



AI 705 Recommendation Systems: Midterm Project

Movie Recommendation System

IMT2021009 Madhav Sood

IMT2021065 Vihan Vashishth

IMT2021066 Yukta Rajapur

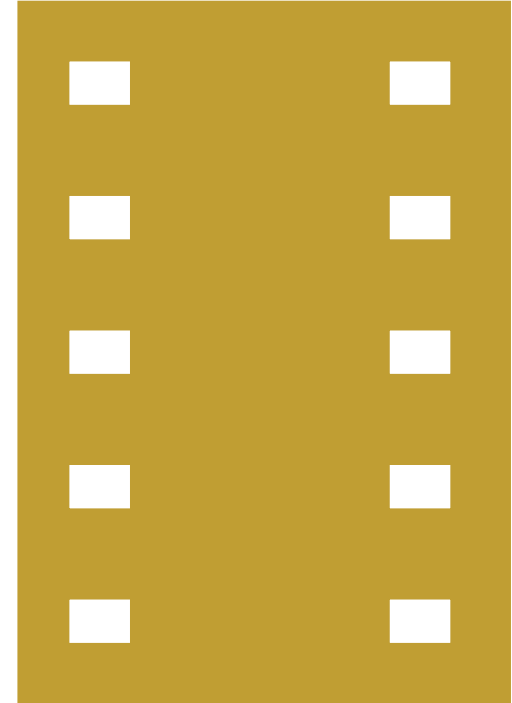
IMT2021067 Brij Desai

Dataset and Problem Statement

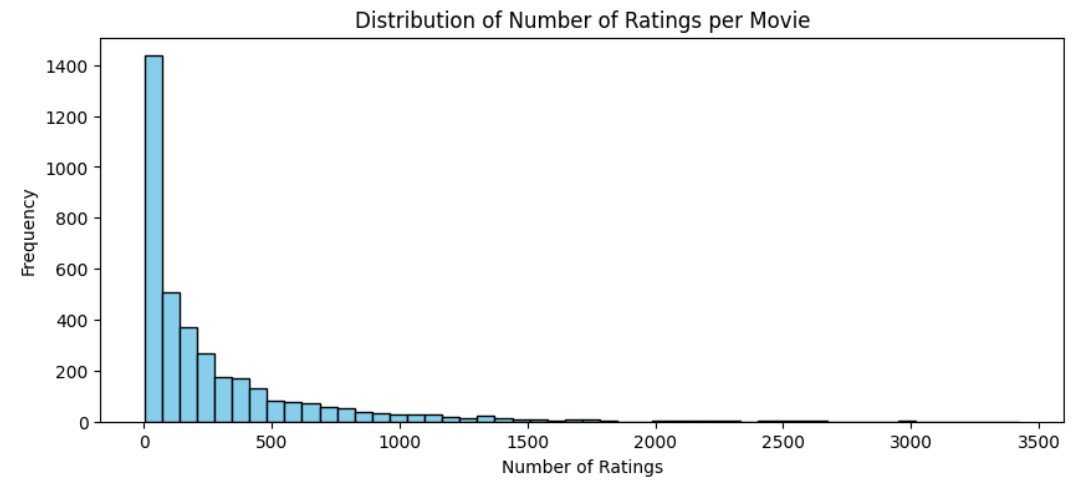
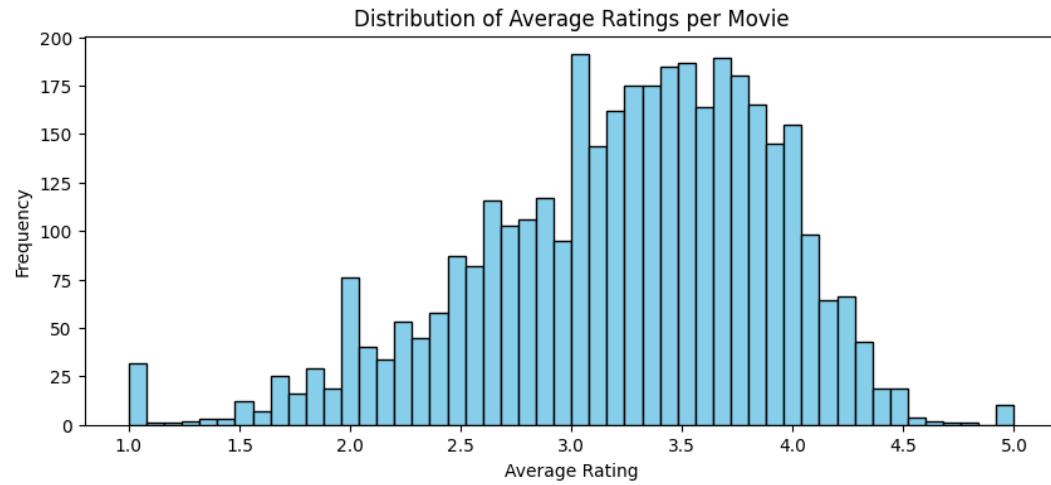
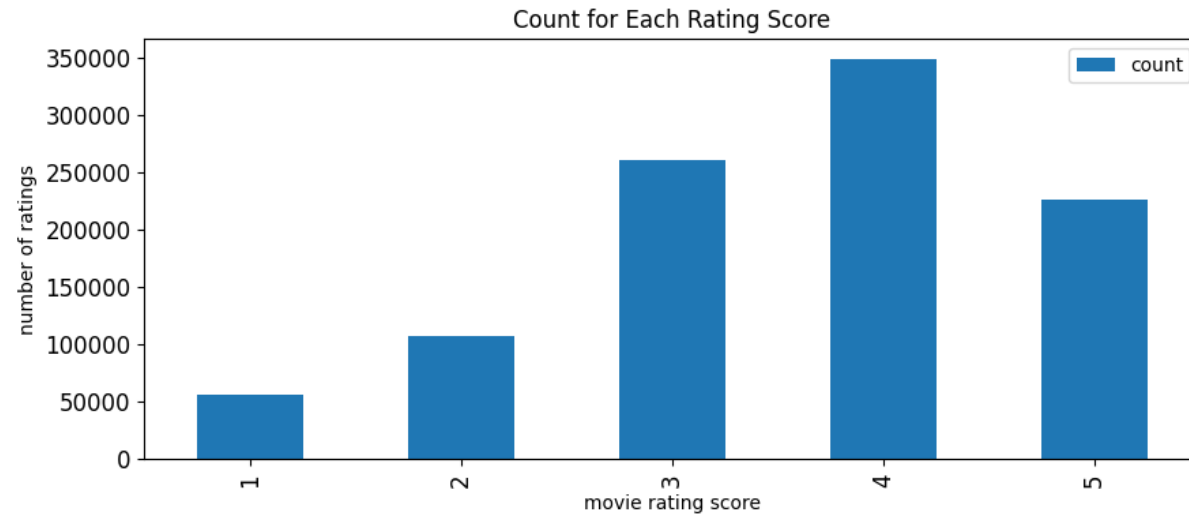
- Dataset used is the MovieLens 1M Dataset with ratings of 6040 users across nearly 4000 movies.

