

Play Store App Review Analysis EDA Capstone Project

Presented By

Abhishek Solanki

Solanki12abhi@gmail.com

Cohort Istanbul

AlmaBetter

WHY ANALYZE THE GOOGLE PLAY STORE?



Mobile App Market is set to grow 20% by 2023



Android Apps comprise 90% of the Mobile App Market



What makes an App popular? Can we predict how popular it's going to be?



What are some interesting patterns in user behavior related to app usage & feedback

Introduction

- Android is the most popular operating system in the world, with over 2.5 billion active users spanning over 190 countries.
- Google Play was launched on March 6, 2012, bringing together Android Market marking a shift in
- Google's digital distribution strategy .
- Android is the dominant mobile operating system today more than 85% of all mobile devices running Google's OS. The Google Play Store is the largest and most popular Android app store.
- There are more than 3.04 million apps found on Google Play Store.
- The Play Store apps data has enormous potential to drive app-making businesses to success.
- Actionable insights can be drawn for developers to work on and capture the Android market. The main goal of our project is-
 - The purpose of our project is to gather and analyze detailed information on apps in the Google Play Store in order to provide insights on app features and the current state of the Android app market.
 - The Objective of the project to Explore and analyze the data to discover key factors responsible for app engagement and success.

PROBLEM STATEMENT

- * Two datasets are provided, one with **basic information** and the other with **user reviews** for the respective app.
- * We must examine and evaluate the data in both datasets in order to identify the important characteristics that influence app engagement and success.

So, what factors influence an app's success?

An app is said to be successful if it has:

- * A high average user rating
- * A good number of positive reviews
- * A good number of monthly average users
- * High revenue per customer and so on.

AGENDA

- ☐ Introduction
- ☐ Category wise play store apps installs
- ☐ Category wise most popular apps
- ☐ Top 10 apps in play store considering all the parameters
- ☐ Average installs, category wise
- ☐ Most installed apps in communication category
- ☐ Average sizes of apps in each category
- ☐ Category wise percentage of paid apps
- ☐ Category wise top installed paid apps
- ☐ Average rating of paid apps
- ☐ Correlation between Rating ,Installs and Price
- ☐ Category wise installed apps with content rating
- ☐ Percentage reviews sentiment distribution

DATASET PREPARATION

- **Loading the data sets:** Two datasets, First Play store app dataset and User Reviews dataset.
- **Import Libraries:** NumPy, Pandas, Seaborn and Matplotlib
- **Data cleaning:** Null values, Finding and removing Outliers, Removing duplicate data.
- **Data Imputation:** Filling the missing categorical values with mode and numerical values with median. Conversion of price, installs, reviews into numerical values.
- **Exploratory Data Analysis:** Analyzing the data sets to summarize their main characteristics using statistical graphics and data visualizations method.

ATTRIBUTES IN PLAY STORE DATA

1. **APP**: THIS COLUMN CONTAINS THE NAME OF THE APP FOR EACH OBSERVATION.
2. **CATEGORY**: THIS COLUMN CONTAINS CATEGORY TO WHICH THE APP BELONGS.
3. **RATING**: THIS COLUMN CONTAINS THE AVERAGE RATING FOR THE APP.
4. **REVIEWS**: THIS COLUMN CONTAINS THE NUMBER OF REVIEWS THAT THE APP HAS RECEIVED ON THE PLAY STORE.
5. **SIZE**: THIS COLUMN CONTAINS THE AMOUNT OF MEMORY THE APP OCCUPIES ON THE DEVICE.
6. **INSTALLS**: THIS COLUMN CONTAINS THE NUMBER OF TIMES THAT THE APP HAS BEEN DOWNLOADED AND INSTALLED FROM THE PLAY STORE.
7. **TYPE**: THIS COLUMN CONTAINS THE INFORMATION WHETHER THE APP IS FREE OR PAID.
8. **PRICE**: IF THE APP IS A PAID APP, THIS COLUMN CONTAINS THE DATA ABOUT ITS PRICE.
9. **CONTENT RATING**: THIS COLUMN CONTAINS THE MATURITY RATING OF THE APP I.E. THE AGE GROUP OF THE AUDIENCE FOR WHICH IT IS SUITABLE.
10. **GENRES**: THIS COLUMN CONTAINS THE DATA ABOUT TO WHICH GENRE THE APP BELONGS. GENRES CAN BE CONSIDERED AS A FURTHER DIVISION OF THE GROUP OF CATEGORY.
11. **LAST UPDATED**: CONTAINS THE DATE ON WHICH THE LATEST UPDATE OF THE APP WAS RELEASED.
12. **CURRENT VERSION**: CONTAINS INFORMATION ON THE CURRENT VERSION OF THE APP AVAILABLE ON THE PLAY STORE.
13. **ANDROID VERSION**: CONTAINS INFORMATION ABOUT THE ANDROID VERSIONS ON WHICH THE APP IS SUPPORTED.

ATTRIBUTES IN USER REVIEWS

1. **App-** Application name
2. **Translated Review-** User review
3. **Sentiment-** Positive/Negative/Neutral
4. **Sentiment Polarity-** Sentiment polarity score
5. **Sentiment Subjectivity-** Sentiment subjectivity score

OVERVIEW OF ANALYSIS

Data Cleaning



Understand the structure of the dataset and clean data before analysis

Data Exploration



Uncover initial patterns, characteristics, and points of interest using visual exploration

