

AIWR ASSIGNMENT – 2

BASIC RECOMMENDER SYSTEM USING CONTENT AND COLLABORATIVE FILTERING

Team Members

PRAJWAL BS- PES1UG20CS287

PRANAV KALWAD - PES1UG20CS291

PURVIK S NUKAL - PES1UG20CS315

RAHUL RANGANATH - PES1UG20CS316

1. Problem Statement

Implement a Recommender System on Movie Dataset using Content and Collaborative Filtering

2. Introduction

Recommender systems are a type of information filtering system that predict and recommend items or products that a user might be interested in based on their past behavior, preferences, and similarities to other users.

In this context, we build a recommender system to suggest movies to users based on their past movie preferences and similar users' movie choices.

The 2 main approaches used for building the recommender system include: content-based filtering and collaborative filtering

Content-based filtering uses the characteristics or attributes of the items being recommended to make suggestions

Wherein Collaborative filtering recommends items based on the preferences and actions of similar users.

3. Dataset

Movie Dataset

<https://www.kaggle.com/datasets/shubhammehta21/movie-lens-small-latest-dataset>

Dataset Summary

This dataset (ml-latest-small) describes 5-star rating and free-text tagging activity from MovieLens, a movie recommendation service. It contains 100836 ratings and 3683 tag applications across 9742 movies. These data were created by 610 users between March 29, 1996 and September 24, 2018. This dataset was generated on September 26, 2018.

Users were selected at random for inclusion. All selected users had rated at least 20 movies. No demographic information is included. Each user is represented by an id, and no other information is provided.

Importing necessary python libraries to build the model

```
1 pip install scikit-surprise

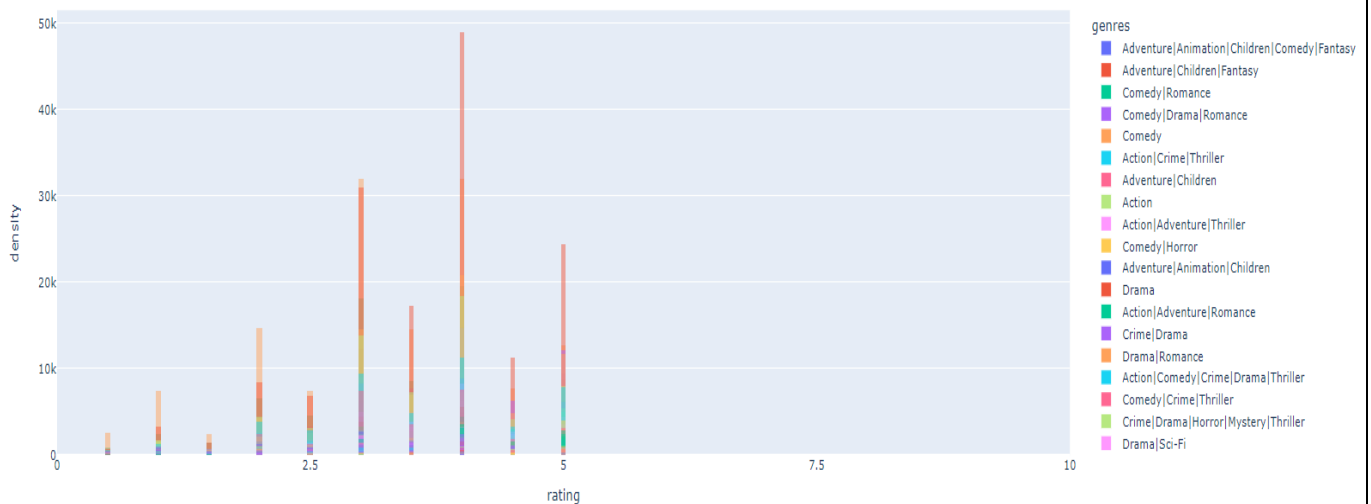
Looking in indexes: https://pypi.org/simple, https://us-python.pkg.dev/colab-wheels/public/simple/
Collecting scikit-surprise
  Downloading scikit-surprise-1.1.3.tar.gz (771 kB)
    772.0/772.0 kB 20.9 MB/s eta 0:00:00
  Preparing metadata (setup.py) ... done
Requirement already satisfied: joblib>=1.0.0 in /usr/local/lib/python3.9/dist-packages (from scikit-surprise) (1.2.0)
Requirement already satisfied: numpy>=1.17.3 in /usr/local/lib/python3.9/dist-packages (from scikit-surprise) (1.22.4)
Requirement already satisfied: scipy>=1.3.2 in /usr/local/lib/python3.9/dist-packages (from scikit-surprise) (1.10.1)
Building wheels for collected packages: scikit-surprise
  Building wheel for scikit-surprise (setup.py) ... done
  Created wheel for scikit-surprise: filename=scikit_surprise-1.1.3-cp39-cp39-linux_x86_64.whl size=3195797 sha256=2dfae1ea2d27286d883b8de86f2b2ebe3865790695ecc3d41b0b3834604707bd
  Stored in directory: /root/.cache/pip/wheels/c6/3a/46/9b17b3512bdf283c6cb84f59929cdd5199d4e754d596d22784
Successfully built scikit-surprise
Installing collected packages: scikit-surprise
Successfully installed scikit-surprise-1.1.3

[3] 1 import pandas as pd
    2 import numpy as np
    3 import matplotlib.pyplot as plt
    4 import seaborn as sns
    5 import plotly.express as px
    6 from collections import defaultdict
    7 from scipy.sparse import csr_matrix
    8 from sklearn.metrics.pairwise import cosine_similarity
    9 from sklearn.feature_extraction.text import TfidfVectorizer
   10 import re
   11 from surprise import SVD, Dataset, Reader
   12
   13
```

