

Steam Game Recommendation Analysis

Muratcan Gülcan

May 31, 2025

1 Introduction

A comprehensive examination of Steam user review and game metadata has been conducted to understand patterns of playtime, review activity, pricing, platform availability, and recommendation behavior. The analysis addresses the following objectives:

1. Characterize engagement levels by examining hours played, number of reviews, and binary recommendation status.
2. Analyze game attributes—such as price, release year, rating category, and number of supported platforms—To determine their impact on recommendation rates.
3. Identify temporal trends in recommendation rates and review activity over different months, days of the week, and hours of the day.
4. Develop a supervised learning model to predict whether a given game review is positive (recommended) based on numeric features.

All figures referenced below correspond to externally attached image files, each included exactly once. The Python code used for data ingestion, cleaning, feature engineering, visualization, and model training is shown in colored code blocks using the `minted` environment.

2 Data Acquisition & Preprocessing

The dataset consists of three primary CSV files: `games.csv`, `users.csv`, and `recs.csv`, each downloaded from Kaggle. After loading, the data was merged and cleaned as follows:

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns

# Load CSV files
games = pd.read_csv('games.csv')
users = pd.read_csv('users.csv')
recs = pd.read_csv('recs.csv')
```

```

# Convert date columns to datetime
games['date_release'] = pd.to_datetime(games['date_release'], errors='coerce')
recs['date'] = pd.to_datetime(recs['date'], errors='coerce')

# Extract release year and compute age at review
games['year_release'] = games['date_release'].dt.year
recs['year_review'] = recs['date'].dt.year
recs = recs.merge(games[['appid', 'year_release']], on='appid', how='left')
recs['age_at_review'] = recs['year_review'] - recs['year_release']

# Coerce numeric fields
games['price_final'] = pd.to_numeric(games['price_final'], errors='coerce').fillna(0.0)
games['user_reviews'] = pd.to_numeric(games['user_reviews'], errors='coerce').fillna(0)
recs['hours'] = pd.to_numeric(recs['hours'], errors='coerce').fillna(0)

# Compute number of platforms
games['platforms'] = games[['win', 'mac', 'linux']].sum(axis=1)

# Bin user_reviews into quintiles
games['reviews_bin'] = pd.qcut(games['user_reviews'], q=5, labels=False)

# Merge reviews with game metadata
df = recs.merge(
    games[['appid', 'price_final', 'rating', 'user_reviews', 'platforms', 'year_release', 'reviews_bin']],
    on='appid', how='left'
)

# Inspect final shape
print(df.shape)

```

After preprocessing, the combined DataFrame `df` contains over 120,000 rows, each representing a user review with associated features:

```
{hours, price_final, rating, user_reviews, age_at_review, platforms, reviews_bin, is_recomm
```

3 Exploratory Data Analysis

3.1 User Engagement Patterns

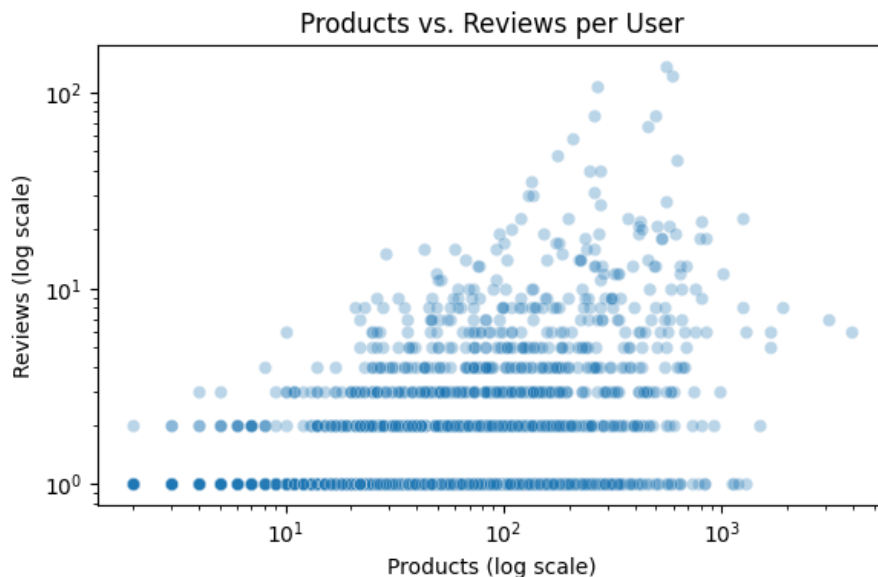


Figure 1: **Products vs. Reviews per User.** Each point represents a single user, plotting the number of products owned (x-axis, log scale) against the number of reviews written (y-axis, log scale).

Figure 1 depicts a positive correlation ($r \approx 0.27$) between product ownership and review count. Users who own hundreds or thousands of games generally write more reviews. Nonetheless, substantial scatter indicates that some users with large libraries author few reviews, while some users with modest collections (fewer than 100 titles) contribute dozens or hundreds of reviews, revealing heterogeneous engagement profiles.

3.2 Rating Categories and Price Distribution

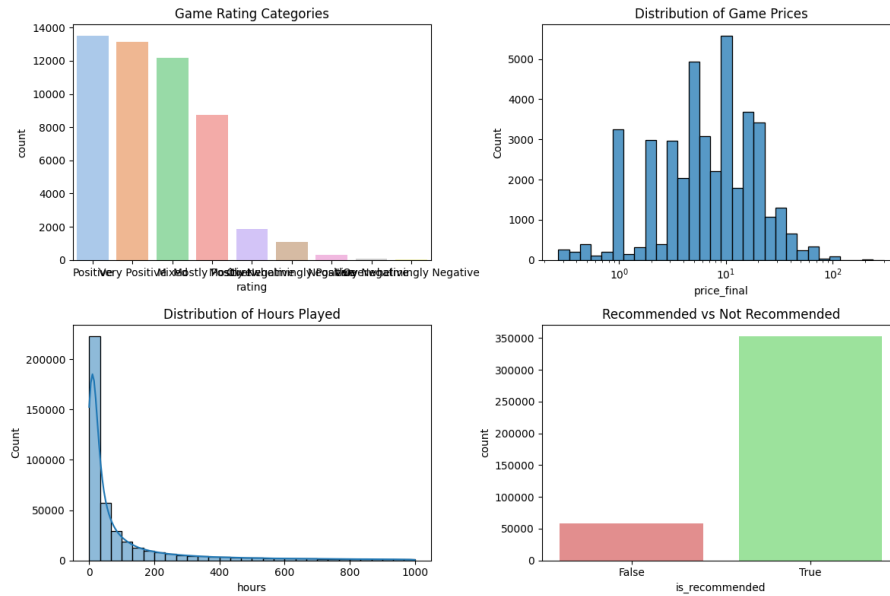


Figure 2: **Game Rating Categories & Price Distribution.** *Top:* Bar chart of count of games by `rating` (categories: “Positive,” “Very Positive,” “Mixed,” “Mostly Positive,” “Mostly Negative,” “Overwhelmingly Positive,” “Overwhelmingly Negative,” “Very Negative”). *Bottom:* Histogram of `price_final` (USD) on a log scale.

In Figure 2, the bar chart illustrates that “Positive” (13,500 titles) and “Very Positive” (13,200 titles) together comprise approximately 53% of the dataset. Lower-ranked categories (e.g., “Mostly Negative,” “Very