# MOVIE LENS RECOMMENDATION SYSTEM

## Business Understanding

## Problem Statement

In today's world where consumers face an overwhelming abundance of choices, recommender systems are no longer optional but essential for businesses to thrive and maintain a competitive edge in a data driven marketplace.

Recommender systems provide personalized guidance, thereby enhancing decision making and user satisfaction by predicting user preferences and surfacing relevant items such as movies, books, music and products from vast inventories through user past behaviour, preference or a combination of both.

Recommendation systems are pivotal in navigating the "Long Tail"—a phenomenon where niche, less-popular items collectively represent significant value but remain undiscovered without intelligent filtering.

By analyzing user behavior and leveraging predictive models, these systems drive customer engagement, loyalty, and profitability. For businesses, they unlock opportunities to monetize the Long Tail, improve retention, and deliver personalized experiences.

**The Key Impact of recommender systems include :**

```
        1.   User Experience: Simplify discovery in a sea of options, ensuring tai
lored recommendations.
        2.   Business Value: Increase sales, reduce choice overload, and capitaliz
e on niche markets.
        3.   Competitive Edge: Transform vast inventories into curated, actionable
insights for users.
```

## Overview

The goal of this project is to build a recommendation model that provides the top 5 personalized movie recommendations to a user based on their past ratings, preferences. By analyzing user preferences and comparing them with patterns in the ratings of other users or movies, the system aims to predict what a user would enjoy watching next.

Primarily, this project will recommend movies using Collaborative Filtering, and hybrid approach which will combine both content based and collaborative filtering to handle cold start cases

## Key Features

1. Dataset: MovieLens latest small having 100,000 ratings
2. Core method: Collaborative filtering approach
3. Evaluation Metrics:
4. Stretch goal: Hybrid approach for cold start issues

## Reasons for Chosing Project

1. **Real World relevance:-** Recommendation systems drive engagement in Platforms such as Amazon,

Instagram, Reddit, Netflix and Google
2. **Academic Benchmark:-** MovieLens is a standard dataset for testing recommendation algorithms from the GroupLens research lab at the University of Minnesota
3. **Is a scalable solution capable of extending with Neural networks or hybrid models**

# Objectives

1. Design a collaborative filtering–based recommendation system to accurately predict user preferences.
2. Mitigate the cold start problem for new users and movies with limited or no ratings.
3. Develop classification models to distinguish between highly liked and disliked movies.
4. Recommend top 5 movies based on user and particular movie watched.

# Success criteria

1. **Model Performance:** Achieve a high precision@K > 0.6 for accurate recommndation to the users
2. **User Satisfaction:** Recommendations should align with user preferences as evaluated via user feedback or online metrics
3. **The system should be scalable to handle new users or movies efficiently**

# Data understanding

**The data seems to have been split between train and test data**

In [ ]:

```python
#import libraries to support in data understanding and cleaning
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns
import warnings
from sklearn.preprocessing import MultiLabelBinarizer, MinMaxScaler, LabelEncoder
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.metrics.pairwise import cosine_similarity
from sklearn.decomposition import TruncatedSVD
from sklearn.model_selection import train_test_split
from math import sqrt
from sklearn.metrics import mean_squared_error
from sklearn.ensemble import RandomForestClassifier
from sklearn.metrics import precision_score, recall_score
from sklearn.tree import DecisionTreeClassifier
from sklearn.ensemble import GradientBoostingClassifier
import tensorflow as tf
from sklearn.preprocessing import StandardScaler
from sklearn.utils import class_weight
from sklearn.metrics import precision_score, recall_score, classification_report


warnings.filterwarnings('ignore')
```

In [ ]:

```python
# load the files
movies= pd.read_csv('movies.csv')
ratings=pd.read_csv('ratings.csv')
tags=pd.read_csv('tags.csv')
```

In [ ]:

```python
# merge the csv files and rename user and timestamp columns as per table
```

```
df_merge = movies.merge(ratings, on="movieId")\
.merge(tags, on=("movieId"))
df_merge
```

Out[ ]:

| | movieId | title | genres | userId_x | rating | timestamp_x | userId_y | tag | ti |
|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | Toy Story (1995) | Adventure\|Animation\|Children\|Comedy\|Fantasy | 1 | 4.0 | 964982703 | 336 | pixar | |
| 1 | 1 | Toy Story (1995) | Adventure\|Animation\|Children\|Comedy\|Fantasy | 1 | 4.0 | 964982703 | 474 | pixar | |
| 2 | 1 | Toy Story (1995) | Adventure\|Animation\|Children\|Comedy\|Fantasy | 1 | 4.0 | 964982703 | 567 | fun | |
| 3 | 1 | Toy Story (1995) | Adventure\|Animation\|Children\|Comedy\|Fantasy | 5 | 4.0 | 847434962 | 336 | pixar | |
| 4 | 1 | Toy Story (1995) | Adventure\|Animation\|Children\|Comedy\|Fantasy | 5 | 4.0 | 847434962 | 474 | pixar | |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 233208 | 187595 | Solo: A Star Wars Story (2018) | Action\|Adventure\|Children\|Sci-Fi | 586 | 5.0 | 1529899556 | 62 | star wars | |
| 233209 | 193565 | Gintama: The Movie (2010) | Action\|Animation\|Comedy\|Sci-Fi | 184 | 3.5 | 1537098554 | 184 | anime | |
| 233210 | 193565 | Gintama: The Movie (2010) | Action\|Animation\|Comedy\|Sci-Fi | 184 | 3.5 | 1537098554 | 184 | comedy | |
| 233211 | 193565 | Gintama: The Movie (2010) | Action\|Animation\|Comedy\|Sci-Fi | 184 | 3.5 | 1537098554 | 184 | gintama | |
| 233212 | 193565 | Gintama: The Movie (2010) | Action\|Animation\|Comedy\|Sci-Fi | 184 | 3.5 | 1537098554 | 184 | remaster | |

**233213 rows × 9 columns**

In [ ]:

```
#renaming merged columns
df_merge = df_merge.rename(columns={
    'userId_x': 'userId_rating',
    'timestamp_x': 'timestamp_rating',
    'userId_y': 'userId_tag',
    'timestamp_y': 'timestamp_tag'
})
```

In [ ]:

```
# checking shape of merged dataframe
df_merge.shape
```

Out[ ]:

```
(233213, 9)
```

```
In [ ]:
```

```
# checking all necessary info related to the dataframe
df_merge.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 233213 entries, 0 to 233212
Data columns (total 9 columns):
 #   Column            Non-Null Count   Dtype
---  ------            --------------   -----
 0   movieId           233213 non-null  int64
 1   title             233213 non-null  object
 2   genres            233213 non-null  object
 3   userId_rating     233213 non-null  int64
 4   rating            233213 non-null  float64
 5   timestamp_rating  233213 non-null  int64
 6   userId_tag        233213 non-null  int64
 7   tag               233213 non-null  object
 8   timestamp_tag     233213 non-null  int64
dtypes: float64(1), int64(5), object(3)
memory usage: 16.0+ MB
```

```
In [ ]:
```

```
# summary statistics
df_merge.describe()
```

```
Out[ ]:
```

|  | movieId | userId_rating | rating | timestamp_rating | userId_tag | timestamp_tag |
|---|---|---|---|---|---|---|
| **count** | 233213.000000 | 233213.000000 | 233213.000000 | 2.332130e+05 | 233213.000000 | 2.332130e+05 |
| **mean** | 12319.999443 | 309.688191 | 3.966535 | 1.213524e+09 | 470.683564 | 1.384774e+09 |
| **std** | 28243.919401 | 178.206387 | 0.968637 | 2.250448e+08 | 153.329632 | 1.534621e+08 |
| **min** | 1.000000 | 1.000000 | 0.500000 | 8.281246e+08 | 2.000000 | 1.137179e+09 |
| **25%** | 296.000000 | 156.000000 | 3.500000 | 1.017365e+09 | 424.000000 | 1.242494e+09 |
| **50%** | 1198.000000 | 309.000000 | 4.000000 | 1.217325e+09 | 477.000000 | 1.457901e+09 |
| **75%** | 4638.000000 | 460.000000 | 5.000000 | 1.443201e+09 | 599.000000 | 1.498457e+09 |
| **max** | 193565.000000 | 610.000000 | 5.000000 | 1.537799e+09 | 610.000000 | 1.537099e+09 |

# Data Cleaning

**Handling null values**

**Dropping duplicates**

**Display and Clean Outliers**

```
In [ ]:
```

```
# checking for missing values
df_merge.isna().sum()
```

```
Out[ ]:
```

|  | 0 |
|---|---|
| **movieId** | 0 |
| **title** | 0 |
| **genres** | 0 |
| **userId_rating** | 0 |

| | rating | 0 |
|---|---|---|
| | | 0 |
| timestamp_rating | 0 |
| userId_tag | 0 |
| tag | 0 |
| timestamp_tag | 0 |

**dtype: int64**

In [ ]:

```python
# checking duplicates
df_merge.duplicated().sum()
```

Out[ ]:

```
np.int64(0)
```

In [ ]:

```python
# checking for outliers in the dataset
Q1 = df_merge['rating'].quantile(0.25)
Q3 = df_merge['rating'].quantile(0.75)
IQR = Q3 - Q1

lower_bound = Q1 - 1.5 * IQR
upper_bound = Q3 + 1.5 * IQR

# Identifying outliers
outliers = df_merge[(df_merge['rating'] < lower_bound) | (df_merge['rating'] > upper_boun
d)]
outliers
```

Out[ ]:

| | movieId | title | genres | userId_rating | rating | timestamp_rating | userId_tag |
|---|---|---|---|---|---|---|---|
| 78 | 1 | Toy Story (1995) | Adventure\|Animation\|Children\|Comedy\|Fantasy | 76 | 0.5 | 1439165548 | 336 |
| 79 | 1 | Toy Story (1995) | Adventure\|Animation\|Children\|Comedy\|Fantasy | 76 | 0.5 | 1439165548 | 474 |
| 80 | 1 | Toy Story (1995) | Adventure\|Animation\|Children\|Comedy\|Fantasy | 76 | 0.5 | 1439165548 | 567 |
| 741 | 2 | Jumanji (1995) | Adventure\|Children\|Fantasy | 149 | 1.0 | 902084874 | 62 |
| 742 | 2 | Jumanji (1995) | Adventure\|Children\|Fantasy | 149 | 1.0 | 902084874 | 62 |
| ... | ... | ... | ... | ... | ... | ... | ... |
| 233155 | 184471 | Tomb Raider (2018) | Action\|Adventure\|Fantasy | 153 | 0.5 | 1525553051 | 62 |
| 233156 | 184471 | Tomb Raider (2018) | Action\|Adventure\|Fantasy | 153 | 0.5 | 1525553051 | 62 |
| 233187 | 187593 | Deadpool 2 (2018) | Action\|Comedy\|Sci-Fi | 338 | 1.0 | 1530148465 | 62 |
| 233188 | 187593 | Deadpool 2 (2018) | Action\|Comedy\|Sci-Fi | 338 | 1.0 | 1530148465 | 62 |
| 233189 | 187593 | Deadpool 2 (2018) | Action\|Comedy\|Sci-Fi | 338 | 1.0 | 1530148465 | 62 |

**5629 rows × 9 columns**

```
In [ ]:
```

```
# plot outliers
plt.figure(figsize=(10, 6))
sns.boxplot(x=df_merge['rating'])
plt.title('Box Plot of Rating')
plt.show()
```



