

# Smart Learning in Document Categorization using Dynamic Learning

S. M. Prabin, N. Selvaganesh, N. Rajesh Pandian, T. Selva Kumar

*Abstract—Clustering is the process of making data groups using similar data items, used for data mining to extract data from available large datasets. A large volume of text documents consisting of personal information is being generated in form of digital libraries and repositories in internet daily. It is conceivable to get to great quality instructive substance and strategies in an increasingly helpful manner. In spite of the fact that a ton of keen instruments have been connected for instructive application, there are just restricted looks into that show the instructive viability of shrewd devices through test contemplations, Clustering organizes large quantity of unordered text documents into small number of meaningful and coherent clusters. A clustering method based on K-Means algorithm is proposed in this paper. K-Means is a unsupervised algorithm based on randomly selected initial centroids used to cluster a highly unstructured and unlabeled document collection. The system will be evaluated using precision as a measure.*

**Keywords:** K-Means, DTM, IDF, Vector Space Model, Document Frequency-Based Selection

## I. INTRODUCTION

Clustering is most important process in data mining for the extraction of information from the larger datasets. It helps in converting an unstructured data into structured data for reducing the searching complexity. Nowadays huge amount of data is available in internet where the problem is managing the information and reducing the time for quick searching. Clusters are formed by group of same categorized information; It collects the similar information from the dissimilar information. Clusters mean collection of same behavioral information. How to make a good clustering? Clusters with information dissimilar from each other should be formed. So the data independency is most important in clustering. Clustering Algorithms helps to make a Unique clusters. Now we have to analyze how the clustering techniques are very important in day to day life. The most wanted sectors are:

**Marketing:** clustering is important in marketing for finding the product selling point in every area. Clustering can make the buyers group based on the product they are purchased already. So the product based clusters are formed from the

**Revised Version Manuscript Received on 20 September, 2019.**

**S.M.Prabin**, Assistant professor, Department of Computer Science and Engineering, PSNA College of Engineering and Technology, Dindigul Tamilnadu, India.

**N.Selvaganesh**, Assistant professor, Department of Computer Science and Engineering, PSNA College of Engineering and Technology, Dindigul, Tamilnadu, India.

**N.Rajesh Pandian**, Assistant professor, Department of Computer Science and Engineering, PSNA College of Engineering and Technology, Dindigul, Tamilnadu, India.

**T.Selvakumar**, Assistant professor, Department of Computer Science and Engineering, PSNA College of Engineering and Technology, Dindigul, Tamilnadu, India.

larger buyers database file.

**Libraries:** here clustering algorithms makes the book clusters based on the nature of each book. The similar Technique Books comes under the single cluster for easy finding of books from the large available books.

**Insurance:** here clustering technique make the group of objects based on the policy holded by the each customers. The unique policy holders are considered as single clusters. It makes the simple customer handling database files. here the clusters are easily identify the fraud users who claims the huger cost.

**Biology:** data classification can be done in two categories. One category is animals and the other category is plants.

**City categorization:** here the clustering technique can form the house clusters. It helps to identify the groups of houses based on their house type, value and geographical locations.

**Earthquake analytics:** clustering technique can observe earthquake scalability to identify dangerous zones for finding a recoverable solution only to that affected peoples.

**WWW:** here clustering can be done the weblog data to discover groups of similar access patterns and making the data classification based on the natural behavior of data's.

### 1.1 Text Classification:

It is very difficult to people to handle the huge amount of text documents from the various resources. Database management systems can access the data's from the different data bases or different data store. But this was only a small part of what could be gained from the internet source. But in real world we want to analyze the huge data's by various techniques, its helps to get the better knowledge about the data explicitly stored to derive and it gives knowledge about the topic. Here we need the sufficient knowledge about data mining or knowledge comes into existence. In regular day to day activity, every people can update their data's via any social medium. So it makes larger amount of data's to be available in Internet itself. It makes growth of information and also quickly increases the number of text and hypertext document in organizational intranets. But already the people having a sufficient knowledge of organization that becomes more and more success in today's information society. It simplifies the managing and analysis of data becomes very important. Nowadays information filtering and information retrieving from the large available information have a great attention to the both domestic and International consumers.<sup>1-2</sup>



Document clustering aims to group the same category documents into a single clusters, it is one of the most important tasks in data mining and artificial intelligence that are having a much importance in recent years. The main aspects is to cluster forming with a high accuracy as needed.

