

Capstone Project – 1

EDA On Play Store App Reviews

By

SHAIK AHMAD BASHA

Problem Statement

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. We have two data sets. One is play store data which contains each app (row) . And the other dataset is user reviews.

Our main objective is to explore and analyze these datasets to discover key factors responsible for app engagement and success

Work Flow:



Data Understanding

After collecting data its very important to understand the data. The App dataset has 10841 rows and 13 columns The Review dataset has 37427 rows and 6 columns:

Features of App dataset:

App : Name of the App

Category : Category of the app

Rating : Rating given for the app

Reviews : Total number of reviews for the app

Size : Size of the app (MB)

Installs : How many installs done for the app

Type : Type of the app

Price : Price of the app

Content Rating : Content rating

Genres : Genre of the app

Last Updated : Date of latest update

Current Ver : Current version of the app

Android Ver : Android version of the app

Features of Reviews dataset:

App : Name of the App

Translated Review : Text of the review

Sentiment : The review is good or not

Sentiment Polarity : Sentiment polarity

Sentiment Subjectivity : Sentiment subjectivity

Data Cleaning

```
# Checking for Null values for app data
app.isnull().sum().sort_values(ascending = False)

[13] # Checking null values for Review Data
review.isnull().sum().sort_values(ascending = False)
```

Rating	1474	Translated_Review	26868
Current Ver	8	Sentiment	26863
Android Ver	3	Sentiment_Polarity	26863
Type	1	Sentiment_Subjectivity	26863
Content Rating	1	App	0
		dtype: int64	

- ❖ The App dataset has 5 features with null values.
- ❖ The Review dataset has 4 features with null values.
- ❖ So we should remove those null values.

Data Manipulation

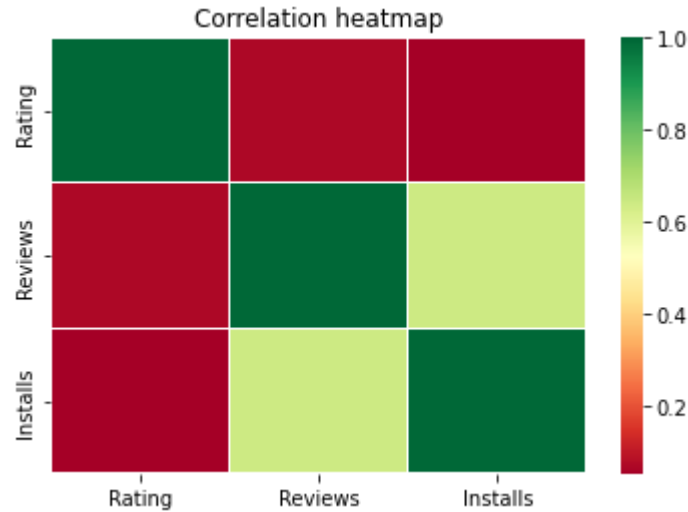
- ❖ Here, we can observe that Reviews, Installs, Price features has Object data type.
- ❖ Reviews, Installs and Price feature values are in string format. So we need to convert them into numerical values for our analysis.
- ❖ Also created a new column named 'Revenue' which is multiplication of price and number of installs for only paid apps.
- ❖ Also created a new column named 'score' which is based and positive and negative reviews for Reviews dataset

```
# Inspecting Reviews column  
new_data_app['Reviews'].dtype  
  
dtype('O')
```

```
# Inspecting Installs Column  
new_data_app['Installs'].dtype  
  
dtype('O')
```

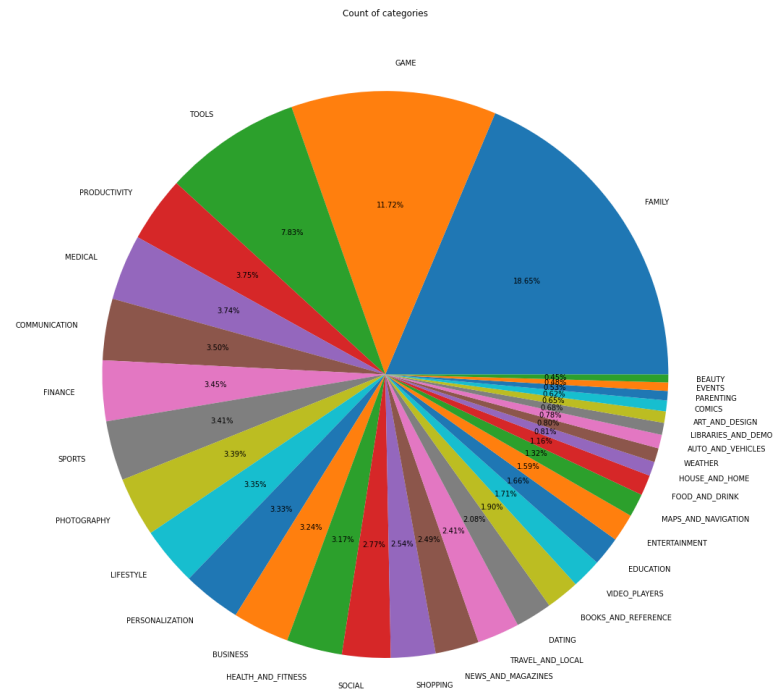
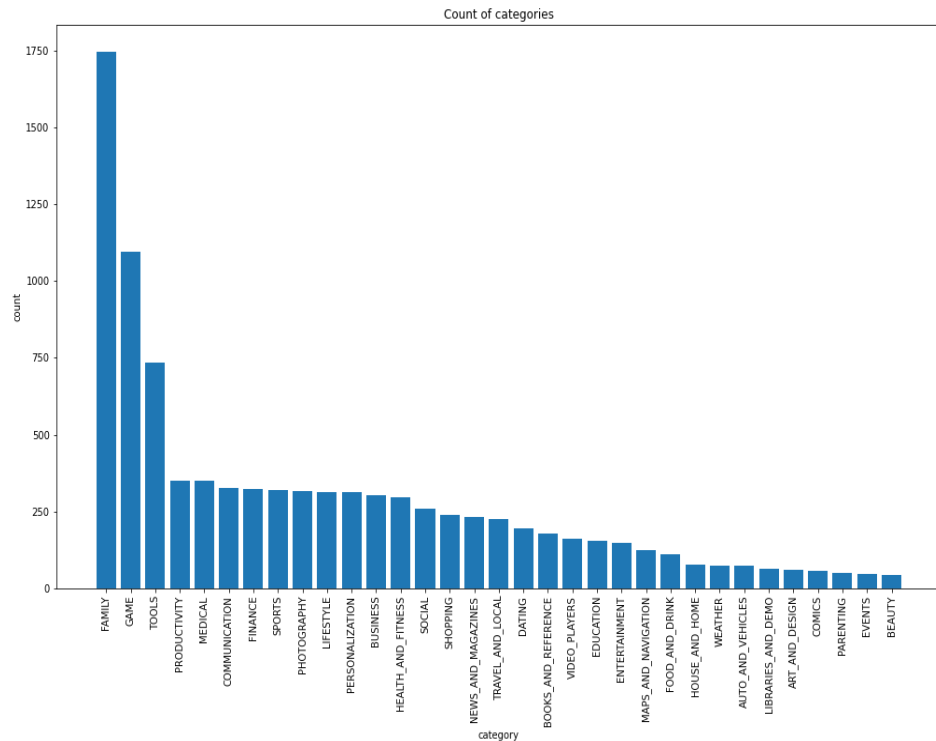
```
# Inspecting price column  
new_data_app['Price'].dtype  
  
dtype('O')
```