Problem Statement:

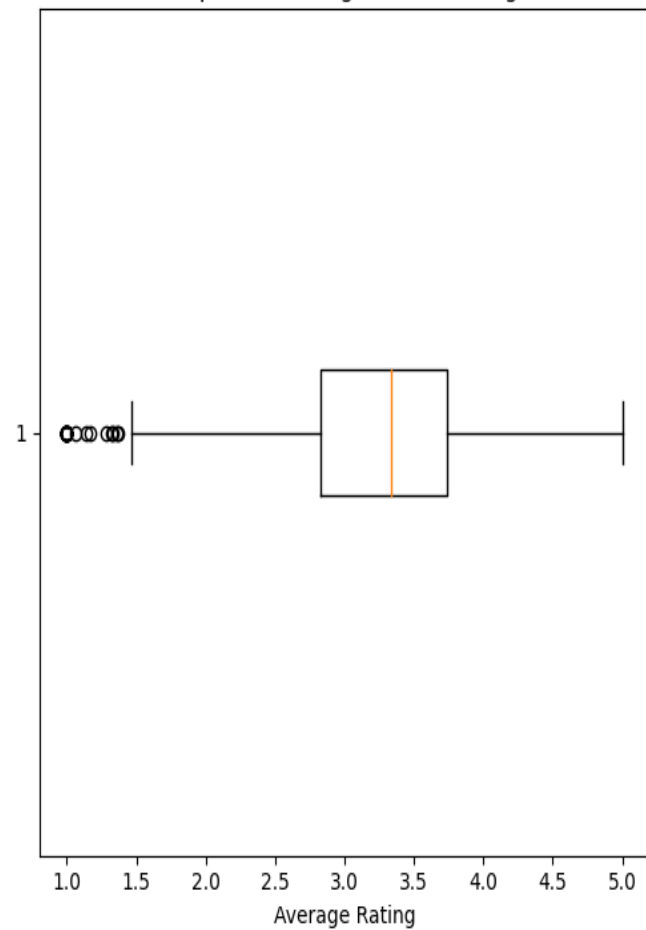
- Build a movie recommendation system, where based on the existing ratings of a user, you should suggest new movies that the user will like.
- Implement necessary methods while adding our own novel ideas.



