

Movie Recommender System Using K-Means Clustering AND K-Nearest Neighbor

Rishabh Ahuja

School of Information and
Communication Technology
Gautam Buddha University
Greater Noida, UP-201308

Email: rishabhahuja279@gmail.com

Arun Solanki

School of Information and
Communication Technology
Gautam Buddha University
Greater Noida, UP-201308

Email: ymca.arun@gmail.com

Anand Nayyar

Graduate School
Duy Tan University
Da Nang, Vietnam

Email: anandnayyar@duytan.edu.vn

Abstract—In the field of Artificial Intelligence Machine learning provides the automatic systems which learn and improve itself from experience without being explicitly programmed. In this research work a movie recommender system is built using the K-Means Clustering and K-Nearest Neighbor algorithms. The movielens dataset is taken from kaggle. The system is implemented in python programming language. The proposed work deals with the introduction of various concepts related to machine learning and recommendation system. In this work, various tools and techniques have been used to build recommender systems. Various algorithms such as K-Means Clustering, KNN, Collaborative Filtering, Content-Based Filtering have been described in detail. Further, after studying different types of machine learning algorithms, there is a clear picture of where to apply which algorithm in different areas of industries such as recommender systems, e-commerce, etc. Then there is an illustration of how implementations and working of the proposed system are used for the implementation of the movie recommender system. Various building blocks of the proposed system such as Architecture, Process Flow, Pseudo Code, Implementation and Working of the System is described in detail. Finally, in this work for different cluster values, different values of Root Mean Squared Error are obtained. In this proposed work as the no of clusters decreases, the value of RMSE also decreases. The best value of RMSE obtained is 1.081648. The results given by the proposed system are better than the existing technique on the basis of RMSE value.

Keywords—Recommender System, k-Means, KNN, Collaborative Filtering, Content-Based Filtering

I. INTRODUCTION

A recommendation system is a type of information filtering system which is used to predict the "rating" or "preference" a user would give to an item. A recommendation system collect data about the user's preferences either implicitly or explicitly on different items like movies, shopping, tourism, TV etc [3], [4], [5], [6], [7]. An implicit acquisition in the development of movie recommendation system uses the user's behavior while watching the movies. On the other hand, a explicit acquisition in the development of movie recommendation system uses the user's previous ratings or history. Collaborative filtering is the technique to filter or calculate the items through the sentiments of other users [8], [9], [10]. Collaborative filtering first collect the movie ratings or preference given by different users and then suggest movies to the different user based on similar tastes and interests in the past. The other supporting

technique that are used in the development of recommendation system is clustering. Clustering[20], [21] is a process to group a set of objects in such a way that objects in the same clusters are more similar to each other than to those in other clusters [11], [12], [13], [14]. K-Means [13], [23], [33] Clustering along with K-Nearest Neighbor [18], [24] is implemented on the movielens dataset in order to obtain the best-optimized result. In existing technique the data is scattered which results in a high number of clusters while in the proposed technique data is gathered and results in a low number of clusters. The process of recommendation of a movie is optimized in the proposed scheme. The proposed recommender system predicts the user's preference of a movie on the basis of different parameters. The recommender system works on the concept that people are having common preference or choice. These user will influence on each other's opinions. This process optimize the process and having lower RMSE.

The work starts with the section I as Introduction section with the basics of recommendation system. Section II discusses the latest work done by recent authors with the details of techniques and tools used by different authors. Section III describe the evolution of the proposed recommendation system. Section IV shows the algorithm of the proposed system. Section V shows the implementation of the proposed system. The section VI discusses the working and results of the system with the help of the snapshot of the system. Section VII is having the conclusion and future work of the proposed system.

II. LITERATURE REVIEW

Content based [40], [41] collaborative [42] and hybrid [43] are the different approaches used by past researcher for the development of recommender system. In 2007 a web-based movie recommendation system using hybrid filtering methods is presented by the authors [35]. In 2011 a movie recommendation system based on genre correlations is proposed by the authors [36]. In 2013 a Bayesian network and Trust model based movie recommendation system is proposed, the Bayesian network is imported for user preference modeling and trust model is used to filter the recommending history data and enable the system to tolerant the noisy data [37]. In 2016, authors proposed Recommender systems to predict the rating for users and items, predominantly from big data to recommend their likes. Movie recommendation systems provide a mechanism to assist users in classifying users with

similar interests. This system (K-mean Cuckoo) has 0.68 MAE [15], [16]. In 2017 authors used a new approach that can solve sparsity problem to a great extent [38]. In 2018, authors built a recommendation engine by analyzing rating data sets collected from Twitter to recommend movies to specific user using R [39].