Exploratory Data Analysis :



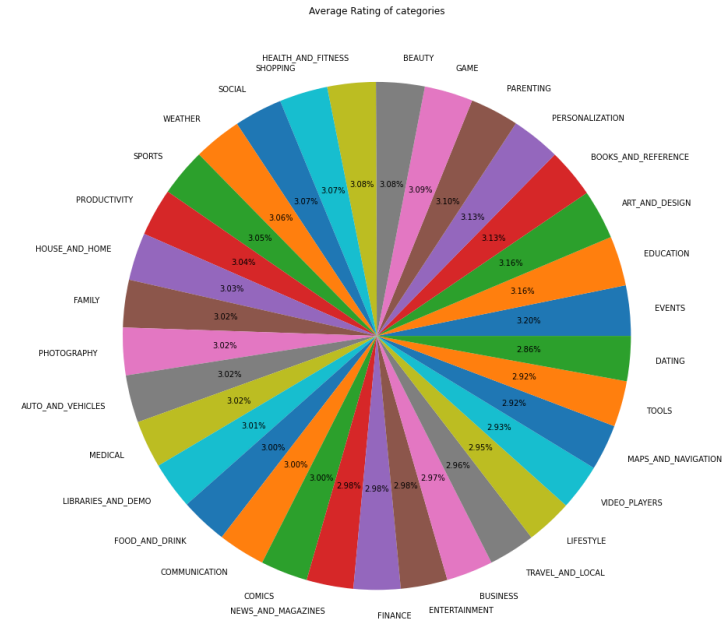
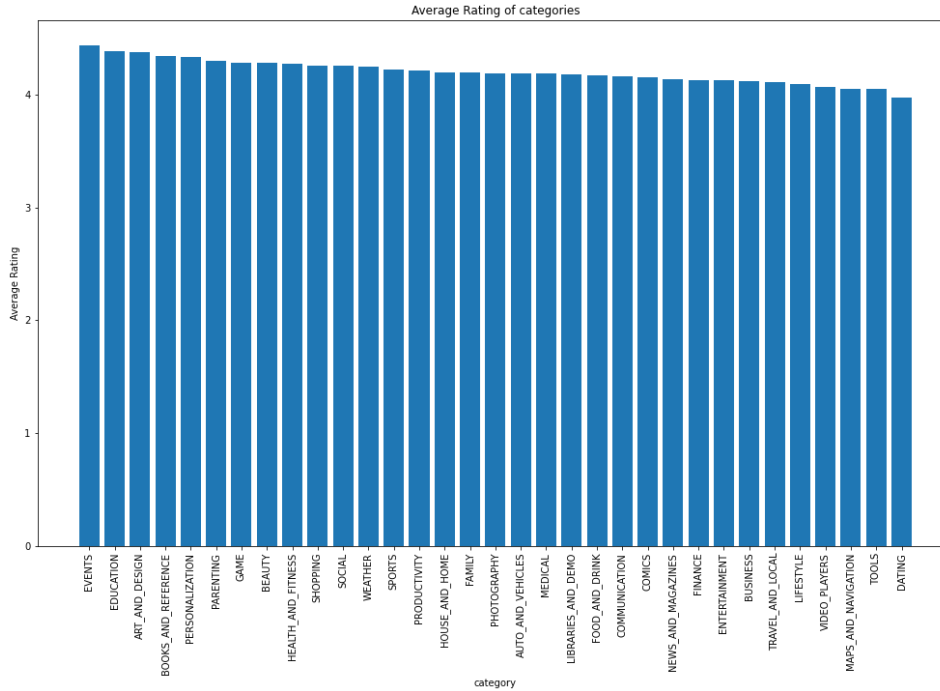
From the above heatmap, we can say that there are no variables with multicollinearity in the data

Category Analysis:



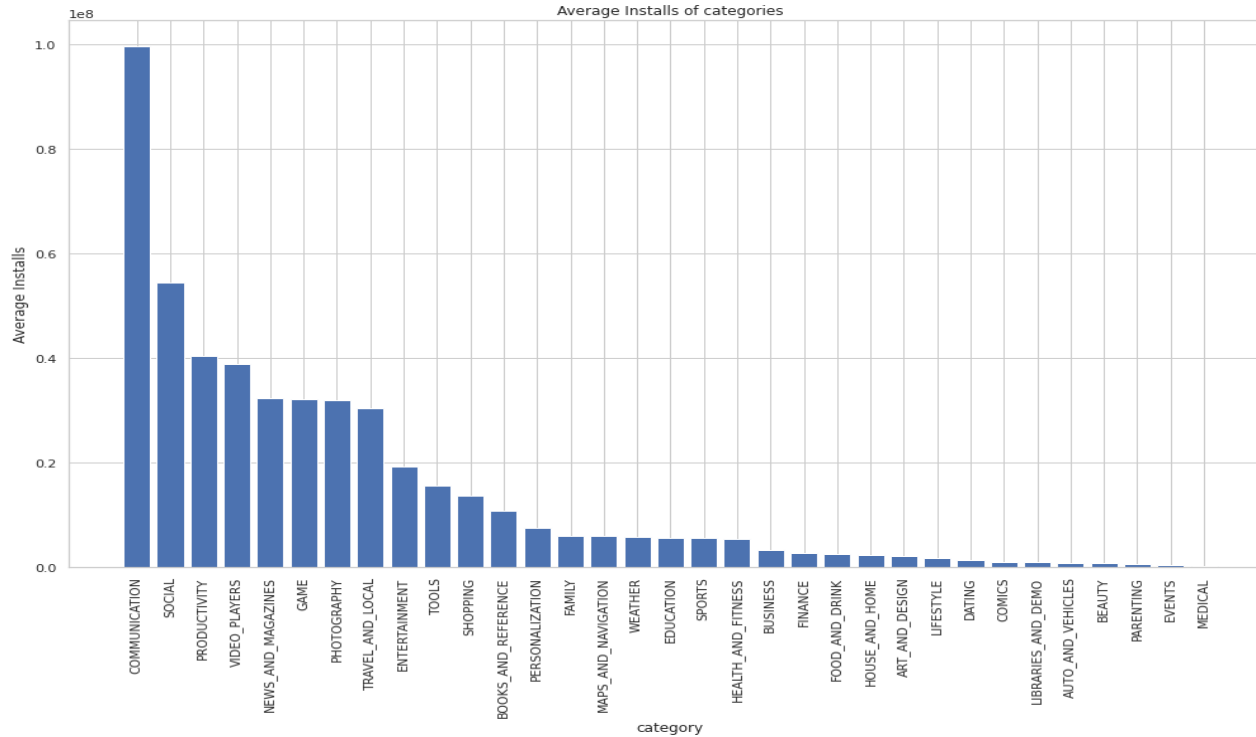
From the above charts, we can say that most apps are from **FAMILY** category followed by **GAME** category and the least is **BEAUTY** category

Category Analysis Based on Ratings :



