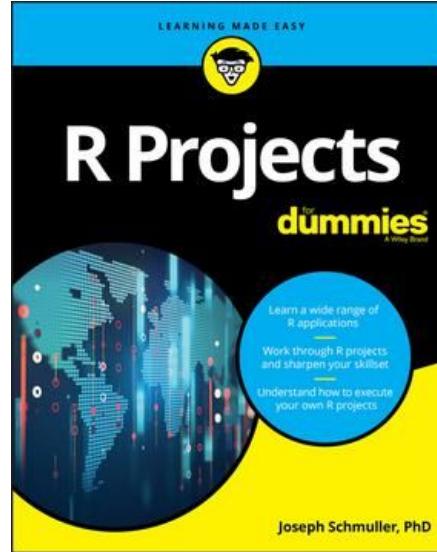
Gender

female

1642
Total Results

464.61
Total Cost

0.28
Average Cost per Result



EDA. WTF.

DEFINE the acronym
Exploratory Data Analysis

What it IS

Learning more about your data without making assumptions
An iterative process

What it IS NOT

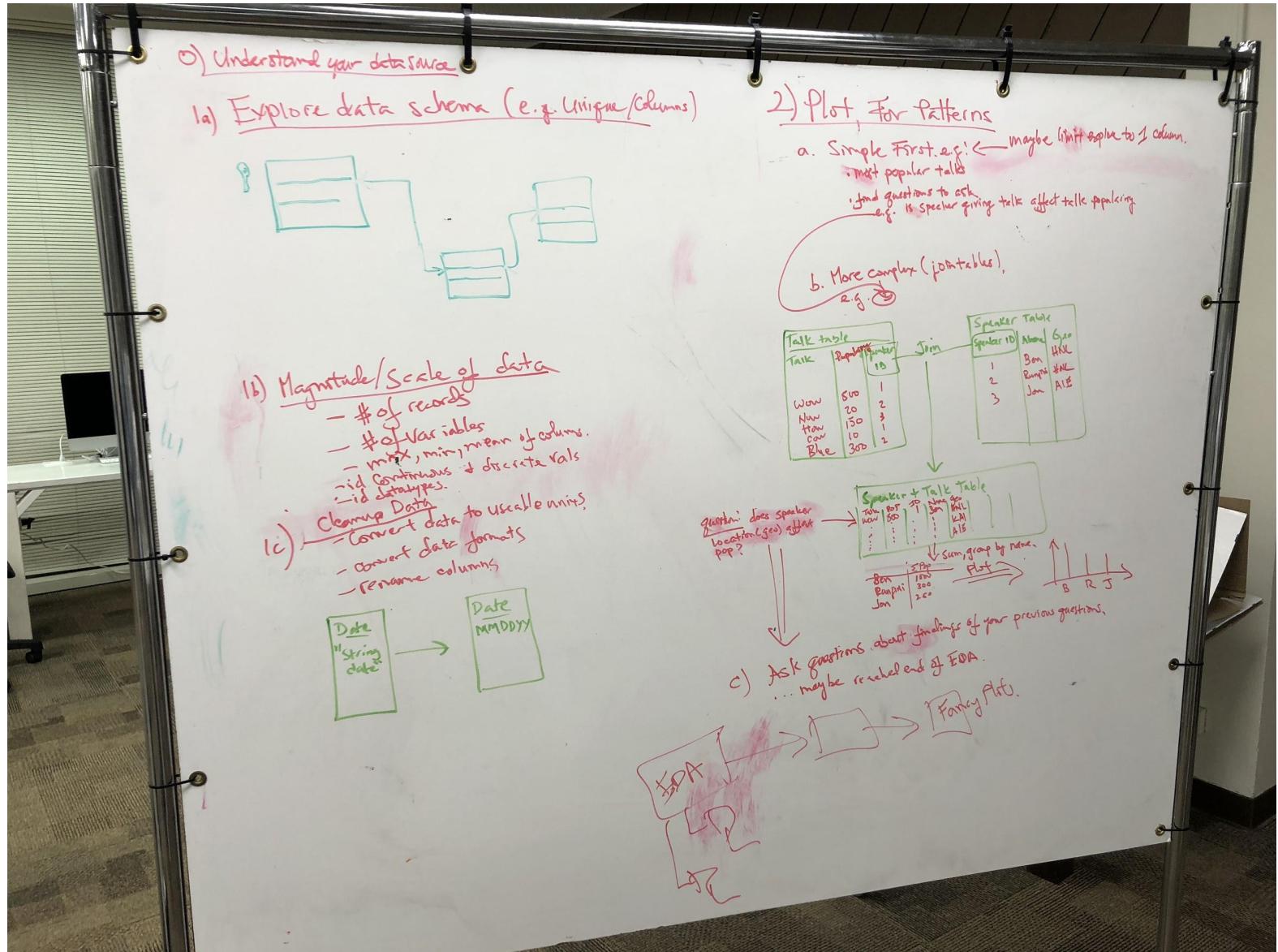
Complete visualizations with fancy aesthetics

HOW it is done

Summary statistics
Quick visualizations (bar plots, histograms, scatterplots)

Example EDA Process

1. Understand data source
2. Explore data schema
3. What is the magnitude/scale of data?
4. Clean up data
5. Plot for patterns: simple → complex
6. Ask more questions
7. Do it again and again



What are ceRcizes?

