

# NETFLIX MOVIE RECOMMENDATION SYSTEM

Netflix is an application that keeps growing bigger and faster with its popularity, shows and content. This is a story telling through its data along with a content-based recommendation system and a wide range of different graphs and visuals

## **Problem Statement :**

Recommendation systems have been around with us for a while now, and they are so powerful. They do have a strong influence on our decisions these days. From movie streaming services to online shopping stores, they are almost everywhere we look. If you are wondering how they know what you might buy after adding an “x” item to your cart, the answer is simple: Power of Data.

Recommendation systems are a very interesting field of machine learning. Recommendation system recommends the movie based on the users movie choices.

## **Data Source :**

MovieLens Dataset : <https://grouplens.org/datasets/movielens/>

This dataset contains 100,000 ratings and 3,600 tag applications applied to 9,000 movies by 600 users.

## **FEATURE DESCRIPTION :**

- MovieId (Quantitative) : Uniquely Identifies the movie
- Title (Qualitative) : Title of the Movie
- Genres (Categorical) : Defines a movie based on narrative elements
- UserId (Quantitative): Uniquely Identifies the user
- Ratings (Quantitative) : Rating of the movie
- Tag (Qualitative): Defines the type of movie
- TimeStamp (Quantitative): Timestamp of the rating

## SAMPLE DATASET :

### Movies.csv

	A	B	C	D	E	F	G	H	I	J
1	movieId	title	genres							
2	1	Toy Story	(Adventure Animation Children Comedy Fantasy							
3	2	Jumanji (1995)	(Adventure Children Fantasy							
4	3	Grumpier (Comedy)	Romance							
5	4	Waiting to (Comedy Drama)	Romance							
6	5	Father of t (Comedy								
7	6	Heat (1995)	Action Crime Thriller							
8	7	Sabrina (1995)	Comedy Romance							
9	8	Tom and H (Adventure Children								
10	9	Sudden De (Action								
11	10	GoldenEye (Action Adventure Thriller								
12	11	American I (Comedy Drama)	Romance							
13	12	Dracula: D (Comedy Horror								
14	13	Balto (1995)	Adventure Animation Children							
15	14	Nixon (1995)	Drama							
16	15	Cutthroat (Action Adventure Romance								
17	16	Casino (1995)	Crime Drama							
18	17	Sense and (Drama)	Romance							
19	18	Four Roos (Comedy								
20	19	Ace Ventu (Comedy								
21	20	Money Tra (Action Comedy Crime Drama Thriller								
22	21	Get Shorty (Comedy Crime Thriller								
23	22	Copycat (1995)	(Crime Drama Horror Mystery Thriller							
24	23	Assassins (Action Crime Thriller								
25	24	Powder (1995)	Drama Sci-Fi							
26	25	Leaving La (Drama)	Romance							
27	26	Othello (1995)	Drama							
28	27	Now and T (Children Drama								
29	28	Persuasior (Drama)	Romance							

### Ratings.csv

	A	B	C	D	E	F
1	userId	movieId	rating	timestamp		
2	1	1	4	9.65E+08		
3	1	3	4	9.65E+08		
4	1	6	4	9.65E+08		
5	1	47	5	9.65E+08		
6	1	50	5	9.65E+08		
7	1	70	3	9.65E+08		
8	1	101	5	9.65E+08		
9	1	110	4	9.65E+08		
10	1	151	5	9.65E+08		
11	1	157	5	9.65E+08		
12	1	163	5	9.65E+08		
13	1	216	5	9.65E+08		
14	1	223	3	9.65E+08		
15	1	231	5	9.65E+08		
16	1	235	4	9.65E+08		
17	1	260	5	9.65E+08		
18	1	296	3	9.65E+08		
19	1	316	3	9.65E+08		
20	1	333	5	9.65E+08		
21	1	349	4	9.65E+08		
22	1	356	4	9.65E+08		
23	1	362	5	9.65E+08		
24	1	367	4	9.65E+08		
25	1	423	3	9.65E+08		

