

# Big Data Project on Visualizing Play Store App Data

by

Avirup Das  
Prosenjit Dey  
Shiuli Subhra Ghosh  
Suman Roy

“Let everything happen to you.  
Beauty and terror.  
Just keep going. No Feeling is Final.”

Rainer Maria Rilke

- From a Visualization's point of view



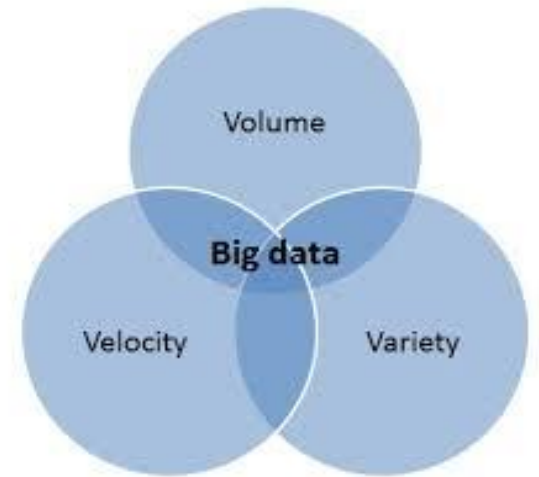
Let's Play

# Why do we need to visualise Google Play store App data?

- Earlier till 2020, the total number of apps on the Google Play Store was 2.87 million which has now **risen to nearly 5 million**.
- In 2020, there were **108.5 billion apps downloaded** by the users on Google Play Store.
- If we analyse and visualise this huge data, we can gain insights on what parameters an app outperforms others.
- Developers might find these insights very useful to before creating an app or making their apps better.

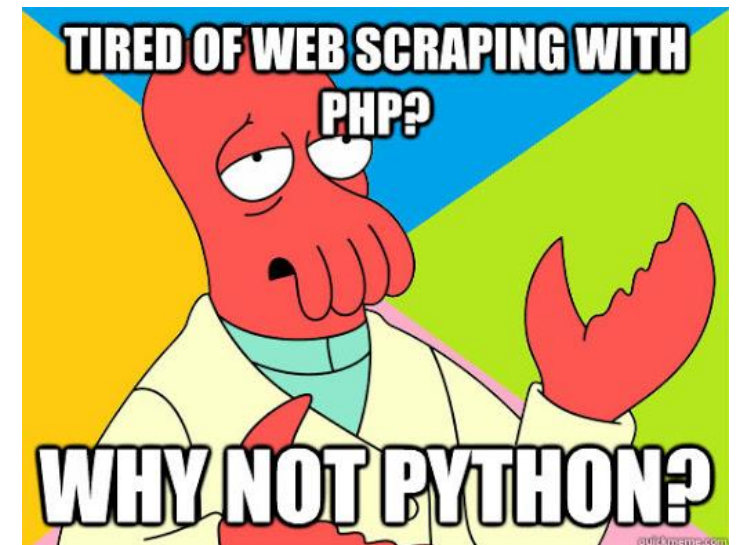
## Why is this a Big Data Project?

- **Big Data** is a collection of data that is huge in volume, yet growing exponentially with time. Google Play Store data is also of huge volume and growing every day.
- Big data is often characterized by the three V's. **Volume, Variety and Velocity**.
- Volume: The number of apps in Google Play Store is almost 5 Million.
- Variety: For each app there are various types of data available like number of installs, score, size, reviews, content rating, price etc.
- Velocity: These data are generated in an enormous speed. New apps, updates of apps, reviews and scores given by users and many other things contribute in generation of data. For effective analysis and visualisation these are to be processed in similar speed.

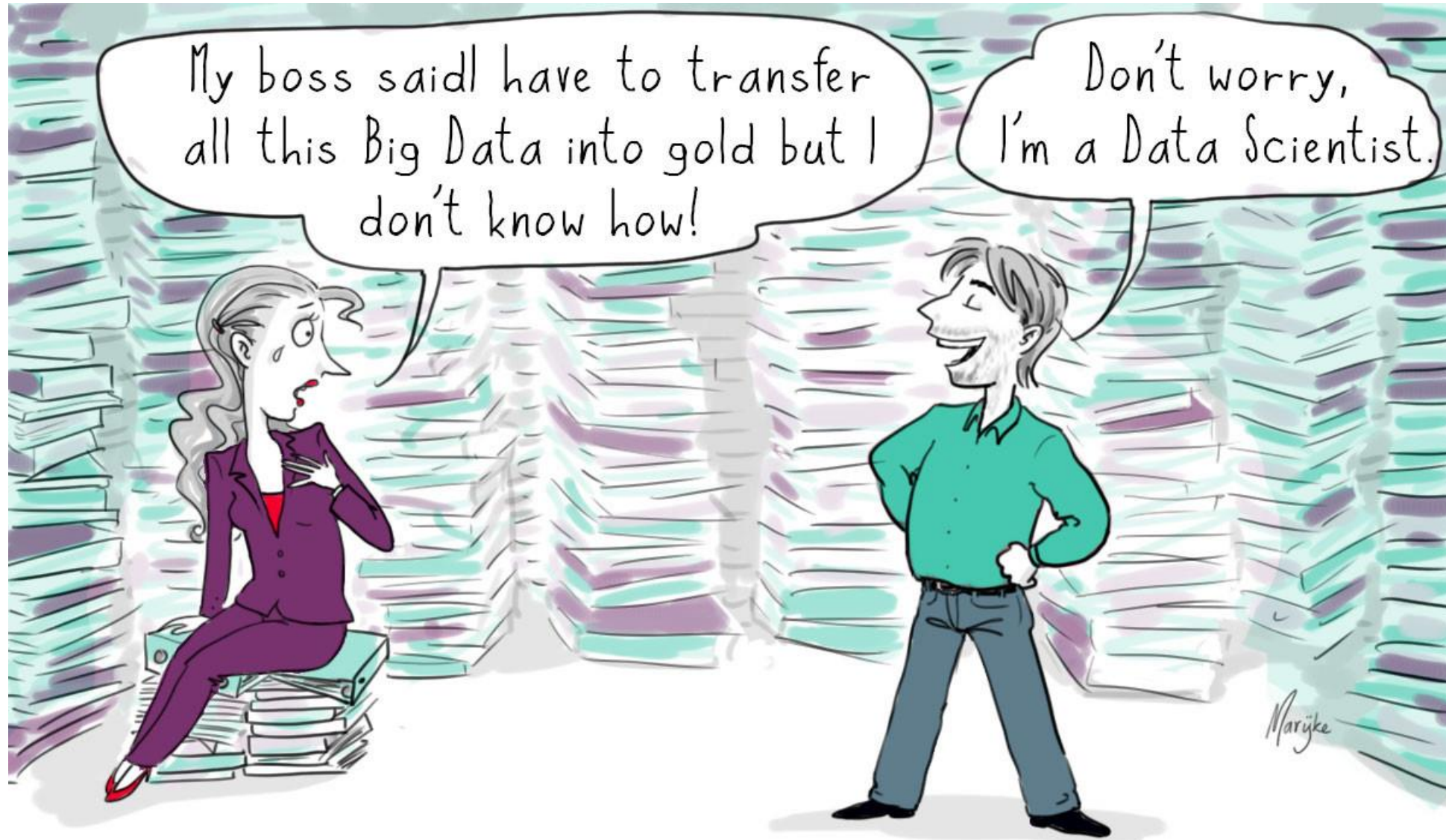


## Web Scraping

- We designed a code which will crawl over the Play-store using the links provided and collect all relevant data.
- The crawler visits each of the links of different categories of app, waits for 5 seconds for the page to load and scrolls down to the bottom of the page.
- Our bot waits for 3 seconds for the new apps to show up and will continue doing this until there are no more apps to show.
- At this stage, it has the app-id's of all the apps it has come across and will visit the app-page to collect the relevant data.
- All this data is then transferred to a CSV file.

