# Movie Recommendation System Using R

A mini project report submitted for

R Programming Lab (Semester VIII)

## By

Harsh Dobariya (B1/814)

Akshay Kalapgar (B2/831)

Mohit Kamble (B2/832)

Under the guidance of

**Prof. Govind Wakure** 

# MANJARA CHARLTABLE TRUST RAJIV GANDHI INSTITUTE OF TECHNOLOGY, MUMBAI

Permanently affiliated to Mumbai University

Department of Information Technology, Mumbai-400615. UNIVERSITY OF MUMBAI 2021



# **CERTIFICATE**

### DEPARTMENT OF INFORMATION TECHNOLOGY

This is to certify that

Harsh Dobariya B 1 / 814

Akshay Kalapgar B 2 / 831

Mohit Kamble B 2 / 832

Have satisfactory completed this synopsis entitled Movie Recommendation System Using R

Towards the partial fulfilment of the FOURTH YEAR OF ENGINEERINGIN (Information Technology)
As laid by University of Mumbai.

Guide **Prof. Govind Wakure** 

H.O.D.

Dr. Sunil Wankhade

Principal

Dr. Sanjay U. Bokade

**Internal Examiner** 

External Examiner



#### **DECLARATION**

We declare that this written submission represents our ideas in our own words and where others ideas or words have been included; we have adequately cited and referenced the original sources. We also declare that we have adhered to all principles of academics honestly and integrity and have not misrepresented or fabricated or falsified any idea/data/fact/sources in our submission. We understand that any violation of the above will be cause for disciplinary action by the institute and can also evoke penal action from the sources which have thus not been properly cited or from whom proper permission has not been taken needed.

(Signature)

1. Harsh Dobariya B1/814

2. Akshay Kalapgar B2/831

3. Mohit Kamble B2/832



#### **ACKNOWLEDGEMENT**

We wish to express our sincere gratitude to Dr. Sanjay U. Bokade, Principle and Dr. Sunil Wankhade, H.O.D of Information Technology Department of RGIT for proving us an opportunity to do our project work on "Movie Recommendation System Using R." This project bears on imprint of many peoples. We sincerely thank our project guide Prof. Govind Wakure for his guidance and encouragement in carrying out this synopsis work. Finally, we would like to thank our colleagues and friends who helped us in completing the Project (synopsis) work successfully.

1. HarshDobariy B1/814

2. Akshay Kalapgar B2/831

3. MohitKamble B2/831



#### **ABSTRACT**

Everyone loves movies irrespective of age, gender, race, color, or geographical location. We all in a way are connected to each other via this amazing medium. Yet what most interesting is the fact that how unique our choices and combinations are in terms of movie preferences. Some people like genre-specific movies are it a thriller, romance, or sci-fi, while others focuses on lead actors and directors. When we take all that into account, it's astoundingly difficult to generalize a movie and say that everyone would like it. But with all that said, it is still seen that similar movies are liked by a specific part of the society. However, to really enhance the user experience through personalized recommendations, we need dedicated recommender systems. From a business standpoint, the more relevant products a user finds on the platform, the higher their engagement. This often results in increased revenue for the platform itself. Various sources say that as much as 35–40% of tech giants' revenue comes from recommendations alone.



# **Table of Content**

Sr.	Topics	Page No.
No.	_	_
1.	Introduction.	1
2.	Existing System	2
3.	Proposed System.	3
4.	Literature survey.	4
5.	Methodology	5
6.	Implementation.	10
7.	Testing and Results.	16
8.	Conclusion.	19
9.	References.	20

# **List of Figures**

Sr.	Figures	Page no.
No.	-	_
1.	Architecture of the system	5
2.	The collaborative filtering process	7
3.	Isolation of the co-rated items and similarity computation.	7
4.	Class Diagram	14
5.	User Case Diagram	15