We created “ceRcizes” as a teaching and learning tool.

Joincersice!!!

```
In [ ]: library(dplyr)
```

```
In [43]: Names <- c("Tori", "Ben", "Hunter", "Runjini")
Macs <- c("MacBook", "Pro2014", "Pro2017", "Pro2011" )
Gender_Mac <- c("Female", "Male", "Male", "Female")
Mac <- data.frame(Names, Macs, Gender_Mac)
```

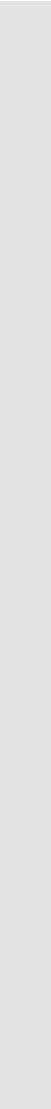
```
In [44]: Names <- c("Jon", "Hunter", "Runjini")
Cars <- c("Fit", "Escape", "Mini")
Gender_Car <- c("Male", "Male", "Female")
Car <- data.frame(Names, Cars, Gender_Car)
```

```
In [49]: Mac
Car
```

Names	Macs	Gender_Mac
Tori	MacBook	Female
Ben	Pro2014	Male
Hunter	Pro2017	Male
Runjini	Pro2011	Female



Exciting Things We Did
and Learned



Runjini: What I Learned

SQL

```
SELECT *  
FROM data_class  
WHERE googled = TRUE
```



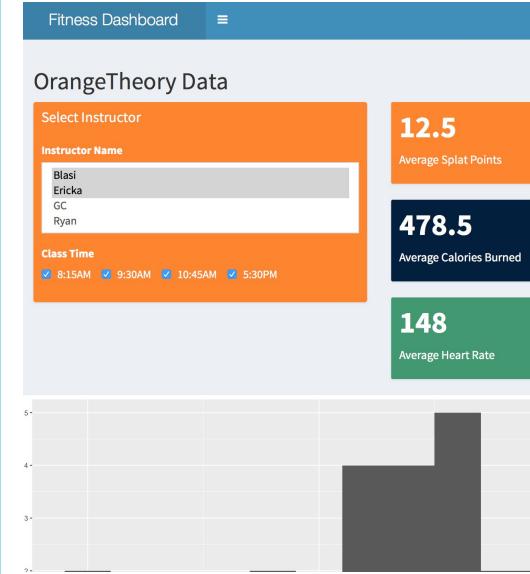
RFacebook



+



Shiny



SQL + Wine

Did you know... 21% of wines on the Kaggle data set of over 130,000 wines are described using the phrases "oak" or "wood"?

```
SELECT COUNT(description)  
FROM wine_reviews  
WHERE description like '%oak%' OR description like '%wood%'
```

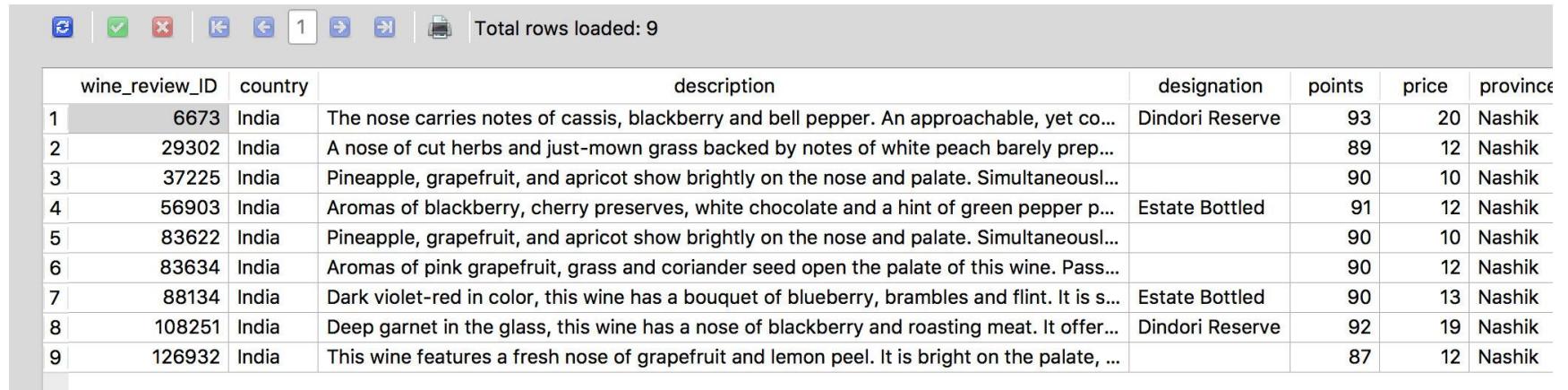
Also researched:

- How many countries are represented on this list? **44**
- How many wines have a perfect score of 100? **19**
- What is the average points rating? **88.44**

SQL + Wine

Provide all information on wines from India.

```
SELECT *
FROM wine_reviews
WHERE country = 'India'
```



The screenshot shows a database query results window with the following details:

- Toolbar icons: Refresh, Checkmark, Cross, Back, Forward, Page Number (1), Print, Total rows loaded: 9.
- Table Headers: wine_review_ID, country, description, designation, points, price, province.
- Table Data (9 rows):

wine_review_ID	country	description	designation	points	price	province
1	6673	India	Dindori Reserve	93	20	Nashik
2	29302	India		89	12	Nashik
3	37225	India		90	10	Nashik
4	56903	India	Estate Bottled	91	12	Nashik
5	83622	India		90	10	Nashik
6	83634	India		90	12	Nashik
7	88134	India	Estate Bottled	90	13	Nashik
8	108251	India	Dindori Reserve	92	19	Nashik
9	126932	India		87	12	Nashik

R + Facebook

All Posts Published

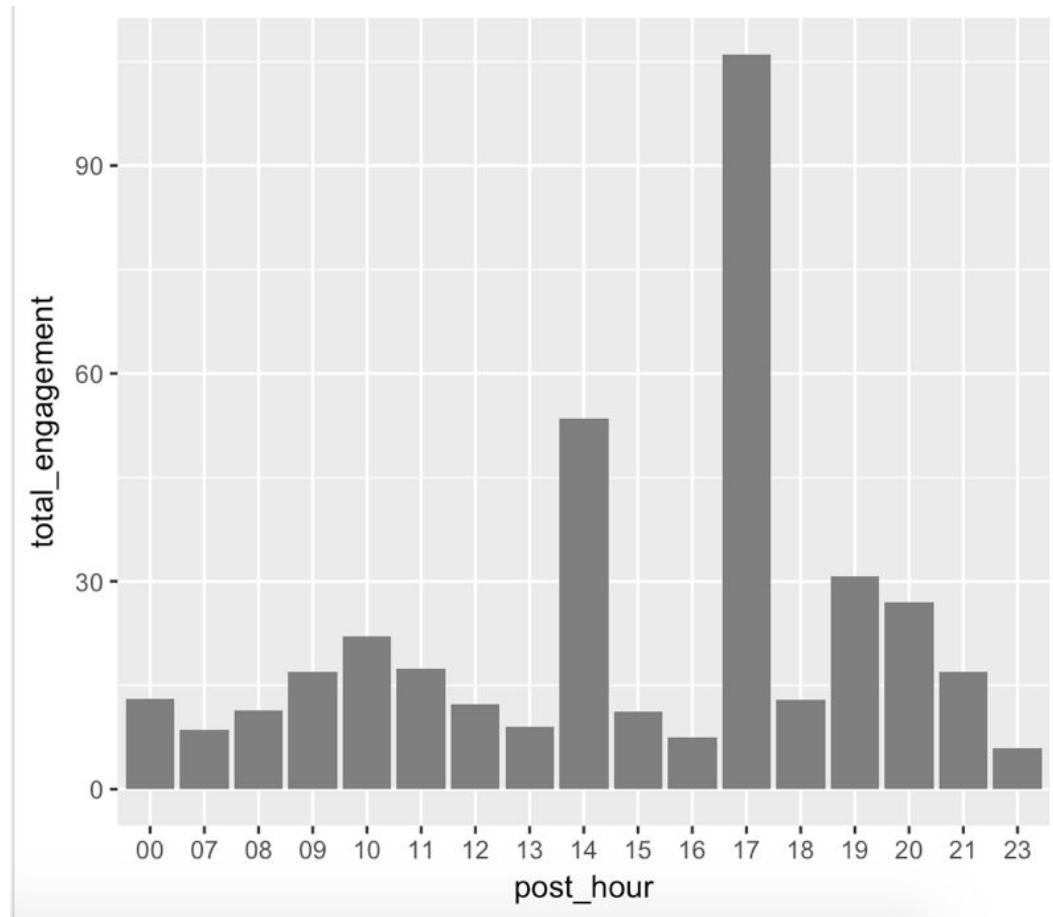
				Reach: Organic / Paid	Reactions	Comments	Shares
Published	Post	Type	Targeting	Reach	Engagement	Promote	
01/10/2018 6:00 pm	 Bright, bold, & unforgettable. For more than five decades			238	0 0 0		
01/10/2018 10:05 am	 Oh what a night! We had an amazing time at the National			1.3K	29 3 1		
01/09/2018 12:00 pm	 Names have power. Learn more about the origins of the			13.7K	210 1 14		
01/08/2018 6:00 pm	 By scouring the seafloor in search of materials for her			616	4 0 0		
01/05/2018 12:00 pm	 Looking to craft a brighter 2018 for yourself but not sure			803	7 0 0		