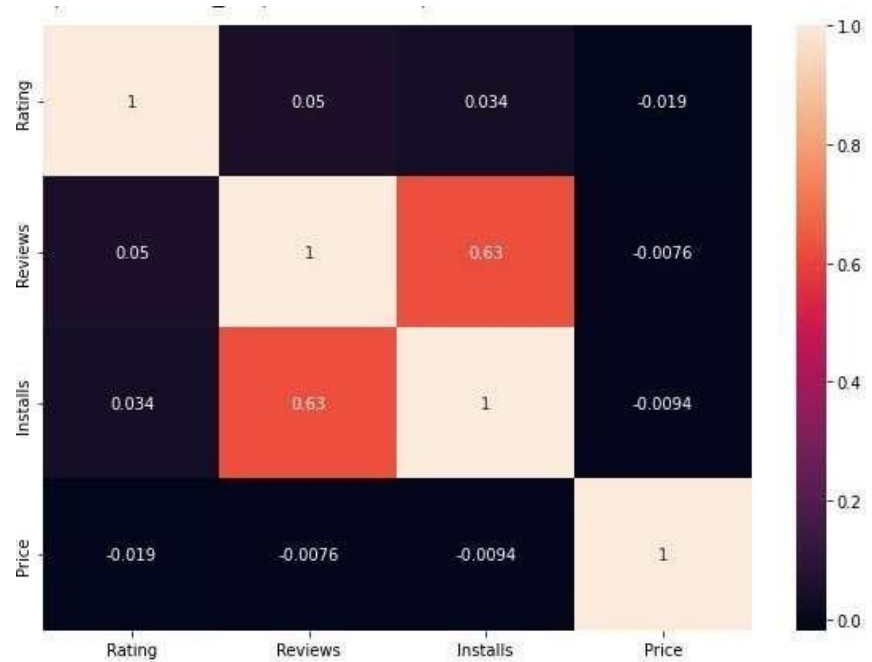
Predictive Modeling



Formulate a statistical model to forecast an outcome using relevant predictors

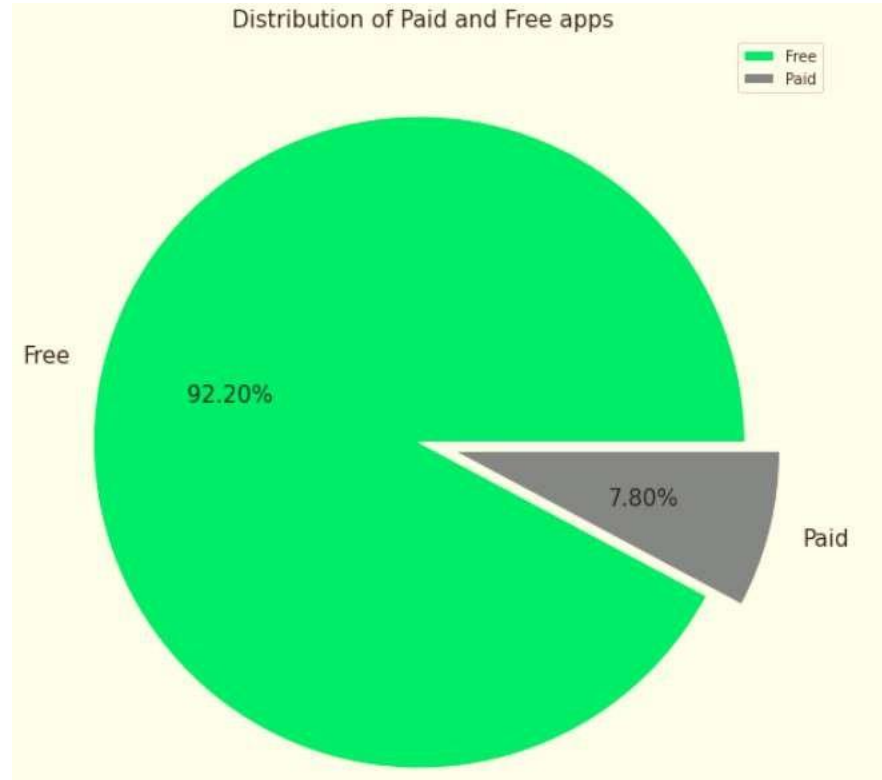
CORRELATION HEATMAP

- There is a strong **positive** correlation between the **Reviews** and **Installs**.
- The **Price** is slightly **negatively** correlated with the **Rating**, **Reviews**, and **Installs**.
- The **Rating** is slightly **positively** correlated with the **Installs** and **Reviews**.



Percentage of Paid apps v/s Free apps

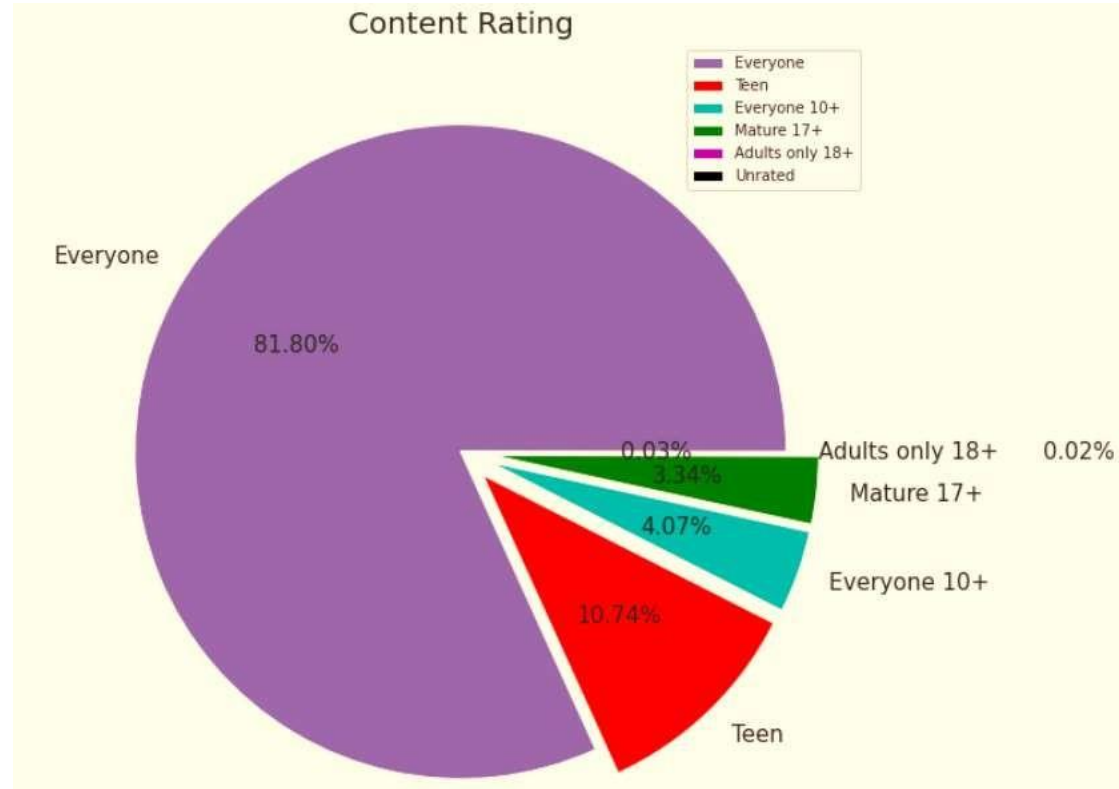
We Observed that **92.20% of Apps are free** and only **7.80% of Apps are paid** in Play store.



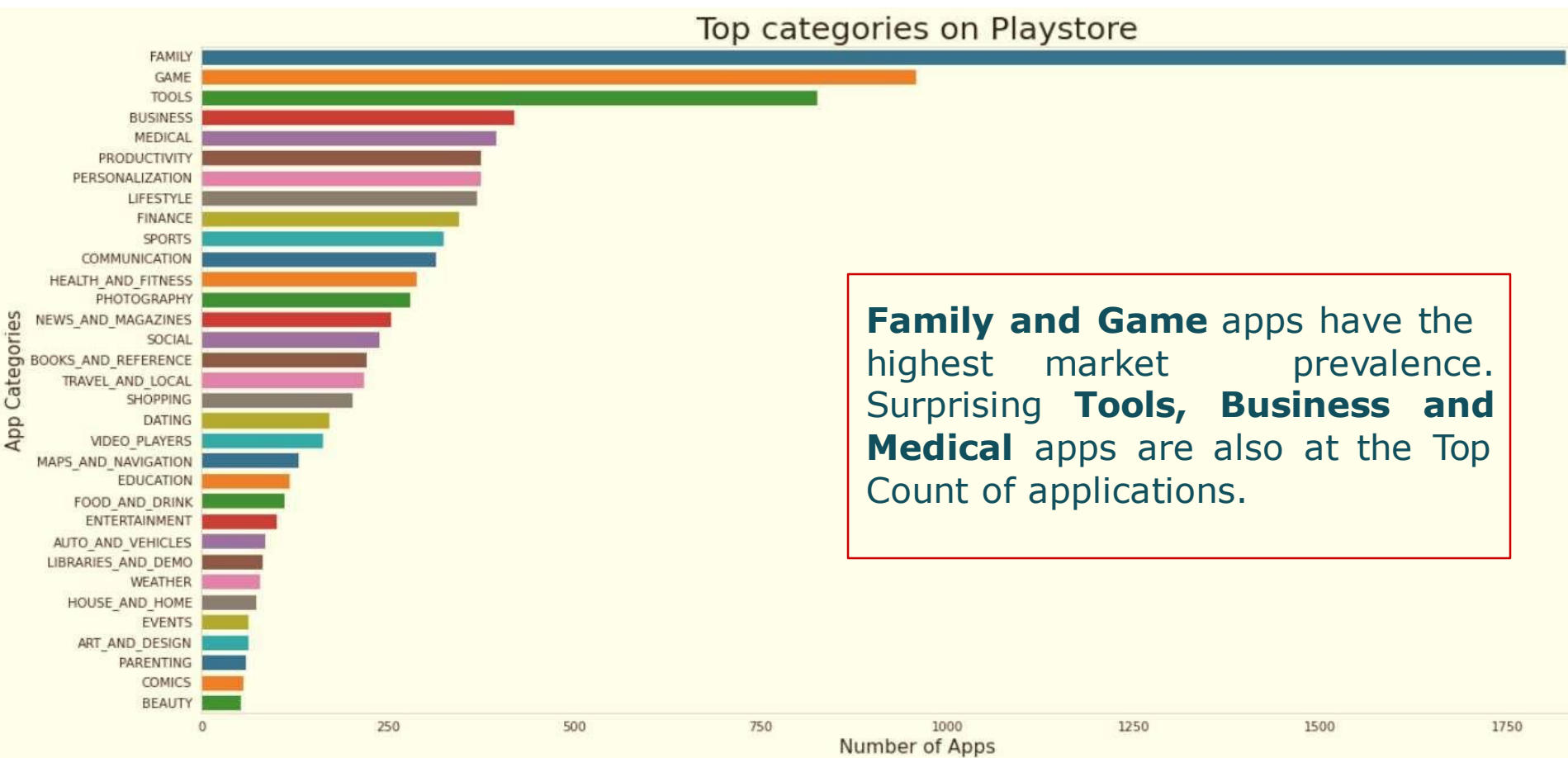
CONTENT RATING

From the above plot we can see that Everyone category having majority of apps count.