Box Plot of Rating
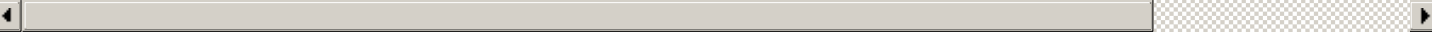
```
In [ ]:
```

```
# removing outliers
df_merge_cleaned = df_merge[(df_merge['rating'] >= lower_bound) & (df_merge['rating'] <=
upper_bound)]
df_merge_cleaned
```

```
Out[ ]:
```

| | movieId | title | genres | userId_rating | rating | timestamp_rating | userId_tag |
|---|---|---|---|---|---|---|---|
| 0 | 1 | Toy Story (1995) | Adventure\|Animation\|Children\|Comedy\|Fantasy | 1 | 4.0 | 964982703 | 336 |
| 1 | 1 | Toy Story (1995) | Adventure\|Animation\|Children\|Comedy\|Fantasy | 1 | 4.0 | 964982703 | 474 |
| 2 | 1 | Toy Story (1995) | Adventure\|Animation\|Children\|Comedy\|Fantasy | 1 | 4.0 | 964982703 | 567 |
| 3 | 1 | Toy Story (1995) | Adventure\|Animation\|Children\|Comedy\|Fantasy | 5 | 4.0 | 847434962 | 336 |
| 4 | 1 | Toy Story (1995) | Adventure\|Animation\|Children\|Comedy\|Fantasy | 5 | 4.0 | 847434962 | 474 |
| ... | ... | ... | | ... | ... | ... | ... |
| | | Solo: A | | | | | |

| | movieId | title | genres | userId_rating | rating | timestamp_rating | userId_tag |
|---|---|---|---|---|---|---|---|
| 233208 | 187595 | Solo: A Star Wars Story (2018) | Action\|Adventure\|Children\|Sci-Fi | 586 | 5.0 | 1529899556 | 62 |
| 233209 | 193565 | Gintama: The Movie (2010) | Action\|Animation\|Comedy\|Sci-Fi | 184 | 3.5 | 1537098554 | 184 |
| 233210 | 193565 | Gintama: The Movie (2010) | Action\|Animation\|Comedy\|Sci-Fi | 184 | 3.5 | 1537098554 | 184 |
| 233211 | 193565 | Gintama: The Movie (2010) | Action\|Animation\|Comedy\|Sci-Fi | 184 | 3.5 | 1537098554 | 184 |
| 233212 | 193565 | Gintama: The Movie (2010) | Action\|Animation\|Comedy\|Sci-Fi | 184 | 3.5 | 1537098554 | 184 |

**227584 rows × 9 columns**

**No major effect on dropping outliers only approx 5000 rows of 220000 were removed**

In [ ]:

```python
# plotting the cleaned dataset
plt.figure(figsize=(10, 6))
sns.boxplot(x=df_merge_cleaned['rating'])
plt.title('Box Plot of Rating')
plt.show()
```

### Box Plot of Rating



In [ ]:

```
# converting timestamp to datetime format
df_merge_cleaned['timestamp_rating'] = pd.to_datetime(df_merge_cleaned['timestamp_rating'], unit='s', errors='coerce')
df_merge_cleaned['timestamp_tag'] = pd.to_datetime(df_merge_cleaned['timestamp_tag'], unit='s', errors='coerce')
df_merge_cleaned
```

Out[ ]:

| | movieId | title | genres | userId_rating | rating | timestamp_rating | userId_tag |
|---|---|---|---|---|---|---|---|
| 0 | 1 | Toy Story (1995) | Adventure\|Animation\|Children\|Comedy\|Fantasy | 1 | 4.0 | 2000-07-30 18:45:03 | 336 |
| 1 | 1 | Toy Story (1995) | Adventure\|Animation\|Children\|Comedy\|Fantasy | 1 | 4.0 | 2000-07-30 18:45:03 | 474 |
| 2 | 1 | Toy Story (1995) | Adventure\|Animation\|Children\|Comedy\|Fantasy | 1 | 4.0 | 2000-07-30 18:45:03 | 567 |
| 3 | 1 | Toy Story (1995) | Adventure\|Animation\|Children\|Comedy\|Fantasy | 5 | 4.0 | 1996-11-08 06:36:02 | 336 |
| 4 | 1 | Toy Story (1995) | Adventure\|Animation\|Children\|Comedy\|Fantasy | 5 | 4.0 | 1996-11-08 06:36:02 | 474 |
| ... | ... | ... | ... | ... | ... | ... | ... |
| 233208 | 187595 | Solo: A Star Wars Story (2018) | Action\|Adventure\|Children\|Sci-Fi | 586 | 5.0 | 2018-06-25 04:05:56 | 62 |
| 233209 | 193565 | Gintama: The Movie (2010) | Action\|Animation\|Comedy\|Sci-Fi | 184 | 3.5 | 2018-09-16 11:49:14 | 184 |
| 233210 | 193565 | Gintama: The Movie (2010) | Action\|Animation\|Comedy\|Sci-Fi | 184 | 3.5 | 2018-09-16 11:49:14 | 184 |
| 233211 | 193565 | Gintama: The Movie (2010) | Action\|Animation\|Comedy\|Sci-Fi | 184 | 3.5 | 2018-09-16 11:49:14 | 184 |
| 233212 | 193565 | Gintama: The Movie (2010) | Action\|Animation\|Comedy\|Sci-Fi | 184 | 3.5 | 2018-09-16 11:49:14 | 184 |

**227584 rows × 9 columns**

In [ ]:

```
# checking total count per rating
df_merge['rating'].value_counts()
```

Out[ ]:

| | count |
|---|---|
| **rating** | |
| 4.0 | 64781 |
| 5.0 | 63845 |
| 4.5 | 31502 |

| rating | count |
|---|---|
| 3.0 | 28550 |
| 2.5 | 22895 |
| 2.0 | 7955 |
| 2.5 | 6488 |
| 1.0 | 3721 |
| 0.5 | 1908 |
| 1.5 | 1568 |

**dtype: int64**

In [ ]:

```python
# creating new column from title by removing year of release from movie name
df_merge_cleaned['year'] = df_merge_cleaned['title'].str.extract(r'\((\d{4})\)')
df_merge_cleaned['title'] = df_merge_cleaned['title'].str.replace(r'\s*\(\d{4}\)$', '',
regex=True)
df_merge_cleaned
```

Out[ ]:

| | movieId | title | genres | userId_rating | rating | timestamp_rating | userId_tag |
|---|---|---|---|---|---|---|---|
| 0 | 1 | Toy Story | Adventure\|Animation\|Children\|Comedy\|Fantasy | 1 | 4.0 | 2000-07-30 18:45:03 | 336 |
| 1 | 1 | Toy Story | Adventure\|Animation\|Children\|Comedy\|Fantasy | 1 | 4.0 | 2000-07-30 18:45:03 | 474 |
| 2 | 1 | Toy Story | Adventure\|Animation\|Children\|Comedy\|Fantasy | 1 | 4.0 | 2000-07-30 18:45:03 | 567 |
| 3 | 1 | Toy Story | Adventure\|Animation\|Children\|Comedy\|Fantasy | 5 | 4.0 | 1996-11-08 06:36:02 | 336 |
| 4 | 1 | Toy Story | Adventure\|Animation\|Children\|Comedy\|Fantasy | 5 | 4.0 | 1996-11-08 06:36:02 | 474 |
| ... | ... | ... | ... | ... | ... | ... | ... |
| 233208 | 187595 | Solo: A Star Wars Story | Action\|Adventure\|Children\|Sci-Fi | 586 | 5.0 | 2018-06-25 04:05:56 | 62 |
| 233209 | 193565 | Gintama: The Movie | Action\|Animation\|Comedy\|Sci-Fi | 184 | 3.5 | 2018-09-16 11:49:14 | 184 |
| 233210 | 193565 | Gintama: The Movie | Action\|Animation\|Comedy\|Sci-Fi | 184 | 3.5 | 2018-09-16 11:49:14 | 184 |
| 233211 | 193565 | Gintama: The Movie | Action\|Animation\|Comedy\|Sci-Fi | 184 | 3.5 | 2018-09-16 11:49:14 | 184 |
| 233212 | 193565 | Gintama: The Movie | Action\|Animation\|Comedy\|Sci-Fi | 184 | 3.5 | 2018-09-16 11:49:14 | 184 |

**227584 rows × 10 columns**

In [ ]:

```python
# Splitting the 'genres' column by '|', then explode
df_exploded = df_merge_cleaned.copy()
df_exploded['genres'] = df_exploded['genres'].str.split('|')
df_exploded = df_exploded.explode('genres').reset_index(drop=True)
df_exploded
```

Out[ ]:

|  | movieId | title | genres | userId_rating | rating | timestamp_rating | userId_tag | tag | timestamp_tag | year |
|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | Toy Story | Adventure | 1 | 4.0 | 2000-07-30 18:45:03 | 336 | pixar | 2006-02-04 09:36:04 | 1995 |
| 1 | 1 | Toy Story | Animation | 1 | 4.0 | 2000-07-30 18:45:03 | 336 | pixar | 2006-02-04 09:36:04 | 1995 |
| 2 | 1 | Toy Story | Children | 1 | 4.0 | 2000-07-30 18:45:03 | 336 | pixar | 2006-02-04 09:36:04 | 1995 |
| 3 | 1 | Toy Story | Comedy | 1 | 4.0 | 2000-07-30 18:45:03 | 336 | pixar | 2006-02-04 09:36:04 | 1995 |
| 4 | 1 | Toy Story | Fantasy | 1 | 4.0 | 2000-07-30 18:45:03 | 336 | pixar | 2006-02-04 09:36:04 | 1995 |
| ... | ... | ... | ... | ... | ... | ... | ... | ... | ... | ... |
| 747116 | 193565 | Gintama: The Movie | Sci-Fi | 184 | 3.5 | 2018-09-16 11:49:14 | 184 | gintama | 2018-09-16 11:50:03 | 2010 |
| 747117 | 193565 | Gintama: The Movie | Action | 184 | 3.5 | 2018-09-16 11:49:14 | 184 | remaster | 2018-09-16 11:49:52 | 2010 |
| 747118 | 193565 | Gintama: The Movie | Animation | 184 | 3.5 | 2018-09-16 11:49:14 | 184 | remaster | 2018-09-16 11:49:52 | 2010 |
| 747119 | 193565 | Gintama: The Movie | Comedy | 184 | 3.5 | 2018-09-16 11:49:14 | 184 | remaster | 2018-09-16 11:49:52 | 2010 |
| 747120 | 193565 | Gintama: The Movie | Sci-Fi | 184 | 3.5 | 2018-09-16 11:49:14 | 184 | remaster | 2018-09-16 11:49:52 | 2010 |

**747121 rows × 10 columns**

# Exploratory Data Analysis

## univariate analysis

In [ ]:

```
#identifying the distribution of movie ratings
sns.histplot(ratings['rating'], bins=9, kde=True)
plt.title("Distribution of Movie Ratings")
plt.xlabel("Rating")
plt.ylabel("Frequency")
plt.show()
```
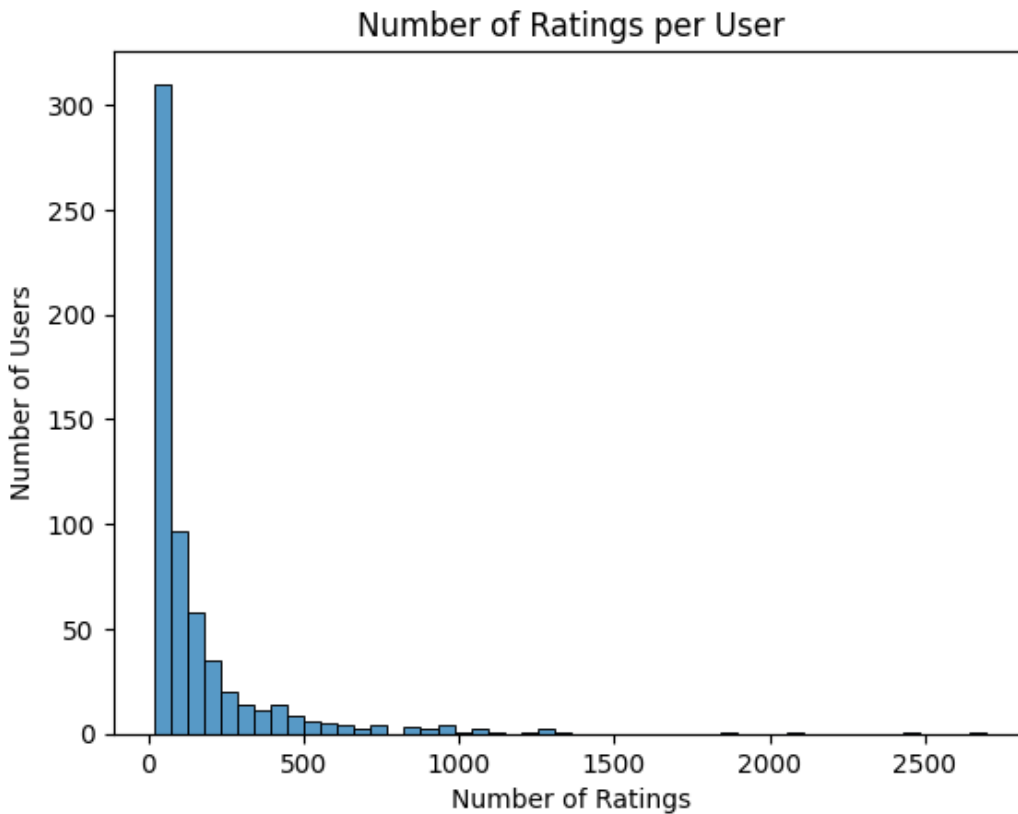