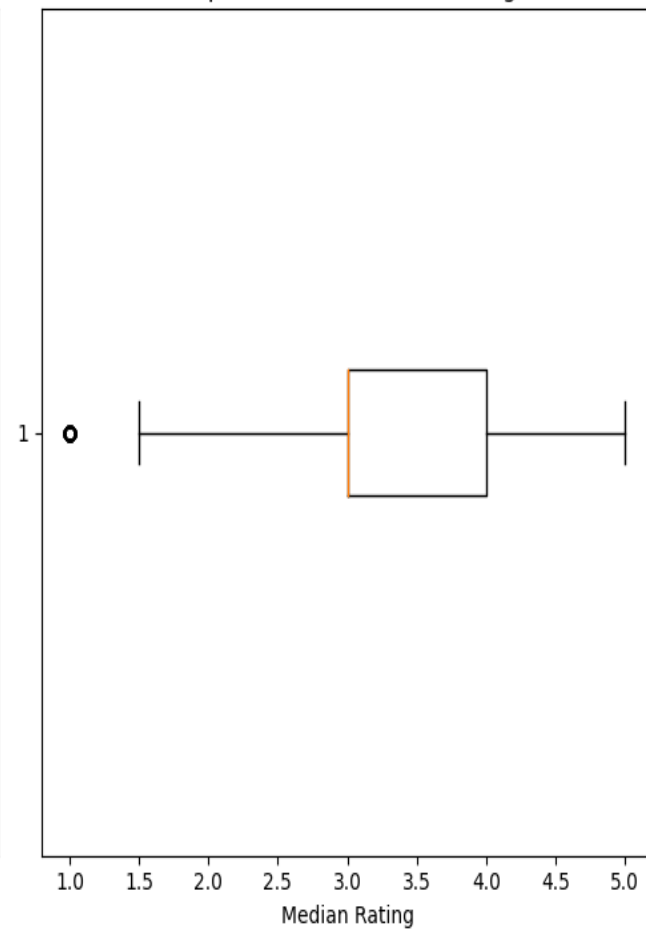
EDA



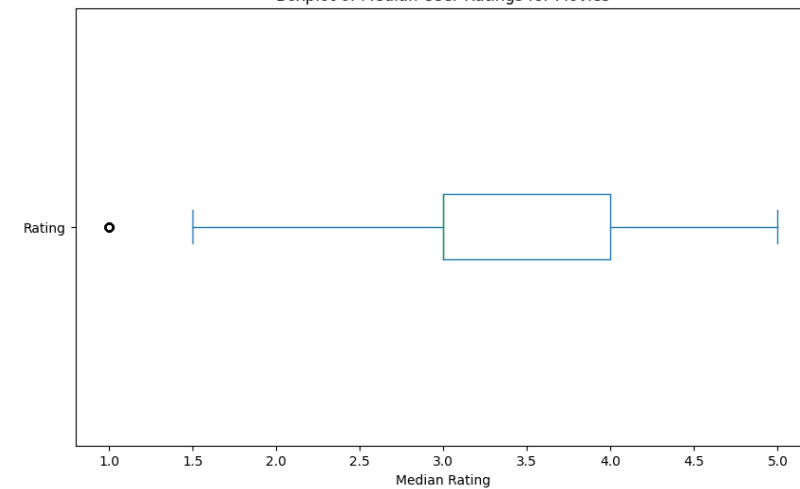
Boxplot of Average Movie Ratings



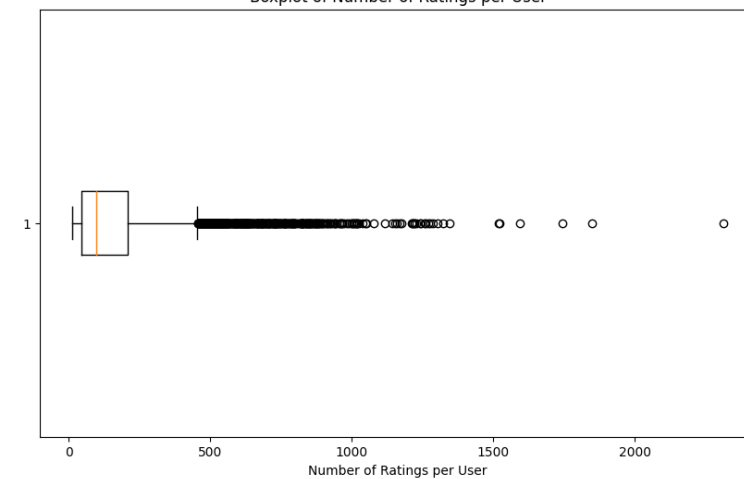
Boxplot of Median Movie Ratings



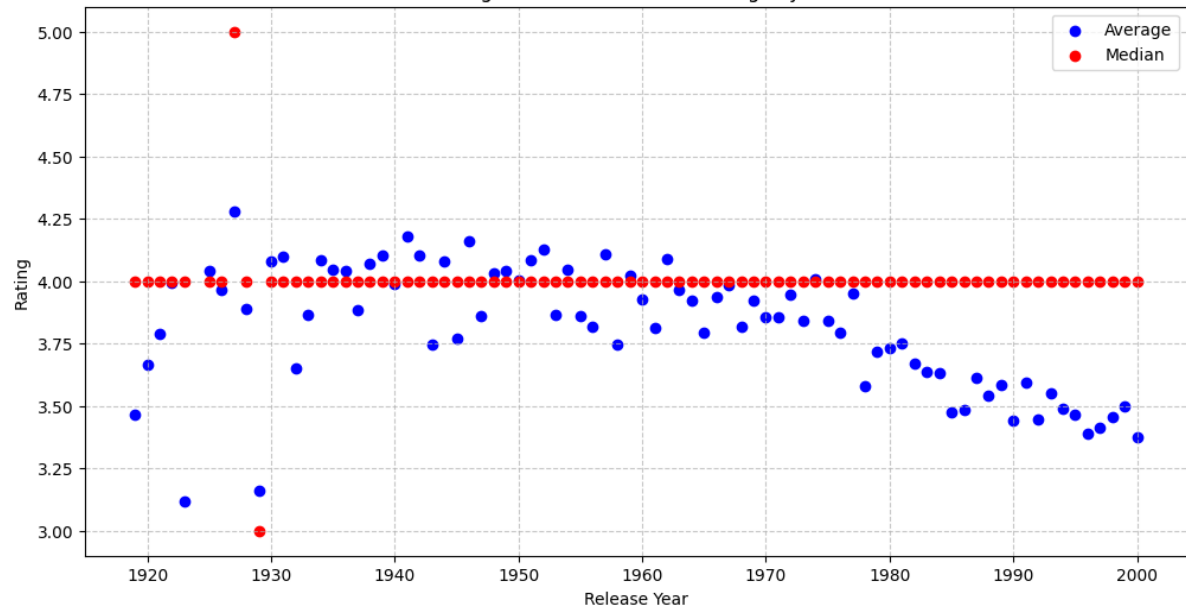
Boxplot of Median User Ratings for Movies



