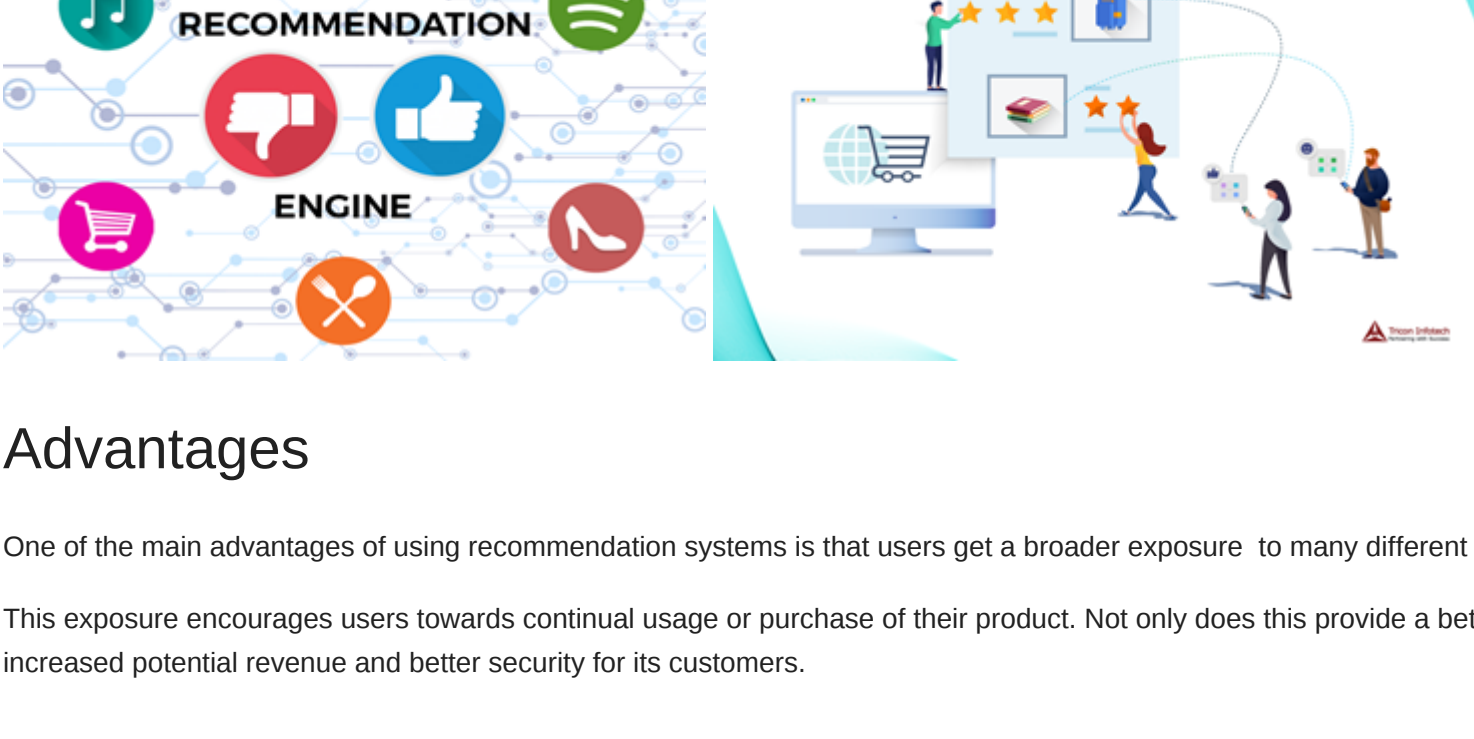


Introduction

What is a Recommendation System?

Simply put a Recommendation System is a filtration program whose prime goal is to predict the "rating" or "preference" of a user towards a domain-specific item or item. In our case, this domain-specific item is a movie, therefore the main focus of our recommendation system is to filter and predict only those movies which a user would prefer given some data about the user him or herself.