#### 1. INTRODUCTION

Recommendation system has become very popular in many aspects in real social networks, such as e-commerce services Amazon.com, movie rating website IMDB, and DVD rental service company Netflix. This project focuses on multiple level link prediction of recommending movies. Movie rating data has become largely accessible via the Internet. Most rating systems are designed to have 5 score choices from 1 to 5. However, in finding a potential interesting movie for a certain user, there is not much difference between score 4 and 5. In fact, This is a classification problem on predicting the category of a movie that has not been watched by the user. So instead of a specific score, we are more concerned with the genre of the movie and based on the user input the code will present the end user with the desirable results. Another important motivation for this project is to build a social recommendation system in the future. The idea comes from recommendations we get from our friends in everyday life. It is quite common that we tend to value more about reviews given from other strangers before watching. In this project, we have used a data set consisting of thousands of movie titles which are saved along with multiple genres along with the ratings. So once the user selects three fields which are genres, the recommendation system will give the user the best results matching the entire genrewhich are having the best rating in a sorting order.



#### 2. EXISTING SYSTEM

The biggest issue facing recommender systems is that they need a lot of data to effectively make recommendations. It's no coincidence that the companies most identified with having excellent recommendations are those with a lot of consumer user data: Google, Amazon, Netflix, Last.fm. The more item and user data a recommender system must work with, the stronger the chances of getting good recommendations. But to get good recommendations, you need a lot of users, so you can get a lot of data for the recommendations. Another problem with search engine would be the rapidly changing data. Clearly an algorithmic approach will find it difficult if not impossible to keep up with fashion trends. Lastly the search engine is unpredictable because itcannot give a perfect match every time.



#### 3. PROPOSED SYSTEM

The Movie Recommendation Search Engine uses a user based collaborative filter. User based CF approach is very like item based CF. But in this one, similarity is computed between users, not item. The underlying assumption of the collaborative filtering approach is that if a person *A* has the same opinion as a person *B* on an issue, A is more likely to have B's opinion on a different issue than that of a randomly chosen person. The search engine will take the user's input on genres and movie name and recommends the user a list of 10 movies that he would also enjoy. The search engine uses k means clustering method to split the object in the dataset into clusters. A class label will be attached to the clusters for matching based on the user input. A hard code on the search engine is that the user's input for movie name should be in the order of 5,4,3 ratings since the search engine would return a movie based on the hard code.



#### 4. LITRATURE SURVEY

The term "collaborative filtering" was introduced in the context of the first commercial recommender system, called Tapestry, which was designed to recommend documents drawn from newsgroups to a collection of users. The motivation was to leverage social collaboration to prevent users from getting inundated by a large volume of streaming documents. Collaborative filtering, which analyzes usage data across users to find well matched user-item pairs, has since been juxtaposed against the older methodology of content filtering which had its original roots in information retrieval. In content filtering, recommendations are not "collaborative" in the sense that suggestions made to a user do not explicitly utilize information across the entire user-base. Some early successes of collaborative filtering on related domains included the Group Lens system.

Collaborative filtering systems have many forms, but many common systems can be reduced to two steps:

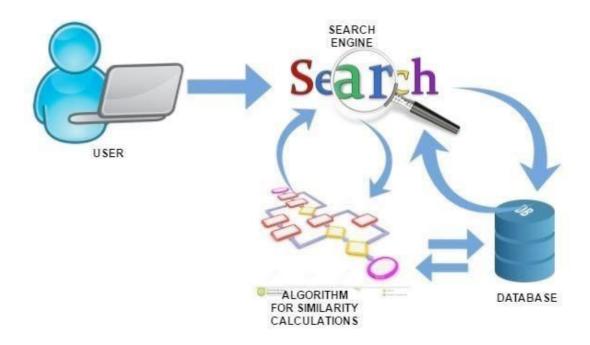
- 1. Look for users who share the same rating patterns with the active user (the userwhom the prediction is for).
- 2. Use the ratings from those like-minded users found in step 1 to calculate aprediction for the active user

Over the last decade, lots of researchers have studied innovative approaches of recommender systems to solve these problems of CF (Collaborative Filter method) and CB (Content Based Filtering method), and to apply them into real world problems. Especially, applications of data mining techniques to recommender systems have been effective to offer personalized information to the user through analysis of his/her preference. In the past, many have used Item Based Collaborative Filtering for their respective search engine. But with an increasing data each day, the search engine should be accurate and the database should be updated often. We use a User Based Collaborative Filter which computes k means algorithm to find the nearest neighbor within the cluster.