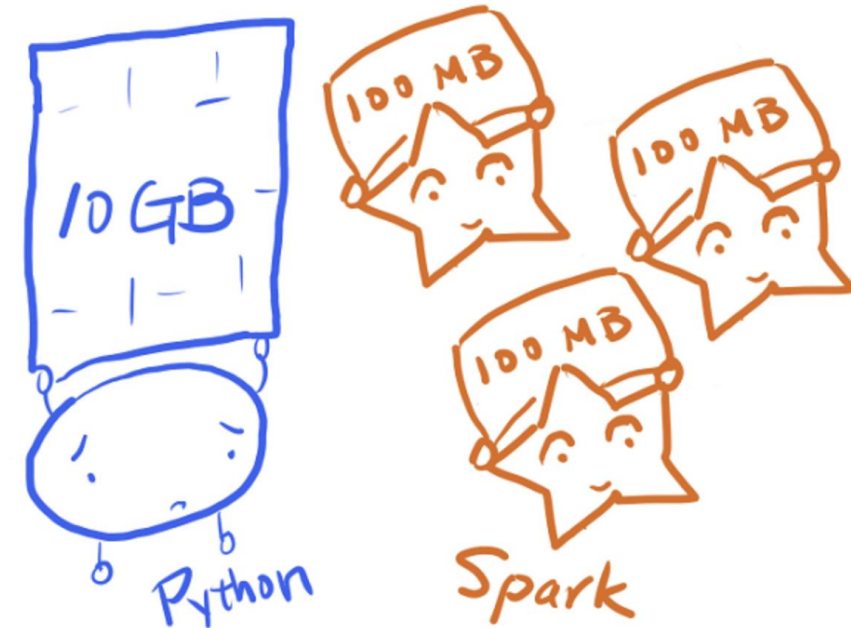




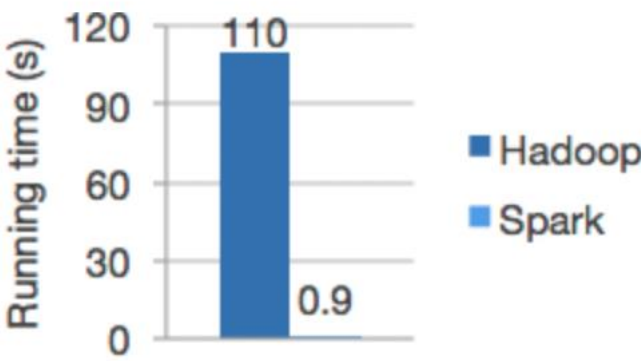
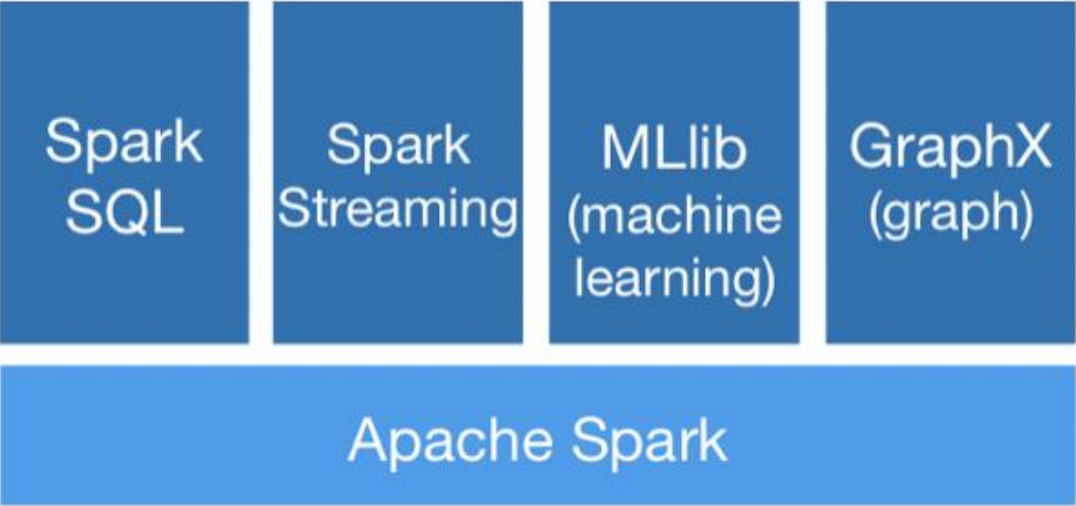


## Why spark?

- Spark has useful features such as Spark SQL, DataFrame, Streaming, MLlib and Spark Core.
- When the dataset is very large, Python will load the data all at once and will perform all operations on the whole data. Most of the time this becomes not possible due to our primary memory constraint.
- Spark re-partitions the data and picks it up one at a time. This way we can handle huge data even with comparatively less resources.
- Apache Spark achieves high performance for both batch and streaming data resulting in running workloads 100 times faster than Hadoop.



