# Soft Clustering Based Recommender System

**Nishant Gupta[1], Deepanshu Jain[2], Shreyansh Jain[3], Achal Kaushik[4]**

[1, 2, 3 ,4]*Department of Computer Science and Engineering, Bhagwan Parshuram Institute of Technology, Delhi*

[4]achalkaushik@gmail.com

*Abstract- Today, the use of internet has been growing and therefore browsing on the internet has also been on an increase. To find the relevant information from the massive information available online is a difficult task and therefore, we aim at providing relevant information to a user by eliminating the irrelevant information. This will save user's time as well as help them getting their needed information quickly without any delay. The web page recommendation is a challenging problem which uses various techniques for solving and has become an intelligent support system for organizations worldwide. Fuzzy C-Mean Clustering algorithm is one such approach which is used to improve the web page recommendations. In our work, we have used sequential ordering aspect for forming cluster and employed soft clustering technique in order to improve the user experience. For each user involved in searching through internet our model provides a better prediction of the web page that a user may visit. We have experimentally evaluated our system and found that it provides better results in comparison with the existing systems.*

*Keywords--- Fuzzy C-Mean Clustering, recommendations system, sequential aspect*

## I. INTRODUCTION

Web page recommendation is specifying which could be the following page a user might visit depending on his foregoing pages. A lot of information is available on the internet and choosing what is relevant is itself a quite difficult task [1]. As e-commerce has been growing at a fast pace and this has allowed the companies to provide user various options on a single platform. These systems are like decision support system that provides users with customized information tailored to their requirements and needs. The recommender systems take into account the various information of customer like social relationships, various buying behaviors, their likes or dislikes, etc. [2]. Most of the e-commerce websites uses recommender systems that are deployed at the back end. Various machine learning algorithms have been employed for improving recommender systems like data mining, heuristic and finding associate patterns, probabilistic models. Various recommender systems have been developed like IMDb for movies, VERSIFI is used for news, amazon for various products etc. The data that are analyzed for making recommendations include various url's searched, various web page content, hyperlinks involved in searching, etc. [1]. The web page recommendations have become very handy popular tool being used by a number of organizations for organizing the recommendations according to the user needs and preferences. The web page recommendation is a challenging problem which uses various techniques for solving and has become an intelligent support system for organizations worldwide [3]. As it makes the website user friendly by accounting a user visit i.e. if a user has visited the webpages contented as entertainment, electronics, footwear, etc. This data will

be used as input for web page recommender system for predicting the next web page the user might be interested in [4].

As the recommendations make user search the desired product quickly and easily and hence increasing the profit of various e-commerce companies [5]. But most of the time, the sequential aspect of recommender system is not considered like in various models like Markov models for making recommendations. Also in some other models like probabilistic models the problem of switching across various probabilities in various categories of web is considered as a priority as it requires knowledge of the expert system in that domain. Also the estimation of exact probability is an ongoing problem and has not been properly taken into consideration [6]. We have used fuzzy c-mean clustering algorithm for our recommender systems which divides the data into various clusters [7]. The outcome of our system is more desirable as some other algorithms like K–means where each datum is assigned to only one cluster as they uses hard clustering whereas we used soft clustering and assigned our datum to more than one cluster centers as they are allotted belongingness to each cluster center [8]. There is increasing development in world wide web, the Internet has modified from a static platform for transferring and sharing information to a social interactive platform with greater participation and communication between people of various regions each having different ideas and choices [9]. We try to evaluate the ordering information by observing the browsing behaviour of a user for making predictions for that user's next visit [10]. The essence of our work is mainly presented as follows:

- We first obtain soft clusters by using fuzzy c-mean clustering methodology.
- Then, we determine the mean value of each cluster centre and selected top n-clusters for evaluation.
- We predict the pages for aimed user that he may want to go to by computing substantial weight of each page.
- Finally, we determined the correctness of our model by using msnbc dataset.

The test is mainly performed on huge data of msnbc where 5000 users were evaluated by considering the six pages they have browsed; we forecast the seventh page he may likely to visit [11]. The rest of the paper is organized as follows; section 3 provides the understanding and backdrop of our experimentation point, section 4 explains our system model and in section 5 we have evaluated the various aspects and calculations of our model and finally section 6 the result section.