## **5. METHODOLOGY**

## **Architecture of the system:**



#### **Algorithm**

This approach uses user rating data to compute the similarity between users or items. This is used for making recommendations. This was an early approach used in many commercial systems. It's effective and easy to implement. Typical examples of this approach are neighborhood-based CF and item-based/user-based top-N recommendations. For example, in user based approaches, the value of ratings user 'u' gives to item 'I' is calculated as an aggregation of some similar users' rating of the item:



$$r_{u,i} = \operatorname{aggr}_{u' \in U} r_{u',i}$$

where 'U' denotes the set of top 'N' users that are most similar to user 'u' who rated item 'i'. Some examples of the aggregation function includes:

$$egin{aligned} r_{u,i} &= rac{1}{N} \sum_{u' \in U} r_{u',i} \ r_{u,i} &= k \sum_{u' \in U} \mathrm{simil}(u,u') r_{u',i} \ r_{u,i} &= ar{r_u} + k \sum_{u' \in U} \mathrm{simil}(u,u') (r_{u',i} - ar{r_{u'}}) \end{aligned}$$

where k is a normalizing factor defined as  $k=1/\sum_{u'\in U}|\operatorname{simil}(u,u')|$ . and  $\bar{r_u}$  is the average rating of user u for all the items rated by u.

The neighborhood-based algorithm calculates the similarity between two users or items produces a prediction for the user by taking the weighted average of all the ratings.

Similarity computation between items or users is an important part of this approach. Multiple measures, such as Pearson correlation and vector cosine based similarity are used for this

The Pearson correlation similarity of two users x, y is defined as

$$ext{simil}(x,y) = rac{\sum\limits_{i \in I_{xy}} (r_{x,i} - ar{r_x}) (r_{y,i} - ar{r_y})}{\sqrt{\sum\limits_{i \in I_{xy}} (r_{x,i} - ar{r_x})^2 \sum\limits_{i \in I_{xy}} (r_{y,i} - ar{r_y})^2}}$$

where Ixy is the set of items rated by both user x and user y.

The cosine-based approach defines the cosine-similarity between two users x and y as:[4]

$$ext{simil}(x,y) = \cos(ec{x},ec{y}) = rac{ec{x} \cdot ec{y}}{||ec{x}|| imes ||ec{y}||} = rac{\sum\limits_{i \in I_{xy}} r_{x,i} r_{y,i}}{\sqrt{\sum\limits_{i \in I_{x}} r_{x,i}^2} \sqrt{\sum\limits_{i \in I_{y}} r_{y,i}^2}}$$

The user based top-N recommendation algorithm uses a similarity-based vector model to identify the k most similar users to an active user. After the k, most similar users are found, their corresponding user-item matrices are aggregated to identify the set of items to be recommended. A popular method to find the similar users is the Locality-sensitive hashing, which implements the nearest neighbor mechanism in linear time.

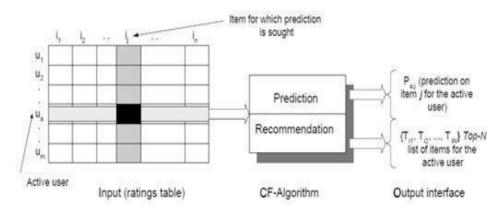
The advantages with this approach include: the explain ability of the results, which is an important aspect of recommendation systems; easy creation and use; easy facilitation of new data; content-independence of the items being recommended; good scaling with co-rated items.