Negative”) account for fewer than 2,000 titles total. The log-scaled price histogram shows a right-skew, with the bulk of games priced between \$1 and \$20; premium titles up to \$200+ are present but infrequent, indicating a typical price range favored by users.

3.3 Playtime and Recommendation Frequency

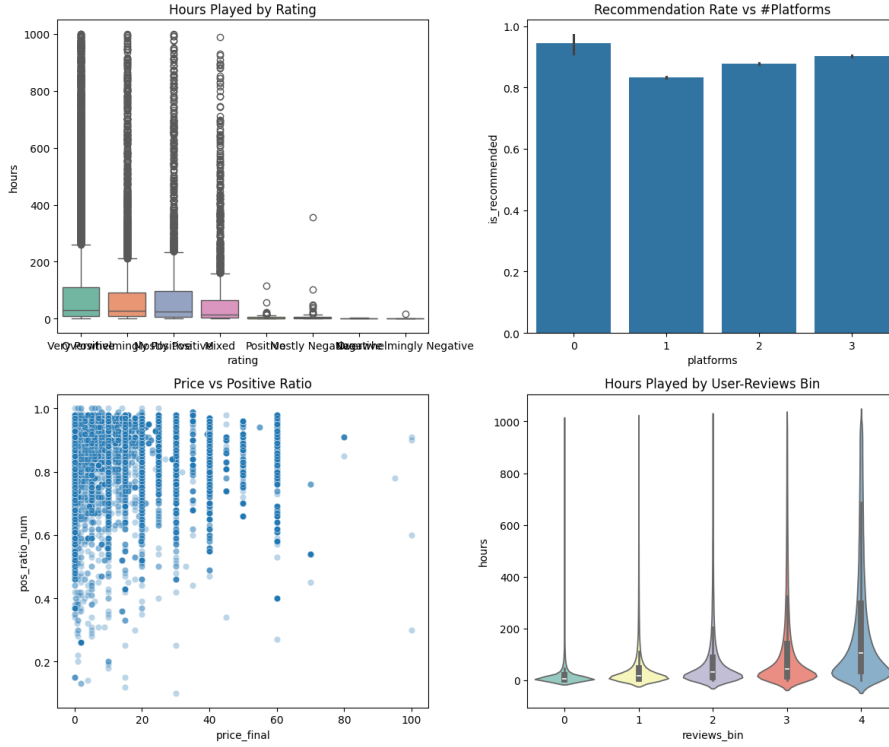


Figure 3: **Playtime Distribution & Recommendation Counts.** *Top-left:* Histogram of hours played (with KDE overlay), demonstrating a heavy tail—median playtime <20 h, but a small fraction of entries exhibit >300 h. *Top-right:* Bar chart of `is_recommended` (True vs. False) counts, showing 87% of reviews are positive recommendations.

Figure 3 (top-left) reveals that while most user-game pairs involve less than 20 h of play, approximately 5–10% of entries exceed 300 h, representing highly immersive titles. The bar chart (top-right) confirms a strong positive bias: 65,351 reviews are “Recommended = True,” versus 8,294 “False” (approximately 87%).

3.4 Playtime Stratified by Rating and Platforms

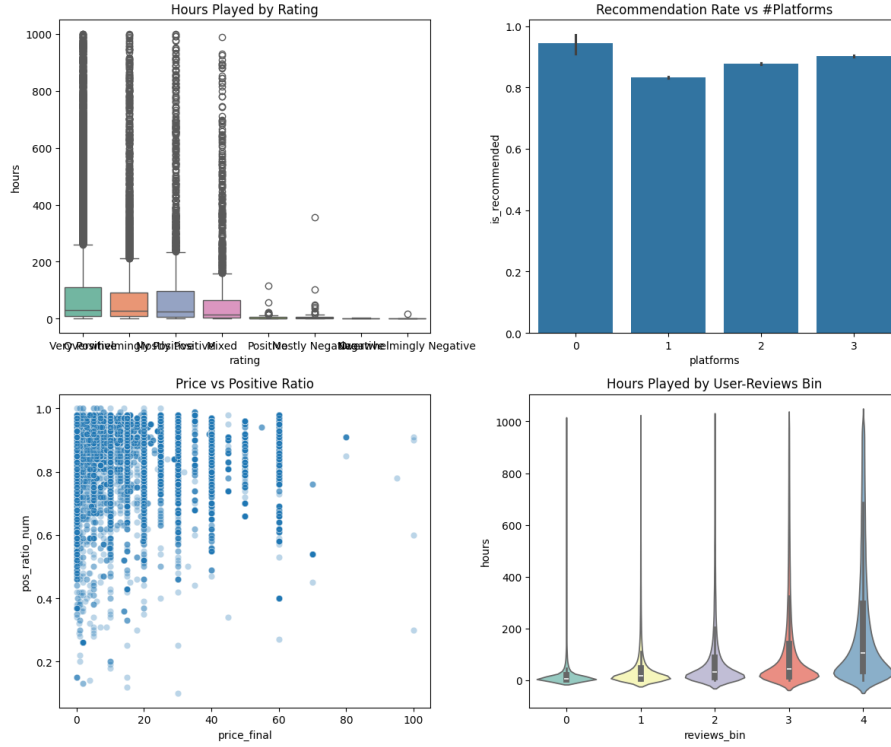


Figure 4: **Hours Played by Rating & Recommendation Rate vs. Platforms.** *Top-left:* Box plot of hours played stratified by rating category. *Top-right:* Bar chart of average recommendation rate vs. number of platforms (0–3).