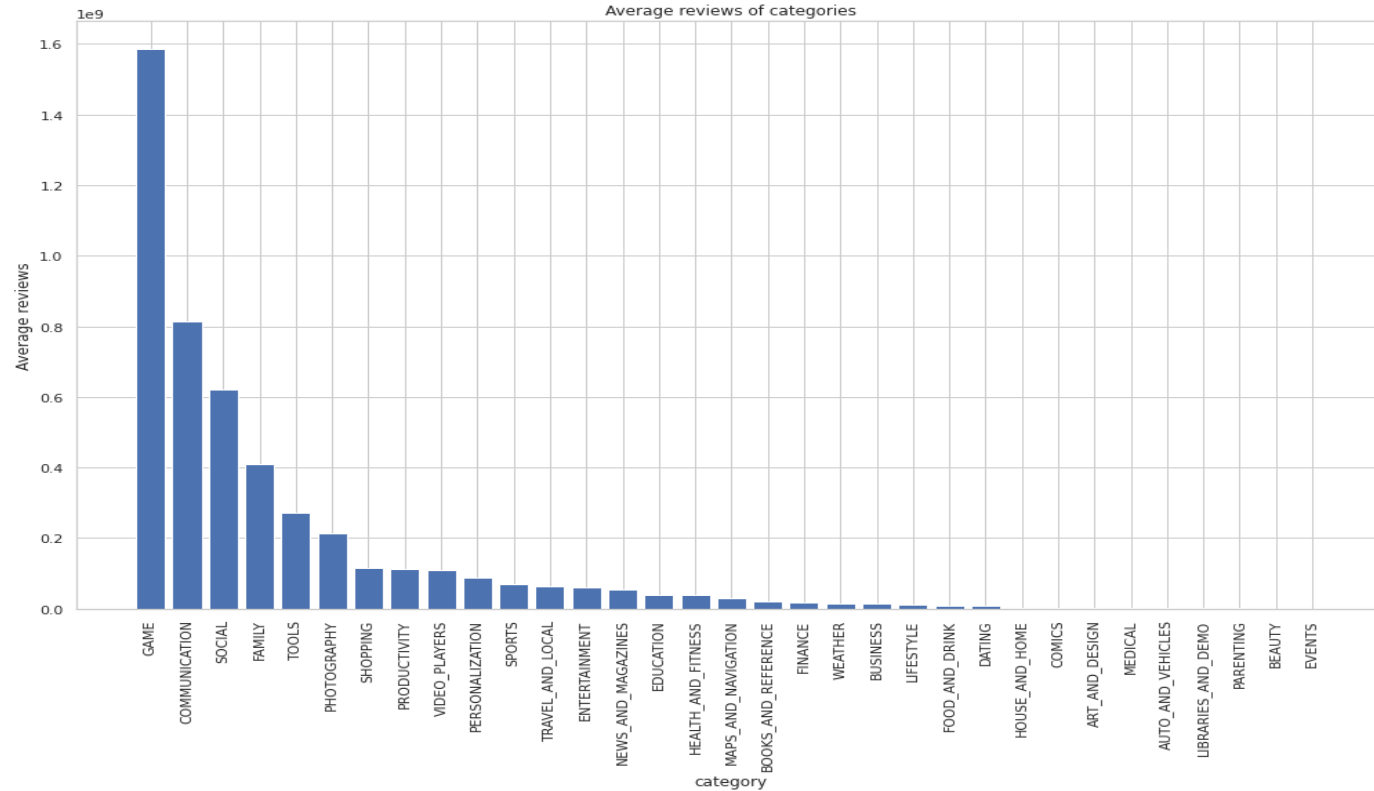
From the above charts, we can say that **EVENTS, EDUCATION, ART_AND_DESIGN, BOOKS_AND_REFERENCE, PERSONALIZATION** are the categories with highest average ratings

Category Analysis Based on Installs :



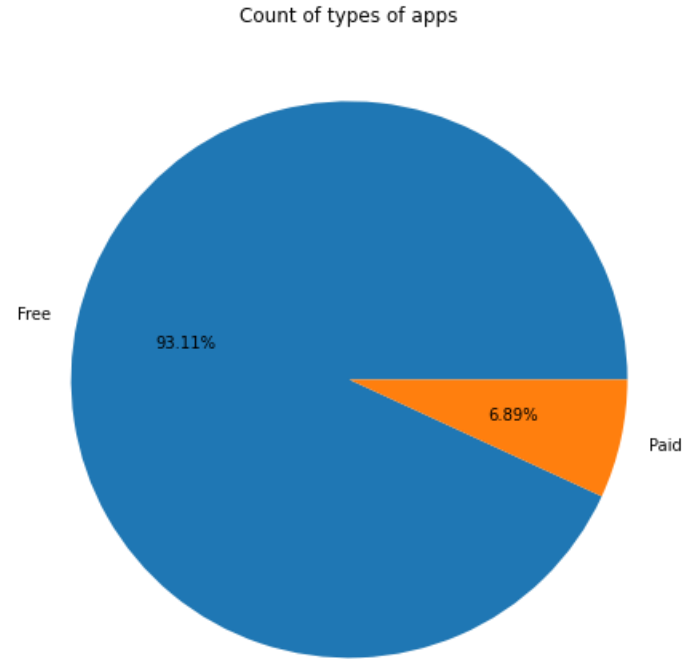
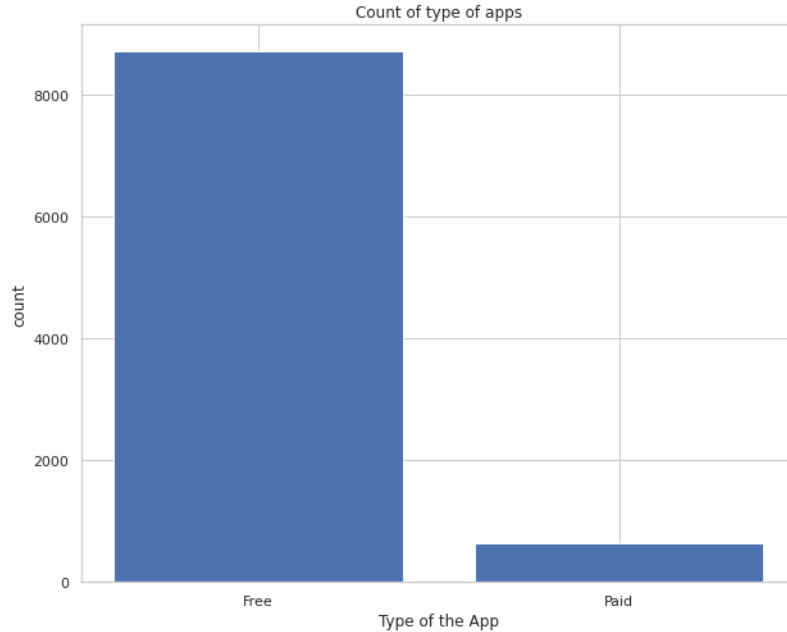
From the above plot we can observe that **COMMUNICATION** category has highest installs followed by **SOCIAL** and the least is **MEDICAL**

Category Analysis Based on Reviews :



