

# **FROM DOWNLOADS TO 5 STARS. THE SCIENCE OF APP VIRALITY**

A Deep Dive into Consumer  
Behavior and Marketplace  
Dynamics

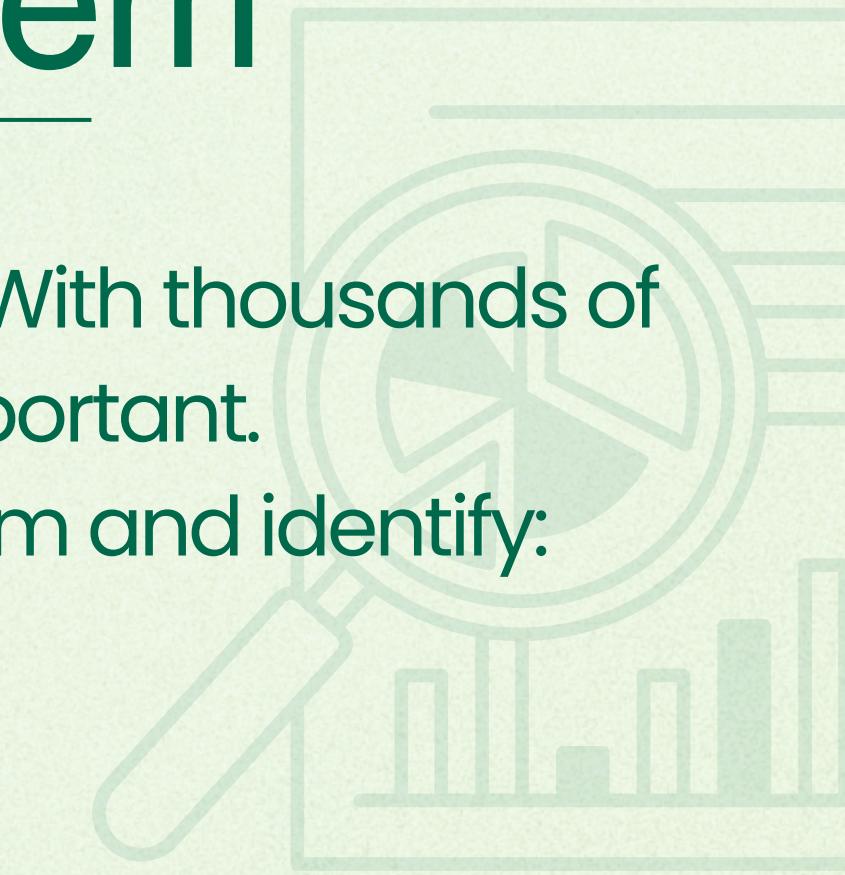
Presented by Team Carpe Diem

Vidhit T S | Padam Rathi | Md Faisal

# Decoding the Digital Ecosystem

## Context & Problem Statement

- Smartphones and apps drive modern communication and entertainment. With thousands of apps on the Play Store, understanding what makes an app successful is important.
- Our objective is to perform a data-driven analysis of the Play Store ecosystem and identify:
  - What factors influence app success,
  - How users rate and respond to apps, and
  - How metadata + user sentiment together reflect real usage trends.



## Datasets Used

We used two datasets from the Google Play Store environment:

### 1. Googleplaystore Apps datasets

Metadata of ~10,000 apps (Category, Rating, Installs, Size, Reviews, Price, Type, Genres)

### 2. Googleplaystore User Reviews

~24,000 user reviews with Sentiment, Polarity, Subjectivity, and translated review text

These datasets allow us to combine app-level features with user opinions for deeper insights.

## Exploratory Themes

1. **Market Dominance:** Which categories get the most installs? Does content rating (Everyone/Teen/Mature) affect popularity?
2. **Success Factors:** Do higher ratings lead to more installs? Do reviews influence visibility and trust?
3. **User Satisfaction:** Are Free apps rated better than Paid apps? How do sentiments reflect issues like ads, crashes, or privacy concerns?

# Methodology used in the analysis

## 1 Preprocessing

Cleaning raw datasets by fixing missing values, correcting formats, removing duplicates, and standardizing fields.

## 2 Engineering

Creating useful features (numeric conversions, one-hot encodings, sentiment fields) to enhance analytical depth.

## 3 Exploration

Analyzing patterns through visualizations (ratings, installs, categories) to understand app behaviour.

## 4 Insights and Inferences

Extracting meaningful findings about category performance, user satisfaction, monetization, and sentiment trends.

## 5 Prediction

Building a simple model (Random Forest) to estimate app success or ratings based on its attributes.

# Preprocessing

- **Classified Missing App Type:**

Filled missing Type values by checking the Price column and labeling apps as Free or Paid.

- **Removed Invalid Content Rating Rows:**

Dropped the 2 rows where Content Rating was missing since they were too few to impact analysis.

- **Dropped Irrelevant Version Columns:**

Removed Current Ver and Android Ver as they had inconsistent formats and were not useful for insights.

- **Imputed Missing Ratings:**

Used a Random Forest Regressor to predict and fill missing Rating values based on other app features.

- **Standardized App Size:**

Converted all Size values (M, k) into megabytes (MB) for uniform analysis.