## II. RELATED WORK

A lot of work has been done in order to optimize the performance of recommender system so that they can produce more accurate predictions. For this purpose, a lot of data mining and clustering methods have been proposed. In [12] authors have proposed the various versions for K-mean clustering that can be used for a finite dataset by constructing various cases for least square error.

In [13] it was mainly researched on the problem of extremely large and redundant information that is available on the internet and tried to use various data mining techniques to obtain some useful information i.e. tried to use pattern recognition for obtaining data relevant for a user. A system was proposed for estimating irregularities that are associated with magnetic image reasoning by employing fuzzy Cartesian and making some change in the

objective function to enumerate for these irregularities [3]. Authors researched a possibilistic fuzzy c-mean algorithm in which we generate the typical values in addition to the membership values. Thus, making the model a combination of fuzzy c-mean (FCM) and possibilistic c-mean algorithms as it resolves some of the limitations of FCM [14]. Authors evaluated one the most essential element of the FCM algorithm which is the weight of the system and determined that the experiment results an outcome are highly variable when we increase or decrease the values of this component [3]. It was proposed that system mainly adds distributed information in the algorithm for eliminating the unevenness generated due to noise and increasing the uniformity of the outcome in comparison to various other prevailing systems by using partitioning of images that are noisy [15].

The spatial function is mainly generated by examining each pixel in the area [15]. In order to deal with the problem of excessive products that are available on e- commerce website and tried to customize it to provide recommendation to a user based on the user, taste and wants of that user thus preventing the system from irrelevant browsing [7]. The user browsing history is evaluated in order to determine the user preference and find out the sites and pages the user has browsed so that using this information we try and calculate the value for each page and the pages having the highest values are the next page that the user might visit [16].
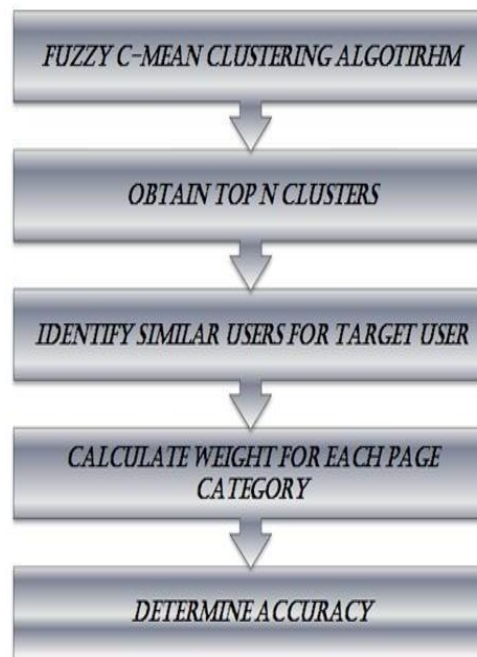


Fig 1: Flow chart of the proposed Soft Clustering based Recommender System

A model is proposed which mainly takes into consideration the semantic connection between various items that are being browsed and it makes use of LDA (Latent Dirichlet Allocation) for this determination [17]. The semantic information will help in discovering the users' social information, his likings, etc. A lot of recommender systems are prevailing but the common problem in all these systems is that they do not consider the ordering information. Therefore, this ordering aspect is mainly considered in our model for forming cluster using various matrices. We employed soft clustering in order to make the prediction of the web page that a user may visit.

### III.     PROPOSED SOFT CLUSTERING BASED RECOMMENDER SYSTEM

There are a lot of models that have been developed involving item rating prediction for websites like amazon, ebay and many others. Fuzzy C-Mean clustering mainly divides the information into various clusters depending on some parameters and then we take only n-top clusters. Each page weight calculation is carried out by finding users having same searching pattern as that of the user for which we are making a recommendation. This calculation helps us determining which page will be most suited for the prediction. So, by using this approach, we predict the seventh page visit for each user by taking into account the six page visit for each user. Thus, our aim is to improve the user experience of each user involved in searching through internet [12].

In web page predictions, we tried to perform pattern recognition by incorporating two aspects, clustering and categorization. For the feasible system outcome, the process of categorizing any information requires the system must itself have the appropriate knowledge and training which can be considered in real situations. Then only the system will make correct predictions for the user and this is proposed in our model.