There are also several disadvantages with this approach. Its performance decreases when data gets sparse, which occurs frequently with web-related items. This hinders the scalability of

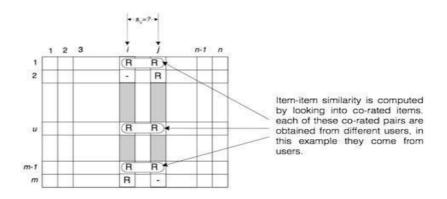


#### RAJIV GANDHI INSTITUTE OF TECHNOLOGY, MUMBAI

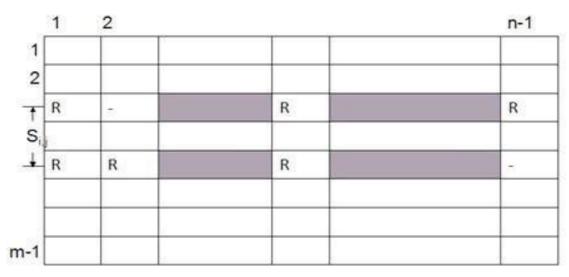
This approach and creates problems with large datasets. Although it can efficiently handle new users because it relies on a data structure, adding new items becomes more complicated since that representation usually relies on a specific vector space. Adding additional items requires inclusion of the additional item and the re-insertion of all the elements in the structure.



"The collaborative filtering process."



<sup>&</sup>quot;Isolation of the co-rated items and similarity computation."



<sup>&</sup>quot;Isolation of the co-rated users and similarity computation."

# K means algorithm:

Algorithm: k-means. The k-means algorithm for partitioning, where each cluster's center

is represented by the mean value of the objects in the cluster.

Input: k: the

number of

clusters,

D: a data set containing n

objects.Output: A set of k

clusters.

Method:

- (1) arbitrarily choose k objects from D as the initial cluster centers;
- (2) repeat



- (3) (re)assign each object to the cluster to which the object is the most similar, based on themean value of the objects in the cluster;
- (4) update the cluster means, that is, calculate the mean value of the objects for each cluster:
- (5) until no change;

#### Libraries used:

- ➤ library(proxy): Distance and Similarity Measures
  Provides an extensible framework for the efficient calculation of auto- and cross- proximities,
  along with implementations of the most popular ones.
- ➤ library(recommenderlab): Lab for Developing and Testing Recommender Algorithms Provides a research infrastructure to test and develop recommender algorithms including UBCF, IBCF, FunkSVD and association rule-based algorithms.
- ➤ library(reshape2): A Reboot of the Reshape Package
  Flexibly restructure and aggregate data using just two functions: melt and 'dcast'(or 'acast').
  library(ggplot2): Create Elegant Data Visualizations Using the Grammar of Graphics A system for 'declaratively' creating graphics, based on "The Grammar of Graphics". You provide the data, tell 'ggplot2' how to map variables to aesthetics, what graphical primitives to use, and it takes care of the details.
- ➤ library(shiny): Web Application Framework for R

  Makes it incredibly easy to build interactive web applications with R. Automatic "reactive" binding between inputs and outputs and extensive pre-built widgets make it possible to build beautiful, responsive, and powerful applications with minimal effort.
- ➤ library(shinythemes):

  Themes for use with Shiny. Includes several Bootstrap themes from, which are packaged for use with Shiny applications



#### 6. IMPLEMENTATION

#### **System Design**

System design is transition from a user oriented document to programmers or database personnel. The design is a solution, how to approach to the creation of a new system. This is composed of several steps. It provides the understanding and procedural details necessary for implementing the system recommended in the feasibility study. Designing goes through logical and physical stages of development, logical design reviews the present physical system, prepare input and output specification, details of implementation plan and prepare a logical design walkthrough.

The database tables are designed by analyzing functions involved in the system and format of the fields is also designed. The fields in the database tables should define their role in the system. The unnecessary fields should be avoided because it affects the storage areas of the system. Then in the input and output screen design, the design should be made user friendly. The menu should be precise and compact.

