

## Lab-6

### MOVIE RECOMMENDATION SYSTEM

Koduri Gokul

19BCD7006

Rating Dataset:

1	userId	movieId	rating	timestamp	
2	1	16	4	1.218E+09	
3	1	24	1.5	1.218E+09	
4	1	32	4	1.218E+09	
5	1	47	4	1.218E+09	
6	1	50	4	1.218E+09	
7	1	110	4	1.218E+09	
8	1	150	3	1.218E+09	
9	1	161	4	1.218E+09	
10	1	165	3	1.218E+09	
11	1	204	0.5	1.218E+09	
12	1	223	4	1.218E+09	
13	1	256	0.5	1.218E+09	
14	1	260	4.5	1.218E+09	
15	1	261	1.5	1.218E+09	
16	1	277	0.5	1.218E+09	
17	1	296	4	1.218E+09	
18	1	318	4	1.218E+09	
19	1	349	4.5	1.218E+09	
20	1	356	3	1.218E+09	
21	1	377	2.5	1.218E+09	
22	1	380	3	1.218E+09	
23	1	457	4	1.218E+09	

## Movie Dataset:

1	movieId	title	genres				
2	1	Toy Story (1995)	Adventure   Animation   Children   Comedy   Fantasy				
3	2	Jumanji (1995)	Adventure   Children   Fantasy				
4	3	Grumpier Old Men	Comedy   Romance				
5	4	Waiting to Exhale	Comedy   Drama   Romance				
6	5	Father of the Bride	Comedy				
7	6	Heat (1995)	Action   Crime   Thriller				
8	7	Sabrina (1995)	Comedy   Romance				
9	8	Tom and Huck	Adventure   Children				
10	9	Sudden Death	Action				
11	10	GoldenEye (1995)	Action   Adventure   Thriller				
12	11	American Pie	Comedy   Drama   Romance				
13	12	Dracula: Dead and Loving It	Comedy   Horror				
14	13	Balto (1995)	Adventure   Animation   Children				
15	14	Nixon (1995)	Drama				
16	15	Cutthroat Island	Action   Adventure   Romance				
17	16	Casino (1995)	Crime   Drama				
18	17	Sense and Sensibility	Drama   Romance				
19	18	Four Rooms	Comedy				
20	19	Ace Ventura: Pet Detective	Comedy				
21	20	Money Train	Action   Comedy   Crime   Drama   Thriller				
22	21	Get Shorty (1995)	Comedy   Crime   Thriller				
23	22	Copycat (1995)	Crime   Drama   Horror   Mystery   Thriller				
24	23	Assassins (1999)	Action   Crime   Thriller				
25	24	Powder (1995)	Drama   Sci-Fi				
26	25	Leaving Las Vegas	Drama   Romance				
27	26	Othello (1995)	Drama				
28	27	Now and Then	Children   Drama				
29	28	Persuasion (1995)	Drama   Romance				
30	29	City of Lost Children	Adventure   Drama   Fantasy   Mystery   Sci-Fi				

```
import pandas as pd
import numpy as np
from scipy.sparse import csr_matrix
from sklearn.neighbors import NearestNeighbors
import matplotlib.pyplot as plt
import seaborn as sns
```

```
movies = pd.read_csv("movies.csv")
ratings = pd.read_csv("ratings.csv")
```

```
movies.head()
```

	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



```
ratings.head()
```



	userId	movieId	rating	timestamp
0	1	1	4.0	964982703
1	1	3	4.0	964981247
2	1	6	4.0	964982224
3	1	47	5.0	964983815
4	1	50	5.0	964982931



```
final_dataset = ratings.pivot(index='movieId',columns='userId',values='rating')
final_dataset.head()
```

userId	1	2	3	4	5	6	7	8	9	10	...	601	602	603	604	605	606	607	608	609	610
movieId																					
1	4.0	NaN	NaN	NaN	4.0	NaN	4.5	NaN	NaN	NaN	...	4.0	NaN	4.0	3.0	4.0	2.5	4.0	2.5	3.0	5.0
2	NaN	NaN	NaN	NaN	NaN	4.0	NaN	4.0	NaN	NaN	...	NaN	4.0	NaN	5.0	3.5	NaN	NaN	2.0	NaN	NaN
3	4.0	NaN	NaN	NaN	NaN	5.0	NaN	NaN	NaN	NaN	...	NaN	NaN	NaN	NaN	NaN	NaN	NaN	2.0	NaN	NaN
4	NaN	NaN	NaN	NaN	NaN	3.0	NaN	NaN	NaN	NaN	...	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN	NaN
5	NaN	NaN	NaN	NaN	NaN	5.0	NaN	NaN	NaN	NaN	...	NaN	NaN	NaN	3.0	NaN	NaN	NaN	NaN	NaN	NaN