### A.  FUZZY C-MEAN CLUSTERING

The main aspect is to first divide the data into clusters in order to eliminate the redundant information that the user has to browse most frequently [9]. The cluster formation is done by training a system about the various likings and interest patterns of the user browsing through the dataset and thus forming the clusters. The process is explained using a flow diagram in figure 1. In our model, we had formed clusters by using the method of soft clustering. Soft clustering is same as that of hard clustering, but the only difference here is that the user or datum is a part of more than one cluster. Further, each of the users has a membership level assigned which emphasizes on the degree to which that data is linked to each of the cluster. Thus, Fuzzy c-mean algorithm is basically a clustering algorithm where the basic idea to categorize the data into clusters for the ease of computation by providing association values to these clusters [13]. Provided a group of data, our system use the FCM algorithm to generate a tally of C-cluster centers $C = \{c1 \dots cc\}$ and then calculating $W = wij \in [0,1], i = 1 \dots n, j = 1 \dots c$ This wij is the partition matrix and are proposed algorithm try to decrease the objective function as is done in other clustering algorithms.

$$argmin_c \sum_{i=1}^{n} \sum_{j=1}^{c} w_{ij}^{m} \left\| x_i - c_j \right\|^2 \qquad (1)$$

where $\quad w_{ij}^{m} = \dfrac{1}{\sum_{k-1}^{c} \left\{ \dfrac{\left\| x_i - c_j \right\|}{\left\| x_i - c_k \right\|} \right\}^{\frac{2}{m-1}}}$

Here, $c_j$ represents each cluster center value, $x_i$ stands for each data element taken into consideration and $w_{ij}$ measures the extent with which each data element $x_i$ is related to the cluster $c_j$, $\sum$ represent summation of all the data elements, m represents the degree of fuzziness in our algorithm. The addition of membership value $w_{ij}$ and the fusilier $m \geq r$ are

the parameters that make our algorithm different from k-means algorithm. The value of m is significantly important in determining membership value $w_{ij}$ and m tells us about the degree of fuzziness in a cluster and if the value of m is larger than $w_{ij}$ values will be small and the clusters fuzziness will increase, But we want to make a set more crisp or decrease the degree of fuzziness therefore we take the value of m to be small. For m=1, $w_{ij}$ may coincide to 0 or 1 therefore, we use m=2, which is most frequently used in experiments [3].

### B. SELECTION ON TOP 'N' CLUSTERS

As the above algorithm has provided us with a number of "c" clusters thus eliminating our need to model the entire data and we can now work on these clusters only. The next step is to determine the clusters where our target user has a belonging to, as there can be more than one cluster. We are using a soft clustering technique where the user mainly has a partial membership which is not the case with hard clustering. Then, we will calculate the mean value of all the cluster centers in order to determine the feasible clusters which will have some influence on our recommendation and will affect the result in a positive manner. Thus, we take only n-clusters which are at the top.

### C. FINDING USERS SIMILAR TO TARGET USER

Based on the target user, we try and determine which all users are alike to our target users as they are the ones who will have higher influence on the predictions and the recommended mode for the target user. For this, we computed and obtained a partition matrix "P" for the alike users and then a matrix NAT is mainly constructed which is the conclusive matrix indicating the list or ordering of the similar users.

### D. WEIGHT CALCULATION FOR EACH PAGE CATEGORY

In the next step, we calculate the significant weight of all the pages and this is the most important aspect of our model. This weight is basically the indication of the possibility of whether that page is relevant for the prediction or not. We have to take a total of seventeen categories of pages, such that each page is assigned a number between 1 and 17 for indicating that page. The categories that are used in our system are as follows: "front page", "news", "tech", "local", "opinion", "on-air", "misc", "weather", "health", "living", "business", "sports", "summary" , "bbs", "travel", "msn-news" and "msnsports". Now, we compute the weight for each page which requires the use of these matrices 'A', 'B', 'C':

- Matrix A indicates the integral value of the number of times a particular page is visited by our target user.

- Matrix B indicates the integral value of the number of times that particular page has been visited by a user similar to our target users.
- Matrix C indicates that if a page has been visited by the target user at least once or not.

Thus, after computing the data obtained through the above matrices, we found out that the two matrices A and B indicates which pages are best for making recommendations where C indicates the web pages which are not suitable for the recommendation.