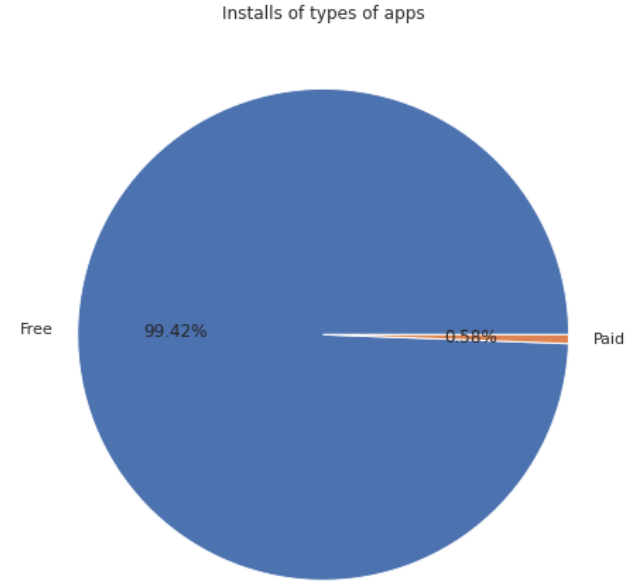
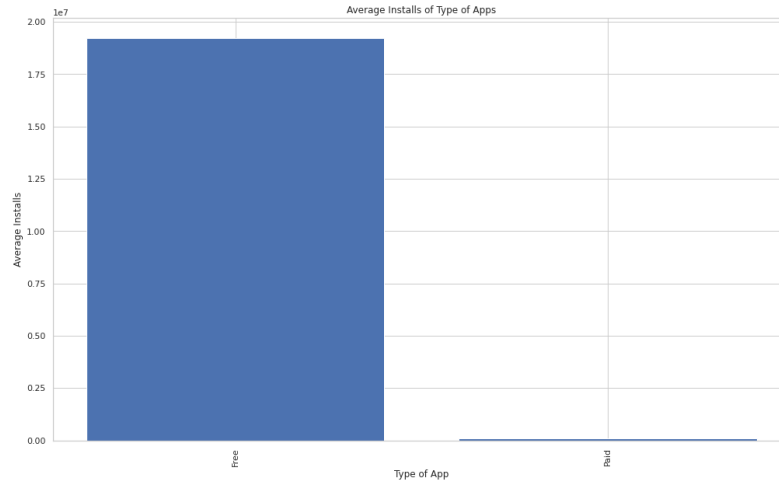
From the above charts we can observe that **GAME** followed by **COMMUNICATION** category has high reviews and least is **EVENTS**.

Type of Apps Analysis :



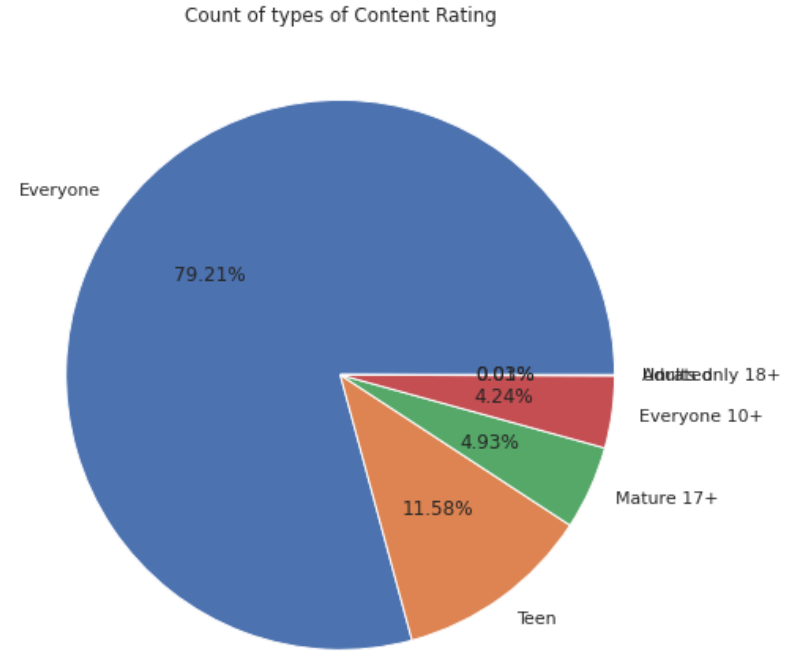
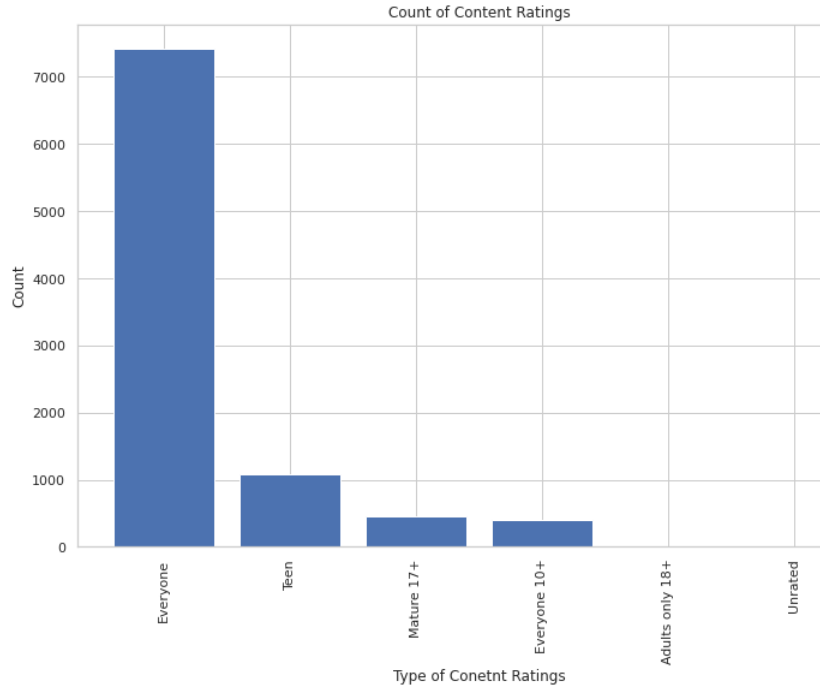
From the above chart we can observe that most of the apps are free (93.11%) and only 6.89% of apps are paid apps.

Type of Apps Analysis Based on Installs:

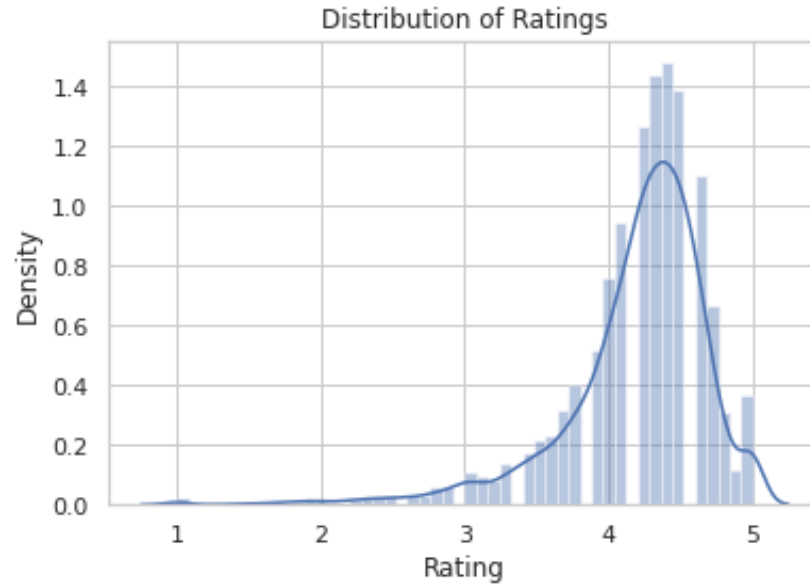


From the above visualizations, we can say that Free Apps has high installs (99.42%) than paid apps (0.58%)

Type of Content Rating Analysis :