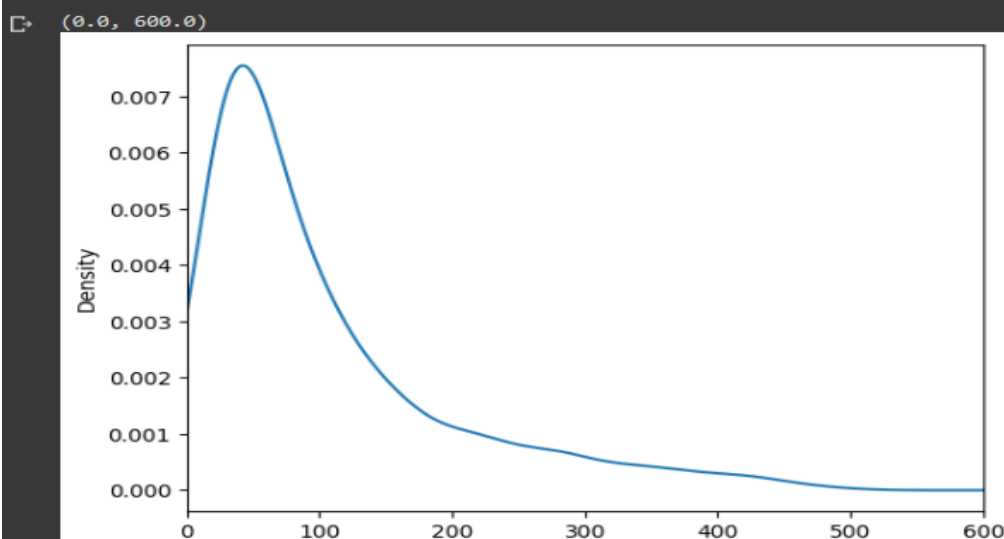
4. EDA

variations of ratings by movie genres

```
1 fig = px.histogram(df.dropna(subset=['rating']), x='rating', color='genres',
2                   histnorm='density', nbins=100, opacity=0.5,
3                   barmode='overlay', range_x=[0, 10])
4 fig.update_layout(xaxis=dict(dtick=2.5))
5
6
7 fig.show()
8
```