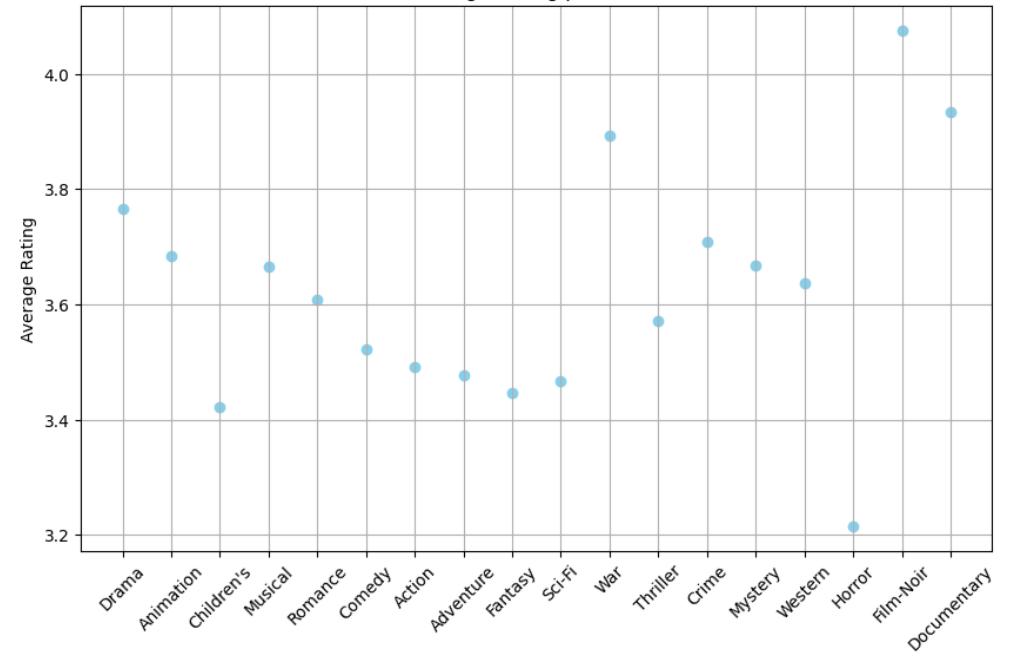
Boxplot of Number of Ratings per User



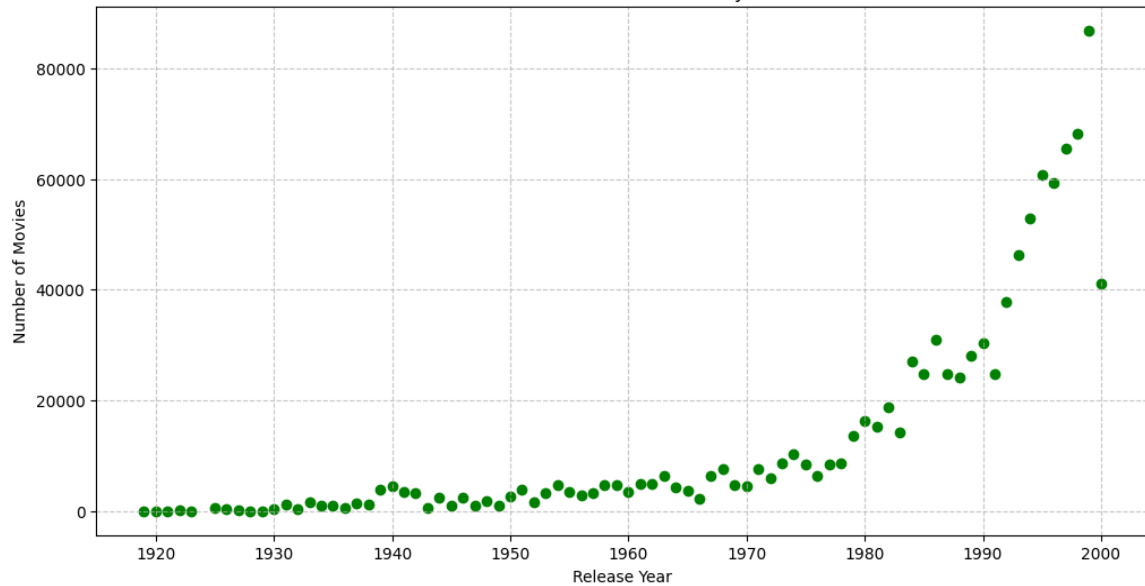
Average and Median Movie Ratings by Year