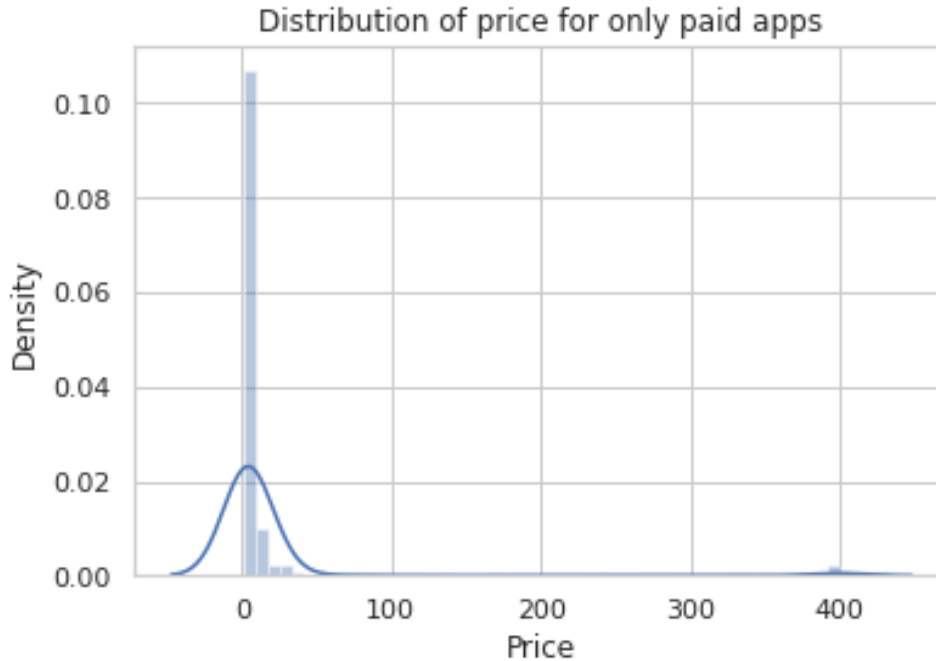
From the above chart we can say that **Everyone** content rating apps are more than others.



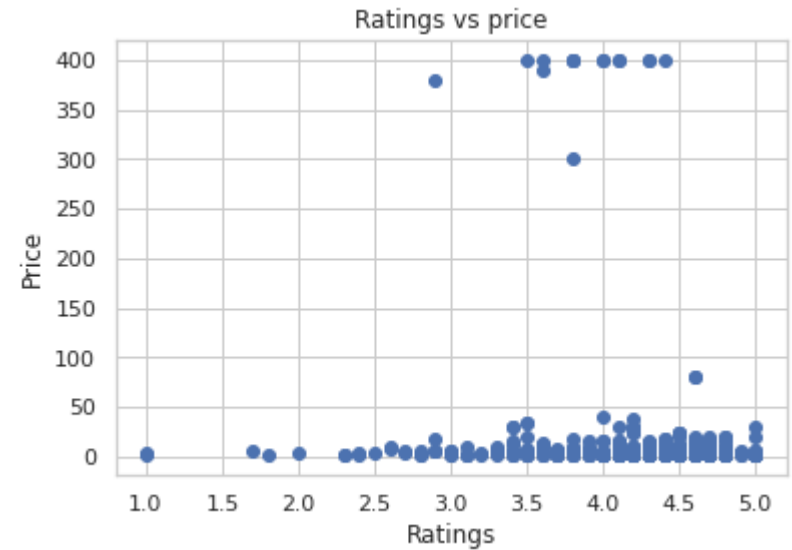
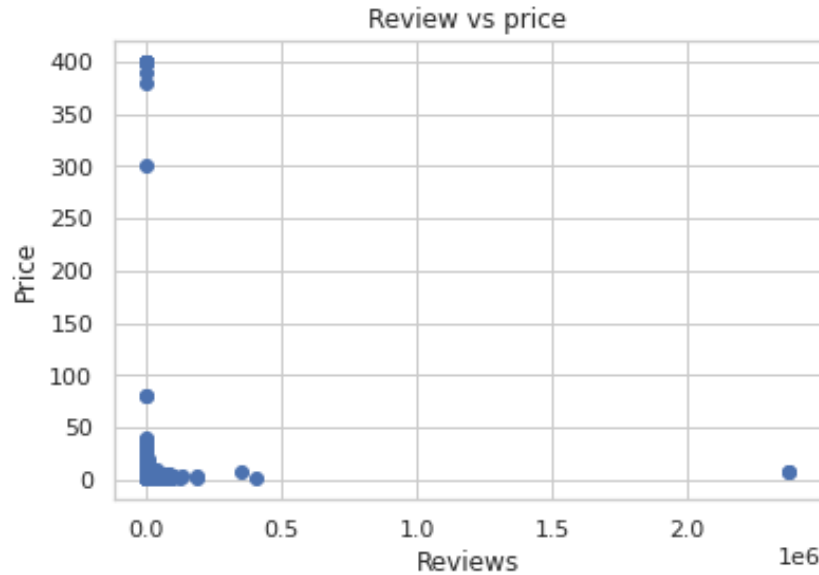
From the above chart we can say that :

- ❖ The above distribution plot is negatively skewed.
- ❖ The most number of ratings are in range of 4 to 5.

Price Analysis :

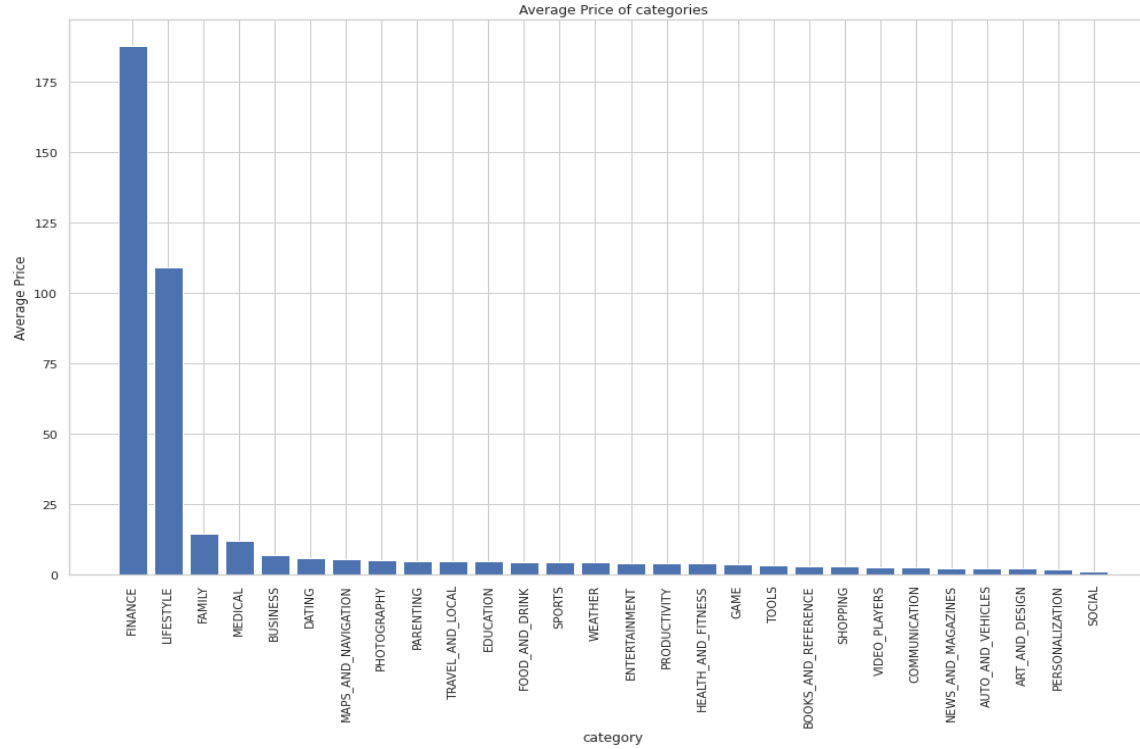


From the above distribution plot, we can say that most prices are in the range of 0 to 50 and a very few are in the range of 50 to 450