In the time of doing the clustering analytics, we first splits the set of data into unique groups based on data similarity between the data. The various types of algorithms are used for clustering the data and to improve the quality to a great extent.<sup>3-4</sup>

### 1.2 Overview

#### 1.2.1 Text Clustering Techniques:

The content bunching (clustering) is a method for gathering an arrangement of comparative information into content subsets (groups) with the end goal that comparative information are characterized like a solitary bunch. Bunches are framed by information closeness. However, every bunch are completely vary from the other bunch Groups. The primary objective of bunching is to break down the disparity between every datum specifically message records. Every bunch or gathering share basic conduct. Gathering framing between the comparable data is the principle part to take note of that the rundown of points is regularly not known preceding the grouping procedure. Likewise the grouping can makes report classification for simple treatment of data. Effective content bunching can give the better data benefits via seeking and dealing with the reports into profitable (justifiable) group orders and gives a superior significant supplement to our typical content looking in web crawlers, when catchphrase based hunt returns an excessive number of archives.

#### 1.2.2 Text Pre-Processing:

Content preprocessing strategy is critical in content mining. Here we need to continue an alternate littler process like tokenizing, evacuating stop words, stemming of terms, limit the term recurrence, distinguishing the weighting terms in report vectors. Tokenization can play out the gathering of information dependent on each single word. Tokenization process can maintains a strategic distance from the accentuation ,full stops, unique characters and the information that are as of now exhibits the essential word reference records and stopwords (is ,was, this ,that) are over and over existed in content archive first we have to evacuate it then no one but we can proficiently discover the term recurrence for shaping the bunches. This procedure is commonly nearness in that season of relational word process. The rundown of precisions process are regularly profit in grouping development, an institutionalize postings of stop words is constantly utilized in numerous looks into articles. The Porter stemmer whenever used to discover the term recurrence (rehashed number of term)in content records. At that point we have to build a vector display for every content archive.<sup>5-6</sup>

#### 1.2.3 Vector Space Model:

This is the frequent model for identifying document weighting terms. Vector Space Model is a process that encodes with “bag-of-words” representation, It explicitly gives the sequential ordering of data is not explicitly

gathered. In Vector Space Model, a document  $d_i$  is represented by a set of terms  $(t_1, t_2, \dots, t_n)$  wherein each  $t_j$  is a word that appears in the text document  $d_i$ , and  $n$  denotes the total number of various words in the index used to identify the meaning of the text document. Word  $t_i$  has a corresponding weight  $w_j$  calculated as a combination of the statistics term frequency  $TF_{t_i, d_j}$  and inverse document frequency  $IDF_{t_i}$ .  $W_t(t_i, d_j) = TF_{t_i, d_j} \times \log(N/DF(t_i))$ , Where  $TF_{t_i, d_j}$  is the frequency of term  $t_i$  in document  $d_j$ ,  $N$  is the total number of documents in the corpus,  $DF(t_i)$  is the number of documents in which term  $t_j$  occurs. Therefore, document  $d_i$  can be represented as a specific  $n$  dimensional vector  $d_i$  as  $d_i = (w_1, w_2, w_3, \dots, w_n)$ . The function encodes the intuitions that: (i) the more often a word occurs in a document, the more it is representative of the content of the text; (ii) The more text the word occurs in, the less distinction it is. In classification, the Inverse Document frequency is a good index of the usefulness of a word. The text document is also subjected to TF and IDF weighting which is used for the training documents. The operation of k-means algorithm is explained as follows:

Algorithm : k-means algorithm

1. Select K points as initial centroids.
2. Repeat(RPT)
3. Form k clusters by assigning each point to its closest centroid.
4. Recompute the centroid of each cluster.
5. until centroids do not change.

#### 1.3 Document Classification:

While cluster is naturally relate unattended learning strategy, it are frequently utilized to help the standard of the winds up in its regulated variation. Most importantly, word-bunches and co-preparing philosophy are frequently used with the end goal to help the arrangement precision of administered applications with the work of cluster systems. we tend to take note of that few classes of recipes like the K-implies calculation, or stratified calculations square measure general methodologies, which may be reached out to any very learning, together with content information. A content archive are regularly depicted either inside the style of double learning, after we utilize the nearness or nonattendance of a word inside the report in order to make a parallel vector. In such cases, it's capability to straightforwardly utilize a scope of downright information bunch calculations on the double representation. An a considerable measure of expanded representation would accompany refined weight procedures bolstered the frequencies of the individual words inside the report still as frequencies of words in an entire arrangement (e.g., TF-IDF weighting).<sup>7-8</sup>