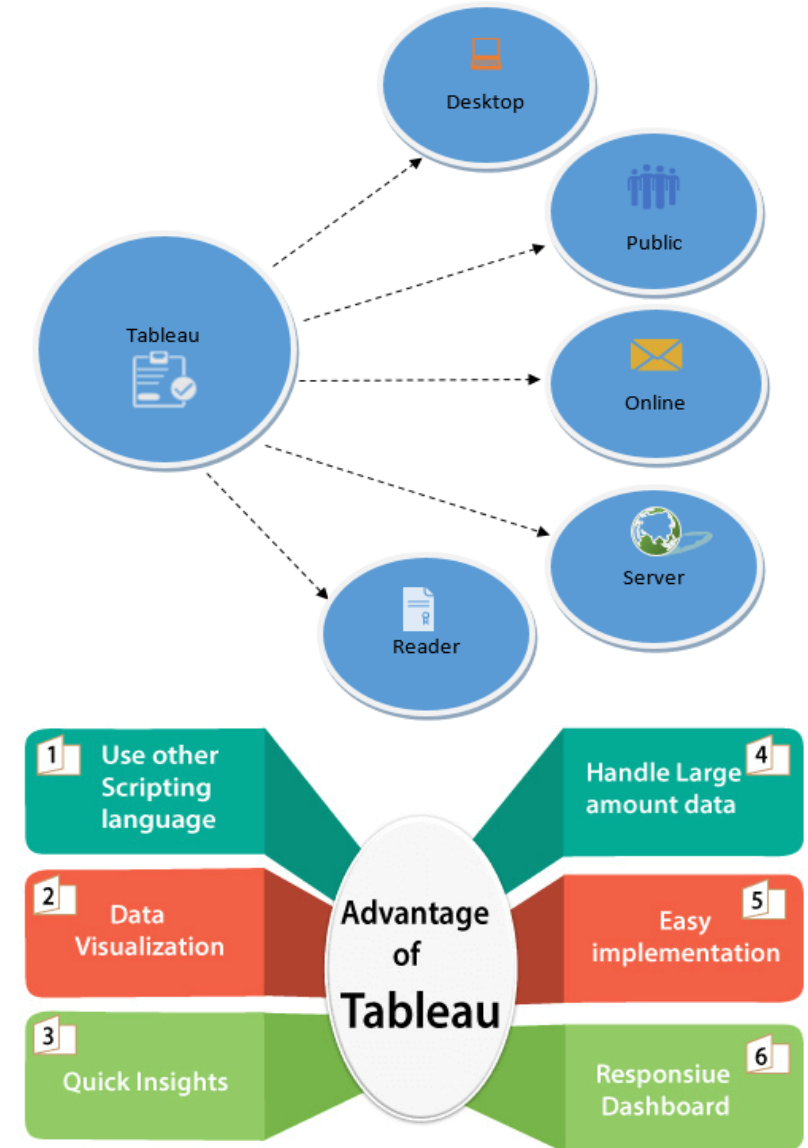




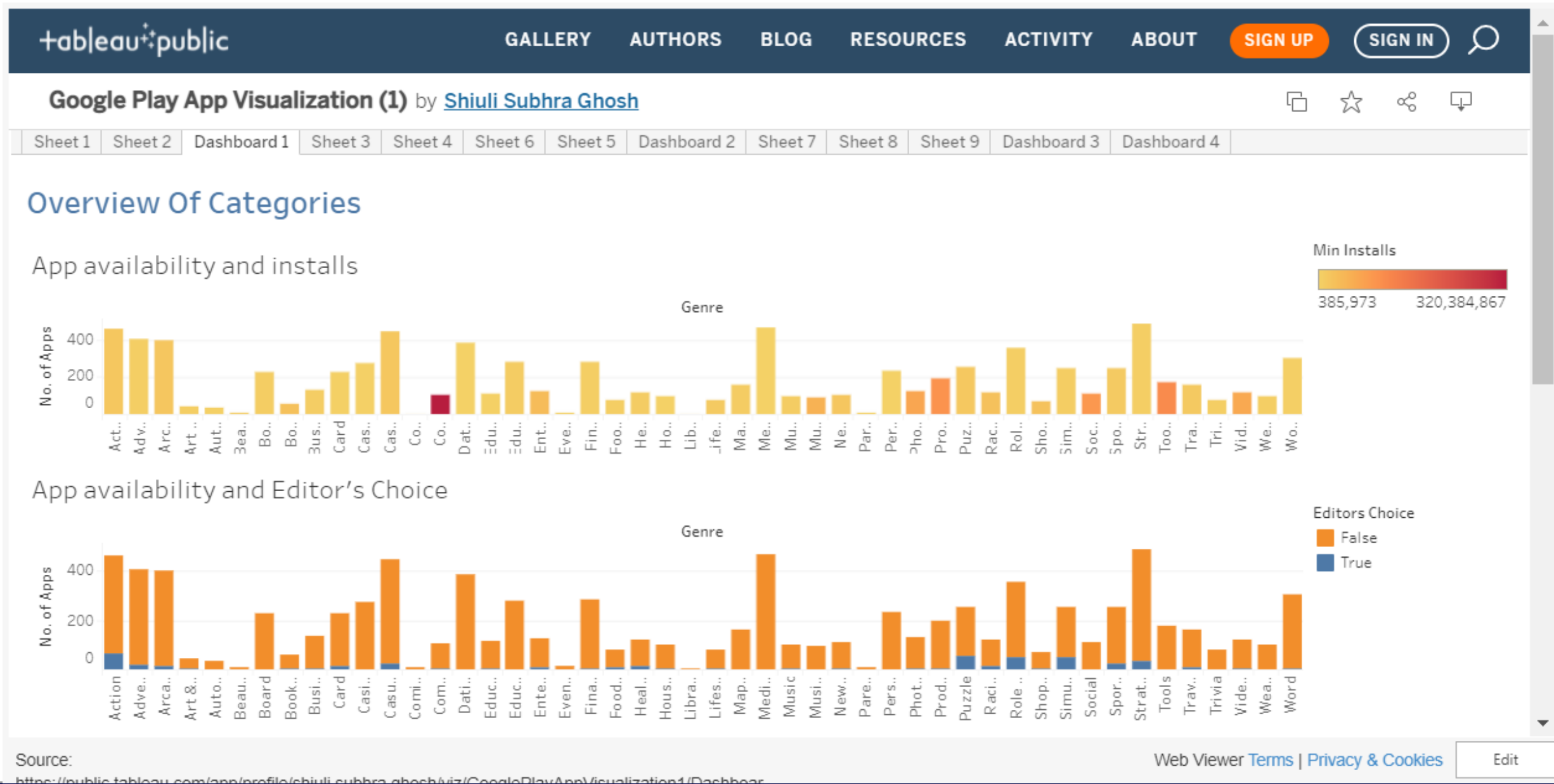
Logistic regression in Hadoop and Spark

## Why Tableau?

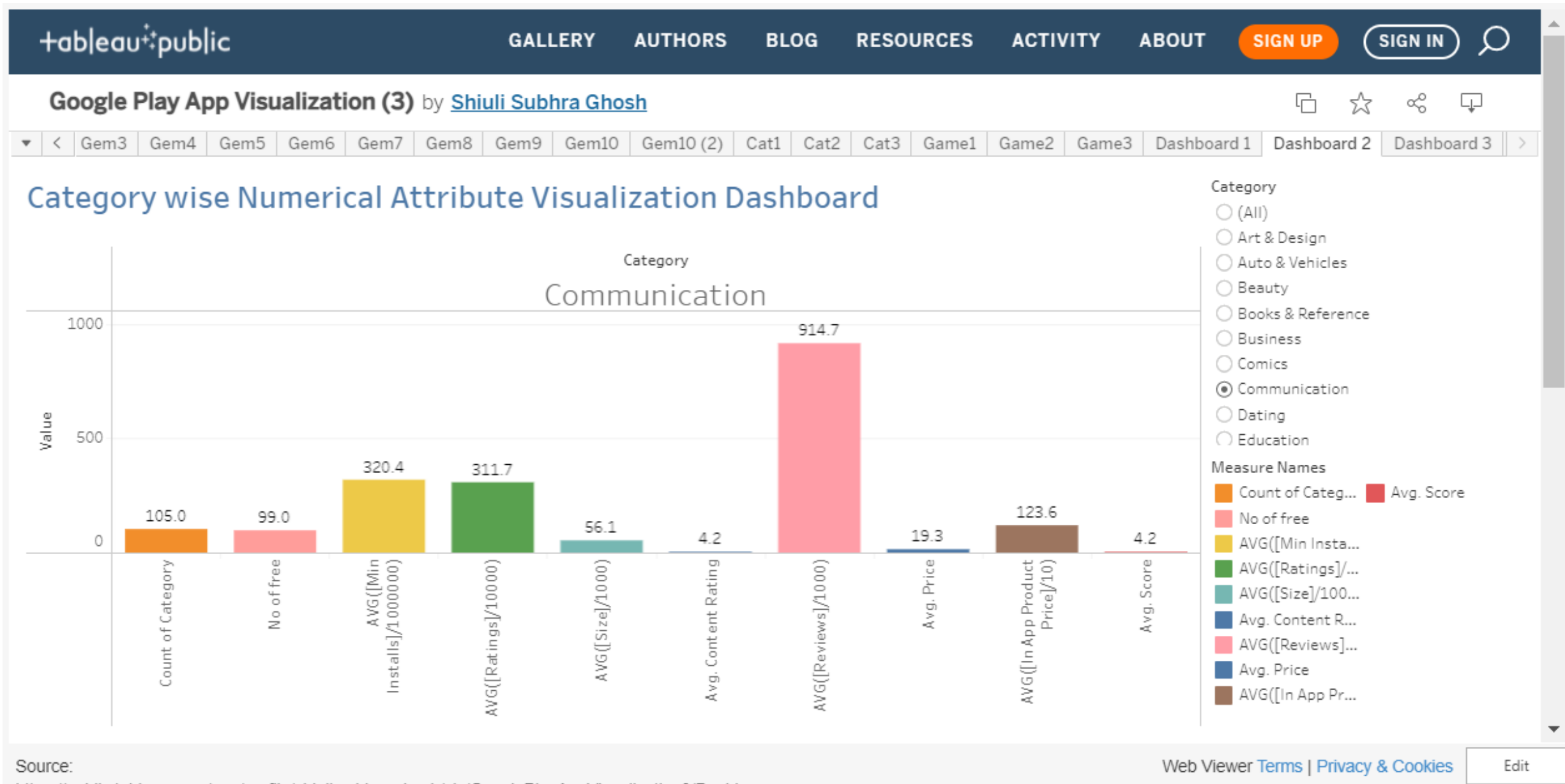
- Big data analytics describes the process of uncovering trends, patterns, and correlations in large amounts of raw data to help make data-informed decisions.
- Tableau Software is a tool that helps make Big Data small, and small data insightful and actionable.
- The main use of tableau software is to help people see and understand their data.
- One of the coolest Tableau uses is performing basic ETL operations fast!
- Different users can collaborate on the same project seamlessly in Tableau.
- Tableau comes with beautiful dashboards for better visualization.



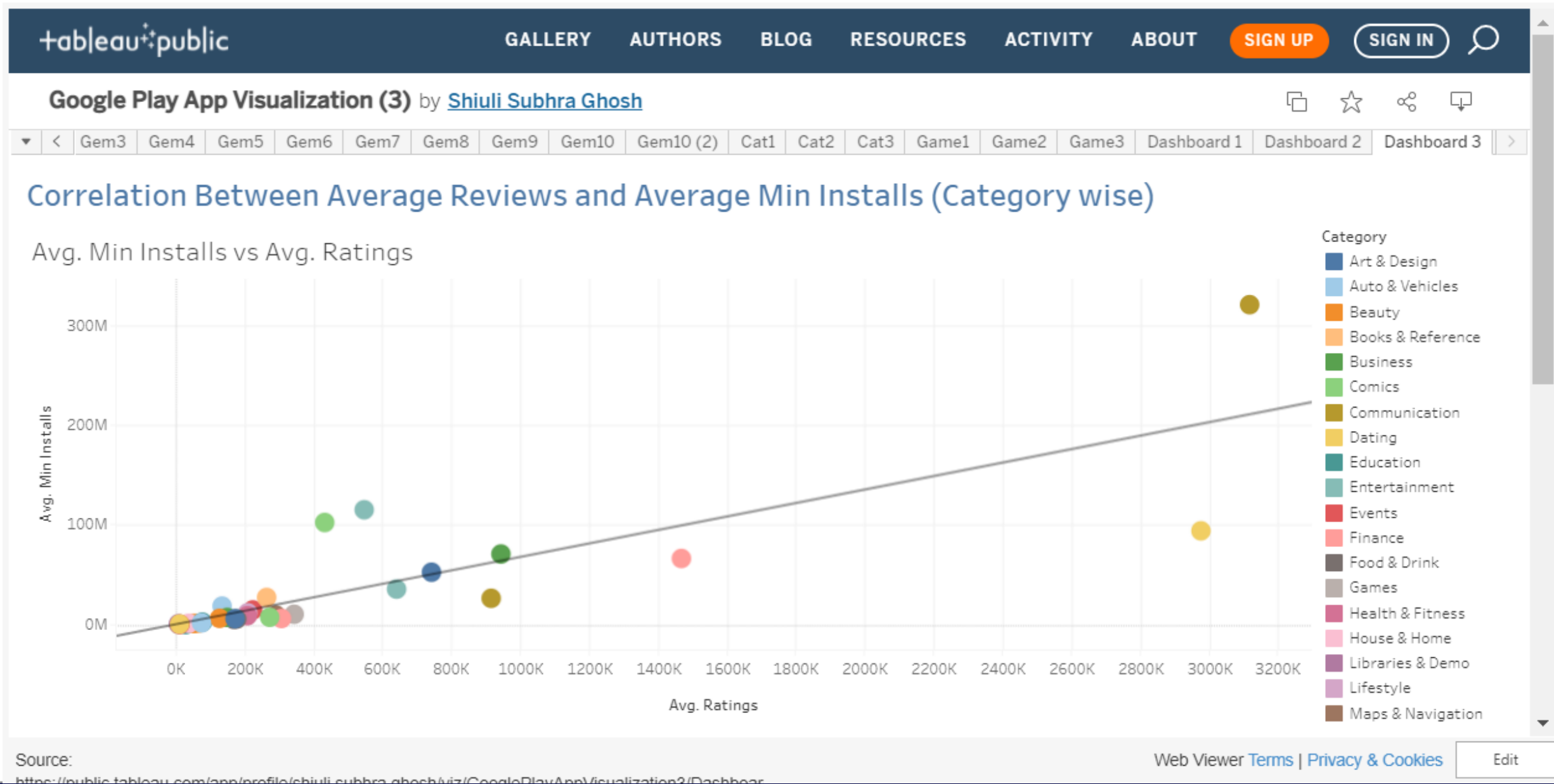
# Overview of Data (Category Wise)



