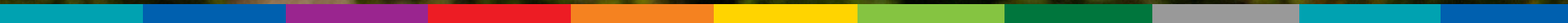




# Moving Beyond Excel and Becoming More Data Science-y


Evolving How We Think About  
and Use Data





*Once upon a time...*



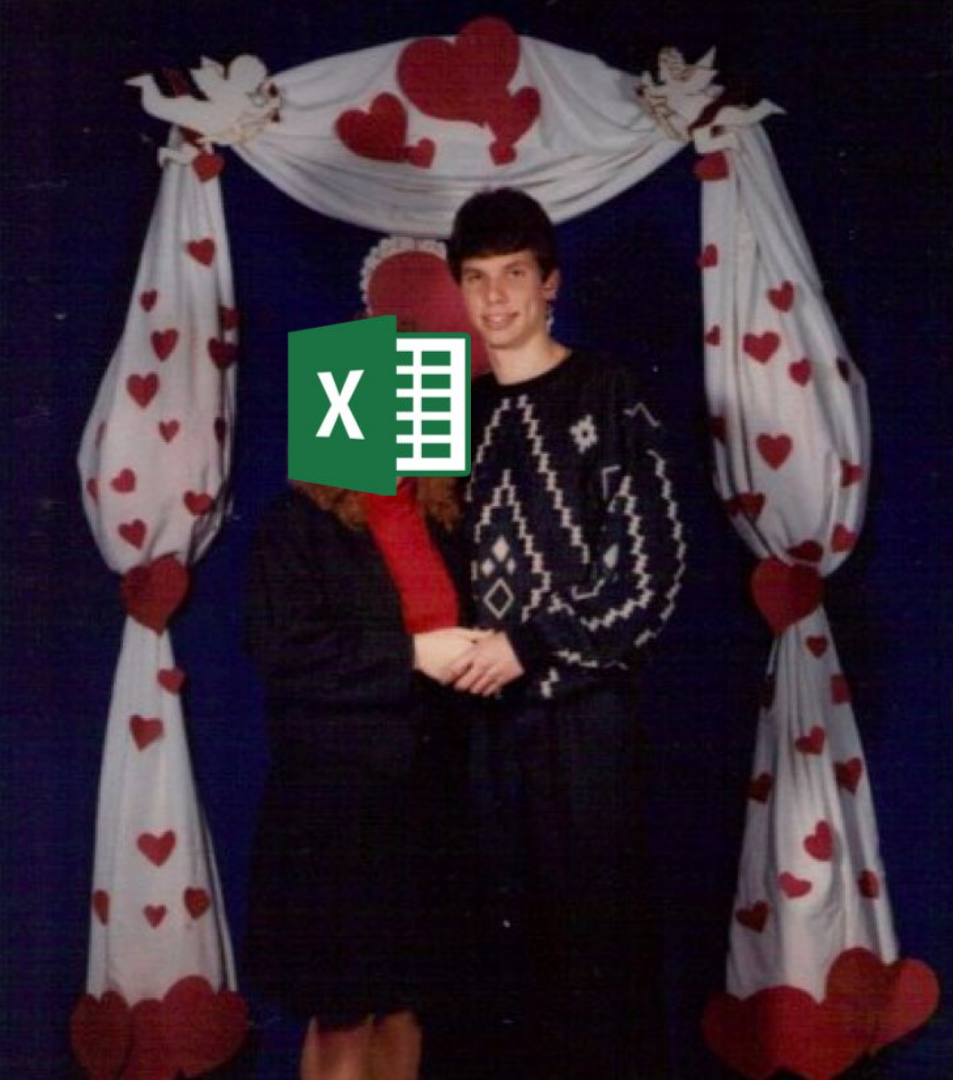


*...there was an analyst.*

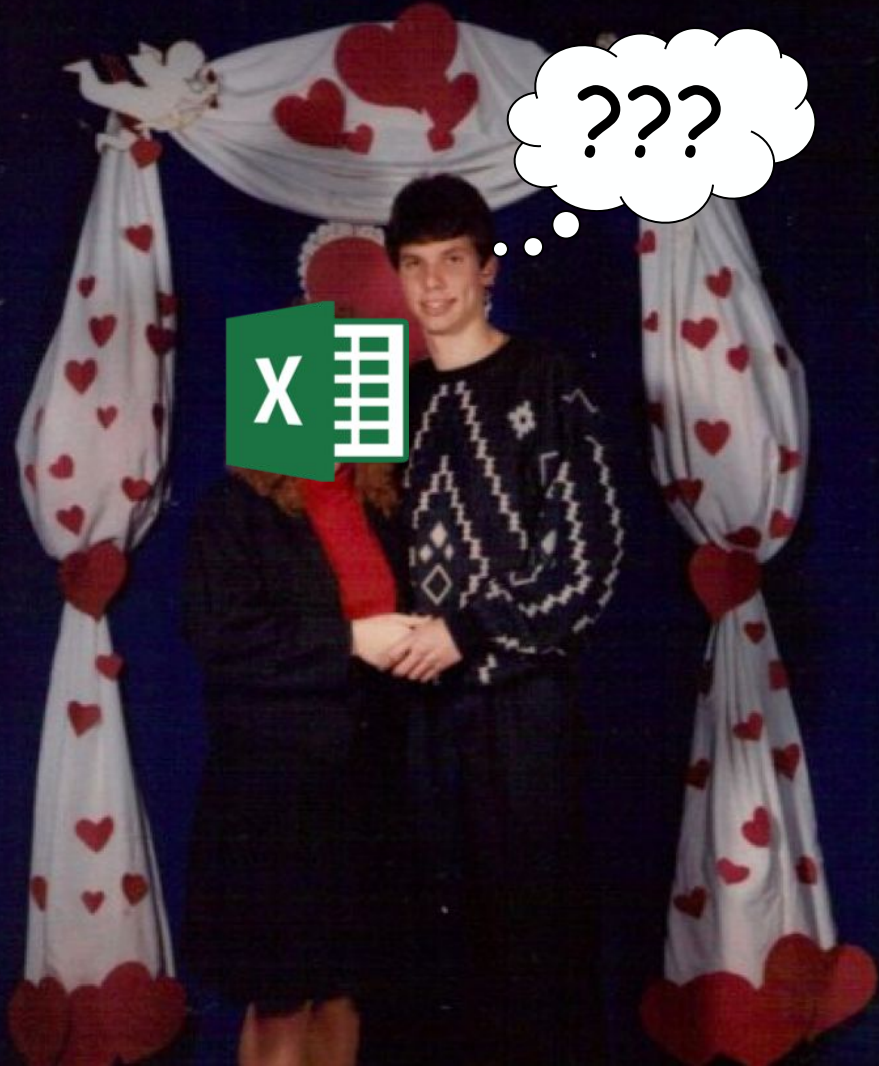


*And that analyst was me.*





And I had a long and  
abiding relationship  
with Microsoft Excel.



But I started to  
question the merits  
of that relationship.



I decided I would  
become a **data  
scientist!**

I would do  
**magical things**  
with the data!

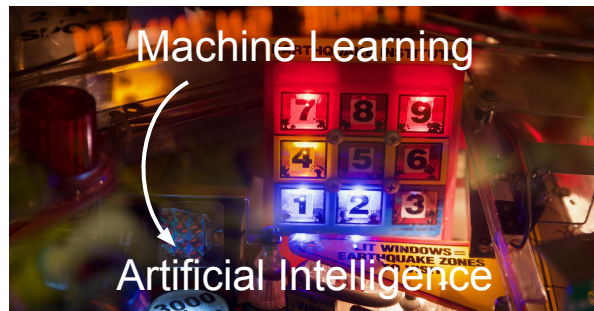
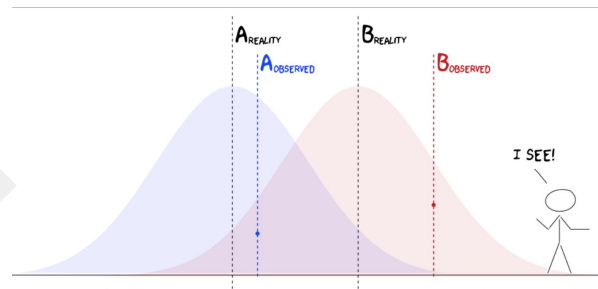


I quickly realized I  
would never be a  
**true** data scientist.





But I was intrigued...

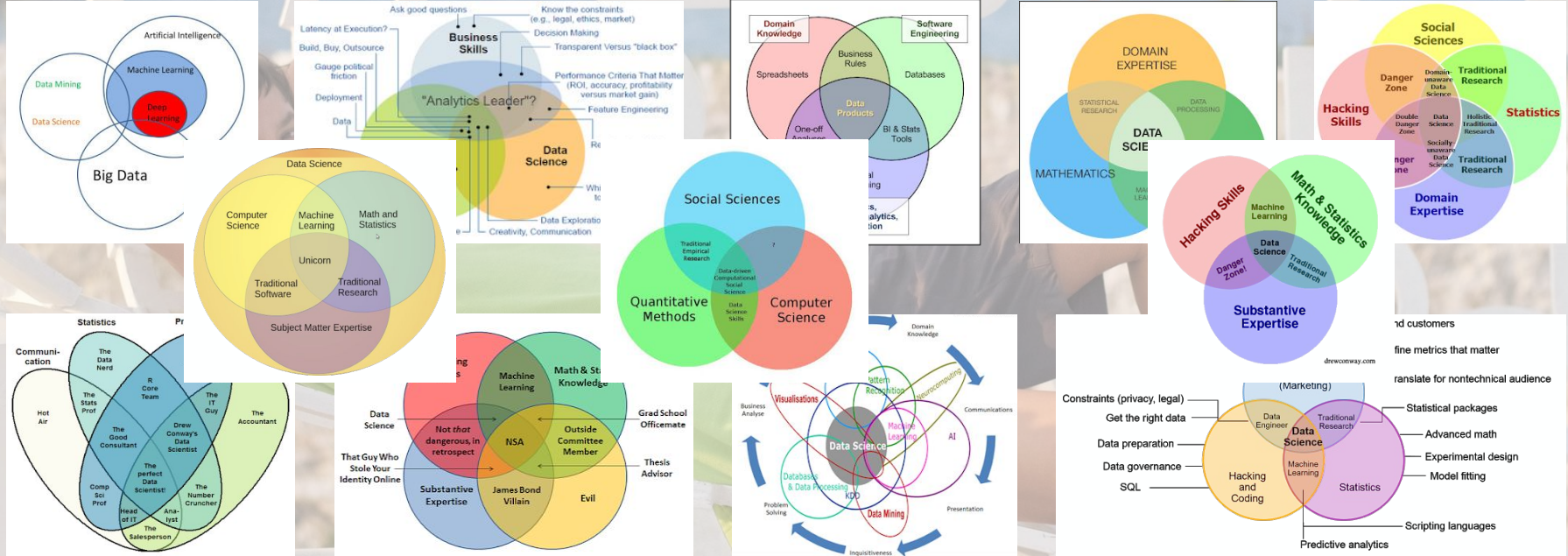


# What is data science?





# What is data science?



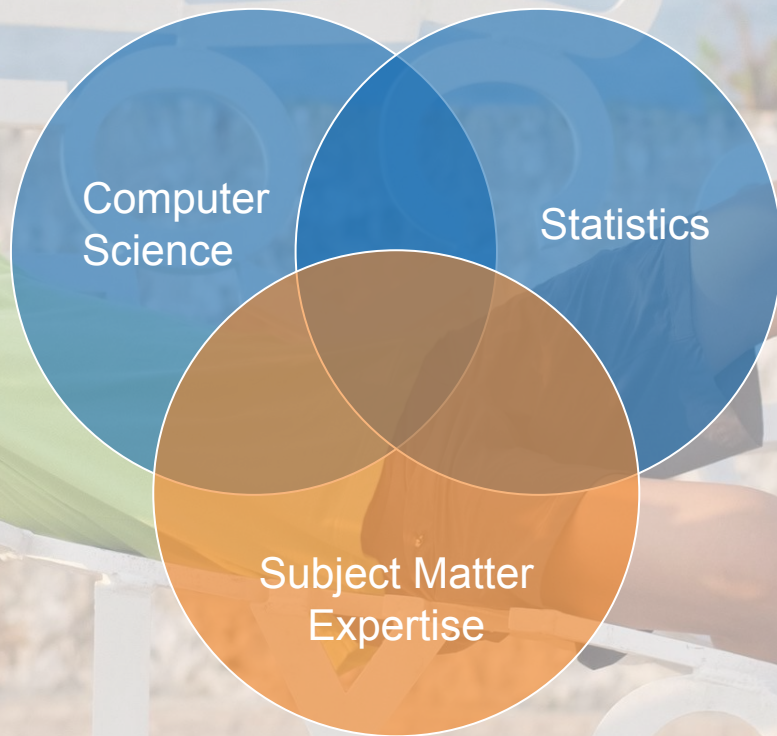
<https://www.kdnuggets.com/2016/10/battle-data-science-venn-diagrams.html>

# An 85% Confidence Venn Diagram





# Let's Start with Subject Matter Expertise



# We've got this! We already...

...have deep knowledge of **how the data is collected**.

...have a deep understanding of **marketing**.

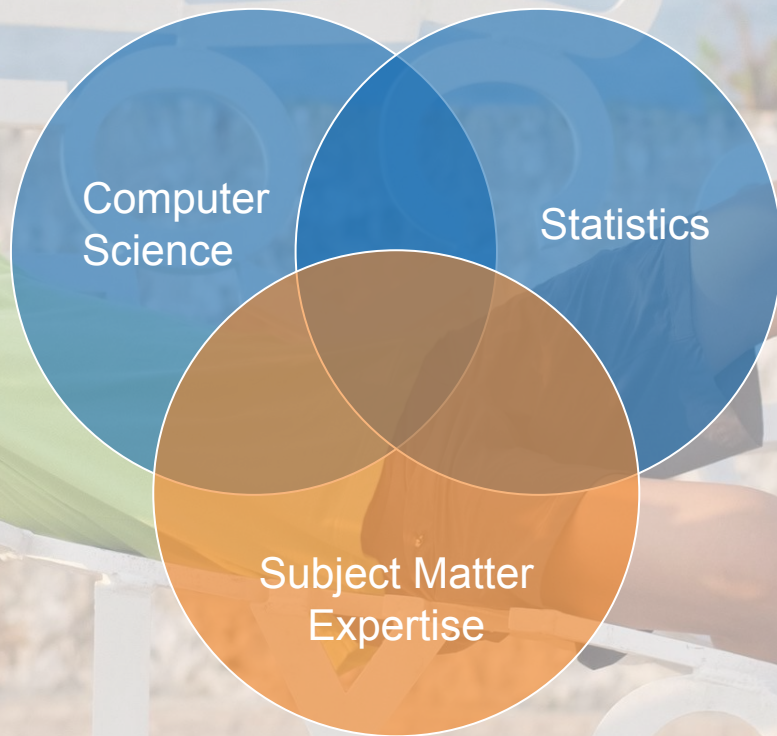
...are able to effectively articulate the **business problems** faced by our stakeholders.

...are able to effectively communicate the **results of analyses** to stakeholders.

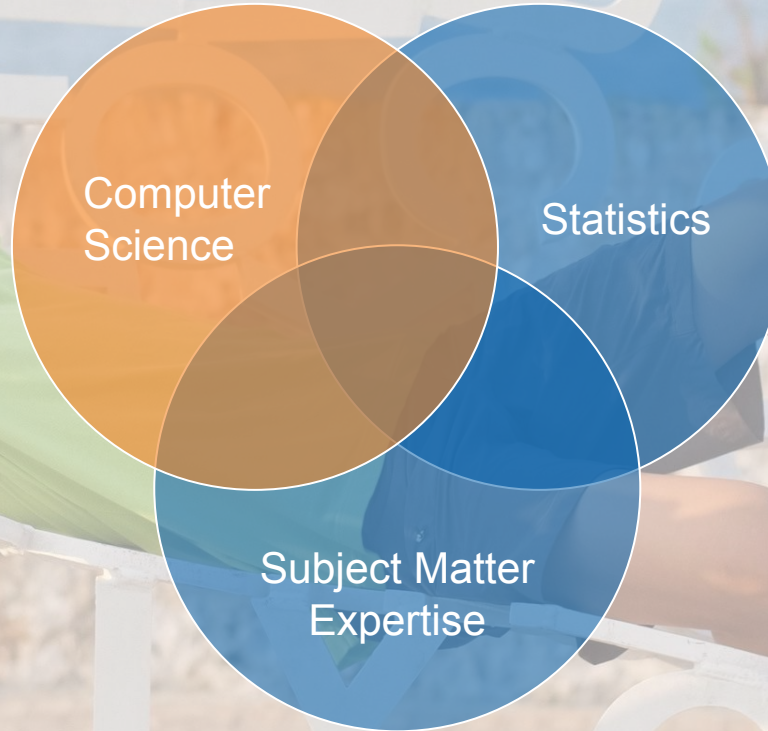




# That's It for Subject Matter Expertise



# Let's Talk...Computer Science!



# Text-Based Programming

```
[ ] ###  
# Generate global vars from Predictor Configurations  
##  
LOCAL_PREDICTOR_CONFIG = COMPANY + '-' + PROPERTY + '-' + PREDICTOR + '-'  
GCS_BUCKET_LOCATION = 'astrologger-2/' + COMPANY + '/' + PROPERTY + '/' +  
GCS_PREDICTOR_CONFIG = GCS_BUCKET_LOCATION + 'predictorConfig.pkl'  
  
###  
# Save configuration settings to file  
##  
def saveConfig(PREDICTOR_CONFIG):  
    f = open(LOCAL_PREDICTOR_CONFIG, 'w')  
    pickle.dump(PREDICTOR_CONFIG, f)  
    f.close()  
  
    # Copy updated file to GCS  
    !gsutil cp {LOCAL_PREDICTOR_CONFIG} {GCS_PREDICTOR_CONFIG}  
  
###  
# Check for .txt file, and if it does not exist  
##  
def getConfig():  
    # Copy file from GCS to local disk if it exists  
    !gsutil cp gs://{GCS_PREDICTOR_CONFIG} {LOCAL_PREDICTOR_CONFIG}  
  
    # Create new config if non exists already  
    if not(os.path.isfile(LOCAL_PREDICTOR_CONFIG)):  
        # Create new Config  
        PREDICTOR_CONFIG = [{  
            'COMPANY': COMPANY,  
            'PROPERTY': PROPERTY,  
            'PREDICTOR': PREDICTOR,  
            'SETTINGS': {  
                'TYPE': TYPE,  
                'CUSTOM_DATA_SOURCE_ID': "",  
                'WEB_PROPERTY_ID': "",  
                'ACCOUNT_ID': ""  
            }  
        }  
    ]  
    saveConfig(PREDICTOR_CONFIG)
```

CONNECTIVITY

AUTOMATION

```
1 # Load the necessary libraries.  
2 if (!require("pacman")) install.packages("pacman")  
3 pacman::p_load(RSiteCatalyst, tidyverse)  
4  
5 # Load the username, shared secret, and report suite ID  
6 username <- Sys.getenv("ADOBE_API_USERNAME")  
7 secret <- Sys.getenv("ADOBE_API_SECRET")  
8  
9 # Authorize Adobe Analytics.  
10 SCAuth(username, secret)  
11  
12 # Set the date range.  
13 start_date <- Sys.Date() - 30  
14 end_date <- Sys.Date()  
15  
16 # Get the Report Suite ID  
17 rsid_df <- GetReportSuiteIDs(start_date, end_date)  
18 numsuites <- nrow(rsid_df)  
19  
20 # Set a counter just to keep track of the progress  
21 i <- 1  
22  
23 # Define a function to pull high-level traffic data  
24 get_traffic <- function(rsid) {  
25     # Pull the traffic data  
26     # Output the RSID to the console so we know what's being processed  
27     cat("Processing", i, "of", numsuites, ":", rsid, "\n")  
28  
29     # Increment counter  
30     i <- i + 1  
31  
32     # Pull the visits and pageviews data  
33     traffic_df <- QueueSummary(rsid, start_date, end_date)  
34  
35     # Arrange by date  
36     traffic_df <- arrange(traffic_df, date)  
37  
38     # Select columns  
39     select(traffic_df, date, visits, pageviews)  
40  
41     # Save the summary to a CSV  
42     write.csv(traffic_df, "rsid_traffic_summary.csv", row.names=FALSE)  
43  
44     return(traffic_df)  
45 }  
46  
47 # Loop through the Report Suite IDs and pull the traffic data  
48 for (i in 1:numsuites) {  
49     get_traffic(rsid_df[i, "rsid"])  
50 }  
51
```

DATA  
VISUALIZATION



A photograph of a public fountain in a park. In the foreground, a white mermaid statue stands on a dark base, spraying a high arc of water into the air. To the right, a swan statue also sprays water upwards. The fountain is surrounded by a black metal fence and a bed of pink and white flowers. In the background, there are trees and a red building.