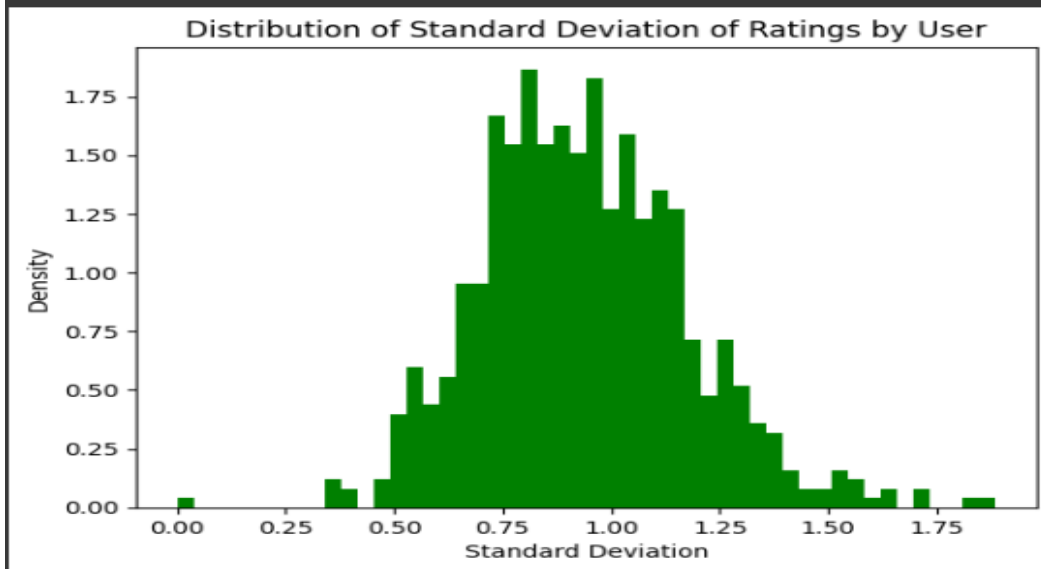


```
1 g = rating.groupby('userId').size().reset_index(name='n')
2 g = g[g['n'] < 500]
3
4 # Density plot using matplotlib
5 fig, ax = plt.subplots()
6 ax = g['n'].plot.density()
7 ax.set_xlim(0, 600)
8
```



Distribution of Standard Deviation of Movie Ratings by User

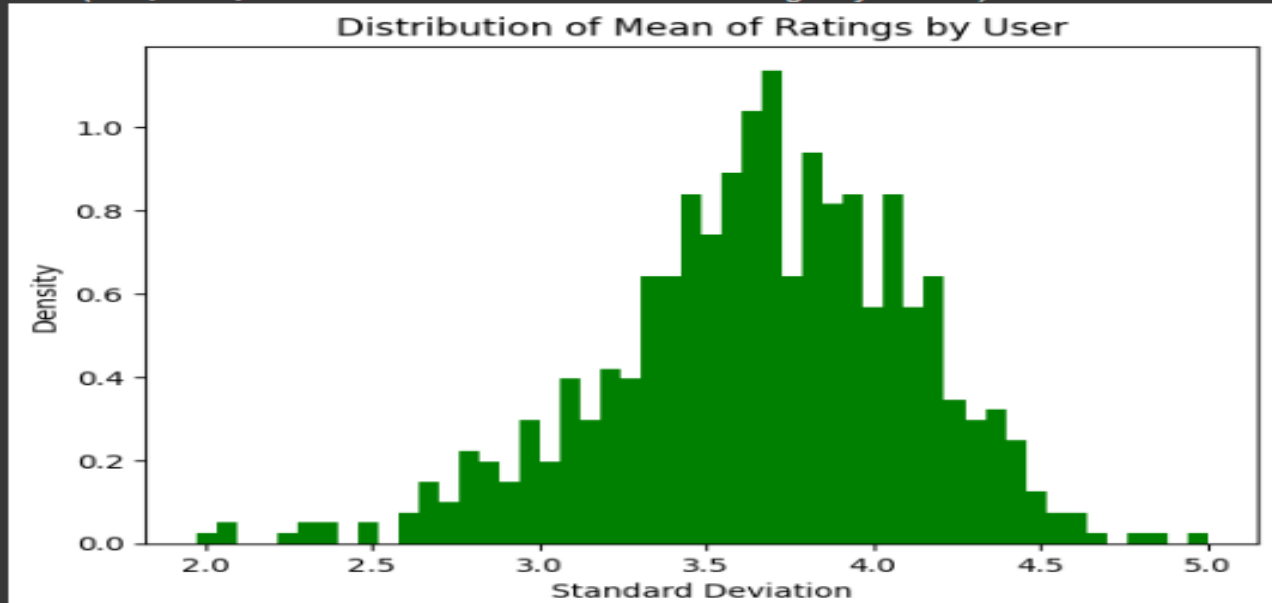
```
1 g = rating.groupby('userId')['rating'].std().reset_index(name='sd')
2 plt.hist(g['sd'], bins=50, density=True,color="green")
3 plt.xlabel('Standard Deviation')
4 plt.ylabel('Density')
5 plt.title('Distribution of Standard Deviation of Ratings by User')
6
7 # Display the plot
8 plt.show()
```