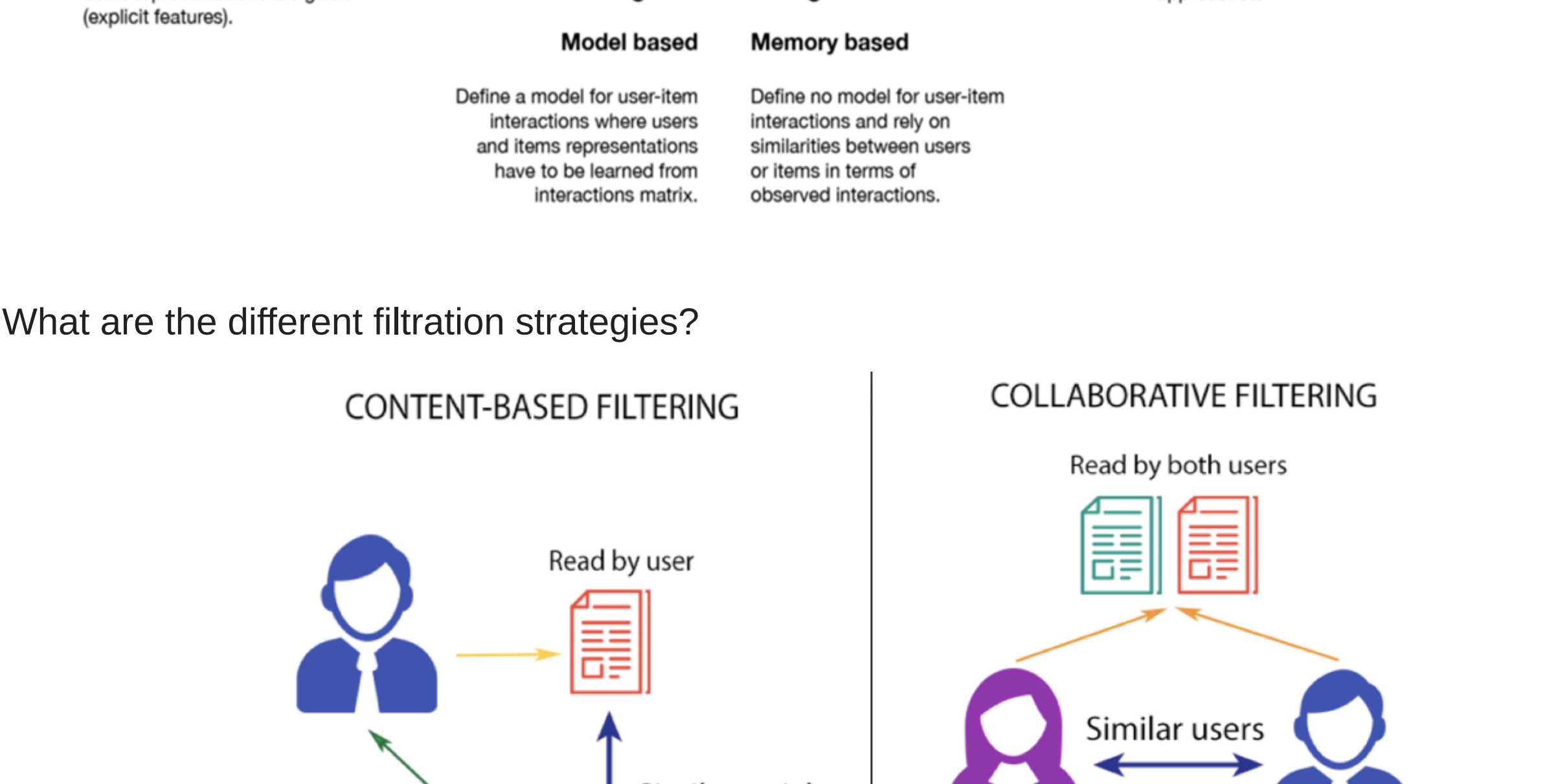


Advantages

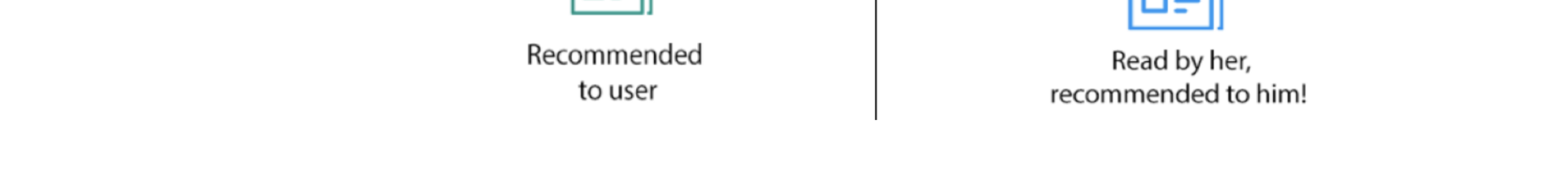
One of the main advantages of using recommendation systems is that users get a broader exposure to many different products they might be interested in.

This exposure encourages users towards continual usage or purchase of their product. Not only does this provide a better experience for the user but it benefits the service provider, as well, with increased potential revenue and better security for its customers.

Types of Recommendation Systems



What are the different filtration strategies?