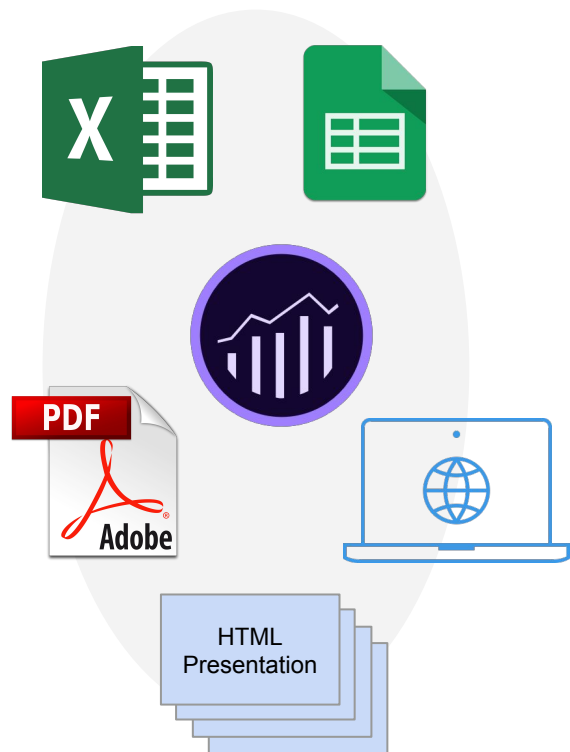
# Connectivity

Inputs and Outputs

## INPUTS\*



## OUTPUTS\*



\* These lists are essentially infinite, so they are limited here just to systems that Tim has pulled/pushed to using R.



A close-up, low-angle shot of a pinball machine. The central focus is a blue devil character with a large, bushy mustache and a goatee, wearing a small crown. He is positioned behind a large, ornate drum with a yellow and red flame pattern. The machine's playfield is visible, showing various targets, bumpers, and a yellow and purple ball. The background is filled with bright, colorful lights, creating a vibrant and dynamic atmosphere.

# Automation

and Reusability



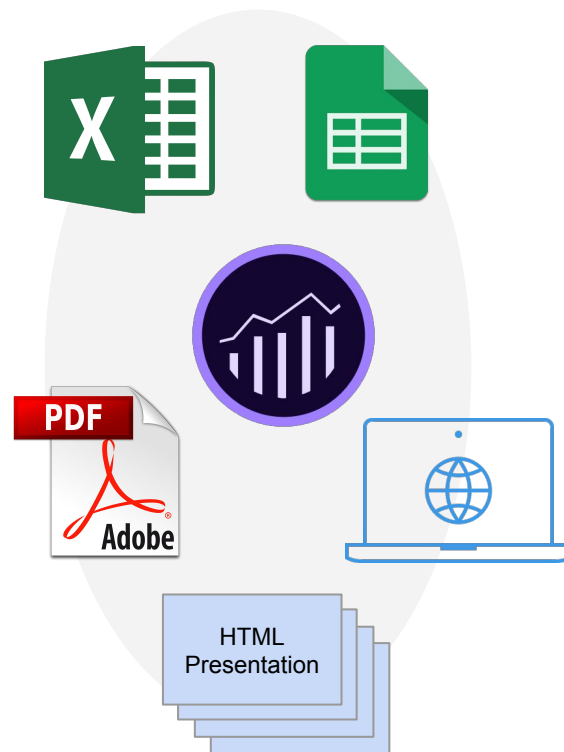


**“How much traffic did we get to each of our sites?”**

## INPUTS\*



## OUTPUTS\*



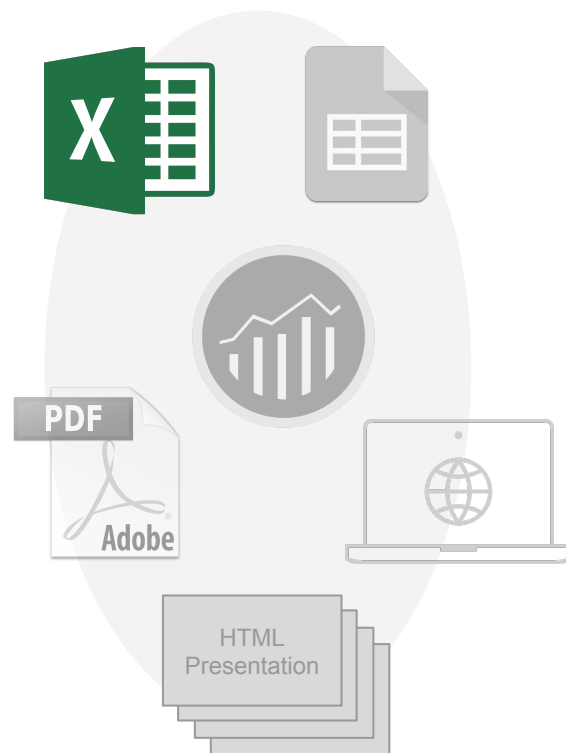
*\* These lists are essentially infinite, so they are limited here just to systems that Tim has pulled/pushed to using R.*



## INPUTS\*



## OUTPUTS\*



\* These lists are essentially infinite, so they are limited here just to systems that Tim has pulled/pushed to using R.

# Visits and Page Views for 475 Report Suites!

```

1 # Load the necessary libraries.
2 if (!require("pacman")) install.packages("pacman")
3 pacman::p_load(RSiteCatalyst, tidyverse)
4
5 # Load the username, shared secret, and report suite ID
6 username <- Sys.getenv("ADOBE_API_USERNAME")
7 secret <- Sys.getenv("ADOBE_API_SECRET")
8
9 # Authorize Adobe Analytics.
10 SCAuth(username, secret)
11
12 # Set the date range.
13 start_date <- Sys.Date() - 31 # 30 days back from yesterday
14 end_date <- Sys.Date() - 1 # Yesterday
15
16 # Get the Report Suites
17 rsid_df <- GetReportSuites()
18 numsuites <- nrow(rsid_df)
19
20 # Set a counter just to report out on the processing count
21 i <- 1
22
23 # Define a function to pull high-level data
24 get_traffic <- function(rsid = "shell-ac-site"){
25
26   # Output the RSID to the console so we'll know what's being processed
27   cat("Processing", i, "of", numsuites, ":", rsid)
28
29   # Increment counter
30   i <- i+1
31
32   # Pull the visits and pageviews data
33   traffic_df <- QueueSummary(rsid,
34     date = "",
35     metrics = c("visits", "pageviews"),
36     date.from = start_date,
37     date.to = end_date)
38 }
39
40 # Pull the summary data for all report suites
41 traffic_summary <- map_dfr(rsid_df$rsid, get_traffic)
42
43 # Add the report suite name back on and sort by visits descending
44 traffic_summary_full <- left_join(traffic_summary, rsid_df, by = c("reportsuite" = "rsid")) %>%
45   mutate(visits = as.numeric(visits), pageviews = as.numeric(pageviews)) %>%
46   arrange(desc(visits)) %>%
47   select(-url)
48
49 # Save the summary to a CSV
50 write.csv(traffic_summary_full, "rsid_traffic_summary.csv", row.names=FALSE)
51

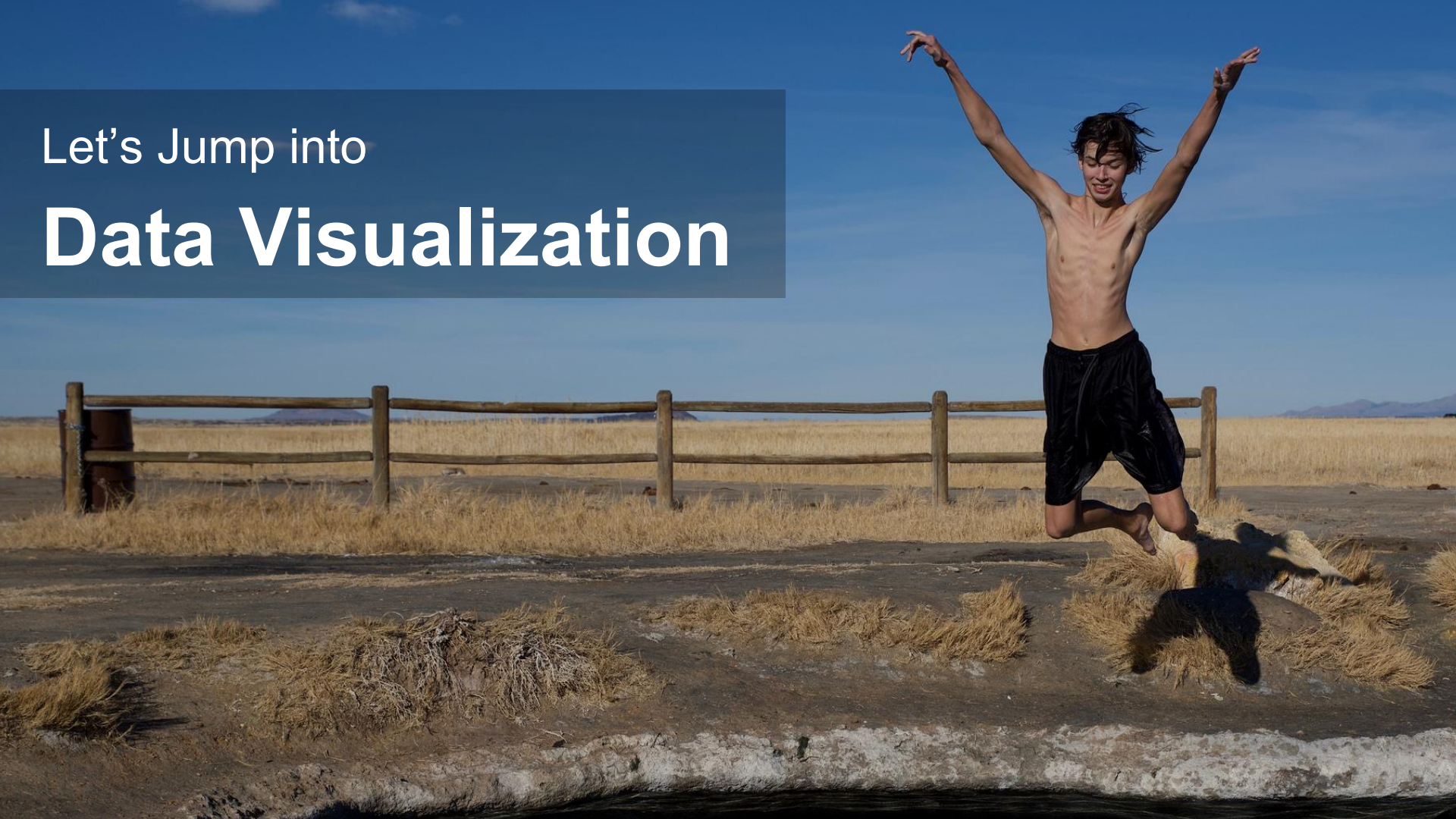
```



	A	B	C	D	E
1	reportsuite	visits	pageviews	site_title	virtual
2		6800684	14385490		NA
3		6298882	13477663		NA
4		3927516	6		NA
5		1044121	2582688		NA
6		808991	3803033		NA
7		679666	1011706		NA
8		563177	845258		NA
9		436928	794315		NA
10		362959	521924		NA
11		332231	532405		NA
12		329702	1463000		NA
13		299037	1604536		TRUE
14		285318	454644		NA
15		269896	795125		NA
16		247984	473322		NA
17		238520	396465		NA
18		237979	486079		NA
19		230299	609724		NA
20		224902	415389		NA
21		224055	244575		NA
22		222245	361846		NA
23		222093	1084912		NA
24		221132	587837		NA
25		160904	277098	s)	NA
26		118195	205674		NA
27		117493	236634		NA
28		99355	258287		NA
29		93757	164877		NA
30		91087	916501		TRUE
31		81557	150017		NA
32		76938	167569		NA
33		75003	157229		NA
34		71960	220363		NA
35		70105	143900		NA
36		69903	141548		NA
<hr/>					
465		0	0		NA
466		0	0		NA
467		0	0		NA
468		0	0		NA
469		0	0		NA
470		0	0		NA
471		0	0		NA
472		0	0		NA
473		0	0		NA
474		0	0		NA
475		0	0		NA
476		0	0		TRUE
477					

Let's Jump into

# Data Visualization







# Mapping Out a Data Ecosystem



	A	B
1	<b>System</b>	
2	Adobe Analytics	
3	Adobe Audience Manager	
4	Domo	
5	Salesforce.com	
6	ExactTarget	
7	Google Analytics	
8	BigQuery	
9		
10		
11		
12		

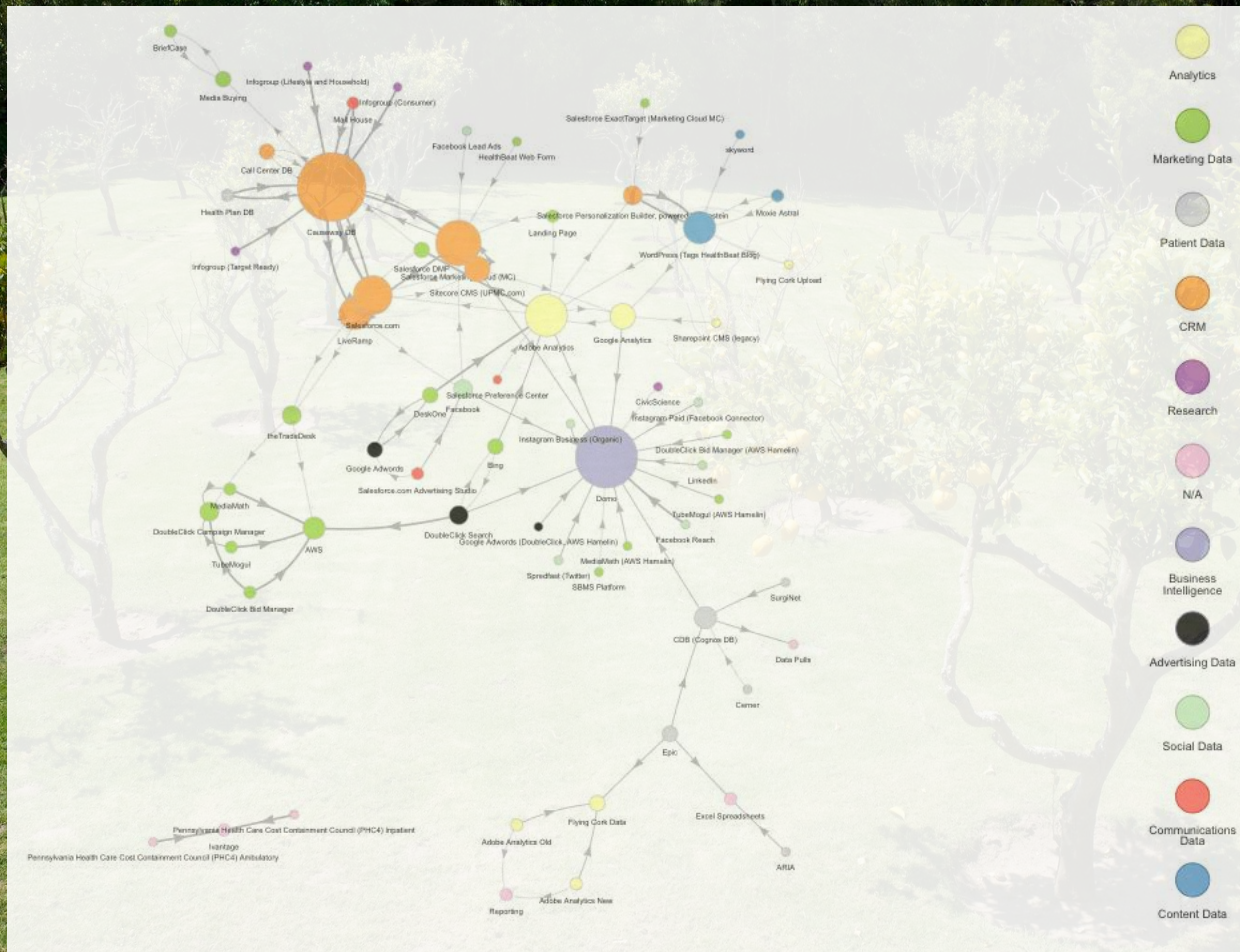
+   Systems ▾   Data Feeds ▾



	A	B
1	<b>From</b>	<b>To</b>
2	Adobe Analytics	Adobe Audience Manager
3	Adobe Analytics	Domo
4	Adobe Analytics	Salesforce.com
5	Adobe Analytics	ExactTarget
6	Adobe Audience Manager	Adobe Analytics
7	Adobe Audience Manager	Domo
8	Salesforce.com	Domo
9	ExactTarget	Domo
10	BigQuery	Domo
11	Google Analytics	BigQuery
12	Google Analytics	Salesforce.com

+   Systems ▾   Data Feeds ▾







A flock of geese is flying in a clear blue sky. The geese are silhouetted against the sky, with some showing lighter underbellies. They are scattered across the frame, with one in the top right, another in the top center, one in the middle left, one in the middle center, one in the bottom center, and one in the bottom right. A dark blue horizontal band runs across the middle of the image, containing the text.

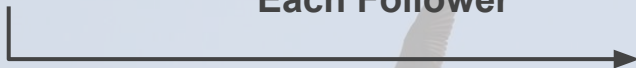
# Exploring Twitter Followers



**Pull a List of All  
Followers of  
@tgwilson**

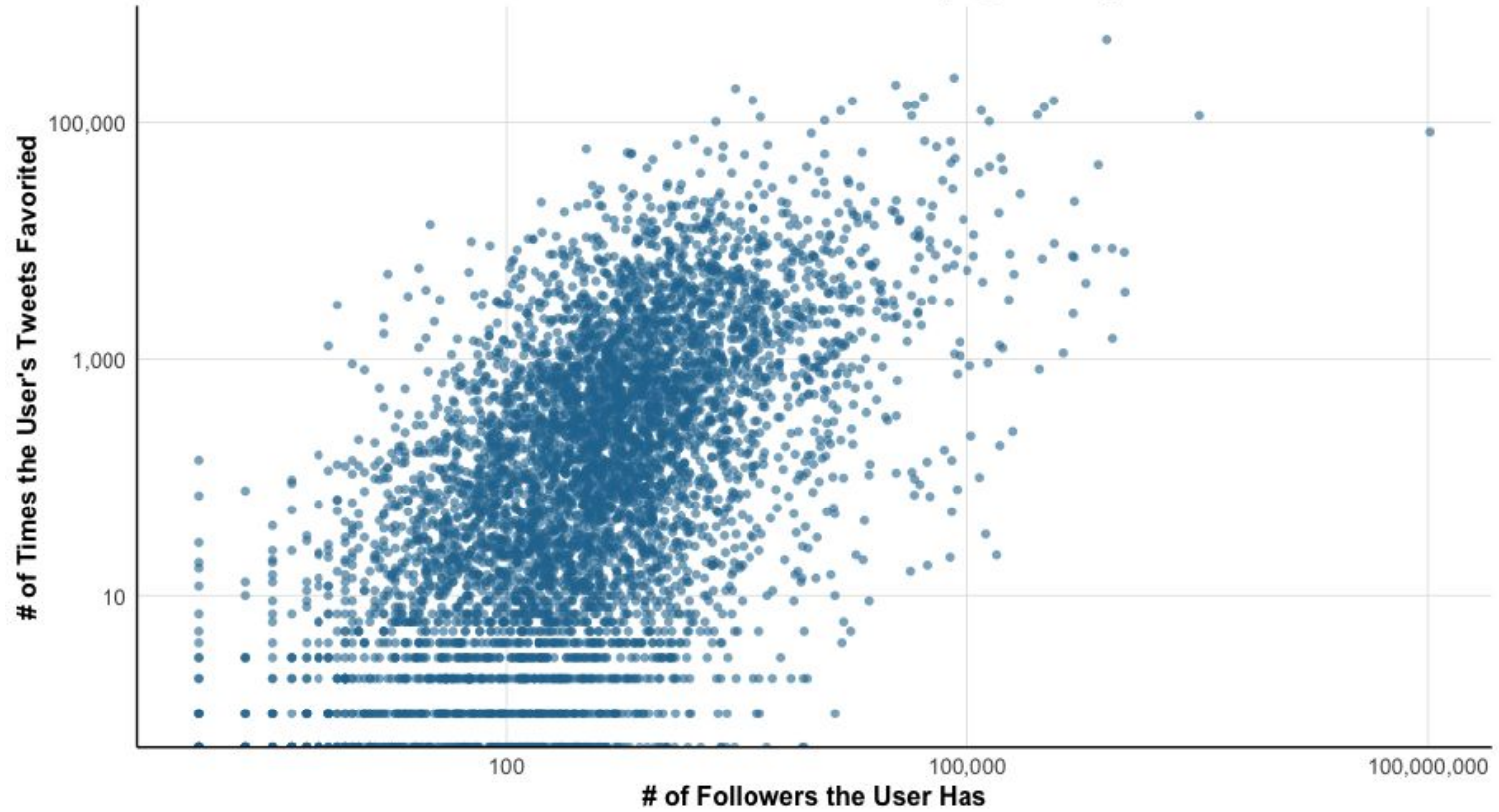


**Pull Basic  
Information About  
Each Follower**



User	Tweets	Flwrs	Faves

All Followers of the Account (Log Scales)







**MIX.MASTA.KING**  
@MixMastaKing Follows you

Actor/Director For New Movie Makaveli  
[#IMDB makavelithepacumentary.com](#)

📍 Los Angeles, CA  
[makavelithepacumentary.com](#)  
📺 View broadcasts  
📅 Joined January 2009

[Tweet to](#) [Message](#)



**Larry Kim** ✓  
@larrykim Follows you

CEO @MobileMonkey\_, Founder  
@WordStream (acquired for \$150M)  
Columnist @Inc, @Medium, @CNBC.  
Startups, AdWords, Chatbots.  
Popularized Unicorns in Marketing.

