



# MOVIE LENS

## MOVIE RATINGS AND METADATA REPORT

This report presents an analysis of the dataset gotten from the movie lens video streaming website focusing on feature engineering and exploratory insights that illuminate patterns in user behaviour, genre performance, and rating dynamics.

### **Engineered features**

Some features were gleaned from the dataset and included as new features to give new perspectives.

### **Dataset summary**

Total unique users: 610, Total ratings: 100836, Total unique movies: 9724, Date range of ratings: 1996 – 2018, Release year range: 1902 – 2018

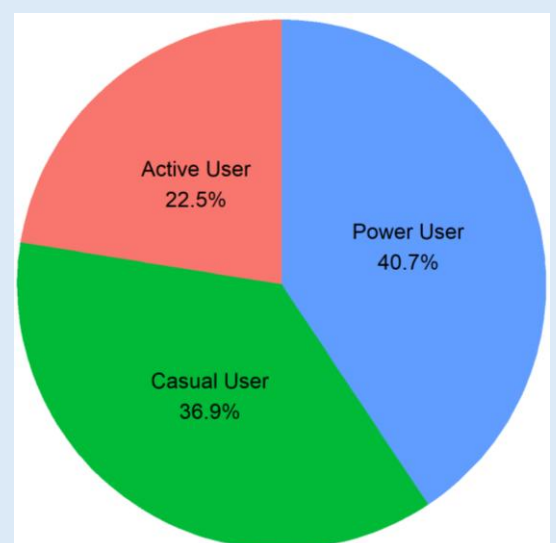
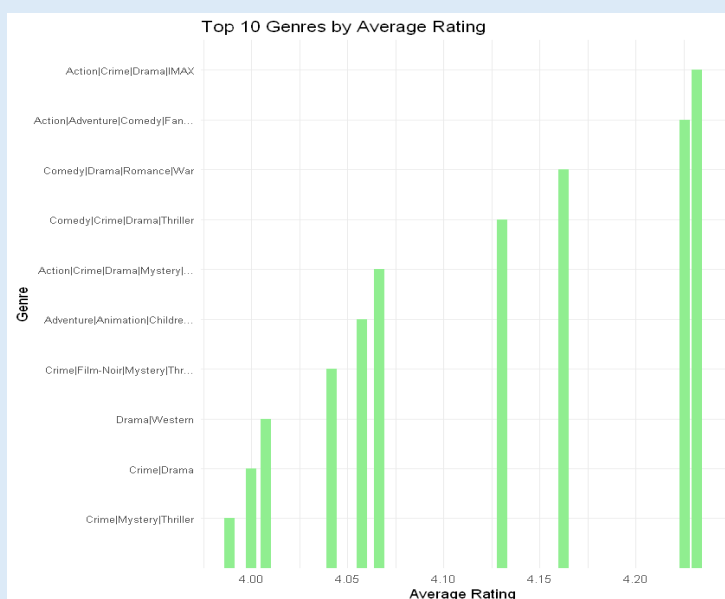
Rating stats: minimum rating: 0.5, median rating: 3.5, mean rating: 3.502, maximum: 5.0