III. EVOLUTION OF PROPOSED MOVIE RECOMMENDATION SYSTEM

This section consist of the architecture and process flow of proposed system.

A. Architecture

Figure 1 shows the architecture of the proposed system. It consists of three modules, namely the input module, a processing module, an output module. Figure 1 gives the clear conceptual idea about the working of the proposed recommendation system.

Next, an illustration of each module is done in detail and explains the architecture of the system. This helps in understanding the architecture of the system in a crisp and clear manner.

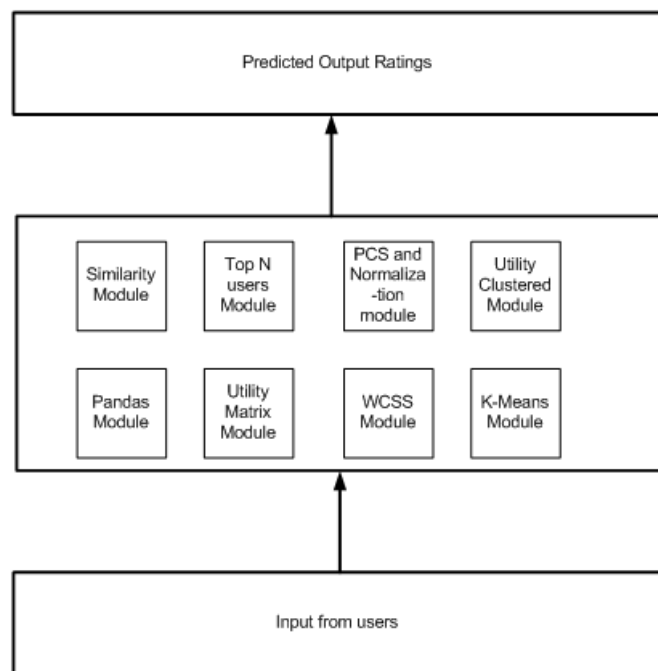


Fig. 1: Architecture of the proposed system

As discussed in Figure 1, the three modules of the proposed in the architecture are as follows:

Input Module

In this module, the user is asked to give the details as input. In input, the user gives the detail about himself by providing details such as userId, age, gender, pin code. This information is further passed to the next module i.e. the processing module.

Processing Module

In this, the panda's module first separates the data from the raw files. It separates the information about the user and movie items into a separate data frame using the panda's library. After separating the data from the raw form, in a utility matrix module a utility matrix is built which defines which user rated which movie. This helps in figuring out how many times each movie is rated by the users. Then based on previous preprocessing of data, separate data frames for the training set and testing set is created. This is done to further evaluate the performance of the system. After getting the utility matrix, K-means clustering is used to build a separate data frame which shows which movie belongs to which genre. The Within-Cluster Sum of Squares (WCSS) is a measure of the variability of the observations within each cluster. In general, a cluster that has a small sum of squares is more compact than a cluster that has a large sum of squares. In WCSS module the right no of clusters is chosen using the technique Within Clustered Sum of Square. Now, for calculating the average rating given by each user given to each cluster, a utility clustered matrix is created. In utility clustered module the utility clustered matrix is used to calculate the similarity between the users. The PCS and normalization module calculates the correlation using the utility clustered matrix. Finally, in the KNN module and similarity module using the K-Nearest Neighbor predictions for movie rating is calculated with the help of the similarity matrix and utility clustered matrix.

Output Module

The output module describes the predicted movies that the input user might like. Further, in the output along with the movies, their predicted ratings are also defined which input user might give to the movies.

B. Process Flow

Figure 2 shows the process of the flow diagram of the Movie Recommendation System. This diagram shows the process flow of the proposed system. Process flow depicts how the system is working, how the system is dealing with the raw data, and how the system predicts the rating for the input userId.