from_id	from_name	message	created_time	type	link	id
3.106683e+14	Ward Village	Where Hawaii name...	2016-10-26T22:45...	link	http://bit.ly/2ePPUyV	310668279073583...
3.106683e+14	Ward Village	Have you heard? The...	2016-10-27T22:45...	link	http://bit.ly/2eA87ST	310668279073583...
3.106683e+14	Ward Village	We know how busy l...	2016-10-28T22:15...	link	http://bit.ly/2fohLJn	310668279073583...
3.106683e+14	Ward Village	Happy Halloween fro...	2016-10-31T22:55...	photo	https://www.facebook.com/WardVillage106683e14/p...	310668279073583...
3.106683e+14	Ward Village	It's official! We've su...	2016-11-01T22:00...	photo	https://www.facebook.com/WardVillage106683e14/p...	310668279073583...
3.106683e+14	Ward Village	Will you be joining u...	2016-11-02T22:45...	photo	https://www.facebook.com/WardVillage106683e14/p...	310668279073583...
3.106683e+14	Ward Village	Coming to Ward Vill...	2016-11-05T05:03...	photo	https://www.facebook.com/WardVillage106683e14/p...	310668279073583...
3.106683e+14	Ward Village	Last year, the Ice Rin...	2016-11-08T04:10...	link	http://bit.ly/2fudoKb	310668279073583...
3.106683e+14	Ward Village	Next week, we'll be ...	2016-11-09T02:54...	photo	https://www.facebook.com/WardVillage106683e14/p...	310668279073583...
3.106683e+14	Ward Village	Grand Opening of S...	2016-11-11T19:34...	video	https://www.facebook.com/WardVillage106683e14/p...	310668279073583...
3.106683e+14	Ward Village	One of our favorite t...	2016-11-16T07:52...	photo	https://www.facebook.com/WardVillage106683e14/p...	310668279073583...
3.106683e+14	Ward Village	Mahalo to everyone ...	2016-11-16T23:31...	video	https://vimeo.com/171111111	310668279073583...
3.106683e+14	Ward Village	Today marks such a ...	2016-11-17T18:51...	video	https://www.facebook.com/WardVillage106683e14/p...	310668279073583...
3.106683e+14	Ward Village	NA	2016-11-19T01:55...	photo	https://www.facebook.com/WardVillage106683e14/p...	310668279073583...
3.106683e+14	Ward Village	This week marked 2 ...	2016-11-19T18:00...	photo	https://www.facebook.com/WardVillage106683e14/p...	310668279073583...

R + Facebook

Using RFacebook, I could determine the best times for us to post in terms of which hours would yield the highest total engagement: 5pm, 2pm, and 7pm.



Shiny

Fitness Dashboard

Quick Stats

Performance Over Time

Frequency

Quick Stats

Select Instructor

Instructor Name

Blasi Ericka GC Ryan Zach

Class Time

8:15AM 9:30AM 10:45AM 5:30PM

15

Average Splat Points

480

Average Calories Burned

148

Average Heart Rate

Selected Workout Results

Instructor	month	day	year	time_adj	splats	calories	averagehr
Blasi	May	02	2018	5:30PM	17	360	151
Blasi	April	24	2018	5:30PM	20	447	152
Blasi	August	10	2018	5:30PM	0	11	NA
Blasi	August	14	2018	5:30PM	17	434	150
Blasi	July	18	2018	5:30PM	18	318	143
Blasi	Anril	04	2018	5:30PM	13	584	152