#### **Sample Codes**

• Dataset for Movie Recommendation Search Engine is in .csv format. So, we have used below code to read data in R studio.

```
search <- read.csv("search.csv", stringsAsFactors=FALSE)
ratings <- read.csv("ratings.csv", header = TRUE)|
search <- search[-which((search$movieId %in% ratings$movieId) == FALSE),]</pre>
```



• Sample dataset screenshot of Movie dataset

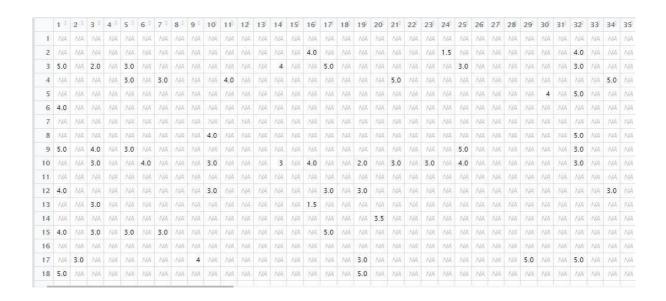
	movield	title	genres
1	1	Toy Story (1995)	Adventure Animation Children Comedy Fantasy
2	2	Jumanji (1995)	Adventure Children Fantasy
3	3	Grumpier Old Men (1995)	Comedy Romance
4	4	Waiting to Exhale (1995)	Comedy Drama Romance
5	5	Father of the Bride Part II (1995)	Comedy
6	6	Heat (1995)	Action Crime Thriller
7	7	Sabrina (1995)	Comedy Romance
8	8	Tom and Huck (1995)	Adventure Children
9	9	Sudden Death (1995)	Action
10	10	GoldenEye (1995)	Action Adventure Thriller

• To use the ratings data for building a recommendation engine with recommender lab, I convert rating matrix into a sparse matrix of type real Rating Matrix.

```
ratingmatrix <- dcast(ratings, userId~movieId, value.var = "rating", na.rm=FALSE)
ratingmatrix <- ratingmatrix[,-1]
colnames(userSelect) <- colnames(ratingmatrix)
ratingmatrix_new <- rbind(userSelect,ratingmatrix)
ratingmatrix_new <- as.matrix(ratingmatr|x_new)

#Convert rating matrix into a sparse matrix
ratingmatrix_new <- as(ratingmatrix_new, "realRatingMatrix")</pre>
```

#### **Output:**



• The recommender lab package contains some options for the recommendationalgorithm:

- We have defined static array to display genres in dropdown.
- Sample shiny code to display output

```
output$table <- renderTable({
    movie_recommendation(input$select, input$select2, input$select3)
})
output$dynamic_value <- renderPrint({
    c(input$select,input$select2,input$select3)
})</pre>
```

#### 4.3 UML Diagrams

#### User Model View:

This view represents the system from the users' perspective.

The analysis representation describes a usage scenario from the end-user's perspective.

#### Structural View

In this model, the data functionality is arrived from the inside the system. This model view models the static structures.

#### Behavioral View

It represents the dynamic of behavioral as parts of the system depicting the interactions of collection between various structural elements described in the user module and structural model view.

#### Implementation Model View

In this structural and behavioral as parts of the system are represented as they are to be built.

MANJARA CHARLTABLE TRUST
RAJIV GANDHI INSTITUTE OF TECHNOLOGY, MUMBAI

#### **Environmental Model View**

In this structural and behavioral aspects of the environment in which the system is to be implemented are represented.