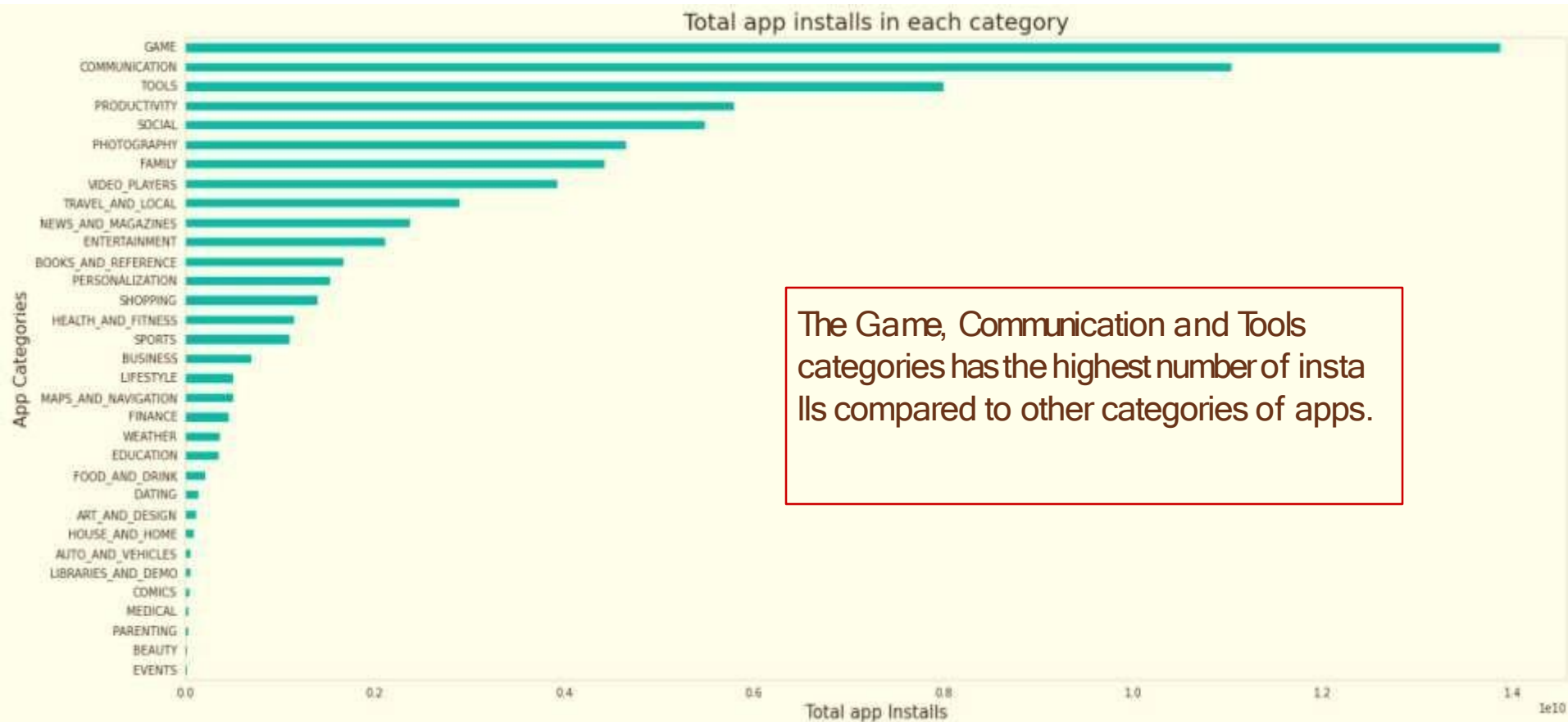
A majority of the apps (**81.80%**) in the play store are can be used by everyone. The remaining apps have various age restrictions to use it.



