**DWDM MINI PROJECT REPORT**

# Project Title: Movie Recommendation and Rating Predictor Using K Nearest Neighbors

*Submitted by*

| Sl# | Full Name | Reg# | Roll Number | Lab Batch |
|---|---|---|---|---|
| 1 | Shreyansh K Gupta | 190911178 | 47 | IT B - 2 |
| 2 | Sumantra Sen | 190911192 | 50 | IT B - 2 |
| 3 | Tarunvir Singh Panesar | 190911196 | 52 | IT B - 2 |

*In partial fulfillment for the award of degree of*

**B. Tech**

**IN**

**INFORMATION TECHNOLOGY**

**MANIPAL INSTITUTE OF TECHNOLOGY**
MANIPAL
*A Constituent Institution of Manipal University*

**Department of Information & Communication Technology**

**May 2022**

**Abstract:**

A recommendation system attempts to predict a user's preference or rating for an item. They are becoming increasingly important as people seek convenience and are always looking for products/services that are best suited to them. As a result, recommendation systems are important because they assist them in making the best decisions. The main goal of this mini project is to provide a personalized experience for each user by showing movies based on the movie of their choice and its imdb ratings.
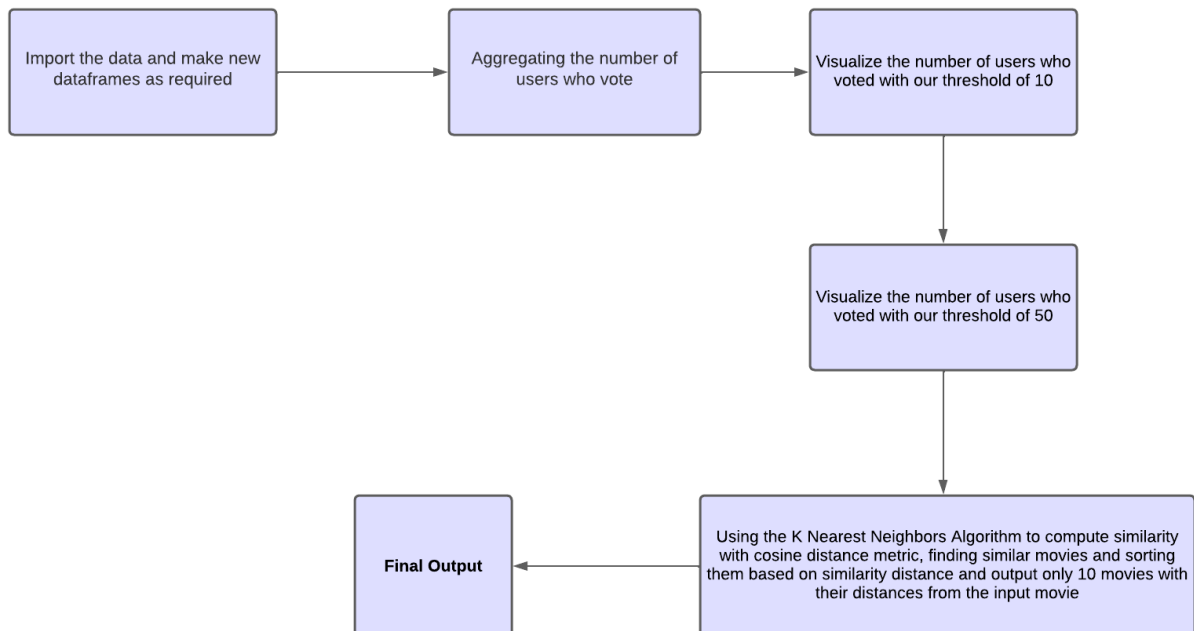
**Introduction:**

This recommender system uses content-based filtering by implementing the K-Nearest Neighbors algorithm. This filtering method is typically based on collecting and analyzing information about a movie's various details such as keywords, actors, genres, and predicting what the user will like based on the movie's similarity to other movies. The main advantage of the content-based filtering approach is that the model doesn't need any data about other users since the recommendations are specific to this user. This makes it easier to scale to many users, so it can accurately recommend complex elements such as movies. Approaches such as collaborative filter algorithms face issues such as scalability, sparsity, and cold start, which can be avoided to some extent in the proposed framework.