The mathematical statements are presented below:

if (msndata(user, i)==n)

$$a(n) = a(n) + 1; \tag{2}$$

if (nat(i, j)==n)

$$b(n)=b(n)+1 \tag{3}$$

if (nat(i, j)==n)

    if (c(n)==0)

$$c(n)=c(n)+1 \tag{4}$$

Here NAT which provides us with the list of similar users for our aimed user, "n" is used to indicate the sequence number assigned to each page and msn data indicates the dataset on which experimental computation is performed.

$$weight(i) \ = (A[i] + B[i])/(size - C[i]) \tag{5}$$

Here, mainly "i" denotes each category of the page and each of the parameter values can be derived from the above equations.

### E. CALCULATION OF ACCURACY

Accuracy is mainly the proportion of count of accurate recommendations to the count of overall recommendations of a system. Accuracy is the ratio of the number of accurate recommendations to the number of overall recommendations. Here, we mainly predict the successive pages that the user may like to visit and then evaluate our prediction with the real

pages that the user visited. So, in our recommendation mainly matches with the real page that the user has visited we define it as a hit otherwise, we define it as a miss.

$$Accuracy = (no\ of\ hits) / (no\ of\ hits + no\ of\ miss) \qquad (6)$$

Thus, the proposed algorithm is defined as follows:

- First, we obtain soft clusters by applying fuzzy c-mean algorithm for clustering msnbc dataset.
- Then, based on the mean value of the cluster centre, we take prime 'N' clusters from all the clusters using eq.1.
- In the next step, we determine the similar users of our target users by first finalizing our target user.
- Then, matrix 'A' values are determined from eq. 2.
- In the same way, the matrix 'B' and 'C' values are accessed from eq. 3 and eq. 4.
- Here, we obtain the matrix page value based on the weight of each page obtained by eq. 5 and then the page is recommended to the user.
- At the end, we compute the accuracy of our proposed work.

We observed the accuracy of our proposed model and are mainly a proportion of admissible responses to all non-admissible and non-trustworthy responses.

### F. DATASET

We mainly worked on the real world dataset of MSNBC for performing our computations and research. In this dataset, there were total 5000 entries of the user and for each user 6 entries were maintained. We partitioned our dataset into two segments as show: Training data=60%&Testing data=40%. Table 1: Accuracy comparison

When the context groups are formed, the textual descriptions are given weight on the basis of their frequency. It is given less weight if it occurs very frequently and vice versa. Each test image $X_0$ is provided with the same conditional probability as context group.

## IV.    MATHEMATICAL MODEL FOR CONTEXT BASED AUTOMATIC IMAGE ANNOTATION

In Suppose there are n number of visual units such as $V_n$ in an image and m number of textual description such as $H_m$. Let there are $CC$ number of context categories where $T \in CC$ corresponds to one context group. Each training image will belong from the one of these T groups.
By picking a context group with conditional probability over test image $X_0$ i.e. $W(H|X_0)$. By selecting a training image $X_t$ within the training set TS with the probability $W(H|X_h)$

for i=1,2,.....,n

2.1 Pick a visual unit $V_i$ having conditional probability $W_R(.|X_h)$

For j=1,2,.....,m

2.2  By selecting a word $h_j$ from conditional probability $W_T(.|X_h)$

The main aim of the proposed approach is to enhance probability metric V and T over the training image $X_h$

$$P(H, V|X_h) = \sum_{H \in CC} P(H|X_h) \sum_{X_h \in HS} P(X_h|H) \prod_{j \in m} w_T(H_j|X_h) \prod_{i \in n} W_R(vI/X_t) \qquad (5)$$

The $W_H(H_j|X_t)$ (Bernoulli distribution) is defined as:

$$W_H(H_j|X_h) = \frac{\mu \delta_{H_j} + N_{H_j}}{\mu + N_H} \qquad (6)$$

 where, $A_{H_j}$ is the members of T with word $T_j$ in their description

$A_{H_j}$ is members of $T_j$

$\delta_{H_j}$ is set to be 1 if description of the image $X_t$ has word $T_j$ in it $\mu$ is empirically selected constant

$W_R(V_i|X_t)$ is the density estimate to generate the visual unit $V_i$ for the training image $X_t$.