#### Relationships in UML

There are four kinds of relationships in the UML:

Dependency

Association

Generalization

Realization

#### Class diagram:

Class diagram, one of the most commonly used diagrams in object-oriented system, models the static design view for a system. The static view mainly supports the functional requirements of a system – the services the system should provide to the end users. We will see from our practical experience that lots of fun comes out when modeling out system with class diagrams. The discussion on different views of class diagrams for the system will be put into emphasis later in this paper. A class diagram shows a set of classes, interfaces, and collaborations and their relationships. Class diagrams involve global system description, such as the system architecture, and detail aspects such as the attributes and operations within a class as well. Themost common contents of a class diagram are:

- Classes
- Interfaces
- Collaborations
- Dependency, generalization, and association relationships
- Notes and constraints



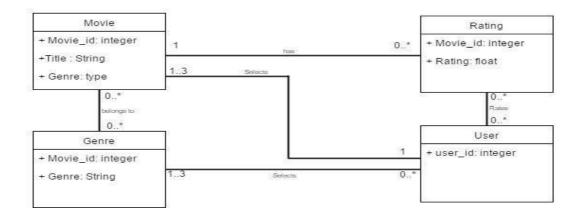
Below is a high-level class diagram for the Movie Recommendation System. This diagram depicts the relationship between users, ratings, Movie, Genre. The multiplicity is also shown to help understand the system better. One user can rate several movies and must select at least one and a maximum of three genre and movies to make the search engine work, a movie can belong to several genres at a time, a movie can be rated by number of users, in the following subsections, four groups of class diagrams are presented and analyzed in detail.

**Movie:** This class contains all the attributes required for identifying a movie uniquely it containsattributes like movie\_id, title, genre.

**Rating:** This class contains the ratings given by different users for different movies.

**User:** This class contains the user\_id attribute and user is responsible for selecting the genre and movies for getting the required movies.

**Genre:** This class contains the genres to which a movie belongs to.

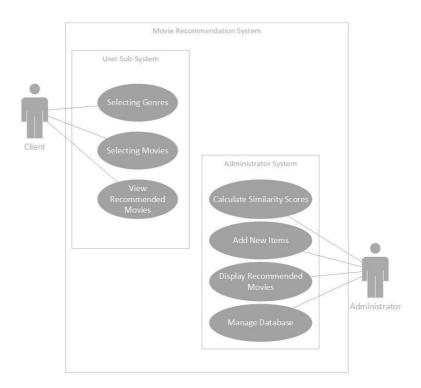


## **Use Case Diagram:**

Use case diagram contains two actors as shown in the below figure. User is an actor who performs the actions such as selecting genres and based on the Genres selected the recommender system will display the movie list drop down. User has an option to select another movie from the dropdown list to get a movie recommendation of his personal favorite. Once user selects all the mandatory genre fields and confirms the movies displayed for each genre the end user will be displayed with the desired movies list based on the Similarity scores and movies will be displayed in a sorted order.

# MANJARA CHARITABLE TRUST RAJIV GANDHI INSTITUTE OF TECHNOLOGY, MUMBAI

Admin is another actor who performs actions like updating the movie database, refactoring the algorithm to use indexing or Hash functionalities to decrease the time to retrieve the search results from database.



# **Software Requirements:**

Operating System: Mac OS, Windows

R studio (V 1.0.143)Cran R (V 3.3.3)

# **Hardware Requirements:**

System: Intel core i5 processor

Ram: 8Gb



#### 7. TESTING AND RESULTS

#### **Testing Methods**

- 1) To ensure that during operation the system will perform as per specification.
- 2) To make sure that system meets the user requirements during operation
- 3) To make sure that during the operation, incorrect input, processing and output will be detected
- 4) To see that when correct inputs are fed to the system the outputs are correct
- 5) To verify that the controls incorporated in the same system as intended
- 6) Testing is a process of executing a program with the intent of finding an error.
- 7) A good test case is one that has a high probability of finding a yet undiscovered error The software developed has been tested successfully using black box testing strategy and any errors that are encountered are corrected and again the part of the program or the procedure.