**Literature Survey:**

o   R. Ahuja, A. Solanki and A. Nayyar, "Movie Recommender System Using K-Means Clustering AND K-Nearest Neighbor," , in this paper we mainly referred it to get an idea as to how to approach a recommendation system and figure out the different steps to the process. They have used K-nearest-neighbors in both content based as well as collaborative filtering. Based on this we decided to go with content based as it was more independent of other users, and only needed one main dataset. They have also used k-means clustering to classify the different movies, which we have avoided, as we felt it was unnecessary for our expected result.

o   G Geetha1, M Safa1, C Fancy1 and D Saranya2 G "A Hybrid Approach using Collaborative filtering and Content based Filtering for Recommender System". In this paper, similar to the previous reference, it goes about implementing the problem using a hybrid approach, using both filtering methods. We have explored only the content based filtering aspect of this paper, but this research planned on further perfecting the obtained results by using collaborative filtering as well.

## Methodology



## Import the python libraries:



For this analysis, we are importing the following python library files:

1. Pandas: This is a software library file written for the Python Programming Language for data manipulation and analysis, it offers data structures and operations for manipulating numerical tables and time series.
2. NumPy: This is a software library file written for the Python Programming Language to provide support for large multidimensional arrays and matrices, along with large collection of high-level mathematical functions to operate on these arrays.
3. SciPy: This Python library file contains modules for linear algebra, integration, interpolation, special functions, used for scientific analysis.
4. Seaborn and Matplotlib: Matplotlib is a plotting library for Python, and Seaborn is a library used for making statistical graphs in python, which builds on top of matplotlib.

**Making a new data frame:** To make sure that each column represents a unique userId and each row represents a unique movieId

## Aggregating the number of users who voted



## Visualize the number of users who voted with our threshold of 10

**Visualize the number of users who voted with our threshold of 50**



Visualize the number of votes by each user with our threshold of 50.

```python
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()
```

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

**Reducing sparsity by using the csr_matrix function from the SciPy library:**



To reduce the sparsity we use the csr_matrix function from the scipy library.

```python
[11] 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)
```
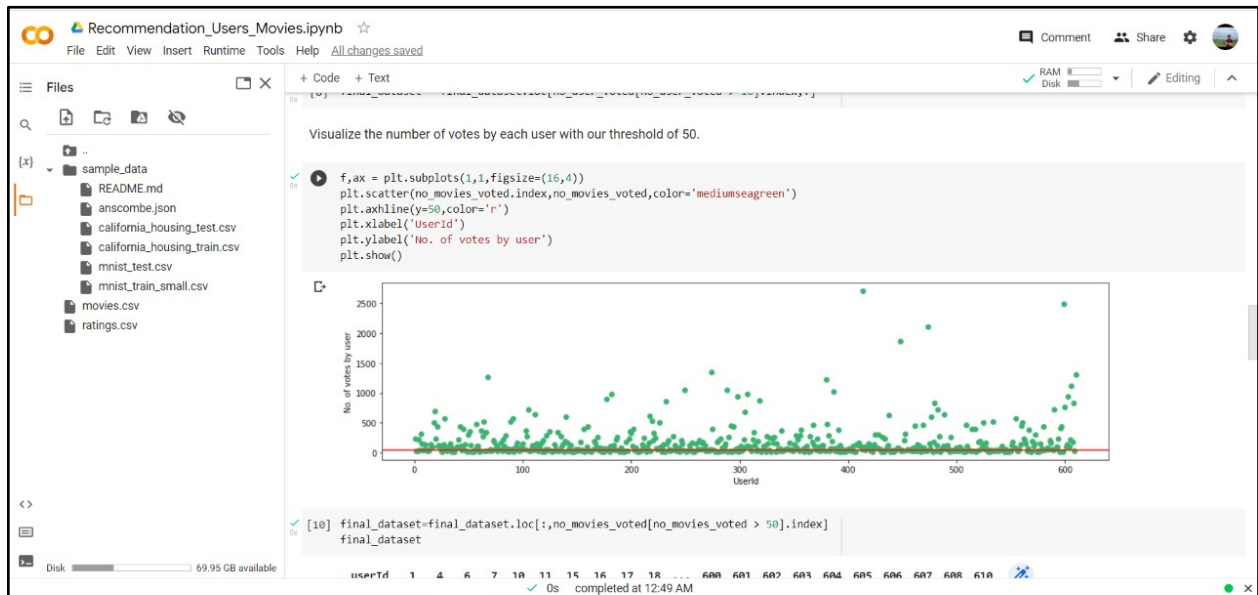
0.7333333333333334

```python
csr_sample = csr_matrix(sample)
print(csr_sample)
```