5 rows × 610 columns

```
final_dataset.fillna(0,inplace=True)
final_dataset.head()
```

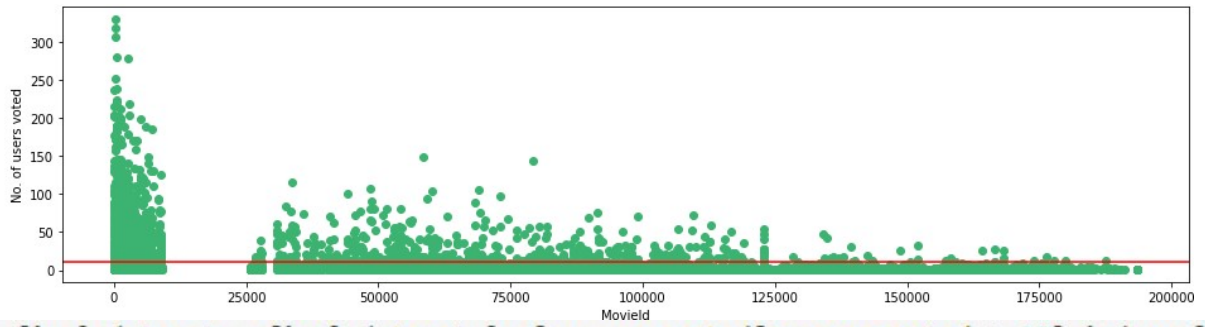
userId	1	2	3	4	5	6	7	8	9	10	...	601	602	603	604	605	606	607	608	609	610
movieId																					
1	4.0	0.0	0.0	0.0	4.0	0.0	4.5	0.0	0.0	0.0	...	4.0	0.0	4.0	3.0	4.0	2.5	4.0	2.5	3.0	5.0
2	0.0	0.0	0.0	0.0	0.0	4.0	0.0	4.0	0.0	0.0	...	0.0	4.0	0.0	5.0	3.5	0.0	0.0	2.0	0.0	0.0
3	4.0	0.0	0.0	0.0	0.0	5.0	0.0	0.0	0.0	0.0	...	0.0	0.0	0.0	0.0	0.0	0.0	0.0	2.0	0.0	0.0
4	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
5	0.0	0.0	0.0	0.0	0.0	5.0	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

5 rows × 610 columns

```
[9] no_user_voted = ratings.groupby('movieId')['rating'].agg('count')
no_movies_voted = ratings.groupby('userId')['rating'].agg('count')
```

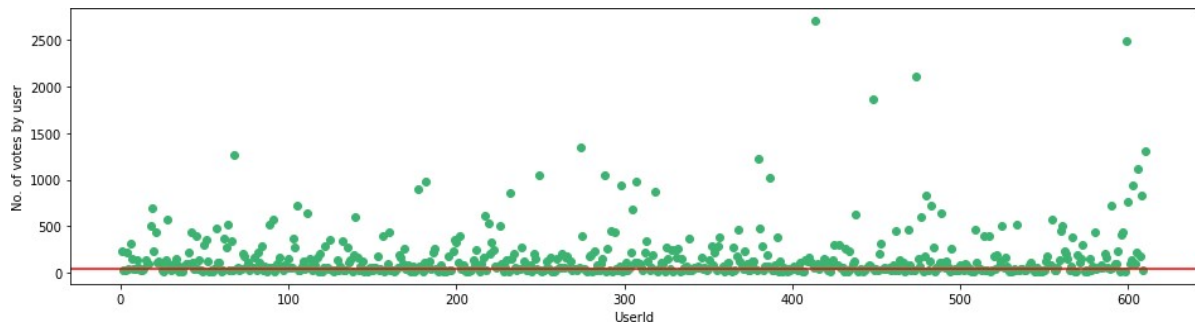
```
f,ax = plt.subplots(1,1,figsize=(16,4))
# ratings['rating'].plot(kind='hist')
plt.scatter(no_user_voted.index,no_user_voted,color='mediumseagreen')
plt.axhline(y=10,color='r')
plt.xlabel('MovieId')
plt.ylabel('No. of users voted')
plt.show()
```

*Removing Noise from Data:*



```
[11] final_dataset = final_dataset.loc[no_user_voted[no_user_voted > 10].index,:]
```

```
f,ax = plt.subplots(1,1,figsize=(16,4))
plt.scatter(no_movies_voted.index,no_movies_voted,color='mediumseagreen')
plt.axhline(y=50,color='r')
plt.xlabel('UserId')
plt.ylabel('No. of votes by user')
plt.show()
```





```

final_dataset=final_dataset.loc[:,no_movies_voted[no_movies_voted > 50].index]
final_dataset