Shiny

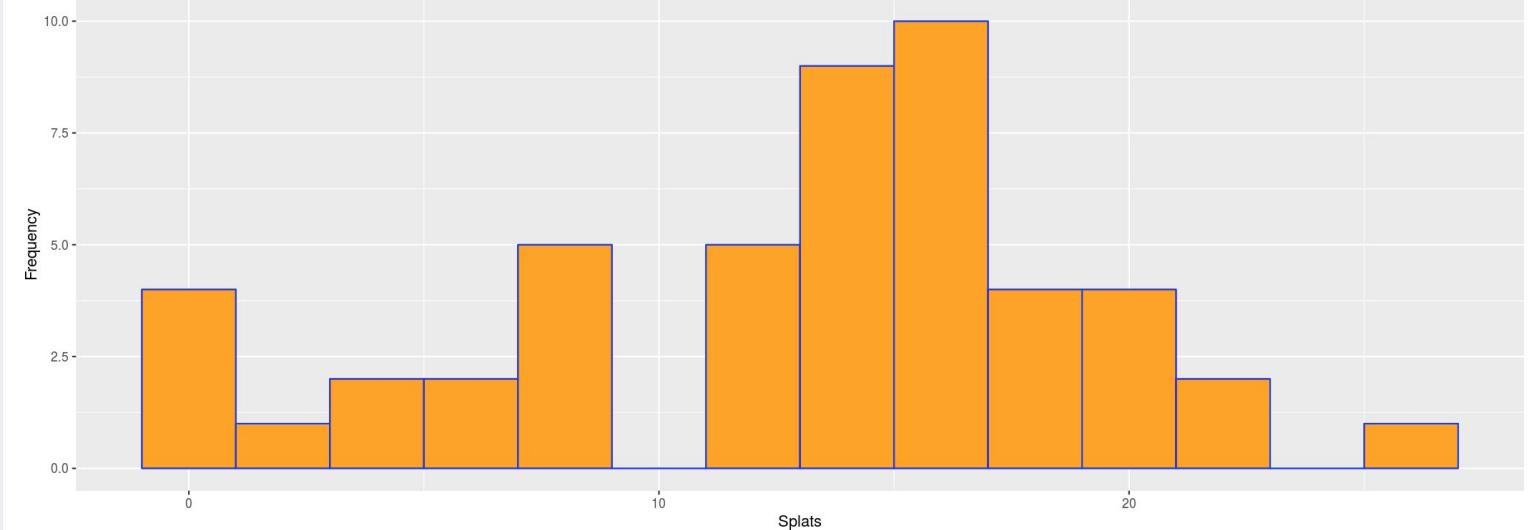
Frequency

Select Instructor

Instructor Name

Blasi Ericka GC Ryan

Frequency of Splats Based on Instructor



Jon: What I Learned

Too Much Pie!



Stories on Jupyter



Auto-Tagging Email



Too much pie!



Stories on Jupyter

	I	P	Q	R	S	T
1	Facility Name (Appendix A1)	NEL_TSS Score AF Wtd Avg	NEL_Turbidity Score AF Wtd Avg	NEL_N Score AF Wtd Avg	NEL_NN Score AF Wtd Avg	NEL_AN Score AF Wtd Avg
2						
3	Ahuimanu Dewatering Facility	0.011	0.011	-	-	-
4	Aiea Fire Station	0.132	0.132	0.011	0.011	0.011
5	Aikahi Fire Station	-	-	-	-	-
6	Aircraft One Fire Station	-	-	-	-	-
7	Airport Vehicle Maintenance Yard	0.031	0.031	-	-	-
8	Ala Moana Regional Park (including Maintenance Baseyard)	0.007	0.007	-	-	-
9	Ala Wai Golf Course	0.472	0.472	0.493	0.493	0.493
10	Alapai Police Station	0.124	0.124	-	-	-
11	Central Fire Station	0.015	0.015	-	-	-
12	Central Oahu Regional Park (including Maintenance Baseyard)	0.085	0.085	-	-	-
13	East Kapolei Fire Station	0.004	0.004	-	-	-
14	Ewa Beach Fire Station	0.004	0.004	-	-	-
15	Ewa Refuse Convenience Center	0.160	0.160	-	-	-

Stories on Jupyter

Jupyter Notebooks offer a multi-purpose solution to tell a story.

Description

Step 4.c: Calculate Age Factor Weighted (AFW) Average Exceedances

We claim that more recent samples better represent current housekeeping and sampling necessary to represent long term trends and to dampen whipsaw effects that may occur. The weighted average exceedance(AFW average exceedance) considers these two viewpoints on older samples. The equation is:

$$AFW_p = \frac{\sum C_p e^{-SRn}}{\sum e^{-SRn}}$$

Coding

```
12 df[AFColName + '_AF'] = np.exp(-df[AFColName + '_DR']) #calculate each year's weight
13 df[AFColName + '_AF*Val'] = df[AFColName + '_AF']*df[AFColName] #apply weight to individual facility values
14 #for each facility, add af weights and wt applied values up. enter to temp datafram
15 tmp = df.groupby(['Facility_Name'])[[AFColName + '_AF', AFColName + '_AF*Val']].sum()
16 columns = {AFColName + '_AF':AFColName + '_AFSUM', AFColName + '_AF*Val':AFColName + '_AF*ValSum'}
17 # display(tmp)
18 df=pd.merge(df,tmp,on='Facility_Name',how='inner') #merge
```

Results

Fiscal Year	Deficient_Count	Total_Count	Deficiency_Rate	Deficiency_Rate_AF
2013	1	38	0.026316	0.018316
2014	0	38	0.000000	0.049787
2015	3	38	0.078947	0.135335
2016	1	38	0.026316	0.367879
2017	1	38	0.026316	1.000000

Categorize Email

Use previously categorized email...