Distribution of Movie Ratings

Rating

In [ ]:

```python
#identifying number of ratings per user
user_rating_counts = ratings['userId'].value_counts()

sns.histplot(user_rating_counts, bins=50, kde=False)
plt.title("Number of Ratings per User")
plt.xlabel("Number of Ratings")
plt.ylabel("Number of Users")
plt.show()
```



In [ ]:

```python
# Removing users who had no ratings
user_rating_counts = df_merge_cleaned['userId_rating'].value_counts()


active_users = user_rating_counts[user_rating_counts > 0].index


ratings_filtered = df_merge_cleaned[df_merge_cleaned['userId_rating'].isin(active_users)
]


print(ratings_filtered.shape)
```

(227584, 10)

In [ ]:

```python
# Top 5 movies that had the most rating by count

movie_ratings_count = df_merge_cleaned.groupby('title')['rating'].count()
top_5_movies = movie_ratings_count.sort_values(ascending=False).head(5)

plt.figure(figsize=(10, 6))
top_5_movies.plot(kind='bar')
```

```
plt.title('Top 5 Movies with Most Ratings')
plt.xlabel('Movie Title')
plt.ylabel('Number of Ratings')
plt.xticks(rotation=45, ha='right')
plt.tight_layout()
plt.show()
```



Top 5 Movies with Most Ratings

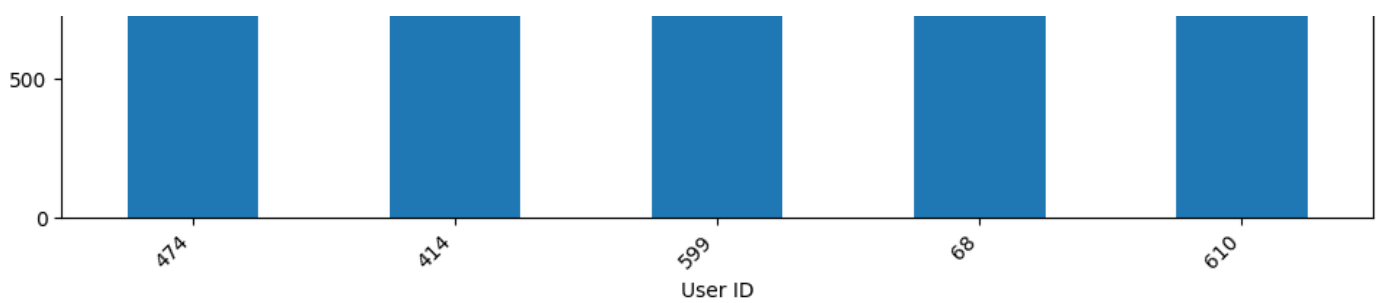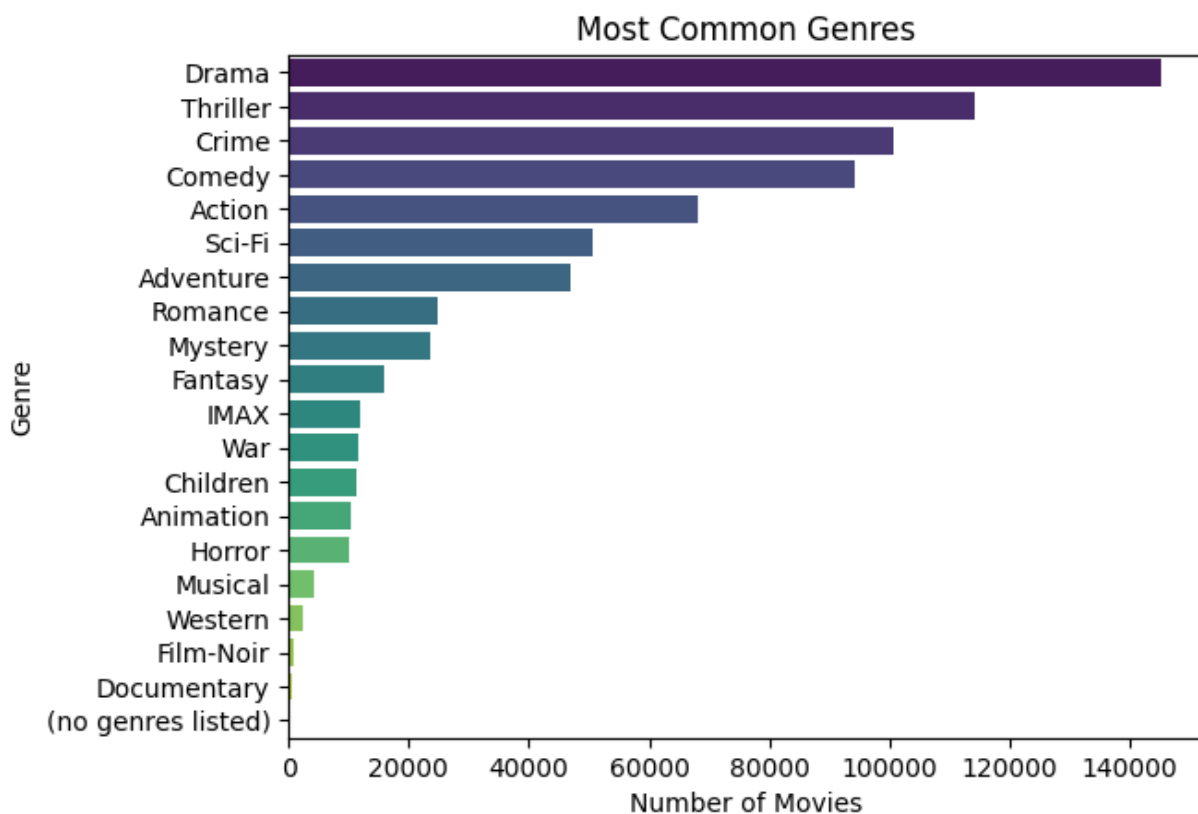**Pulp Fiction seems to be the highest rated movie. This would be a good movie to draw insights from**

In [ ]:

```
# Top 5 users who have rated the most movies by count

user_rating_counts = df_merge_cleaned.groupby('userId_rating')['rating'].count()
top_5_users = user_rating_counts.sort_values(ascending=False).head(5)

plt.figure(figsize=(10, 6))
top_5_users.plot(kind='bar')
plt.title('Top 5 Users with Most Ratings')
plt.xlabel('User ID')
plt.ylabel('Number of Ratings')
plt.xticks(rotation=45, ha='right')
plt.tight_layout()
plt.show()
```



Top 5 Users with Most Ratings

**Top five users with the most count of ratings will help us draw further insisghts. User 68 is a good selection as he/she is not the highest.**

In [ ]:

```python
#genre with the most number of movies from df_exploded
genre_counts = df_exploded['genres'].str.split('|').explode().value_counts()

sns.barplot(y=genre_counts.index, x=genre_counts.values, palette='viridis')
plt.title("Most Common Genres")
plt.xlabel("Number of Movies")
plt.ylabel("Genre")
plt.show()
```
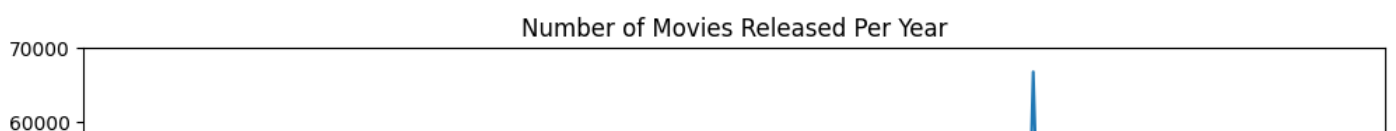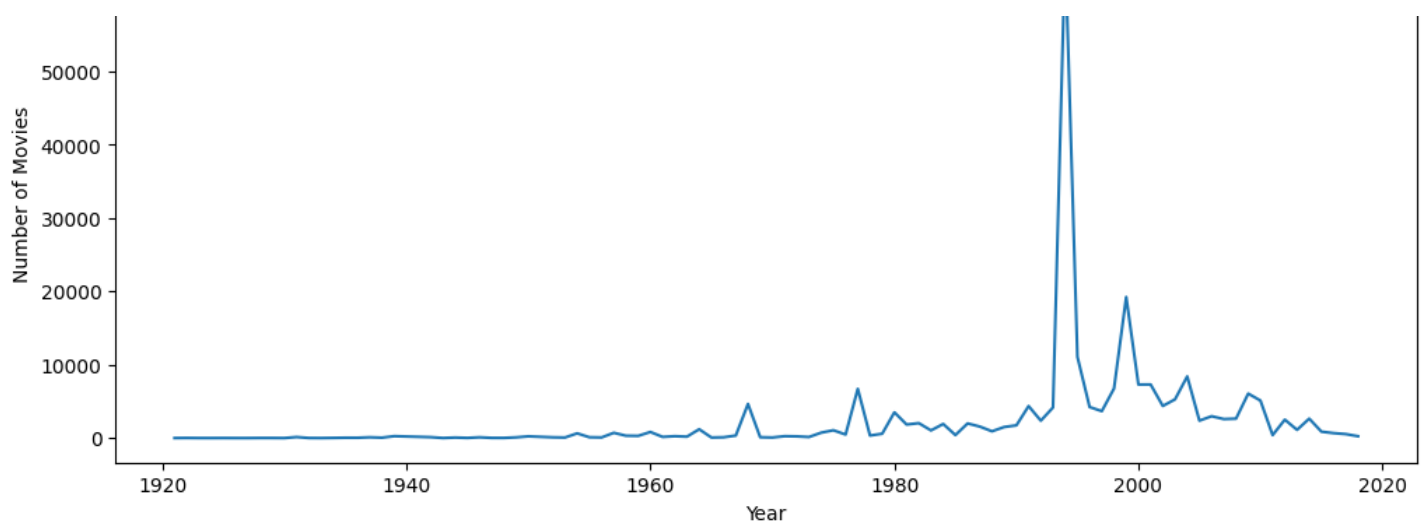