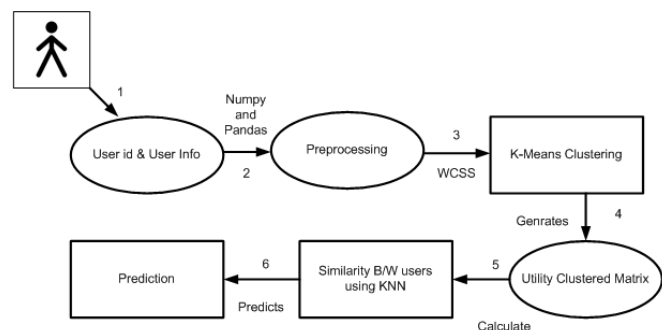


Fig. 2: Process Flow Diagram

Step 1: The user gives the userId and information such as

gender, age, pin code.

Step 2: Using the numpy and pandas library the raw data is preprocessed into separate data frames.

Step 3: Within Clustered Sum of the Squared method is used to find the right no clusters so that K-means clustering can be applied to the movie.

Step 4: After applying K-means clustering a utility clustered matrix is build which defines average rating the user gives to each cluster

Step 5: Using the utility clustered matrix and Pearson correlation similarity between the users are calculated.

Step 6: Finally KNN uses the utility clustered matrix and similarity to predict the movies for input user.

IV. ALGORITHM OF PROPOSED SYSTEM

Algorithm for the proposed algorithm is as follows:

Step1: Import the python libraries: Numpy, Pandas, Matplotlib, sklearn

Step2: Read the csv information as data frames in user and item variable.

Step3: Split the data into the training set and test set as data frame into the variables rating and rating test.

Step4: Create a utility matrix name utility which tells which user rated which movie.

Step5: Using the WCSS method choose the right number of clusters so that the K-means Clustering technique can be applied to classify the movies according to the number of clusters.

Step6: Define the utility clustered matrix after applying the K-means clustering algorithm.

Step7: Apply Pearson Correlation metric on utility clustered matrix to calculate the similarity matrix between the users.

Step8: Normalize the values stored in utility matrix.

Step9: Guess() function takes two parameters as input userID and topN users which is used by KNN to predict the movie ratings for topN similar users.

Step10: ratingTest data frame ratings are used for comparison while using the guess function for predicting the ratings of test users.

Step11: RMSE is calculated to evaluate the accuracy of the model.

V. IMPLEMENTATION

The system has been implemented in python programming language using K-Means clustering library and K-Nearest Neighbor. The implementation of the system consists of many sub-sections which are standard processes to be followed while solving any machine learning [17], [19], [22], [27], [28], [29], [30], [31], [32], [34] problem. These are as follows:

- 1) Data Collection
- 2) Data Preparation
- 3) Model Creation
- 4) Model Training

5) Model Testing

A. Data Collection

The first step in the process of implementation is the data collection step. In this step, the right dataset is chosen so as to perform further computations. In the case of movie recommendation system movielens, the dataset is taken from the kaggle website. The dataset consists of 100,000 movie rating from (1-5). Further, there are 943 users and 1682 no movies. With this information, further computations are done using the Python programming language.

B. Data Preparation

The second step in the process of implementation is the data preparation step. In this step data preprocessing is done. It represents the utility matrix which tells which user rated which movie. This is done by first separating the user data and movie data into the separate data frames. Then, using both the data frames, utility matrix is created.

C. Data Creation

The third step in the process of implementation is the data creation step. In this step, the K-Means clustering model is applied. The right number of clusters is chosen using the WCSS method. After choosing the right no of cluster movies are divided into clusters by applying the K-Means Clustering model. This leads to the creation of utility clustered matrix.

D. Data Training

The fourth step in the process of implementation is the data training. In this step normalization of utility clustered matrix is done. Then the similarity between the users is calculated using the Pearson Correlation Matrix. Then, using the KNN [18] prediction for the movie ratings for top N users is done.

E. Data Testing

The fifth step in the process of implementation is the data training. In his step prediction for the movie, the rating is done for the test users. This is done for the evaluation of our model, by using some evaluation metric.

VI. WORKING AND RESULTS OF PROPOSED SYSTEM

The proposed system working is discussed using the following steps:

Step 1: In this step the user information is taken as input the userID and information such as gender, age, pin code as shown in Figure 3.

```
In [7]: input("Enter the user Information")
Enter the user Information:944, 21, 'M', 'student', 110018
```

Fig. 3: Utility Matrix

Step 2: Then using the numpy and pandas library the raw data is preprocessed into separate data frames as shown in

figure 4.