```

	userId	1	4	6	7	10	11	15	16	17	18	...	600	601	602	603	604	605	606	607	608	610
movieId																						
1	4.0	0.0	0.0	4.5	0.0	0.0	2.5	0.0	4.5	3.5	...	2.5	4.0	0.0	4.0	3.0	4.0	2.5	4.0	2.5	5.0	
2	0.0	0.0	4.0	0.0	0.0	0.0	0.0	0.0	0.0	3.0	...	4.0	0.0	4.0	0.0	5.0	3.5	0.0	0.0	2.0	0.0	
3	4.0	0.0	5.0	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	2.0	0.0	
5	0.0	0.0	5.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	2.5	0.0	0.0	0.0	3.0	0.0	0.0	0.0	0.0	0.0	
6	4.0	0.0	4.0	0.0	0.0	5.0	0.0	0.0	0.0	4.0	...	0.0	0.0	3.0	4.0	3.0	0.0	0.0	0.0	0.0	5.0	
...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	...	
174055	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	
176371	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	
177765	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	...	0.0	4.5	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	
179819	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	0.0	0.0	0.0	
187593	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	0.0	0.0	0.0	

2121 rows x 378 columns

```

sample = np.array([[0,0,3,0,0],[4,0,0,0,2],[0,0,0,0,1]])
sparsity = 1.0 - ( np.count_nonzero(sample) / float(sample.size) )
print(sparsity)

```

0.7333333333333334

```

csr_sample = csr_matrix(sample)
print(csr_sample)

```

```

(0, 2)      3
(1, 0)      4
(1, 4)      2
(2, 4)      1

```

```

[16] csr_data = csr_matrix(final_dataset.values)
final_dataset.reset_index(inplace=True)

```

*Removing Sparsity:*

**Making the movie recommendation system model:**

```

knn = NearestNeighbors(metric='cosine', algorithm='brute', n_neighbors=20, n_jobs=-1)
knn.fit(csr_data)

```

NearestNeighbors(algorithm='brute', metric='cosine', n\_jobs=-1, n\_neighbors=20)

## Making the recommendation function:

```
def get_movie_recommendation(movie_name):
    n_movies_to_reccomend = 10
    movie_list = movies[movies['title'].str.contains(movie_name)]
    if len(movie_list):
        movie_idx= movie_list.iloc[0]['movieId']
        movie_idx = final_dataset[final_dataset['movieId'] == movie_idx].index[0]
        distances , indices = knn.kneighbors(csr_data[movie_idx],n_neighbors=n_movies_to_reccomend+1)
        rec_movie_indices = sorted(list(zip(indices.squeeze().tolist(),distances.squeeze().tolist())),key=lambda x: x[1])[:0:-1])
        recommend_frame = []
        for val in rec_movie_indices:
            movie_idx = final_dataset.iloc[val[0]]['movieId']
            idx = movies[movies['movieId'] == movie_idx].index
            recommend_frame.append({'Title':movies.iloc[idx]['title'].values[0],'Distance':val[1]})
        df = pd.DataFrame(recommend_frame,index=range(1,n_movies_to_reccomend+1))
        return df
    else:
        return "No movies found. Please check your input"
```

## Final Outputs:

```
get_movie_recommendation('Iron Man')
```

	Title	Distance
1	Up (2009)	0.368857
2	Guardians of the Galaxy (2014)	0.368758
3	Watchmen (2009)	0.368558
4	Star Trek (2009)	0.366029
5	Batman Begins (2005)	0.362759
6	Avatar (2009)	0.310893
7	Iron Man 2 (2010)	0.307492
8	WALL·E (2008)	0.298138
9	Dark Knight, The (2008)	0.285835
10	Avengers, The (2012)	0.285319



```
get_movie_recommendation('Memento')
```

	Title	Distance
1	American Beauty (1999)	0.389346
2	American History X (1998)	0.388615
3	Pulp Fiction (1994)	0.386235
4	Lord of the Rings: The Return of the King, The...	0.371622
5	Kill Bill: Vol. 1 (2003)	0.350167
6	Lord of the Rings: The Two Towers, The (2002)	0.348358
7	Eternal Sunshine of the Spotless Mind (2004)	0.346196
8	Matrix, The (1999)	0.326215
9	Lord of the Rings: The Fellowship of the Ring,...	0.316777
10	Fight Club (1999)	0.272380