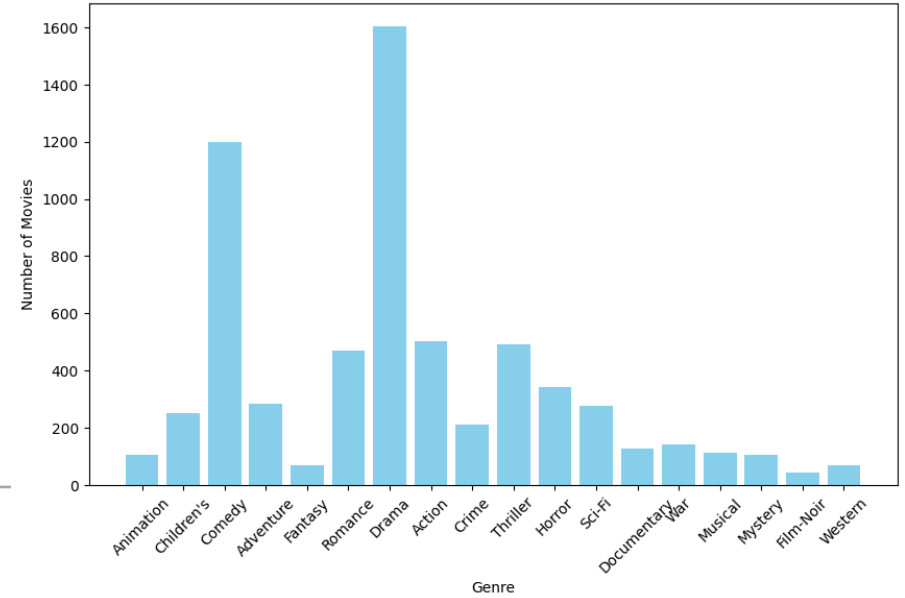
Average Rating per Genre

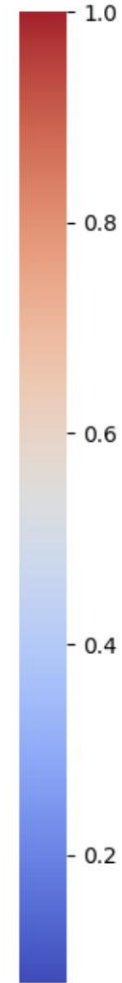
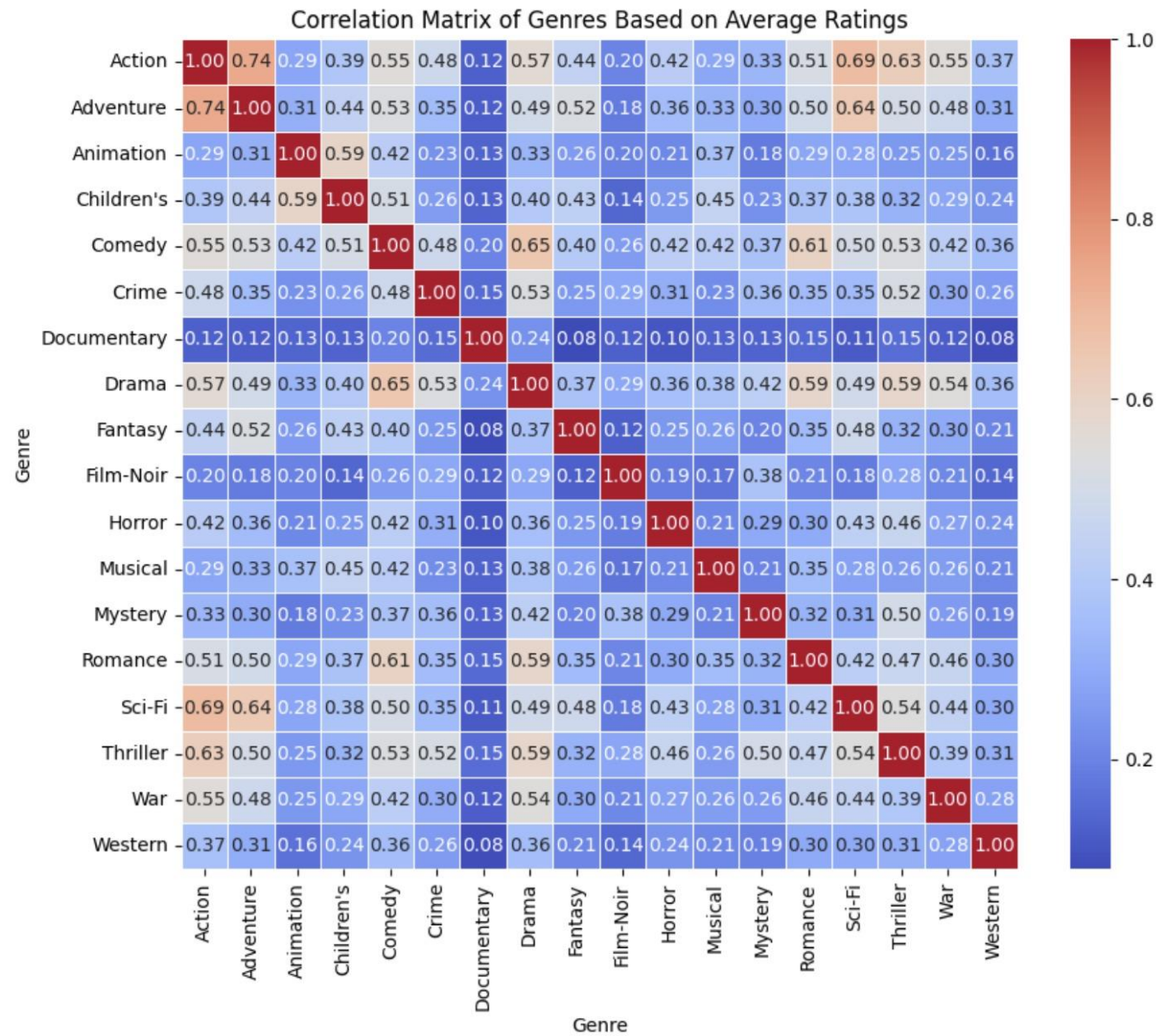