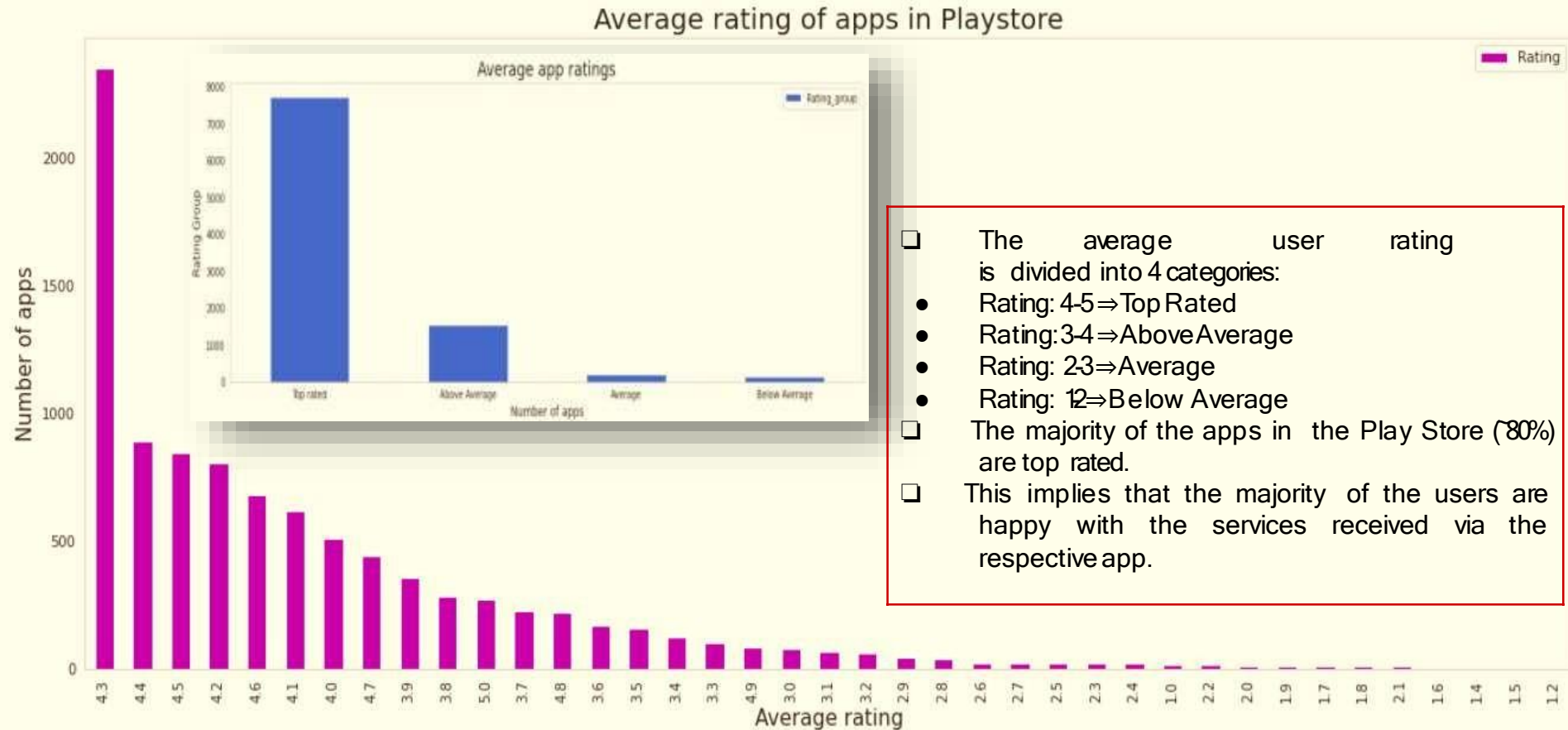
COUNT OF APPLICATION'S IN EACH CATEGORY



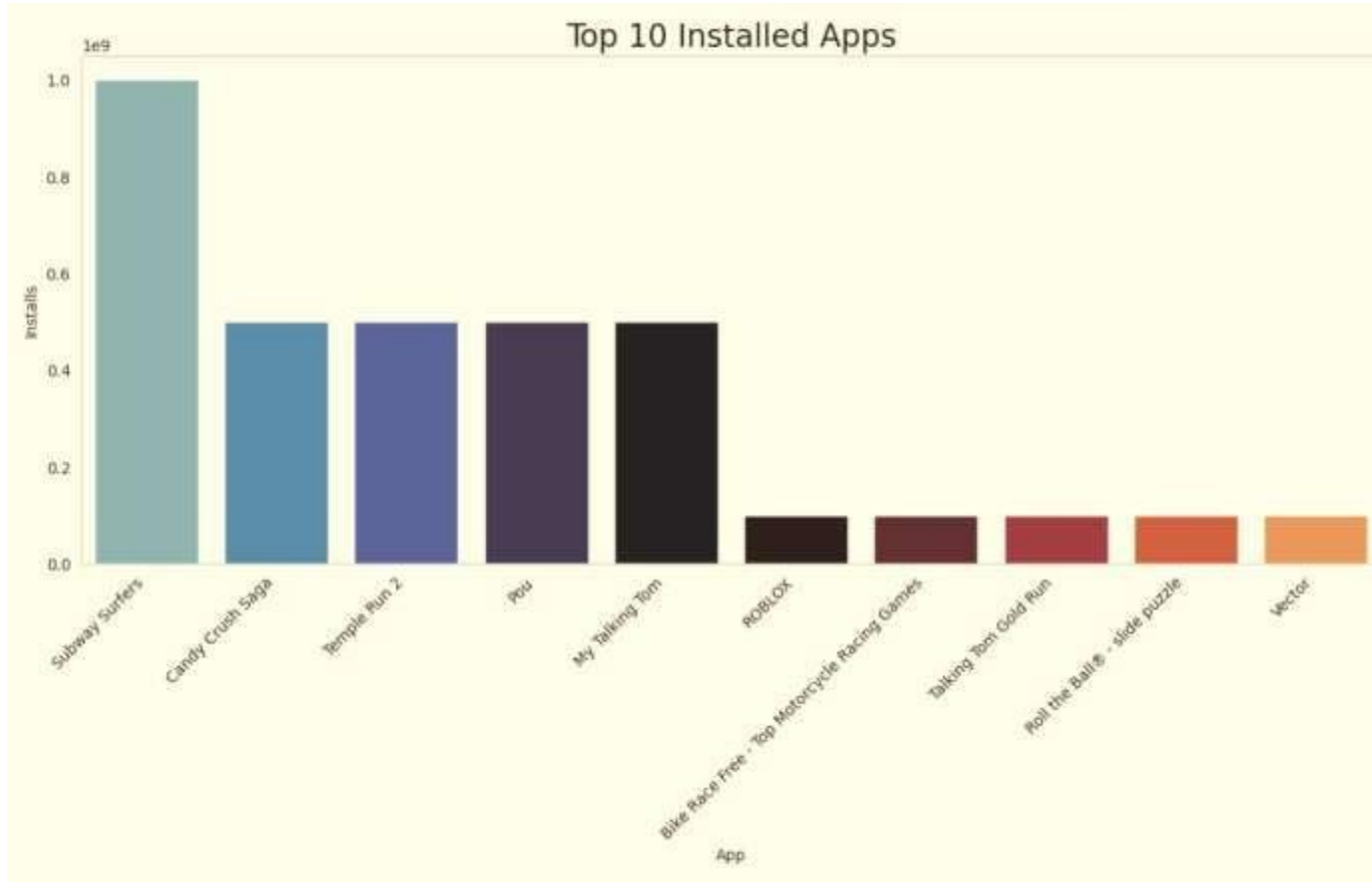
CATEGORY OF APP'S HAVE MOST NUMBER OF INSTALLS



AVERAGE RATING OF APPS IN PLAY STORE



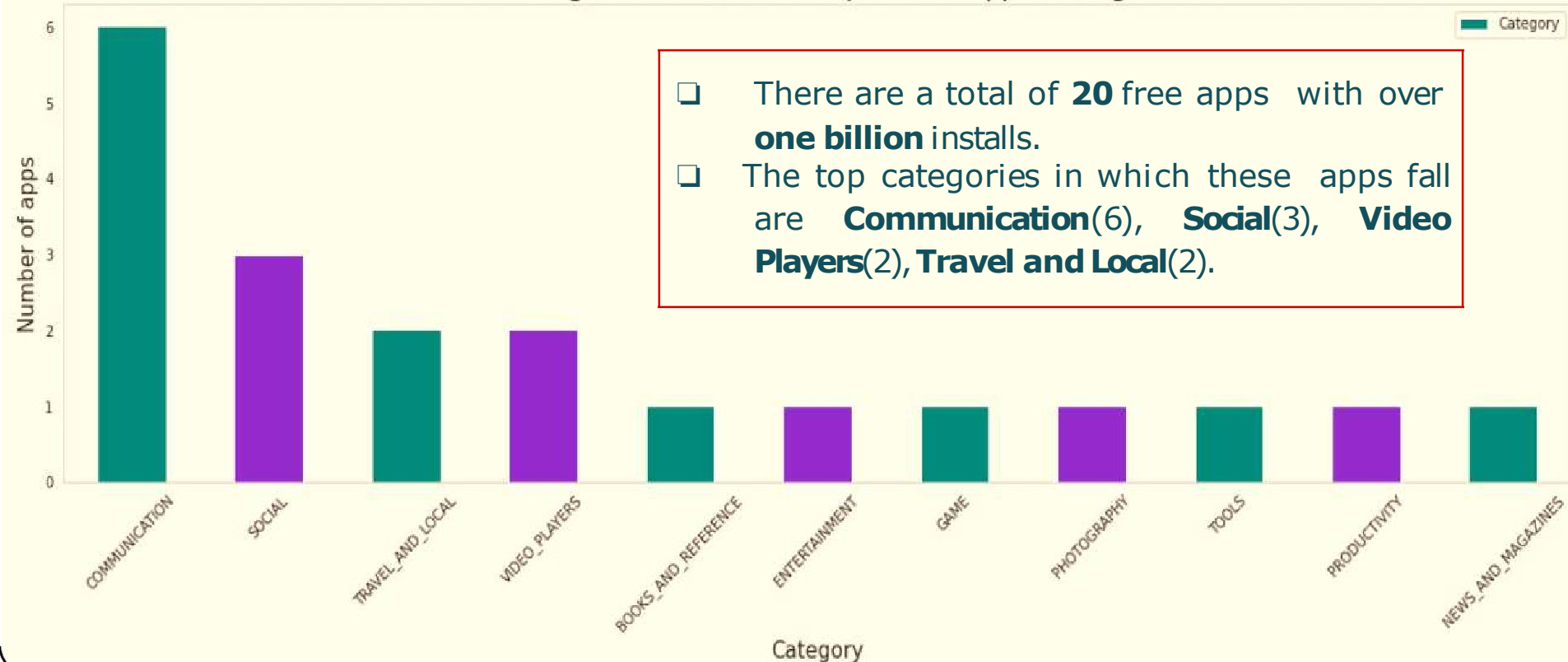
TOP 10 INSTALLS APP'S IN ANY CATEGORY



This graph shows the top installed apps in the '**Games**' category. Further looking into the play store reveals that these apps are light, casual, single player games.