Gaussian kernel is employed for this density estimate. Suppose if the visual units of the training image $X_t$ are $\{VT_1, VT_2, ... ..., VT_n\}$ then

$$(V_i|X_h) = \frac{e^{(-(V_i - VH_n)^T (\Sigma V_i - VH_n)^{-1})}}{\sqrt{2\pi|\Sigma|}} \qquad (7)$$

where $\sum$ is the covariance matrix.

## V.    RESULTS AND DISCUSSION

The proposed algorithm is tested on Corel-10k dataset. In this paper results are compared without using filter and with using filter. In both the cases the tucker decomposition level is taken as 3. To check the effectiveness of the algorithm the comparison has been made between precision, recall, and accuracy. Since the dataset contains 10k image therefore complete dataset are divided into small context groups for easy and fast analysis purposes. We have also checked the efficiency of the algorithm by taking different percentage combination of training and testing data i.e. 70% and 30%, 50% and 50% etc.  The results thus obtained are as in the table I-IV.

- Some tags are correct and retrieved
- Some tags are incorrect and retrieved
- Some tags are not retrieved giving some random value.

Table I: Results analysis of different context groups

| Group's name (Training images-Testing images) | No. of correct and retrieved tags in proposed approach | No. of correct and retrieved tags in base approach | Total tags that should be correct and retrieved |
|---|---|---|---|
| New 17(70-30) | 59 | 45 | 60 |
| New 19(70-30) | 51 | 44 | 60 |
| New 21(50-50) | 101 | 83 | 113 |
| New 12(70-30) | 57 | 56 | 60 |
| New 18(50-50) | 92 | 86 | 100 |
| New 32(50-50) | 106 | 94 | 131 |
| New 23(5-5) | 11 | 11 | 12 |
| New 24(5-5) | 10 | 10 | 10 |

- **Precision:** Precision is calculated as fraction of relevant tags among retrieved tags as in equation (8).

  Precision P=D/E            (8)

  where D is number of relevant images retrieved whereas E is total number of images retrieved.

Table II: Precision table for both the cases

| Group name(Training images-Testing images) | Proposed approach Precision in % | Base approach Precision in % |
|---|---|---|
| New 17(70-30) | 98.3333 | 83.3333 |
| New 19(70-30) | 86.4407 | 81.4815 |
| New 21(50-50) | 90.9910 | 74.778 |
| New 12(50-50) | 98.2759 | 98.2456 |
| New 18(50-50) | 92.9293 | 88.6598 |
| New 32(50-50) | 79.6992 | 71.557 |

- **Recall:** Recall is calculated as fraction of total relevant tags that are retrieved equation (9).

  Recall=D/F            (9)

  where D is number of relevant images retrieved. F is the number of images that are relevant in the dataset.

As it can be seen from the Table I to IV that all the analysis parameters such as precision, recall, and accuracy has been improved many fold in the proposed approach as compared to the base approach.

Table III: Recall table for both the cases

| Group 'name (Training images-testing images) | Proposed approach recall in % | Base approach Recall in % |
|---|---|---|
| New 17(70-30) | 100 | 90 |
| New 19(70-30) | 98.0769 | 95.6522 |
| New 21(50-50) | 98.0583 | 90.2714 |
| New 12(50-50) | 96.6102 | 94.9153 |
| New 18(50-50) | 98.9247 | 97.7273 |
| New 32(50-50) | 85.4839 | 71.557 |

Accuracy is calculated as total no of correct observation divided by total no of observations. The model proposed by author Tariq et.al [2] have used single level tucker decomposition while in this paper three level tucker decomposition is used to model the base case. Even in the base case the results are better as compared to the results obtained by Tariq et. al [2]. Since all the images contain Gaussian noise by default therefore adopting the Gaussian filtering technique improves the results.