📍 Boston, MA  
[bit.ly/LarryKim-Linke...](#)  
📅 Joined December 2008

[Tweet to](#) [Message](#)



**Google Analytics** ✓  
@googleanalytics Follows you

Get the latest news and product updates  
on Google Analytics, Data Studio,  
Optimize, Surveys, and Tag Manager.  
Learn more at [g.co/marketingplatf...](#)

📍 Most places  
[g.co/marketingplatf...](#)  
📅 Joined June 2009

[Tweet to](#) [Message](#)







**MIX.MASTA.KING**  
@MixMastaKing Follows you

Actor/Director For New Movie Makaveli  
[#IMDB makavelithethepacumentary.com](#)

📍 Los Angeles, CA

[makavelithethepacumentary.com](#)

📺 View broadcasts

📅 Joined January 2009

[Tweet to](#) [Message](#)



**Larry Kim** ✓  
@larrykim Follows you

CEO @MobileMonkey\_, Founder  
@WordStream (acquired for \$150M)  
Columnist @Inc, @Medium, @CNBC.  
Startups, AdWords, Chatbots.  
Popularized Unicorns in Marketing.

📍 Boston, MA

[bit.ly/LarryKim-Linke...](#)

📅 Joined December 2008

[Tweet to](#) [Message](#)



**Google Analytics** ✓  
@googleanalytics Follows you

Get the latest news and product updates  
on Google Analytics, Data Studio,  
Optimize, Surveys, and Tag Manager.  
Learn more at [g.co/marketingplatf...](#)

📍 Most places

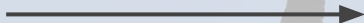
[g.co/marketingplatf...](#)

📅 Joined June 2009

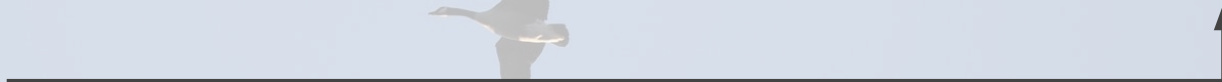
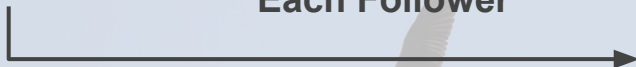
[Tweet to](#) [Message](#)



**Pull a List of All  
Followers of  
@tgwilson**



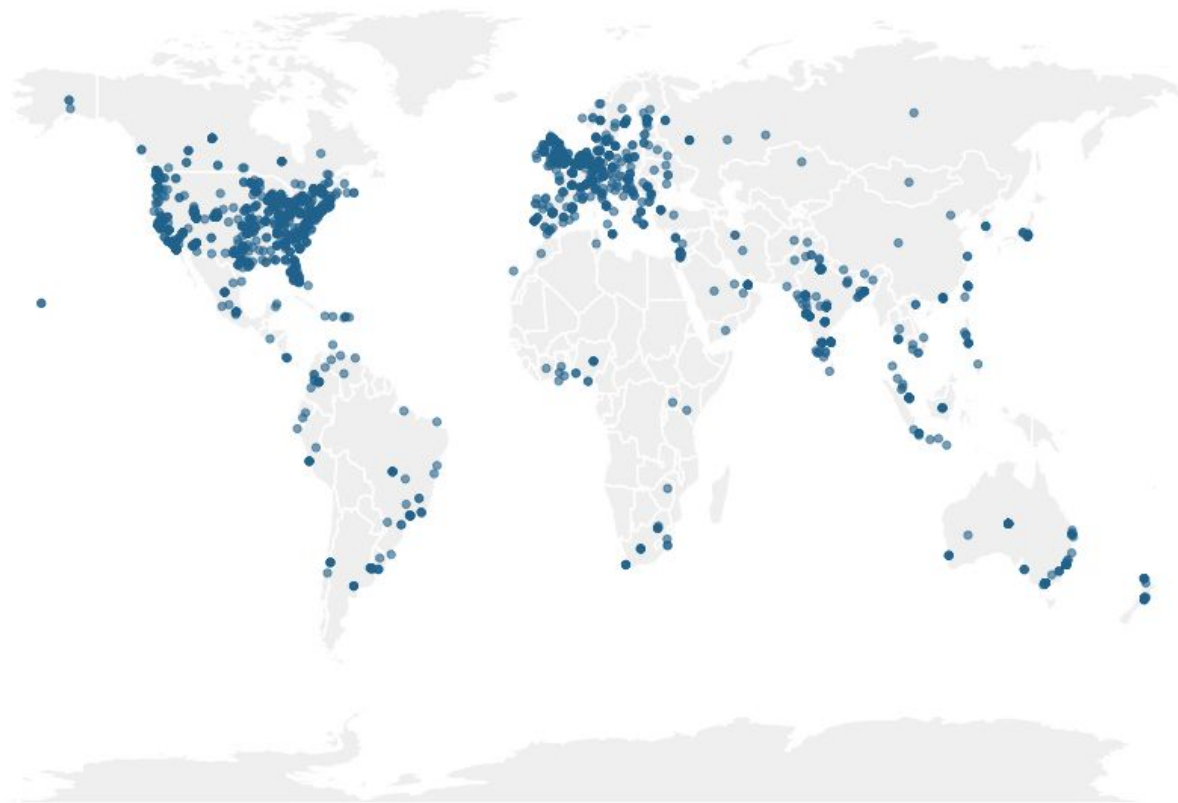
**Pull Basic  
Information About  
Each Follower**