#### 1.4 Document Frequency-Based Selection:

The most straightforward feasible strategy for highlight decision in report bunch is that of the usage of record recurrence to strain insignificant alternatives. Where as the use of reverse archive frequencies diminishes the significance of such words, this may not the only one be extra

proportional back the clamor impacts of horrendously visit words. In elective words, words that zone unit excessively visit inside the corpus are regularly expelled because of they are by and large basic words like "an", "a", "the", or "of" that aren't discriminative from a cluster point of view. Such words additionally are raised as stop words. a scope of techniques region unit generally reachable for stop-word evacuation. By and large normally possible stop word arrangements of concerning three hundred to four hundred words zone unit utilized for the recovery technique. Furthermore, words that happen once in a while might be off from the social occasion. this is regularly because of such words don't add something to the likeness calculations that region unit utilized in most bunch procedures. Now and again, such words is likewise incorrect spellings or typographic mistakes in records. crying content accumulations that territory unit got from the net, online journals or interpersonal organizations region unit extra without a doubt to contain such terms. we tend to take note of that a few lines of examination plot report recurrence fundamentally constructed decision entirely with respect to the preface of horrendously rare terms, because of these terms contribute the littlest add up to the similitude counts. Be that as it may, it should be focused on that appallingly visit words should even be evacuated, especially in the event that they're not discriminative between bunches. Note that the TF-IDF weight philosophy may normally strain very basic words in a "delicate" approach. Unmistakably, the quality arrangement of stop words offers a real arrangement of words to prune. Withal, we'd kind of a methodology of measuring the significance of a term on to the cluster technique that is essential for extra forceful pruning.

### 1.5 Objectives:

The goal of this project is to cluster the unlabeled data or the text document collection that is highly unstructured and, it further improves the quality of clustering solutions and accuracy.

## II. LITERATURE SURVEY

### I. An Efficient Approach for Text Clustering Based on Frequent Itemsets

As of late, the gigantic amount of issue data open in electronic sort is developing at stunning rate. This expanding scope of issue data has diode to the errand of mining accommodating or entrancing incessant thing sets (words/terms) from frightfully gigantic content databases and still it has all the earmarks of being very troublesome. the use of such regular thing sets for content agglomeration has gotten a decent arrangement of consideration in examination network since the strip-mined successive thing sets downsize the spatiality of the reports radically. inside the anticipated investigation, we have conceived Associate in Nursing prudent methodology for content agglomeration upheld the regular thing sets. A praised system, known as Apriori algorithmic program is utilized for mining the successive thing sets. The strip-mined continuous thing sets region unit at that point utilized for getting the parcel, wherever the archives territory unit at first bunched while not covering. Additionally, the resultant bunches region unit successfully gotten by gathering the archives inside the parcel by implies

that of determined watchwords.

### Perception:

Because of the exponential increment inside the volume of content report accumulations and in this manner the might want for investigating content records, numerous procedures are produced for mining the incessant relationship from content archives. Inside the content mining environment, content bundle connotes one among the premier successful ways to deal with bunch reports in partner degree unsupervised way. Amid this paper, we have a tendency to build up a decent methodology for content pack as per the continuous thing sets that has critical spatiality decrease. We have a tendency to get an accumulation of non-covering allotments exploitation this successive thing sets and along these lines the resultant bunch is produced inside the segment for the record accumulations.<sup>7-8</sup>

### II Text agglomeration with Extended User Feedback

Content agglomeration is most normally regarded as a thoroughly machine-controlled assignment while not client criticism. Nonetheless, a spread of scientists has investigated blended activity agglomeration ways which allow a client to move with and exhort the agglomeration run the show. This blended activity approach expressly is exceptionally captivating for content agglomeration errands wherever the client is making an endeavor to set up a corpus of reports into bunches for a couple of specific purposes. This paper acquaints a supplanting approach with blended activity agglomeration that handles numerous regular assortments of client criticism. we tend to beginning present a substitution probabilistic generative model for content agglomeration (the Spe-Clustering model) and demonstrate that it beats the unremarkably utilized blend of multinomial's agglomeration display, even once utilized in completely self-ruling mode with no client input. we tend to then depict the best approach to join four unmistakable assortments of client input into the agglomeration run, and supply exploratory confirmation demonstrating generous improvements in content agglomeration once this client criticism is consolidated.