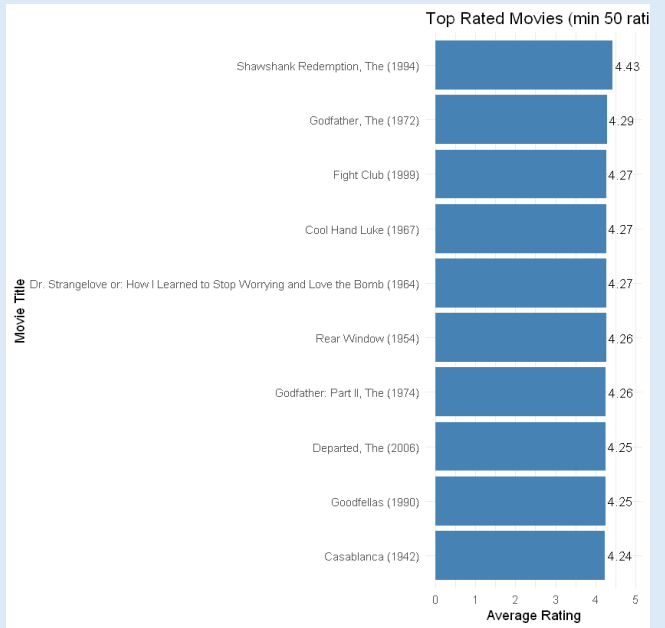
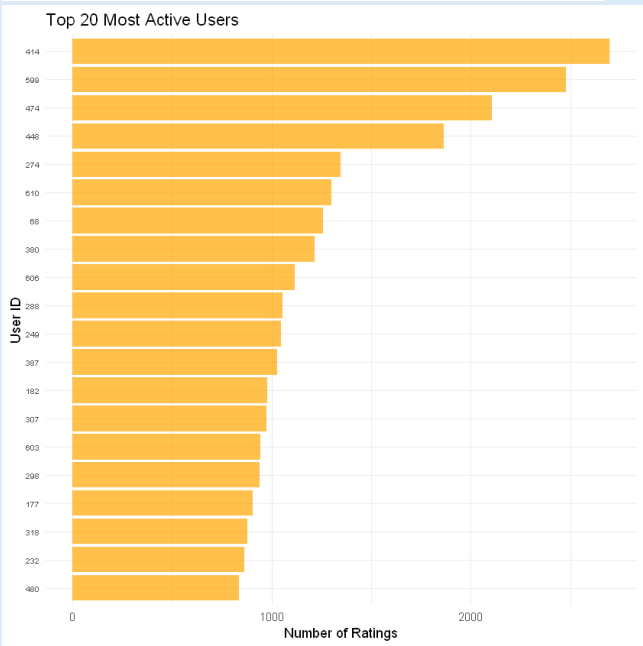
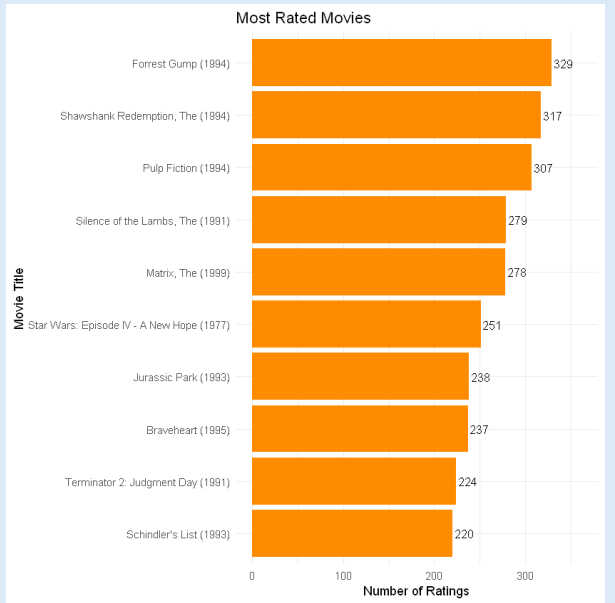
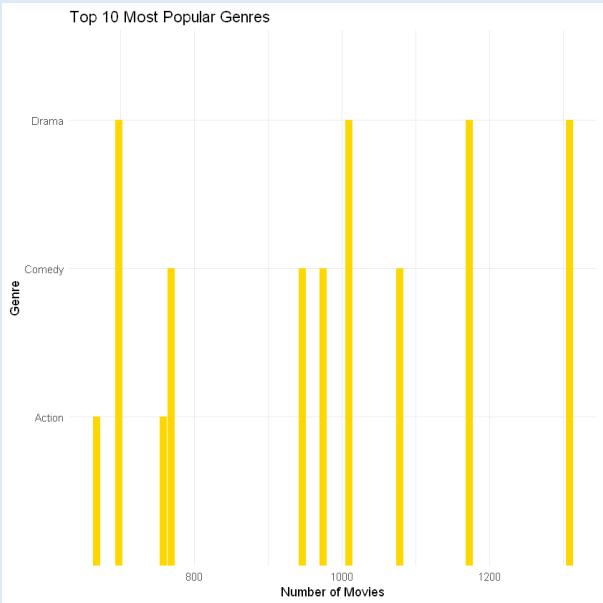
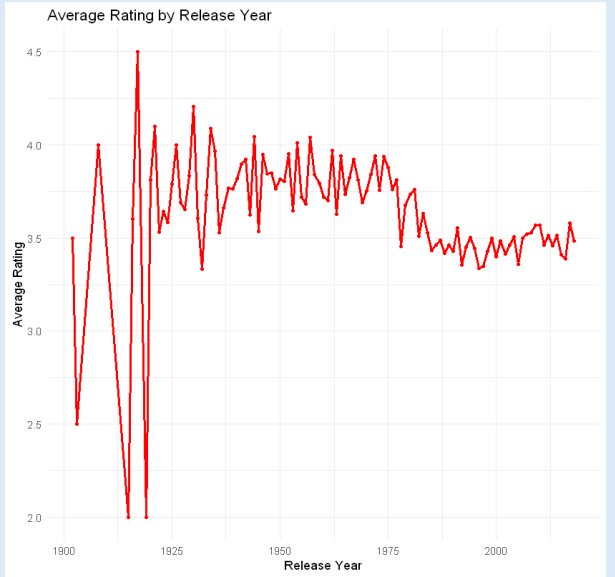
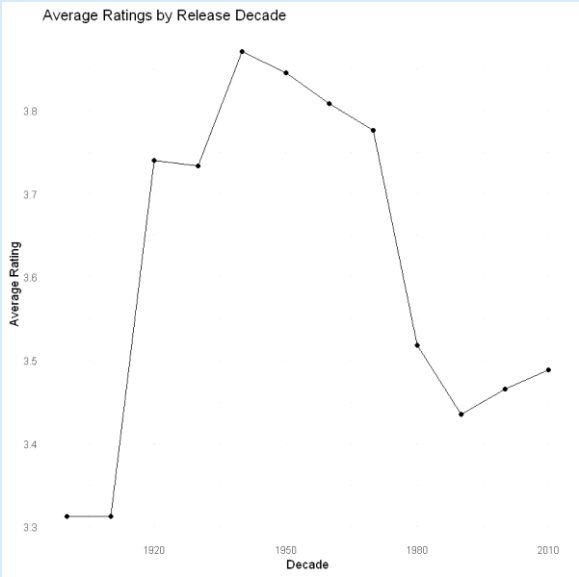
## Engineered features