Content-based Filtering

This filtration strategy is based on the data provided about the items. The algorithm recommends products that are similar to the ones that a user has liked in the past. This similarity (generally cosine similarity) is computed from the data we have about the items as well as the user's past preferences. For example, if a user likes movies such as "The Prestige" then we can recommend him the movies of "Christian Bale" or movies with the genre "Thriller" or maybe even movies directed by "Christopher Nolan".

Disadvantages

- > Different products do not get much exposure to the user.
- > Businesses cannot be expanded as the user does not try different types of products.

Colabarative Filtering

This filtration strategy is based on the combination of the user's behavior and comparing and contrasting that with other users' behavior in the database. The history of all users plays an important role in this algorithm. The main difference between content-based filtering and collaborative filtering is that in the latter, the interaction of all users with the items influences the recommendation algorithm while for content-based filtering only the concerned user's data is taken into account.

There are multiple ways to implement collaborative filtering but the main concept to be grasped is that in collaborative filtering multiple user's data influences the outcome of the recommendation, and doesn't depend on only one user's data for modeling.

There are 2 types of collaborative filtering algorithms:

User-based Collaborative filtering The basic idea here is to find users that have similar past preference patterns as the user 'A' has had and then recommending him or her items liked by those similar users which 'A' has not encountered yet. This is achieved by making a matrix of items each user has rated/viewed/liked/clicked depending upon the task at hand, and then computing the similarity score between the users and finally recommending items that the concerned user isn't aware of but users similar to him/her are and liked it.

Disadvantages

People are fickle-minded i.e their taste change from time to time and as this algorithm is based on user similarity it may pick up initial similarity patterns between 2 users who after a while may have completely different preferences. There are many more users than items therefore it becomes very difficult to maintain such large matrices and therefore needs to be recomputed very regularly. This algorithm is very susceptible to shilling attacks where fake users profiles consisting of biased preference patterns are used to manipulate key decisions.

Item-based Collaborative Filtering

The concept in this case is to find similar movies instead of similar users and then recommending similar movies to that 'A' has had in his/her past preferences. This is executed by finding every pair of items that were rated/viewed/liked/clicked by the same user, then measuring the similarity of those rated/viewed/liked/clicked across all user who rated/viewed/liked/clicked both, and finally recommending them based on similarity scores.

Here, for example, we take 2 movies 'A' and 'B' and check their ratings by all users who have rated both the movies and based on the similarity of these ratings, and based on this rating similarity by users who have rated both we find similar movies. So if most common users have rated 'A' and 'B' both similarly and it is highly probable that 'A' and 'B' are similar, therefore if someone has watched and liked 'A' they should be recommended 'B' and vice versa.

Advantages over User-based Collaborative Filtering

- > Unlike people's taste, movies don't change.
- > There are usually a lot fewer items than people, therefore easier to maintain and compute the matrices.
- > Shilling attacks are much harder because items cannot be faked.

Let's start coding up our own Movie recommendation system

In this implementation, when the user searches for a movie we will recommend the top 10 similar movies using our movie recommendation system. We will be using an **item-based collaborative filtering** algorithm for our purpose. The dataset used in this demonstration is the movielens-small dataset.

```
In [2]: 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_excel(r'H:\Users\BHOAKARKAR\Downloads\Study Material\My Stuff\Courses\Building Recommendation Systems In Python\movies_exe.xlsx')
ratings = pd.read_csv(r'H:\Users\BHOAKARKAR\Downloads\Study Material\My Stuff\Courses\Building Recommendation Systems In Python\user_ratings.csv')
```

```
In [3]: # First 5 Rows
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

Movie dataset has -

movieId – once the recommendation is done, we get a list of all similar movieId and get the title for each movie from this dataset.

genres – which is not required for this filtering approach.

```
In [4]: # Top 5 Rows of Ratings
ratings.head()
```

	userid	movieid	rating	timestamp	title	genres
0	1	1	4.0	964982703	Toy Story (1995)	Adventure Animation Children Comedy Fantasy
1	5	1	4.0	847434962	Toy Story (1995)	Adventure Animation Children Comedy Fantasy
2	7	1	4.5	1106635946	Toy Story (1995)	Adventure Animation Children Comedy Fantasy
3	15	1	2.5	1510577970	Toy Story (1995)	Adventure Animation Children Comedy Fantasy
4	17	1	4.5	1305696483	Toy Story (1995)	Adventure Animation Children Comedy Fantasy

Ratings dataset has-

userId – unique for each user.

movieId – using this feature, we take the title of the movie from the movies dataset.

rating – Ratings given by each user to all the movies using this we are going to predict the top 10 similar movies.

Here, we can see that user1d 1 has watched movieId 1 & 3 and rated both of them 4.0 but has not rated movieId 2 at all. This interpretation is harder to extract from this dataframe. Therefore, to make things easier to understand and work with, we are going to make a new dataframe where each column would represent each unique userId and each row represents each unique movieId.