In [ ]:

```python
# visualize the year with the most produced movies
df_merge_cleaned['year'] = df_merge_cleaned['year'].dropna().astype(int)

yearly_counts = df_merge_cleaned['year'].value_counts().sort_index()

plt.figure(figsize=(12, 5))
sns.lineplot(x=yearly_counts.index, y=yearly_counts.values)
plt.title("Number of Movies Released Per Year")
plt.xlabel("Year")
plt.ylabel("Number of Movies")
plt.show()
```
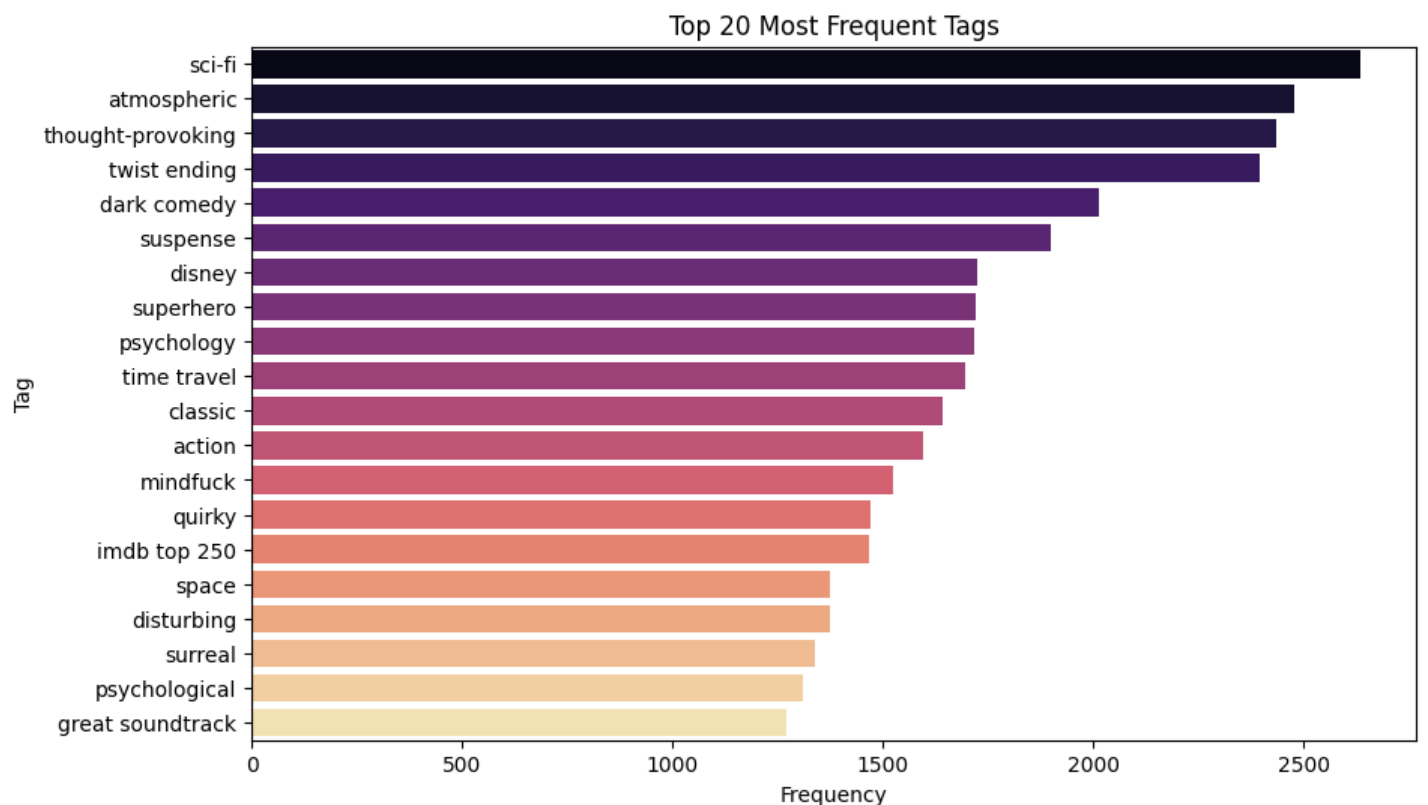
```
In [ ]:

#Top performing tags
top_tags = df_merge_cleaned['tag'].str.lower().value_counts().head(20)

plt.figure(figsize=(10, 6))
sns.barplot(y=top_tags.index, x=top_tags.values, palette='magma')
plt.title("Top 20 Most Frequent Tags")
plt.xlabel("Frequency")
plt.ylabel("Tag")
plt.show()
```
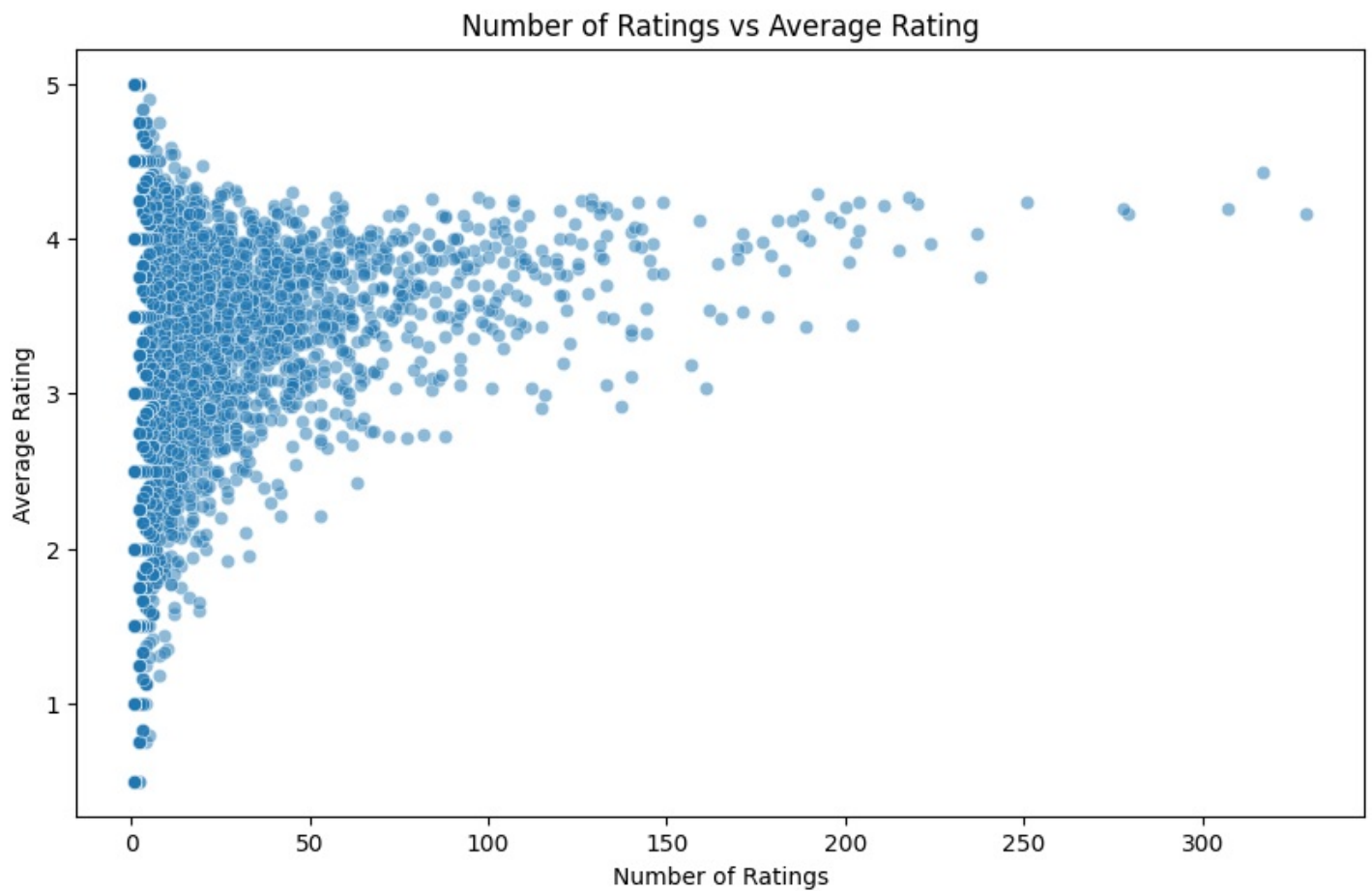
## Bivariate analysis

```
In [ ]:

#number of ratings vs average ratings
movie_stats = ratings.groupby('movieId')['rating'].agg(['count', 'mean']).reset_index()

plt.figure(figsize=(10,6))
sns.scatterplot(data=movie_stats, x='count', y='mean', alpha=0.5)
plt.title("Number of Ratings vs Average Rating")
plt.xlabel("Number of Ratings")
plt.ylabel("Average Rating")
```
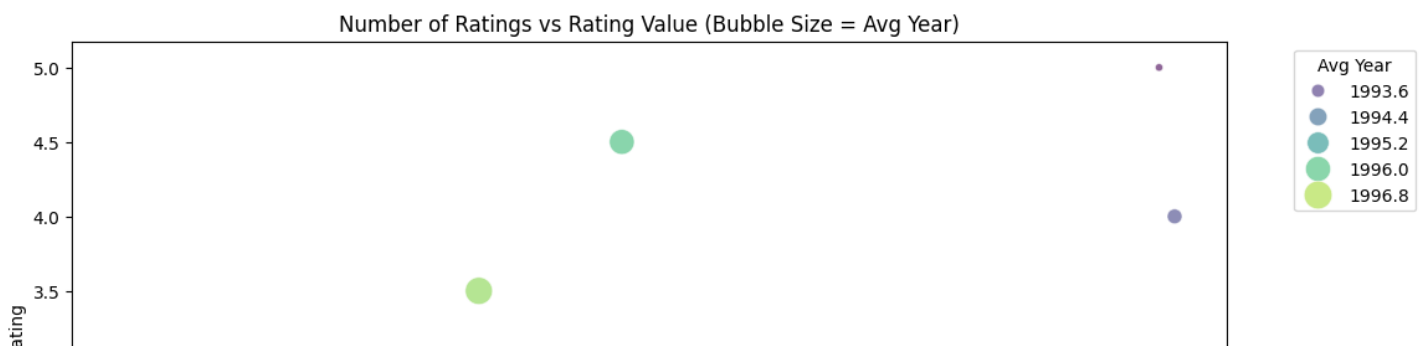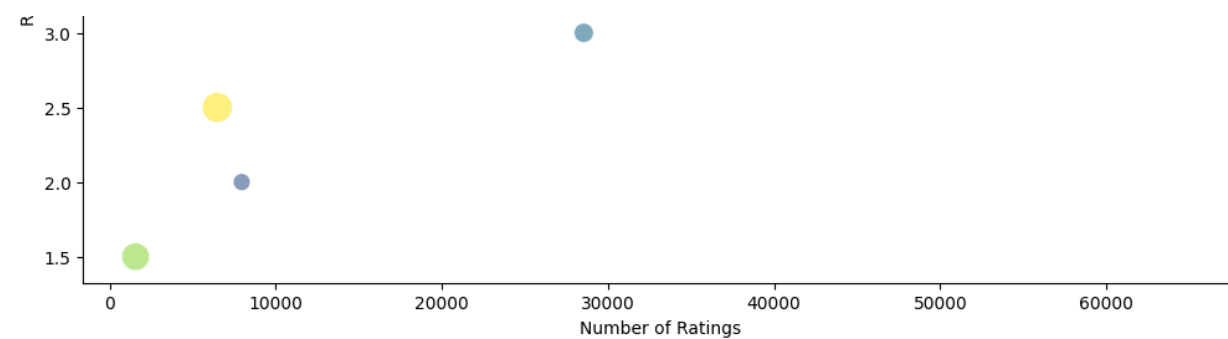
```
plt.show()
```



Number of Ratings vs Average Rating

## Mutivariate Analysis

In [ ]:

```
#count of ratings vs average rating vs average year
rating_summary = df_merge_cleaned.groupby('rating').agg(
    count=('rating', 'count'),
    mean_year=('year', 'mean')
).reset_index()


plt.figure(figsize=(12, 6))
sns.scatterplot(
    data=rating_summary,
    x='count',
    y='rating',
    size='mean_year',
    hue='mean_year',
    sizes=(20, 300), palette='viridis', alpha=0.6
)
plt.title("Number of Ratings vs Rating Value (Bubble Size = Avg Year)")
plt.xlabel("Number of Ratings")
plt.ylabel("Rating")
plt.legend(title='Avg Year', bbox_to_anchor=(1.05, 1), loc='upper left')
plt.show()
```



Number of Ratings vs Rating Value (Bubble Size = Avg Year)
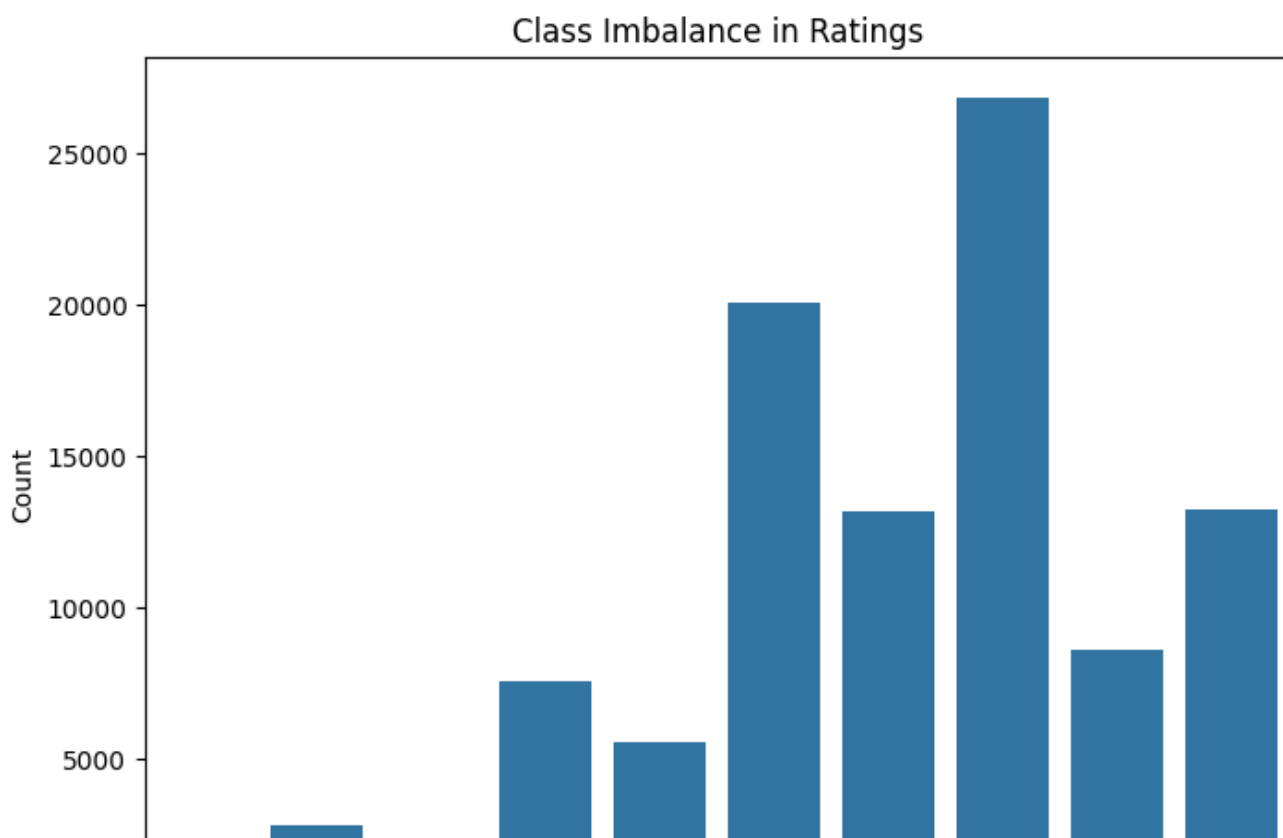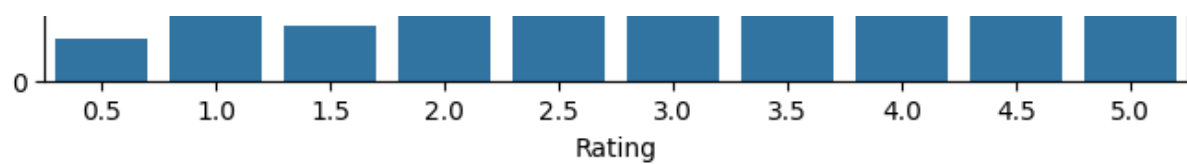
# Preprocessing