	0	1	2	3	4	5	6	7	8	9	10	11	12
0	5.000	3.000	4.000	3.000	5.000	5.000	4.000	1.000	5.000	3.000	2.000	5.000	5.000
1	4.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	2.000	0.000	0.000	0.000
2	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
3	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	4.000	0.000	0.000
4	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
5	4.000	0.000	0.000	0.000	0.000	0.000	2.000	4.000	4.000	0.000	0.000	4.000	2.000
6	0.000	0.000	0.000	5.000	0.000	0.000	5.000	5.000	5.000	4.000	3.000	5.000	0.000
7	0.000	0.000	0.000	0.000	0.000	0.000	3.000	0.000	0.000	0.000	3.000	0.000	0.000
8	0.000	0.000	0.000	0.000	0.000	0.000	4.000	0.000	0.000	0.000	0.000	0.000	0.000
9	4.000	0.000	0.000	4.000	0.000	0.000	0.000	0.000	4.000	0.000	4.000	5.000	3.000
10	0.000	0.000	0.000	0.000	0.000	0.000	0.000	4.000	5.000	0.000	2.000	2.000	0.000
11	0.000	0.000	0.000	5.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000	0.000
12	3.000	3.000	0.000	5.000	1.000	0.000	2.000	4.000	3.000	0.000	1.000	5.000	5.000
13	0.000	0.000	0.000	0.000	0.000	0.000	0.000	5.000	0.000	4.000	0.000	0.000	4.000
14	1.000	0.000	0.000	0.000	0.000	0.000	1.000	0.000	4.000	0.000	0.000	0.000	1.000
15	5.000	0.000	0.000	5.000	0.000	0.000	0.000	5.000	0.000	5.000	3.000	0.000	0.000
16	0.000	0.000	0.000	0.000	0.000	0.000	0.000	4.000	0.000	0.000	0.000	0.000	0.000
17	5.000	0.000	0.000	3.000	0.000	5.000	0.000	5.000	5.000	0.000	0.000	5.000	5.000
18	0.000	0.000	0.000	0.000	0.000	0.000	0.000	5.000	0.000	0.000	0.000	0.000	0.000

Fig. 4: Utility Matrix

Step 3:The method Within Clustered Sum of Squared is used to find the right no clusters so that K-means clustering can be applied in the movie as shown in figure 5.

	0	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18
0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
1	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	0	0	0	0	0	0	0	0	0	0	0	0	0	0
3	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
4	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
5	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
6	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
7	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
8	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
9	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
10	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
11	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
12	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
13	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
14	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
15	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
16	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
17	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
18	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1

Fig. 5: Utility Matrix

Step 4:After applying K-means clustering a utility clustered matrix is build which defines average rating the user gives to each cluster as shown in figure 6.

	0	1
0	3.392	3.923
1	3.571	3.871
2	2.680	3.211
3	4.182	5.000
4	2.043	2.720
5	3.558	3.711
6	3.815	4.280
7	3.657	4.357
8	4.429	3.880
9	4.153	4.283
10	3.217	3.747
11	4.227	4.474
12	2.885	3.528
13	4.495	4.171
14	2.487	3.261
15	4.380	4.367
16	3.857	3.304
17	3.667	4.480
18	3.580	3.758

Fig. 6: Utility Clustered Matrix

Step 5:Using the utility clustered matrix and Pearson correlation similarity between the users are calculated as

shown in figure 7.