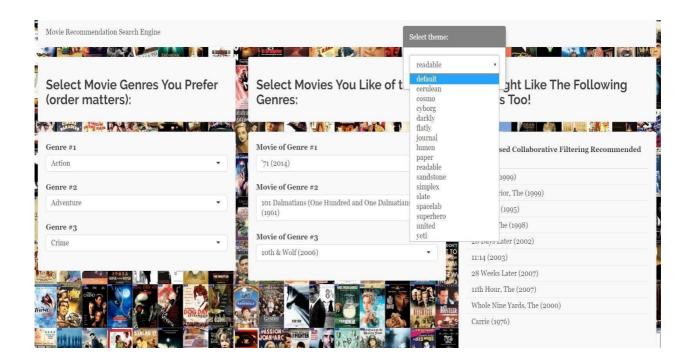
Black box testing strategy is a case design method that uses the control structure of the Transaction data set and k out of N strategy is used and designed to derive test cases. All independents path in a module are exercised at least once, all logical decisions are exercised at once, execute all loops at boundaries and within their operational bounds exercise internal datastructure to ensure their validity.

In this application, it does not watch the internal variables during testing. This gives clear idea about what is going on during execution. The points at which the bug occurs were all clear and were removed.

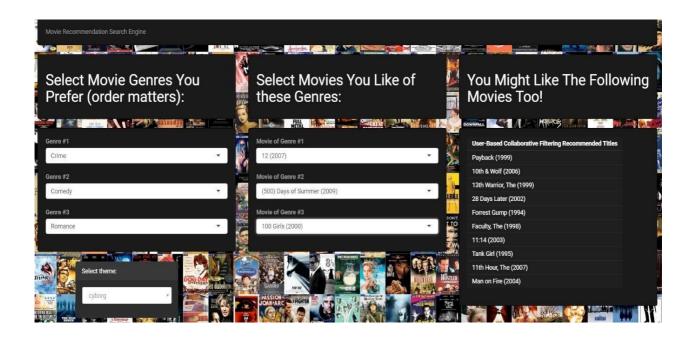
Note that the result of the system testing will prove that the system is working correctly. It will give confidence to system designer, users of the system; prevent frustration during implementation process etc.



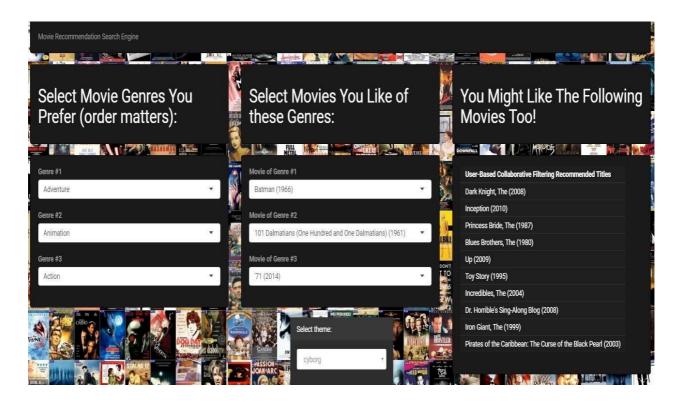
#### **Results**



User has selected Genre Crime, Comedy and Romance and movies that are displayed by default are as Below:









# 8. CONCLUSION AND FUTURE WORK

The movie recommendation search engine displays the result list of 10 movies which users may also be interested in. Item based collaborative filtering had less error than user based collaborative filtering. In addition, its less-dynamic model was computed less often and stored in a smaller matrix, so item- based system performance was better than user based systems. Soon, we would like to improve the accuracy of prediction by using Slope One algorithm. Slope One is a family of algorithms used for collaborative filtering, introduced in a 2005 paper by Daniel Lemire and Anna Maclachlan. Arguably, it is the simplest form of non-trivial item- based collaborative filtering based on ratings. Their simplicity makes it especially easy to implement them efficiently while their accuracy is often on par with more complicated and computationally expensive algorithms. They have also been used as building blocks to improve other algorithms. They are part of major open-source libraries such as Apache Mahout and Easyrec.



# **REFERENCES**

- (1) https://rpubs.com/tarashnot/recommender\_comparison
- (2) https://link.springer.com/article/10.1007/s10479-016-2367-1#Sec15
- (3) https://shiny.rstudio.com/tutorial/
- (4) https://muffynomster.wordpress.com/2015/06/07/building-a-movie-recommendation-engine-with-r/