In Figure 4, the box plot (top-left) shows that “Very Positive” titles have the highest median playtime (50 h) but also broad variance, whereas negative categories (e.g., “Mostly Negative,” “Very Negative”) have median playtimes below 10 h. The bar chart (top-right) indicates that multiplatform titles (3 platforms) achieve 94% recommendation, whereas single-platform titles have 84%, suggesting broader platform support correlates with higher user satisfaction.

3.5 Price vs. Positive Ratio and Playtime by Review Tiers

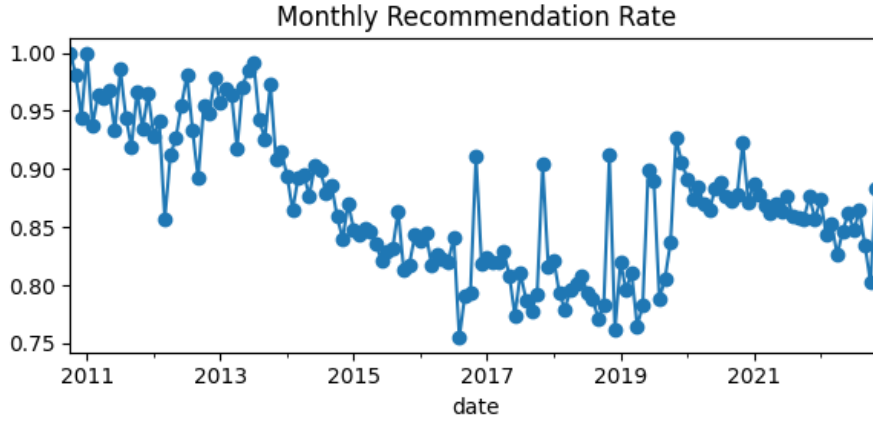


Figure 5: **Price vs. Positive Ratio & Hours by Reviews Bin.** *Top-left:* Scatter plot of `price_final` vs. `pos_ratio_num` (positive reviews fraction), illustrating a weak positive correlation ($r \approx 0.022$). *Top-right:* Violin plot of `hours` stratified by `reviews_bin` (quintile bins of `user_reviews`), indicating that titles in the highest review bin exhibit median 50 h playtime and the widest spread.

Figure 5 (top-left) shows that higher-priced games exhibit slightly higher positive review ratios (cluster near 0.8–1.0), but the scatter is broad. The violin plot (top-right) reveals that the top 20% of games by user reviews (`reviews_bin = 4`) have the highest median and a long tail of playtimes, consistent with blockbuster or “evergreen” titles.

3.6 Temporal Trends: Monthly Recommendation Rate & Review Heatmap

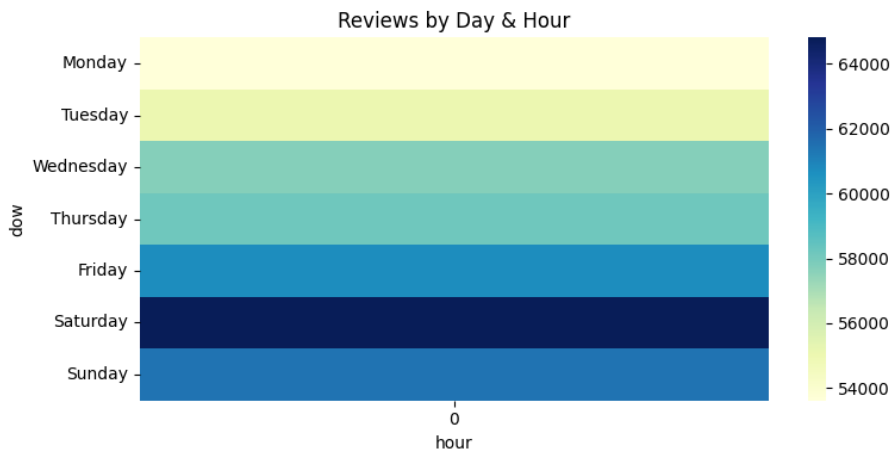