## Dataset Initially(Missing values)

• Rating	1474
• Type	1
• Content Rating	1
• Current Ver	8
• Android Ver	3

# Feature Engineering

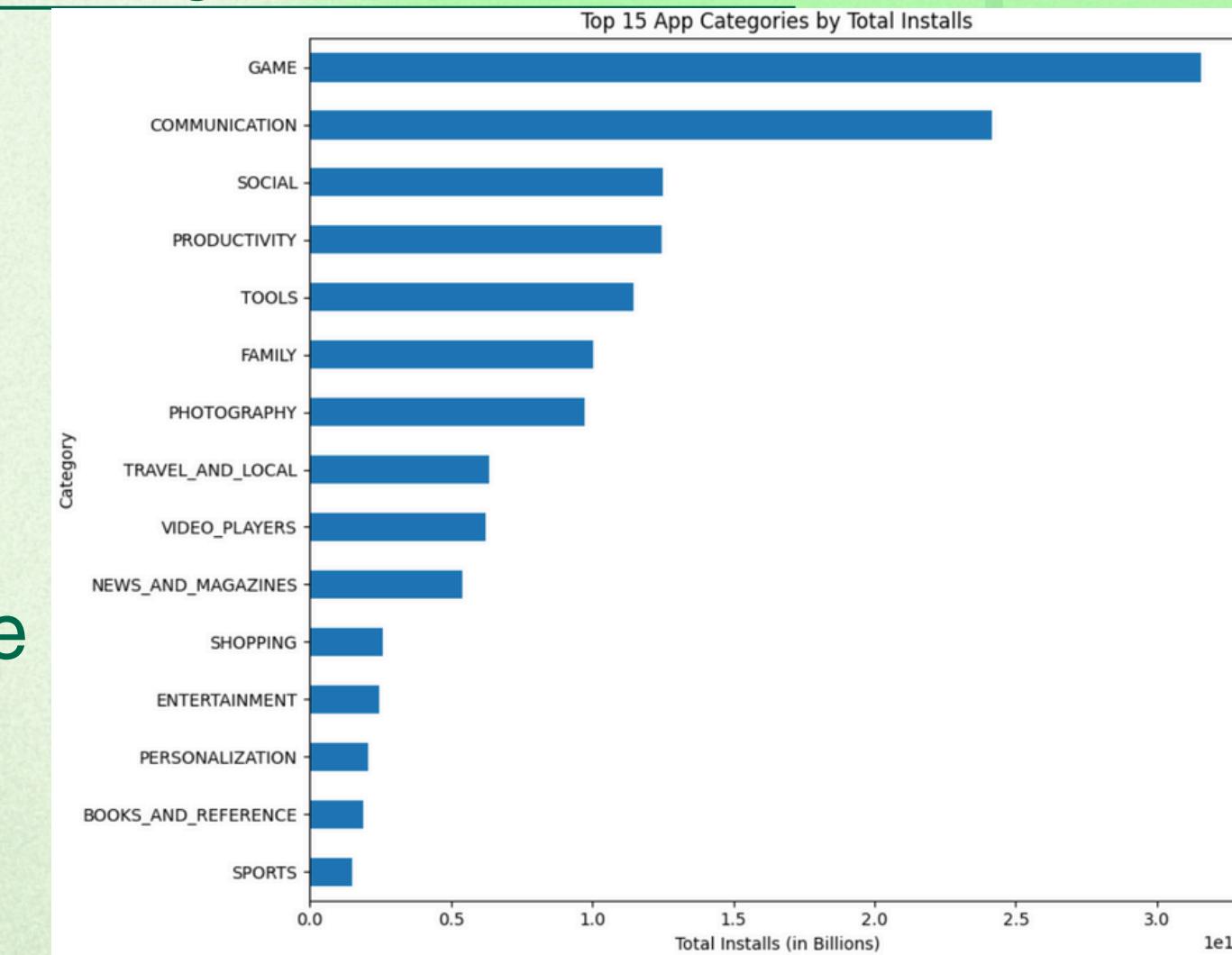
- 1. Extracted Update Information:** Converted Last Updated into datetime and created new fields: Update Year, Update Month, and Days Since Update.
- 2. Created Engagement Metric:** Added an Engagement Rate feature calculated as  $(\text{Reviews} \div \text{Installs})$  to measure user interaction.
- 3. Categorized Install Counts:** Binned Installs into Install Tiers (Low → Massive) for better comparison across apps.
- 4. Separated Genres:** Split the Genres column into Primary Genre and Secondary Genre for more detailed analysis.



# Google playstore Dataset analysis

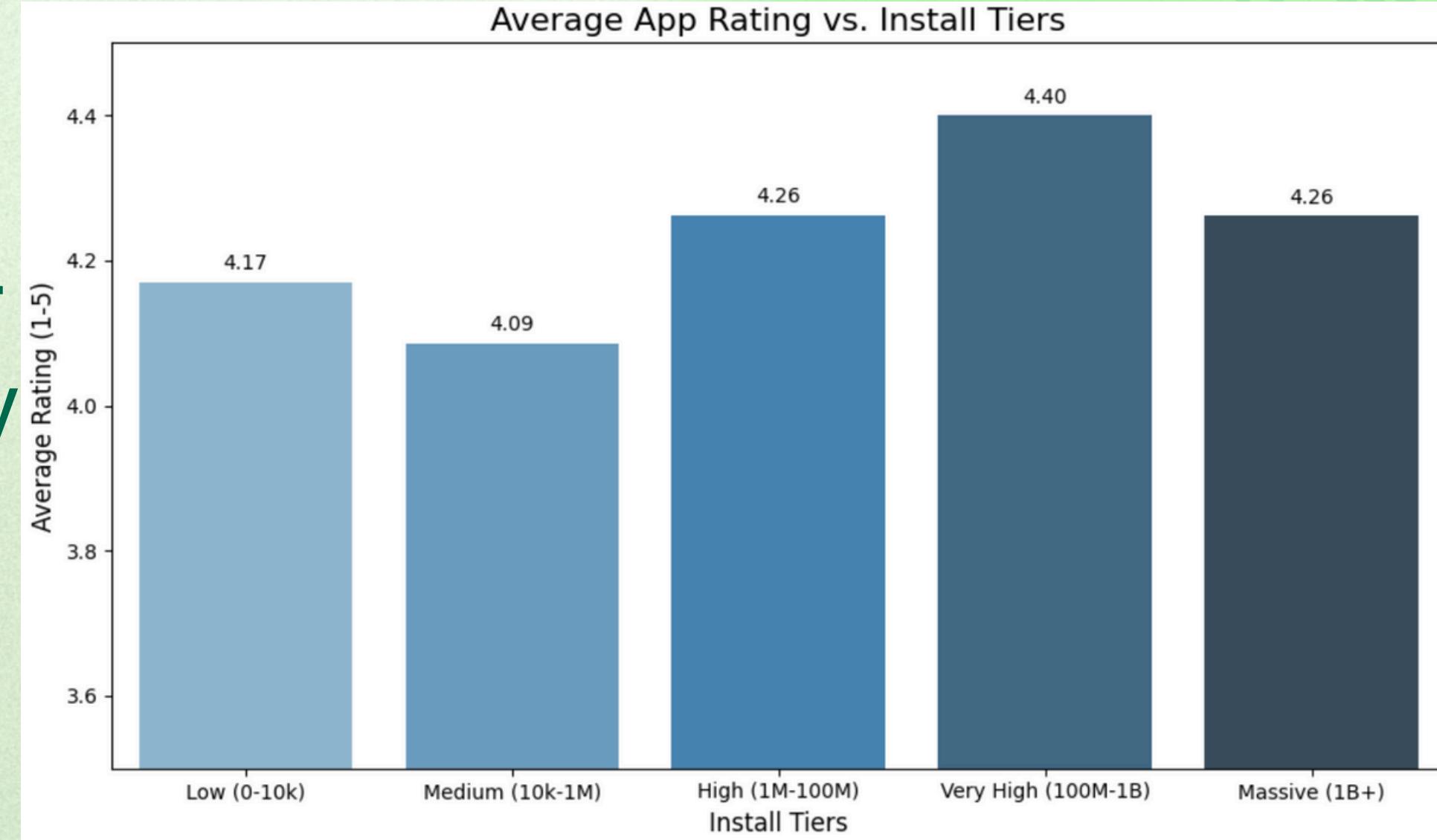
Which app categories lead the Play Store in total installs, and how does their popularity compare across the platform?