Feature name	Description	Purpose
<b>release_year</b>	Extracted from movie titles	Enables trend detection capture era effects (temporal preferences). Useful for time-aware recommendations
<b>genre_count</b>	Number of genres per movie	Captures genre diversity and complexity, proxy for 'broad appeal' vs niche; informs content-based filtering and similarity measures.
<b>movie_avg_rating</b>	Average rating per movie	Measures overall movie popularity and quality
<b>genre_avg_rating</b>	Average rating per genre	Assesses genre performance and user preferences
<b>user_rating_count</b>	Number of ratings per user	Identifies user engagement levels, proxy for user activity
<b>is_classic</b>	Boolean flag for movies released before 2000	Supports segmentation by era and nostalgia-based recommendations

## Key Insights Discovered

- ❖ **Top Movies:** Shawshank Redemption, Godfather, and Fight Club lead in ratings with substantial review counts.
- ❖ **Genre Popularity:** Drama, Comedy, and Action dominate in volume.
- ❖ **Genre Performance:** Genres like Action|Crime|Drama|IMAX and Action|Adventure|Comedy|Fantasy received the highest average ratings.
- ❖ **Average Ratings per Decade** reveals that movies in the 1900s are generally more liked than movies in 2000s, the 1940's were the golden era of movies (based on this dataset)
- ❖ **Genre Count vs Rating:** there is no definitive indication that movies with more genres tend to receive higher ratings
- ❖ **Genre Popularity:** Drama, Comedy, and Action dominate in volume.
- ❖ **Engagement by Movie Type:** Classic movies had slightly lower engagement but similar median ratings.
- ❖ **User Segmentation:** Users were classified into Power(>100), Active(50-100), and Casual (<50) based on rating frequency.





# How Do The Features Support a Future Recommendation System?

- ❖ **User Profiling:** `user_rating_count` and rating behavior help tailor suggestions based on each user's engagement level.
- ❖ **Content-Based Filtering:** `genre_count`, `genres`, and `release_year` enable matching users to similar movies. Users that prefer certain genres will get recommendations of those genres.
- ❖ **Collaborative Filtering:** `movie_avg_rating` and `user_type` support clustering users with similar tastes and movie preferences from the clusters can inform recommendations for similar users.
- ❖ **Temporal Dynamics:** `rating_year` and `release_year` allow for trend-aware recommendations. Ratings often spike around holidays or after major events (e.g., Oscar wins, streaming releases). `rating_year` allows the system to factor in these temporal spikes and suggest trending content.
- ❖ **Cold Start Solutions:** Genre and tag-based features help recommend new or unrated movies. When a movie has no ratings yet, genre and tag-based features allow the system to:
  - Compare the new movie to others with similar genres or tags and recommend it to users who have shown interest in those genres or themes.
  - New users will answer onboarding questions on preferred genres, this will be used for recommendations
- ❖ **Diversity and Novelty:** Insights into genre performance and popularity help balance familiar vs. novel suggestions. By analysing genre performance, we can introduce diverse genres that are still well-rated, expanding user horizons without sacrificing quality.

## Conclusion

This analysis of the MovieLens dataset provided a strong foundation for understanding user–item interactions and feature dynamics. Through feature engineering and visual exploration, we identified patterns such as consistent user rating tendencies, genre-based preferences, and temporal trends in movie popularity. These findings support the development of a robust recommendation system by enabling:

- User segmentation for personalized recommendations.
- Trend-aware filtering using temporal features like `release_year` and `rating_year`.
- Cold start solutions through genre and tag-based similarity.
- Balanced suggestions that promote both diversity and novelty.

*By integrating collaborative filtering with temporal features, future models can deliver more accurate, engaging, and personalized recommendations.*