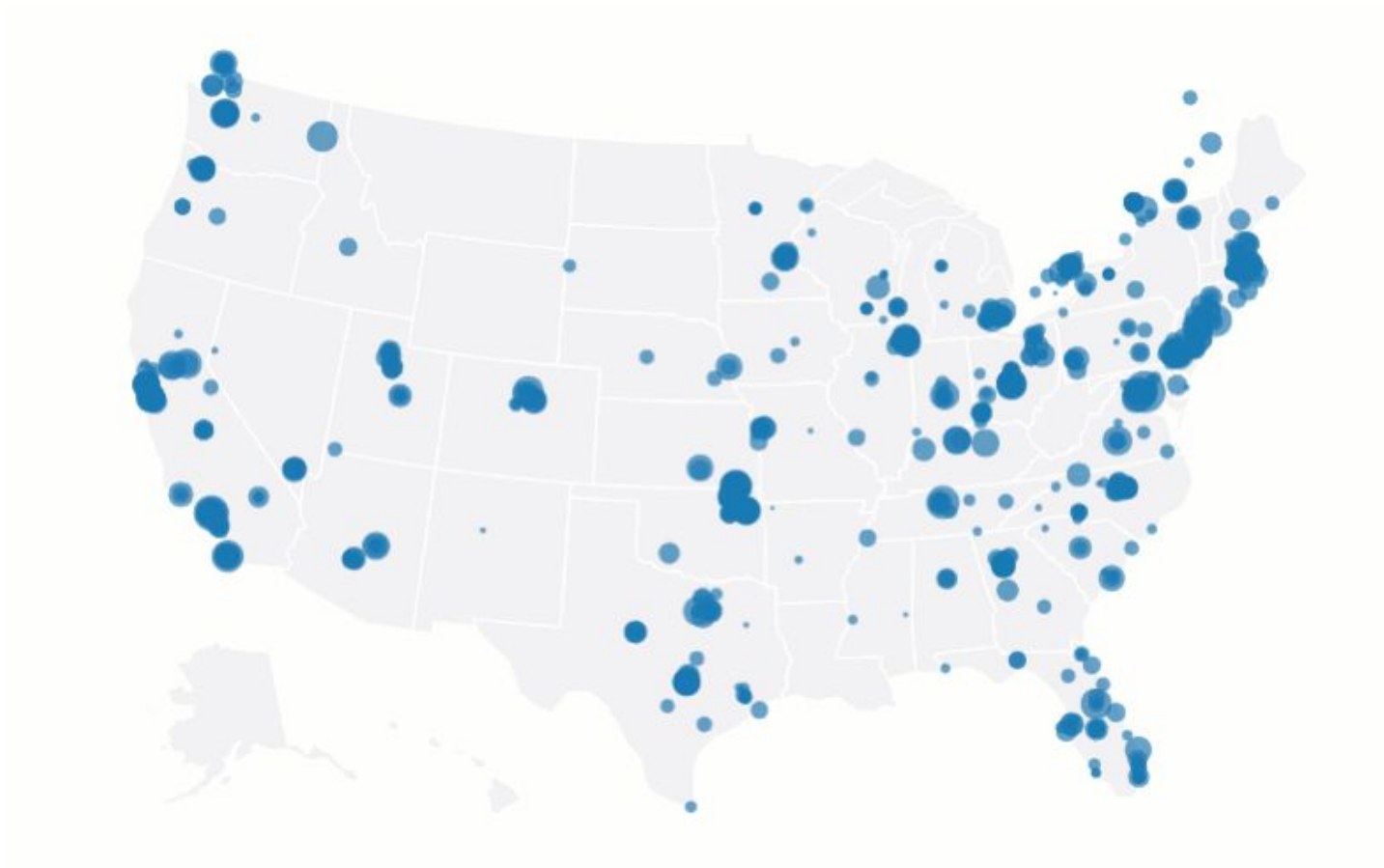


User	Tweets	Flwrs	Faves	Lat	Long



### Top Followers of the Account (Based on Followers, When Location Available)



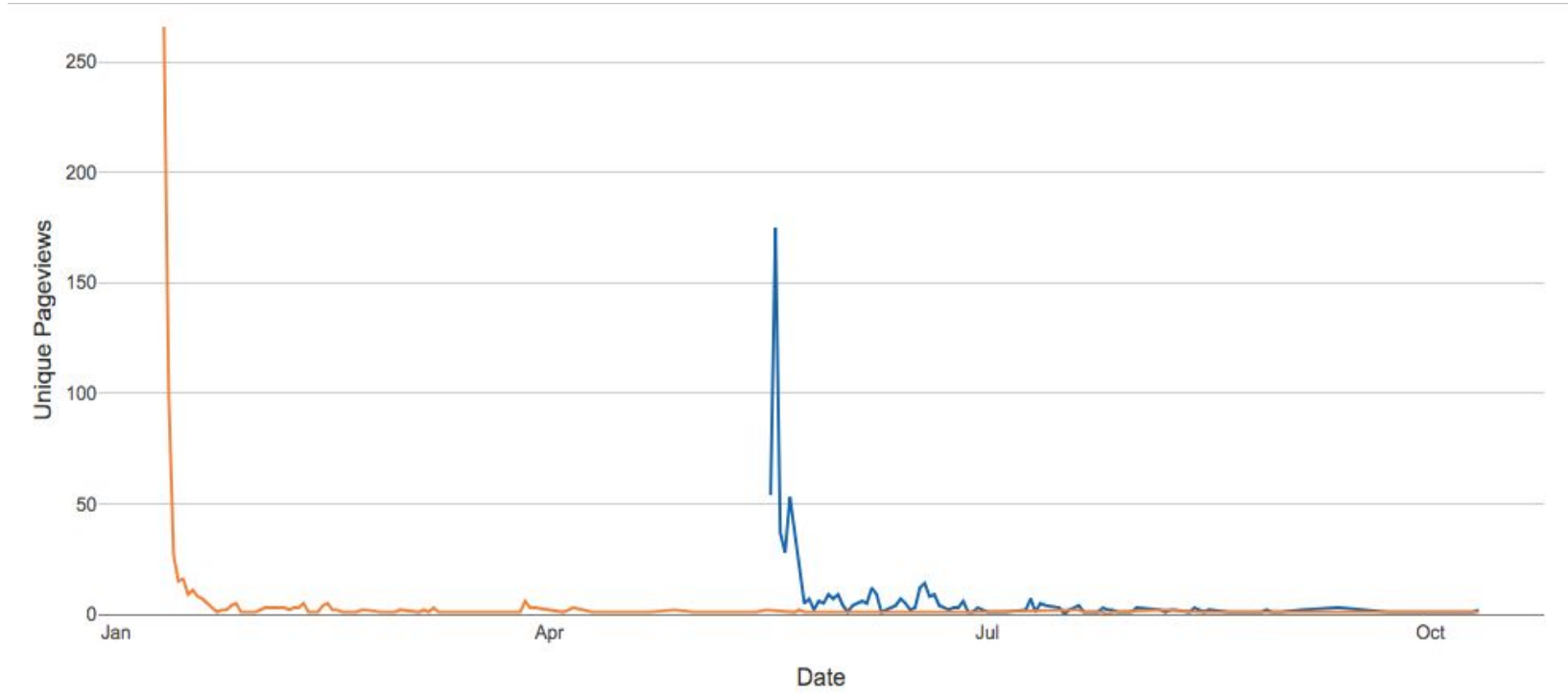




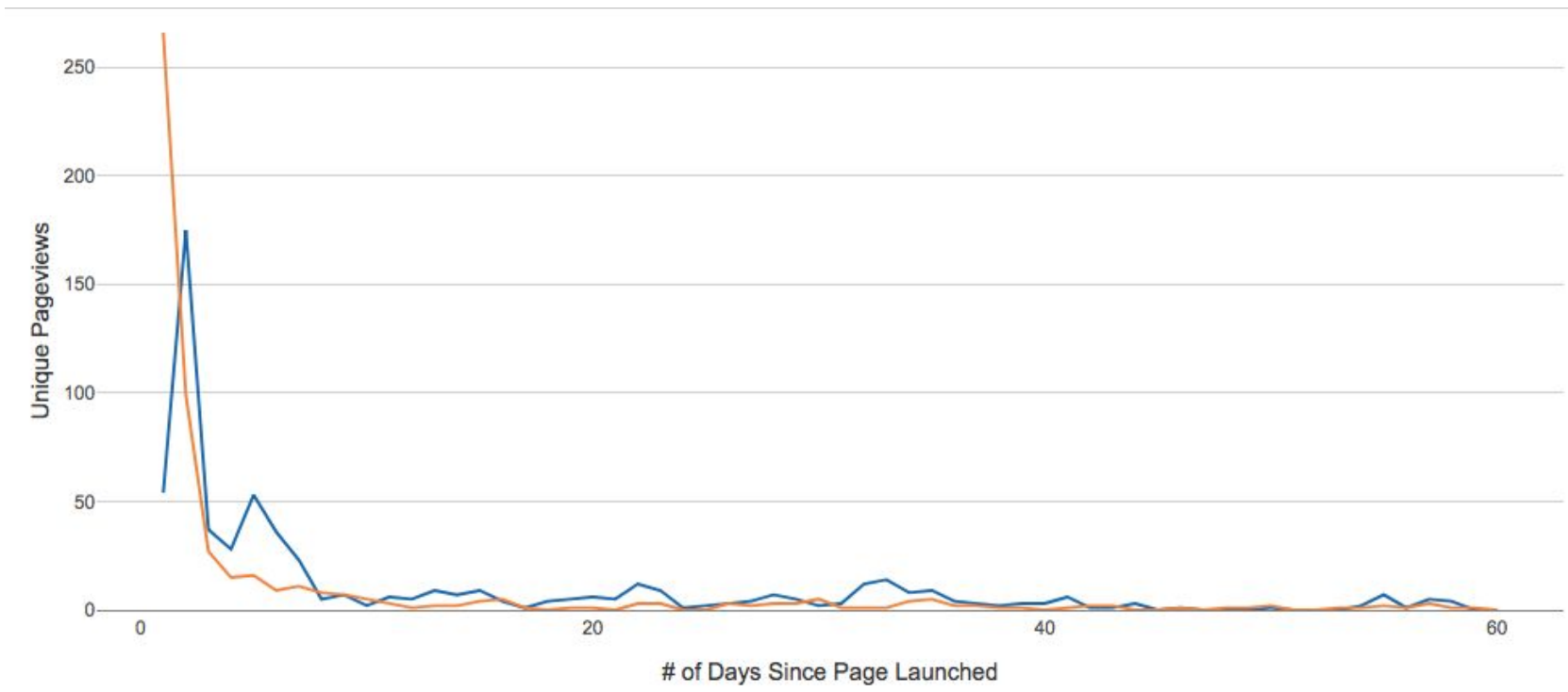
**Page Traffic: Launch vs. Lifecycle**



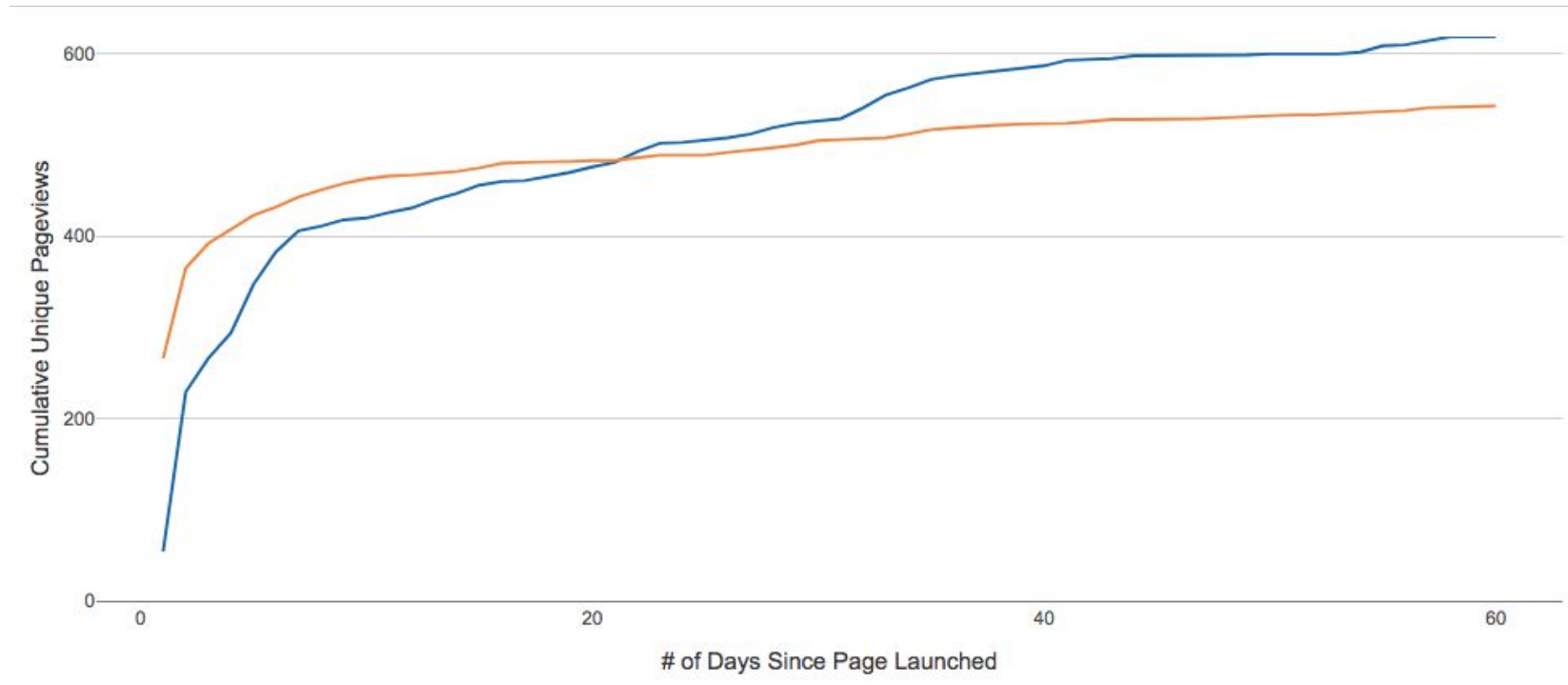
# Consider Two Blog Posts



# We Can “Time-Normalize” Them to Launch Date

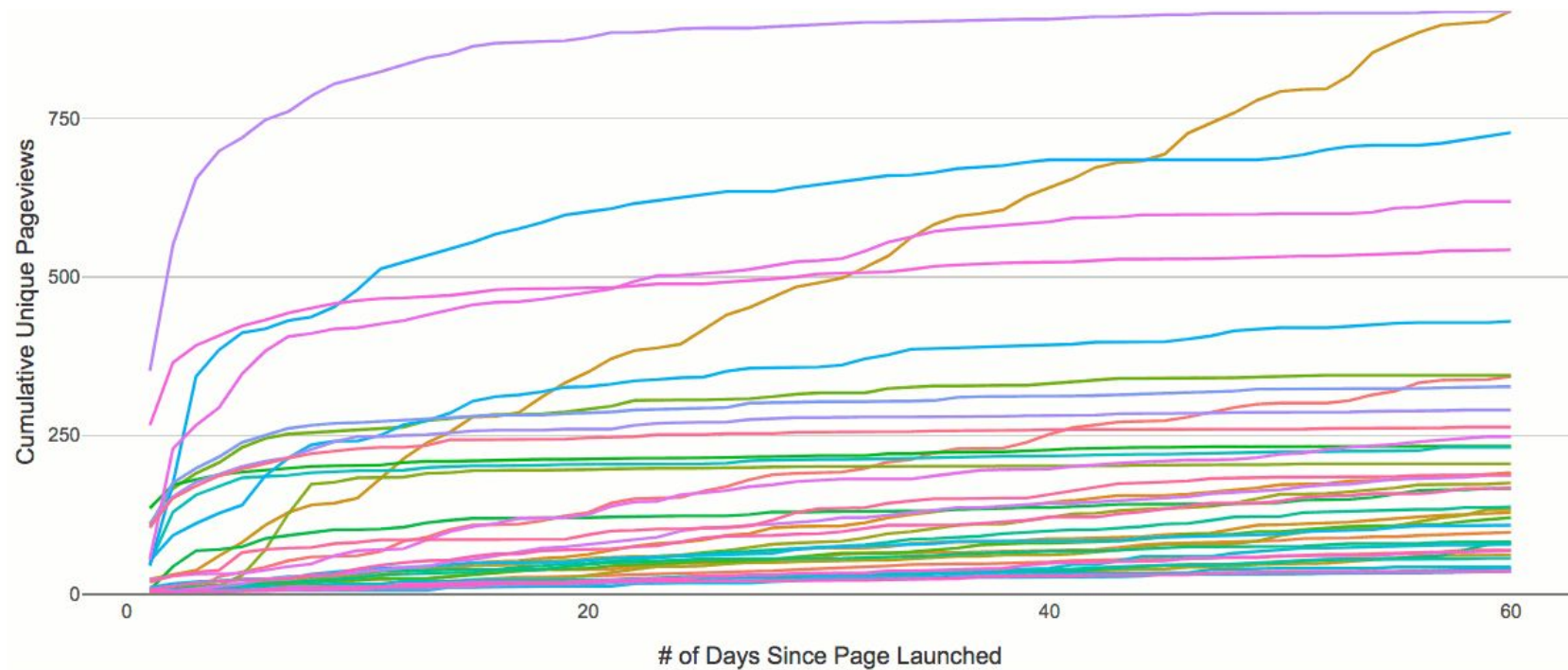


# We Can “Time-Normalize” Them to Launch Date





# We Can “Time-Normalize” Them to Launch Date



# So...should I learn



or





So...should I learn

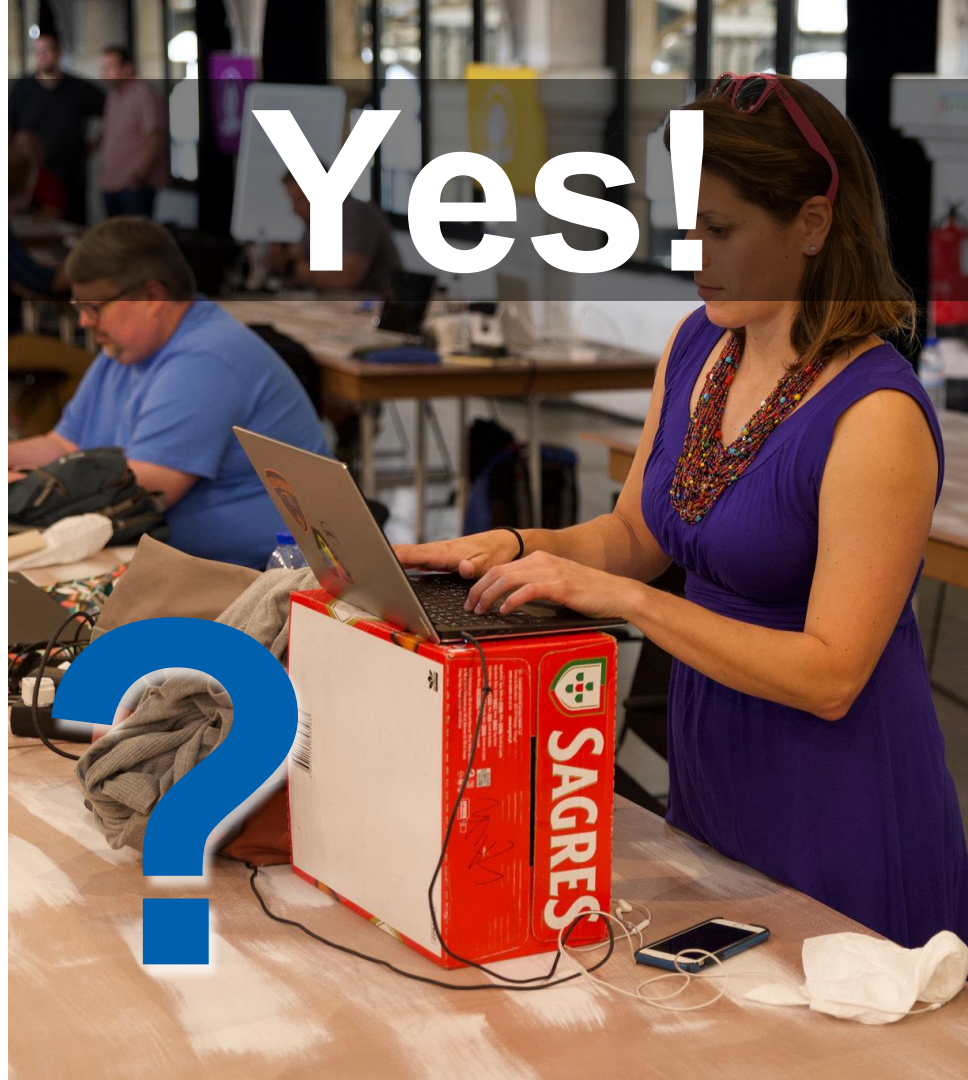


or



python™

Yes!



So...should I learn



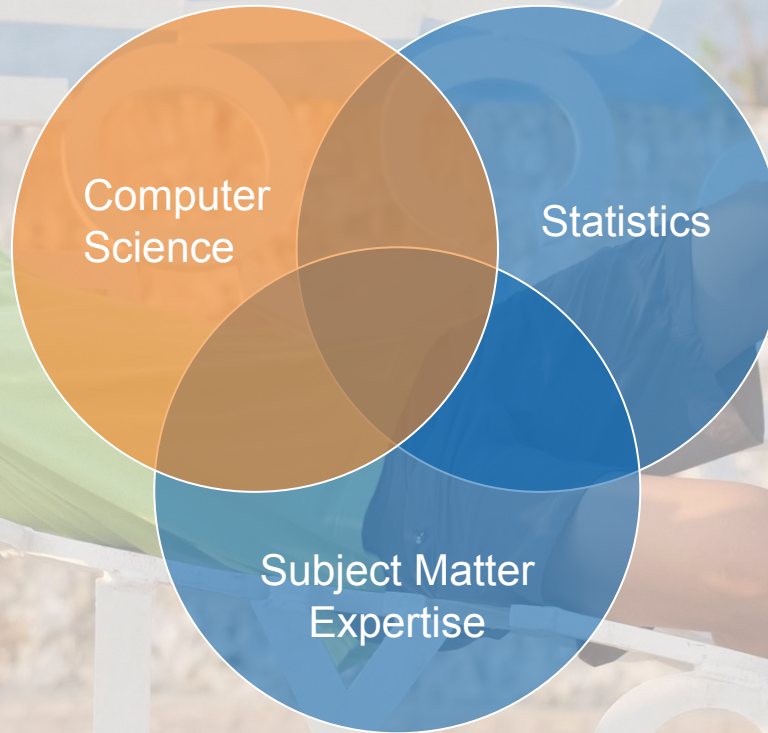
or



TBD!

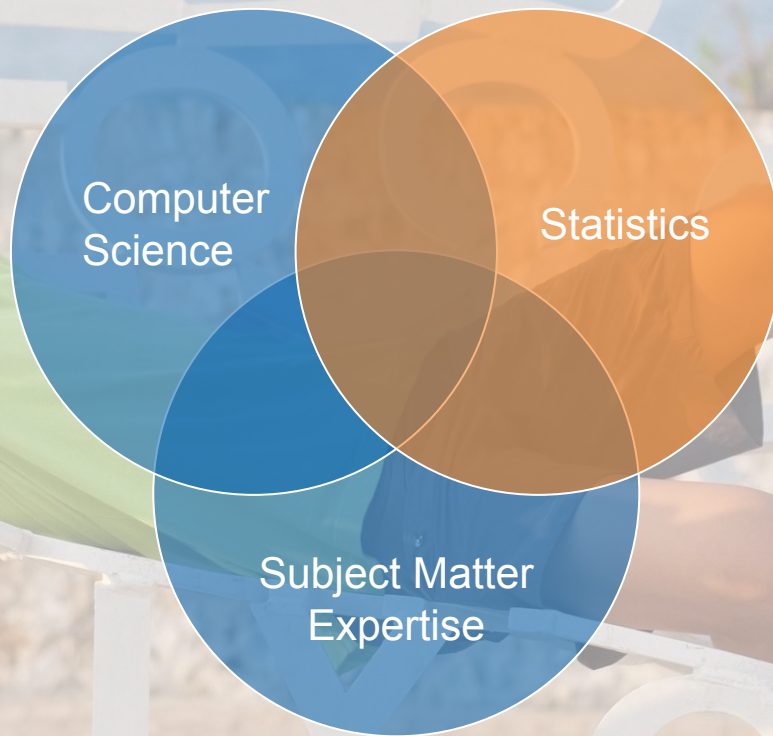


# So...That Was Computer Science!





# I'm 100% Confident We Should Talk Statistics



“We have a problem  
in America with  
thinking  
probabilistically.”

- Annie Duke



Source: Annie Duke



“We have a problem  
in [marketing] with  
thinking  
probabilistically.”

- Annie Duke



Source: Annie Duke





How much traffic  
came to the site last  
month?

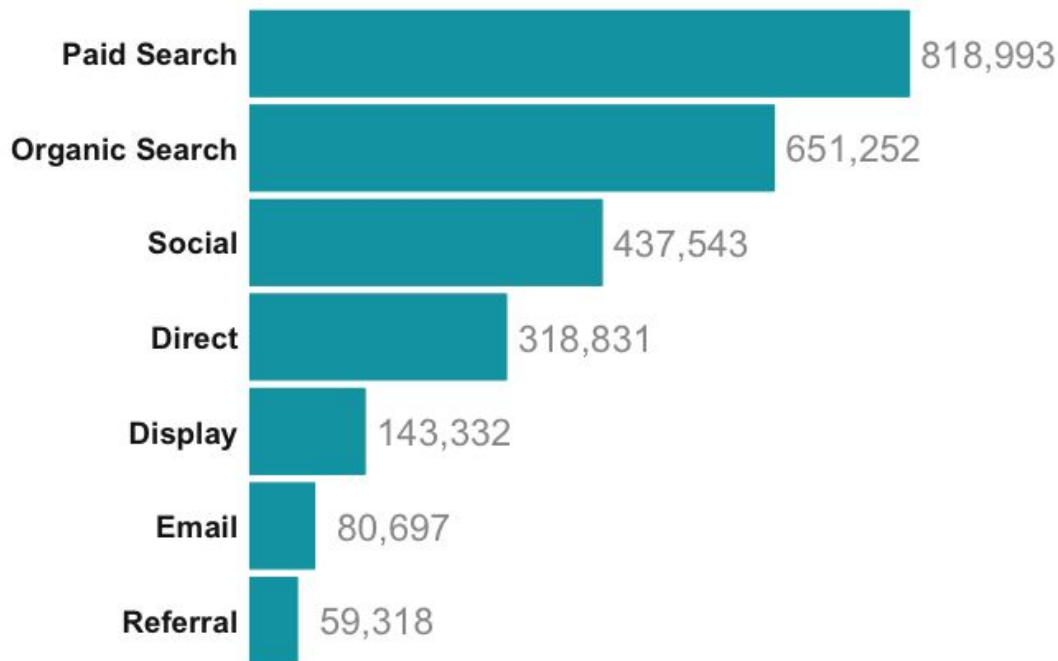
2,509,966  
*sessions*

# Trend it!

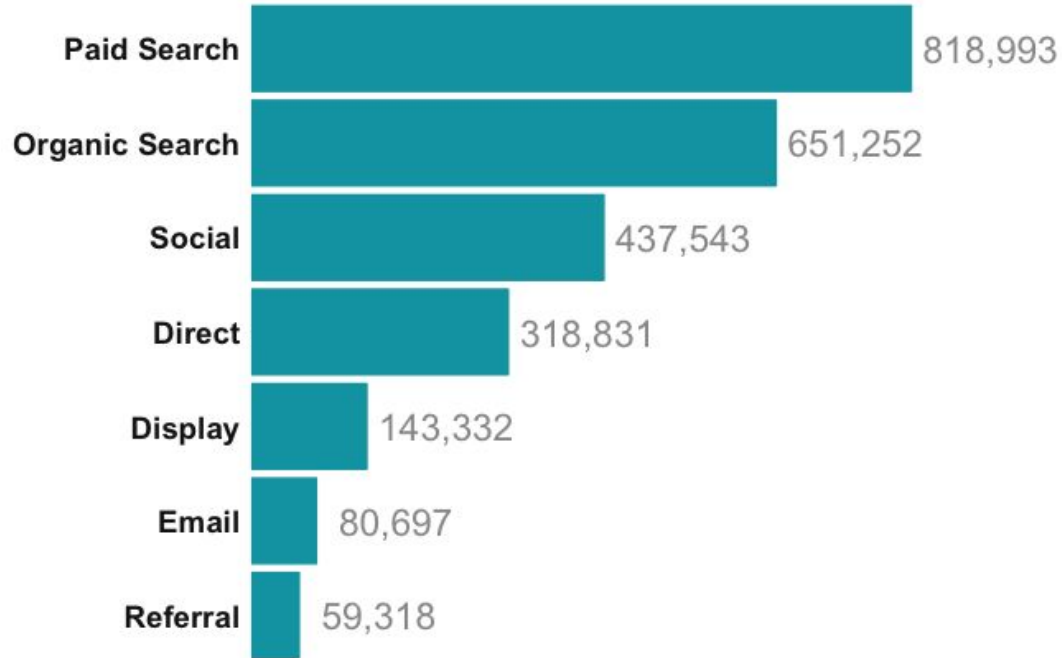




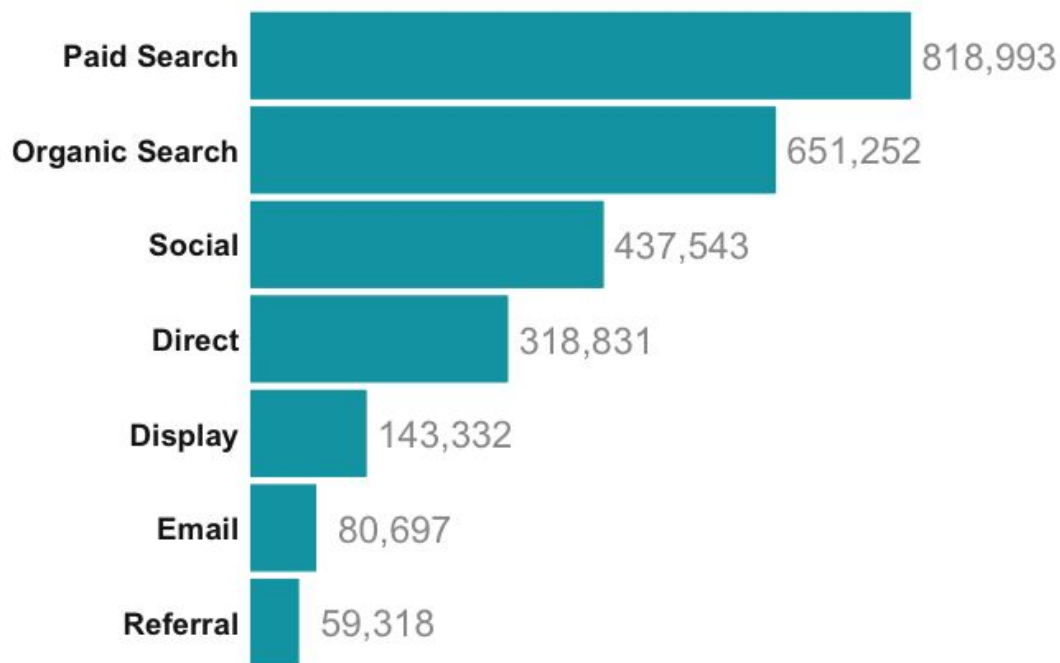
# Or break it down by channel!



# Let's ponder this breakdown a bit.



# How many rows of data went into this table?

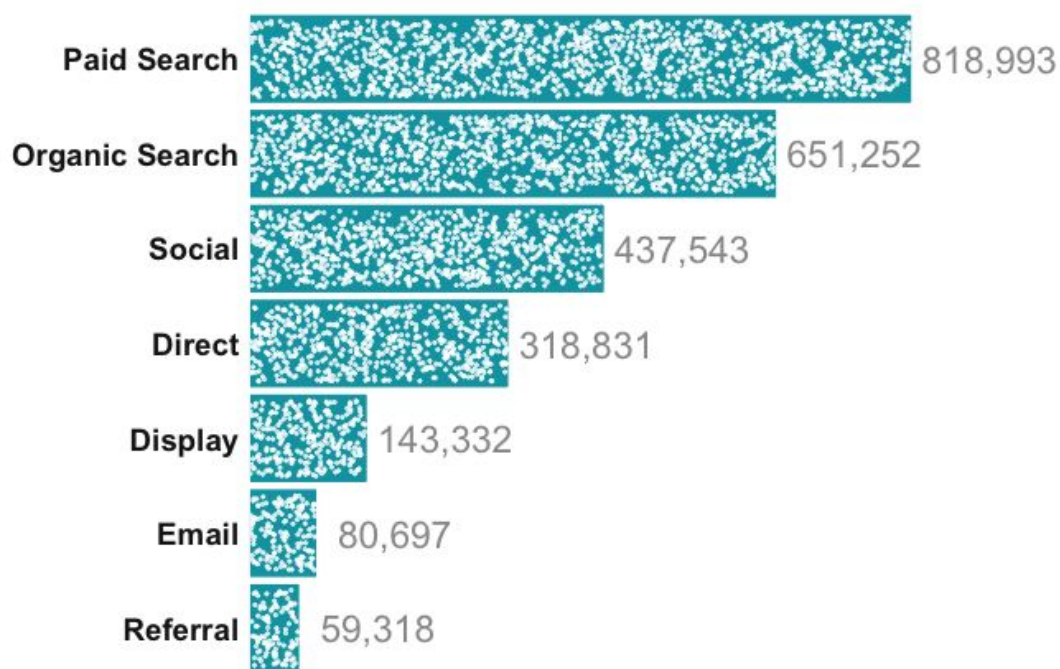


Channel	Sessions
Paid Search	818,993
Organic Search	651,252
Social	437,543
Direct	318,831
Display	143,332
Email	80,697
Referral	59,318

7!



# How would a data scientist answer the question?



Visit ID	Channel
7521516054.5930117975	Organic Search
6639584132.5334606522	Organic Search
2072704531.5757008375	Organic Search
3700666076.4855492580	Social
3901278585.2803798674	Social
9493011014.8450893410	Referral
4257955221.2341693619	Organic Search
5653594021.2053548613	Paid Search
2103008782.3786630871	Organic Search
2156814867.9852847980	Direct
3969242874.8434979641	Paid Search
8975736469.8764719209	Social
3127960728.4142998058	Social
3863953834.9545306506	Organic Search
1769847183.6232068180	Organic Search

2.5  
million

# This shift also helps with some basic stats concepts

## What is the variable?

## What are the values of the variable?

Channel	Sessions
Paid Search	818,993
Organic Search	651,252
Social	437,543
Direct	318,831
Display	143,332
Email	80,697
Referral	59,318

?

Visit ID	Channel
7521516054.5930117975	Organic Search
6639584132.5334606522	Organic Search
2072704531.5757008375	Organic Search
3700666076.4855492580	Social
3901278585.2803798674	Social
9493011014.8450893410	Referral
4257955221.2341693619	Organic Search
5653594021.2053548613	Paid Search
2103008782.3786630871	Organic Search
2156814867.9852847980	Direct
3969242874.8434979641	Paid Search
8975736469.8764719209	Social
3127960728.4142998058	Social
3863953834.9545306506	Organic Search
1769847183.6232068180	Organic Search

# This shift also helps with some basic stats concepts

## What is the variable?

## What are the values of the variable?

Channel	Sessions
Paid Search	818,993
Organic Search	651,252
Social	437,543
Direct	318,831
Display	143,332
Email	80,697
Referral	59,318

?

Visit ID	Channel
7521516054.5930117975	Organic Search
6639584132.5334606522	Organic Search
2072704531.5757008375	Organic Search
3700666076.4855492580	Social
3901278585.2803798674	Social
9493011014.8450893410	Referral
4257955221.2341693619	Organic Search
5653594021.2053548613	Paid Search
2103008782.3786630871	Organic Search
2156814867.9852847980	Direct
3969242874.8434979641	Paid Search
8975736469.8764719209	Social
3127960728.4142998058	Social
3863953834.9545306506	Organic Search
1769847183.6232068180	Organic Search



# This shift also helps with some basic stats concepts

## What is the variable?

## What are the values of the variable?

Channel	Sessions
Paid Search	818,993
Organic Search	651,252
Social	437,543
Direct	318,831
Display	143,332
Email	80,697
Referral	59,318

?

Visit ID	Channel
7521516054.5930117975	Organic Search
6639584132.5334606522	Organic Search
2072704531.5757008375	Organic Search
3700666076.4855492580	Social
3901278585.2803798674	Social
9493011014.8450893410	Referral
4257955221.2341693619	Organic Search
5653594021.2053548613	Paid Search
2103008782.3786630871	Organic Search
2156814867.9852847980	Direct
3969242874.8434979641	Paid Search
8975736469.8764719209	Social
3127960728.4142998058	Social
3863953834.9545306506	Organic Search
1769847183.6232068180	Organic Search

# This shift also helps with some basic stats concepts

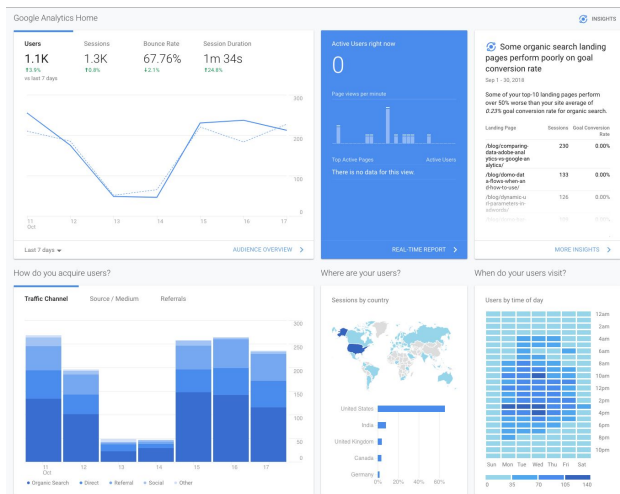
What is the variable?