- Games dominate the Play Store ecosystem, accounting for the highest total installs by a large margin. This shows gaming apps drive the majority of user engagement on mobile devices.
- Communication and Social apps follow closely, highlighting that messaging, calling, and social networking remain core smartphone activities.
- Productivity, Tools, and Family categories also show very high install volumes, indicating that practical, everyday-use apps have broad, consistent demand across users.
- Categories like Sports, Books & Reference, and Personalization show significantly lower total installs, suggesting that these app types cater to more niche audiences and have limited mass-market demand compared to communication, gaming, or utility apps.



# Rating Trends Across Install Tiers

Is there a strong relationship between an app's rating and its install count, or are installs more heavily influenced by external factors like promotion and platform visibility?

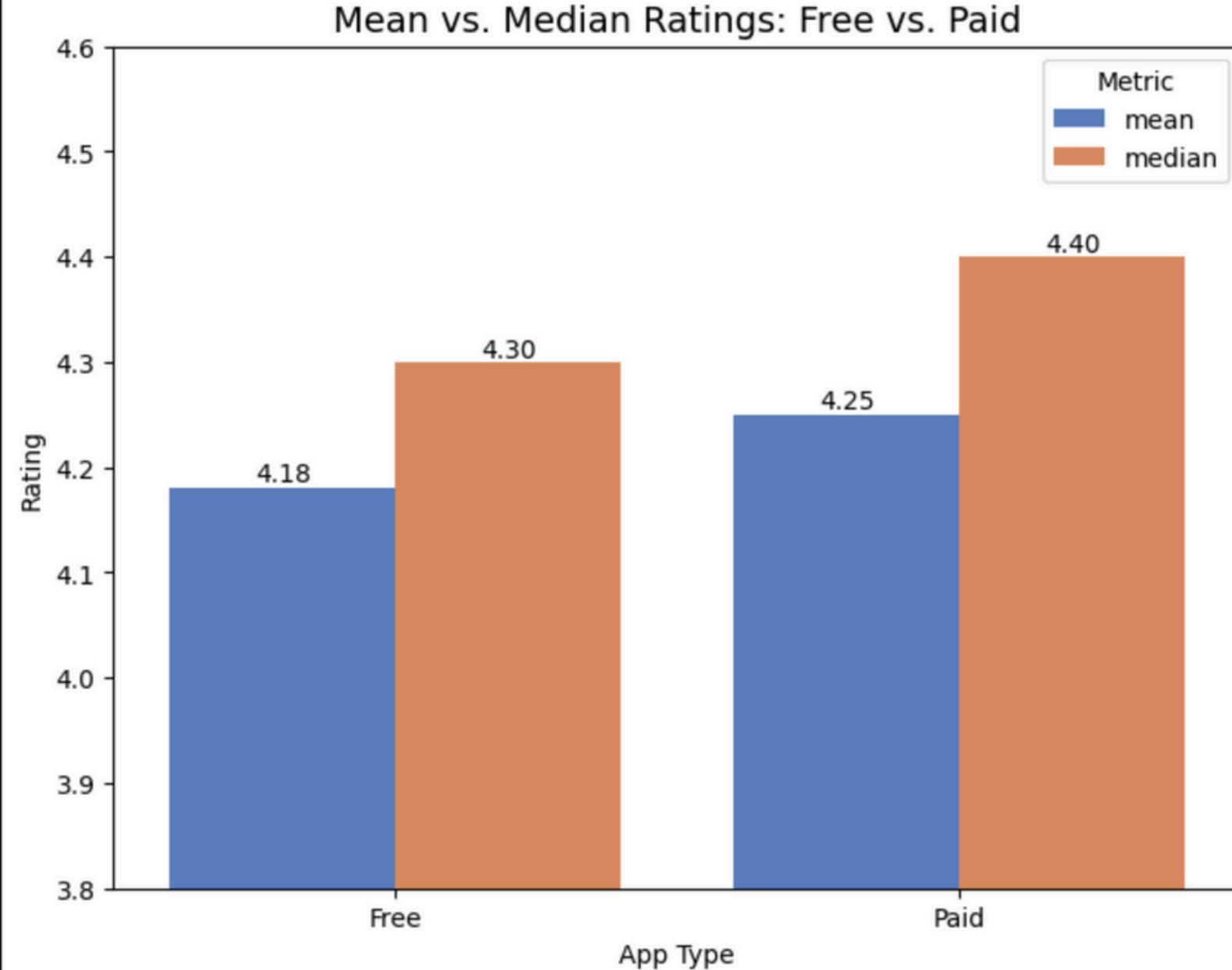


- Ratings rise steadily from the Medium to Very High tiers, showing that only consistently high-quality apps manage to scale to tens or hundreds of millions of installs.
- Low-install apps often show high initial ratings because early reviews usually come from supportive or biased early users; ratings normalize once the app reaches a larger audience (Medium tier).
- At the 1B+ tier, ratings dip slightly because widely used apps receive harsher criticism — when an app becomes universal, even small issues generate large volumes of negative feedback.