In [ ]:

```python
# check class imbalance
class_counts = ratings['rating'].value_counts()
print(class_counts)
```

```
rating
4.0    26818
3.0    20047
5.0    13211
3.5    13136
4.5     8551
2.0     7551
2.5     5550
1.0     2811
1.5     1791
0.5     1370
Name: count, dtype: int64
```

In [ ]:

```python
# visualize class imbalance
plt.figure(figsize=(8, 6))
sns.barplot(x=class_counts.index, y=class_counts.values)
plt.title("Class Imbalance in Ratings")
plt.xlabel("Rating")
plt.ylabel("Count")
plt.show()
```

## Encoding

In [ ]:

```python
#One-hot encoding
mlb = MultiLabelBinarizer()
genre_encoded = mlb.fit_transform(df_merge_cleaned['genres'].astype(str).str.split('|'))
genre_df = pd.DataFrame(genre_encoded, columns=mlb.classes_)


df_merge_cleaned = pd.concat([df_merge_cleaned, genre_df], axis=1)


df_merge_cleaned.head()
```

Out[ ]:

| | movieId | title | genres | userId_rating | rating | timestamp_rating | userId_tag | tag | tim |
|---|---|---|---|---|---|---|---|---|---|
| 0 | 1.0 | Toy Story | Adventure\|Animation\|Children\|Comedy\|Fantasy | 1.0 | 4.0 | 2000-07-30 18:45:03 | 336.0 | pixar | |
| 1 | 1.0 | Toy Story | Adventure\|Animation\|Children\|Comedy\|Fantasy | 1.0 | 4.0 | 2000-07-30 18:45:03 | 474.0 | pixar | |
| 2 | 1.0 | Toy Story | Adventure\|Animation\|Children\|Comedy\|Fantasy | 1.0 | 4.0 | 2000-07-30 18:45:03 | 567.0 | fun | |
| 3 | 1.0 | Toy Story | Adventure\|Animation\|Children\|Comedy\|Fantasy | 5.0 | 4.0 | 1996-11-08 06:36:02 | 336.0 | pixar | |
| 4 | 1.0 | Toy Story | Adventure\|Animation\|Children\|Comedy\|Fantasy | 5.0 | 4.0 | 1996-11-08 06:36:02 | 474.0 | pixar | |

**5 rows × 30 columns**

One hot encoding was done for the genre column since it is categorical data

In [ ]:

```python
df_merge_cleaned.shape
```

Out[ ]:

(233016, 30)

An increase in number of columns resulted but not so high and workable therefore no need for dimensionality reduction

In [ ]:

```python
# LabelEncoding
user_encoder = LabelEncoder()
movie_encoder = LabelEncoder()
userid_encoder = LabelEncoder()


df_merge_cleaned['userId_rating'] = user_encoder.fit_transform(df_merge_cleaned['userId_rating'])


df_merge_cleaned['movieId'] = movie_encoder.fit_transform(df_merge_cleaned['movieId'])
```

```
df_merge_cleaned['userId_tag'] = userid_encoder.fit_transform(df_merge_cleaned['userId_ta
g'])

df_merge_cleaned.head()
```

Out[ ]:

| | movieId | title | genres | userId_rating | rating | timestamp_rating | userId_tag | tag | tim |
|---|---|---|---|---|---|---|---|---|---|
| 0 | 0 | Toy Story | Adventure\|Animation\|Children\|Comedy\|Fantasy | 0 | 4.0 | 2000-07-30 18:45:03 | 32 | pixar | |
| 1 | 0 | Toy Story | Adventure\|Animation\|Children\|Comedy\|Fantasy | 0 | 4.0 | 2000-07-30 18:45:03 | 41 | pixar | |
| 2 | 0 | Toy Story | Adventure\|Animation\|Children\|Comedy\|Fantasy | 0 | 4.0 | 2000-07-30 18:45:03 | 50 | fun | |
| 3 | 0 | Toy Story | Adventure\|Animation\|Children\|Comedy\|Fantasy | 3 | 4.0 | 1996-11-08 06:36:02 | 32 | pixar | |
| 4 | 0 | Toy Story | Adventure\|Animation\|Children\|Comedy\|Fantasy | 3 | 4.0 | 1996-11-08 06:36:02 | 41 | pixar | |

**5 rows × 30 columns**

```
# This is formatted as code
```

**Label encoding was done for the userId_tag, movieId and userId_rating since it is ordinal data to be used in models such as neural net and matrix vectorisation**

## Utility Matrix

In [ ]:

```
# Utility matrix creation
utility_matrix = ratings.pivot_table(index='userId', columns='movieId', values='rating')
utility_matrix.head()
```

Out[ ]:

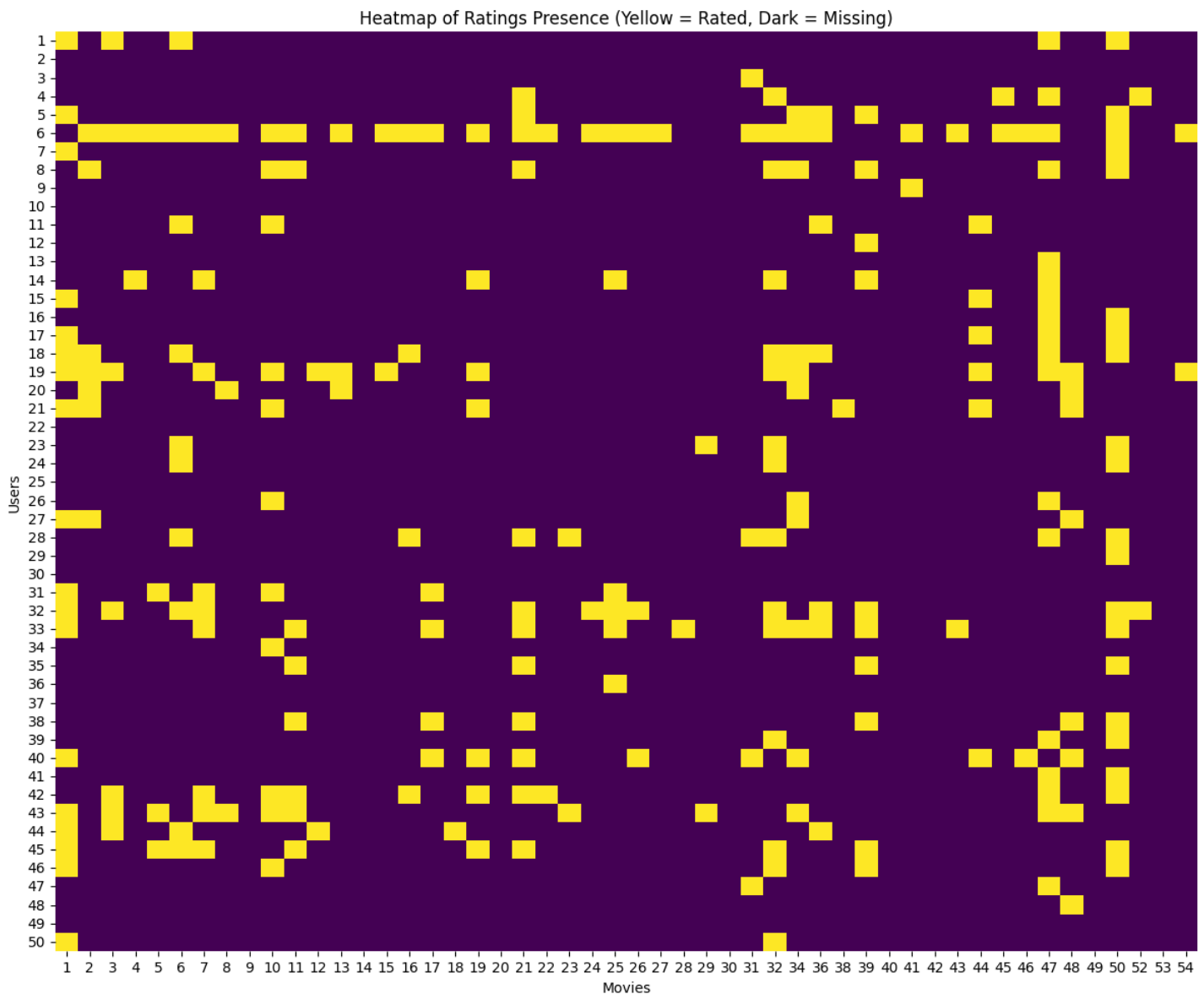| movieId | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | ... | 193565 | 193567 | 193571 | 193573 | 193579 | 193581 | 193 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| userId | | | | | | | | | | | | | | | | | | |
| 1 | 4.0 | NaN | 4.0 | NaN | NaN | 4.0 | NaN | NaN | NaN | NaN | ... | NaN | NaN | NaN | NaN | NaN | NaN | |
| 2 | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | ... | NaN | NaN | NaN | NaN | NaN | NaN | |
| 3 | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | ... | NaN | NaN | NaN | NaN | NaN | NaN | |
| 4 | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | ... | NaN | NaN | NaN | NaN | NaN | NaN | |
| 5 | 4.0 | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | NaN | ... | NaN | NaN | NaN | NaN | NaN | NaN | |

**5 rows × 9724 columns**

In [ ]:

```
# Spasity percentage
total_possible_ratings = utility_matrix.shape[0] * utility_matrix.shape[1]
actual_ratings = utility_matrix.count().sum()
sparsity = 1 - (actual_ratings / total_possible_ratings)
print(f"Matrix Sparsity: {sparsity:.4f} ({sparsity * 100:.2f}% missing)")
```

Matrix Sparsity: 0.9830 (98.30% missing)

```python
# Small subset for visualization
sampled_matrix = utility_matrix.iloc[:50, :50]
plt.figure(figsize=(12, 10))
sns.heatmap(sampled_matrix.notna(), cbar=False, cmap='viridis')
plt.title('Heatmap of Ratings Presence (Yellow = Rated, Dark = Missing)')
plt.xlabel('Movies')
plt.ylabel('Users')
plt.tight_layout()
plt.show()
```



Heatmap of Ratings Presence (Yellow = Rated, Dark = Missing)

**A utility matrix is important for user based collaborative filtering**

**A lot of null values identified meaning few users rated most movies**

## TF-IDF Vectorisation

In [ ]:

```python
tfidf = TfidfVectorizer(stop_words='english')
tfidf_matrix = tfidf.fit_transform(df_merge_cleaned['tag'].fillna(''))
tfidf_df = pd.DataFrame(tfidf_matrix.toarray(), columns=tfidf.get_feature_names_out())
tfidf_df.head()
```

Out [ ]:

| | 06 | 1900s | 1920s | 1950s | 1960s | 1970s | 1980s | 1990s | 2001 | 250 | ... | wrongful | wry | york | younger | zellweger | zither | zoe |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | ... | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |

| | 1 0.0<br>06 | 0.0<br>1900s | 0.0<br>1920s | 0.0<br>1950s | 0.0<br>1960s | 0.0<br>1970s | 0.0<br>1980s | 0.0<br>1990s | 0.0<br>2001 | 0.0<br>250 | ... | 0.0<br>wrongful | 0.0<br>wry | 0.0<br>york | 0.0<br>younger | 0.0<br>zellweger | 0.0<br>zither | 0.0<br>zoe |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 2 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | ... | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| 3 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | ... | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| 4 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | ... | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |

**5 rows × 1670 columns**

**Tf-IDF will be useful during content based filtering with special focus on the tags column**

# Scaling

In [ ]:

```
# performing minmax scaler for rating necessary for uniformity
scaler = MinMaxScaler()
df_exploded['rating'] = scaler.fit_transform(df_exploded[['rating']])
df_exploded.head()
```