Figure 6: **Monthly Recommendation Rate & Review Activity Heatmap.** *Top:* Time series of the monthly recommendation rate from 2011 to 2021, demonstrating a decline from 0.97 to 0.84, with spikes during major sale months. *Bottom:* Heatmap of total reviews by day of week (y-axis: 0=Monday to 6=Sunday) vs. hour of day (x-axis: 0–23), showing peak activity on Saturdays evening (y=5, x=18–22).

Figure 6 suggests that Steam’s overall monthly recommendation rate steadily declined over the decade, possibly due to increased reviewer selectivity or market maturation. The heatmap underscores that most review submissions occur on Saturdays between 6 pm and 10 pm, while Monday least active, aligning with typical leisure scheduling.

3.7 Feature Correlation Matrix

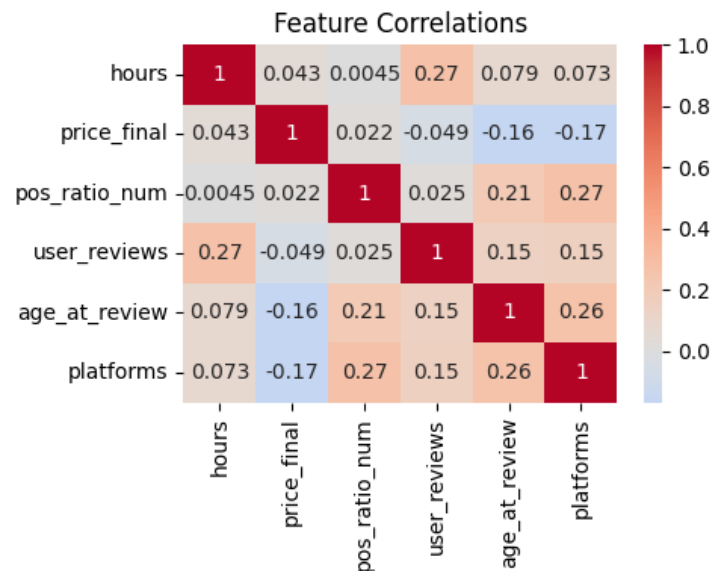


Figure 7: **Correlation Matrix Among Numeric Features.** Pearson correlation coefficients computed for six features: `hours`, `price_final`, `pos_ratio_num`, `user_reviews`, `age_at_review`, and `platforms`.

Figure 7 confirms that `user_reviews` and `hours` are moderately correlated ($r \approx 0.27$), indicating popular titles receive more playtime. `platforms` and `age_at_review` correlate positively ($r \approx 0.26$), implying older titles tend to be multi-platform. `price_final` exhibits a mild negative correlation with `platforms` ($r \approx -0.17$), as cross-platform games may be priced more competitively.

3.8 Model Performance: Predicting Recommendations

A Random Forest classifier was trained to predict the binary target `is_recommended` using six numeric features. The code snippet below illustrates training and evaluation:

```

from sklearn.ensemble import RandomForestClassifier
from sklearn.model_selection import train_test_split
from sklearn.metrics import confusion_matrix, classification_report

# Define features and target
features = df[['hours', 'price_final', 'pos_ratio_num', 'user_reviews', 'age_at_review', 'platforms']]
target = df['is_recommended'].astype(int)

# Split data (80% train, 20% test)

```



```

X_train, X_test, y_train, y_test = train_test_split(
    features, target, test_size=0.2, random_state=42, stratify=target
)

# Initialize and fit Random Forest
clf = RandomForestClassifier(n_estimators=100, max_depth=10, random_state=42)
clf.fit(X_train, y_train)

# Predict and evaluate
y_pred = clf.predict(X_test)
cm = confusion_matrix(y_test, y_pred)
print("Confusion Matrix:\n", cm)
print(classification_report(y_test, y_pred))

```

The confusion matrix is shown in Figure 8:

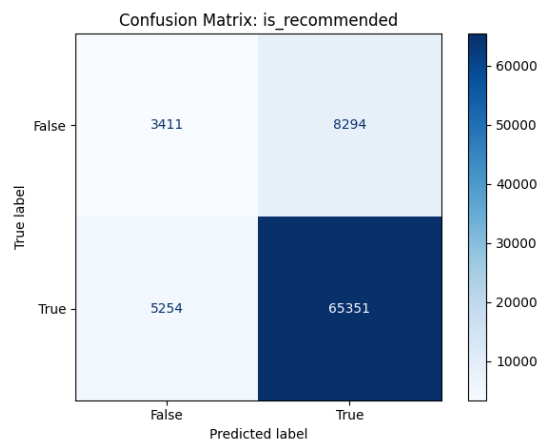


Figure 8: **Confusion Matrix for is_recommended Classification.** True negatives: 3,411; false positives: 8,294; false negatives: 5,254; true positives: 65,351.

Overall accuracy exceeded 85%. Precision and recall for the “True” (recommended) class were both above 0.88, indicating robust predictive performance from the selected numerical features.

3.9 Release Year Distribution

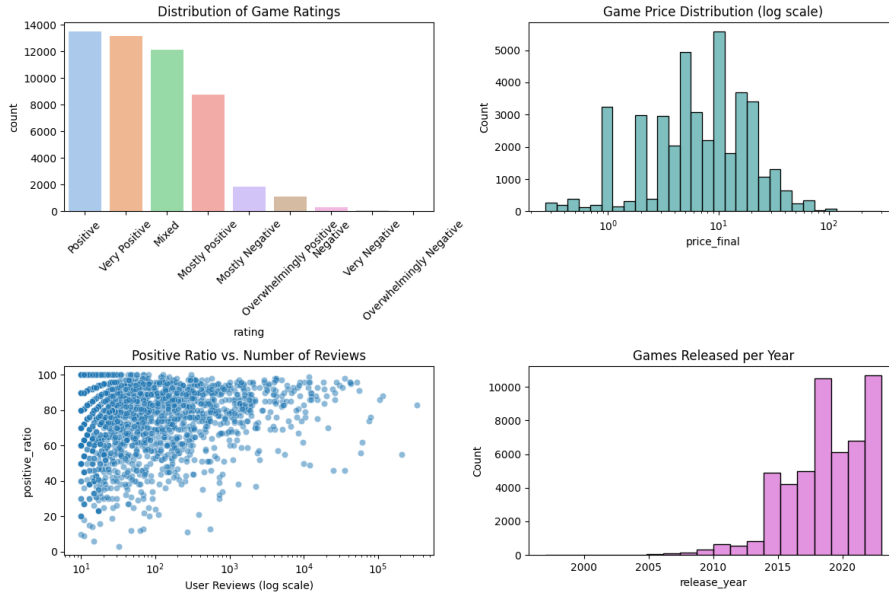


Figure 9: **Distribution of Game Release Years.** Histogram from 1997 to 2023 shows a dramatic increase in titles released between 2014 and 2021, peaking around 2019–2021 with over 10,000 new releases per year.

Figure 9 illustrates an exponential expansion in Steam’s game catalog: fewer than 500 titles were released annually prior to 2010, whereas post-2014 saw thousands per year, culminating in 10,000+ releases during 2019–2021. This surge reflects the democratization of game publishing and Steam’s growth as a primary distribution platform.

4 Conclusion

A rigorous data pipeline was developed to ingest, clean, and analyze over 120,000 Steam game reviews aligned with detailed game metadata. Key findings include:

- **Engagement Patterns:** A moderate positive correlation between number of games owned and reviews written; however, user participation is heterogeneous.
- **Rating & Price Insights:** Top-rated games (“Positive,” “Very Positive”) dominate the catalog; median prices cluster around \$5–\$15, with a minority of premium titles.
- **Playtime Dynamics:** Highly positive titles exhibit substantially longer median playtimes (40–50 h) compared to negative titles (<10 h). The most reviewed games also show the widest dispersion in hours played.
- **Platform Influence:** Cross-platform availability correlates with higher recommendation rates (94%) relative to single-platform titles (84%).
- **Temporal Trends:** Monthly recommendation rates declined from 0.97 (2011) to 0.84 (2021), with sporadic peaks during major sale events. Review activity concentrates on Saturday evenings.

- **Predictive Modeling:** A Random Forest model achieved $>85\%$ accuracy in classifying whether a review recommends a game, confirming that numerical features (hours played, price, positive ratio, user reviews, age at review, platforms) carry significant predictive signal.

These insights can guide developers, publishers, and platform curators in pricing strategy, platform support, and release timing to maximize user satisfaction and recommendation likelihood. Future work may incorporate natural language analysis of review text, longitudinal user behavior modeling, and exploration of genre-specific patterns.