## Rating Distribution: Free vs. Paid Apps

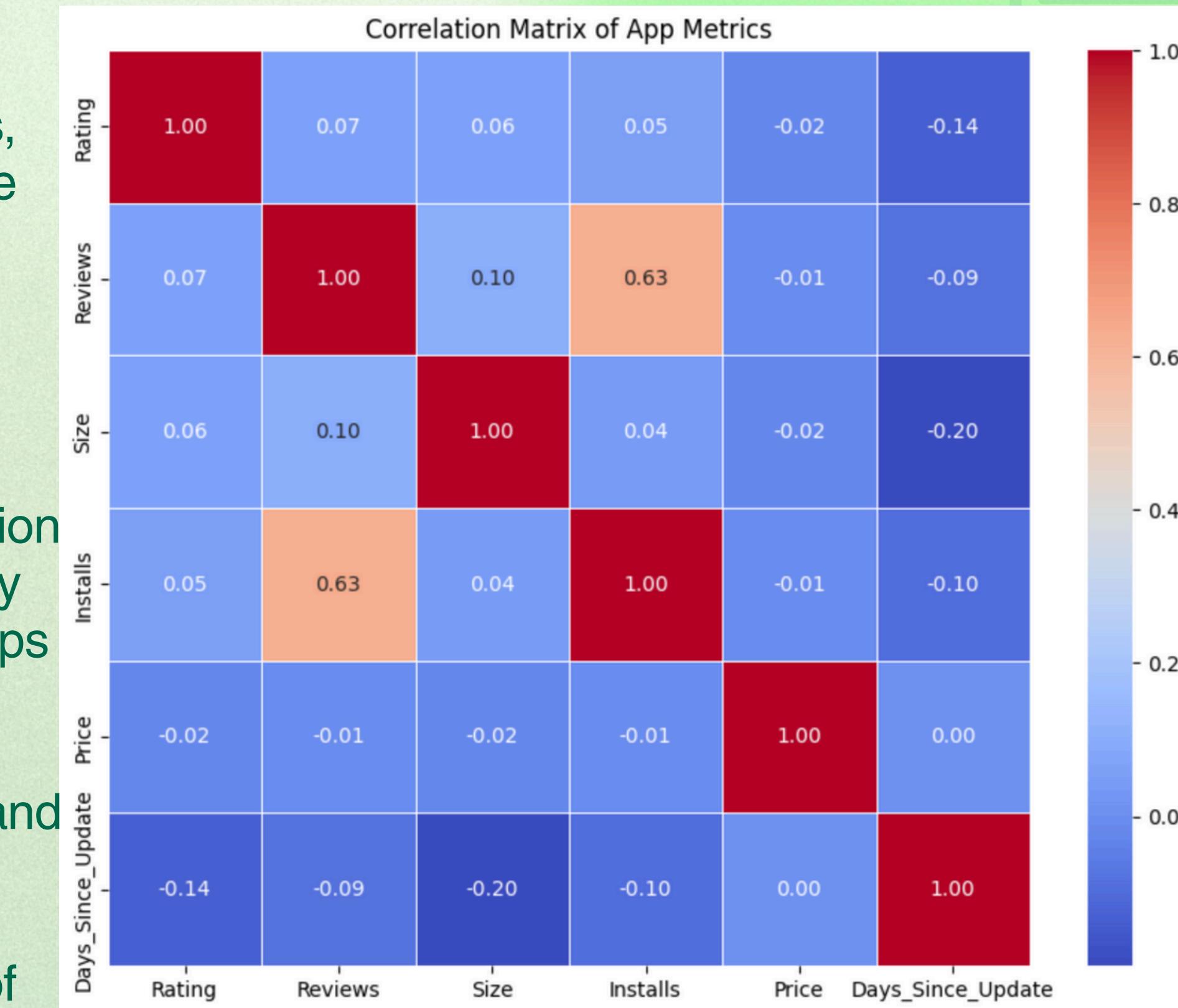
Are paid apps rated better or worse than free apps, and what does this reveal about how monetization affects user experience?

- **Premium Quality Advantage:** Paid apps receive higher ratings (Mean 4.26, Median 4.40) than free apps (Mean 4.18, Median 4.30), showing users are more satisfied with paid apps—likely because they offer ad-free experiences, better features, and more polished quality to justify their price.
- **Negative-Tail Skew in Both Categories:** For both free and paid apps, the Median is higher than the Mean, revealing a left-skewed distribution where overall satisfaction is high, but a small number of extremely low (1-star) reviews disproportionately drags down the average rating.
- **Expectations vs Reality Gap:** The rating difference between free and paid apps is small (~0.08), suggesting that while paid apps do perform better, paying users expect more and rate more critically so even small bugs or shortcomings prevent paid apps from achieving much higher rating gaps.



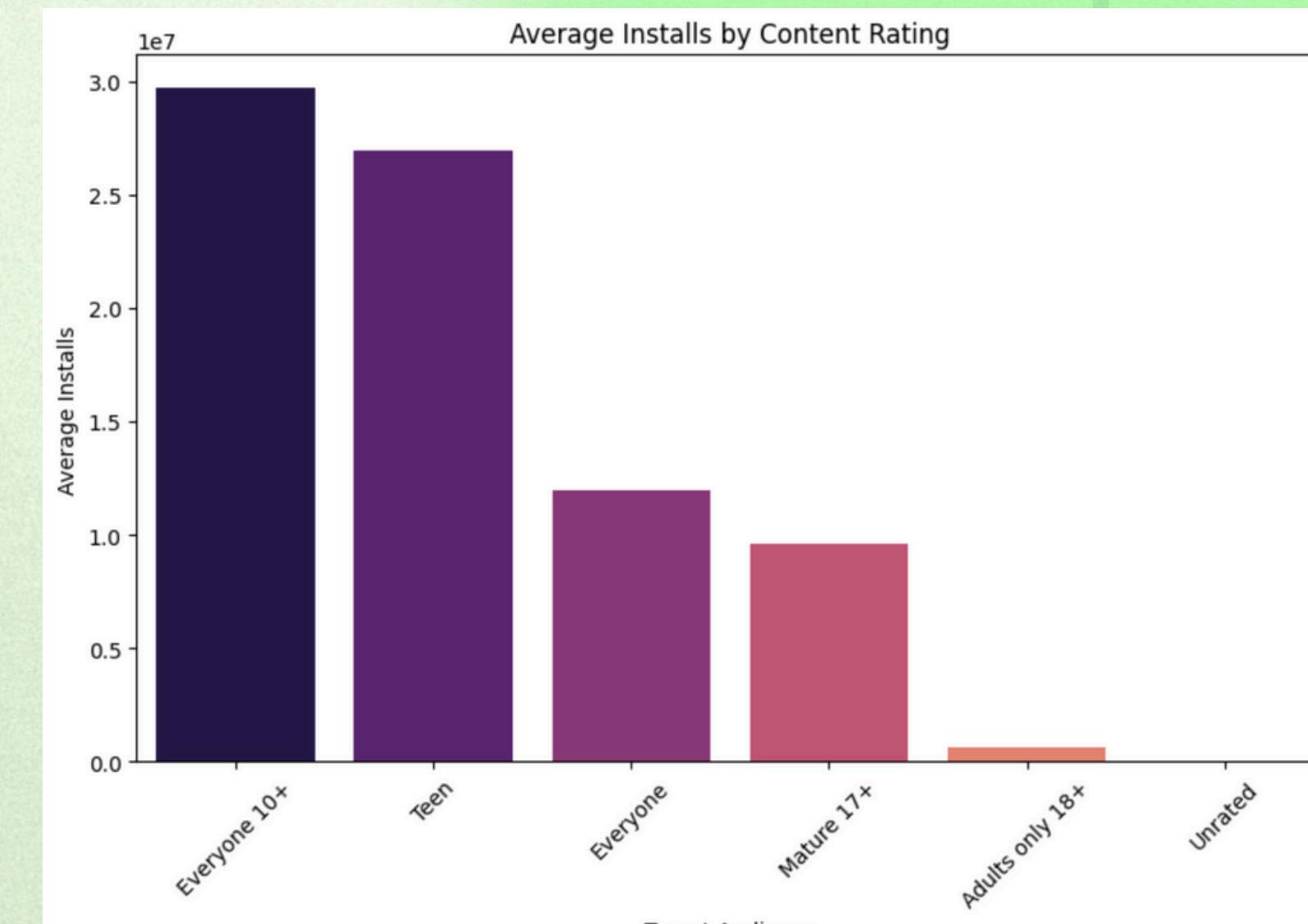
# Correlation Analysis: What Actually Drives App Downloads?