Distribution of Mean of Movie Ratings by User

```
1 g = rating.groupby('userId')['rating'].mean().reset_index(name='m')
2 plt.hist(g['m'], bins=50, density=True,color="green")
3 plt.xlabel('Standard Deviation')
4 plt.ylabel('Density')
5 plt.title('Distribution of Mean of Ratings by User')
```

Text(0.5, 1.0, 'Distribution of Mean of Ratings by User')



```
1 combine_movie_rating = df.dropna(axis = 0, subset = ['title'])
```

```
1 combine_movie_rating
```

	movieId	title	genres	userId	rating	timestamp
0	1	Toy Story (1995)	Adventure Animation Children Comedy Fantasy	2	5.0	859046895
1	1	Toy Story (1995)	Adventure Animation Children Comedy Fantasy	5	4.0	1303501039
2	1	Toy Story (1995)	Adventure Animation Children Comedy Fantasy	8	5.0	858610933
3	1	Toy Story (1995)	Adventure Animation Children Comedy Fantasy	11	4.0	850815810
4	1	Toy Story (1995)	Adventure Animation Children Comedy Fantasy	14	4.0	851766286
...
105334	148238	A Very Murray Christmas (2015)	Comedy	475	3.0	1451213043
105335	148626	The Big Short (2015)	Drama	458	4.0	1452014749
105336	148626	The Big Short (2015)	Drama	576	4.5	1451687664
105337	148626	The Big Short (2015)	Drama	668	4.5	1451148148
105338	149532	Marco Polo: One Hundred Eyes (2015)	(no genres listed)	475	4.0	1451223429

105339 rows × 6 columns

5. Preprocessing

```
1 movie=pd.read_csv('/content/drive/MyDrive/AIWR_A2_DATASET/movies.csv')
2 rating=pd.read_csv('/content/drive/MyDrive/AIWR_A2_DATASET/ratings.csv')
3
```

Double-click (or enter) to edit

```
[5] 1 movie
```


	movieId	title	genres
0	1	Toy Story (1995)	Adventure Animation Children Comedy Fantasy
1	2	Jumanji (1995)	Adventure Children Fantasy
2	3	Grumpier Old Men (1995)	Comedy Romance
3	4	Waiting to Exhale (1995)	Comedy Drama Romance
4	5	Father of the Bride Part II (1995)	Comedy
...
10324	146684	Cosmic Scrat-tastrophe (2015)	Animation Children Comedy
10325	146878	Le Grand Restaurant (1966)	Comedy
10326	148238	A Very Murray Christmas (2015)	Comedy
10327	148626	The Big Short (2015)	Drama
10328	149532	Marco Polo: One Hundred Eyes (2015)	(no genres listed)

10329 rows × 3 columns

```
1 #dropping the null values in movie df
2 #replacing ratings which are -1 to nan in rating column in rating df
3 movie=movie.dropna()
4 rating.rating.replace(-1, np.NaN,inplace=True)
5
```

```
[9] 1 for text in movie['title']:
2     text = re.sub(r'&quot;', '', text)
3     text = re.sub(r'.hack//', '', text)
4     text = re.sub(r'&#039;', '', text)
5     text = re.sub(r'&A&#039;s', '', text)
6     text = re.sub(r'I&#039;', 'I\'' , text)
7     text = re.sub(r'&amp;', 'and', text)
8
```

```
[10] 1 #merging movie dataframe and rating df
2 df=pd.merge(movie,rating,on='movieId',suffixes= ['', '_user'])
3 df = df.rename(columns={'rating_user': 'user_rating'})
4 df.head(10)
```