	0	1	2	3	4	5	6	7	8	9	10	11	12
0	0.000	1.000	1.000	1.000	-1.000	1.000	1.000	1.000	-1.000	1.000	1.000	1.000	1.000
1	1.000	0.000	1.000	1.000	-1.000	1.000	1.000	1.000	-1.000	1.000	1.000	1.000	1.000
2	1.000	1.000	0.000	1.000	-1.000	1.000	1.000	1.000	-1.000	1.000	1.000	1.000	1.000
3	1.000	1.000	1.000	0.000	-1.000	1.000	1.000	1.000	-1.000	1.000	1.000	1.000	1.000
4	-1.000	-1.000	-1.000	-1.000	0.000	-1.000	-1.000	-1.000	1.000	-1.000	-1.000	-1.000	-1.000
5	1.000	1.000	1.000	1.000	-1.000	0.000	1.000	1.000	-1.000	1.000	1.000	1.000	1.000
6	1.000	1.000	1.000	1.000	-1.000	0.000	0.000	1.000	-1.000	1.000	1.000	1.000	1.000
7	1.000	1.000	1.000	1.000	-1.000	1.000	1.000	0.000	-1.000	1.000	1.000	1.000	1.000
8	-1.000	-1.000	-1.000	-1.000	0.000	-1.000	-1.000	0.000	1.000	-1.000	-1.000	-1.000	-1.000
9	1.000	1.000	1.000	1.000	-1.000	1.000	1.000	1.000	-1.000	0.000	1.000	1.000	1.000
10	1.000	1.000	1.000	1.000	-1.000	1.000	1.000	1.000	-1.000	1.000	0.000	1.000	1.000
11	1.000	1.000	1.000	1.000	-1.000	1.000	1.000	1.000	-1.000	1.000	1.000	0.000	1.000
12	1.000	1.000	1.000	1.000	-1.000	1.000	1.000	1.000	-1.000	1.000	1.000	1.000	0.000
13	1.000	1.000	1.000	1.000	-1.000	1.000	1.000	1.000	-1.000	1.000	1.000	1.000	1.000
14	1.000	1.000	1.000	1.000	-1.000	1.000	1.000	1.000	-1.000	1.000	1.000	1.000	1.000
15	1.000	1.000	1.000	1.000	-1.000	1.000	1.000	1.000	-1.000	1.000	1.000	1.000	1.000
16	1.000	1.000	1.000	1.000	-1.000	1.000	1.000	1.000	-1.000	1.000	1.000	1.000	1.000
17	1.000	1.000	1.000	1.000	-1.000	1.000	1.000	1.000	-1.000	1.000	1.000	1.000	1.000
18	1.000	1.000	1.000	1.000	-1.000	1.000	1.000	1.000	-1.000	1.000	1.000	1.000	1.000

Fig. 7: Similarity Matrix

Step 6:Finally KNN uses the utility clustered matrix and similarity to predict the movies for input user as shown in figure 8.

```

.....
Finding the similar types of movies
Men With Guns (1997)
Star Trek III: The Search for Spock (1984)
Blood For Dracula (Andy Warhol's Dracula) (1974)
Nutty Professor, The (1996)
Terminator 2: Judgment Day (1991)
MURDER and murder (1996)

```

Fig. 8: Output

Root Mean Square Error (RMSE) [26] can be calculated from the Equation (1).

$$RMSE = \sqrt{\frac{\sum_{i=1}^n (\hat{y}_i - y_i)^2}{n}} \quad (1)$$

Table 1 shows the results obtained after using different no of clusters for the proposed technique.

TABLE I: Results

K-means + KNN	
Number of clusters	Root Mean Squared Error
19	2.504990
18	2.375555
17	2.337194
16	2.416212
15	2.256299
14	2.080751
13	1.994332
12	1.928682
11	1.861167
10	1.820095
9	1.625027
8	1.493939
7	1.441855
6	1.439451
5	1.269583
4	1.166091
3	1.141065
2	1.081648

A. Comparison with Existing Technology

The table 2 and 3 compares the result of the proposed system with the existing technique. These tables shows a comparison of RMSE with the existing technique i.e. cuckoo search. It is seen from the tables that for the existing technique the RMSE value is 1.23154 for cluster equal to 68, RMSE value using proposed technique is 1.233 to 19 clusters and RMSE value using proposed technique is 1.081648 to 2 clusters.

TABLE II: RMSE in Proposed Technique

Root Mean Squared Error	No. of Cluster
1.23154	68

TABLE III: RMSE in Existing Technique

Root Mean Squared Error	No. of Cluster
1.2333	19
1.081648	2

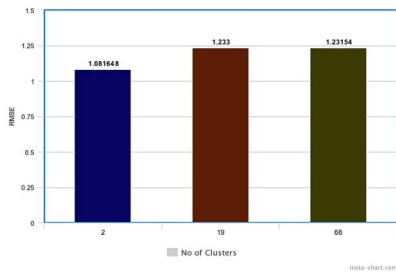


Fig. 9: Comparison Graph with the Existing Technique

Figure 9 compares the RMSE value for existing technique with the RMSE value of the proposed technique. The X-axis represents the No of Clusters and Y-axis represents the RMSE values. It is seen from the graph that for the existing technique the RMSE value is 1.23154 for cluster equal to 68, RMSE value using proposed technique is 1.233 to 19 clusters and RMSE value using proposed technique is 1.081648 to 2 clusters.