Number of Movies Released by Year



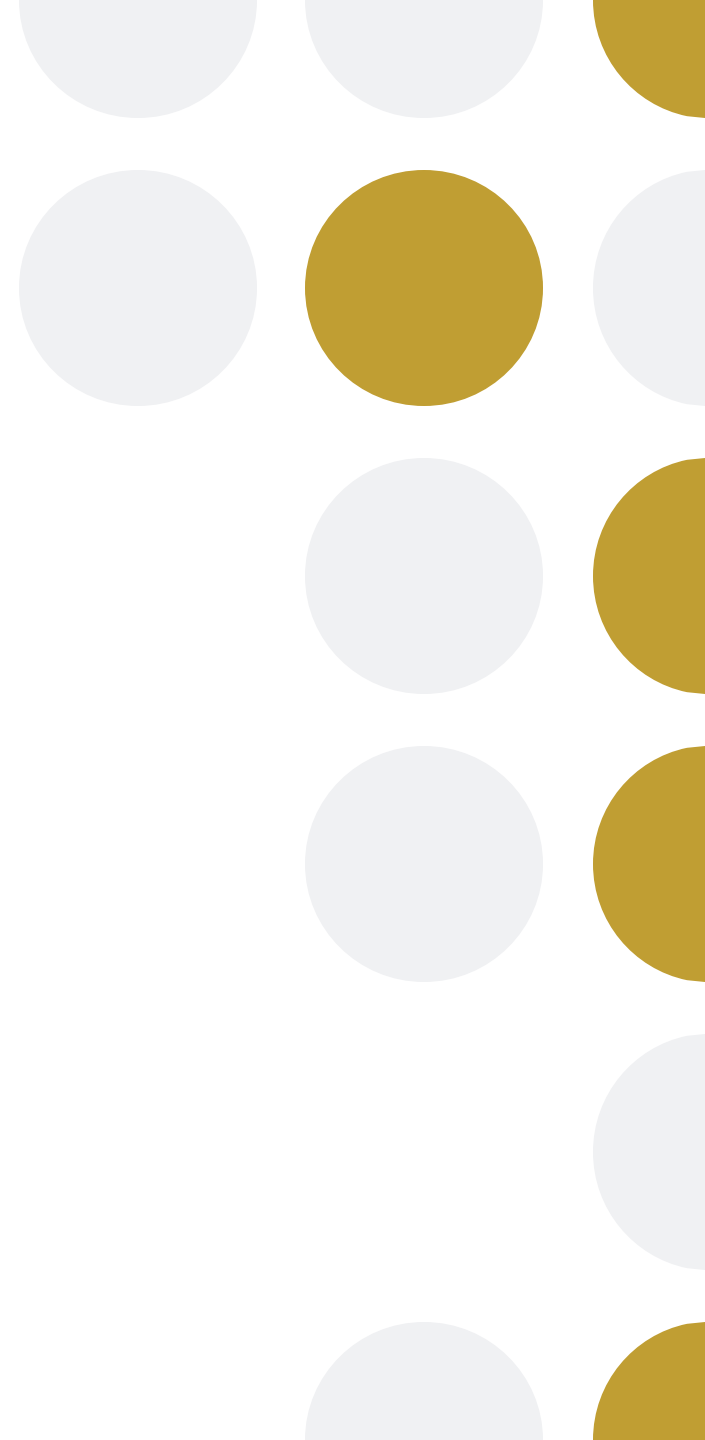
Number of Movies per Genre





Preprocessing

1. User Bias:
 - Calculated user bias as the mean of all their ratings. This is necessary to consider in a user-user setting.
2. Defined function to create the movie-user interaction matrix



Our Approach

1. Collaborative Filtering

- User-User based using Pearson's correlation
 - Unweighted correlation: Vanilla approach with only similarity index and user bias.
 - Weighted correlation: Added weights to movies based on their frequency.
- Item-Item based, using K-Nearest Neighbours (KNN)

2. SVD followed by KMeans++

Naive Approach: User-User Collaborative Filtering using Pearson's correlation

1. Find valid user pairs:

- Group users by those that have watched **at least one** movie the target user has watched.
- Give priority to users that have more movies in common (by sorting based on size of group)

2. Compute Similarity

- Calculate the Pearson correlation coefficient between each pair of users based on their ratings.
- Return the similarity index for each valid user-user pair

3. Find Closest users

- Select top 50 users that are closest to target user.