TOP FREE APP'S

Categories in which the top 20 free apps belong

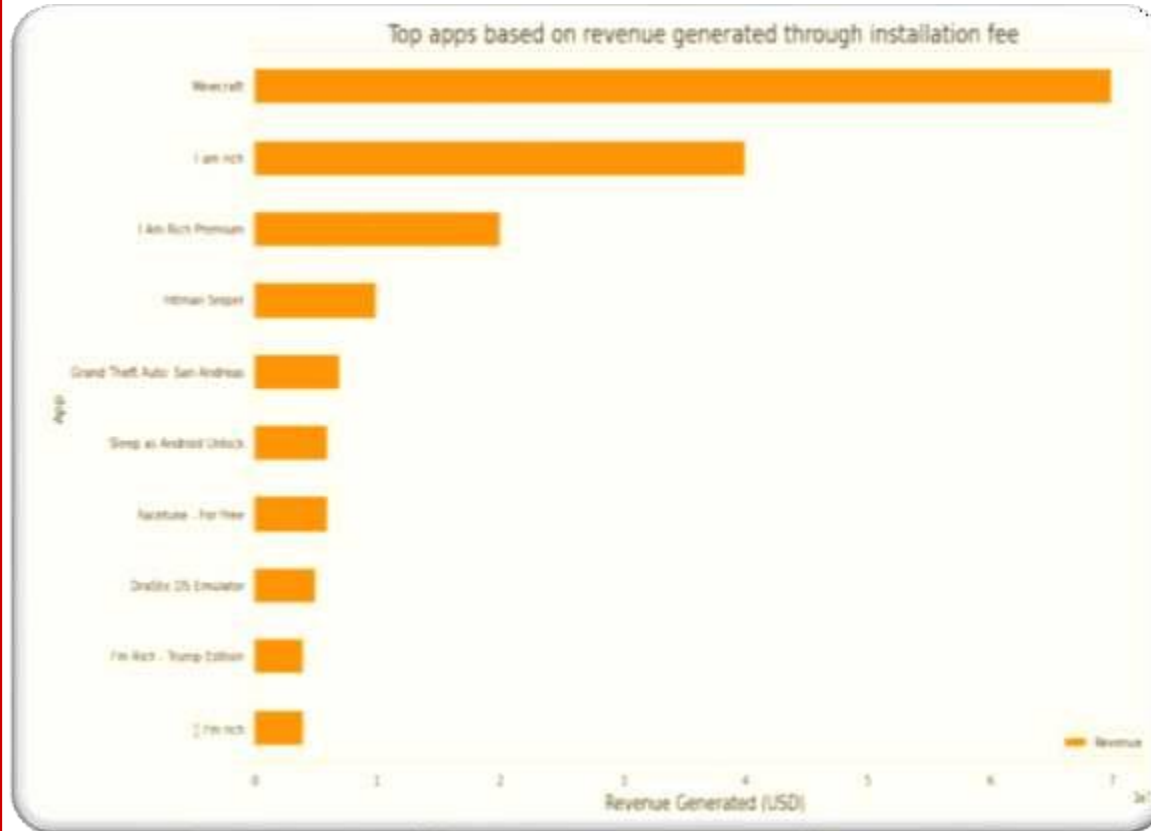


TOP PAID APP'S BASED ON REVENUE GENERATED

- Revenue generated is given by the formula:

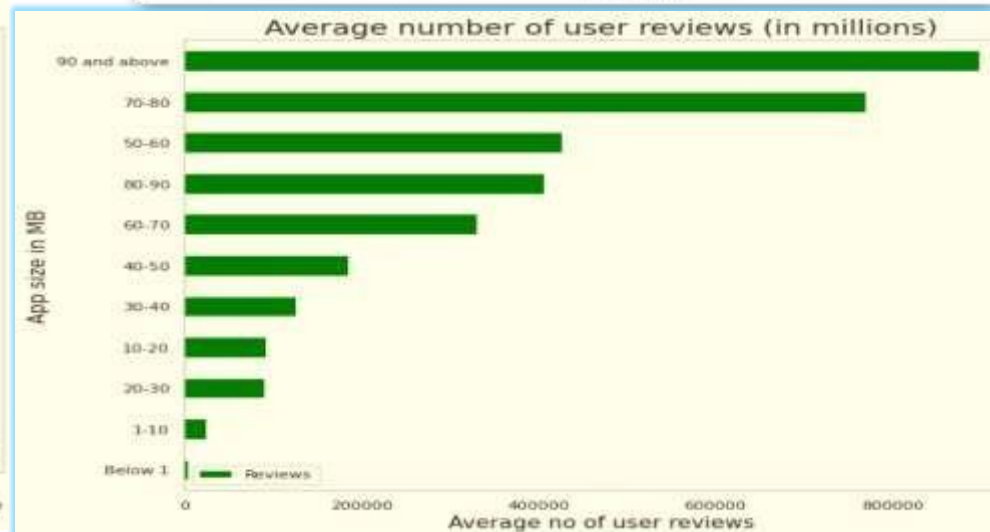
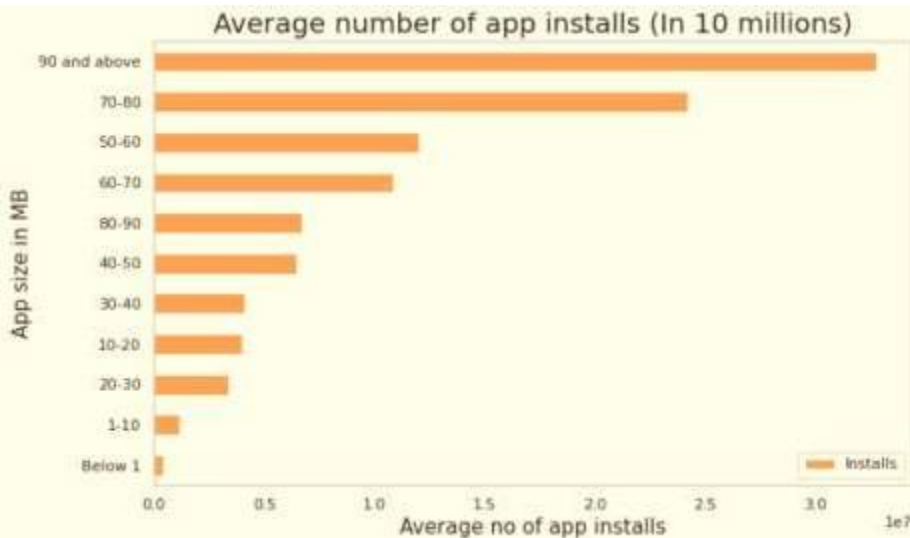
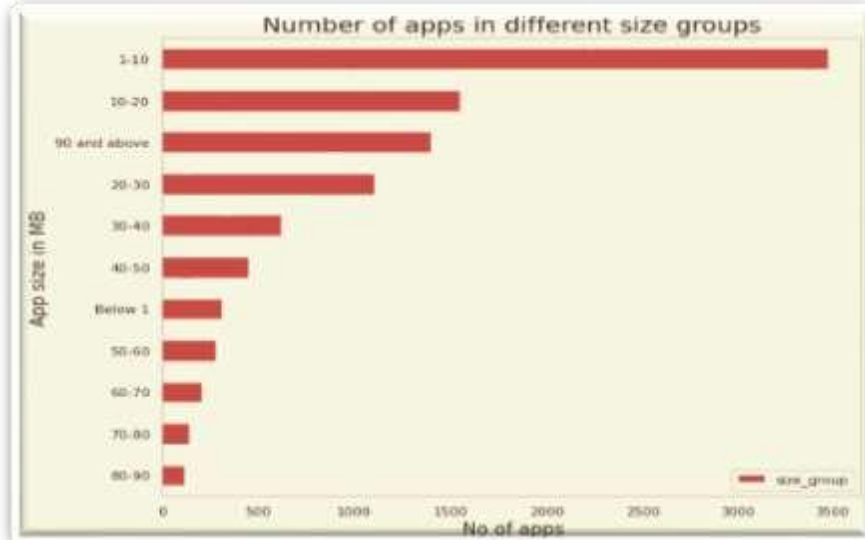
$$\text{Revenue} = \text{Installs} * \text{Price}$$

- Note that in this case, revenue refers to the money earned only from paid appinstalls.
- The top categories in which these apps fall are **Lifestyle**(5), **Family**(5), and **Game**(4).
- Minecraft**, **I am rich**, and **I am rich premium** are the top paid apps based on revenue generated.
- Minecraft** is the only app that has over **10M** installs.