Tags.csv

	A	B	C	D	E
1	userId	movieId	tag	timestamp	
2	2	60756	funny	1.45E+09	
3	2	60756	Highly quot	1.45E+09	
4	2	60756	will ferrell	1.45E+09	
5	2	89774	Boxing sto	1.45E+09	
6	2	89774	MMA	1.45E+09	
7	2	89774	Tom Hardy	1.45E+09	
8	2	106782	drugs	1.45E+09	
9	2	106782	Leonardo	1.45E+09	
10	2	106782	Martin Sc	1.45E+09	
11	7	48516	way too lo	1.17E+09	
12	18	431	Al Pacino	1.46E+09	
13	18	431	gangster	1.46E+09	
14	18	431	mafia	1.46E+09	
15	18	1221	Al Pacino	1.46E+09	
16	18	1221	Mafia	1.46E+09	
17	18	5995	holocaust	1.46E+09	
18	18	5995	true story	1.46E+09	
19	18	44665	twist endir	1.46E+09	
20	18	52604	Anthony H	1.46E+09	
21	18	52604	courtroom	1.46E+09	
22	18	52604	twist endir	1.46E+09	
23	18	88094	britpop	1.46E+09	
24	18	88094	indie recor	1.46E+09	
25	18	88094	music	1.46E+09	
26	18	144210	dumpster c	1.46E+09	
27	18	144210	Sustainabil	1.46E+09	
28	21	1569	romantic c	1.42E+09	
29	21	1569	wedding	1.42E+09	

All these datasets are merged through and to be used for our recommendation system.

## TOOLS USED :

**Python** : Python is an interpreted, high-level, general-purpose programming language used for performing the statistical analysis. When applying the technique of Web Scraping, Python coding will scrap the internet for selected data.

**Open CV** : OpenCV is a library of programming functions mainly aimed at real-time computer vision. Originally developed by Intel, it was later supported by Willow Garage then It sees. The library is cross-platform and free for use under the open-source BSD license.

**Pandas** : Pandas is a software library written for the Python programming language for data manipulation and analysis. In particular, it offers data structures and operations for manipulating numerical tables and time series.

**Numpy** : NumPy is a python library used for working with arrays. It also has functions for working in domain of linear algebra, fourier transform, and matrices.

**Seaborn** : Seaborn is a Python data visualization library based on matplotlib. It provides a high-level interface for drawing attractive and informative statistical graphics

## Methodology:

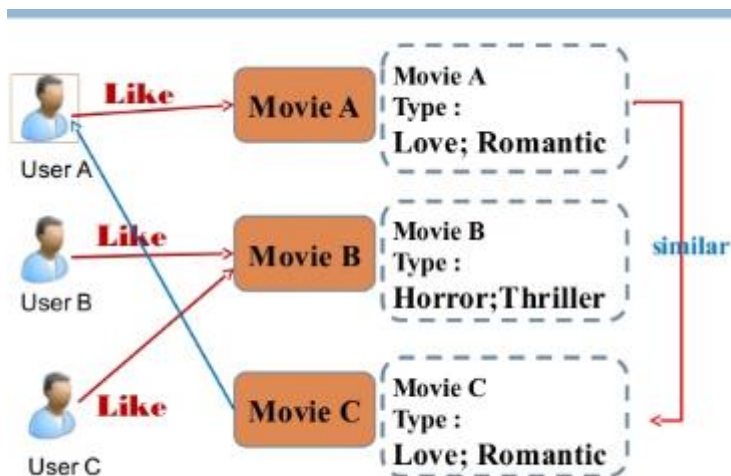
There are 5 major steps involved in the building a ML model for Movie Recommendation System. This encapsulates the following steps:

- Data loading
- Data cleaning
- Data Analysis
- Recommendation model

## Demographic Filtering :

## Content Based Filtering :

In this recommender system the content of the movie (overview, cast, crew, keyword, tagline etc) is used to find its similarity with other movies. Then the movies that are most likely to be similar are recommended.



## Collaborative Filtering :

### Collaborative Filtering

Our content based engine suffers from some severe limitations. It is only capable of suggesting movies which are close to a certain movie. That is, it is not capable of capturing tastes and providing recommendations across genres.

Also, the engine that we built is not really personal in that it doesn't capture the personal tastes and biases of a user. Anyone querying our engine for recommendations based on a movie will receive the same recommendations for that movie, regardless of who she/he is.

Therefore, in this section, we will use a technique called Collaborative Filtering to make recommendations to Movie Watchers. It is basically of two types:-

**User based filtering**- These systems recommend products to a user that similar users have liked. For measuring the similarity between two users we can either use pearson correlation or cosine similarity. This filtering technique can be illustrated with an example. In the following matrixes, each row represents a user, while the columns correspond to different movies except the last one which records the similarity between that user and the target user. Each cell represents the rating that the user gives to that movie. Assume user E is the target.