Table IV: Accuracy table for both the cases

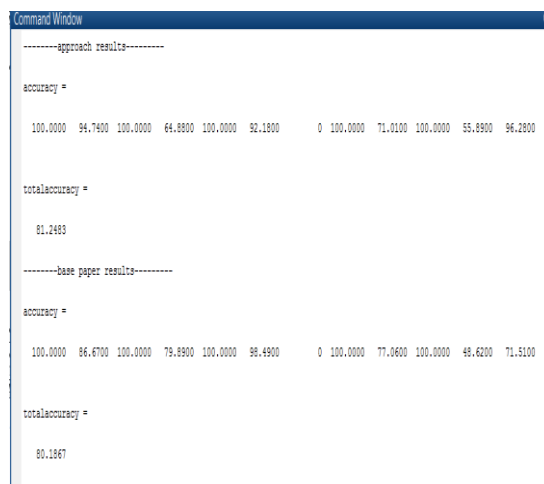| Group's name(Training images-Testing images) | Proposed approach accuracy in% | Base approach accuracy in% |
|---|---|---|
| New 17(70-30) | 88.2602 | 74.9217 |
| New 19(70-30) | 74.6538 | 70.4888 |
| New 21(50-50) | 78.1253 | 72.8023 |
| New 12(50-50) | 87.3810 | 84.7413 |
| New 18(50-50) | 83.9296 | 81.4908 |
| New 32(50-50) | 71.9161 | 64.4773 |

Fig 4: Snapshot of results obtained in MATLAB

The graphical representation of the above table can be seen in the Fig. 5 to Fig. 7. MATLAB software is used for the simulation and analysis purposes. The results obtained during simulation are as in Fig. 4.
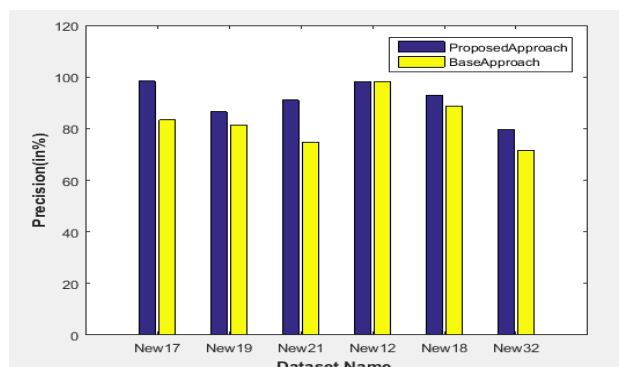


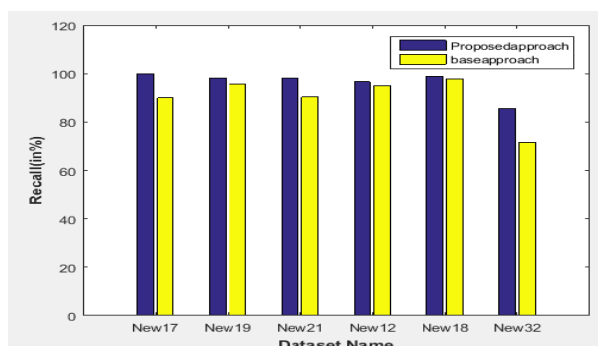Fig 5: Precision graph for different context groups of dataset



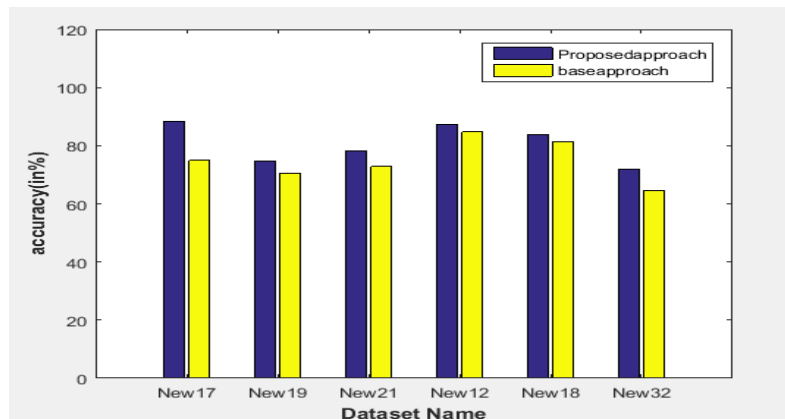Fig 6: Recall graph for different context groups of dataset

Fig 7: Accuracy graph for different context groups of dataset

## VI. CONCLUSION AND FUTURE SCOPE

Gaussian filtering based novel context estimation and tensor decomposition system is proposed. Tensor formation and context estimation is used in this research paper to minimize the semantic gap while the tensor decomposition is used to find the best correlation between the context group images. Due to minimization of semantic gap the accuracy improves significantly.3-level tucker decomposition is adopted to model the framework for better correlation among the context groups. The results are compared between filtered context groups and unfiltered context groups. The evidence of the effectiveness of the proposed algorithm can be seen from the results and discussion section. In future, deep learning technique can be used for better accuracy as well as to reduce the searching time.