Out [ ]:

| | movieId | title | genres | userId_rating | rating | timestamp_rating | userId_tag | tag | timestamp_tag | year |
|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | Toy Story | Adventure | 1 | 0.714286 | 2000-07-30 18:45:03 | 336 | pixar | 2006-02-04 09:36:04 | 1995 |
| 1 | 1 | Toy Story | Animation | 1 | 0.714286 | 2000-07-30 18:45:03 | 336 | pixar | 2006-02-04 09:36:04 | 1995 |
| 2 | 1 | Toy Story | Children | 1 | 0.714286 | 2000-07-30 18:45:03 | 336 | pixar | 2006-02-04 09:36:04 | 1995 |
| 3 | 1 | Toy Story | Comedy | 1 | 0.714286 | 2000-07-30 18:45:03 | 336 | pixar | 2006-02-04 09:36:04 | 1995 |
| 4 | 1 | Toy Story | Fantasy | 1 | 0.714286 | 2000-07-30 18:45:03 | 336 | pixar | 2006-02-04 09:36:04 | 1995 |

**Scaling is done to provide uniformity on the ratings column values especially for sensitive models such as SVD (Matrix Vectorisation)**

# Modeling

## Filtering

### User Based Collaborative Filtering

In [ ]:

```
#checking for the user id for pulp fiction to use it in the collaborative filtering input
s
def get_movie_id_by_title(movies_df, title):

  movie = movies_df[movies_df['title'] == title]
  if not movie.empty:
    return movie['movieId'].iloc[0]
  else:
    return None


movies = pd.read_csv('movies.csv')
movie_id = get_movie_id_by_title(movies, 'Pulp Fiction (1994)')

if movie_id:
  print(f"The movie ID for 'Pulp Fiction' is: {movie_id}")
else:
  print("Movie not found in the dataset.")
```

The movie ID for 'Pulp Fiction' is: 296

In [ ]:

```python
# Predicting top 5 movies from user based collaborative filtering using KNN
from sklearn.neighbors import NearestNeighbors
def get_knn_recommendations(user_id, movie_id, n=5):
    utility_matrix_filled = utility_matrix.fillna(0)

    knn_model = NearestNeighbors(metric='cosine', algorithm='brute', n_neighbors=6)
    knn_model.fit(utility_matrix_filled)

    distances, indices = knn_model.kneighbors(utility_matrix_filled.iloc[user_id - 1].va
lues.reshape(1, -1), n_neighbors=6)

    similar_user_ids = indices.flatten()[1:]
    similar_user_distances = distances.flatten()[1:]

    weighted_rating = 0
    for i, user in enumerate(similar_user_ids):
      rating = utility_matrix.iloc[user].get(movie_id)
      if pd.notna(rating):
        weighted_rating += rating * (1 - similar_user_distances[i])

    if similar_user_distances.sum() > 0:
        predicted_rating = weighted_rating / sum(1 - similar_user_distances)
    else:
        predicted_rating = 0

    weighted_ratings = {}
    for i, user in enumerate(similar_user_ids):
        for movie, rating in utility_matrix.iloc[user].items():
            if pd.notna(rating) and movie not in utility_matrix.iloc[user_id-1].dropna()
.index:
                weighted_ratings.setdefault(movie, 0)
                weighted_ratings[movie] += rating * (1 - similar_user_distances[i])

    for movie, rating in weighted_ratings.items():
        weighted_ratings[movie] = rating / sum(1 - similar_user_distances) if sum(1 - si
milar_user_distances) > 0 else 0


    recommendations = sorted(weighted_ratings.items(), key=lambda x: x[1], reverse=True)
[:n]

    return recommendations

recommendations_user_based = get_knn_recommendations(user_id=68, movie_id=296)
print("\nRecommendations for user 68 and movie 296:")
recommendations_user_based
```

Recommendations for user 68 and movie 296:

Out[ ]:

```
[(750, np.float64(4.209720908281964)),
 (5956, np.float64(3.799335285590348)),
 (32, np.float64(3.5283076613659614)),
 (5903, np.float64(3.516133680984722)),
 (778, np.float64(3.3313719270314213))]
```

In [ ]:

```python
#getting movie titles from movieId
def get_movie_titles(df, movie_ids):

  if not isinstance(df, pd.DataFrame) or not {'movieId', 'title'}.issubset(df.columns):
    print("Error: df must have 'movieId' and 'title' columns.")
    return []

  if not isinstance(movie_ids, list) or not all(isinstance(movie_id, (int, float)) for m
ovie_id in movie_ids):
```

```
            print("Error:'movie_ids' must be a list of numbers.")
            return []

    movie_titles = []
    for movie_id in movie_ids:
        title = df[df['movieId'] == movie_id]['title'].values
        if len(title) > 0:
            movie_titles.append(title[0])
        else:
            print(f"Warning: No movie found with ID {movie_id}")
            movie_titles.append(None)

    return movie_titles


movie_ids_to_find = [750, 5956,32 , 5903, 778]

movie_titles = get_movie_titles(df_merge_cleaned, movie_ids_to_find)
movie_titles
```

Warning: No movie found with ID 5956
Warning: No movie found with ID 5903

Out[ ]:

```
['Official Story, The (La historia oficial)',
 None,
 'In the Bleak Midwinter',
 None,
 'Evil Dead, The']
```

In [ ]:

```
#identifying precision@K and recall@K metrics
recommended_movie_ids = [movie[0] for movie in recommendations_user_based]
liked_threshold = 3.5
user_actual_ratings = utility_matrix.loc[68]
relevant_movies = user_actual_ratings[user_actual_ratings >= liked_threshold].index.tolis
t()


def precision_recall_at_k(recommended, relevant, k):
    recommended_k = recommended[:k]
    relevant_set = set(relevant)
    hits = len([movie for movie in recommended_k if movie in relevant_set])

    precision = hits / k
    recall = hits / len(relevant_set) if relevant_set else 0.0
    return precision, recall


k = 5
precision, recall = precision_recall_at_k(recommended_movie_ids, relevant_movies, k)


print(f"Precision@{k}: {precision:.3f}")
print(f"Recall@{k}: {recall:.3f}")
```

Precision@5: 0.000
Recall@5: 0.000

**We have poor metrics this might be due to the cold start problem since most of the users have not rated the movies. It is clear from the utility matrix that there are a lot of null values**

**Item Based Collaborative Filtering**

In [ ]:

```
# Predicting top 5 movies using item based filtering using cosine similarity
```

```python
from sklearn.neighbors import NearestNeighbors

def get_item_based_recommendations(user_id, movie_id, n=5):

    utility_matrix_filled = utility_matrix.fillna(0)

    item_similarity = cosine_similarity(utility_matrix_filled.T)
    item_similarity_df = pd.DataFrame(item_similarity, index=utility_matrix.columns, col
umns=utility_matrix.columns)

    user_ratings = utility_matrix.loc[user_id].dropna()

    scores = item_similarity_df[movie_id].drop(movie_id)
    scores = scores * user_ratings.mean()


    recommendations = scores.sort_values(ascending=False).head(n).index.tolist()

    return recommendations

from sklearn.metrics.pairwise import cosine_similarity
recommendations_item_based = get_item_based_recommendations(user_id=68, movie_id=296)
print("\nItem-Based CF Recommendations for user 68 and movie 296:")
print(recommendations_item_based)
movie_titles = get_movie_titles(movies, recommendations_item_based)
movie_titles
```

```
Item-Based CF Recommendations for user 68 and movie 296:
[593, 318, 47, 356, 50]
```

Out[ ]:

```
['Silence of the Lambs, The (1991)',
 'Shawshank Redemption, The (1994)',
 'Seven (a.k.a. Se7en) (1995)',
 'Forrest Gump (1994)',
 'Usual Suspects, The (1995)']
```

In [ ]:

```python
#identifying precision@K and recall@K metrics
recommended_movie_ids = recommendations_item_based
liked_threshold = 3.5
user_actual_ratings = utility_matrix.loc[68]
relevant_movies = user_actual_ratings[user_actual_ratings >= liked_threshold].index.tolis
t()


def precision_recall_at_k(recommended, relevant, k):
    recommended_k = recommended[:k]
    relevant_set = set(relevant)
    hits = len([movie for movie in recommended_k if movie in relevant_set])

    precision = hits / k
    recall = hits / len(relevant_set) if relevant_set else 0.0
    return precision, recall

k = 5
precision, recall = precision_recall_at_k(recommended_movie_ids, relevant_movies, k)


print(f"Precision@{k}: {precision:.3f}")
print(f"Recall@{k}: {recall:.3f}")
```

```
Precision@5: 0.600
Recall@5: 0.005
```

The metrics above show that i was able to recommend 3 movies out of 5 already liked by the user. The recall shows that I have recommended very few movies out of the large pool of liked movies

There is significant improvement on the precision metrics since there was a selection based on the item as opposed to the user. This has reduced the cold start problem

**Matrix Vectorisation using SVD**

In [ ]:

```
#predicting top 5 movies using matrix vectorisation

def get_svd_recommendations_for_movie(user_id, movie_id, n=5):

    svd = TruncatedSVD(n_components=50, random_state=42)
    matrix_svd = svd.fit_transform(utility_matrix.fillna(utility_matrix.mean(axis=0)))

    user_ratings = matrix_svd[user_id - 1]

    try:
        movie_index = utility_matrix.columns.get_loc(movie_id)
    except KeyError:
        print(f"Movie ID {movie_id} not found in the utility matrix.")
        return []

    movie_latent_vector = svd.components_.T[movie_index]

    predicted_rating = np.dot(user_ratings, movie_latent_vector)

    movie_similarities = cosine_similarity([movie_latent_vector], svd.components_.T)

    similar_movies_indices = movie_similarities.argsort()[0][::-1][1:n+1]
    similar_movies = [utility_matrix.columns[i] for i in similar_movies_indices]

    rated_movies = utility_matrix.loc[user_id - 1].dropna().index
    similar_movies = [mid for mid in similar_movies if mid not in rated_movies]

    recommendations = []
    for movie_id in similar_movies:
        pred_rating = np.dot(user_ratings, svd.components_.T[utility_matrix.columns.get_
loc(movie_id)])
        recommendations.append((movie_id, pred_rating))

    return recommendations

recommendations = get_svd_recommendations_for_movie(user_id=68, movie_id=296)
print("\nSVD Recommendations for user 68 and movie 296:")
print(recommendations)

movie_ids_to_find = [movie[0] for movie in recommendations]
movie_titles = get_movie_titles(movies, recommendations)
movie_titles
```

```
SVD Recommendations for user 68 and movie 296:
[(np.int64(509), np.float64(3.5057600391478747)), (np.int64(32), np.float64(4.06991647189
5481)), (np.int64(107141), np.float64(3.6872740599711578)), (np.int64(1966), np.float64(3
.649465984405974)), (np.int64(171), np.float64(3.833297973290365))]
Error:'movie_ids' must be a list of numbers.
```

Out[ ]:

```
[]
```

In [ ]:

```
# Matching MovieId to title
def get_movie_titles(df, movie_ids):

  if not isinstance(df, pd.DataFrame) or not {'movieId', 'title'}.issubset(df.columns):
    print("Error: Input 'df' must be a DataFrame with 'movieId' and 'title' columns.")
    return []

  if not isinstance(movie_ids, list) or not all(isinstance(movie_id, (int, float)) for m
ovie_id in movie_ids):
```

```
        print("Error: Input 'movie_ids' must be a list of numbers.")
        return []

    movie_titles = []
    for movie_id in movie_ids:
        title = df[df['movieId'] == movie_id]['title'].values
        if len(title) > 0:
            movie_titles.append(title[0])
        else:
            print(f"Warning: No movie found with ID {movie_id}")
            movie_titles.append(None)

    return movie_titles

movie_ids_to_find = [509, 32, 107141, 1966, 171]

movie_titles = get_movie_titles(df_merge_cleaned, movie_ids_to_find)
movie_titles
```

```
Warning: No movie found with ID 107141
Warning: No movie found with ID 1966
```

Out[ ]:

```
['Addams Family, The',
 'In the Bleak Midwinter',
 None,
 None,
 'Independence Day (a.k.a. ID4)']
```