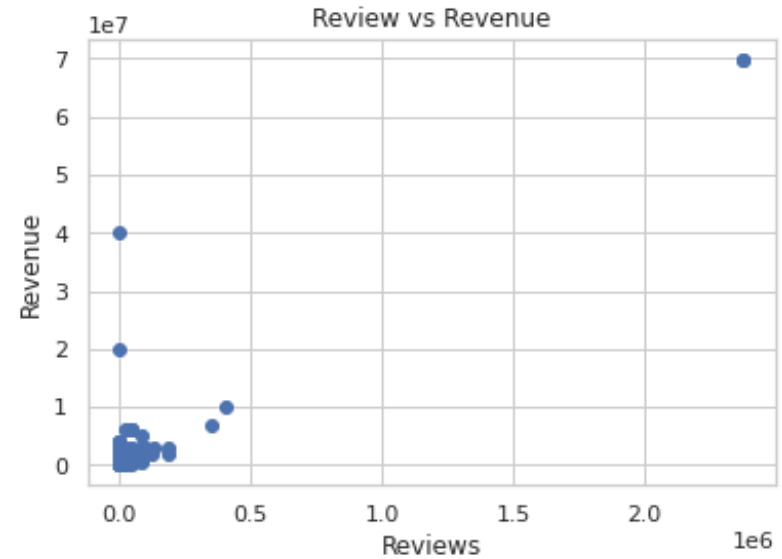
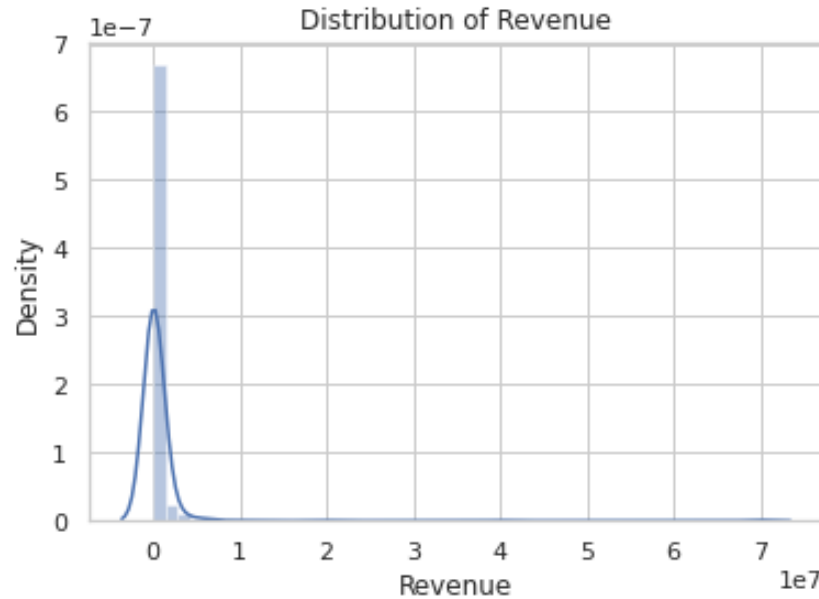
From the above scatter plots, we can observe that

- ❖ Most of the apps has prices of less than 50 USD
- ❖ And very few apps has high prices.

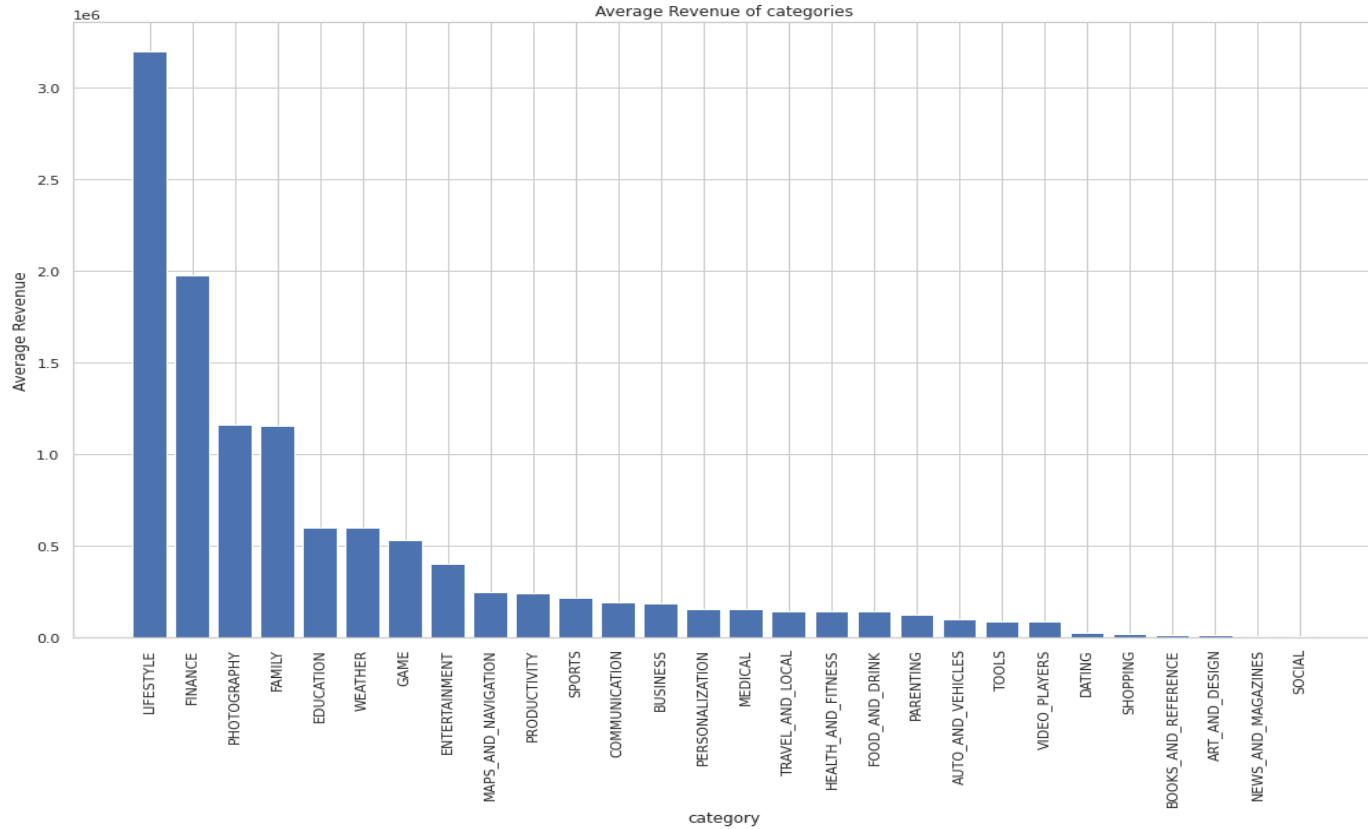


From the above visualization, **FINANCE, LIFESTYLE** category has high average prices and the least is **SOCIAL**.