1. **Reviews Drive Growth:** A strong correlation (0.63) shows that apps with more reviews gain more installs, confirming that user engagement and social proof are the primary engines of app growth.
2. **Ratings Don't Influence Installs:** With almost no correlation (0.05), ratings barely affect download counts—apps can be hugely popular with mediocre ratings, proving that popularity ≠ quality.
3. **Frequent Updates Boost Installs:** A negative correlation (-0.10) with Days Since Update indicates that recently updated apps attract more users, while neglected apps rarely scale.
4. **App Size Doesn't Affect Popularity:** A near-zero correlation (0.05) shows users download both large and small apps, meaning app size is not a meaningful barrier to adoption.
5. **Pricing Has No Predictive Power:** With a correlation of -0.01, price shows no real link to installs; utility and visibility outweigh cost when users choose apps.



# Target Audience vs App Popularity Insights

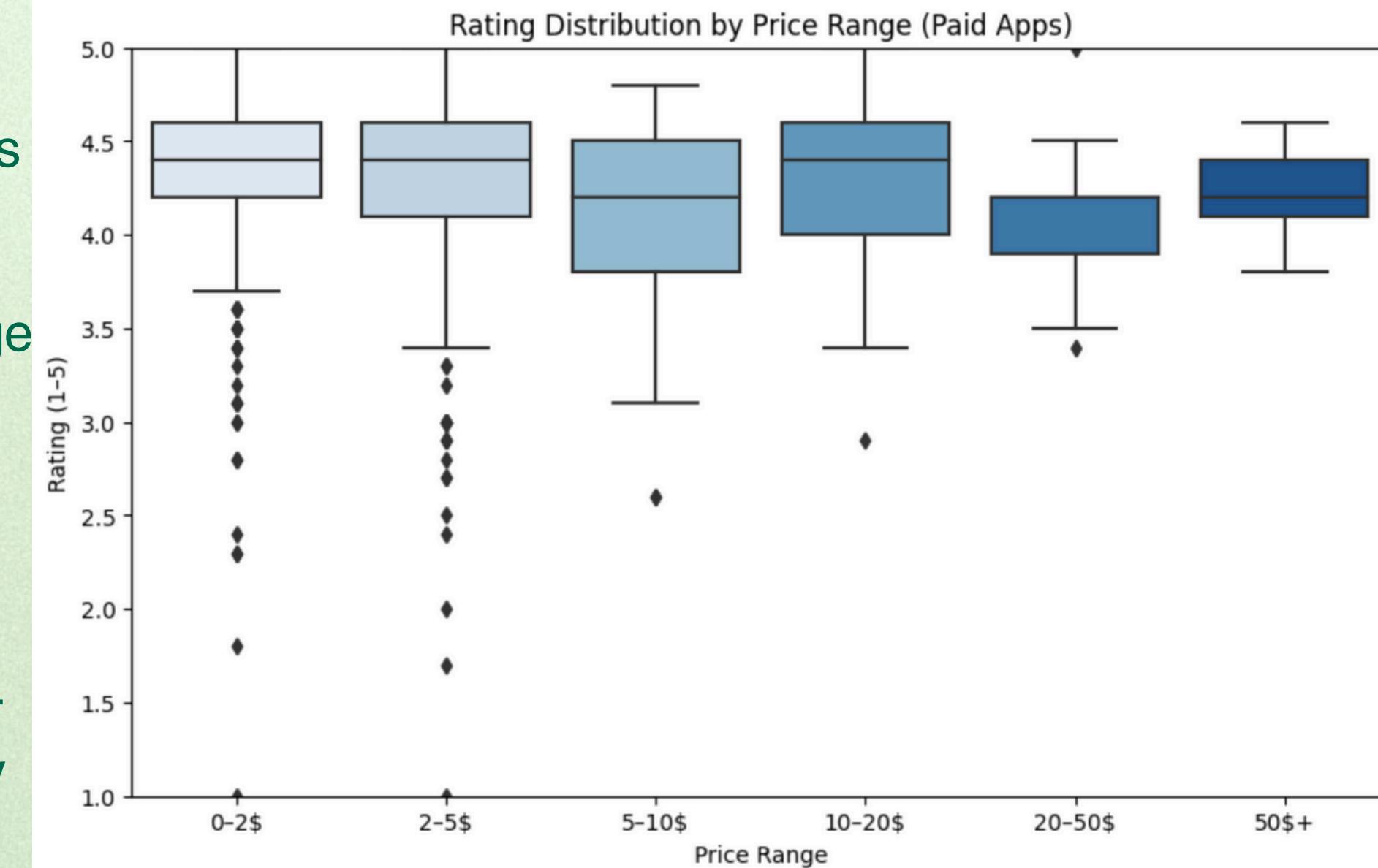
- **Older Kids & Teens Drive the Most Installs:** Apps rated “Everyone 10+” and “Teen” show the highest average installs (~30M and ~27M), indicating that pre-teens and teenagers are the largest drivers of app virality. These segments download games, social apps, and entertainment content at a much higher rate, making them prime targets for apps aiming for rapid growth.



- **“Everyone” Category is Highly Crowded:** Although the “Everyone” rating contains the majority of apps (8,382 apps), its average installs (~12M) are far lower than teen-focused categories. This shows that the general audience segment is oversaturated, making it harder for individual apps to stand out or achieve viral-scale adoption.
- **Mature Apps Struggle to Gain Traction:** Apps rated “Mature 17+” consistently show much lower install averages. This suggests that adult-oriented apps face smaller audiences, stricter content filters, and lower discovery potential, leading to limited mass-market reach compared to teen and youth categories.