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

```python
[13] csr_data = csr_matrix(final_dataset.values)
     final_dataset.reset_index(inplace=True)
```

We will be using the knn algorithm to compute similarity with cosine distance metric which is very fast

```python
[17] 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)
```

**Using the K Nearest Neighbors Algorithm to compute similarity with cosine distance metric, finding similar movies and sorting them based on similarity distance and output only 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
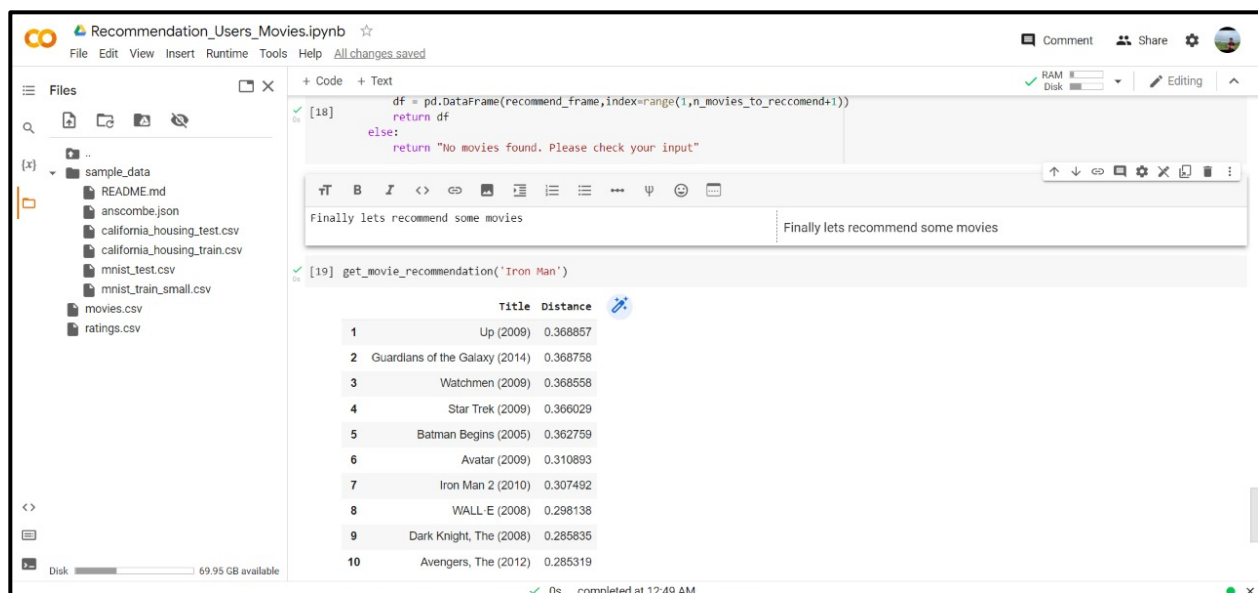
```
[17] knn = NearestNeighbors(metric='cosine', algorithm='brute', n_neighbors=20, n_jobs=-1)
     knn.fit(csr_data)

     NearestNeighbors(algorithm='brute', metric='cosine', n_jobs=-1, n_neighbors=20)
```

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.

```
[18] 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 Output**

Finally lets recommend some movies

```
[19] 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 |

**Conclusions:**

The project explores the workings of recommendation systems and touches upon the topic of content-based filtering and K nearest neighbors. It gives us insight into how the various systems in our life make tasks much easier for us. We also can see that these systems can be surprisingly accurate without training. In conclusion, there are various categories of systems that can be implemented in different ways using varied algorithms, all to make tasks easier and life better.

**References:**

o  G Geetha1, M Safa1, C Fancy1 and D Saranya2 G *"A Hybrid Approach using Collaborative filtering and Content based Filtering for Recommender System"* Geetha et al 2018 J. Phys.: Conf. Ser. 1000 012101

o  R. Ahuja, A. Solanki and A. Nayyar, "Movie Recommender System Using K-Means Clustering AND K-Nearest Neighbor," *2019 9th International Conference on Cloud Computing, Data Science & Engineering (Confluence)*, 2019, pp. 263-268, doi: 10.1109/CONFLUENCE.2019.8776969.

o  B. Chikhaoui, M. Chiazzaro and S. Wang, "An Improved Hybrid Recommender System By Combining Predictions, " IEEE Workshops of International Conference on Advanced Information Networking and Applications, pp. 644-649, 2011.