4. Predict Ratings

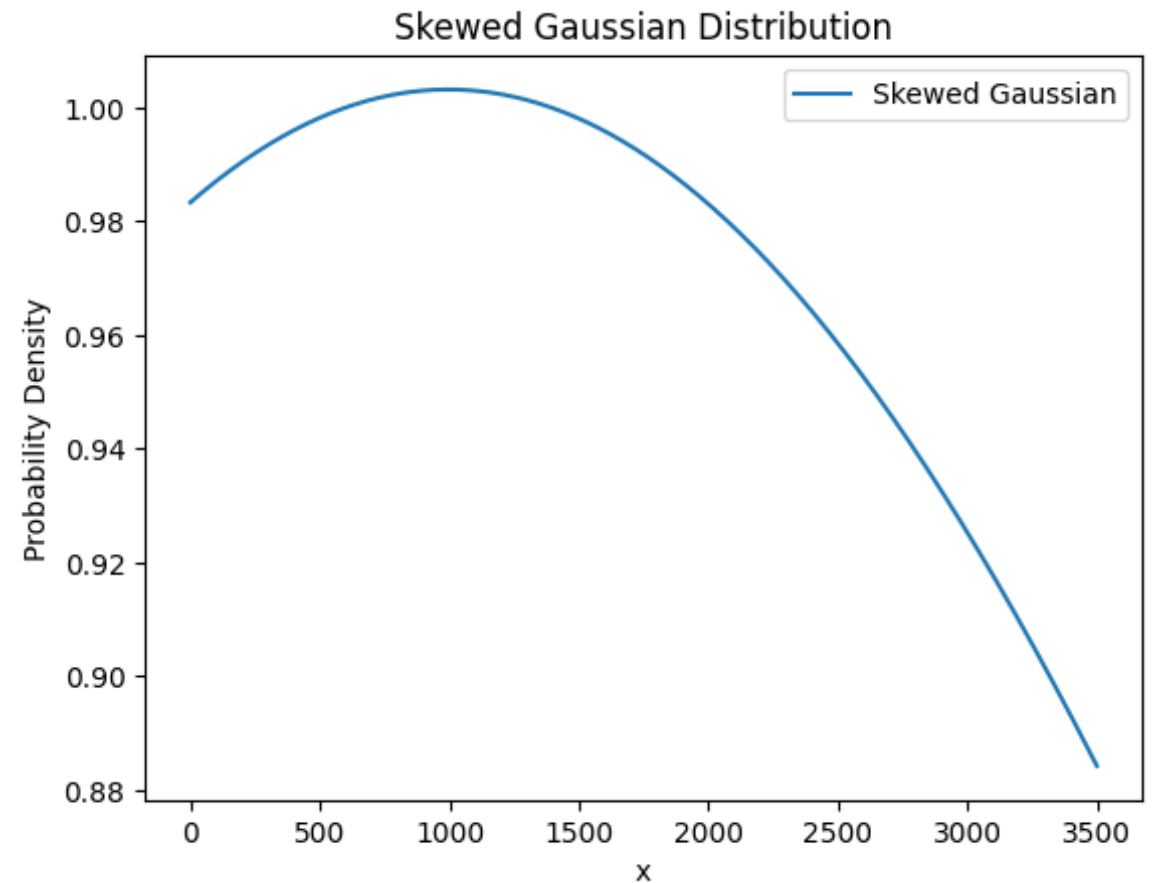
- Weight each neighbor's rating by their correlation coefficient with the target user as per formula.
- Normalize the ratings by subtracting the user's mean rating to handle biases.

5. Final Recommendations

- Return top k movies based on reverse sort of predicted ratings
-

Adding **Novelty**: User-User Collaborative Filtering using **Weighted** Pearson's correlation

- In class we were introduced to the concept of giving higher weights to movies that were rated by fewer people similar to the Inverse Document Frequency (IDF) approach used in NLP.
- Instead of a IDF graph that is strictly decreasing we make use of a skewed gaussian with maximum value near 1000.
- Our hypothesis is that extremely low frequency movies don't say much about a user's taste (for eg. The users that have rated them might have done so because their friends were involved in the production)
- The skewed gaussian preserves the key ideals of the IDF approach while also taking into account, the possible vagueness of extremely low frequency movies.



Comparison (Unweighted vs Weighted)

Input:

*New target user with 11 movies with main genres as **Drama, Mystery and Romance**. Average rating of user was 4.3. We used both methods to compare the top 10 recommended movies.*

Naïve Pearson's Correlation

- Along with the same 2-3 genres we got few new genres as well like Comedy, Horror.
- This seemed out of place compared to the input genres

Weighted Pearson's Correlation

- Got related 2-3 genres of movies as per target user data along with 1-2 Comedy movies.
- Broadly Matched with overall target user profile.

Given the results, we decided to move ahead with the weighted approach as final recommendations for user-user method.

Item-Item Collaborative Filtering using KNN

1. Calculate similarity between movies using KNN.

- In KNN, we use cosine similarity to find the closest set of movies to each movie by calculating the distances of all other movies and taking the nearest K.

2. Predict rating for unwatched movies

- Find the closest movies based on similarity from previous step and calculate the weighted average of the ratings.
- The weight used will be inverse of the distance of how close one movie is to the other.

3. Recommend movies which have highest predicted rating.

Comparision (User-User vs Item-Item)

Input: New target user with 11 movies with main genres as Drama, Mystery and Romance. Avg rating was 4.3. We used both methods to compare the top 10 recommended movies.

User-User Weighted

- Along with the same 2-3 genres we got few new genres as well like Comedy, Sci-Fi.
- Explainable and overall decent results, but lacked some aspect of Accuracy.

Item-Item KNN

- Got related 2-3 genres of movies as per target user data along with 1-2 Comedy movies.
- More accurate and desired results. For eg. On input of Toy Story, Godfather, recommendations contained Toy Story 2 and Godfather III.