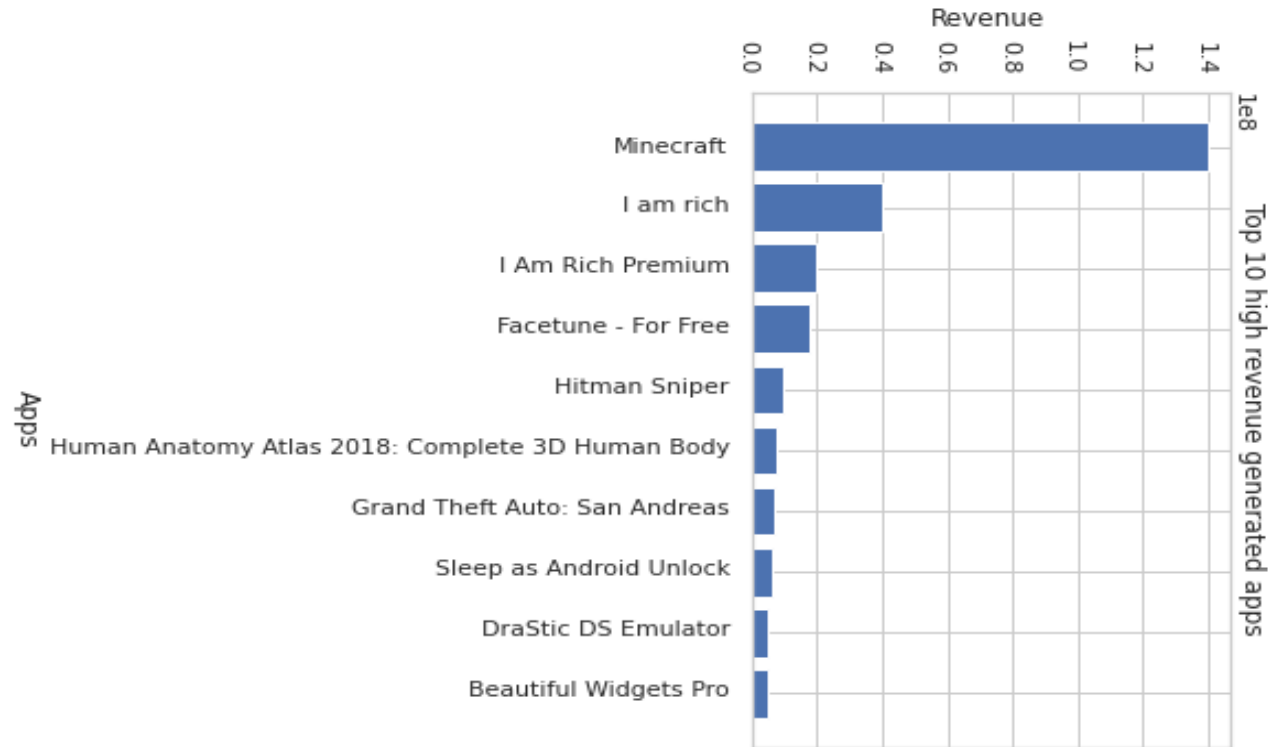
Revenue Analysis Of Paid Apps:



From the above plots we can observe that very few paid apps have highest revenue. Most of the paid apps have some decent amount of revenue.

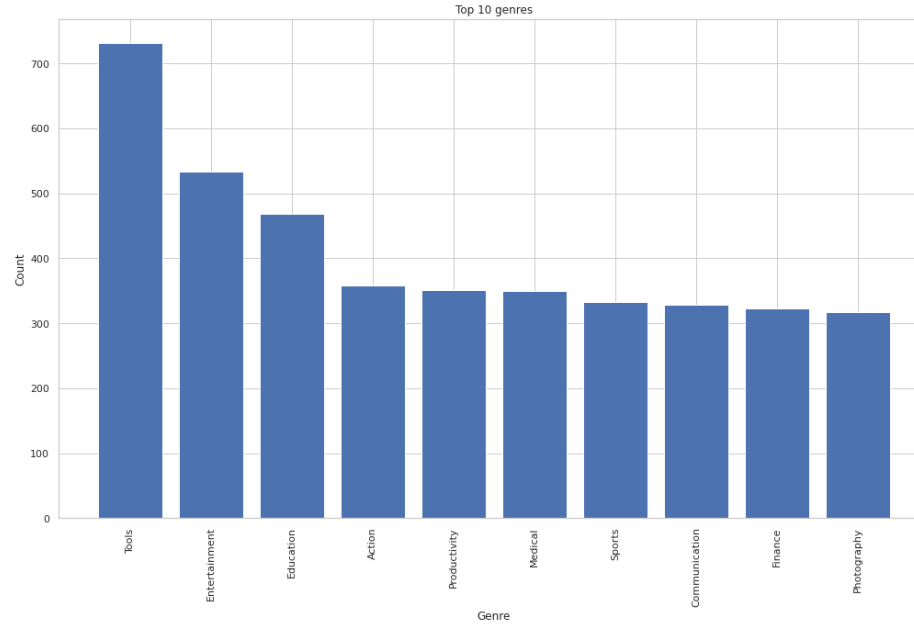


From the above visualizations, LIFESTYLE category has highest average revenue.



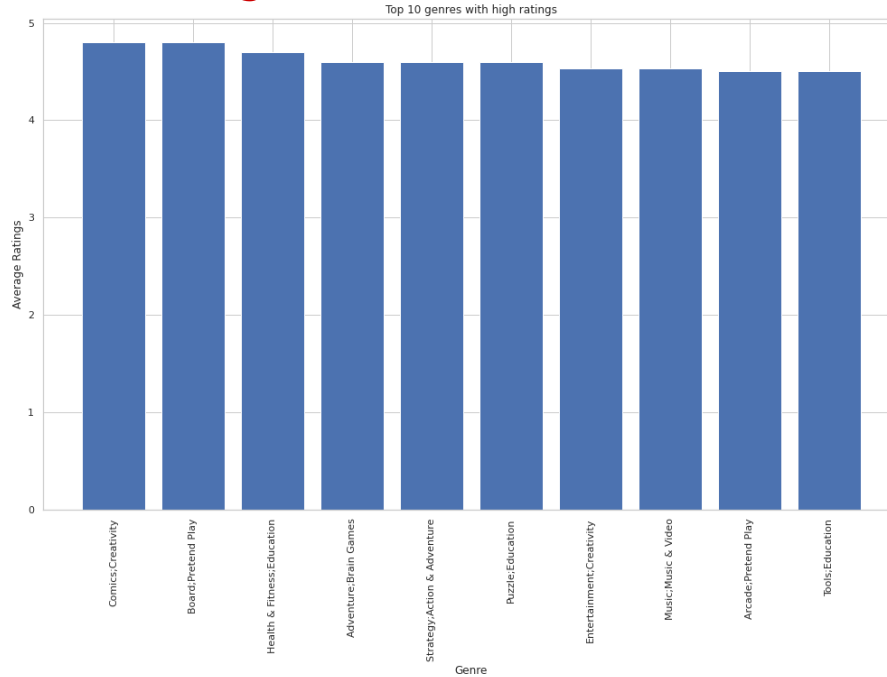
- ❖ From the above visualization, The top 5 high revenue apps are Minecraft, I am rich, I Am Rich Premium, Facetune - For Free, Hitman Sniper.