FROM	SUBJECT	CATEGORIES
Underground Solutions	Comparing HDD to Open-Cut	Meeting
EnviroCert International, Inc.	Hurricane Florence: EnviroCert East Coast Office Closure	Flyers
Storm Water Solutions	Earn PDH credits without the travel expenses SWS Webinar Fest	Flyers
Cortana	Heads-up	Meeting

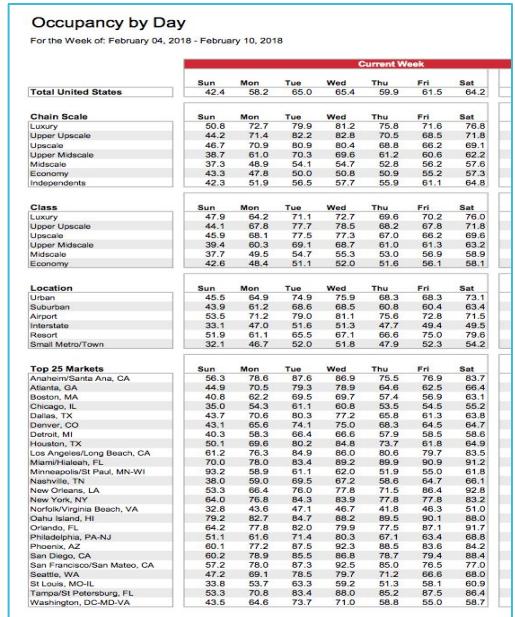


To categorize new email.

FROM	SUBJECT	CATEGORIES
Lauren Baltas, Storm Water...	Hurricane Florence threatens the East Coast	<input type="checkbox"/>
Jeremy Huckaby	Statement of Clarification	<input type="checkbox"/>
Weeklysafety.com, LLC	NEW OFFER! Fast 100 Now Includes Construction AND General I...	<input type="checkbox"/>