**Item Based Collaborative Filtering** - Instead of measuring the similarity between users, the item-based CF recommends items based on their similarity with the items that the target user rated. Likewise, the similarity can be computed with Pearson Correlation or Cosine Similarity. The major difference is that, with item-based collaborative filtering, we fill in the blank vertically, as oppose to the horizontal manner that user-based CF does.

## **EVALUATION METRIC :**

### **Personalization :**

Personalization is a great way to assess if a model recommends many of the same items to different users. It is the dissimilarity ( $1 - \text{cosine similarity}$ ) between user's lists of recommendations.

### **Intra-list Similarity:**

Intra-list similarity is the average cosine similarity of all items in a list of recommendations. This calculation uses features of the recommended items (such as movie genre) to calculate the similarity.

### **Coverage:**

Coverage is the percent of items in the training data the model is able to recommend on a test set. In this example, the popularity recommender has only 0.05% coverage, since it only ever recommends 10 items. The random recommender has nearly 100% coverage as expected. Surprisingly, the collaborative filter is only able to recommend 8.42% of the items it was trained on.

## EXPLORATORY DATA ANALYSIS :

Table :

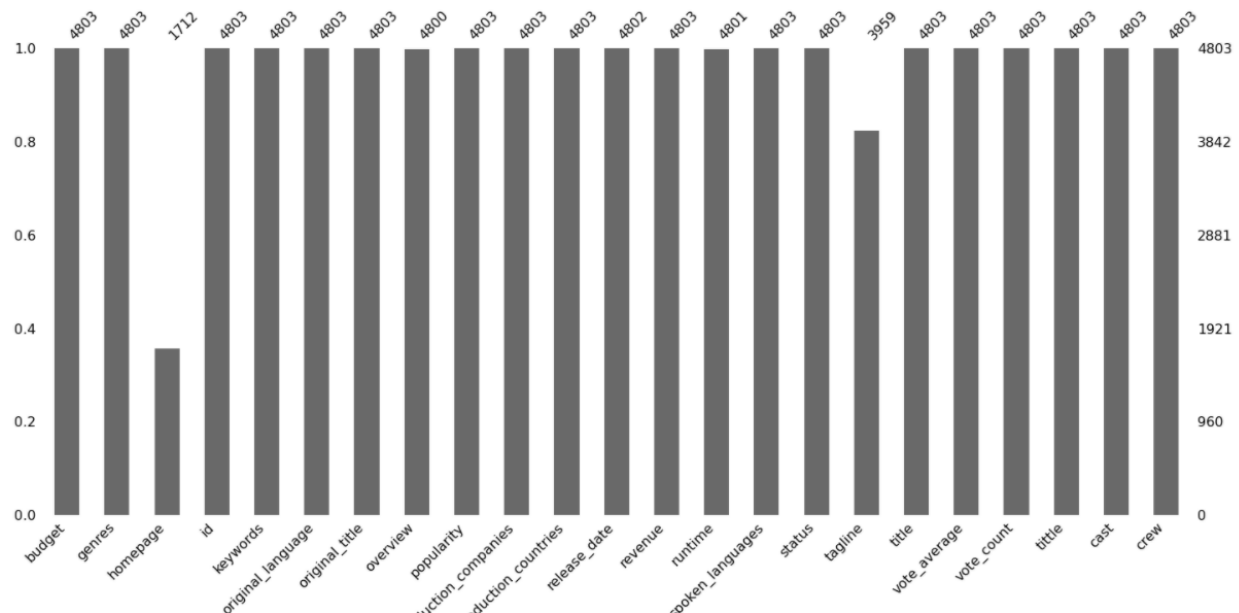
```
df_movies.describe()
```

	budget	id	popularity	revenue	runtime	vote_average	vote_count
count	4.803000e+03	4803.000000	4803.000000	4.803000e+03	4801.000000	4803.000000	4803.000000
mean	2.904504e+07	57165.484281	21.492301	8.226064e+07	106.875859	6.092172	690.217989
std	4.072239e+07	88694.614033	31.816650	1.628571e+08	22.611935	1.194612	1234.585891
min	0.000000e+00	5.000000	0.000000	0.000000e+00	0.000000	0.000000	0.000000
25%	7.900000e+05	9014.500000	4.668070	0.000000e+00	94.000000	5.600000	54.000000
50%	1.500000e+07	14629.000000	12.921594	1.917000e+07	103.000000	6.200000	235.000000
75%	4.000000e+07	58610.500000	28.313505	9.291719e+07	118.000000	6.800000	737.000000
max	3.800000e+08	459488.000000	875.581305	2.787965e+09	338.000000	10.000000	13752.000000

The above table gives the complete description of our data . We can check the mean , standard deviation , Min and Max of all the columns in our data.