### Perception:

Our trial results demonstrate our unattended Spe-Clustering principle beats the unremarkably utilized multinomial innocent mathematician agglomeration manage for every one of the content data sets we tend to contemplated. More finished, when given client criticism, the Spe-Clustering model additions essential change in an extremely close to home email dataset and inside the newsgroup dataset once the agglomeration results is shrieking anyway huge.

The proposed methodology joins the upside of the machine's procedure capacity to explore mammoth data sets, with the advantages of a human's comprehension of classes of intrigue. The outcomes demonstrate that participation among PCs and people might be a promising bearing for future work. There square measure a few future difficulties, similar to exploitation dynamic learning standards to



streamline the rundown of a group, and building extra unobtrusive models to allow extra common assortments of client input.<sup>8-9</sup>

### III. A Study of graded agglomeration rule

Bunching is that the strategy for gathering the data into classifications or groups, so questions inside a group have high similitude contrasted with 1 another anyway these articles square measure unpleasantly not at all like the items that square measure in various groups. Agglomeration ways square measure principally isolated into 2 gatherings: evaluated and apportioning ways. Evaluated agglomeration blend data objects into bunches, those groups into bigger bunches, at that point forward, making a chain of command of bunches. In parceling agglomeration ways various segments square measure made so assessments of those segments square measure performed by some basis. This paper presents expand dialog on some enhanced reviewed agglomeration calculations. Moreover to the current, creator have given a few criteria on that one may likewise check the best among these specified calculations.

#### Perception:

Various leveled agglomeration might be a technique of group investigation that tries to make a progressive system of bunches. The standard of an unadulterated reviewed agglomeration procedure experiences its powerlessness to perform change, once a consolidation or choice has been dead. This consolidation or choice, if not well picked at some progression, could result in some-what gauge groups. One promising bearing for rising the agglomeration nature of reviewed routes is to incorporate evaluated agglomeration with various procedures for different area agglomeration.

### III. EXISTING SYSTEM

#### Agglomerative Hierarchical Clustering Algorithms:

In the bunching writing for records of various types that incorporates multidimensional numerical information, content information and unmitigated information, Hierarchical grouping information have been widely examined. To help an assortment of seeking techniques, agglomerative various leveled bunching is especially utilized, in light of the fact that it makes a tree-like chain of command, which can be utilized for the inquiry procedure. In light of the likeness with each other, it progressively consolidates reports into bunches and it is the general idea of agglomerative various leveled grouping. Between the gatherings of reports dependent on the best match insightful likeness, Hierarchical bunching calculations progressively consolidate gatherings.<sup>8-10</sup>

The fundamental contrast between the classes of techniques is the means by which the calculation of combine savvy closeness is done between the archives of various gatherings. For instance, from the sets of gatherings, can figure the similitude between the sets of gatherings as best-case closeness, normal case comparability and most pessimistic scenario likeness between archives.

Theoretically, a bunch progressive system is made because of the way toward agglomerating archives into bunches of larger amounts, in which singular records are considered as leaf hubs, and the consolidated gatherings of

groups are considered as inner hubs. On the off chance that converging of two gatherings occur, there will exist another hub in the tree which relates to the bigger blended gathering.

The following are the steps for hierarchical clustering:

Step1: Start by assigning each item to a cluster, so that if there are N items, then there will be N clusters, each containing just one item. Let the distances (similarities) between the clusters the same as the distances (similarities) between the items they contain.

Step2: Find the closest (most similar) pair of clusters and merge them into a single cluster, so that, can have one cluster less.

Step3: Compute distances (similarities) between the new cluster and each of the old clusters.

Step4: Repeat steps 2 and 3 until all items are clustered into K number of clusters The two categorizations of Hierarchical clustering are Agglomerative approach and Divisive approach.

#### 3.1 Disadvantage:

- Computational time is quite slow.
- Identifying the correct number of clusters in dendrogram will be difficult, sometimes.
- Handling different size of clusters is difficult in hierarchical clustering algorithm.
- Efficiency and Accuracy is quite difficult to attain high level in hierarchical clustering.