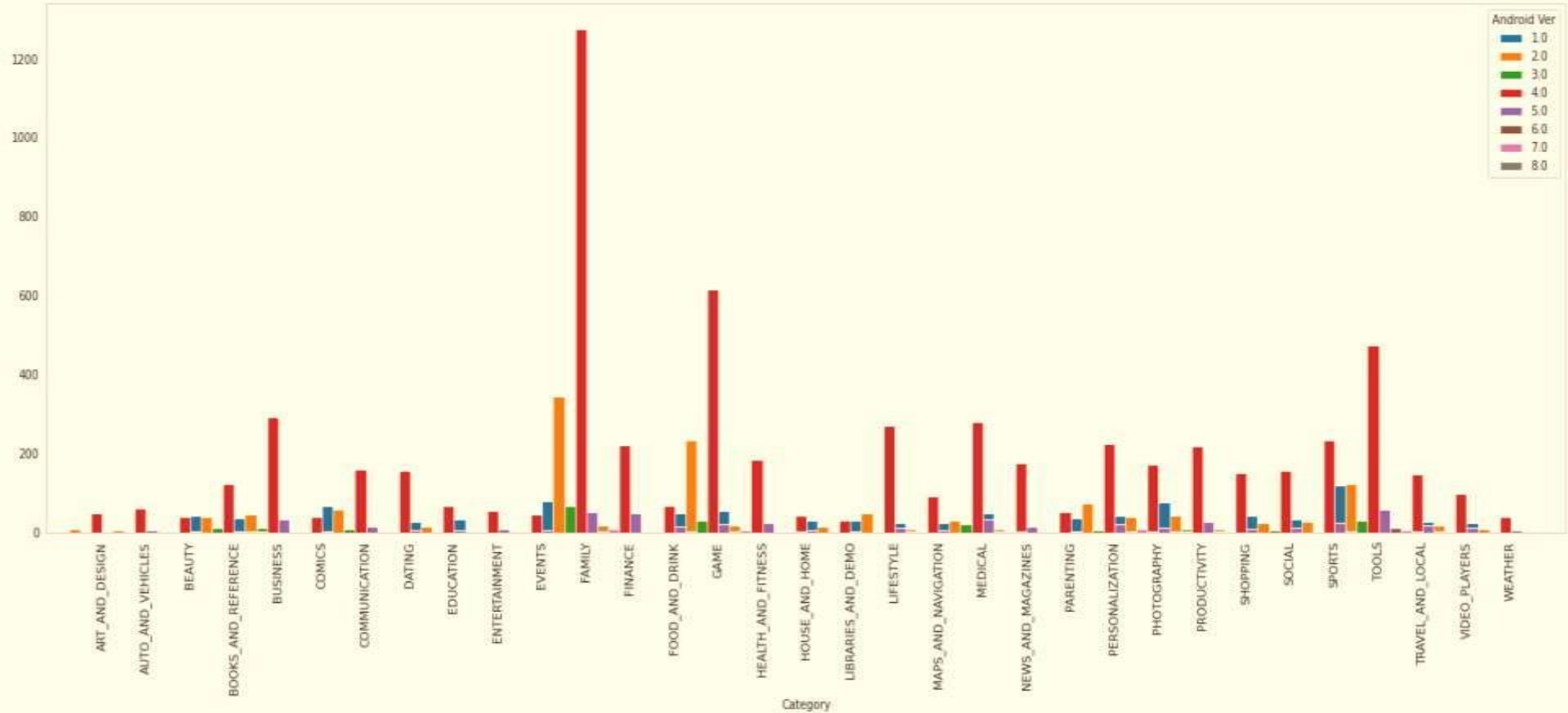


APP SIZE ANALYSIS

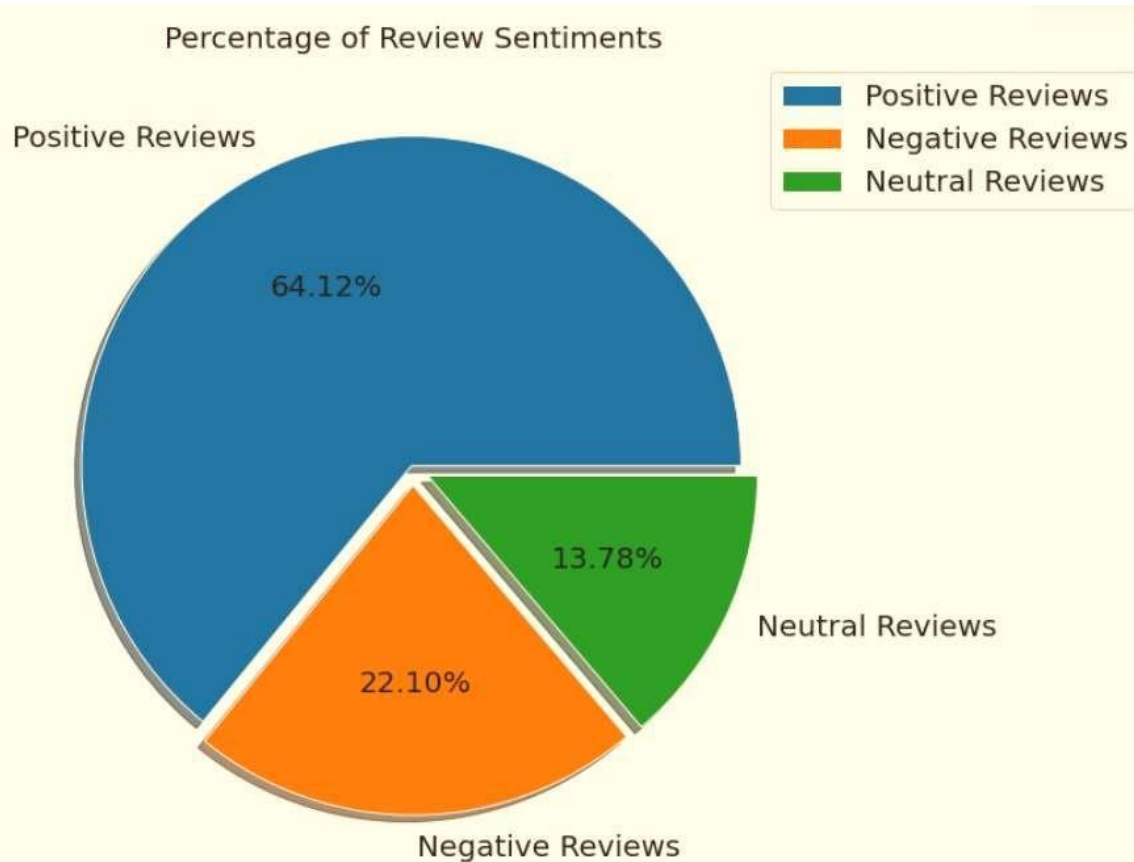
- ❑ The apps are categorized based on its size between ~0 to 100 MB in the intervals of 10 MB each.
- ❑ The total number of apps in each size category indicates the **competition**.
- ❑ Average number of **user reviews** and **average app installs** in each size category indicates the **popularity** of the respective app.



ANDROID VERSION BASED ON EACH CATEGORY



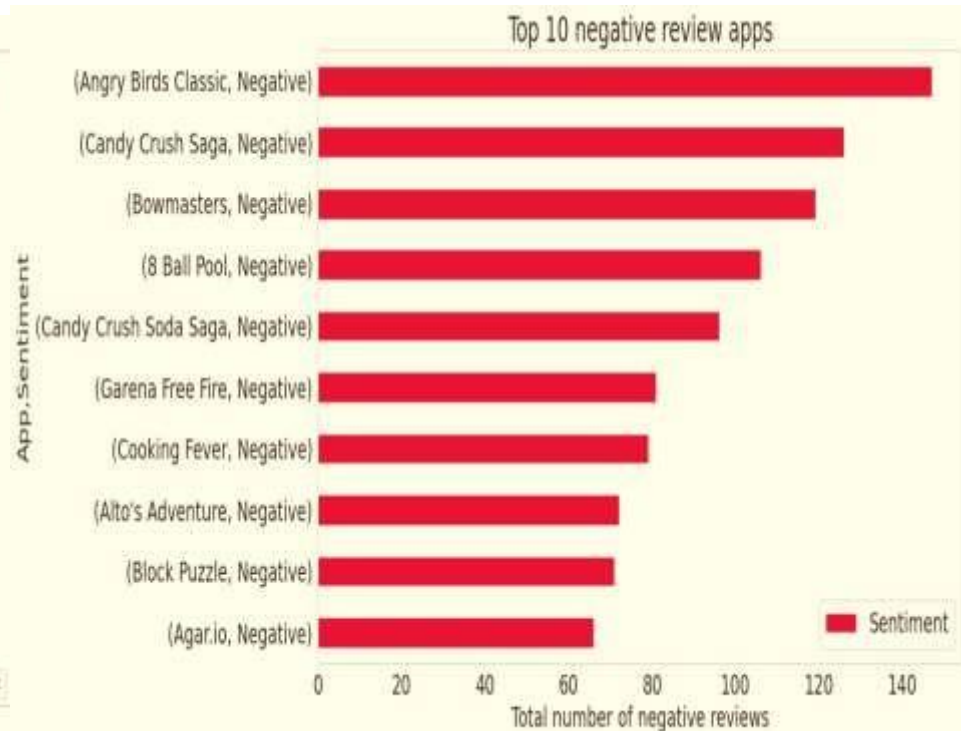
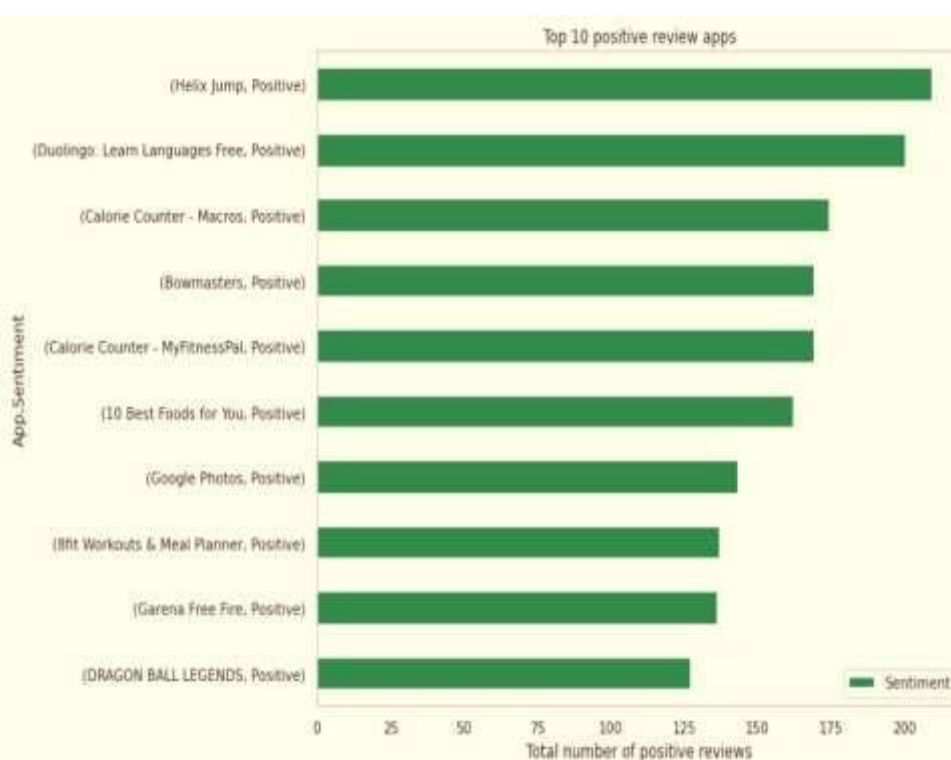
PERCENTAGE OF REVIEW SENTIMENTS