# Price Sensitivity Analysis: Does Cost Correlate with Quality?

- **Stable Quality Across Affordable Apps (\$0–\$20):** Across all price tiers up to \$20, the median rating stays consistently high ( $\approx 4.4$ ), showing that price does not determine user satisfaction—a \$1 app delivers nearly the same perceived quality as a \$15 app. Paid apps in this range are generally polished and expectations remain manageable.
- **Quality–Expectation Mismatch in Mid-High Prices (\$20–\$50):** Apps priced between \$20–\$50 show a dip in median rating (~4.2), reflecting a danger zone where expectations exceed delivery. These apps appeal to neither casual users nor true professionals, and high prices amplify user criticism for even small flaws.



- **Professional Tools Shine in Premium Tier (\$50+):** Apps above \$50 have the highest median rating (~4.6), proving that specialized professional apps command both high prices and high satisfaction. Small, targeted user groups purchase these tools for specific needs—leading to fewer downloads but extremely positive reviews when the app performs its intended function well.
- **Cheapest Apps Show the Widest Quality Spread (\$0–\$2):** The lowest price tier has the largest variance, with many low-effort or low-quality apps pulling down the lower quartile, meaning cheap apps include both great gems and extremely poor products, making the category highly inconsistent despite a strong median.

# Analysis of User Sentiment & App Reviews

- Sentiment Analysis & Opinion Mining of Google Play Store Reviews

# Data Methodology (Preprocessing)

Data Source: googleplaystore\_user\_reviews.csv

Initial State: ~64,000 rows (Dirty, Noisy Data).

The Cleaning Pipeline:

Null Handling: Removed 42% of rows containing missing reviews/sentiment.

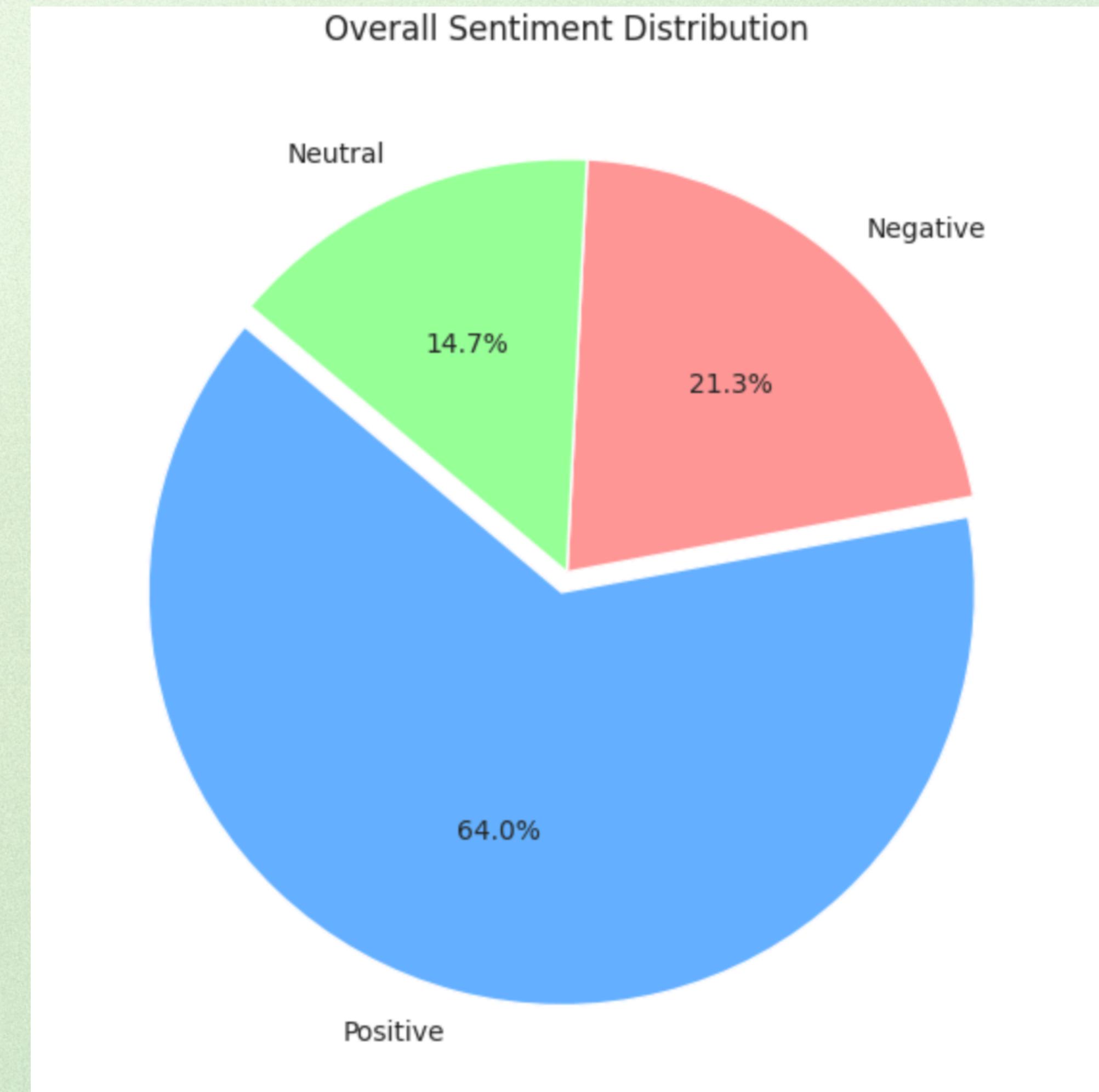
De-Duplication: Removed duplicate entries to prevent "spam bias" in sentiment counts.

Text Normalization: Converted text to lowercase; removed special characters & emojis.

Final Dataset: ~37,000 unique, high-quality reviews.