# Overview of Data (Category wise)

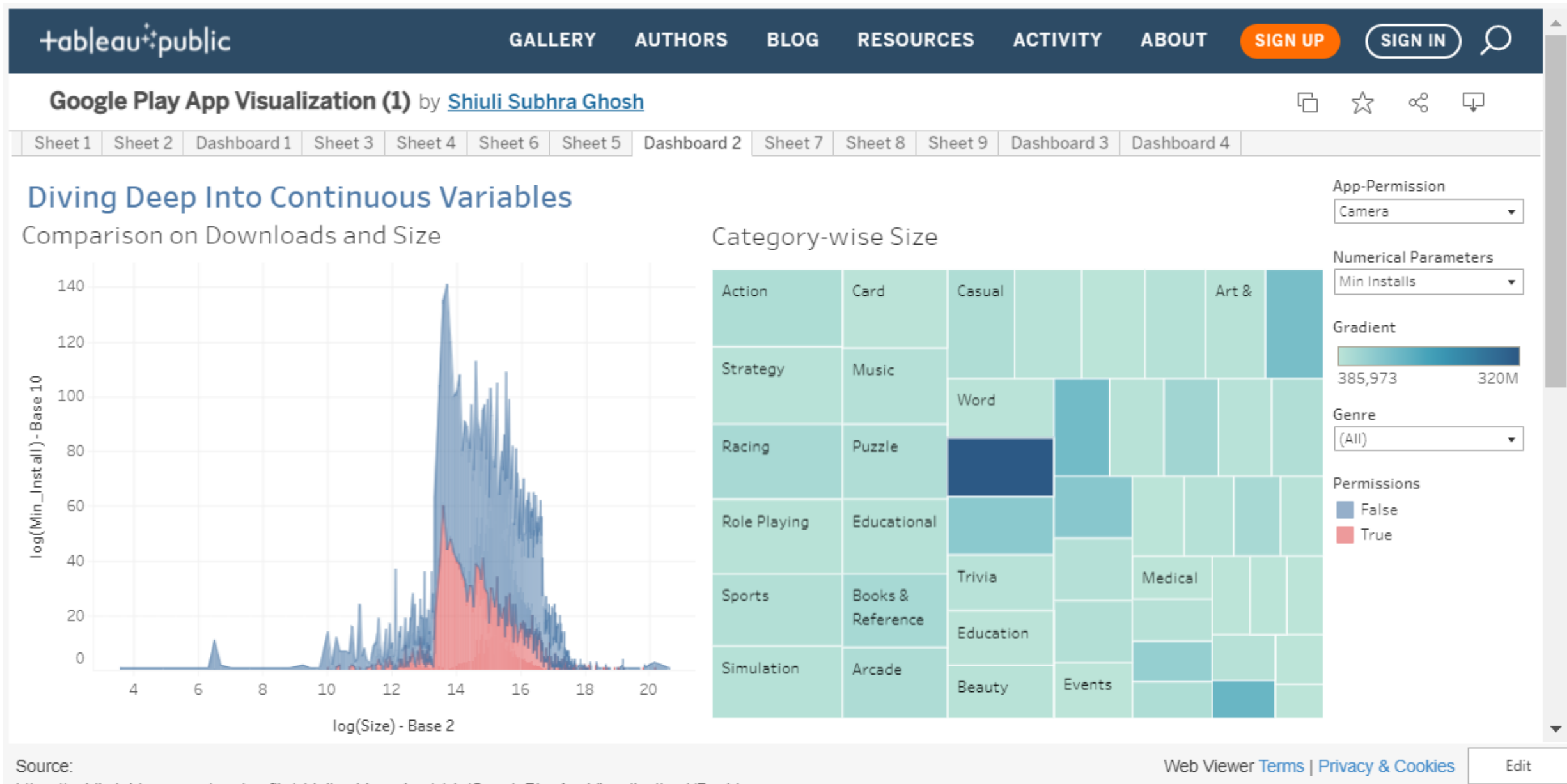


# Overview of Data (Category wise)

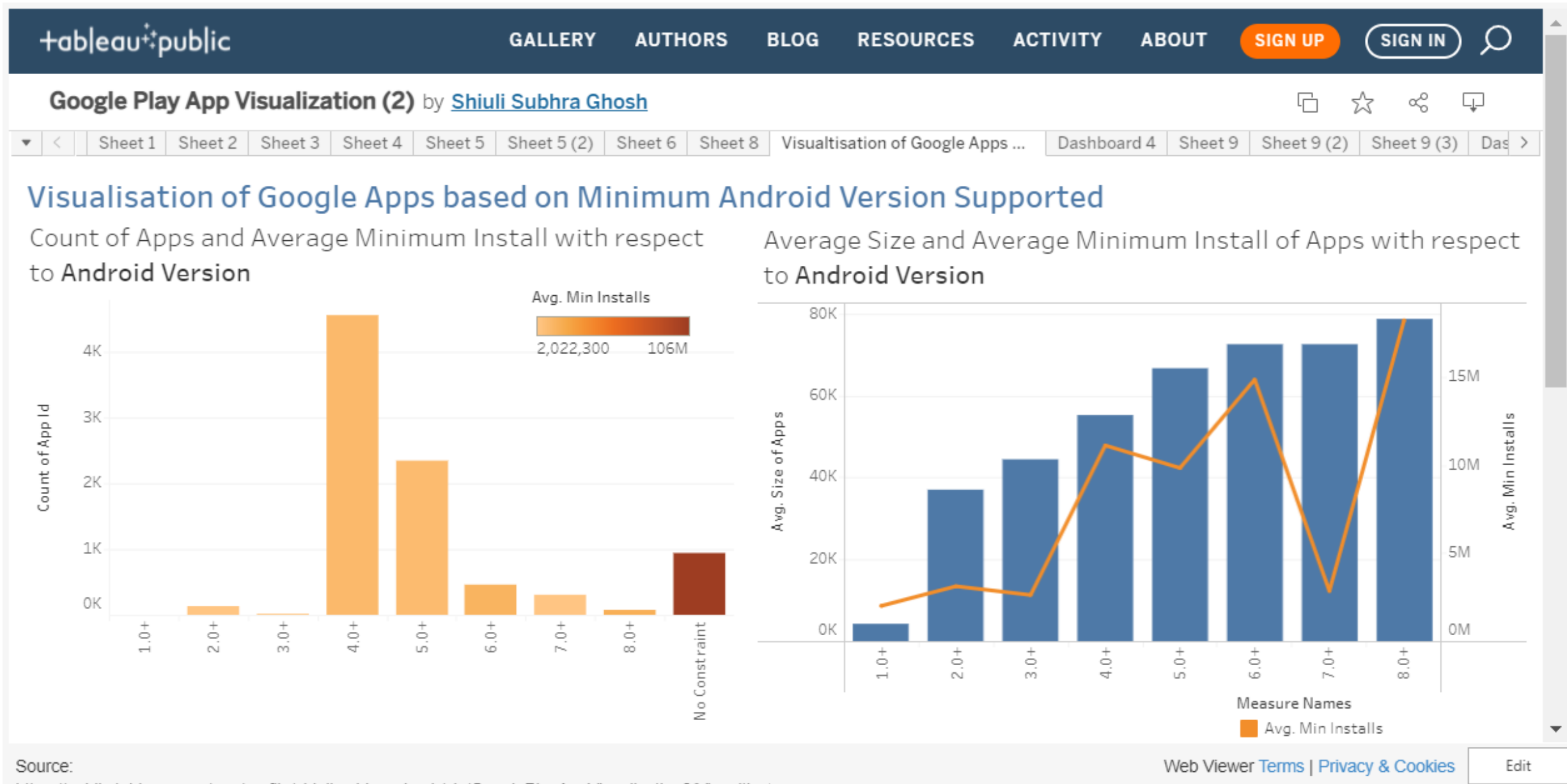




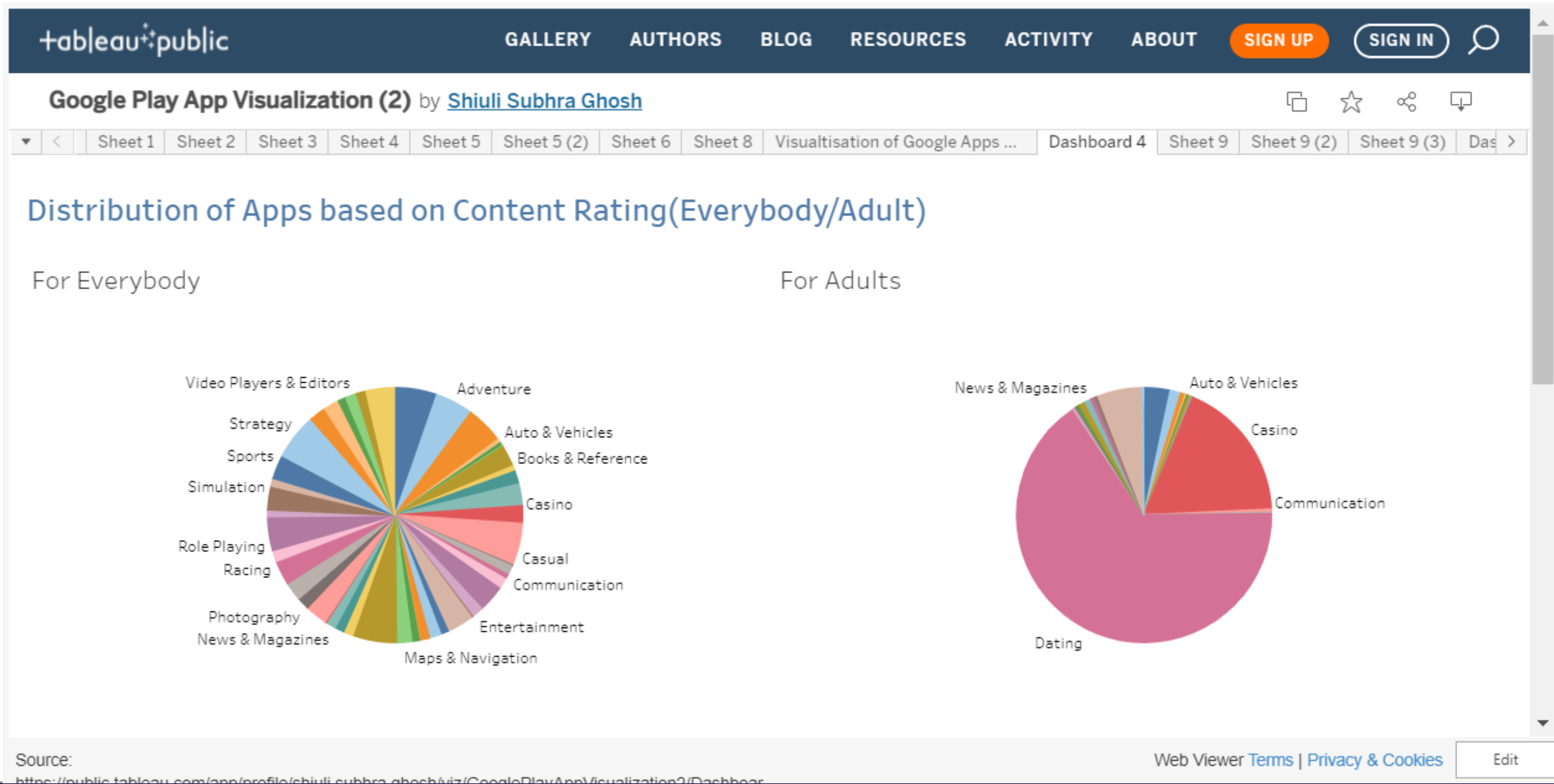