What are the values of the variable?

Channel	Sessions
Paid Search	818,993
Organic Search	651,252
Social	437,543
Direct	318,831
Display	143,332
Email	80,697
Referral	59,318

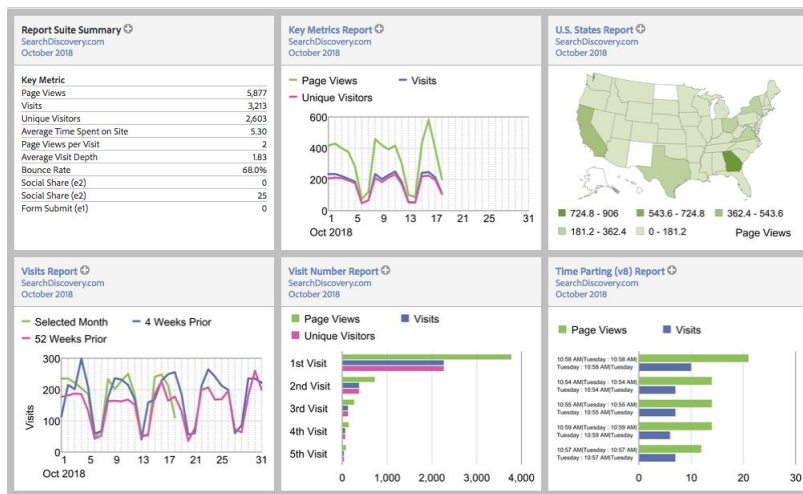
Visit ID	Channel
7521516054.5930117975	Organic Search
6639584132.5334606522	Organic Search
2072704531.5757008375	Organic Search
3700666076.4855492580	Social
3901278585.2803798674	Social
9493011014.8450893410	Referral
4257955221.2341693619	Organic Search
5653594021.2053548613	Paid Search
2103008782.3786630871	Organic Search
2156814867.9852847980	Direct
3969242874.8434979641	Paid Search
8975736469.8764719209	Social
3127960728.4142998058	Social
3863953834.9545306506	Organic Search
1769847183.6232068180	Organic Search

# We are used to working with aggregated data!



+

Google Data Studio



+

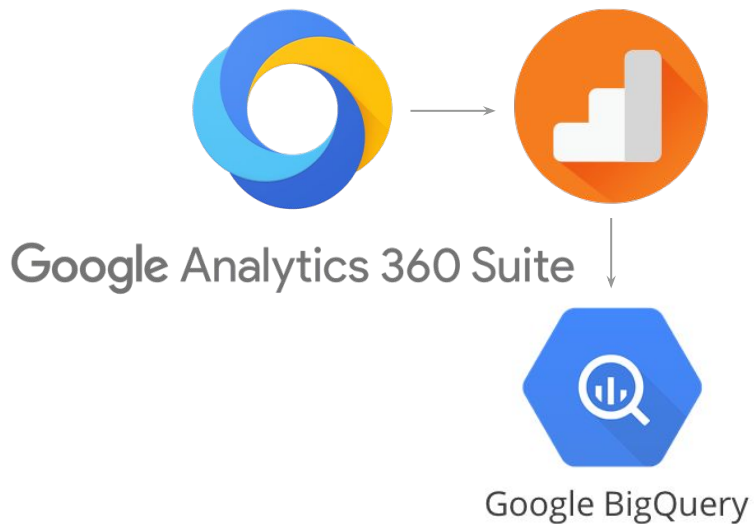
Analysis Workspace

+

Reporting APIs for Either Platform



# Google and Adobe know this!



Data Feed



Should we  
abandon all  
hope if we  
don't have  
session-level  
detail?

No!





A red bicycle with black handlebars and wheels is leaning against a weathered wooden fence. On the fence, there is a large, colorful mural of the Pokémon Pikachu, which is yellow with red cheeks and a black outline. To the right of the Pikachu mural, there is graffiti that reads "Catch me if you can" in black and white. Other graffiti includes "F.T." in red and "OTT" in pink. The fence is made of vertical wooden planks and is surrounded by green grass and weeds at the base. A semi-transparent black banner with white text is overlaid across the middle of the image.

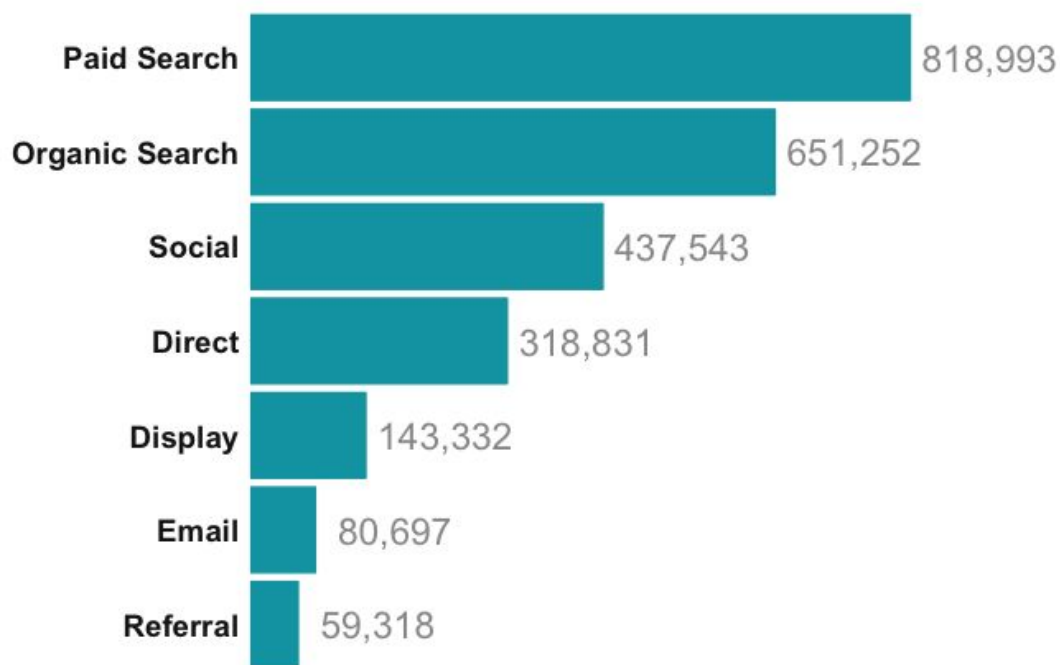
We just have to (carefully) cheat a little bit.



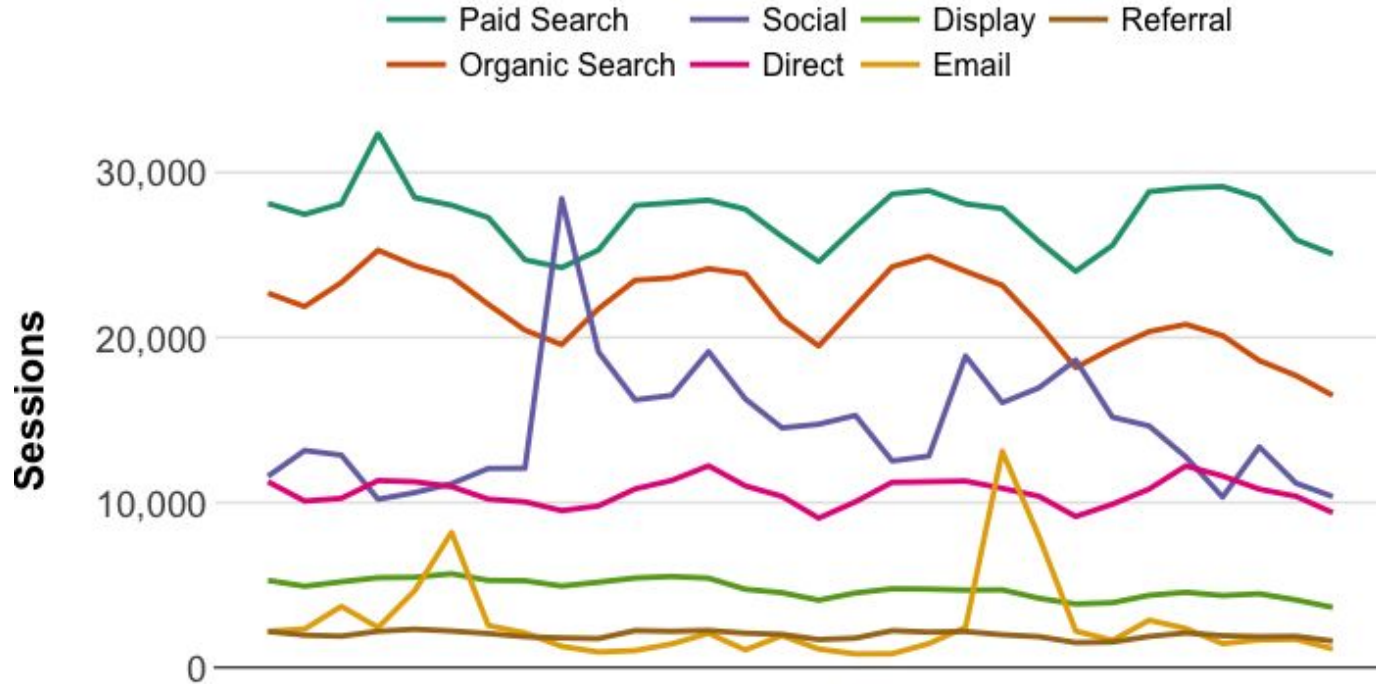


Let's ask a  
different  
question.

# Do sessions **really** differ by channel?



# But if we eyeball it by day, things get murkier.





# Deaggregation

(if not ideal detail)

30 Days

x 7 channels

---

**210 observations**





A red bicycle with black handlebars and a black seat is leaning against a weathered wooden fence. The fence is covered in graffiti, including a large yellow Pikachu character with its mouth open, and the words "Catch me if you can" in black. Other graffiti includes "F.T." in red and "OTT" in pink. The ground is covered in green grass and weeds.

Here's where we're going to cheat a little bit.



A red bicycle with black handlebars and wheels is leaning against a weathered wooden fence. The fence is covered in graffiti, including a large yellow Pikachu character on the left and a speech bubble in the center that says "Catch me if you can". Other graffiti includes "F.T." and "OTT" in red and pink. The scene is outdoors with green grass and weeds at the base of the fence.

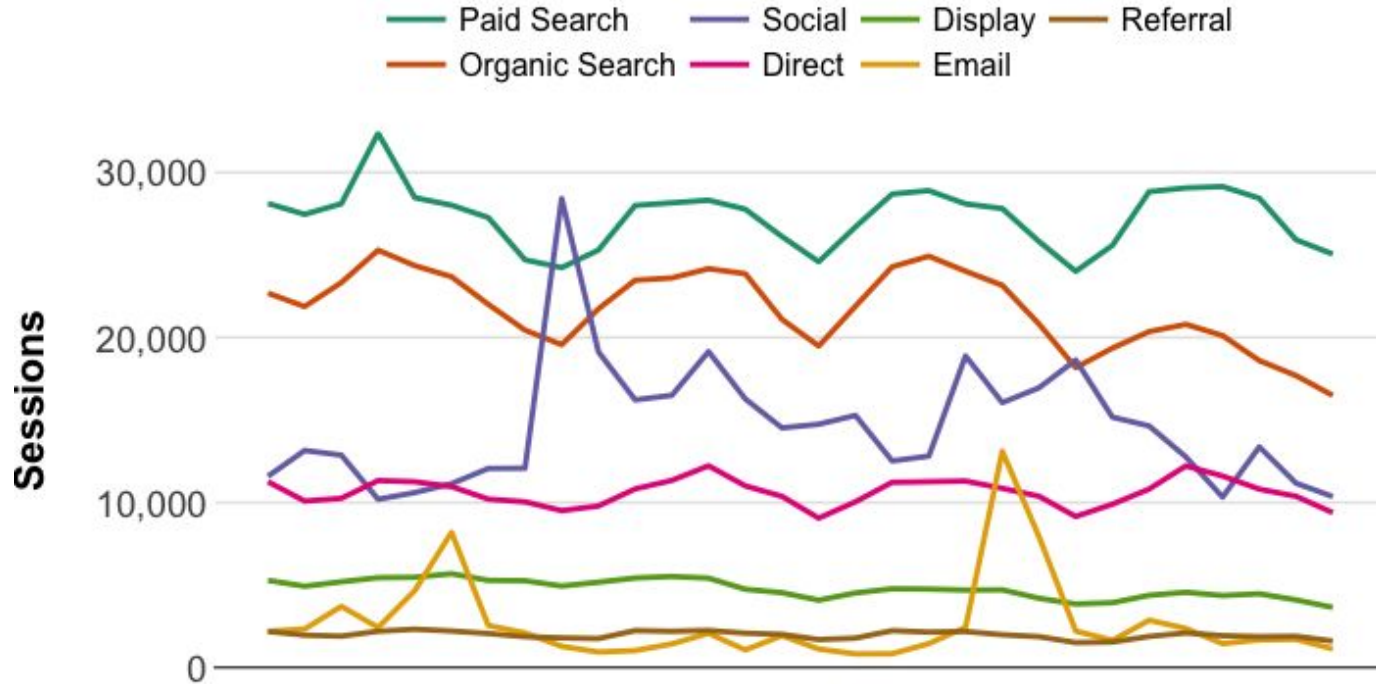
We're going to use **day**...





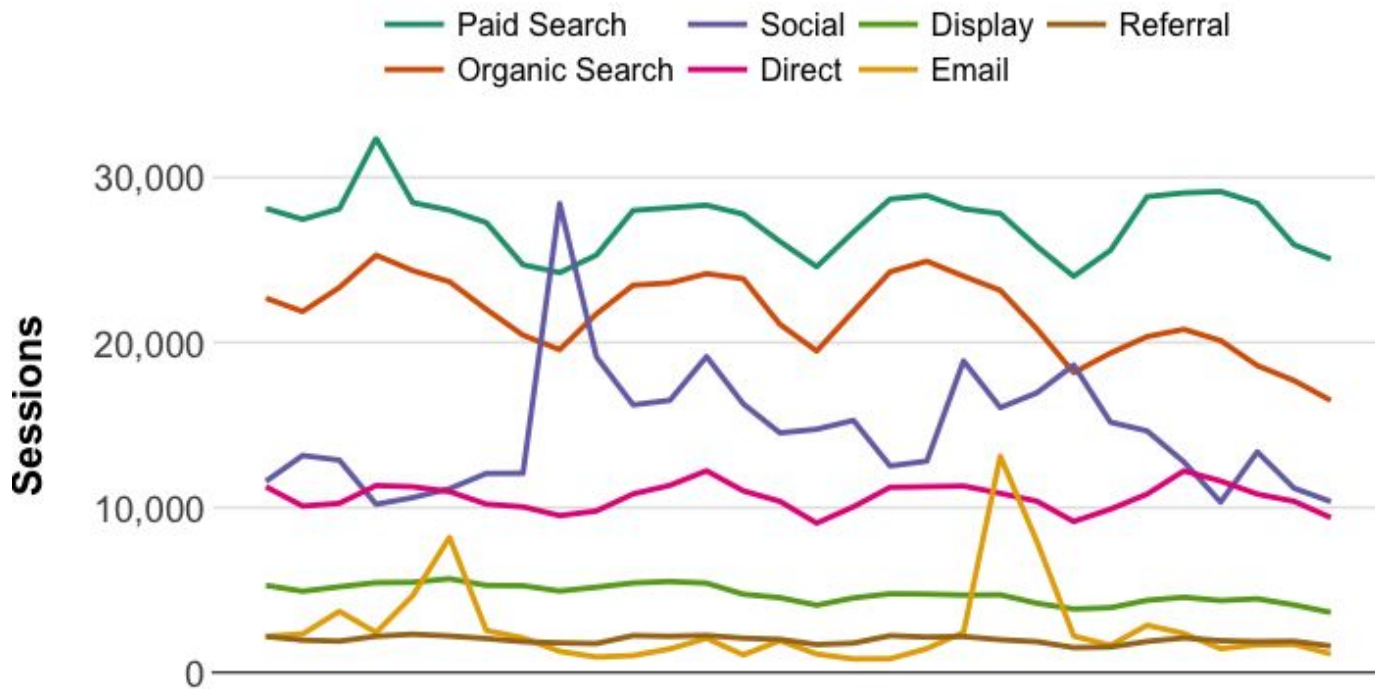
...as the **observational unit** for our analysis.

“Day” seems like an okay way to deaggregate.