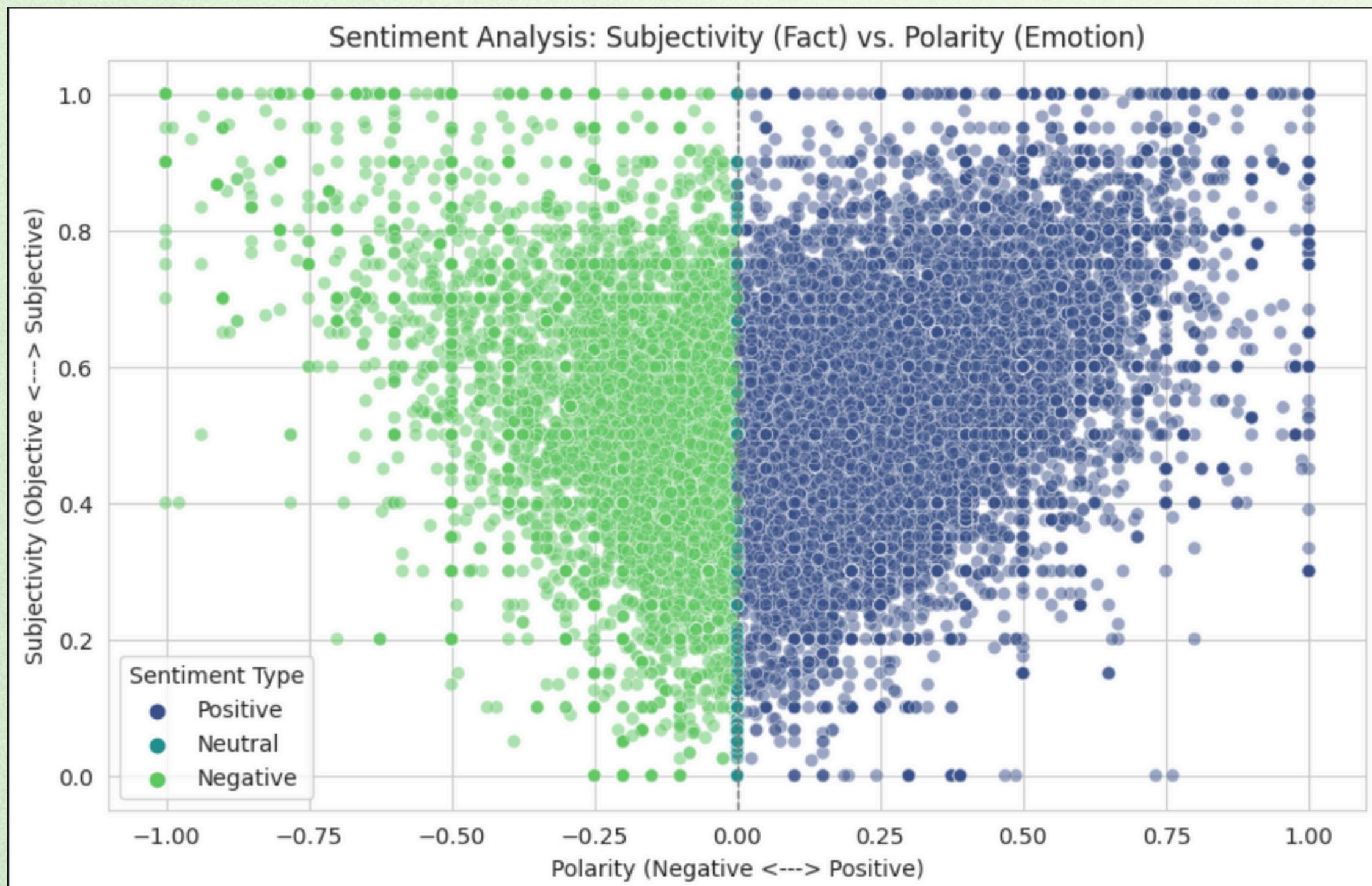
# The Sentiment Landscape

- Positivity Bias: 64% of all reviews are Positive.
- The Minority: Only 22% of reviews are Negative.
- Implication: Users are generally satisfied; however, the minority (Negative) contains the most critical feedback for developers.



# Subjectivity vs. Polarity

- Metric 1: Polarity: Measures Emotion (-1 Negative to +1 Positive).
- Metric 2: Subjectivity: Measures Fact (0) vs. Opinion (1).
- The Correlation: Data exhibits a "V-Shape" distribution.
- Insight:
- Neutral Reviews => Objective Facts (e.g., "App requires login").
- Strong Emotions => Subjective Opinions (e.g., "I hate this interface").



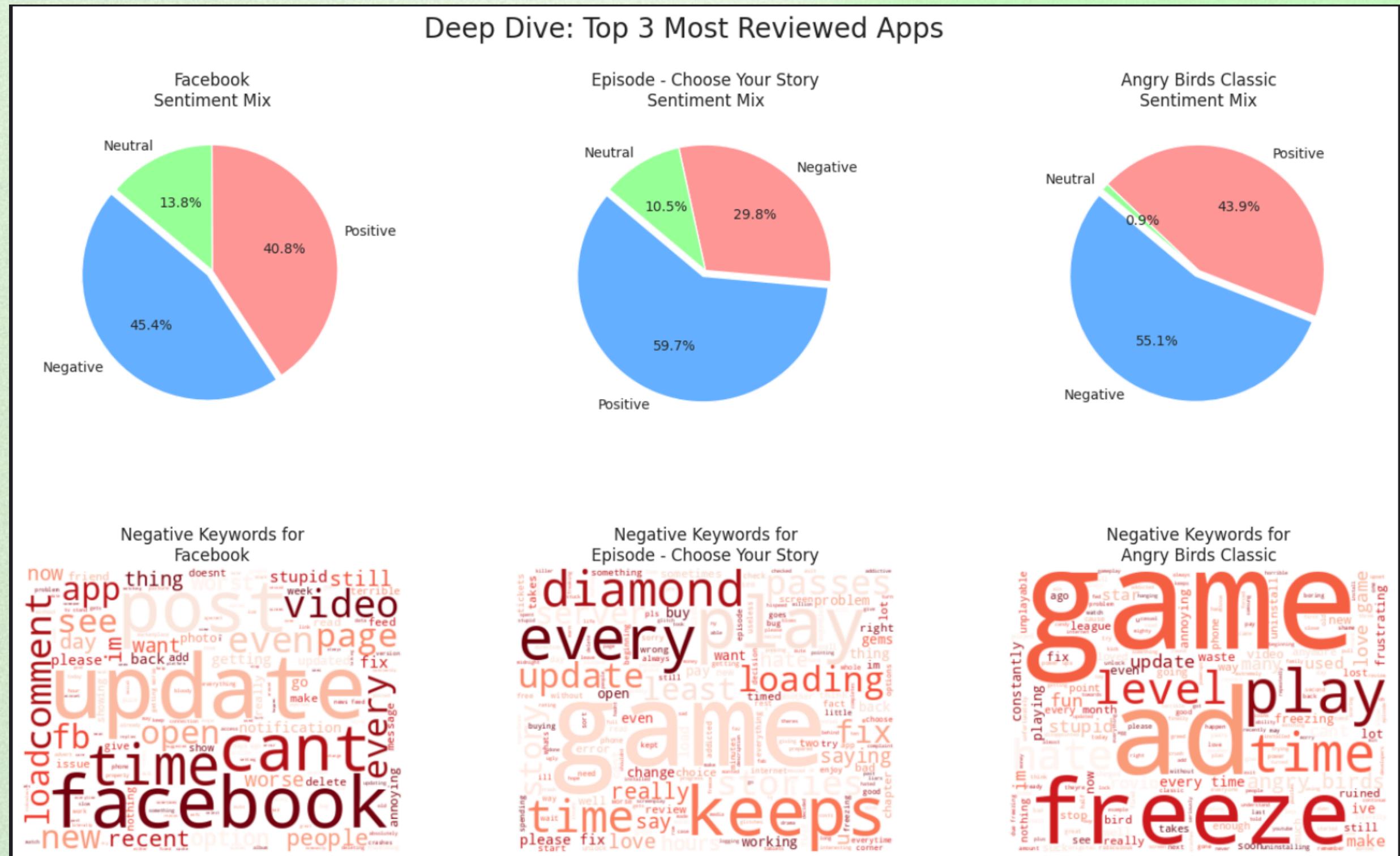
# Mining the "Why" (Negative Keywords)