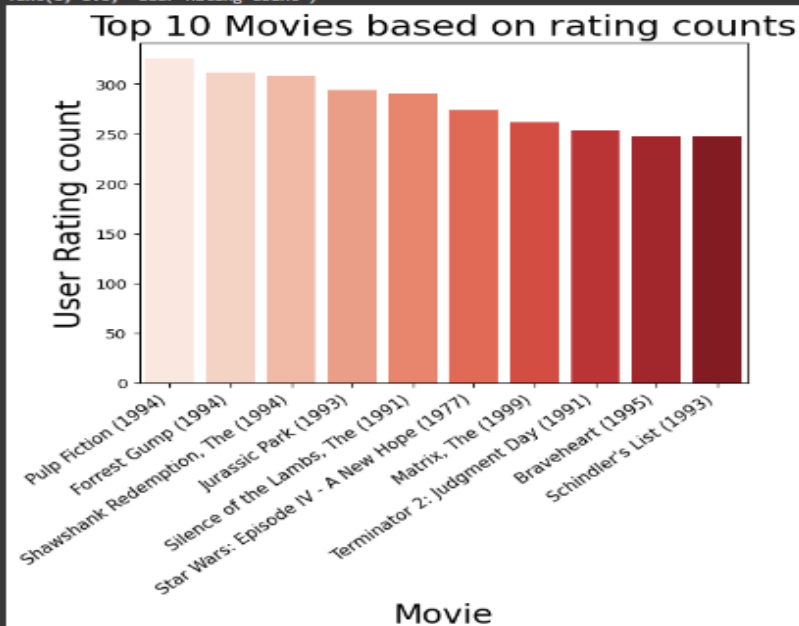
movieId		title	genres	userId	rating	timestamp	
0	1	Toy Story (1995)	Adventure Animation Children Comedy Fantasy	2	5.0	859046895	
1	1	Toy Story (1995)	Adventure Animation Children Comedy Fantasy	5	4.0	1303501039	
2	1	Toy Story (1995)	Adventure Animation Children Comedy Fantasy	8	5.0	858610933	
3	1	Toy Story (1995)	Adventure Animation Children Comedy Fantasy	11	4.0	850815810	
4	1	Toy Story (1995)	Adventure Animation Children Comedy Fantasy	14	4.0	851766286	
5	1	Toy Story (1995)	Adventure Animation Children Comedy Fantasy	17	5.0	1350206819	
6	1	Toy Story (1995)	Adventure Animation Children Comedy Fantasy	28	3.0	884098350	
7	1	Toy Story (1995)	Adventure Animation Children Comedy Fantasy	29	4.0	846942580	
8	1	Toy Story (1995)	Adventure Animation Children Comedy Fantasy	30	4.5	1292690069	
9	1	Toy Story (1995)	Adventure Animation Children Comedy Fantasy	31	4.0	832397496	

Ranking top 10 movies based on movies

Top movies based on rating

```
1 combine_movie_rating = df.dropna(axis=0, subset=['title'])
2
3 movie_ratingCount = (combine_movie_rating
4                       .groupby(by=['title'])['rating']
5                       .count()
6                       .reset_index()
7                       .rename(columns={'rating': 'totalRatingCount'})[['title', 'totalRatingCount']]
8
9
10 top10_movieRating = movie_ratingCount[['title', 'totalRatingCount']].sort_values(by='totalRatingCount', ascending=False).head(10)
11
12 ax = sns.barplot(x="title", y="totalRatingCount", data=top10_movieRating, palette="Reds")
13 ax.set_xticklabels(ax.get_xticklabels(), fontsize=11, rotation=40, ha="right")
14 ax.set_title('Top 10 Movies based on rating counts', fontsize=22)
15 ax.set_xlabel('Movie', fontsize=20)
16 ax.set_ylabel('User Rating count', fontsize=20)
17
```

Text(0, 0.5, 'User Rating count')



6. Content Based Filtering

Content Based recommendation

```
1 movie_features=movie[['genres']].astype('str')
2 movie_features["genres"]
3 movie_features["genres"]=movie_features["genres"].str.split(',').str.join(' ')
4 movie_features
```

	genres
0	Adventure Animation Children Comedy Fantasy
1	Adventure Children Fantasy
2	Comedy Romance
3	Comedy Drama Romance
4	Comedy
...	...
10324	Animation Children Comedy
10325	Comedy
10326	Comedy
10327	Drama
10328	(no genres listed)