Item-Item gave overall better predictions while user-user captured more novelty

SVD + K-Means++

1. Formed a User-Genre Rating Matrix:

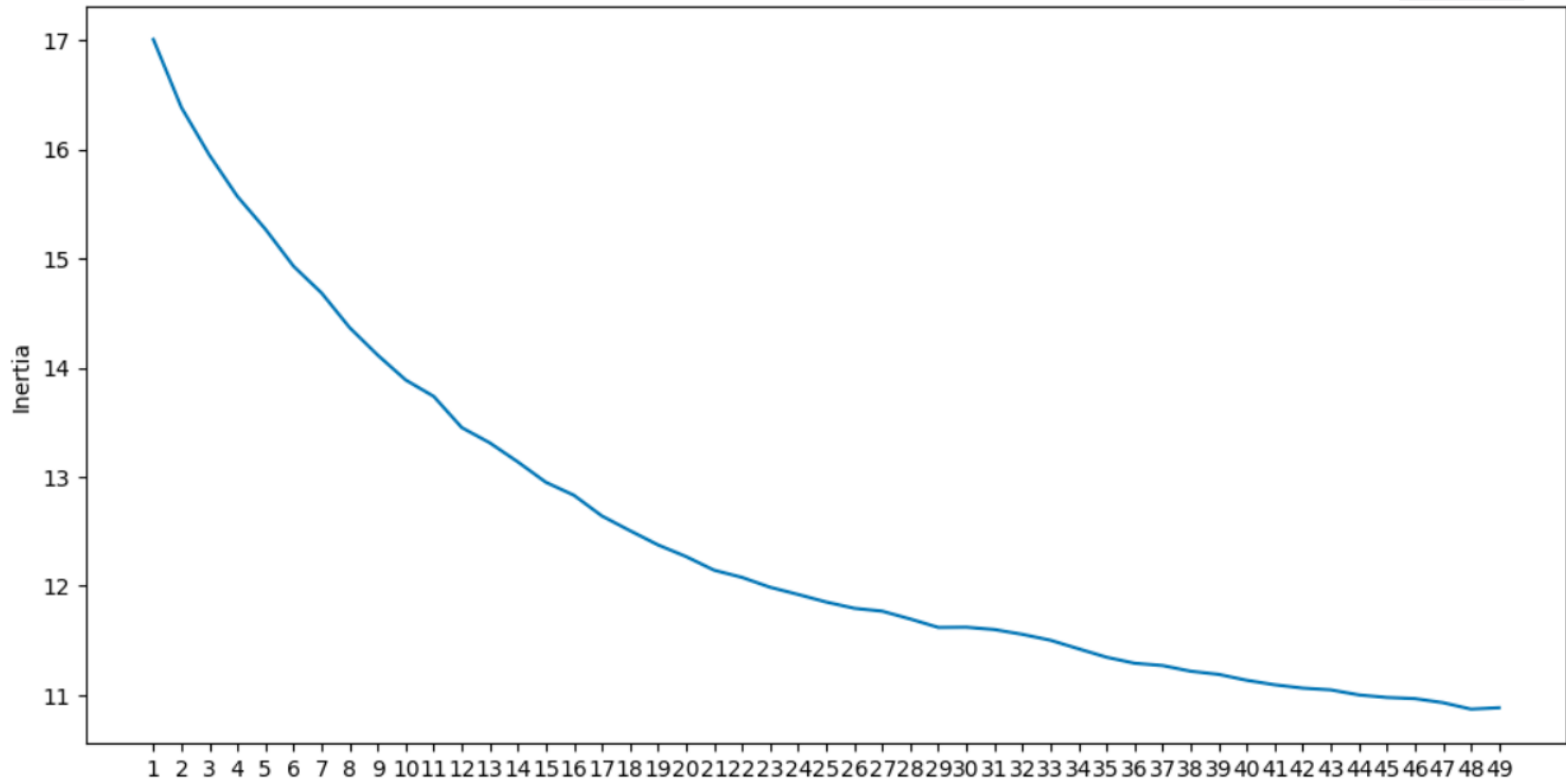
- The ratings for each user in a specific genre is computed by averaging the ratings of the movies that the user has watched within that genre
- This process is performed for all users across all genres, and the results are stored in a pandas dataframe.
- Null Values are filled by the average of the column (which is average of the genre)

2. Computing SVD of the User-Genre Matrix:

- SVD of the user-genre matrix is calculated using the power iteration method. This involves 2 steps:
- First, the covariance matrix ($X^T.X = V.\Sigma^2.V^T$) is used to compute the V matrix and Σ matrix (where $X = U\Sigma V^T$)
- Second, the equation $XV=U\Sigma$ is used to find the U matrix.
- After performing SVD on the user-genre matrix, we get the U, Σ and the V matrix.

3. Performing K-Means++ on the 'U' matrix:

- Now we perform K-Means++ on the U matrix, which consists of user representations.
 - Inertia vs K Plot: (Next Slide)
-



SVD + K-Means++

3. Performing K-Means++ on the 'U' Matrix:

- The Number of clusters (k) is chosen to be 16 while clustering the users.

4. Predict Ratings:

- After a new user has supplied ratings for several movies, a suitable user representation is constructed based on the U matrix. This user representation is subsequently inputted into the k-means model to predict the cluster to which the new user is assigned.
 - Now, the ratings of this new user are predicted by passing the user's ratings through the weighted collaborative filtering, with the added constraint that similarity is only calculated with the users that are in the same cluster as the new user.
 - To compute the ratings, only the top-30 most similar users were considered.
 - Recommend movies based on the highest value of predicted ratings.
-

Comparision (User-User vs SVD + K-Means++)

Input: New target user with 11 movies with main genres as Drama, Mystery and Romance. Avg rating was 4.3. We used both methods to compare the top 10 recommended movies.

User-User Weighted

- Along with the same 2-3 genres we got few new genres as well like Comedy, Sci-Fi.
- Explainable and overall decent results, but lacked some aspect of Accuracy.

SVD+K-Means++

- Got movies with related genres as the input movies as well as children's genre and comedy genre
- Got comparable accuracy in terms of average ratings of movies that were suggested

SVD+K-Means++ was comparable in terms of accuracy and recommendations with the user-user collaborative filtering approach

Thank You