- Analysis: Filtered for Negative Sentiment + Low Subjectivity.
  - Top Recurring Keywords:
  - "Update", "Fix", "Crash" (Stability Issues)
  - "Screen", "Login", "Slow" (UI/UX Issues)
  - "Ads", "Money" (Monetization Issues)



# Case Studies (Top 3 Apps)

- Comparison: Analyzed the 3 most reviewed apps in the dataset.
- Volume != Satisfaction: The most reviewed app does not necessarily have the highest polarity score.
- Category Behavior:
- Games: Complaints focus on "Levels" and "Ads".
- Social: Complaints focus on "Updates" and "Privacy".



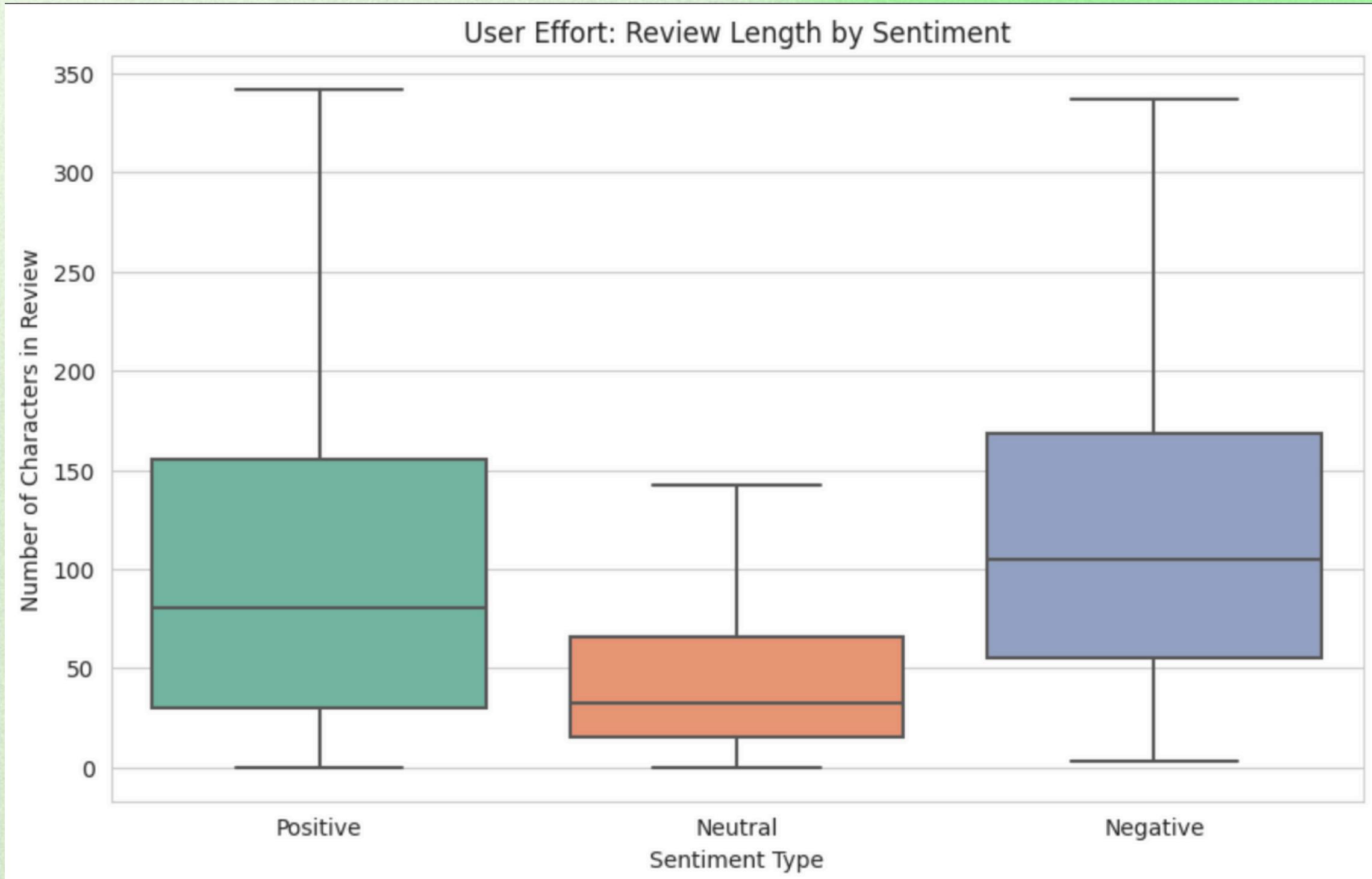
# User Effort Analysis

Hypothesis: Do angry users write more?

Data: Calculated character count per review.

Result: Negative reviews are statistically longer than Positive reviews.

Interpretation: Happy users leave short praise ("Good app"). Dissatisfied users write detailed reports to explain their frustration.



# Future Scope

- PHASE 3 (NEXT STEP): MERGING SENTIMENT DATA WITH METADATA.
- GOAL: DETERMINE WHICH APP CATEGORIES ARE THE MOST "TOXIC."
- GOAL: CORRELATION BETWEEN PRICE AND SENTIMENT (DO PAID APPS HAVE HAPPIER USERS?).
- TECHNIQUE: WILL USE APP NAME AS THE PRIMARY KEY TO JOIN THE TWO DATASETS.

Presented by Juliana Silva

# THANK YOU SO MUCH!

Presentations Templates are communication tools that  
can be used as lectures, reports, and more.

[www.reallygreatsite.com](http://www.reallygreatsite.com)