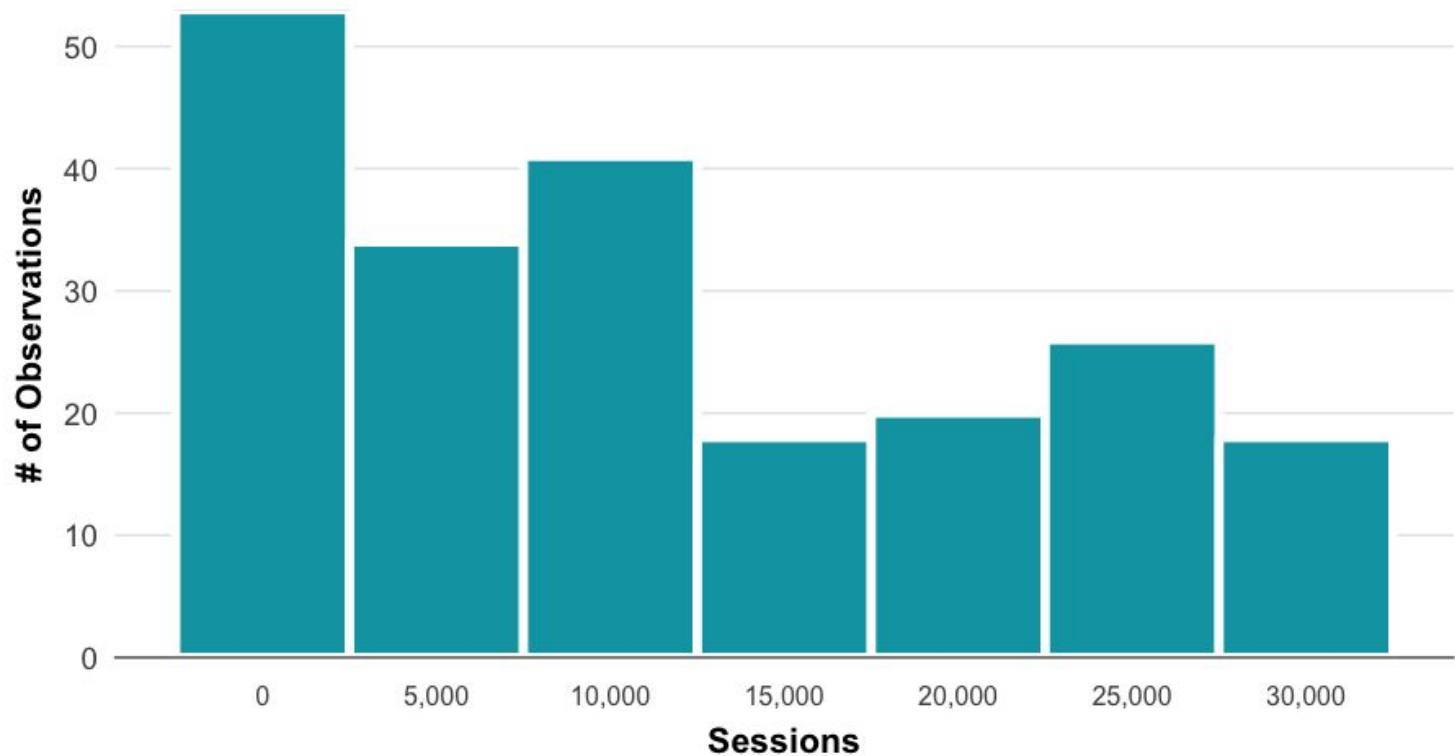


Let's ignore the **time aspect of date.**

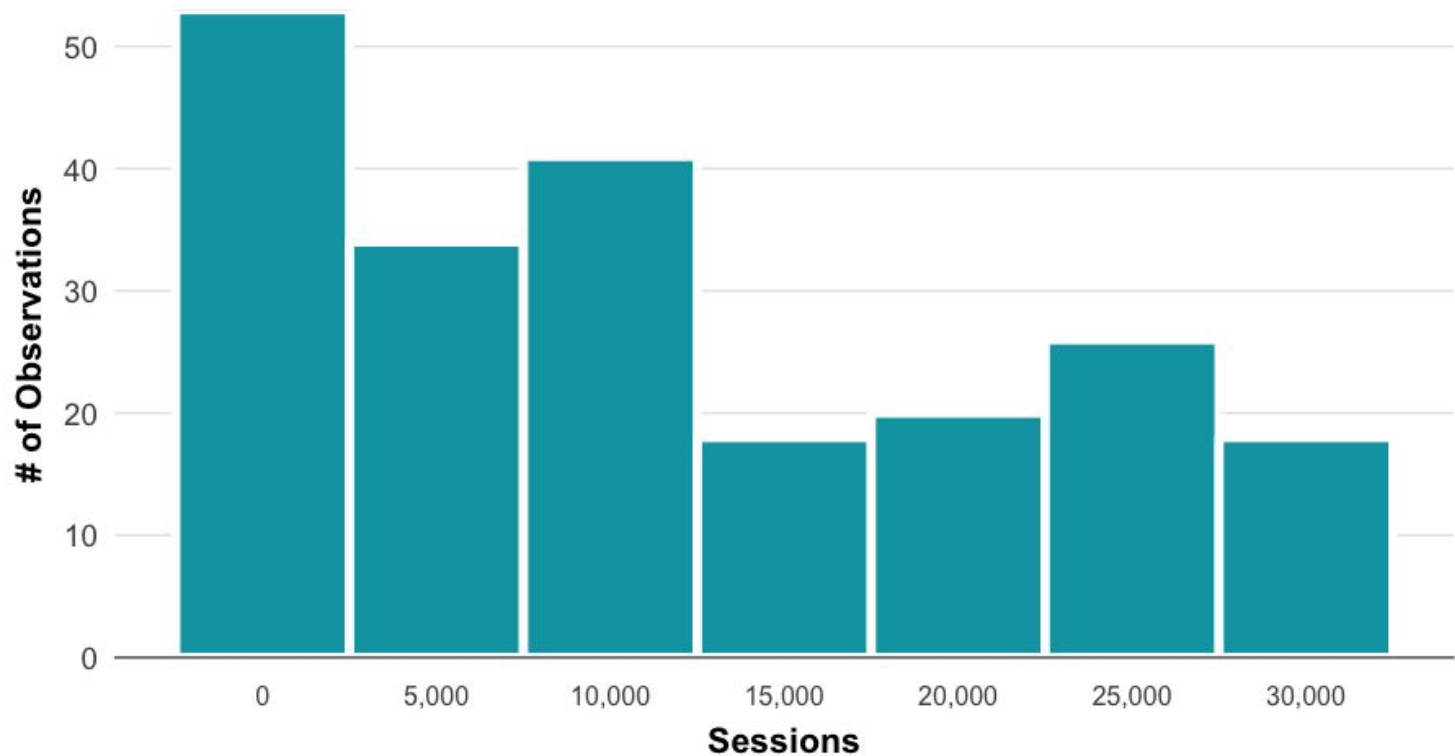




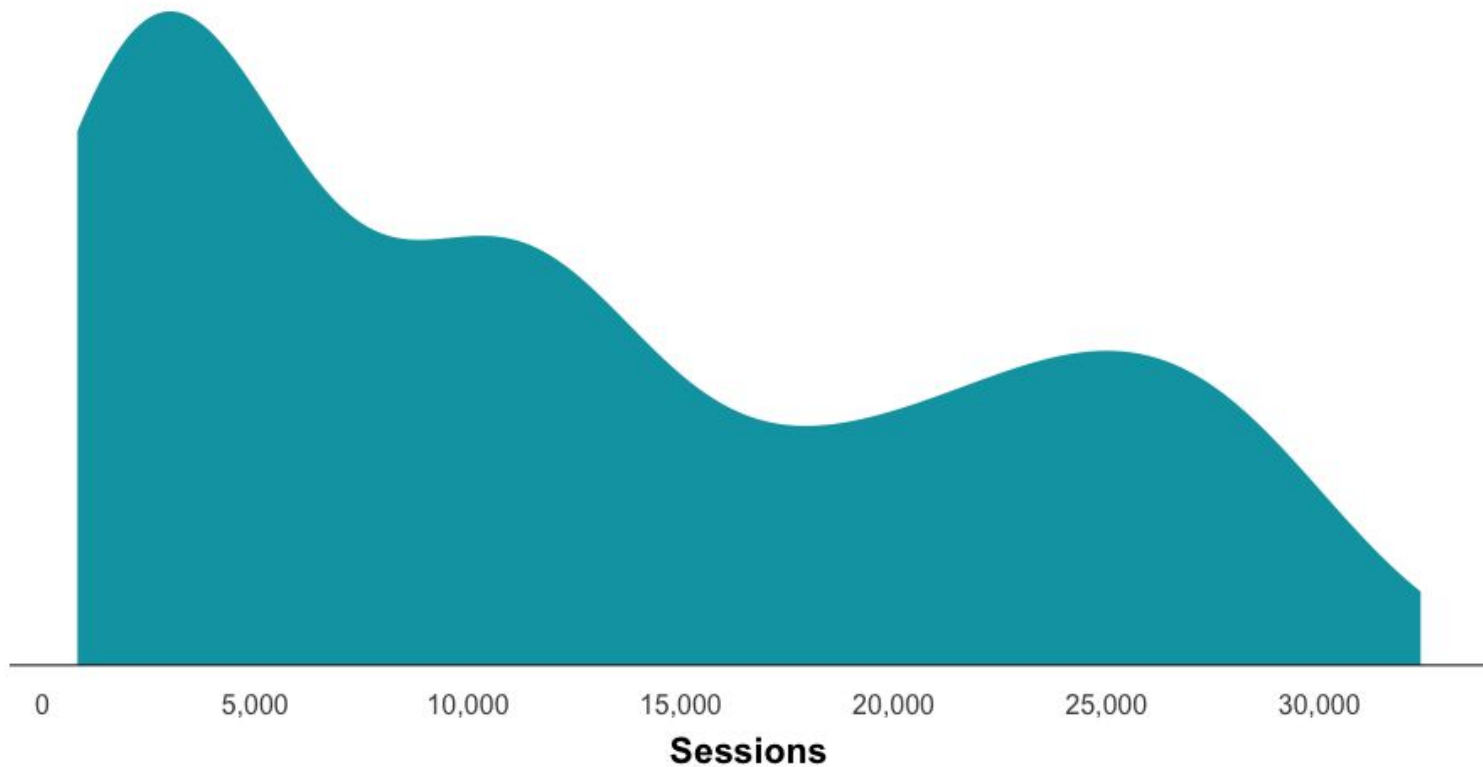
We can look at the data as a histogram (n = 210).



It's the **distribution** of the **observations**.

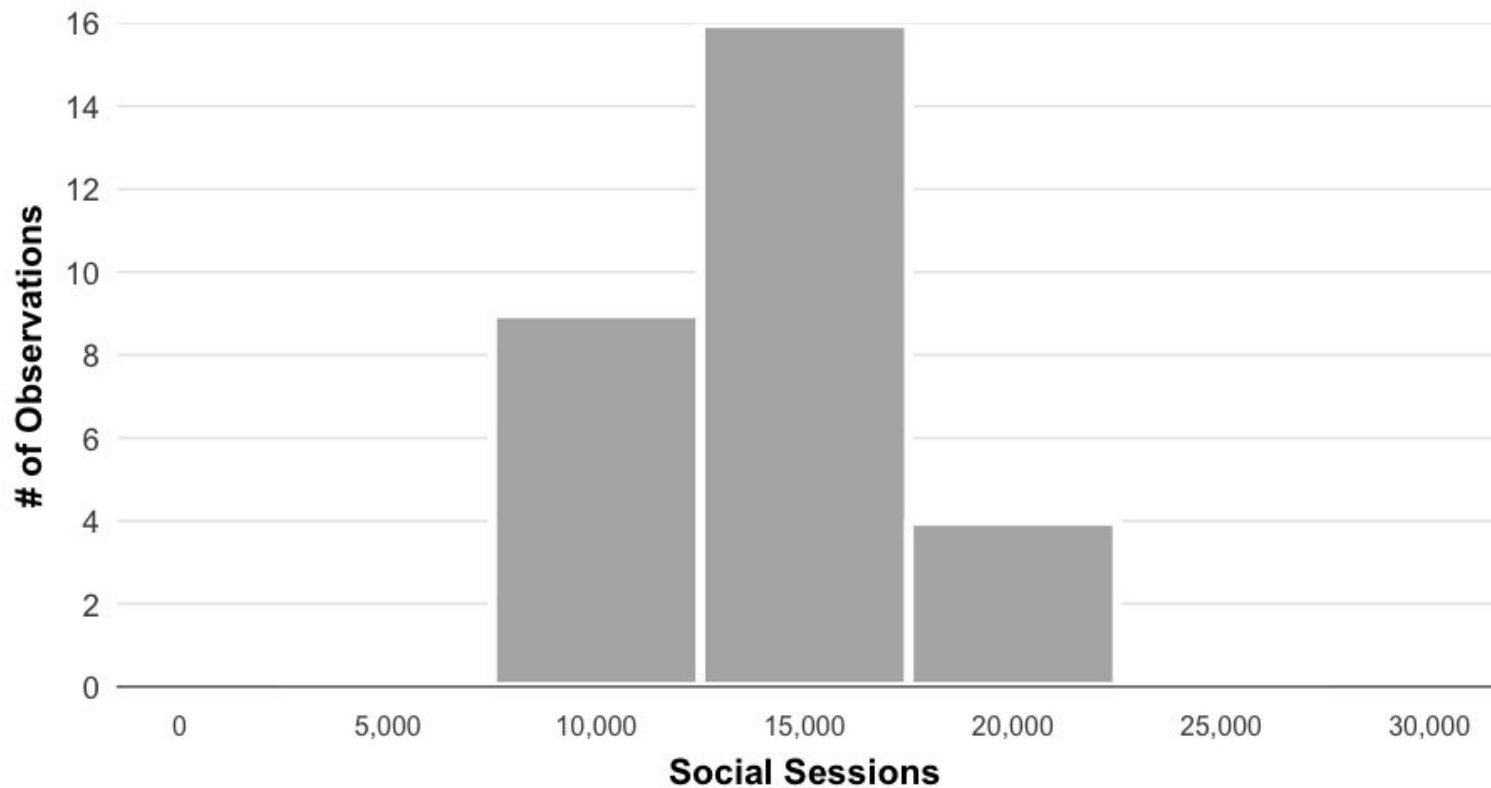


A **density plot** provides a smooth distribution.

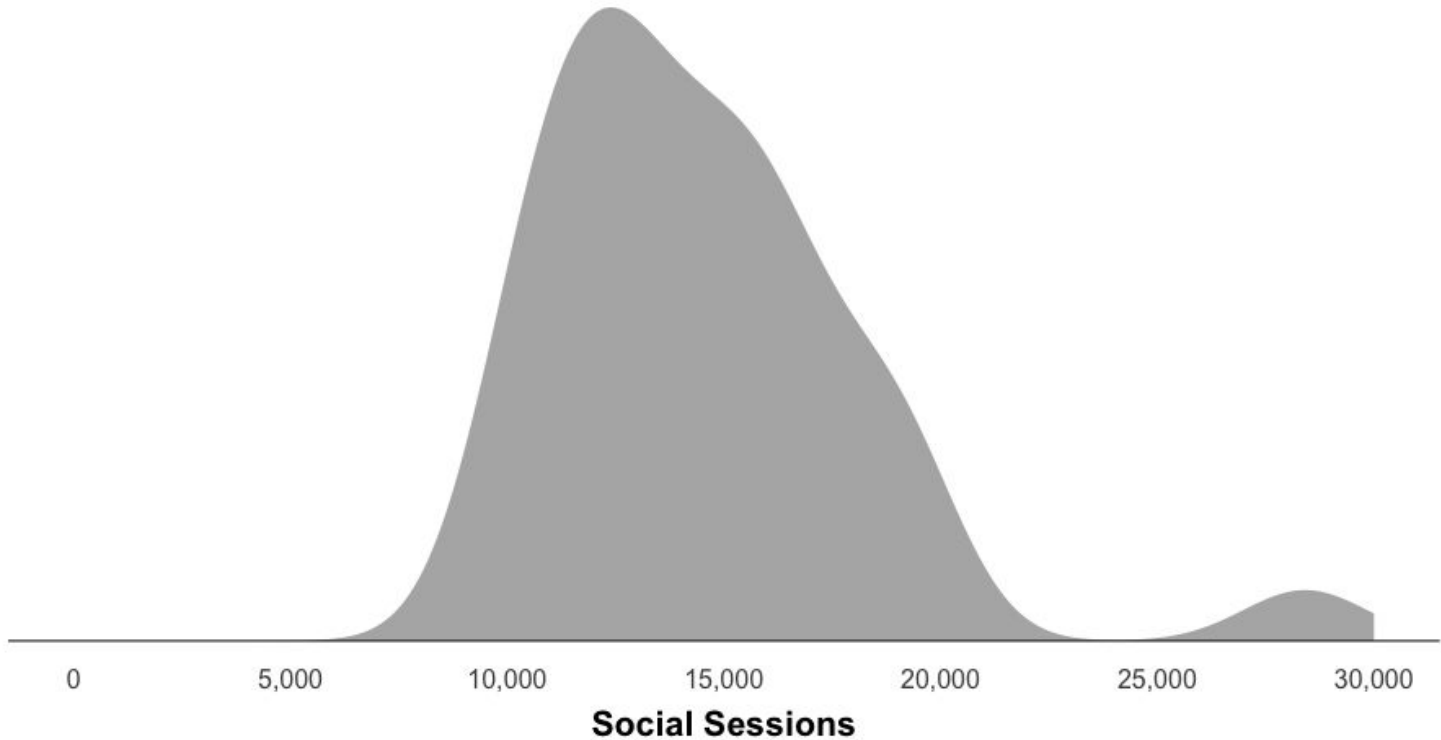




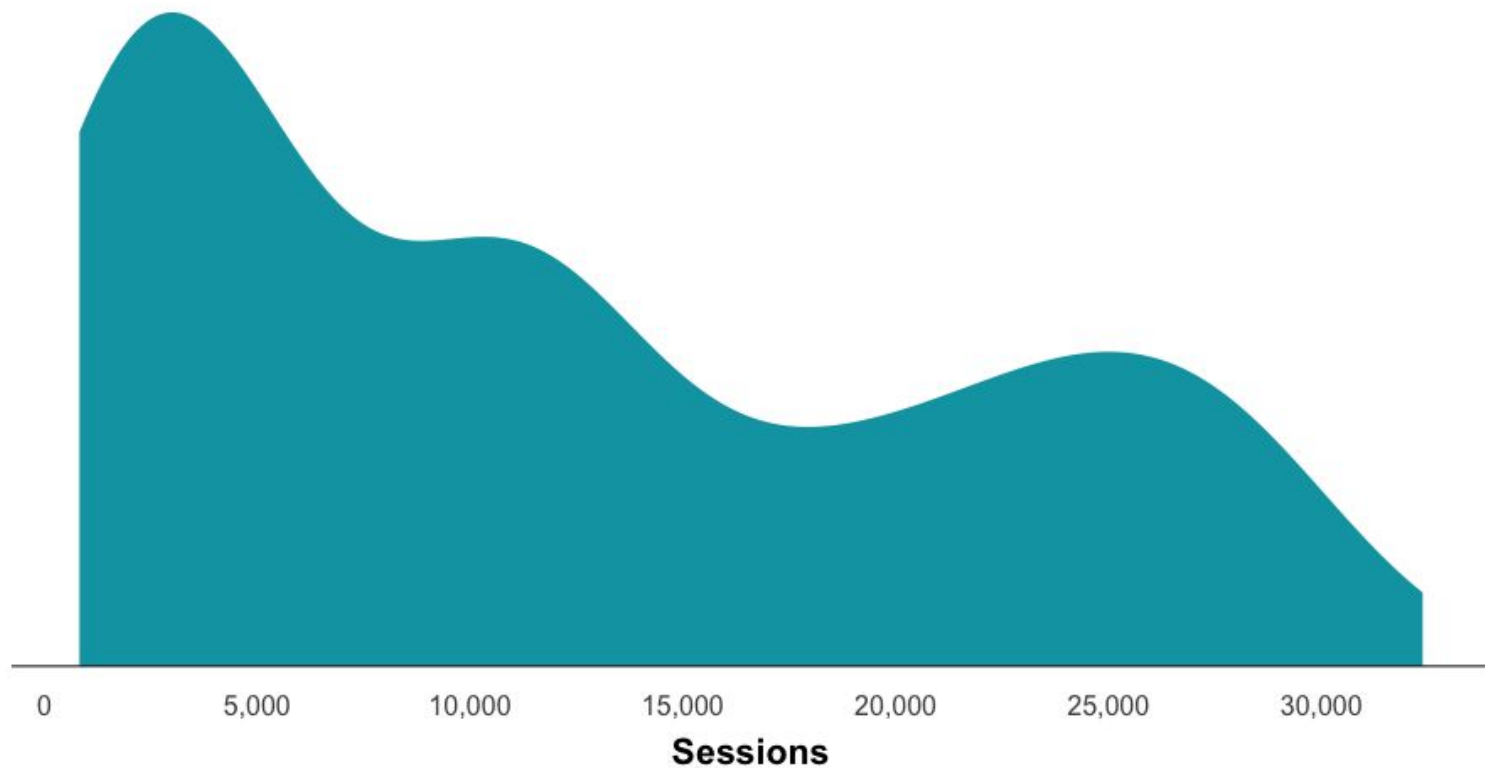
# Let's look at just one channel (n = 30)



Let's look at just one channel (n = 30)