10329 rows x 1 columns

```
[26] 1 #creating tf-idf
2 tfidf = TfidfVectorizer(stop_words='english')
3 movie_matrix = tfidf.fit_transform(movie_features.apply(lambda x: ' '.join(x), axis=1))
4 cosine_similarities = cosine_similarity(movie_matrix) #creating cosine similarity
5 #creating a list which contains all movie names
6 movie_names=movie["title"]
7 indices = pd.Series(movie.index, index=movie['title'])
```



```

1 def get_recommendations(movie_name,k):
2     idx = indices[movie_name]
3     # Get the pairwise similarity scores
4     sim_scores = list(enumerate(cosine_similarities[idx]))
5     # Sort the movie based on the similarity scores
6     sim_scores = sorted(sim_scores, key=lambda x: x[1], reverse=True)
7     # Get the scores of the k most similar movie
8     sim_scores = sim_scores[1:k+1]
9     # Get the movie indices
10    movie_indices = [i[0] for i in sim_scores]
11    # Return the top k most similar movie
12    return movie_names.iloc[movie_indices]

1 get_recommendations("Toy Story (1995)",10)

1815                                Antz (1998)
2496                                Toy Story 2 (1999)
2967    Adventures of Rocky and Bullwinkle, The (2000)
3166    Emperor's New Groove, The (2000)
3811                                Monsters, Inc. (2001)
6617    DuckTales: The Movie - Treasure of the Lost La...
6997                                Wild, The (2006)
7382                                Shrek the Third (2007)
7987                                Tale of Despereaux, The (2008)
9215    Asterix and the Vikings (Astérix et les Viking...
Name: title, dtype: object

```

7. Collaborative Filtering

Collaborative Filtering is done by creating a pivot table (which is normalized later)

Cosine similarity is the metric used to compute the similarity

Collaborative Filtering

Collaborative filtering can be done using creating a pivot table, since the data is huge it takes lot of time to create the pivot table

```

[21] 1 df_collab=df[['userId', 'title', 'rating']]
      2 df_collab=df_collab.dropna()
      3 df_collab= df_collab[df_collab.userId <= 20000]
      4

[22] 1 df_collab=df[['userId', 'title', 'rating']]
      2 # df_collab=df_collab.dropna()
      3 df_collab= df_collab[df_collab.userId <= 20000]
      4 #creating pivot table
      5 pivot_table= df_collab.pivot_table(index=['userId'], columns=['title'], values='rating')
      6
      7 #normalizing
      8 normalised_pivot_table = pivot_table.apply(lambda x: (x-np.mean(x))/(np.max(x)-np.min(x)), axis=1)
      9 # Drop all columns containing only zeros representing users who did not rate
     10 normalised_pivot_table.fillna(0, inplace=True)
     11 normalised_pivot_table = normalised_pivot_table.T
     12 normalised_pivot_table = normalised_pivot_table.loc[:, (normalised_pivot_table != 0).any(axis=0)]
     13 #converting pivot table into sparse matrix
     14 sparse_matrix = csr_matrix(normalised_pivot_table.values)
     15 #getting item_similarity and user similarity
     16 item_similarity = cosine_similarity(sparse_matrix)
     17 user_similarity = cosine_similarity(sparse_matrix.T)
     18
     19 #make item and user similarity matrices into df
     20 item_sim_df = pd.DataFrame(item_similarity, index = normalised_pivot_table.index, columns = normalised_pivot_table.index)
     21 user_sim_df = pd.DataFrame(user_similarity, index = normalised_pivot_table.columns, columns = normalised_pivot_table.columns)
     22

```

```
[23] 1 def collaborative_filtering_recommendation(movie_name,top_n):
      2     count = 1
      3     print('Similar shows to {} include:\n'.format(movie_name))
      4     for item in item_sim_df.sort_values(by = movie_name, ascending = False).index[1:top_n+1]:
      5         print('No. {}: {}'.format(count, item))
      6         count +=1

[24] 1 collaborative_filtering_recommendation('The Big Short (2015)',10)
      2
```