### IV. PROPOSED SYSTEM& RESULTS

#### 4.1 K-Means Algorithm:

In particular, for document clustering there are two main approaches and they are K-Means and agglomerative hierarchical clustering. From a huge collection of document which is in unstructured format, retrieving important information from it is a very difficult and time consuming task. Often, Hierarchical clustering is the better quality clustering approach, but because of its increased time complexity it is limited. K-Means thereby reduces the time complexity and to classify a given data set through a certain number of clusters, it follows a simple and easy way.

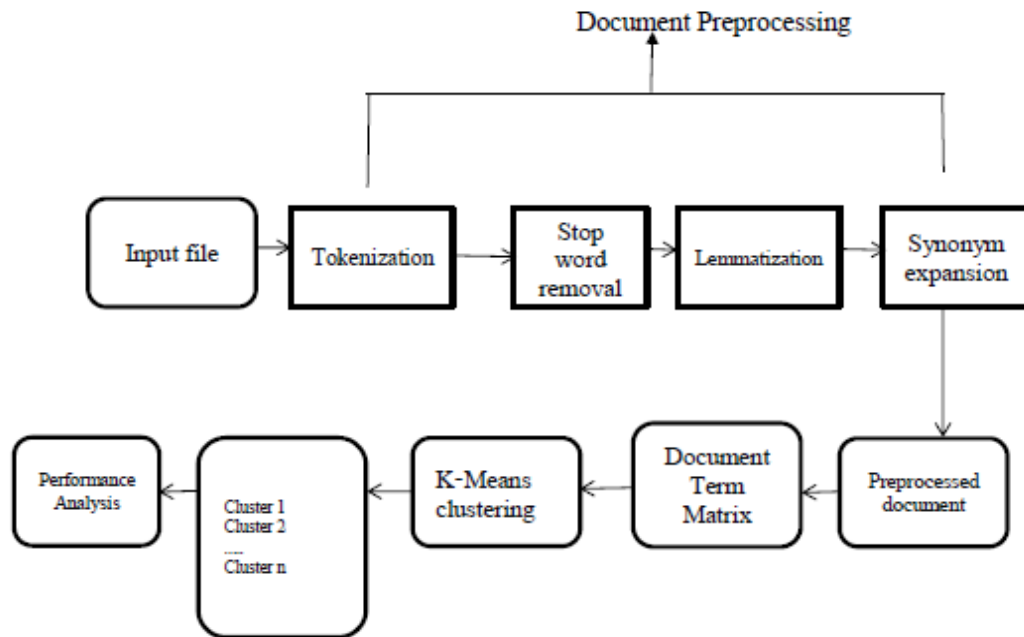
### V. SYSTEM ARCHITECTURE

The input file is first fed to the document preprocessing techniques, which includes tokenization, lemmatization, stopword removal and synonym expansion. After this process, can get preprocessed document, from that preprocessed document, document term matrix can be produced, in which apply K-Means algorithm for creating clusters, and finally performance evaluation is done by using precision recall method.



To “get the best of both worlds” agglomerative hierarchical approaches and K-Means can be combined. However, rather than hierarchical approaches, K-Means is a

good or better approach. Tighter clusters are produced in K-Means rather than hierarchical clustering.



### 5.1 Module Description

#### Input File:

Text File is given to document preprocessing. It is used for condensing a huge number of text documents into well-organized form for getting the desired results in less time.

#### Document Pre-Processing:

It consists of steps that take input as a plain text document and gives output as set of tokens. The document pre-processing process includes the following steps:

#### Tokenization

A sequence of strings is breaking up into pieces such as words, phrases and other elements called tokens and this process is called tokenization. Some characters like punctuation marks are discarded in the process of tokenization.

#### Stop Word Removal

Sometimes extremely common words which would appear to be of little value in helping select documents matching a user's need are excluded from the vocabulary entirely. These words are stop words and the process is called stop word removal. Stop words such as a, an, the, and prepositions are discarded.

#### Lemmatization

After the text is tokenized, to analyze as a single term, each inflected term is reduced to its base form. Lemmatization goes through the morphological analysis of the word in the text, so it is a preferred method to stemming.

#### Synonym Expansion

Searching each token in the dictionary and transforming each word to the base words is the process of synonym expansion. The dictionary consists of a list of words and all of their synonyms.

#### Pre-Processed Document:

The document which comes after applying all the preprocessing techniques is called pre-processed document.

#### Document-Term Matrix:

A document-term matrix (DTM) or term-document matrix that describes the frequency of terms that occur in a collection of document.