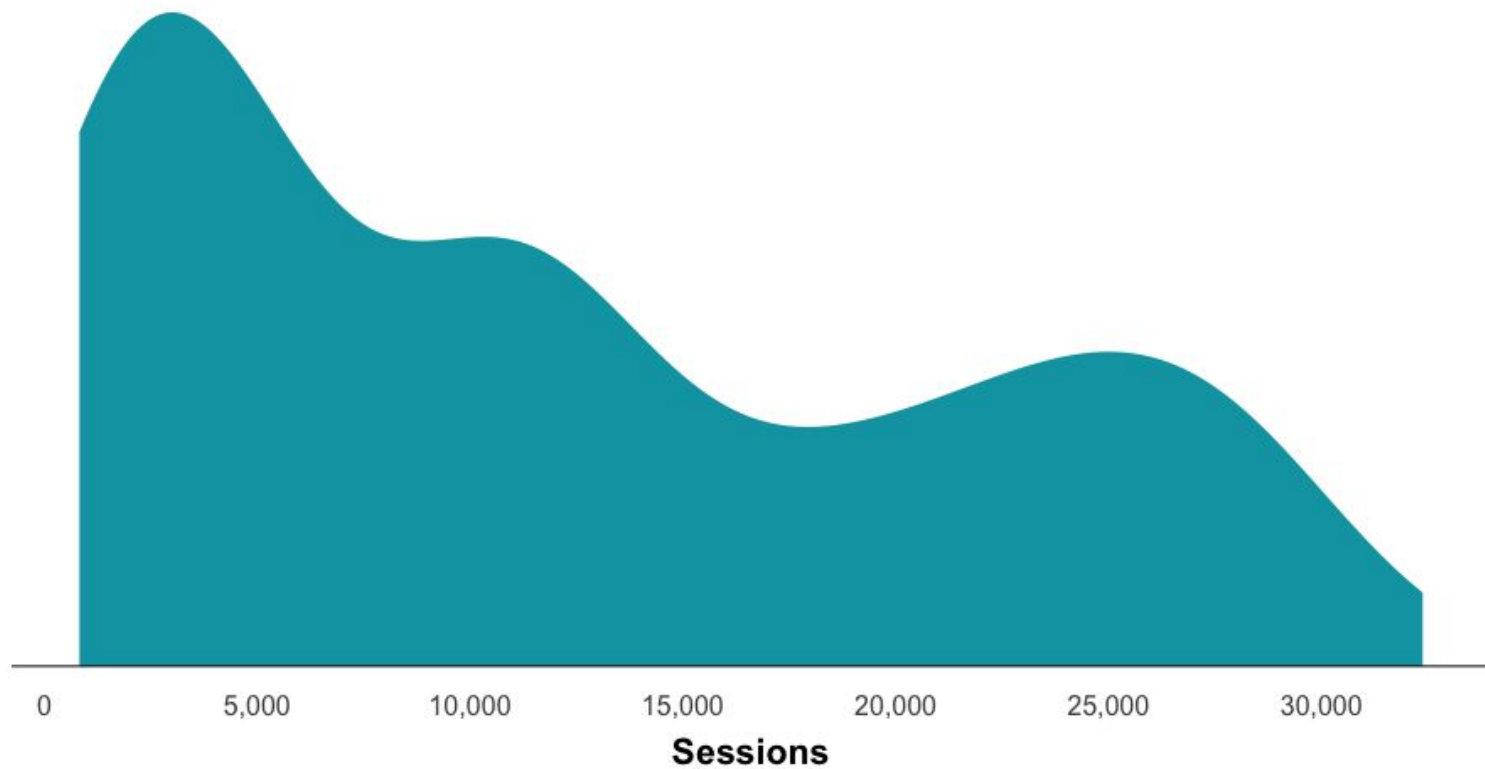


# Start with our overall distribution (n = 210)

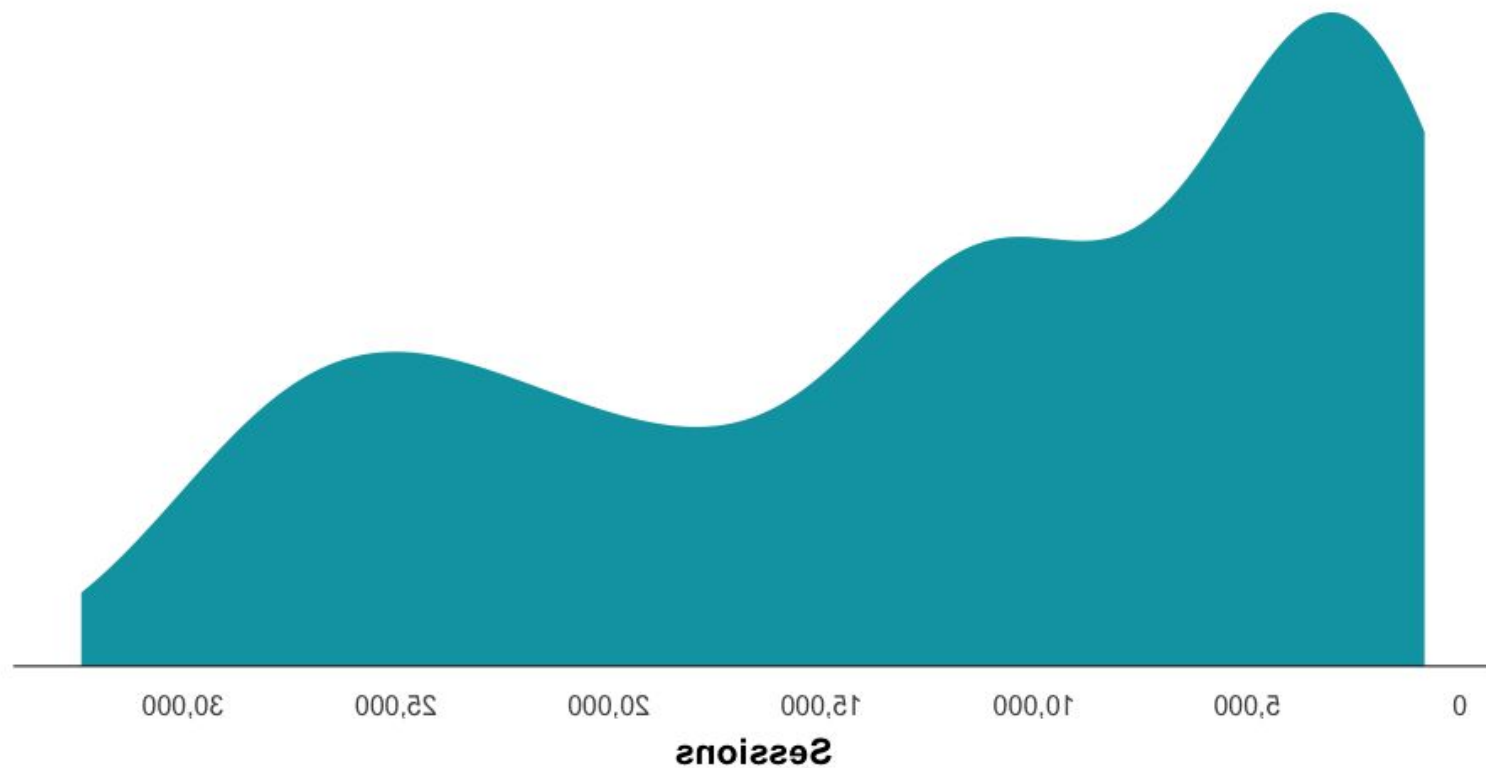




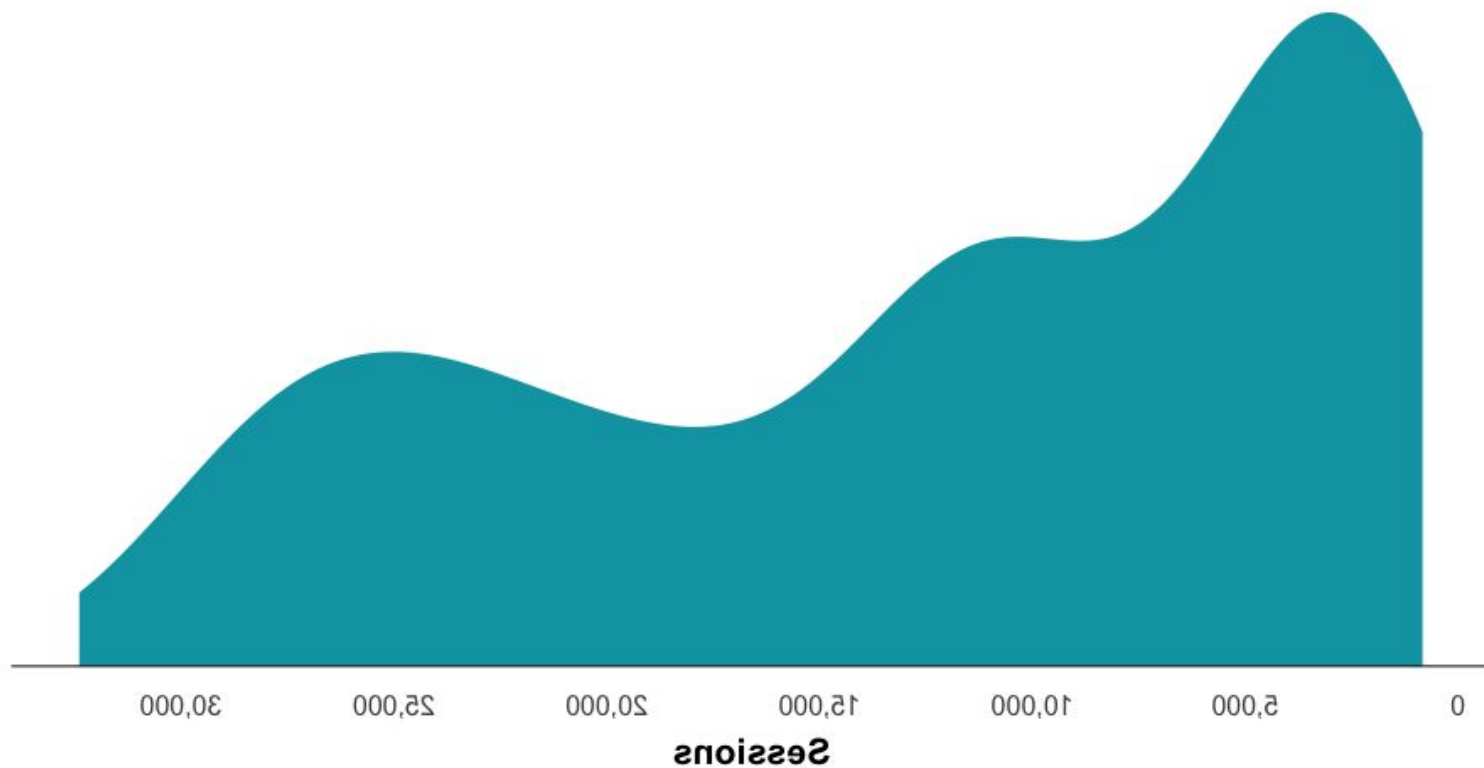
Flip it horizontally.



Flip it horizontally.

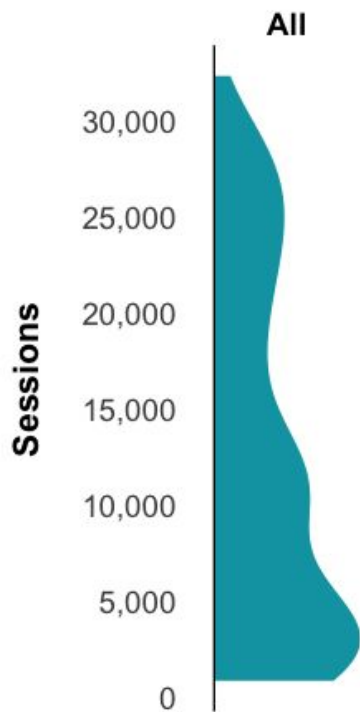


And rotate it and squish it.

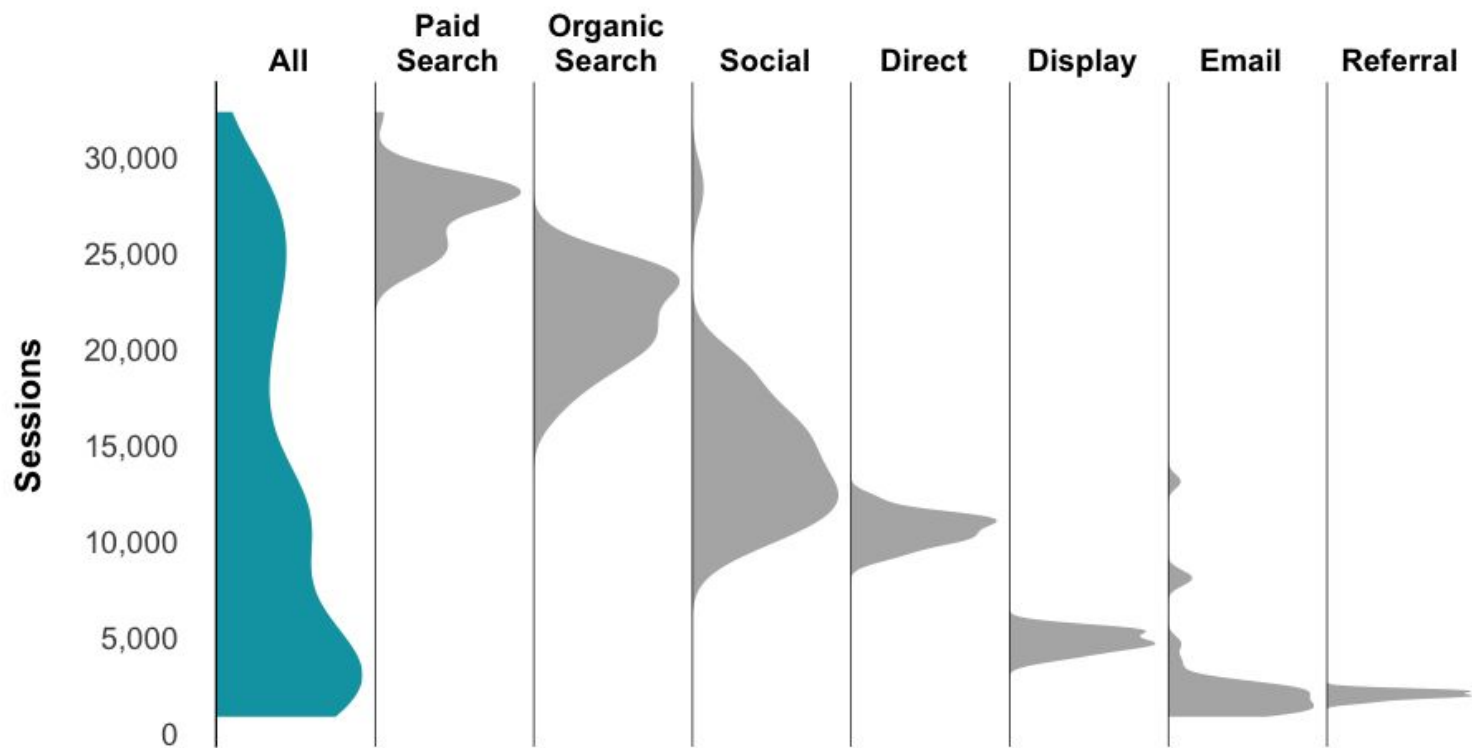




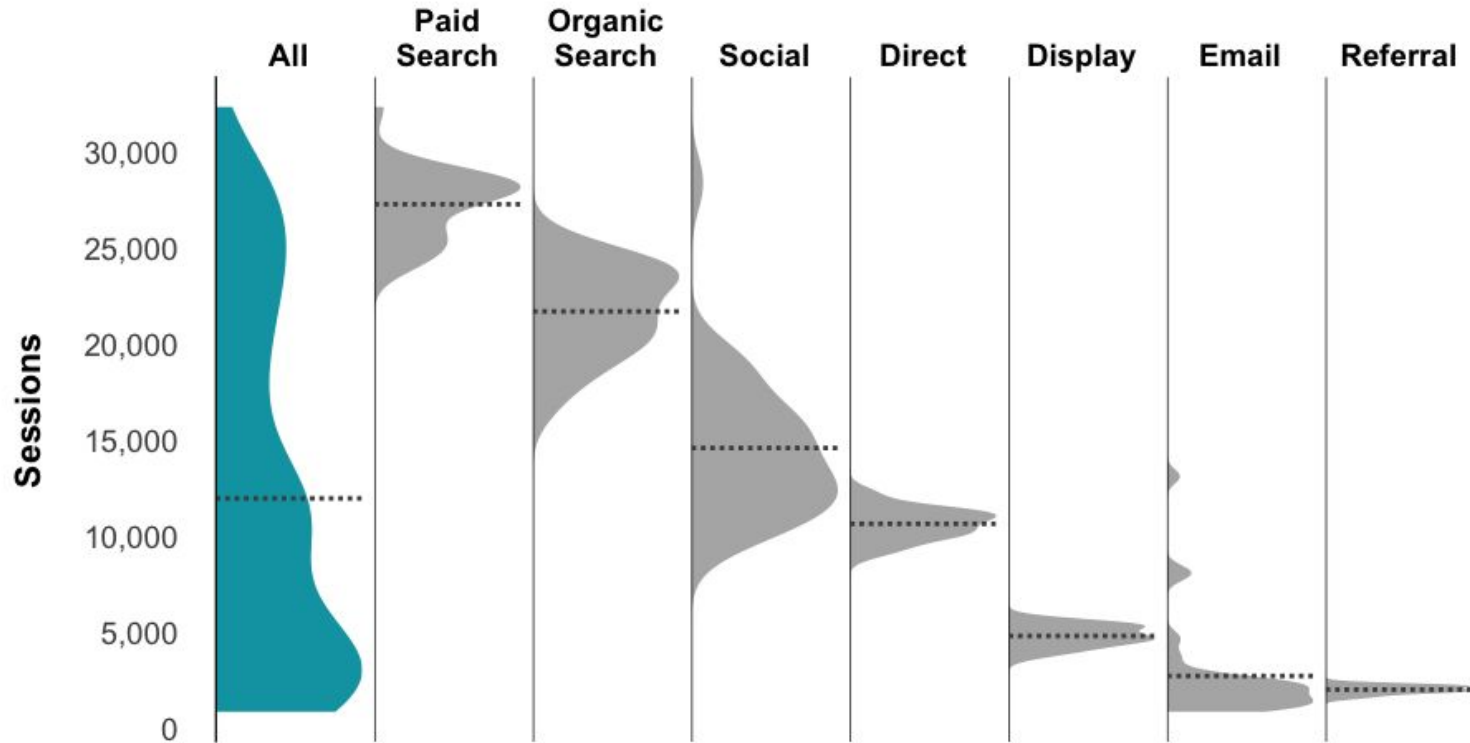
And rotate it and squish it.



“All” is the sum of the channels.

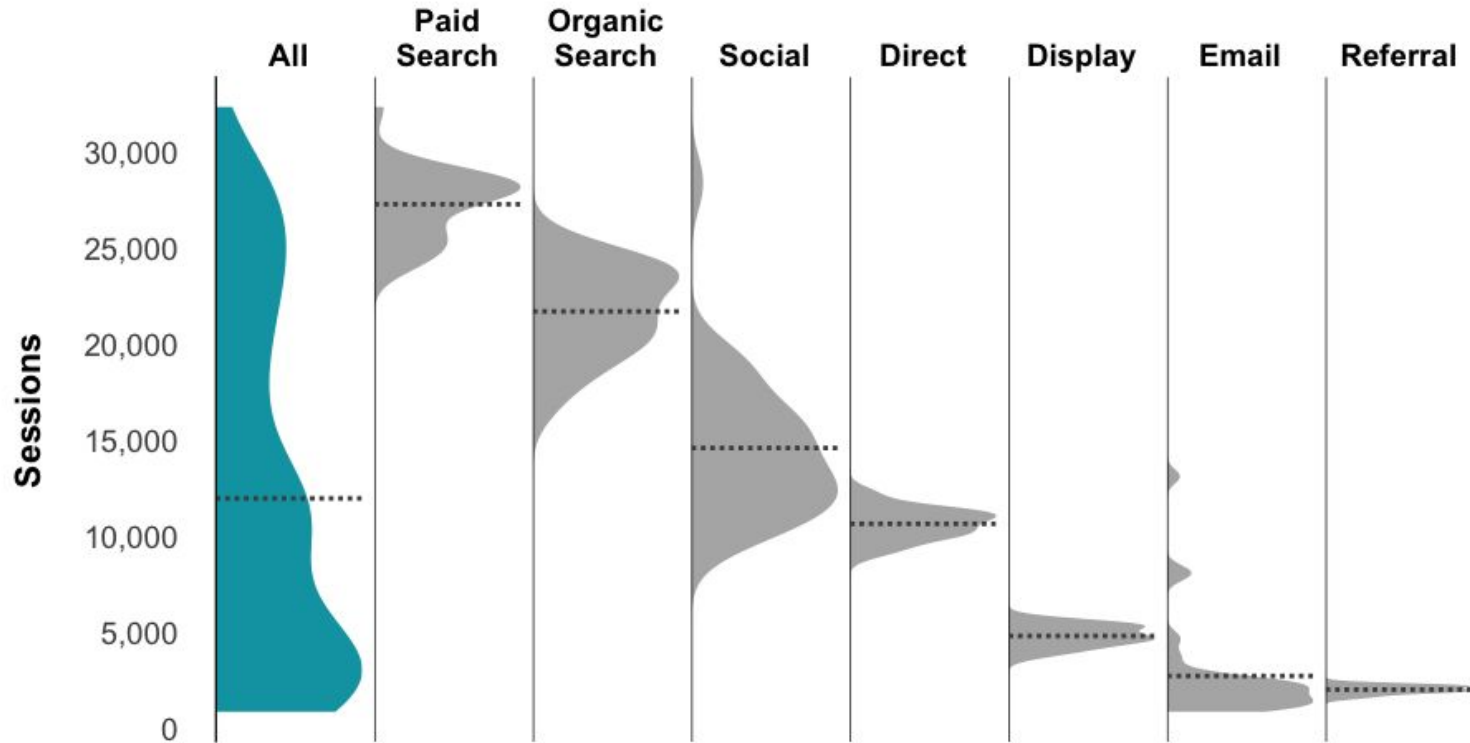


And we can add the **mean** for each distribution.