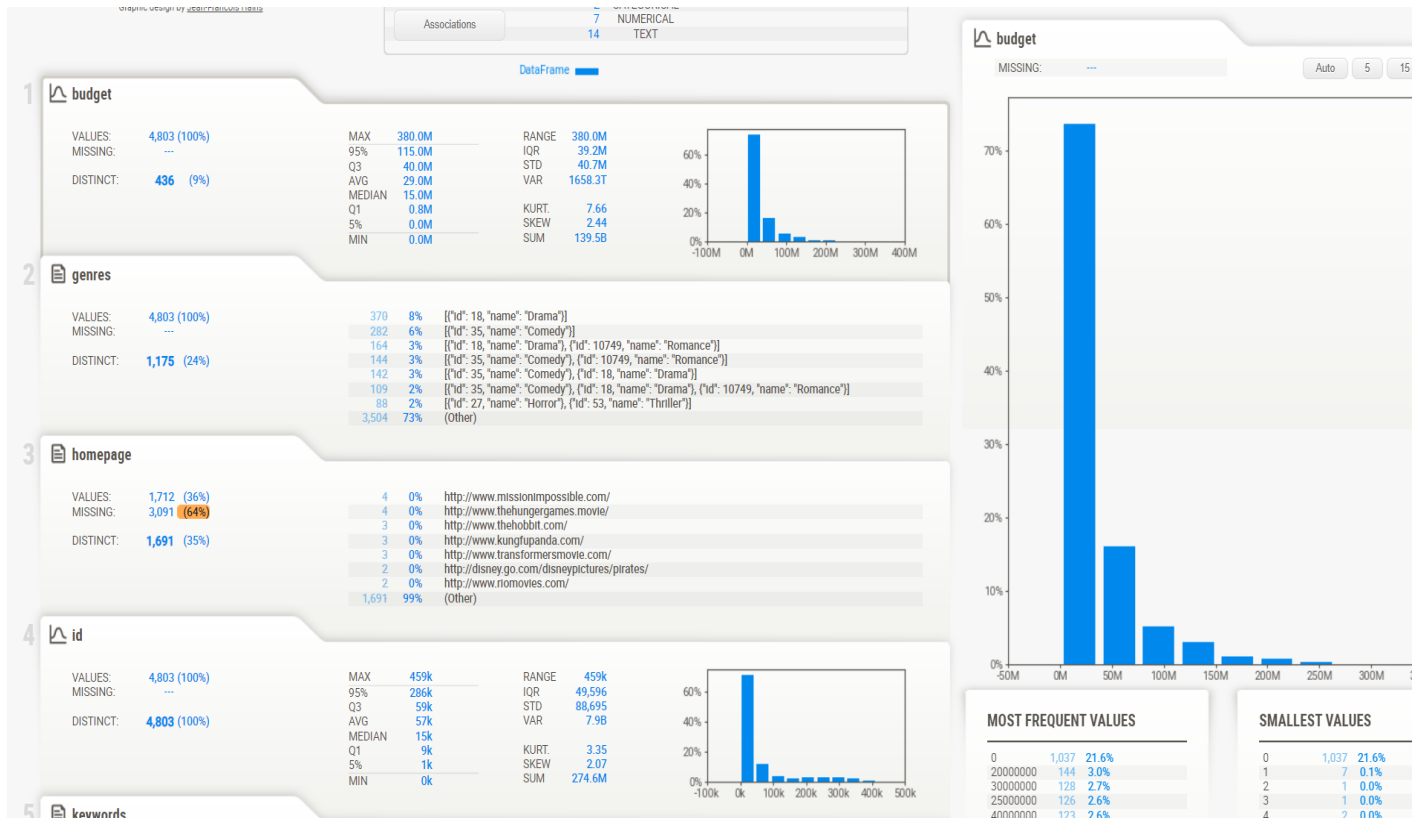
Bar Graph :

```
import missingno as msno
msno.bar(df_movies);
```



The above bar chart gives Information about the missing values in all the columns of the data. We can observe that homepage and tagline column has more missing values than any other column.