The number of **Unique** Apps from Play store and User reviews merged dataset are **816**.

From Sentiment column, **64%** are **Positive**, **22%** are **Negative** and **14%** are **Neutral** values.

POSITIVE AND NEGATIVE REVIEWS

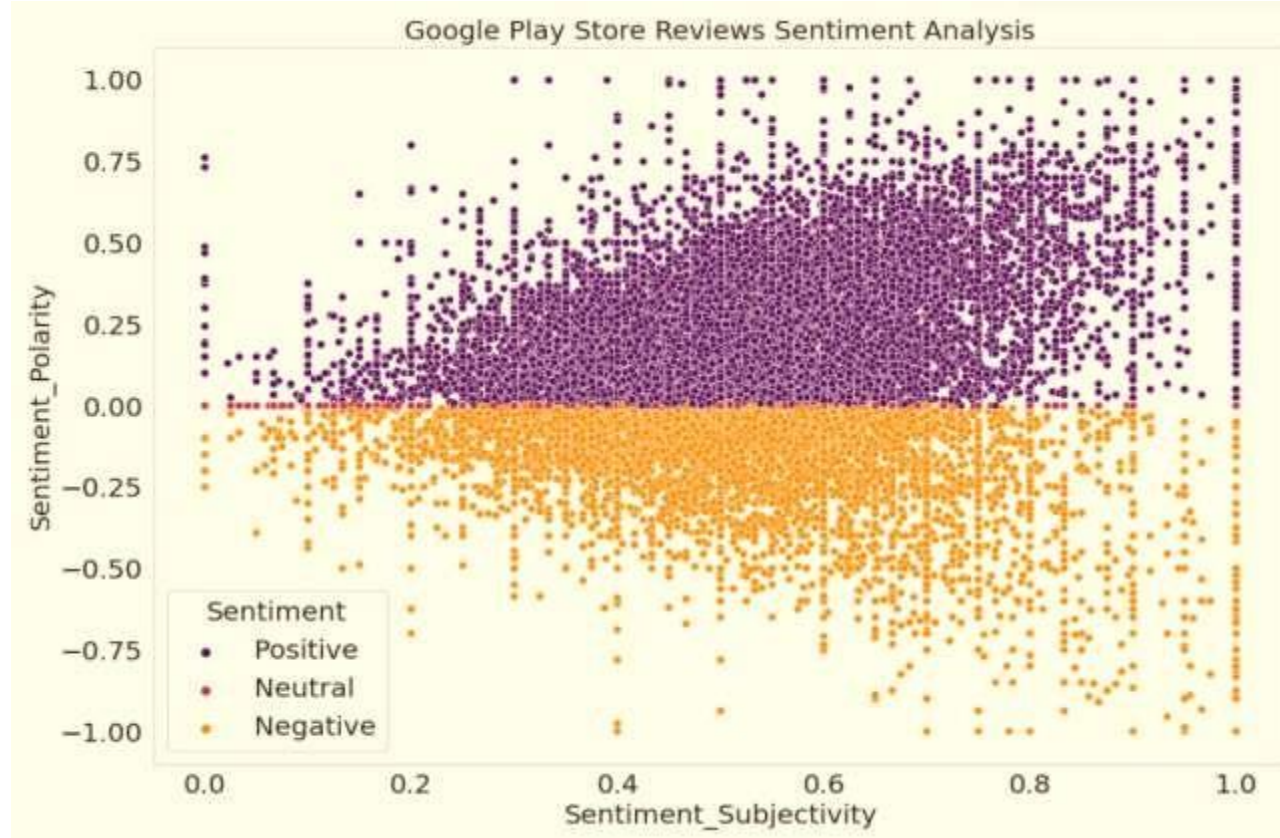


Helix Jump is a App from merged dataset has highest 209 Positive sentimentcount.

Angry Bird Classic is a app from merged dataset has highest **147 Negative** sentimentcount.

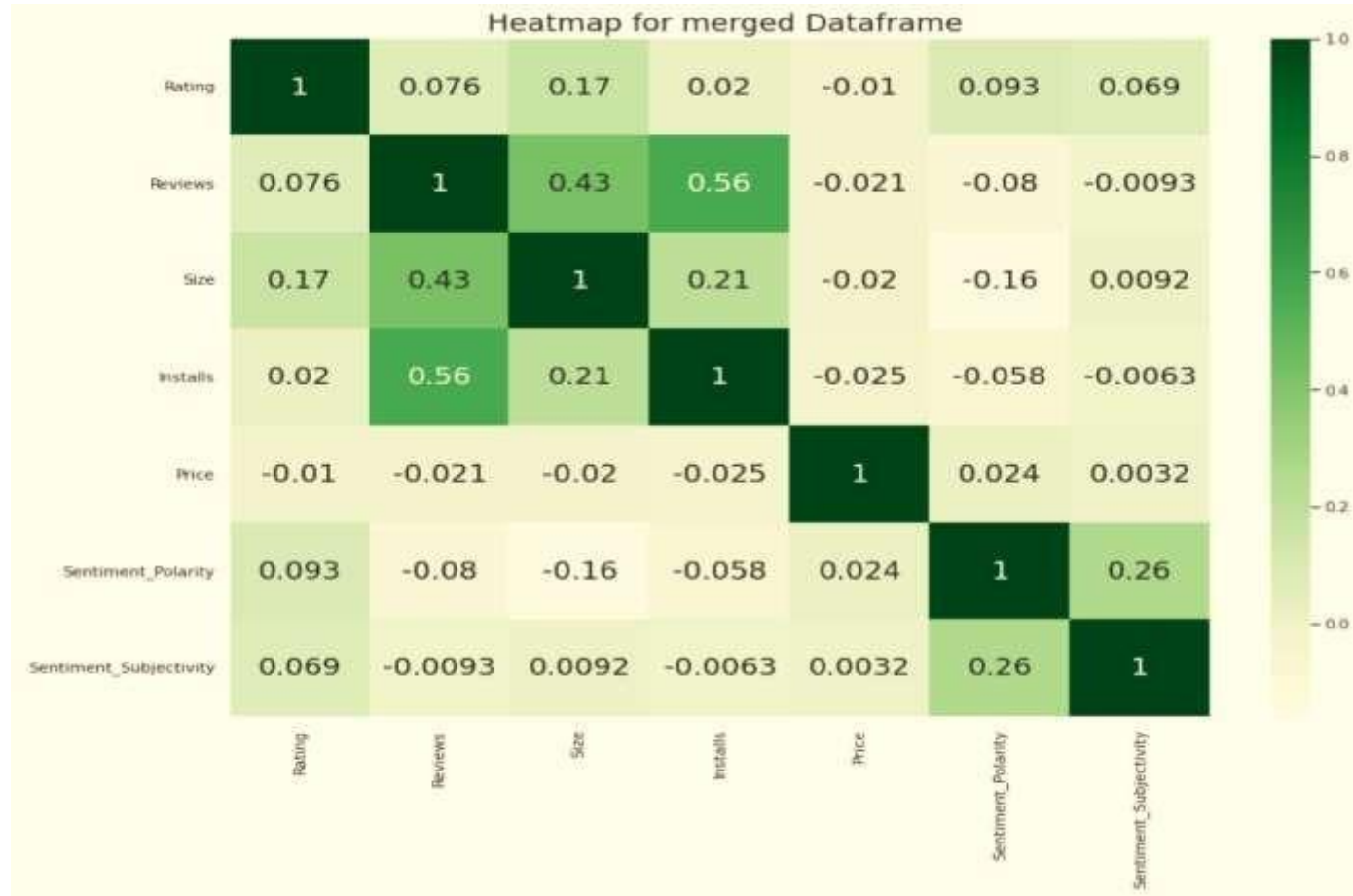
IS SENTIMENTAL_SUBJECTIVITY PROPORTIONAL TO SENTIMENTAL_POLARITY