VII. CONCLUSION

Machine learning is a method of data analysis that automates analytical model building. It is a branch of artificial intelligence based on the idea that systems can learn from data, identify patterns and make decisions with minimal human intervention[25]. In this proposed system a movie recommender system is built using the K-Means Clustering and K-Nearest Neighbor algorithms. The data are taken from movielens data set. The system is implemented in python programming language. It is seen that after implementing the system in the python programming language the RMSE value of the proposed technique is better than the existing technique. It is also seen that the RMSE value of the proposed system is achieving the same value as the existing technique but with less no of clusters. The proposed work can be improved using more

data sets. Sentimental Analysis concept can be used in future to enhance the efficiency of movie recommendation system, so the model can be tuned to accommodate more situations. In future, individual characteristic may be removed which is hidden in the recommendation of the users.

REFERENCES

- [1] Goel A., Khandelwal D., Mundhra J., Tiwari R. (2018) Intelligent and Integrated Book Recommendation and Best Price Identifier System Using Machine Learning. In: Bhateja V., Coello Coello C., Satapathy S., Pattanaik P. (eds) Intelligent Engineering Informatics. Advances in Intelligent Systems and Computing, vol 695. Springer, Singapore
- [2] Bao J., Zheng Y. (2017) Location-Based Recommendation Systems. In: Shekhar S., Xiong H., Zhou X. (eds) Encyclopedia of GIS. Springer, Cham
- [3] Chavarriaga O., Florian-Gaviria B., Solarte O. (2014) A Recommender System for Students Based on Social Knowledge and Assessment Data of Competences. In: Rensing C., de Freitas S., Ley T., Muñoz-Merino P.J. (eds) Open Learning and Teaching in Educational Communities. EC-TEL 2014. Lecture Notes in Computer Science, vol 8719. Springer, Cham
- [4] F.O.Isinkaye et. al, Recommendation systems: Principles, methods and evaluation, Egyptian Informatics Journal Volume 16, Issue 3, November 2015, Pages 261-273
- [5] H. Drachsler, T. Bogers, R. Vuorikari, K. Verbert, E. Duval, N. Manouselis, G. Beham, S. Lindstaedt, H. Stern, M. Friedrich, et al. Issues and considerations regarding sharable data sets for recommender systems in technology enhanced learning. Procedia Computer Science, 1(2):2849–2858, 2010.
- [6] ÁlvaroTejeda-Lorente, A quality based recommender system to disseminate information in a university digital library,Information Sciences Volume 261, 10 March 2014, Pages 52-69
- [7] Trang Tran, T.N., Atas, M., Felfernig, A. et al. J Intell Inf Syst (2018) 50: 501. <https://doi.org/10.1007/s10844-017-0469-0>
- [8] Xin Luo, Mengchu Zhou, Yunni Xia, and Qingsheng Zhu, An Efficient Non-Negative Matrix-Factorization-Based Approach to Collaborative Filtering for Recommender Systems, IEEE Transactions on Industrial Informatics (Volume: 10 , Issue: 2 , May 2014)
- [9] Badsha, S., Yi, X. Khalil, I. Data Sci. Eng. (2016) 1: 161. <https://doi.org/10.1007/s41019-016-0020-2>
- [10] Farman Ullah, Ghulam Sarwar, Sung Chang Lee, Yun Kyung Park, Kyeong Deok Moon, Jin Tae Kim, Hybrid Recommender System with Temporal Information, The International Conference on Information Network, 2012, DOI: 10.1109/ICIN.2012.6164413
- [11] Jing Jiang, Jie Lu, Guangquan Zhang, Guodong Long, Scaling-up Item-based Collaborative Filtering Recommendation Algorithm based on hadoop, 2011 IEEE World Congress on Services, 4-9 July 2011, 10.1109/SERVICES.2011.66
- [12] Vibhor Kanta , Kamal K. Bharadwaj, Enhancing Recommendation Quality of Content-based, Filtering through Collaborative Predictions and Fuzzy Similarity Measures, Procedia Engineering Volume 38, 2012, Pages 939-942.
- [13] Jiang Z., Zang W., Liu X. (2016) Research of K-means Clustering Method Based on DNA Genetic Algorithm and P System. In: Zu Q., Hu B. (eds) Human Centered Computing. HCC 2016. Lecture Notes in Computer Science, vol 9567. Springer, Cham
- [14] Sanjoy K. Sinha a, Nan M. Lairdb, Garrett M. Fitzmaurice, Multivariate logistic regression with incomplete covariate and auxiliary information, Elsevier, 2010
- [15] D.A. Adeniyi, Z. Wei, Y. Yongquan, Automated web usage data mining and recommendation system using K-Nearest Neighbor (KNN) classification method, Saudi Computer Society, King Saud University, October 2014
- [16] Rahul Kataria , Om Prakash Verma, An effective collaborative movie recommender system with cuckoo search, Egyptian Informatics Journal, 2016, Volume 18, Issue 2, July 2017, Pages 105-112 <https://doi.org/10.1016/j.eij.2016.10.002>