In [ ]:

```
#identifying metrics for precision@K and recall@K
recommended_movie_ids = recommendations
liked_threshold = 3.5
user_actual_ratings = utility_matrix.loc[68]
relevant_movies = user_actual_ratings[user_actual_ratings >= liked_threshold].index.tolis
t()


def precision_recall_at_k(recommended, relevant, k):
    recommended_k = recommended[:k]
    relevant_set = set(relevant)
    hits = len([movie for movie in recommended_k if movie in relevant_set])

    precision = hits / k
    recall = hits / len(relevant_set) if relevant_set else 0.0
    return precision, recall

k = 5
precision, recall = precision_recall_at_k(recommended_movie_ids, relevant_movies, k)


print(f"Precision@{k}: {precision:.3f}")
print(f"Recall@{k}: {recall:.3f}")
```

```
Precision@5: 0.000
Recall@5: 0.000
```

**Matrix vectorisation has poor perfomance and may not be performing as well due to a matrix factorisation which
has reduced the total data shape and also due to imputing mean which may not give the actual values**


**Content Based Filtering**

In [ ]:

```
# predicting top 5 movies using content based fitering while using tags

def get_content_based_recommendations(movie_title, n=5, df_exploded=None):
```

```python
    if df_exploded is None:
        return ["Error: DataFrame not provided."]

    df_grouped = df_merge_cleaned.groupby(['movieId', 'title'])['tag'].apply(lambda x: '
'.join(x.dropna())).reset_index()

    if movie_title not in df_grouped['title'].values:
        print(f"Movie '{movie_title}' not found in the dataset.")
        return []

    tfidf = TfidfVectorizer(stop_words='english')
    tfidf_matrix = tfidf.fit_transform(df_grouped['tag'].fillna(''))

    idx = df_grouped[df_grouped['title'] == movie_title].index[0]

    cosine_sim = cosine_similarity(tfidf_matrix, tfidf_matrix)

    sim_scores = list(enumerate(cosine_sim[idx]))

    sim_scores = sorted(sim_scores, key=lambda x: x[1], reverse=True)
    sim_scores = sim_scores[1:n+1]

    movie_indices = [i[0] for i in sim_scores]
    recommended_titles = df_grouped['title'].iloc[movie_indices].tolist()

    return recommended_titles

sample_data = df_exploded

df_exploded = pd.DataFrame(sample_data)

recommendations = get_content_based_recommendations("Pulp Fiction", n=5, df_exploded=df_e
xploded)

print("Top 5 content-based recommendations for 'Pulp Fiction':")
for i, title in enumerate(recommendations, 1):
    print(f"{i}. {title}")
```

```
Top 5 content-based recommendations for 'Pulp Fiction':
1. Reservoir Dogs
2. Big Lebowski, The
3. Django Unchained
4. Kiss Kiss Bang Bang
5. Fight Club
```

In [ ]:

```python
#identifying precision@K and recall@K metrics
test_movies = ["Pulp Fiction"]
k = 5
liked_threshold = 3.5

precision_list = []
recall_list = []

for movie_title in test_movies:
    recommended_titles = get_content_based_recommendations(movie_title, n=k)
    recommended_ids = movies[movies['title'].isin(recommended_titles)]['movieId'].tolist
()

    movie_id = movies[Movies['title'].str.contains(movie_title)]['movieId'].values[0]
    users_who_liked_input = utility_matrix[utility_matrix[movie_id] >= liked_threshold].
index

    relevant_movies = set()
    for user in users_who_liked_input:
        liked_movies = utility_matrix.loc[user]
        liked_movie_ids = liked_movies[liked_movies >= liked_threshold].index.tolist()
        liked_movie_ids = [mid for mid in liked_movie_ids if mid != movie_id]
        relevant_movies.update(liked_movie_ids)

    hits = len([mid for mid in recommended_ids if mid in relevant_movies])
```

```
    precision = hits / k
    recall = hits / len(relevant_movies) if relevant_movies else 0.0

    precision_list.append(precision)
    recall_list.append(recall)

    print(f"\nInput Movie: {movie_title}")
    print(f"Precision@{k}: {precision:.3f}")
    print(f"Recall@{k}: {recall:.3f}")
```

```
Input Movie: Pulp Fiction
Precision@5: 0.000
Recall@5: 0.000
```

For content based filtering we may not also have great metrics since the result may be affected by the sparsity of the tf-idf matrix

**Hybrid: Content based + User Based Collaborative filtering**

In [ ]:

```
# Predicting top 5 movies based on the hybrid model using userid 68 and movie pulp fiction

def hybrid_recommendations(user_id, movie_title, n_recommendations=5):

    content_based_recs = get_content_based_recommendations(movie_title, n=n_recommendations) # Changed 'n' to 'n_recommendations'

    content_based_ids = movies[movies['title'].isin(content_based_recs)]['movieId'].tolist()

    movie_id = movies[movies['title'].str.contains(movie_title)]['movieId'].iloc[0]
    collaborative_recs = get_item_based_recommendations(user_id, movie_id, n=n_recommendations)

    hybrid_recs = content_based_ids + collaborative_recs
    hybrid_recs = list(set(hybrid_recs))

    hybrid_movie_titles = get_movie_titles(movies, hybrid_recs)

    return hybrid_movie_titles[:n_recommendations]

user_id = 68
movie_title = "Pulp Fiction"
recs = hybrid_recommendations(user_id, movie_title, n_recommendations=5)

print(f"\nHybrid Recommendations for User {user_id} based on '{movie_title}':")
for i, title in enumerate(recs, 1):
    print(f"{i}. {title}")
```

```
Hybrid Recommendations for User 68 based on 'Pulp Fiction':
1. Forrest Gump (1994)
2. Seven (a.k.a. Se7en) (1995)
3. Silence of the Lambs, The (1991)
4. Usual Suspects, The (1995)
5. Shawshank Redemption, The (1994)
```

In [ ]:

```
#identifying precision@K and recall@k metrics
def precision_recall_at_k_hybrid(user_id, movie_title, k, ratings_df, movies_df, threshold=3.5):

    user_ratings = ratings_df[ratings_df['userId'] == user_id]
    liked_movie_ids = user_ratings[user_ratings['rating'] >= threshold]['movieId'].tolist()
    relevant_set = set(liked_movie_ids)
```

```
    recommended_titles = hybrid_recommendations(user_id, movie_title, n_recommendations=
k)

    recommended_ids = movies_df[movies_df['title'].isin(recommended_titles)]['movieId'].
tolist()

    hits = len([mid for mid in recommended_ids if mid in relevant_set])

    precision = hits / k

    recall = hits / len(relevant_set) if relevant_set else 0.0

    return precision, recall

user_id = 68
movie_title = "Pulp Fiction"
k = 5

precision, recall = precision_recall_at_k_hybrid(user_id, movie_title, k, ratings, movie
s)

print(f"Precision@{k}: {precision:.3f}")
print(f"Recall@{k}: {recall:.3f}")
```

```
Precision@5: 0.600
Recall@5: 0.005
```

**From the above it is clear that a combination of both content based and collaboration filtering have significantly improved the prediction to 3 out of 5 liked movies similar to item based filtering**

**Recall is still not of key focus since customers will rate multiple movies but only want a recommendation of a few**

## Classification

### Random Forest Classifier

In [ ]:

```
# classification with random forest model for movie ratings <3.5 to be liked and >3.5 not
liked and giving metrics for precision and recall

df_merge_cleaned['liked'] = df_merge_cleaned['rating'] < 3.5


X = df_merge_cleaned.drop(columns=['rating','year','timestamp_rating','genres','title','
tag','timestamp_tag', 'liked'])
y = df_merge_cleaned['liked']

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42
)

rf_classifier = RandomForestClassifier(n_estimators=100, random_state=42)
rf_classifier.fit(X_train, y_train)

y_pred = rf_classifier.predict(X_test)

precision = precision_score(y_test, y_pred)
recall = recall_score(y_test, y_pred)

print(f"Precision: {precision}")
print(f"Recall: {recall}")
```

```
Precision: 0.8609767334715503
Recall: 0.8293575945855985
```

**This is a good performance on both metrics. More exploration needs to be done on other models to find out whether we can get better performance**

## Decision Tree Classifier

In [ ]:

```
# classification with decision tree  model for movie ratings <3.5 to be liked and >3.5 not liked and giving metrics for precision and recall

dt_classifier = DecisionTreeClassifier(random_state=42)
dt_classifier.fit(X_train, y_train)

y_pred = dt_classifier.predict(X_test)

precision = precision_score(y_test, y_pred)
recall = recall_score(y_test, y_pred)

print(f"Precision: {precision}")
print(f"Recall: {recall}")
```

```
Precision: 0.8394993962015589
Recall: 0.8484411405747254
```

**This is also good performance but slightly lower than random forest model. More exploration required**

## Gradient Boost Classifier

In [ ]:

```
# classification with gradient boost model for movie ratings <3.5 to be liked and >3.5 not liked and giving metrics for precision and recall
from sklearn.impute import SimpleImputer
imputer = SimpleImputer(strategy='mean')
X_train_imputed = imputer.fit_transform(X_train)
X_test_imputed = imputer.transform(X_test)

gb_classifier = GradientBoostingClassifier(n_estimators=100, learning_rate=0.1, random_state=42)
gb_classifier.fit(X_train_imputed, y_train)

y_pred = gb_classifier.predict(X_test_imputed)

precision = precision_score(y_test, y_pred)
recall = recall_score(y_test, y_pred)

print(f"Precision: {precision}")
print(f"Recall: {recall}")
```

```
Precision: 0.7982954545454546
Recall: 0.031177188505492066
```

**This model performed poorly especially on recall and may not be reliable for the outcome we would like to achieve**

## Neural Net Model

In [ ]:

```
#Train test split with validation

X = df_merge_cleaned.drop(columns=['rating','year','timestamp_rating','genres','title','tag','timestamp_tag'])
y = df_merge_cleaned['rating']

X_train, X_temp, y_train, y_temp = train_test_split(X, y, test_size=0.2, random_state=42)
```

```
X_val, X_test, y_val, y_test = train_test_split(X_temp, y_temp, test_size=0.5, random_st
ate=42)
```

In [ ]:

```python
# Developing neural net model with classification for movie ratings <3.5 to be liked and
>3.5 not liked and giving metrics for precision and recall
import tensorflow as tf
X_train = X_train.astype(np.float32)
X_val = X_val.astype(np.float32)

y_train_binary = (y_train > 3.5).astype(int)
y_val_binary = (y_val > 3.5).astype(int)

model = tf.keras.Sequential([
    tf.keras.layers.Dense(64, activation='relu', input_shape=(X_train.shape[1],)),
    tf.keras.layers.Dense(32, activation='relu'),
    tf.keras.layers.Dense(1, activation='sigmoid')
 ])

model.compile(optimizer='adam', loss='binary_crossentropy', metrics=['accuracy'])

history = model.fit(X_train, y_train_binary, epochs=5, batch_size=32)

y_pred_probs = model.predict(X_val)
y_pred_binary = (y_pred_probs > 0.5).astype(int)

precision = precision_score(y_val_binary, y_pred_binary)
recall = recall_score(y_val_binary, y_pred_binary)

print(f"Precision: {precision:.3f}")
print(f"Recall: {recall:.3f}")
```

```
Epoch 1/5
5826/5826 ──────────────────── 29s 4ms/step - accuracy: 0.3134 - loss: nan
Epoch 2/5
5826/5826 ──────────────────── 18s 3ms/step - accuracy: 0.3117 - loss: nan
Epoch 3/5
5826/5826 ──────────────────── 17s 2ms/step - accuracy: 0.3131 - loss: nan
Epoch 4/5
5826/5826 ──────────────────── 27s 4ms/step - accuracy: 0.3107 - loss: nan
Epoch 5/5
5826/5826 ──────────────────── 14s 2ms/step - accuracy: 0.3133 - loss: nan
729/729 ──────────────── 1s 2ms/step
Precision: 0.000
Recall: 0.000
```