# Diving Deep Into Continuous Variables



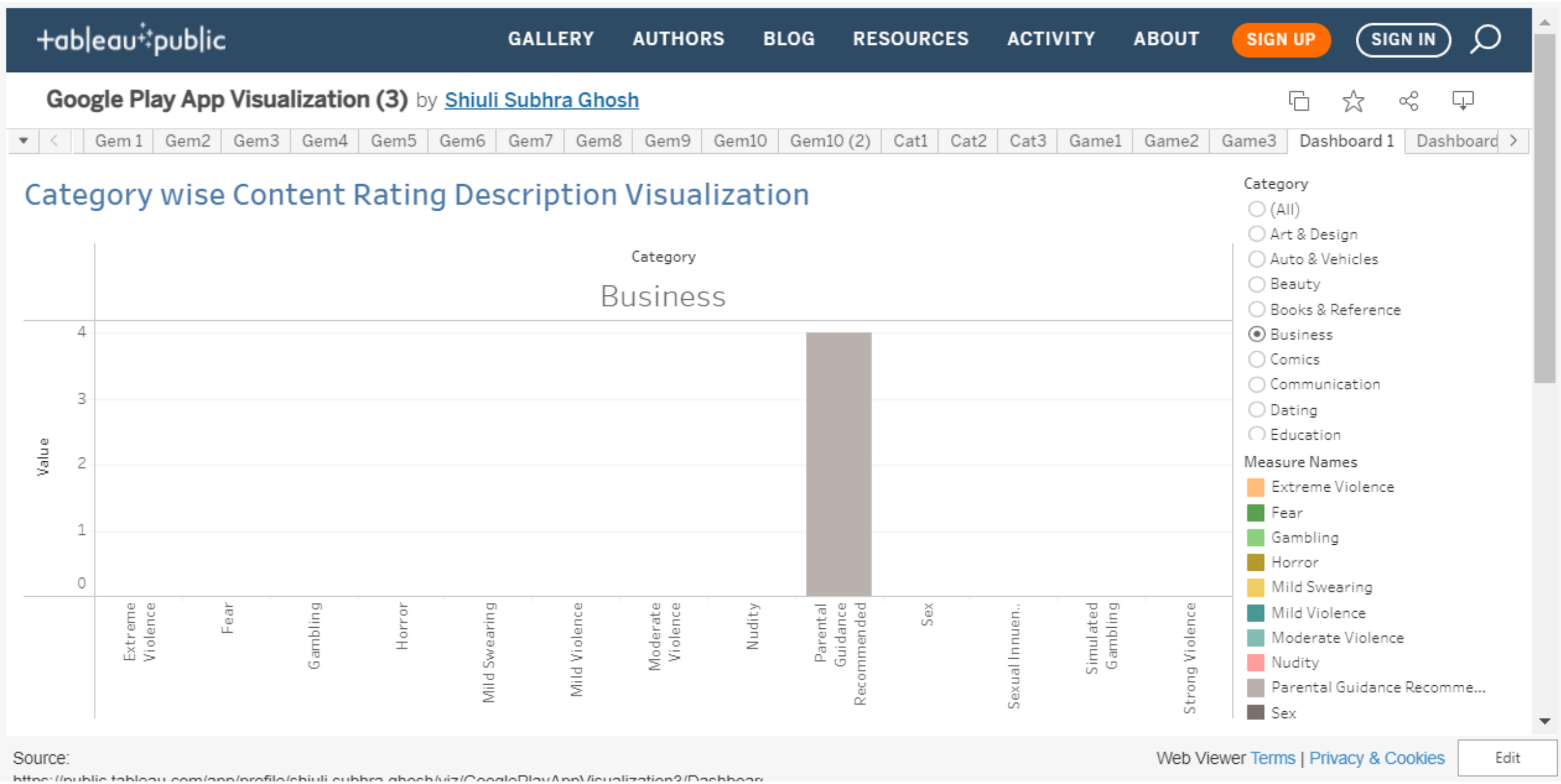
# Visualization of Google Apps based on Minimum Android Version Supported



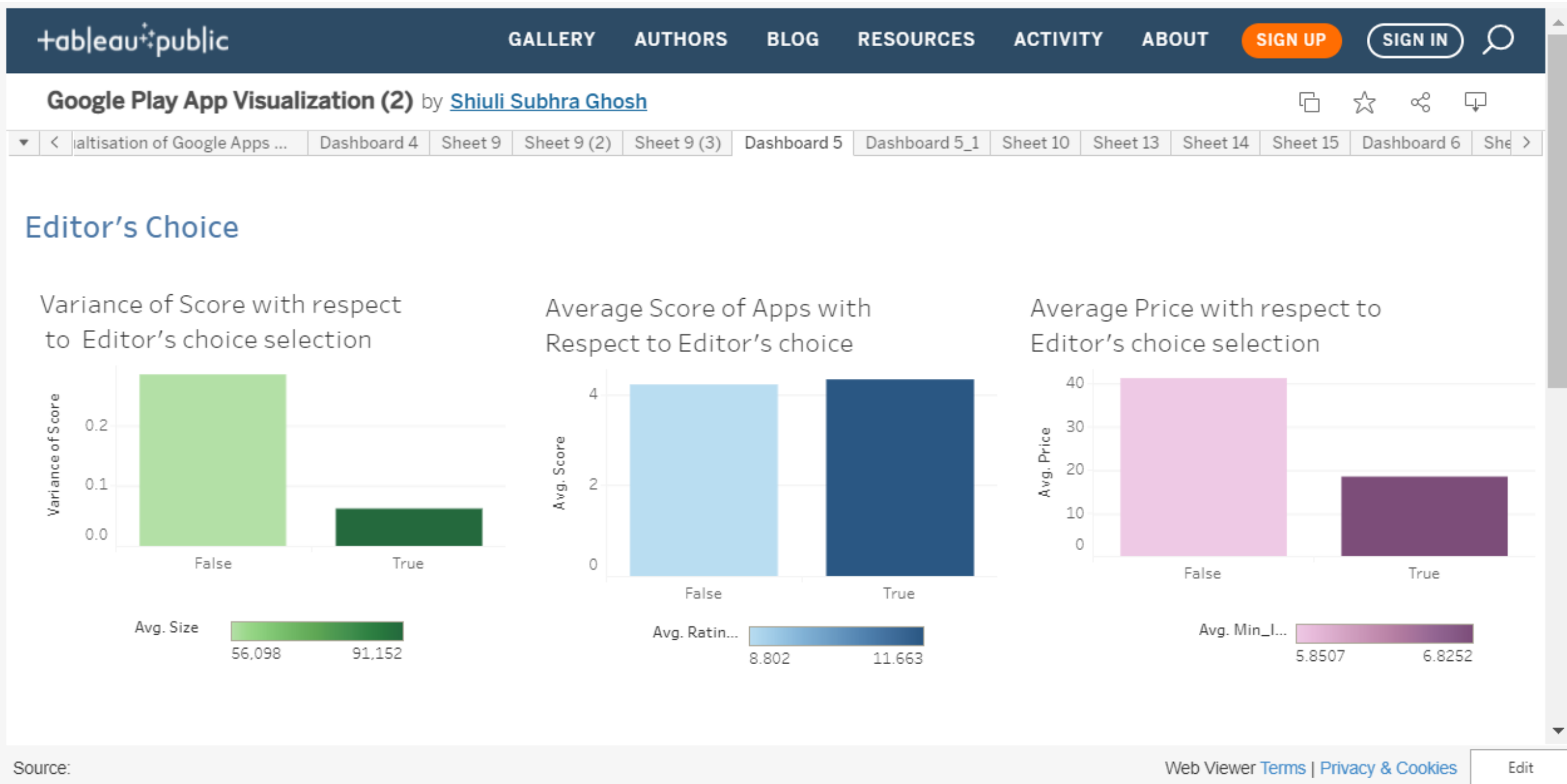
# Distribution of Apps based on Content rating (Everybody/Adult)