Genre Analysis :



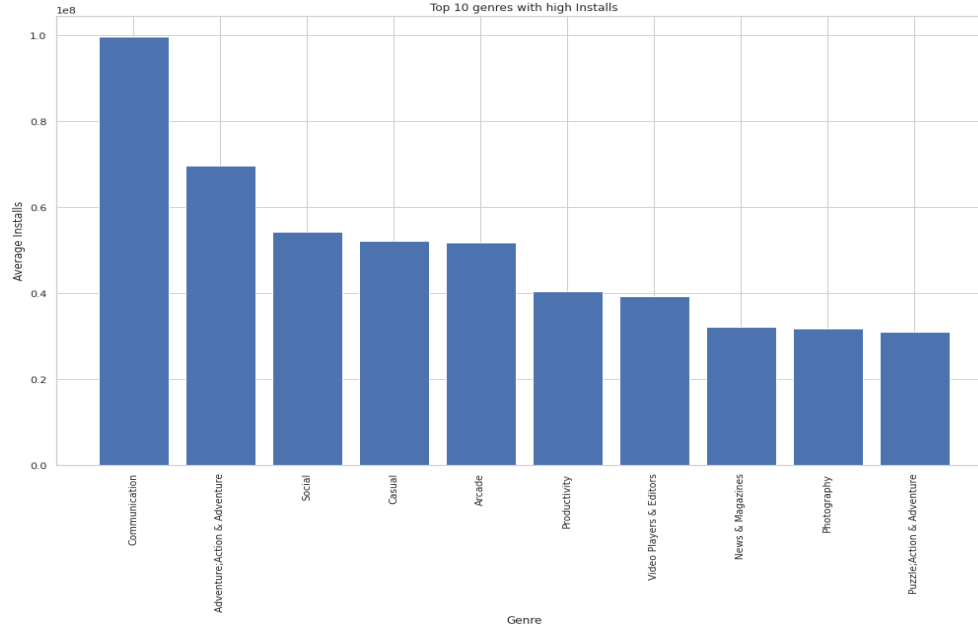
- ❖ There are total 115 genres exists in the data
- ❖ The above bar graph displays top 10 genres with most number of apps. And they are Tools, Entertainment, Education, Action, Productivity, Medical, Sports, Communication, Finance, Photography.

Genre Analysis Based On Ratings :



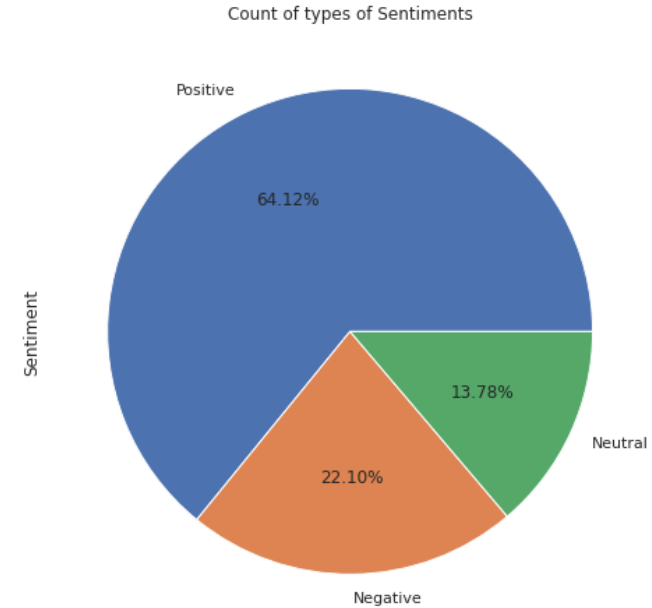
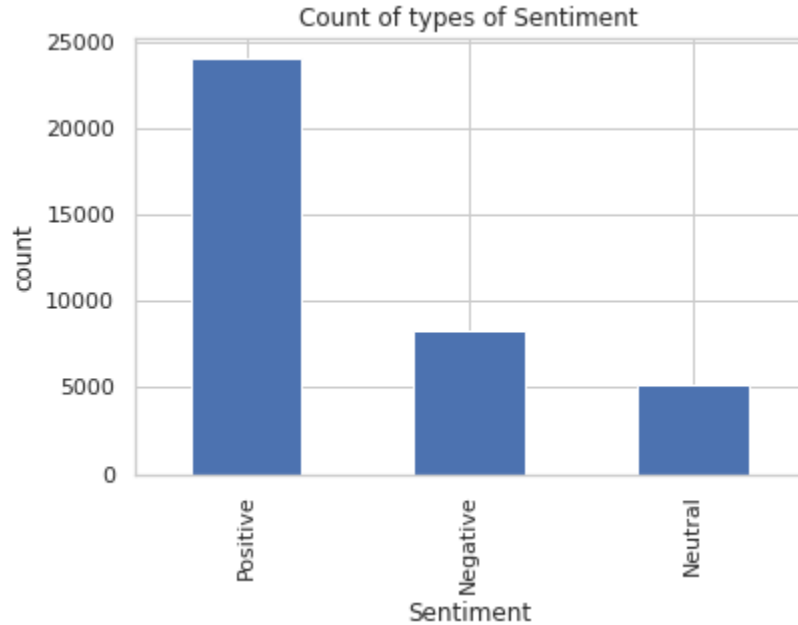
- ❖ From the above visualization, we can see the genres with highest ratings.
- ❖ They are Comics;Creativity, Board;Pretend Play, Health & Fitness;Education, Adventure;Brain Games, Strategy;Action & Adventure, Puzzle;Education, Entertainment;Creativity, Music;Music & Video, Arcade;Pretend Play, Tools;Education

Genre Analysis Based On Installs :



- ❖ From the above visualization, we can see the genres with highest installs.
- ❖ The top 10 genres with high reviews are Communication, Adventure;Action & Adventure, Social, Casual, Arcade, Productivity, Video Players & Editors, News & Magazines, Photography, Puzzle;Action & Adventure.

Sentiment of Reviews Analysis :



- ❖ From the above visualizations, most reviews are positive i.e., 64.12%, 22.10% reviews are Negative, and 13.78% reviews are Neutral.

Top 10 Popular Apps :



- ❖ From the above visualization, we can observe the apps with more positive reviews.
- ❖ Duolingo: Learn Language Free App has high positive reviews.

Final Conclusions :

- ❖ Most apps are from FAMILY category followed by GAME category and the least apps are from BEAUTY category.
- ❖ The top 5 Categories with high ratings are EVENTS, EDUCATION, ART_AND_DESIGN, BOOKS_AND_REFERENCE, PERSONALIZATION
- ❖ COMMUNICATION category has highest installs and the least is MEDICAL
- ❖ The top 5 Categories with high installs are COMMUNICATION, SOCIAL, PRODUCTIVITY, VIDEO_PLAYERS, NEWS_AND_MAGAZINES
- ❖ The top 5 Categories with high reviews are GAME, COMMUNICATION, SOCIAL, FAMILY, TOOLS.
- ❖ Most of the apps are free (93.11%) and only 6.89% of apps are paid apps.

- ❖ Free Apps has high installs (99.42%) than paid apps (0.58%)
- ❖ Most number of ratings are in range of 4 to 5 and most prices of paid apps are in the range of 0 to 50 USD *and a very few are in the range of 50 to 450 USD.*
- ❖ A very few paid apps have highest revenue. Most of the paid apps have some decent amount of revenue.
- ❖ LIFESTYLE category has highest average revenue.
- ❖ FINANCE, LIFESTYLE category has high average prices and the least is SOCIAL.
- ❖ The top 5 revenue apps are Minecraft, I am rich, I Am Rich Premium, Facetune - For Free, Hitman Sniper
- ❖ There are total 115 genres exists in the data. The top 10 genres with most number of apps are Tools, Entertainment, Education, Action, Productivity, Medical, Sports, Communication, France, Photography.

- ❖ **Comics;Creativity genre has highest rating and the Parenting;Brain Games genre has least ratings.**
- ❖ **Communication genre has high installs and Board; Pretend Play genre has very few installs**
- ❖ **Most of the reviews are positive i.e., 64.12%, 22.10% reviews are Negative, and 13.78% reviews are Neutral.**
- ❖ **The most popular apps (more positive reviews) are Duolingo: Learn Languages Free, Calorie Counter - Macros, 10 Best Foods for You, Helix Jump, 8fit Workouts & Meal Planner, Calorie Counter - MyFitnessPal, Calorie Counter - MyNetDiary, ColorNote Notepad Notes, Google Photos, Calorie Counter & Diet Tracker.**

THANK YOU