- [17] Czarnowski I., Jdrzejowicz P. (2008) Data Reduction Algorithm for Machine Learning and Data Mining. In: Nguyen N.T., Borzowski L., Grzech A., Ali M. (eds) New Frontiers in Applied Artificial Intelligence. IEA/AIE 2008. Lecture Notes in Computer Science, vol 5027. Springer, Berlin, Heidelberg
- [18] Vejmelka M., Hlaváčková-Schindler K. (2007) Mutual Information Estimation in Higher Dimensions: A Speed-Up of a k-Nearest Neighbor Based Estimator. In: Beliczynski B., Dzieliński A., Iwanowski M., Ribeiro B. (eds) Adaptive and Natural Computing Algorithms. ICANNGA 2007. Lecture Notes in Computer Science, vol 4431. Springer, Berlin, Heidelberg
- [19] Duarte D., Ståhl N. (2019) Machine Learning: A Concise Overview. In: Said A., Torra V. (eds) Data Science in Practice. Studies in Big Data, vol 46. Springer, Cham
- [20] Fan Y., Dong L., Sun X., Wang D., Qin W., Aizeng C. (2018) Research on Auto-Generating Test-Paper Model Based on Spatial-Temporal Clustering Analysis. In: Huang DS., Jo KH., Zhang XL. (eds) Intelligent Computing Theories and Application. ICIC 2018. Lecture Notes in Computer Science, vol 10955. Springer, Cham
- [21] Kushwaha N., Pant M. (2019) A Teaching-Learning-Based Particle Swarm Optimization for Data Clustering. In: Tanveer M., Pachori R. (eds) Machine Intelligence and Signal Analysis. Advances in Intelligent Systems and Computing, vol 748. Springer, Singapore
- [22] Howley T., Madden M.G., O'Connell ML., Ryder A.G. (2006) The Effect of Principal Component Analysis on Machine Learning Accuracy with High Dimensional Spectral Data. In: Macintosh A., Ellis R., Allen T. (eds) Applications and Innovations in Intelligent Systems XIII. SGAI 2005. Springer, London
- [23] Hartigan, J.A., Wong, M.A.: Algorithm as 136: A k-means clustering algorithm. Journal of the Royal Statistical Society. Series C 28(1), 100–108 (1979)
- [24] J. Laaksonen and E. Oja, Classification with learning k-nearest neighbors, Proceedings of International Conference on Neural Networks (ICNN'96), 3-6 June 1996, 10.1109/ICNN.1996.549118
- [25] Anna L. Buczak et al, A Survey of Data Mining and Machine Learning Methods for Cyber Security Intrusion Detection, IEEE Communications Surveys Tutorials (Volume: 18 , Issue: 2 , Secondquarter 2016)DOI: 10.1109/COMST.2015.2494502
- [26] Author(s) C. Kilgus et.al, Root-Mean-Square Error in Encoded Digital Telemetry, IEEE Transactions on Communications,(Volume: 20 , Issue: 3 , Jun 1972),DOI: 10.1109/TCOM.1972.1091174
- [27] Alexandra L'Heureux et.al, Machine Learning With Big Data: Challenges and Approaches, IEEE Access (Volume: 5) Page(s): 7776 - 7797,DOI: 10.1109/ACCESS.2017.2696365
- [28] Bernard Marr,What Is The Difference Between Artificial Intelligence And Machine Learning,6 December 2016,8 Feb 2018 accessed from:<https://www.forbes.com/sites/bernardmarr/2016/12/06/what-is-the-difference-between-artificial-intelligence-and-machine-learning/7e94a51a2742>
- [29] Nick Mccrea,An Introduction to Machine Learning Theory and Its Applications: A Visual Tutorial with Examples,Aug 2016,Feb 2018 accessed from: <https://www.toptal.com/machine-learning/machine-learning-theory-an-introductory-primer>