# Category wise Content Rating Description Visualization

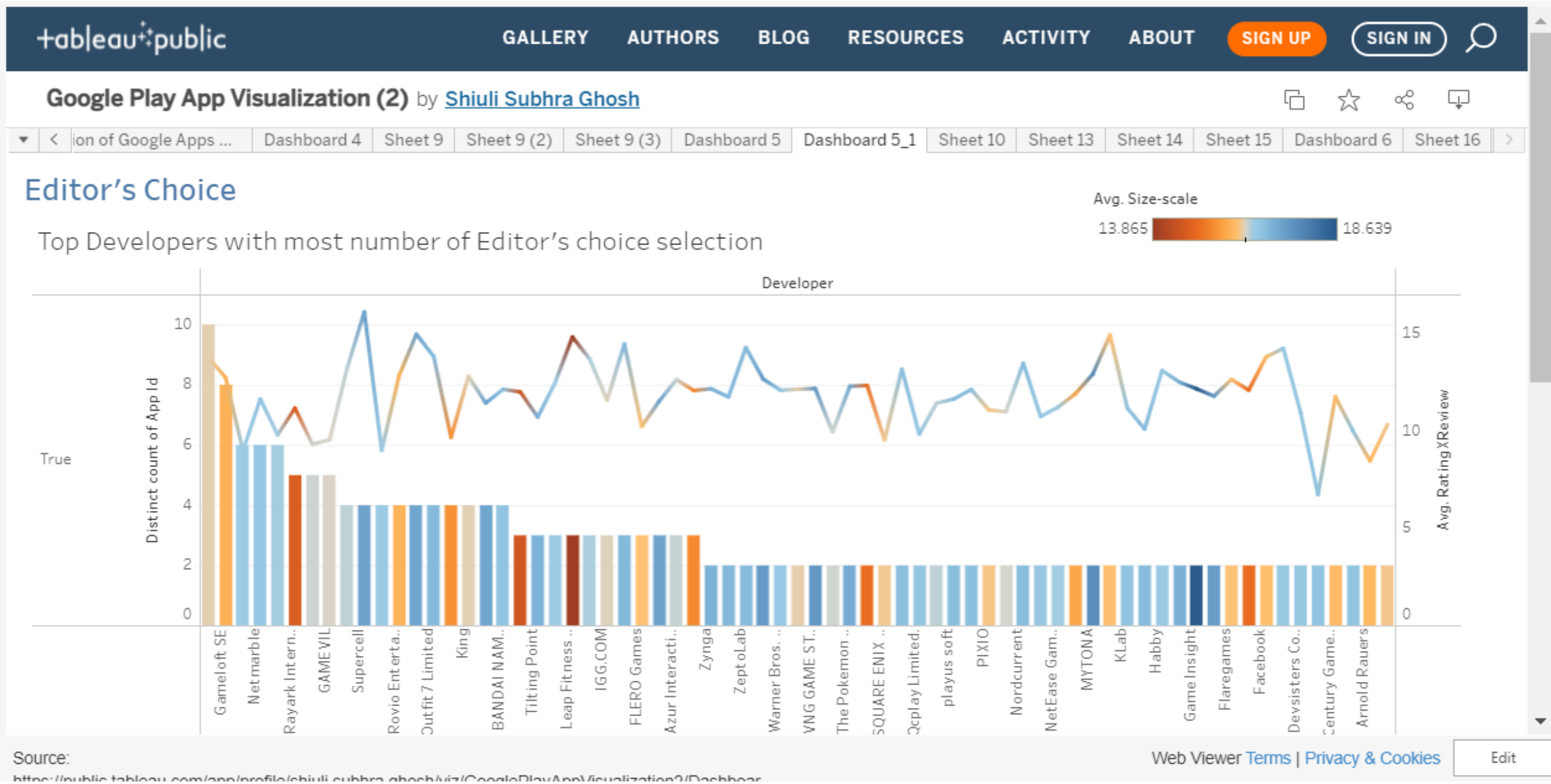


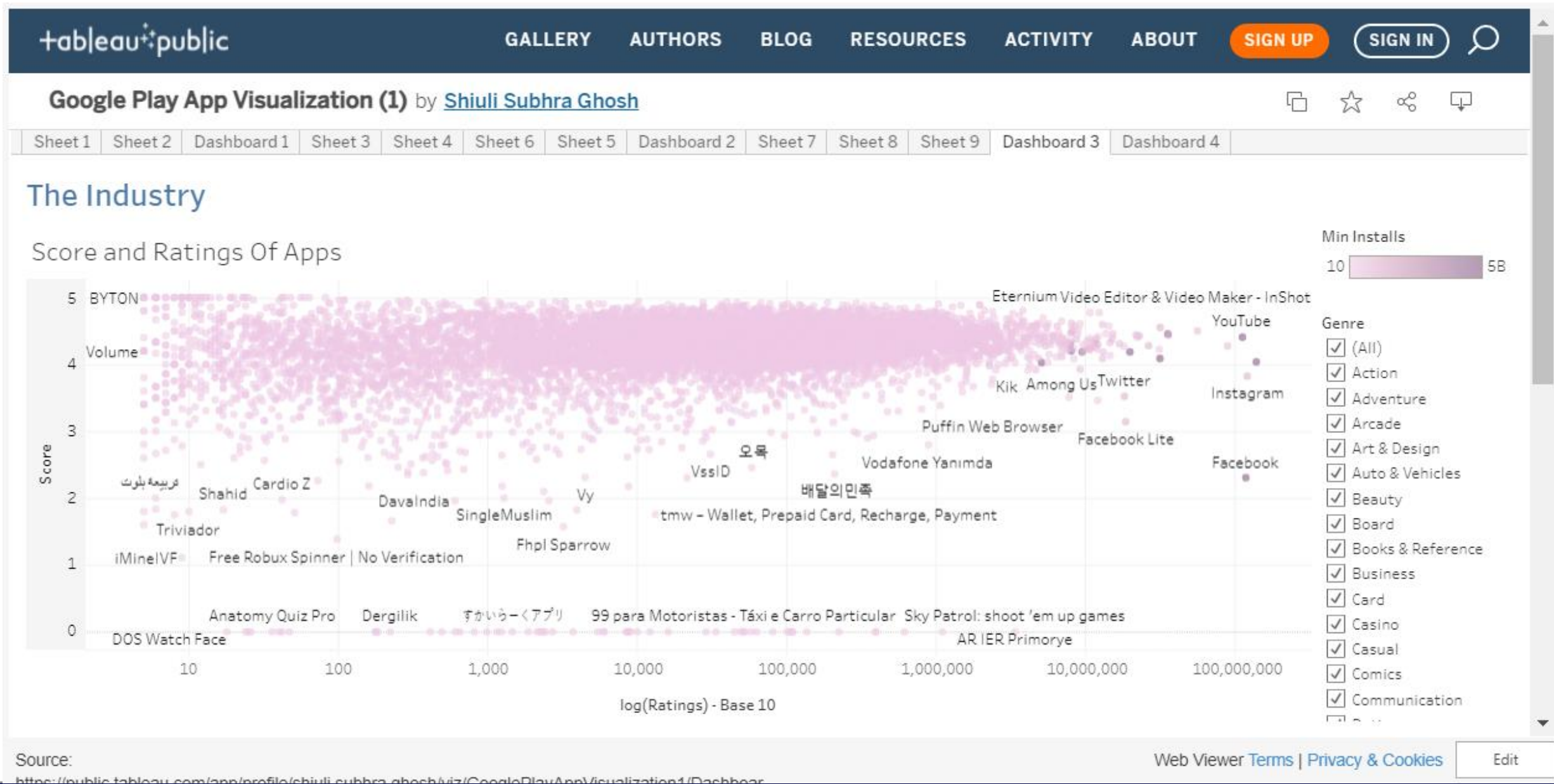
# Visualizing Editor's Choice

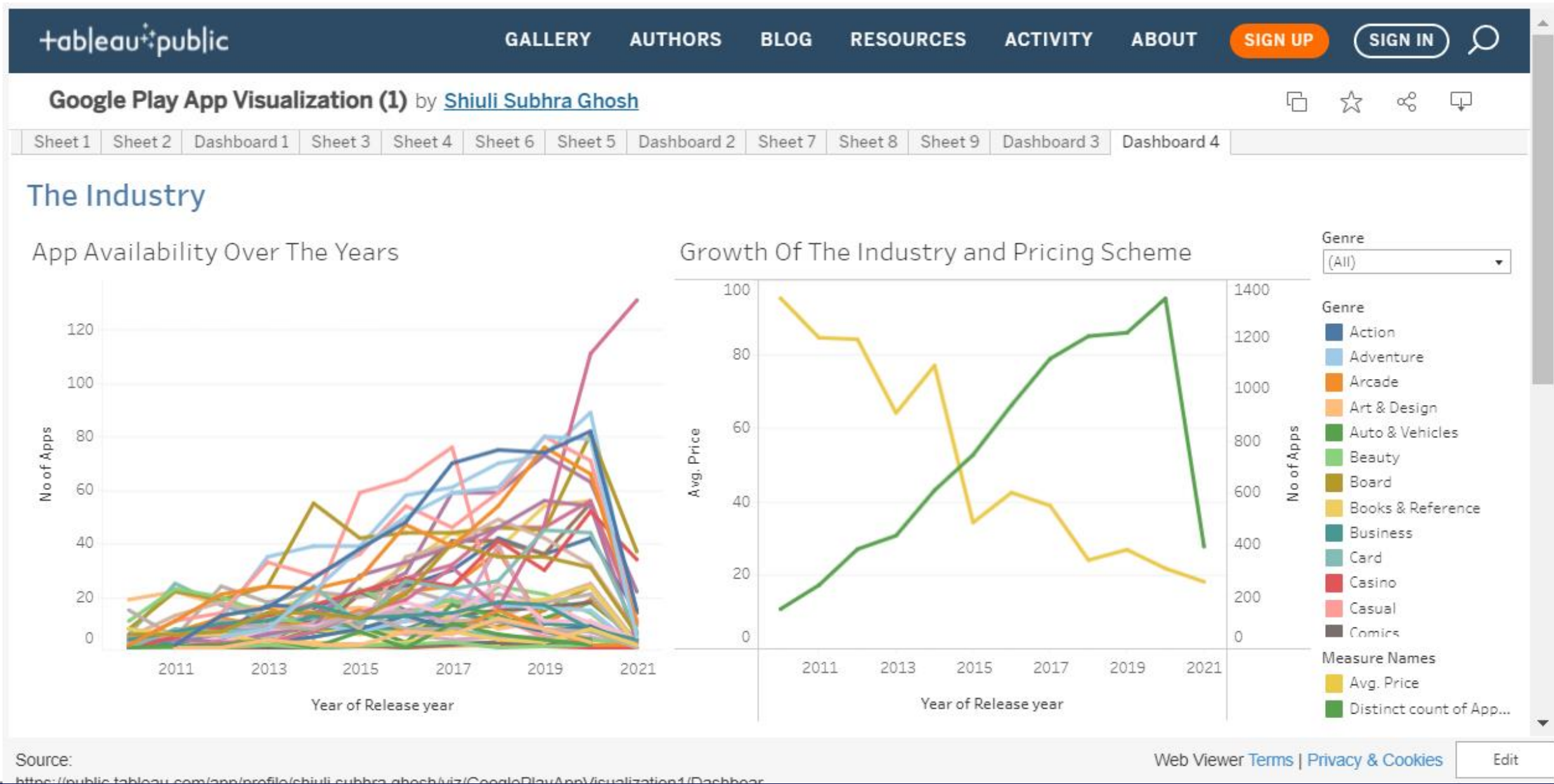




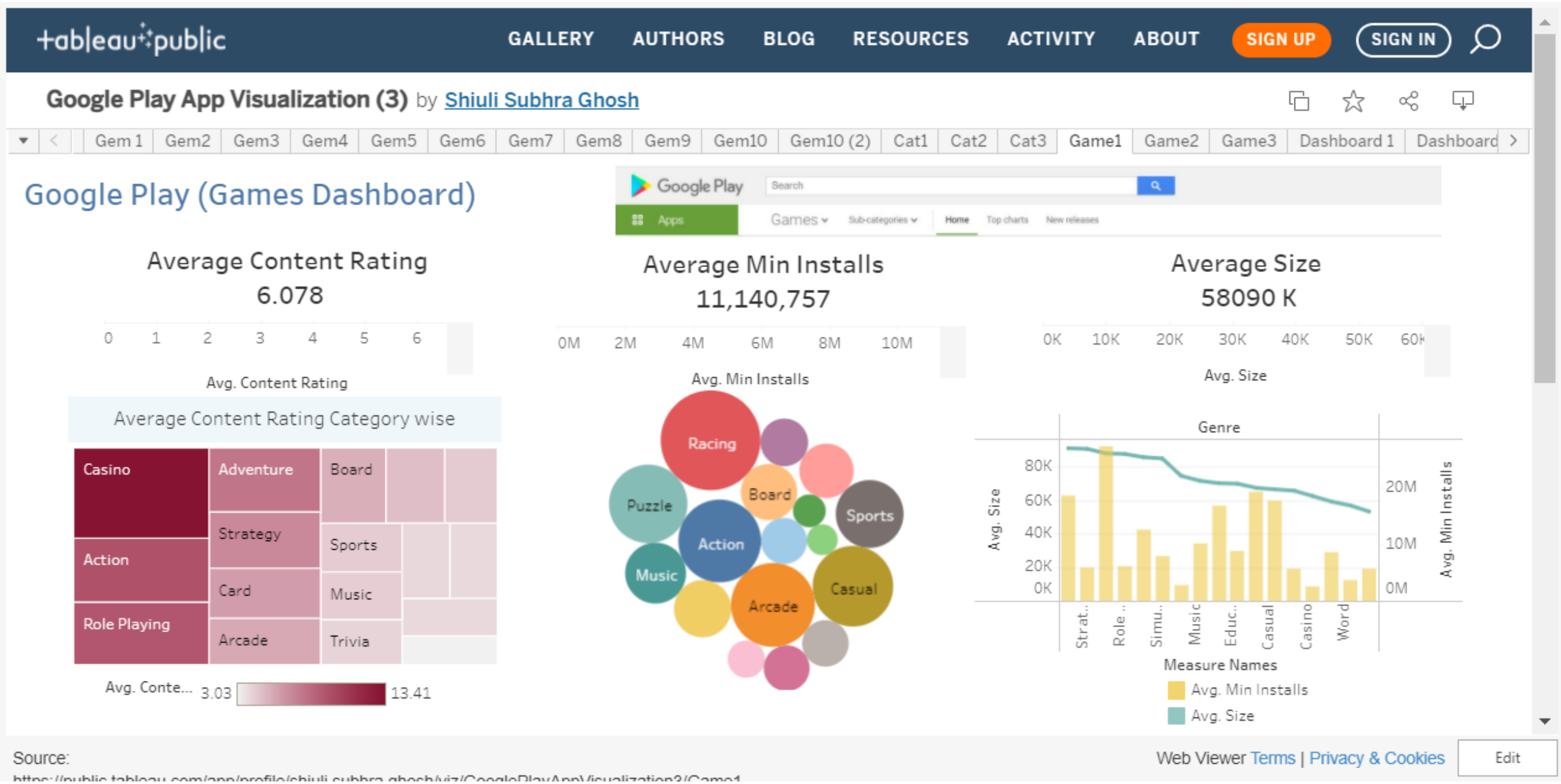
# Visualizing Editor's Choice



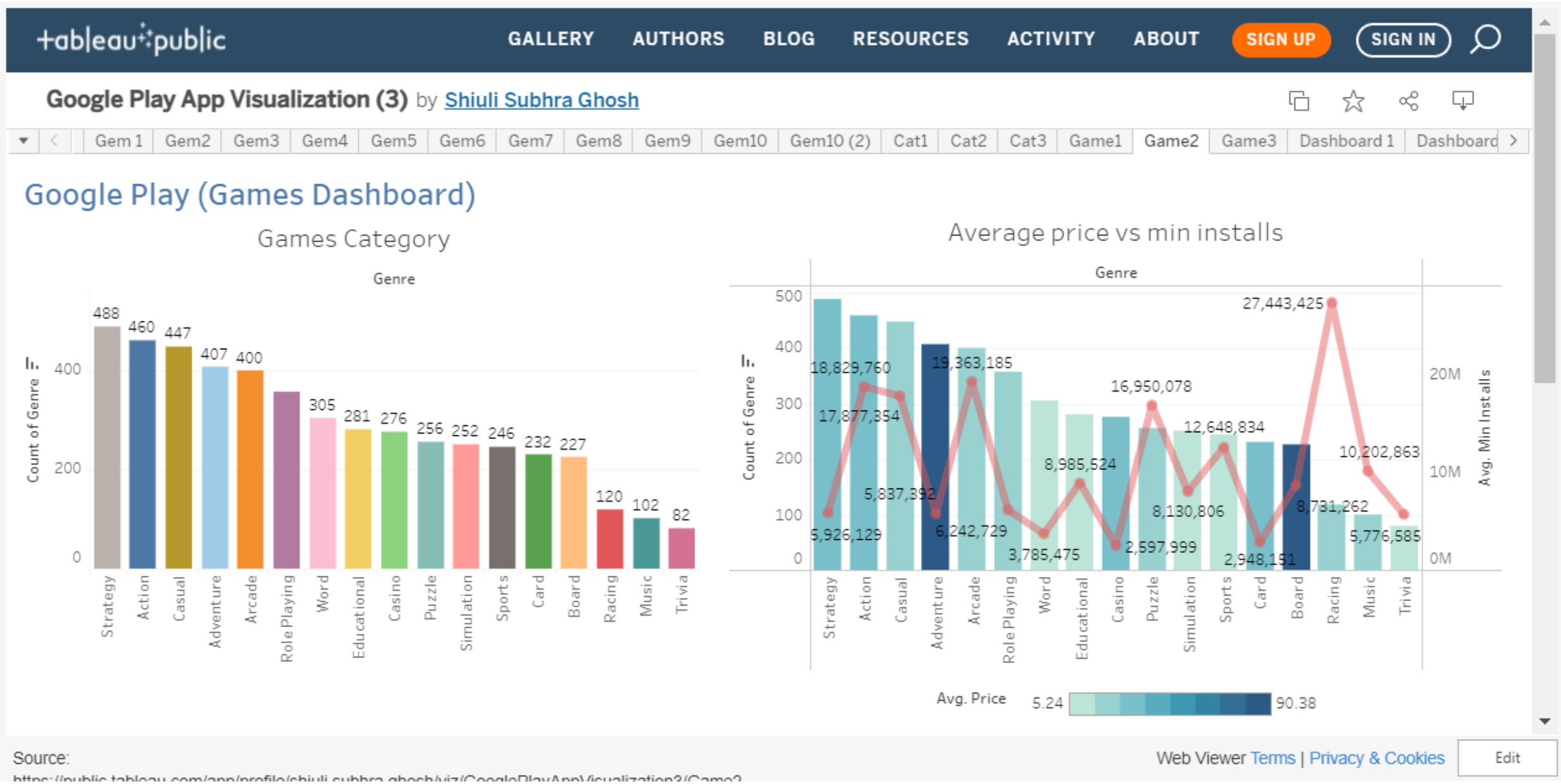




# Comprehensive Game Dashboard



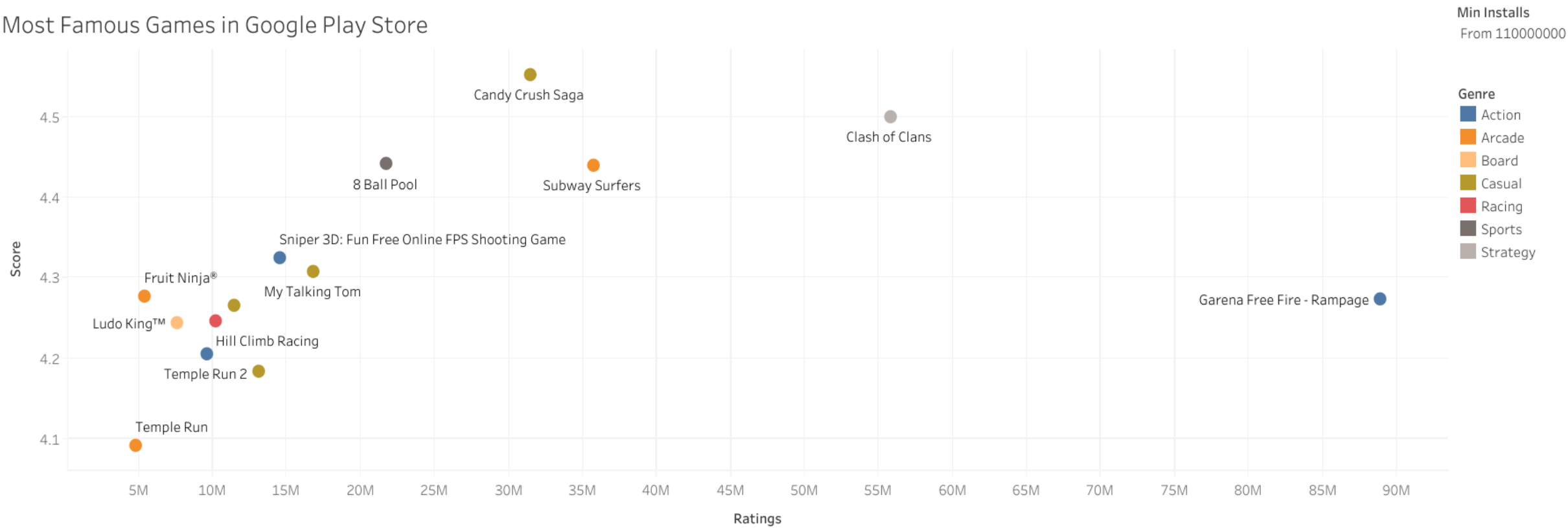
# Game Dashboard Category wise Visualization



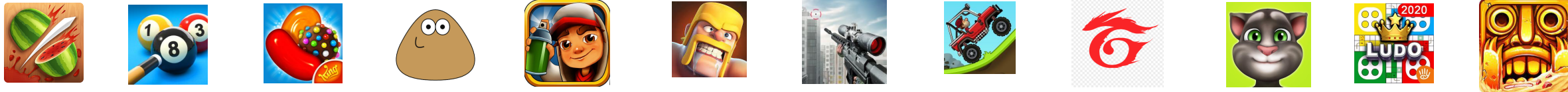


## Famous games

Most Famous Games in Google Play Store



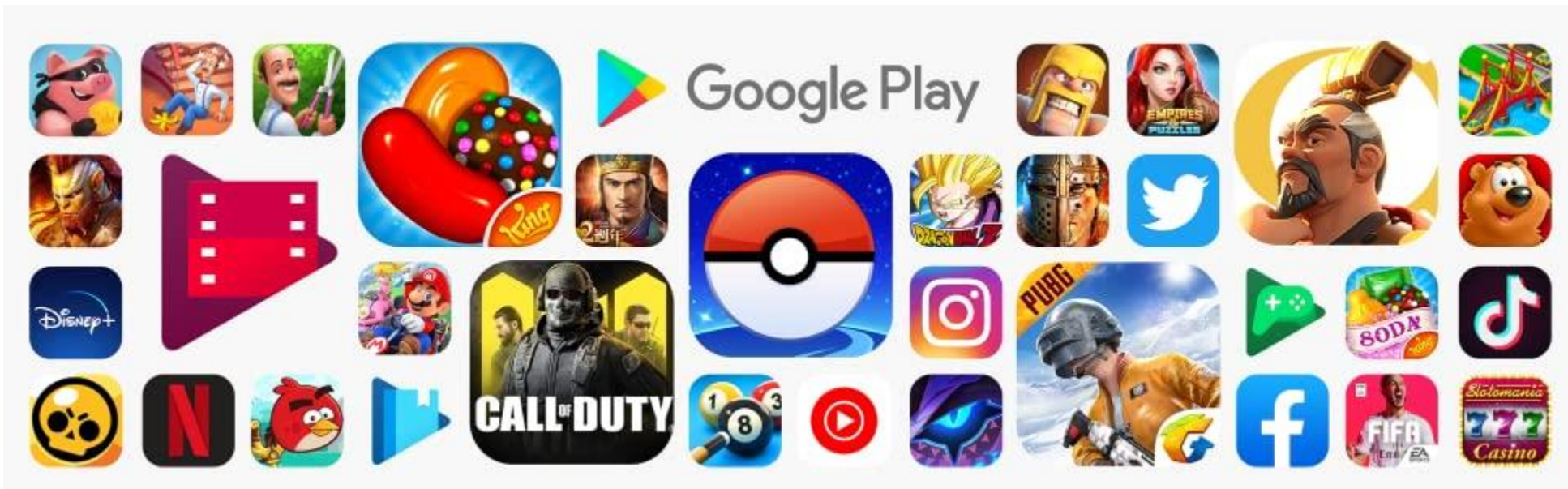
