



**INTERNATIONAL JOURNAL OF ENGINEERING SCIENCES & RESEARCH
TECHNOLOGY**

TEXT CLUSTERING TECHNIQUES: A SURVEY

Mohit *

* Department of Computer Science, The NorthCap University, Gurugram, India

DOI: 10.5281/zenodo.573535

ABSTRACT

The advancements in the fields of mobile computing, grid computing, cloud computing, Internet of Things and primarily due to the availability of internet in the hand-held devices were the vital key factors in the growth of large amounts of data. The main challenge is to organize this big data in a structured manner that helps to derive new insights, predictive analysis and to find trends, patterns and their correlations. One of the solutions is to cluster the text, a significant technique of data mining. This paper investigates various techniques experimented in text clustering. It also describes the process of text clustering along with various similarity measures.

KEYWORDS: Data Mining, k-means, k-medoids, Text Clustering Techniques, Document Clustering.

INTRODUCTION

1.1 Data Mining

Data mining [1] is developed as an area worried with separation of valuable information from the data. Data mining approaches are connected to comprehend an extensive variety of genuine issues and is the computational procedure of searching patterns in extensive data sets. The general goal of the data mining is to concentrate data from a data set and change it into a reasonable structure for further analysis.

Clustering is the assignment of separating the data objects into various clusters with the end goal that data objects in similar clusters are more like other data objects in a similar cluster than those in different clusters. In basic words, the point is to isolate groups with comparable attributes and appoint them into groups.

We can understand by taking an example. Suppose, you are the head of a grocery store and wants to examine the preferences of your customers to tune up your business. Is it possible to see the details of every customer and formulate a business plan for all of them separately? Obviously not. For this, you can do clustering for your customers, like to cluster them into 5 groups on the basis of their purchasing routine. This is what we can call it the process of Clustering.

1.2 Text Mining

The Text Mining can be pointed out as data mining innovation which determines profitable and significant data from unstructured text information. The text mining innovation permits end users to extricate important data from a boundless measure of data and to distinguish the associations with the other data. It likewise includes text categorization, data recovery, and so on. High-limit linguistic assets and complex measurable pattern learning techniques are utilized to make the system do a top to bottom investigation of the data written in human expression and to find the concealed data from the given data. As information turn out to be huge information and social networks like SNS and online journals are extended, Text Mining is generally being utilized for publicizing, promoting, law case investigation, data recovery, etc.

Text mining is a prospering new field that endeavors to gather significant data from the natural language text. It might be approximately portrayed as the way toward dissecting text to concentrate data that is valuable for specific purposes. Contrasted and the sort of information put away in databases, text is unstructured, and hard to manage algorithmically. In any case, in present day culture, text is the most widely recognized vehicle for the formal trade of data. The area of text mining ordinarily manages texts whose capacity is the correspondence of authentic data or conclusions, and the inspiration for attempting to similar data from such text consequently is convincing, regardless of the possibility that achievement is just fractional. Similarly as text mining can be roughly depicted as searching for patterns in text, text mining is about searching for patterns in text. Be that as it may, the trivial correlation between the two conceals actual divergence. Data mining can be all the more completely described as the extraction of verifiable, already unidentified, and probably helpful information from

details. The information is certain in the aid data: it is concealed, obscure, and could barely be obtained without plan of action to programmed methods of data mining. With text mining, in spite of the data to be removed is obviously and expressly expressed in the text.

1.3 Text Clustering

Clustering [8] is an unsupervised technique in data mining where names of the articles are unknown. The idea of what constitutes a valuable cluster relies on upon the application and there are numerous techniques for discovering clusters subject to different criteria, both specially appointed and orderly. Clustering can be connected to various types of data including text. By managing the text data, articles can be files, sections, or words. Text Clustering [5] is defined as the way of combining identical kind of text documents. The issue can be defined as: given the collection of documents it is needed to split them into different clusters, with the end goal that documents in the similar cluster are more like to each other than to documents in different clusters. Clustering is especially helpful for sorting out documents to enhance recovery and browsing. Conventional Clustering algorithms can be reached out to manage textual information. However, there are many difficulties in grouping textual information. The text is typically expressed in high magnitude area even when it is in reality little. In addition, relationship between words showing up in text should be treated in clustering assignment. The variety in file sizes is another test which influences the representation. Thusly, the standardization of text representation is needed.