#### Clustering Using K-Means Algorithm:

After the construction of the Document Term Matrix, the process of clustering is carried out. The K-means clustering algorithm is used to meet the purpose of this project.

The basic algorithm of K-means used for the project is as following:

K-means Algorithm where each cluster's center is represented by the mean value of the objects in the cluster.

#### Input:

k: the number of clusters,

#### Output:

A set of k clusters.

### 5.2 Performance Analysis:

Performance Analysis is done to determine the accuracy of the clusters formed due to the algorithm (K-Means clustering). Precision Recall method is used here.

## VI. CONCLUSION AND FUTUREWORK

### 6.1 Conclusion:

The project study concerned the nice deal of labor on numerous areas of knowledge retrieval and text mining and



targeted on the assorted strategies for document pre-processing and document agglomeration. Text mining and agglomeration techniques are extremely powerful. The greater part of the examines are related with m-learning (portable adapting); as of late, the investigates on keen learning are expanding, the greater part of the investigates are done for advanced education, and fundamental recipients of each exploration are the two instructors and understudies. Furthermore, the principle apparatus which is connected to instructive condition of each paper is a cell phone, and PC is used frequently likewise for other instructive situations. The instructive use of shrewd devices, the vast majority of the creators have positive sentiments. As it were, they imagine that the presentation of brilliant instruments gives a beneficial outcome on the instructive condition. This paper was utterly support these techniques. The system was created for locating the similarities among the text document. For making ready the corpus of a pre-processed document, numerous techniques are applied. Lastly, the k-means agglomeration algorithmic program was used for making the similar clusters of the text document. The similar text documents were classified into one cluster. The \$64000 world application of the project study would facilitate folks to seek out the similar text document on completely different document portal from one platform. This may not be doable while not the utilization of text mining and agglomeration techniques. In general, it's not possible to manually rummage around for similar text in every of the portals and so compare every of them to seek out similarities between them.

### 6.2 Future Work:

The projected methodology but produces higher cluster solutions and is computationally quicker, however the accuracy provided here is kind of consistent not attain high accuracy and finding shortest path to cluster or create cluster is kind of tough. K-Means bunch rule usually doesn't work well for prime dimensions; thus potency must be improved. In Future, to beat the issues of projected system, the ant colony rule is employed. Once an ant finds a brief path from the colony to a food supply, alternative ant's square measure additional probably to follow that path, eventually ends up in all the ants following one path. The thought of the ant colony rule represents the matter to resolve, thereby applying this rule can increase the accuracy and potency usually works well for prime dimensions.

## REFERENCES

- 1 Vinod S.Badgujar, Asha H.Pawar, "Search Engine Using Clustering and Text Mining" in International Research Journal of Engineering and Technology (IRJET) Nov-2015
- 2 S.C.Punitha, P.Ranjith, JebaThangaiah and M.Punithavalli, "Performance Analysis of Clustering using Partitioning and Hierarchical Clustering Techniques" in Research Scholar, Department of Computer Science and Engineering, Karunya University, Coimbatore, India., 2014
- 3 K.Premalatha and A.M.Natarajan, "A Literature review on Document Clustering" in Information Technology Journal, 2010.
- 4 NoumanAzam, JingTao Yao, "Comparison of term frequency and document frequency based feature selection metrics in text categorization" in Department of Computer Science, University of Regina, Canada ,2012
- 5 J.Sathyapriya, S.Priyadharshini, "Clustering Techniques in Data Mining for Text Documents" in International Journal of Computer Science and Information Technologies, 2012
- 6 K.Sathiyakumari, G.Manimekalai, "A Survey on Various Approaches in Document Clustering" in IJCTA, sept-oct 2011
- 7 M.Thangamani and P.Thangaraj, "Integrated Clustering and Feature Selection Scheme for Text Documents" in Journal of Computer Science, 2010
- 8 SepidehSeifzadeh, Ahmad K.Farahath, Mohamed S.Kamel, "Short-Text Clustering using Statistical Semantics" in International World Wide Web Conference Committee, 2015
- 9 Dr.S.Vijayarani, Ms.P.Jothi, "An Efficient Clustering Algorithm for Outlier Detection in Data Streams" in International Journal of Advanced Research in Computer and Communication Engineering , September 2013
- 10 AkshayKrishnamurthy, SivaramanBalakrishnan, Aarti Singh, "Efficient Active Algorithms for Hierarchical Clustering" in International Conference on Machine Learning, Edinburgh, Scotland, UK, 2012