# Famous Games



## Most Famous Games Statistics

Genre	Title	Min Installs	Content Rating	Price	Score	Size
Action	Temple Run 2	500,000,000	7	0	4	133,120
	Sniper 3D: Fun Free Online FPS Shooting G..	500,000,000	18	0	4	139,264
	Garena Free Fire - Rampage	500,000,000	12	0	4	47,104
Arcade	Subway Surfers	1,000,000,000	7	0	4	191,488
	Temple Run	500,000,000	3	0	4	47,104
	Fruit Ninja®	500,000,000	3	0	4	128,000
Board	Ludo King™	500,000,000	3	0	4	50,176
Casual	Candy Crush Saga	1,000,000,000	3	0	5	46,080
	Pou	500,000,000	3	0	4	22,528
	My Talking Tom	500,000,000	3	0	4	116,736
	My Talking Angela	500,000,000	3	0	4	117,760
Racing	Hill Climb Racing	500,000,000	3	0	4	60,416
Sports	8 Ball Pool	500,000,000	12	0	4	68,608
Strategy	Clash of Clans	500,000,000	7	0	4	178,176

- What are the top categories on Google Play Store?
- Percentage of free and paid apps.
- Distribution of the content ratings of the apps.
- How does the size of an app influence installation of an app?
- Genre wise analysis of different numerical factors.
- How is an app selected in "Editor's Choice" different from the rest?
- How did the count of apps progress yearly for different genres?



Thank You