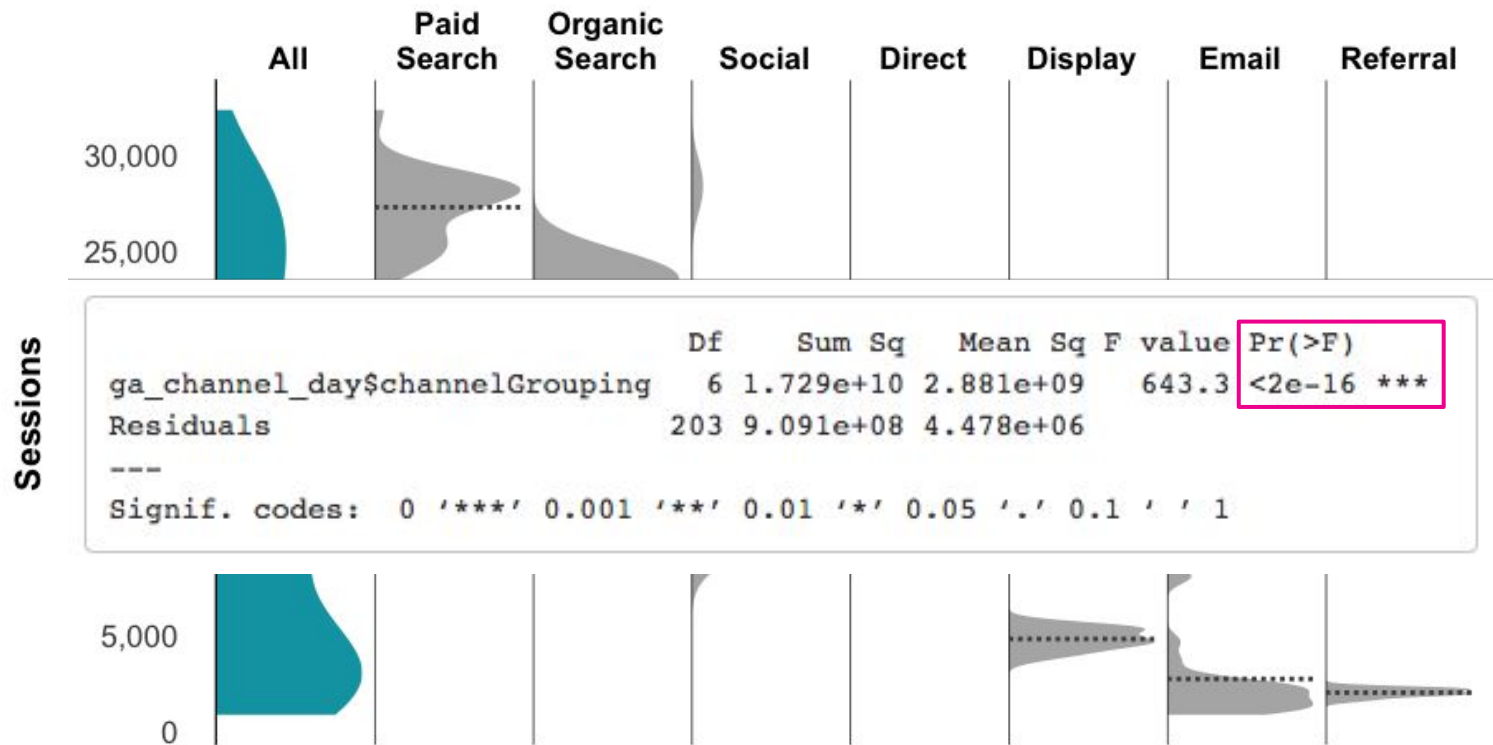




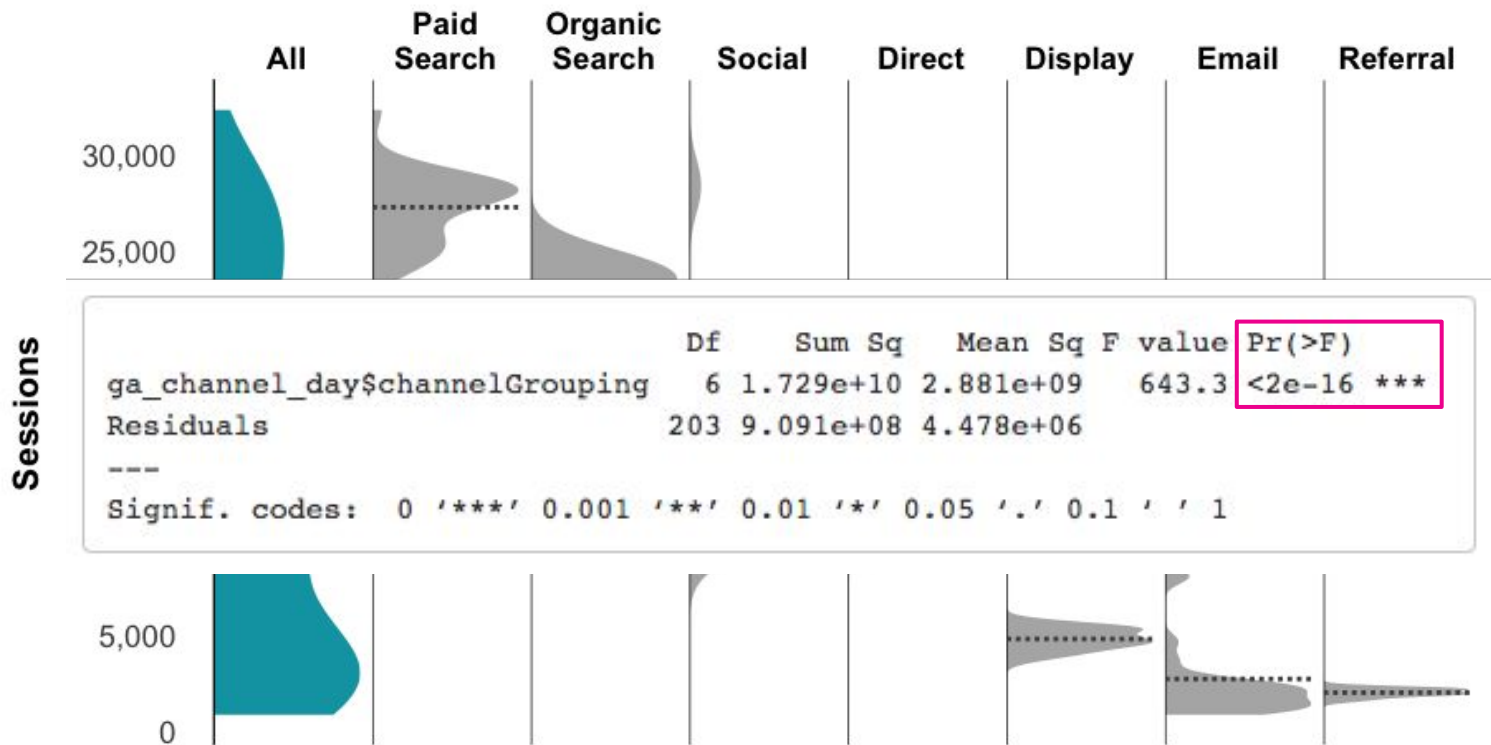
# We're set to do an **AN**alysis **Of** **VA**riance



# We're set to perform a 1-Way ANOVA.



WE can reject the null hypothesis.





*The Null Hypothesis ( $H_0$ ):*  
*Sessions do not really (“significantly”)*  
*differ by channel.*

*So...what we're saying:  
Sessions **do** appear to differ  
by channel.*

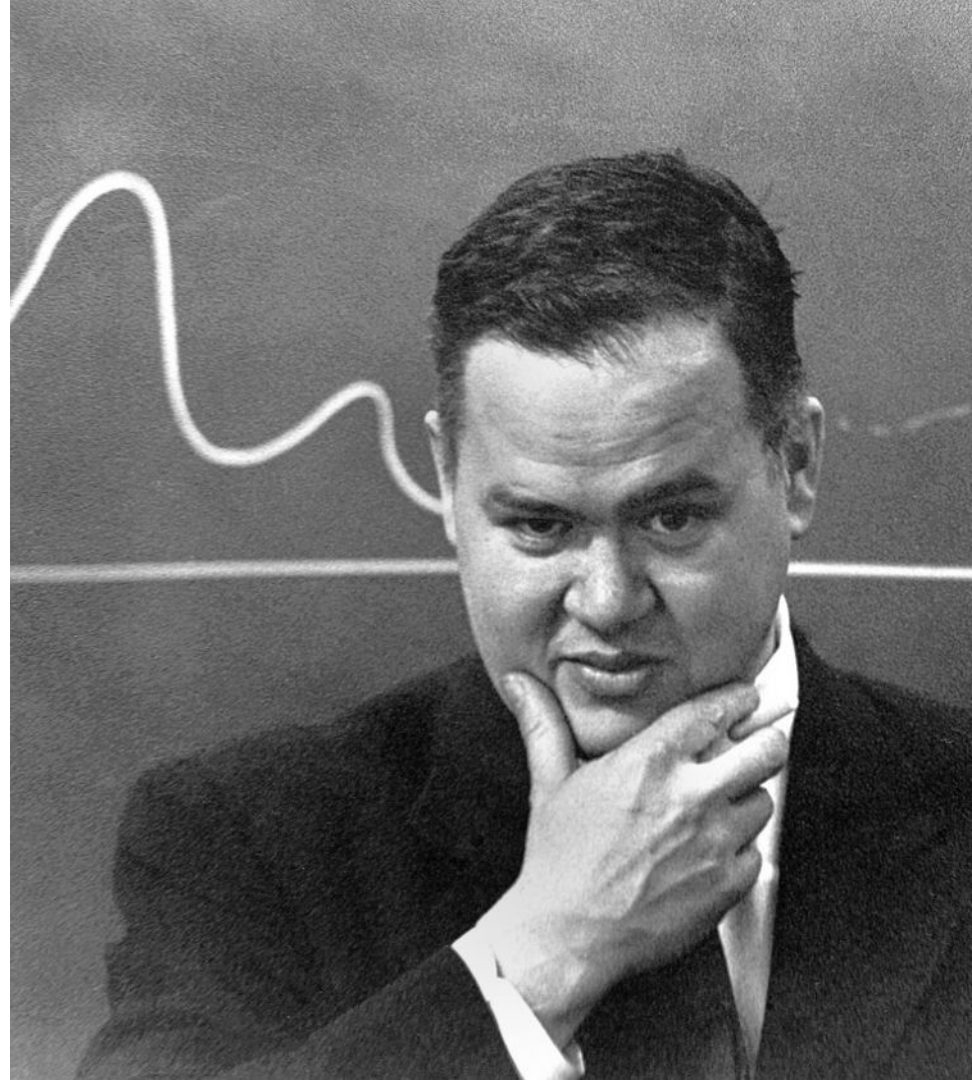
*An ANOVA does not tell us  
**which** channels differ.*



*But a **post hoc analysis**  
can do that.*

# John Tukey

...brought us the  
**Tukey Post Hoc Test.**

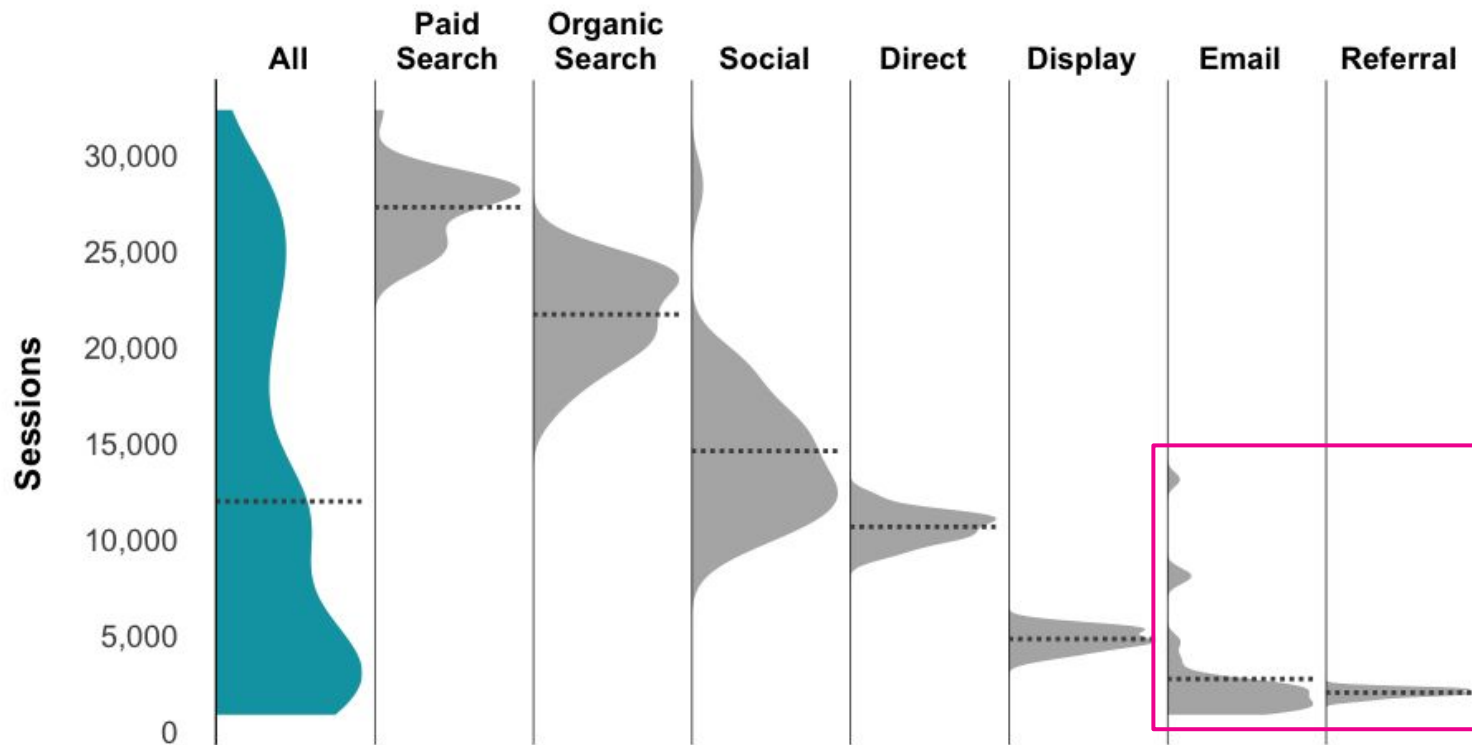


# Tukey Sayeth: Most Channels Differ

Email	0.84938					
Display	0.00001	0.00330				
Direct	0.00000	0.00000	0.00000			
Social	0.00000	0.00000	0.00000	0.00000		
Organic Search	0.00000	0.00000	0.00000	0.00000	0.00000	
Paid Search	0.00000	0.00000	0.00000	0.00000	0.00000	0.00000
	Referral	Email	Display	Direct	Social	Organic Search



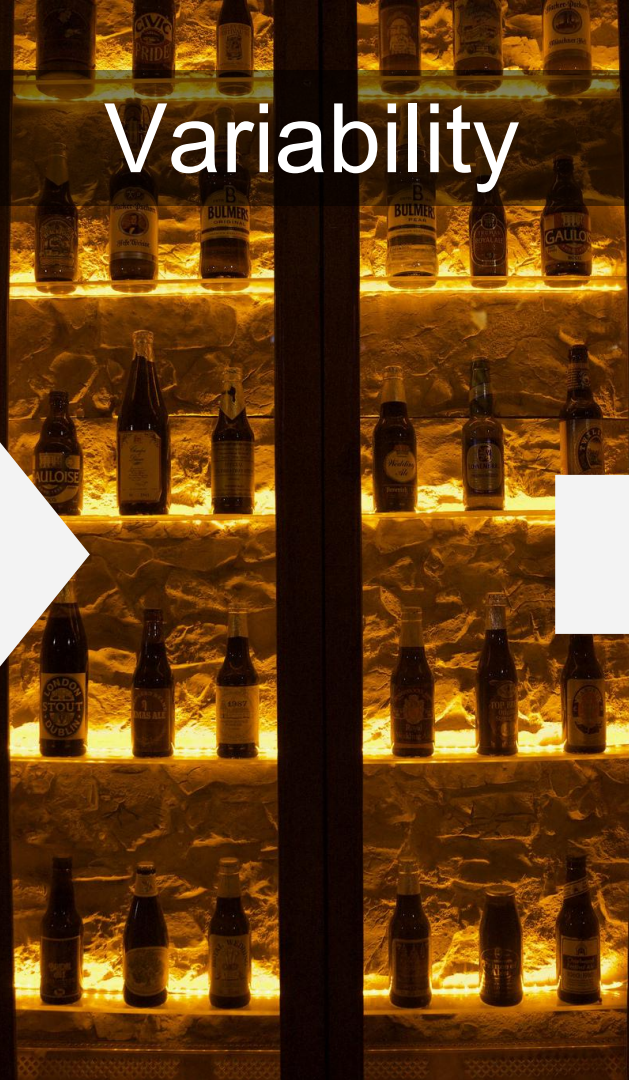
# The result matches our intuition!



# Distributions



# Variability



# Uncertainty



“We have a problem  
in [marketing] with  
thinking  
probabilistically.”

- Annie Duke



Source: Annie Duke



# This is really important!

All marketers inherently operate under **conditions of uncertainty**.

There is a **cost to reduce uncertainty**.

Uncertainty **can't be eliminated**

I've been **#mattgershuffed**

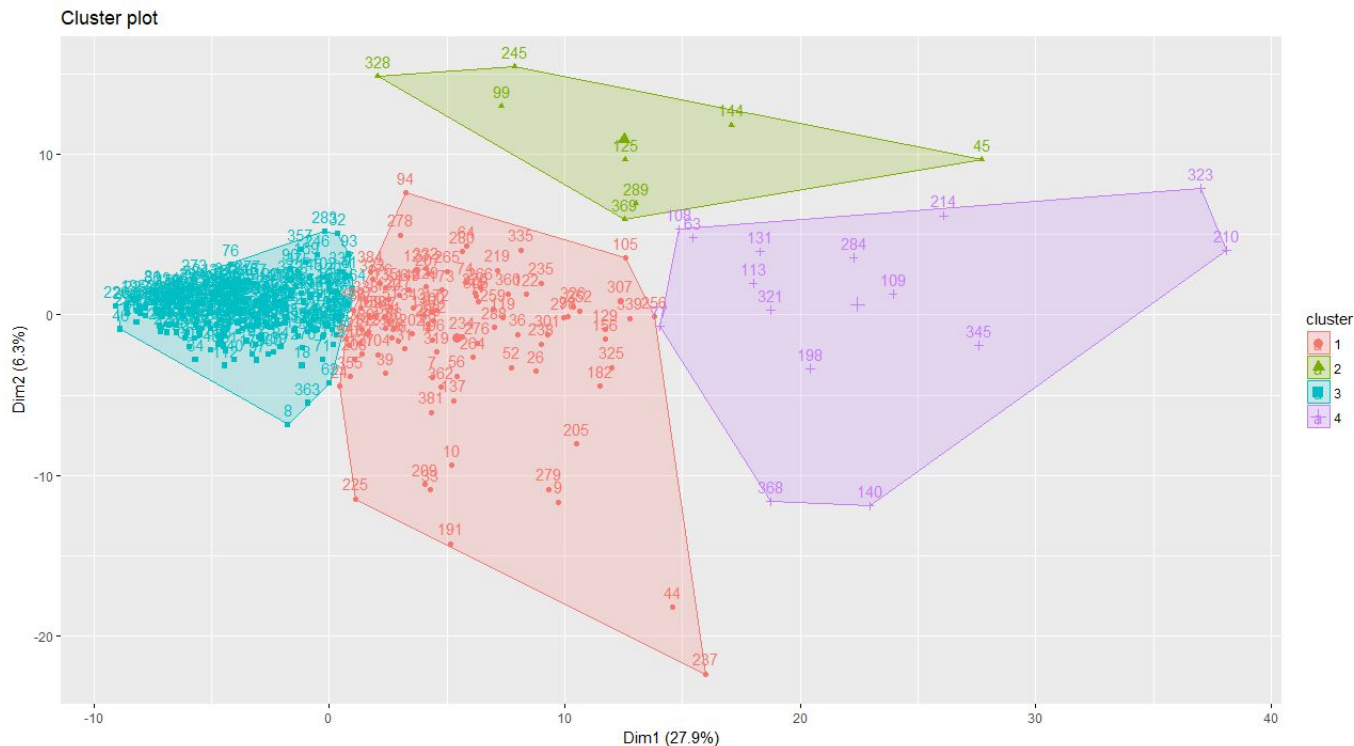


Matt Gershoff  
CEO, Conductrics

# So...Data Science

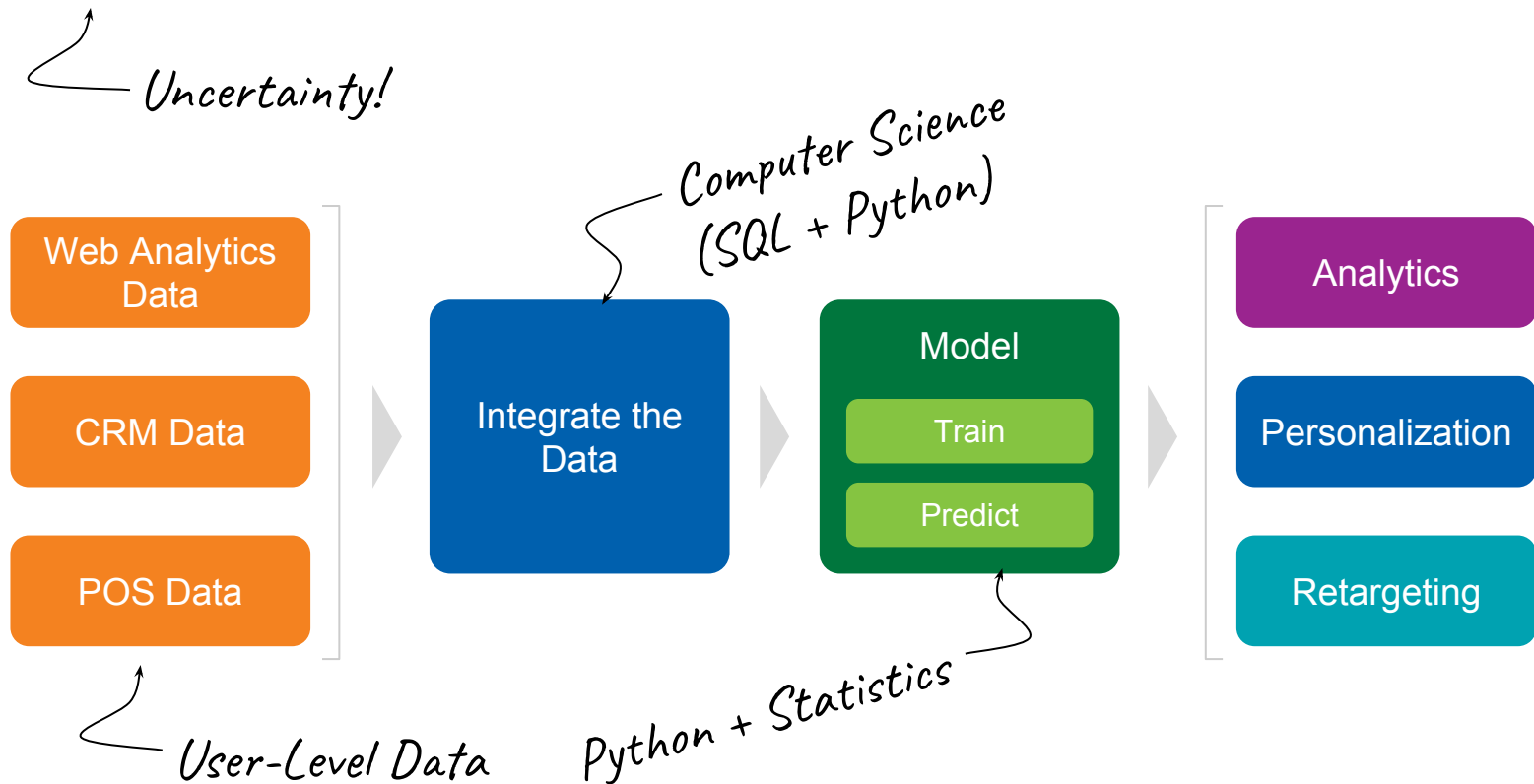


# *k*-means Clustering of Franchises





# Propensity Modeling



# So...where to begin?





[bit.ly/data-science-y](https://bit.ly/data-science-y)