Similar shows to The Big Short (2015) include:

No. 1: Calvary (2014)
 No. 2: You Don't Know Jack (2010)
 No. 3: Arbitrage (2012)
 No. 4: Rush (2013)
 No. 5: Sea Inside, The (Mar adentro) (2004)
 No. 6: Stranger, The (1946)
 No. 7: Duellists, The (1977)
 No. 8: Somebody Up There Likes Me (1956)
 No. 9: Naked City, The (1948)
 No. 10: Sleepless Night (Nuit blanche) (2011)

8. Results

Movie Recommendations for a given user id

```
1 algo = SVD()
2 algo.fit(trainset)
3
4 # Ask for user ID
5 user_id = input("Enter user ID:")
6
7 # Get the top 10 recommendations for the user
8 user_items = ratings[ratings['userId'] == int(user_id)]['movieId']
9 user_items = list(set(user_items))
10 other_items = [i for i in ratings['movieId'].unique() if i not in user_items]
11 random.shuffle(other_items)
12 test_items = other_items[:100]
13 test_data = [(int(user_id), i, 4) for i in test_items]
14 predictions = algo.test(test_data)
15 predictions = sorted(predictions, key=lambda x: x.est, reverse=True)[:10]
16
17 movies = pd.read_csv('/content/drive/MyDrive/AIWR_A2_DATASET/movies.csv')
18 print("Top 10 movie recommendations for user", user_id)
19 for prediction in predictions:
20     movie_id = prediction.iid
21     movie_title = movies[movies['movieId'] == movie_id]['title'].values[0]
22     print(movie_id, movie_title)
```

Enter user ID:1
 Top 10 movie recommendations for user 1
 541 Blade Runner (1982)
 46578 Little Miss Sunshine (2006)
 1283 High Noon (1952)
 111 Taxi Driver (1976)
 353 Crow, The (1994)
 82459 True Grit (2010)
 2791 Airplane! (1980)
 67255 Girl with the Dragon Tattoo, The (Män som hatar kvinnor) (2009)
 161582 Hell or High Water (2016)
 5377 About a Boy (2002)

The evaluation metrics used include MAE, RMSE, Precision and Recall

Evaluation results for Content-based recommendation

```
MAP value: 0.9830508474576272
```

```
Precision for Pulp Fiction (1994): 0.1  
Recall for Pulp Fiction (1994): 0.004310344827586207
```

Evaluation results for Collaborative recommendation

```
Computing the cosine similarity matrix...  
Done computing similarity matrix.  
MAE: 0.7496  
RMSE: 0.9743  
Computing the cosine similarity matrix...  
Done computing similarity matrix.  
MAE: 0.7476  
RMSE: 0.9699  
Computing the cosine similarity matrix...  
Done computing similarity matrix.  
MAE: 0.7476  
RMSE: 0.9696  
Computing the cosine similarity matrix...  
Done computing similarity matrix.  
MAE: 0.7565  
RMSE: 0.9817  
Computing the cosine similarity matrix...  
Done computing similarity matrix.  
MAE: 0.7421  
RMSE: 0.9669  
MAE: 0.6734  
RMSE: 0.8774  
MAE: 0.6705  
RMSE: 0.8740  
MAE: 0.6683  
RMSE: 0.8680  
MAE: 0.6733  
RMSE: 0.8753  
MAE: 0.6678  
RMSE: 0.8695  
KNN: MAE = 0.749, RMSE = 0.972  
SVD: MAE = 0.671, RMSE = 0.873
```

```
KNN-based Collaborative Filtering:  
Precision@10 = 0.75  
Recall@10 = 0.52  
SVD-based Collaborative Filtering:  
Precision@10 = 0.65  
Recall@10 = 0.62
```

9. Conclusion

In conclusion, the inclusion of content-based and collaborative filtering methods in building a movie recommender system has given significantly accurate results.

We also realize the potential of these 2 filtering models that provide the user with an effective and personalized recommendation of movies and can infer that several factors such as the quality and quantity of data, the choice of similarity metrics, and the user-item interaction patterns affect the performance of a recommender system.

These factors should be carefully considered when designing and optimizing recommender systems.