Tori: What I Learned

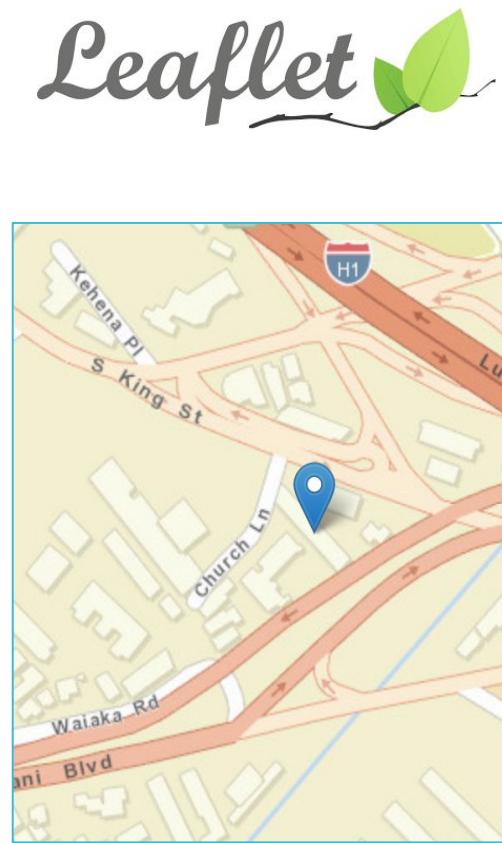
PDF Tools



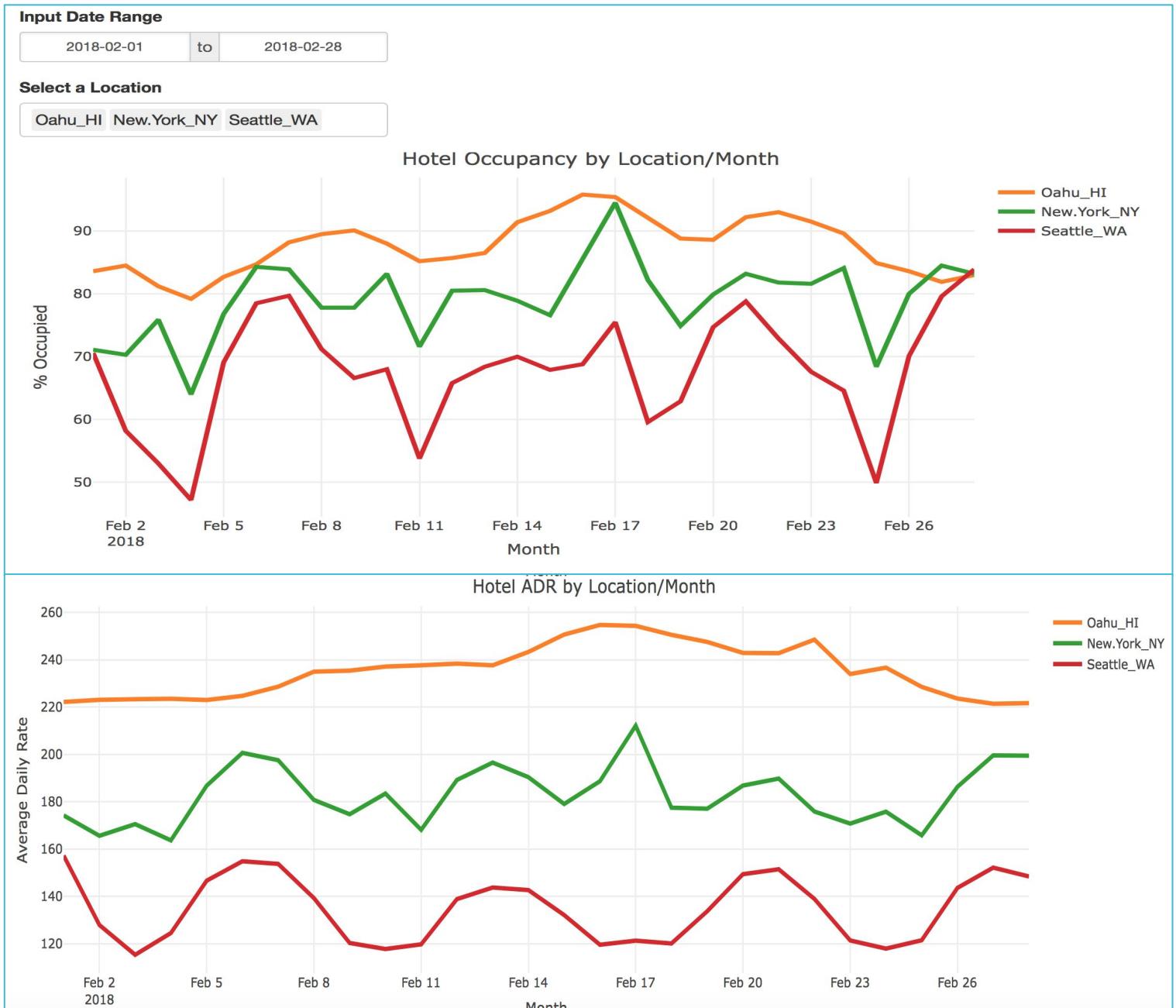
Shiny Dashboards



Leaflet

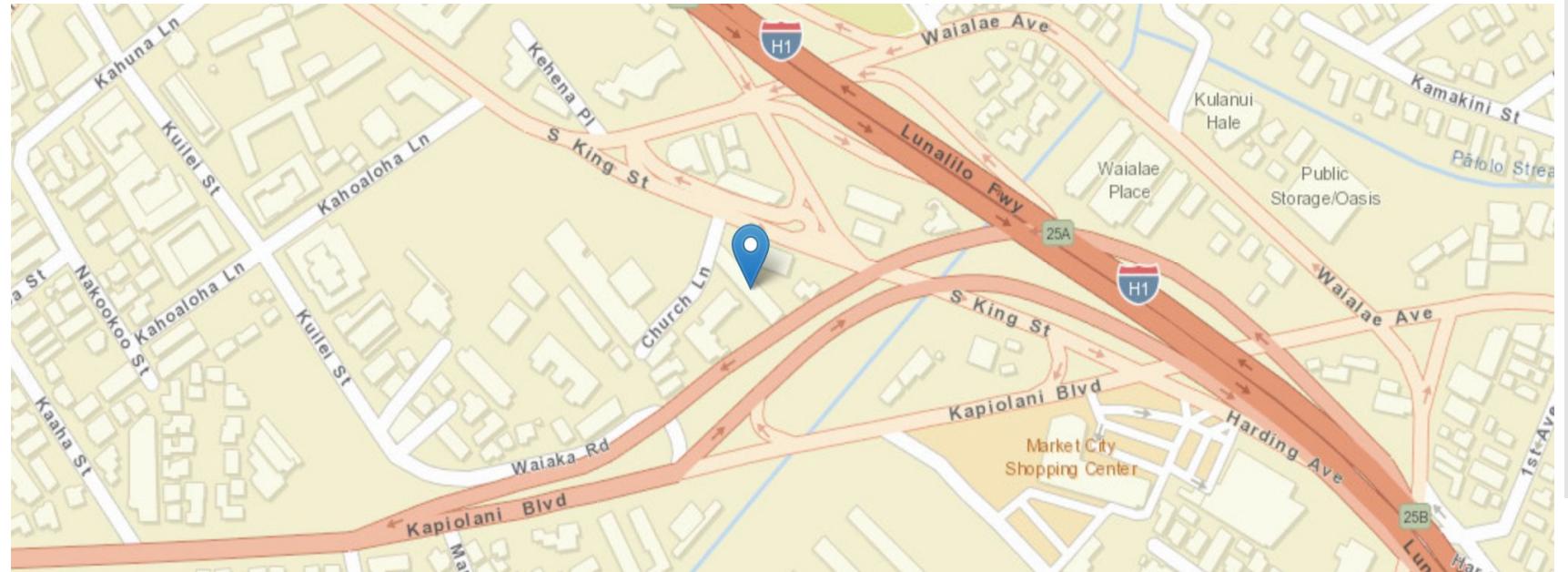


Shiny Dashboards

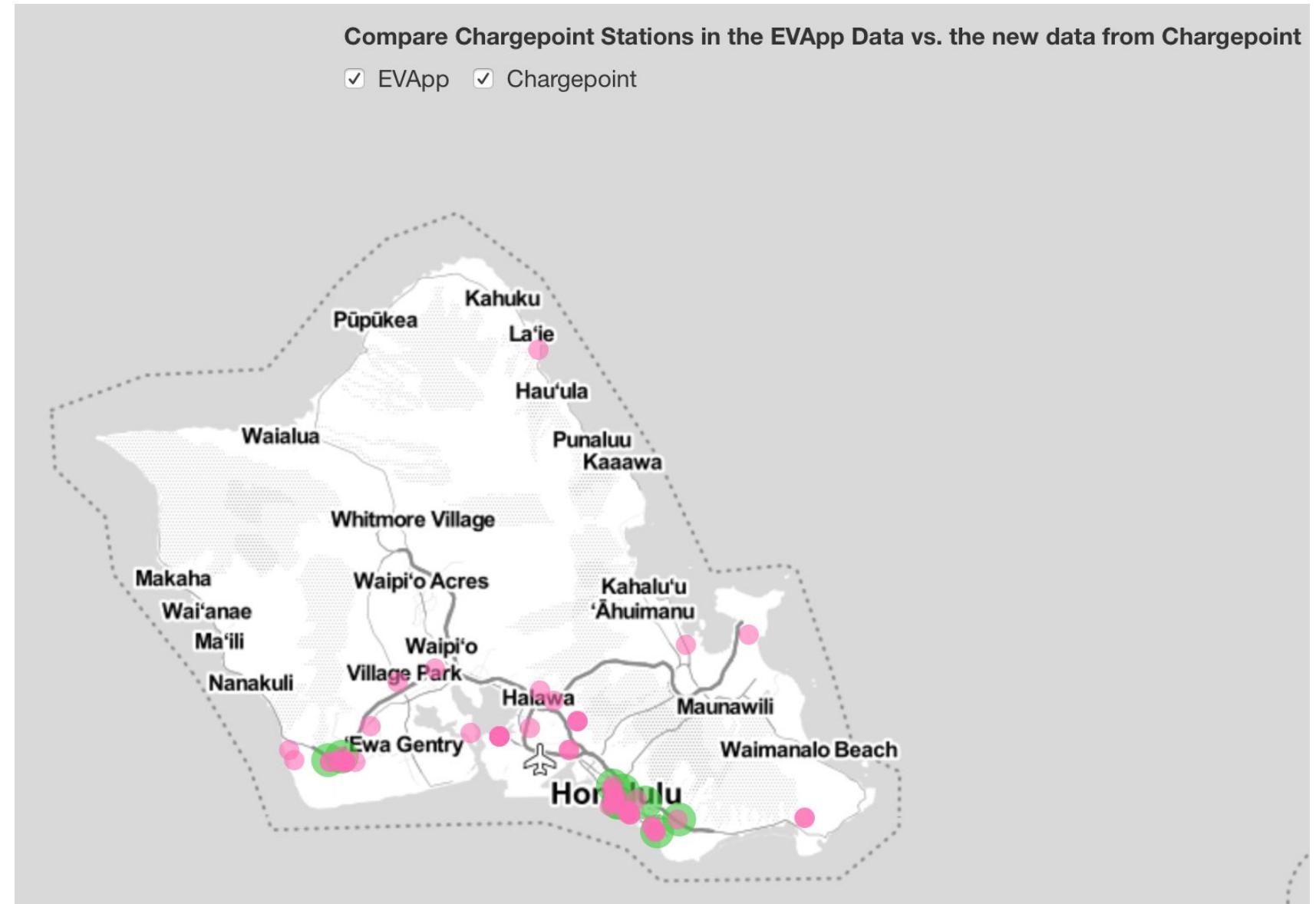


Leaflet

```
1 #install.packages("leaflet")
2 library(leaflet)
3 library(RColorBrewer)
4 library(shiny)
5
6 m <- leaflet()
7 m <- addTiles(m)
8 m <- addMarkers(m, lng = -157.817508, lat = 21.289207)
9
10 m %>% addProviderTiles(providers$Stamen.Watercolor)
```



EV Charging App



The Future

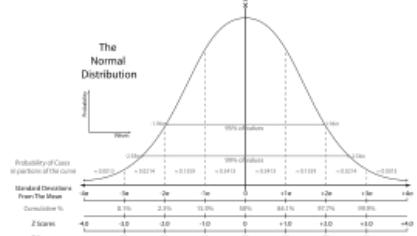
What's Next?

Python?

Practice!!

```
37 import pygame
38
39 def text_to_screen(scr
40     color = (2
41
42     try:
43
44         text = str(tex
45         font = pygame.
46         text = font.re
```

Runjini



Statistics Review

DataCamp



Jon

Kaggle
Competitions



Stock Valuation

Tori



Leaflet



DataCamp

Google those error
messages

A photograph of Leonardo DiCaprio in a black tuxedo and bow tie, smiling and holding a glass of champagne towards the camera. He is standing in front of a dark background with blurred lights and confetti falling around him.

WE DID IT

TIME TO CELEBRATE

The end.

Mahalo!