```
In [5]: final_dataset = ratings.pivot(index='movieId', columns='userId', values='rating')
final_dataset.head()
```

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 rows × 610 columns

Now, it's much easier to interpret that user1d 1 has rated movieId 1 & 3 4.0 but has not rated movieId 2 at all (therefore they are represented as **NaN**) and therefore their rating data is missing.

Let's fix this and impute NaN with 0 to make things understandable for the algorithm and also making the data more eye-soothing.

```
In [6]: final_dataset.fillna(0,inplace=True)
final_dataset.head()
```

To qualify a user, a minimum of 50 movies should have voted by the user.

Let's visualize how these filters look like

Aggregating the number of users who voted and the number of movies that were voted.

5 rows × 610 columns

Removing Noise from the data

In the real-world, ratings are very **sparse** and data points are mostly collected from very popular movies and highly engaged users. We wouldn't want movies that were rated by a small number of users because it's **not credible** enough. Similarly, users who have rated only **a handful of movies** should also not be taken into account.

So with all that taken into account and some trial and error experiments, we will reduce the noise by adding some filters for the final dataset.

To qualify a movie, a minimum of **10** users should have voted a movie.

To qualify a user, a minimum of **50** movies should have voted by the user.

Let's visualize how these filters look like

Aggregating the number of users who voted and the number of movies that were voted.

```
In [7]: no_user_voted = ratings.groupby('movieId')['rating'].agg('count')
no_movies_voted = ratings.groupby('userId')['rating'].agg('count')
```

Let's visualize the number of users who voted with our threshold of 10.

```
In [8]: f,ax = plt.subplots(1,1,figsize=(16,4))
# ratings['rating'].plot(kind='hist')
plt.scatter(no_user_voted.index,no_movies_voted,color='mediumseagreen')
plt.axhline(y=10,color='r')
plt.xlabel('MovieId')
plt.ylabel('No. of users voted')
plt.show()
```



Making the necessary modifications as per the threshold set.

```
In [9]: final_dataset = final_dataset.loc[no_user_voted[no_user_voted > 10].index,:]
```

Let's visualize the number of votes by each user with our threshold of 50.

```
In [10]: 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()
```



Making the necessary modifications as per the threshold set.

```
In [11]: final_dataset=final_dataset.loc[:,no_movies_voted[no_movies_voted > 50].index]
final_dataset
```

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	0.0
--------	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----	-----

2121 rows x 378 columns

Removing sparsity

174055 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

176371 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

177765 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 0.0 ... 0.0 4.5 0.0 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 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 0.0 0.0

2121 rows × 378 columns

Removing sparsity

Our final dataset has dimensions of **2121 × 378** where most of the values are sparse. We are using only a small dataset but for the original large dataset of movie lens which has more than **100000** features, our system may run out of computational resources when that is feed to the model. To reduce the sparsity we use the **csr_matrix** function from the **scipy** library.

```
In [12]: csr_data = csr_matrix(final_dataset.values)
final_dataset.reset_index(inplace=True)
```

Making the movie recommendation system model

The working principle is very simple. We first check if the movie name input is in the database and if it is we use our recommendation system to **find similar movies** and sort them based on their similarity distance and output only the **top 10** movies with their distances from the input movie.

We will be using the **KNN algorithm** to compute similarity with **cosine** distance metric which is **very fast** and more preferable than **pearson coefficient**.

```
In [15]: knn = NearestNeighbors(metric='cosine', algorithm='brute', n_neighbors=20, n_jobs=-1)
knn.fit(csr_data)
```

```
Out[15]: NearestNeighbors(algorithm='brute', metric='cosine', n_jobs=-1, n_neighbors=20)
```

```
In [16]: 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(zip(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"
```

```
In [17]: # Lets get recommendation for Iron Man Movie
get_movie_recommendation('Memento')
```

	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 (The (2008)	0.298138
9	Dark Knight, The (2008)	0.285835
10	Avengers, The (2012)	0.285319

I personally think the results are pretty good. All the movies at the top are superhero or animation movies which are ideal for kids as is the input movie "Iron Man".

```
In [18]: # Let's try another one :
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

All the movies in the top 10 are serious and mindful movies just like "Memento" itself, therefore I think the result, in this case, is also good.

```
In [ ]:
```