- [30] Miroslav Kubat, An Introduction to Machine Learning,<https://doi.org/10.1007/978-3-319-63913-0>, Springer International Publishing AG 2017,Print ISBN 978-3-319-63912-3, Online ISBN 978-3-319-63913-0
- [31] Priyadharshini,Machine Learning: What it is and Why it Matters,March 2018,April 2018 accessed from: <https://www.decipher.com/machine-learning-matters/>
- [32] Fabien Dubosson et.al,A Python Framework for Exhaustive Machine Learning Algorithms and Features Evaluations,2016 IEEE 30th International Conference on Advanced Information Networking and Applications (AINA),23-25 March 2016,ISSN: 1550-445X,DOI: 10.1109/AINA.2016.160
- [33] Jianpeng Qi et. al, K*-Means: An Effective and Efficient K-Means Clustering Algorithm, 2016 IEEE International Conferences on Big Data and Cloud Computing (BDCloud), Social Computing and Networking (SocialCom), Sustainable Computing and Communications (SustainCom) (BDCloud-SocialCom-SustainCom),8-10 Oct. 2016,DOI: 10.1109/BDCloud-SocialCom-SustainCom.2016.46
- [34] Anna L. Buczak,A Survey of Data Mining and Machine Learning Methods for Cyber Security Intrusion Detection, IEEE COMMUNICATIONS SURVEYS TUTORIALS, VOL. 18, NO. 2, SECOND QUARTER 2016.
- [35] Nguyen N.T., Rakowski M., Rusin M., Sobecki J., Jain L.C. (2007) Hybrid Filtering Methods Applied in Web-Based Movie Recommendation System. In: Apolloni B., Howlett R.J., Jain L. (eds) Knowledge-Based Intelligent Information and Engineering Systems. KES 2007. Lecture Notes in Computer Science, vol 4692. Springer, Berlin, Heidelberg
- [36] Ko SK. et al. (2011) A Smart Movie Recommendation System. In: Smith M.J., Salvendy G. (eds) Human Interface and the Management of Information. Interacting with Information. Human Interface 2011. Lecture Notes in Computer Science, vol 6771. Springer, Berlin, Heidelberg
- [37] Wei D., Junliang C. (2013) The Bayesian Network and Trust Model Based Movie Recommendation System. In: Du Z. (eds) Intelligence Computation and Evolutionary Computation. Advances in Intelligent Systems and Computing, vol 180. Springer, Berlin, Heidelberg
- [38] Mishra N., Chaturvedi S., Mishra V., Srivastava R., Bargah P. (2017) Solving Sparsity Problem in Rating-Based Movie Recommendation System. In: Behera H., Mohapatra D. (eds) Computational Intelligence in Data Mining. Advances in Intelligent Systems and Computing, vol 556. Springer, Singapore
- [39] Das D., Chidananda H.T., Sahoo L. (2018) Personalized Movie Recommendation System Using Twitter Data. In: Pattnaik P., Rautaray S., Das H., Nayak J. (eds) Progress in Computing, Analytics and Networking. Advances in Intelligent Systems and Computing, vol 710. Springer, Singapore
- [40] Lops P., de Gemmis M., Semeraro G. (2011) Content-based Recommender Systems: State of the Art and Trends. In: Ricci F., Rokach L., Shapira B., Kantor P. (eds) Recommender Systems Handbook. Springer, Boston, MA
- [41] Hatami, M., Pashazadeh, S.(2014)Improving results and performance of collaborative filtering-based recommender systems using cuckoo optimization algorithm,Int J Comput Appl Volume 88, Pages 46-51
- [42] Z. Huang, D. Zeng, H. Chen (2007) A comparison of collaborative-filtering algorithms for e-commerce IEEE Intell Syst, 22, pp. 68-78
- [43] R. Burke (2007) Hybrid web recommender systems, Adapt Web, pp. 377-408, 10.1007/978-3-540-72079-9_12