From the above scatter plot it can be concluded that sentiment subjectivity is not always proportional to sentiment polarity but in maximum number of case, shows a proportional behavior, when variance is too high or low

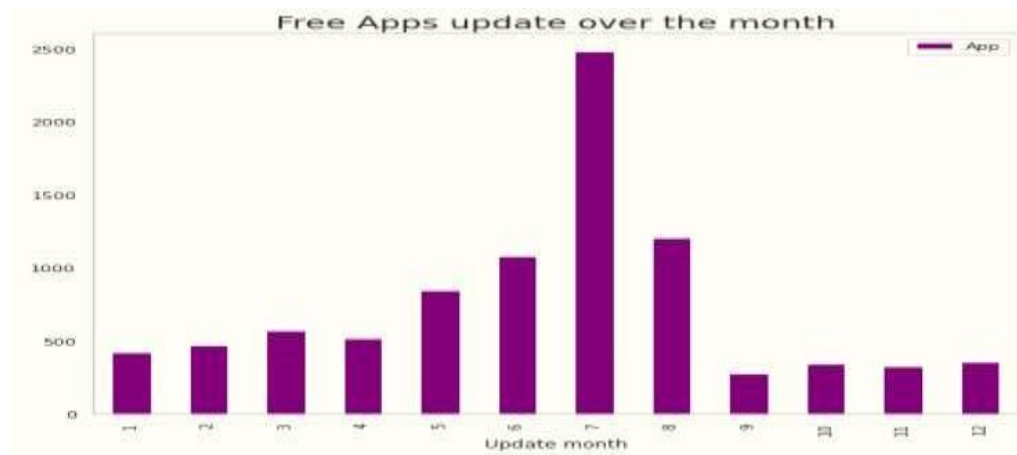
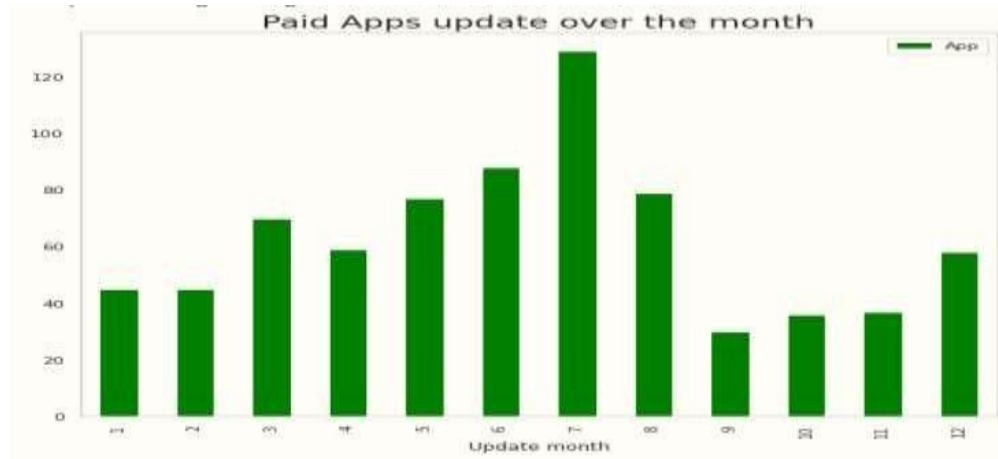
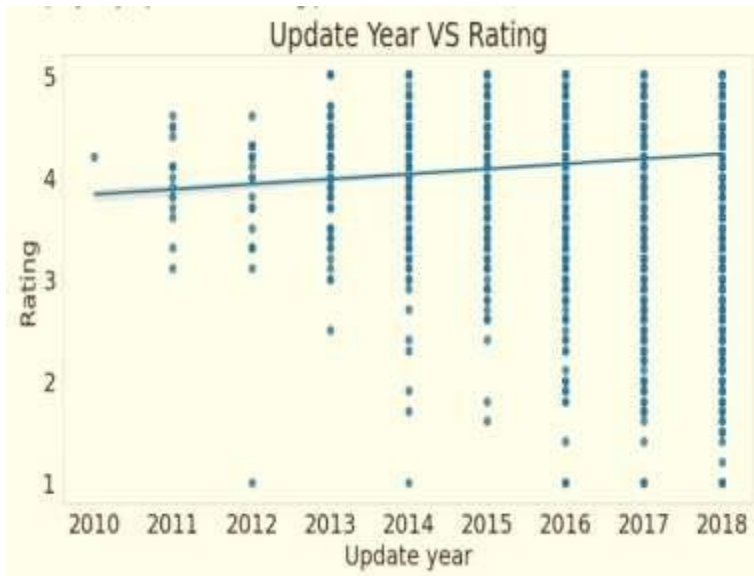


CORRELATION IN MERGED DATAFRAME

In this correlation matrix, There is not a significant relationship between Rating, Reviews, Size and Installs with respect to the Sentiment polarity and Sentiment subjectivity.



DISTRIBUTION OF APPS UPDATED OVER THE YEAR AND MONTH



CHALLENGES FACED

- ❑ Reading the dataset and comprehending the problem statement.
- ❑ Examining the business KPIs for app development and devising a solution to the problem.
- ❑ Handling the error, duplicate and NaN values in the dataset.
- ❑ Designing multiple visualizations to summarize the information in the dataset and successfully communicate the results and trends to the reader.

CONCLUSION'S

92.19% apps are Free and 7.81% apps are paid in type.

81.80% apps have Everyone content rating.

Events category has a highest mean rating of 4.39 and Dating category has lowest 4.05 rating.

Family, Game and Tools are top three categories having 1906, 926 and 829 app count.

Most competitive category: Family

Category with the highest number of installs: Game

Tools, Entertainment, Education, Business and Medical are top Genres.

8783 Apps are having size less than or equal to 50 MB.

7749 Apps has rating more than 4.0 including both type of app.

Overall sentiment count of merged dataset in which Positive sentiment count is 64%, Negative 22% and Neutral 14%.

CONCUSSION'S

It's good to develop a Free type app and having a content rating for Everyone.

Percentage of apps that are top rated = 81.80%

There are 20 free apps that have been installed over a billion times

Minecraft is the only app in the paid category with over 10M installs, and also has produced the most revenue only from installation fee.

Price, Rating, Size has no or very less correlation with Sentiment Polarity.

The median size of the apps in the play store is 12 MB

The apps whose size varies with device has the highest number average app installs.

The apps whose size is greater than 90 MB has the highest number of average user reviews, i.e., they are more popular than the rest.

Helix Jump has the highest number of positive reviews and Angry Birds Classic has the highest number of negative reviews.

THANK YOU