## Sweetviz :

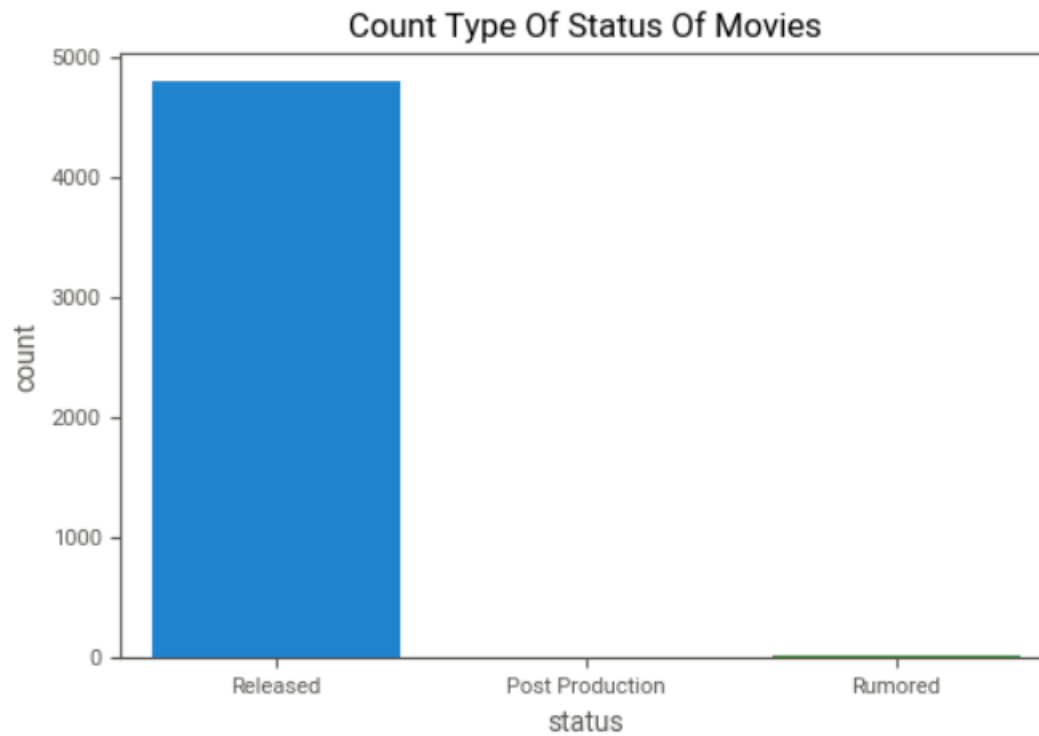


The above chart using the sweetviz gives the detailed information about each column in the dataframe . It gives the the count of each value and maximum,minimum values of each column.

## Count Plot :

```
d = sns.countplot(x='status',data=df_movies)  
d.axes.set_title("Count Type Of Status Of Movies")
```

```
Text(0.5, 1.0, 'Count Type Of Status Of Movies')
```



The above count plot gives the count of Different status of Movies . Most of the movies have Released status and very few movies have Rumored status.

Table :



Title	Budget
Chasing Papi	0
The Savages	0
Breakin' All the Rules	0
Jab Tak Hai Jaan	0
Red Dog	0
Kites	0
Richard III	0
Good Intentions	0
Partition	0
It's a Wonderful Afterlife	0
Jungle Shuffle	0
Fifty Dead Men Walking	0
White Noise 2: The Light	0
Eulogy	0
The Reef	0
Free Style	0

The above table gives information about the lowest budget movies . There are movies that are made with zero budget.

Title	Budget
Pirates of the Caribbean: At World's End	300000000
Avengers: Age of Ultron	280000000
Superman Returns	270000000
John Carter	260000000
Tangled	260000000
Spider-Man 3	258000000
The Lone Ranger	255000000
X-Men: Days of Future Past	250000000
The Hobbit: The Desolation of Smaug	250000000
Captain America: Civil War	250000000
Batman v Superman: Dawn of Justice	250000000
The Hobbit: An Unexpected Journey	250000000
The Hobbit: The Battle of the Five Armies	250000000
Harry Potter and the Half-Blood Prince	250000000
The Dark Knight Rises	250000000
Spectre	245000000

We can observe from the table that the maximum budget of the movie is 3000000000 and next highest is 28