## REFERENCES

[1]. BahramiS,Abadeh M. S, Automatic Image Annotation Using an Evolutionary Algorithm (IAGA),7th International Symposium on Telecommunication, pp. 320–325, 2014.

[2]. Tariq A, Foroosh H, Feature-Independent Context Estimation for Automatic Image Annotation, IEEE,pp.1958-1965,2015.

[3]. S. A. Zhu, C.-W. Ngo, and Y.-G. Jiang, "Sampling and ontologically pooling web images for visual concept learning," IEEE Trans. on MM, vol. 14, no. 4, pp. 1068–1078, 2012.

[4]. X.-R. Li, C. G. M. Snoek, and M. Worring, "Learning social tag relevance by neighbor voting," IEEE Trans. on MM, vol. 11, no. 7, pp. 1310–1322, 2009.

[5]. M. Guillaumin, T. Mensink, J. J. Verbeek, and C. Schmid, "TagProp: Discriminative metric learning in nearest neighbor models for image auto-annotation," in ICCV, 2009.

[6]. J. Jeon, V. Lavrenko, and R. Manmatha. Automatic image annotation and retrieval using cross-media relevance models. In Proceedings of the 26th annual international ACM SIGIR conference on Research and development in information retrieval, pp.119-126,2003.

[7]. V. Lavrenko, R. Manmatha, and J. Jeon. A model for learning the semantics of pictures. In Advances in neural information processing systems, 2003.

[8]. S. Feng, R. Manmatha, and V. Lavrenko. Multiple bernoulli relevance models for image and video annotation. In IEEE Computer Society Conference on Computer Vision and Pattern Recognition, pp. 1-8,2004.

[9]. S. Moran and V. Lavrenko. Optimal tag sets for automatic image annotation. In Proceedings of the British Machine Vision Conference,pp. 1-11, 2011.D. Liu, X.-S. Hua, L.-J. Yang, M. Wang, and H.-J. Zhang, "Tag ranking," in WWW,pp. 351-360, 2009. A.Rae, B.Sigurbj¨ornsson,andR.-V.Zwol, "Improving tag recommendation using social networks," in RIAO, 2010.

[10]. http://www.ci.gxnu.edu.cn/cbir/Dataset.aspx

[11]. Tran, T.-H.; Tran, X.-H.; Nguyen, V.-T.; Nguyen-An, K. Building an Automatic Image Tagger with DenseNet and Transfer Learning; IEEE: Piscataway, NJ, USA, pp. 34–41,2019

[12]. Chu, Y., et al.: Automatic image captioning based on ResNet50 and LSTM with soft attention. In: Wireless Communications and Mobile Computing, 2020

[13]. R. Subash November 2019 Journal of Physics Conference Series 1362:012096 : Automatic Image Captioning Using Convolution Neural Networks and LSTM.

[14]. Bai, S., An, S.: A Survey on automatic image caption generation. Neurocomputing 311, 291–304 (2018)

[15]. Tanti, M., Gatt, A., Camilleri, K.: What is the Role of Recurrent Neural Networks (rnns) in an Image Caption Generator? in: proceedings of the 10th International Conference on Natura l Language Generation, pp. 51–60 (2017).

[16]. Kamal, A.h., Jishan, M.a., Mansoor, N.: Textmage: The Automated Bangla Caption Generator Based on Deep Learning. in: 2020 International Conference on Decision aid Sciences and Application (DASA), pp. 822–826. IEEE (2020)

[17]. Teera Siriteerakul , Kunlabut Suriyakanon ,Sofia Sarideh , (2018 ) " Automatic Restaurant Image Tagging " , International Journal of Electrical, Electronics and Data Communication (IJEEDC) , pp. 1-4, volume-6,issue-4

[18]. Marielet Guillermo, Robert Kerwin Billones, Argel Bandala, Ryan Rhay Vicerra, Edwin Sybingco, Elmer P. Dadios, Alexis Fillone, "Implementation of Automated Annotation through Mask RCNN Object Detection Model in Cvat using AWS ec2 instance", Region 10 Conference (TENCON) 2020 IEEE, pp. 708-713, 2020.