In [ ]:

```python
#improving the above neural net model by standardisation, introduction of class balance a
nd prediction threshold
scaler = StandardScaler()
X_train_scaled = scaler.fit_transform(X_train)
X_val_scaled = scaler.transform(X_val)


y_train_binary = (y_train > 3.5).astype(int)
y_val_binary = (y_val > 3.5).astype(int)

class_weights = class_weight.compute_class_weight(
    class_weight='balanced',
    classes=np.unique(y_train_binary),
    y=y_train_binary
)
class_weights_dict = dict(enumerate(class_weights))

model = tf.keras.Sequential([
    tf.keras.layers.Dense(128, activation='relu', input_shape=(X_train_scaled.shape[1],)
),
    tf.keras.layers.Dense(64, activation='relu'),
    tf.keras.layers.Dense(32, activation='relu'),
```

```
    tf.keras.layers.Dense(1, activation='sigmoid')
])

optimizer = tf.keras.optimizers.Adam(learning_rate=0.0001)

model.compile(optimizer=optimizer, loss='binary_crossentropy', metrics=['accuracy'])

history = model.fit(X_train_scaled, y_train_binary, epochs=10, batch_size=32, class_weig
ht=class_weights_dict)

y_pred_probs = model.predict(X_val_scaled)
y_pred_binary = (y_pred_probs > 0.3).astype(int)

print(classification_report(y_val_binary, y_pred_binary))
```

```
Epoch 1/10
5826/5826 ──────────────── 19s 3ms/step - accuracy: 0.3141 - loss: nan
Epoch 2/10
5826/5826 ──────────────── 34s 5ms/step - accuracy: 0.3115 - loss: nan
Epoch 3/10
5826/5826 ──────────────── 18s 3ms/step - accuracy: 0.3146 - loss: nan
Epoch 4/10
5826/5826 ──────────────── 25s 4ms/step - accuracy: 0.3128 - loss: nan
Epoch 5/10
5826/5826 ──────────────── 36s 3ms/step - accuracy: 0.3126 - loss: nan
Epoch 6/10
5826/5826 ──────────────── 17s 3ms/step - accuracy: 0.3125 - loss: nan
Epoch 7/10
5826/5826 ──────────────── 20s 3ms/step - accuracy: 0.3140 - loss: nan
Epoch 8/10
5826/5826 ──────────────── 20s 3ms/step - accuracy: 0.3143 - loss: nan
Epoch 9/10
5826/5826 ──────────────── 16s 3ms/step - accuracy: 0.3128 - loss: nan
Epoch 10/10
5826/5826 ──────────────── 21s 3ms/step - accuracy: 0.3123 - loss: nan
729/729 ──────────────── 1s 1ms/step
              precision    recall  f1-score   support

           0       0.31      1.00      0.48      7310
           1       0.00      0.00      0.00     15992

    accuracy                           0.31     23302
   macro avg       0.16      0.50      0.24     23302
weighted avg       0.10      0.31      0.15     23302
```

**The above code was introduced to tune the previous model by introducing scaling, class balance to ensure class 0 and 1 are balance and reducing the prediction threshold to ensure class 1 gets better metrics however not much improvement was noticed and further future investigations need to be done on this model**

## Evaluation

In [ ]:

```
#results of filtering models
model_results = {
    'User Based Collaborative Filtering': {'Precision@K': 0.0, 'Recall@K': 0.0},
    'Item Based Collaborative Filtering': {'Precision@K': 0.6, 'Recall@K': 0.005},
    'Matrix Vectorization SVD': {'Precision@K': 0.0, 'Recall@K': 0.0},
    'Content Based Filtering': {'Precision@K': 0.0, 'Recall@K': 0.0},
    'Hybrid: Content Based and Collaborative Filtering': {'Precision@K': 0.6, 'Recall@K':
0.005}
}

results_df = pd.DataFrame.from_dict(model_results, orient='index')

results_df
```

Out[ ]:

|  | Precision@K | Recall@K |
|---|---|---|
| User Based Collaborative Filtering | 0.0 | 0.000 |
| Item Based Collaborative Filtering | 0.6 | 0.005 |
| Matrix Vectorization SVD | 0.0 | 0.000 |
| Content Based Filtering | 0.0 | 0.000 |
| Hybrid: Content Based and Collaborative Filtering | 0.6 | 0.005 |

User-Based Collaborative Filtering showed poor performance, largely due to the cold start problem—only a few users have rated many movies, while most movies lack sufficient ratings. This is also evident from the high number of null values in the utility matrix.

SVD (Singular Value Decomposition) underperformed. While it's a powerful dimensionality reduction technique, it can limit the scope of the dataset, especially when imputation (e.g., mean filling) is used, potentially introducing bias or or noise.

Item-Based Collaborative Filtering performed significantly better, as it focuses on item similarities rather than user behavior. This approach is less affected by sparsity and cold-start issues.

The Hybrid Model delivered the best results by combining collaborative and content-based filtering, leveraging the strengths of both approaches. Precision was prioritized over recall in this case to better assess the accuracy of the recommended movies in matching user preferences. The system was able to accurately recommend 3 out of 5 movies on average.

While recall provides the percentage of relevant items retrieved from all possible relevant ones, it was less informative here due to the imbalance in rating distribution.

In [ ]:

```
#Results from classification models

model_results = {
    'Random Forest Classifier': {'precision': 0.875, 'recall': 0.822},
    'Gradient Boosting Classifier': {'precision': 0.747, 'recall': 0.04},
    'Neural Network Classifier': {'precision': 0.0, 'recall': 0.0},
    'Decision Tree Classifier': {'precision': 0.838, 'recall': 0.8408}
}

results_df = pd.DataFrame.from_dict(model_results, orient='index')

results_df
```

Out[ ]:

|  | precision | recall |
|---|---|---|
| Random Forest Classifier | 0.875 | 0.8220 |
| Gradient Boosting Classifier | 0.747 | 0.0400 |
| Neural Network Classifier | 0.000 | 0.0000 |
| Decision Tree Classifier | 0.838 | 0.8408 |

In [ ]:

```
# visualization for the models used in the project

model_results_filtering = {
    'User Based Collaborative Filtering': {'Precision@K': 0.0, 'Recall@K': 0.0},
    'Item Based Collaborative Filtering': {'Precision@K': 0.6, 'Recall@K': 0.005},
    'Matrix Vectorization SVD': {'Precision@K': 0.0, 'Recall@K': 0.0},
    'Content Based Filtering': {'Precision@K': 0.0, 'Recall@K': 0.0},
    'Hybrid: Content Based and Collaborative Filtering': {'Precision@K': 0.6, 'Recall@K':
0.005}
}
```

```python
model_results_classification = {
    'Random Forest Classifier': {'precision': 0.875, 'recall': 0.822},
    'Gradient Boosting Classifier': {'precision': 0.747, 'recall': 0.04},
    'Neural Network Classifier': {'precision': 0.0, 'recall': 0.0},
    'Decision Tree Classifier': {'precision': 0.838, 'recall': 0.8408}
}


results_df_filtering = pd.DataFrame.from_dict(model_results_filtering, orient='index')
results_df_classification = pd.DataFrame.from_dict(model_results_classification, orient='index')


plt.figure(figsize=(12, 6))

plt.subplot(1, 2, 1)
results_df_filtering.plot(kind='bar', ax=plt.gca())
plt.title('Filtering Models Evaluation Metrics')
plt.ylabel('Score')
plt.xticks(rotation=45, ha='right')

plt.subplot(1, 2, 2)
results_df_classification.plot(kind='bar', ax=plt.gca())
plt.title('Classification Models Evaluation Metrics')
plt.ylabel('Score')
plt.xticks(rotation=45, ha='right')

plt.tight_layout()
plt.show()
```
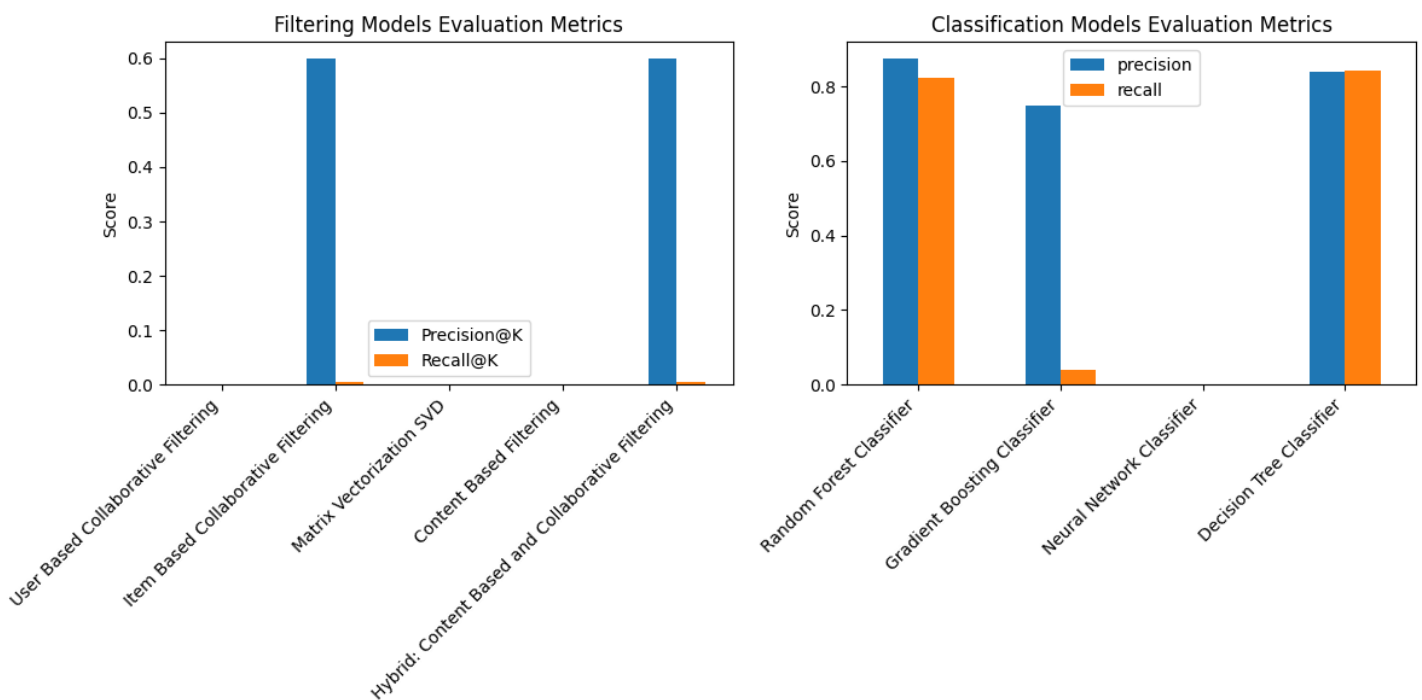


The Neural Network model was the least effective, primarily due to class imbalance and low confidence in predicting positive (class 1) outcomes, which are critical for identifying movies users are likely to enjoy. Further hyperparameter tuning and class balancing techniques may improve its performance.

The Random Forest model, on the other hand, performed best—especially in terms of precision, achieving 87%, which makes it highly reliable for accurately recommending good movies. Precision was chosen as the primary metric because it better captures the quality of recommendations by focusing on how many of the recommended movies are actually liked by users. Since users typically prefer a smaller set of high-quality suggestions, precision is a more meaningful metric than recall in this context. This model goes beyond traditional filtering methods (e.g., collaborative or content-based) by introducing a broader variety of highly rated movies, even those that may not match a user's historical preferences or content similarity patterns.

Ultimately, this approach can expand users' viewing experience, encouraging exploration of diverse genres and styles, rather than just reinforcing past behaviors.

# Summary

**Project Objectives & Outcomes**

**1. Design a collaborative filtering–based recommendation system to accurately predict user preferences.**

We implemented both user-based and item-based collaborative filtering techniques. While they showed potential, we observed that integrating content-based filtering or building a hybrid model significantly improved precision. One key challenge was the cold start problem, which hindered performance when ratings were sparse.

**2. Mitigate the cold start problem for new users and movies with limited or no ratings.**

Due to the limited number of user ratings, the cold start problem was quite evident. To address this, we incorporated content-based filtering and item similarity–based methods, which rely on metadata rather than user history. The hybrid model—which blends collaborative and content-based methods—proved especially effective in mitigating cold start scenarios.

**3. Develop classification models to distinguish between highly liked and disliked movies.**

By introducing classification models, we provided a mechanism for users to explore movies beyond their usual preferences. These models help uncover hidden gems and challenge user biases, potentially increasing their exposure to a broader variety of content. This strategy also helps address the long tail problem by recommending lesser-known but highly rated movies.

**4. Recommend top 5 movies based on user and particular movie watched.**

Top 5 movies based on user 68 and movie 'Pulp Fiction' were;

1. Forrest Gump
2. Ferris Bueller's Day Off
3. Seven (a.k.a. Se7en)
4. Silence of the Lambs
5. The Usual Suspects

**The success factors were also achieved:**

**1. Model Performance: Achieve a high precision@K > 0.6 for accurate recommndation to the users**

A high precision was achieved of 0.6 which will now have 3 out of 5 movies recommended as per the user preference

**2. User Satisfaction: Recommendations should align with user preferences as evaluated via user feedback or online metrics**

From the metrics output the comparison between the models outcome and the actual user preference indicates that users will be satisfied however more tuning of the models needs to be done to increase precision

**3. The system should be scalable to handle new users or movies efficiently**

There is provision to include new datasets for prediction

# Limitations and Future Recommendations

1. Further hyperparameter tuning could have enhanced the models' performance and improved evaluation metrics.
2. The models were computationally intensive, resulting in long training times and occasional runtime crashes due to resource constraints.
3. The dataset had a limited number of movie ratings, which led to data sparsity and impacted model effectiveness.