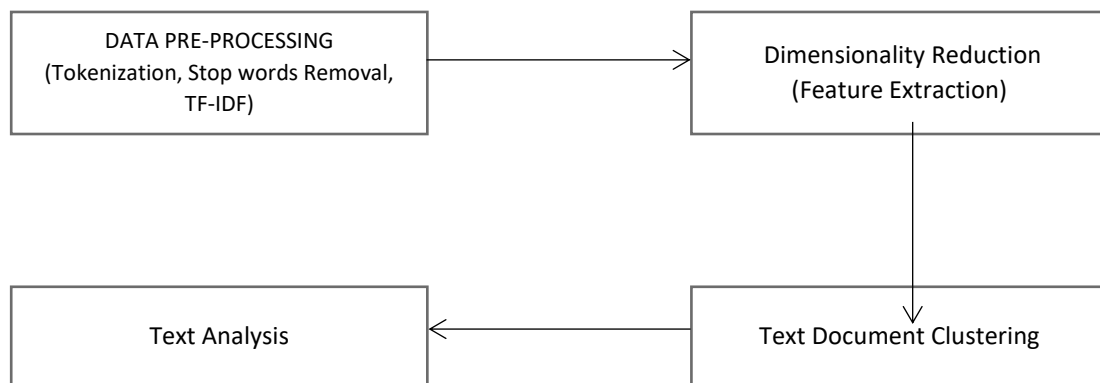


Figure 1: Basic steps of Text Clustering

In Figure 1, the steps which are generally followed to perform the clustering are represented. At initial, the dataset is taken on which we want to apply clustering, the first step is of data-preprocessing in which data is preprocessed like all the data is converted to such a form which can be used for making the analysis easy and more accurate. Then, in the second step the dimensions of data which are not so much important for making the analysis are removed from the dataset. The next step is to make the clusters by using some clustering algorithm as per our need, we cannot say that any of the clustering algorithms is best, it may varies to the size and type of datasets.

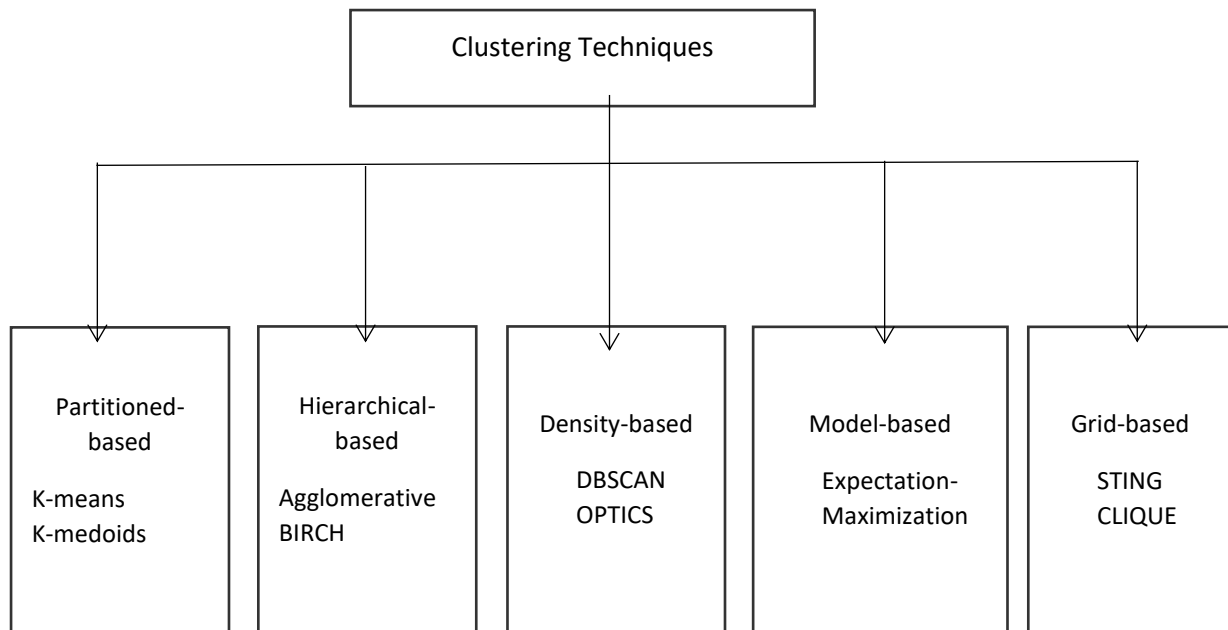


Figure 2: Categories of Different types of Clustering Techniques

Figure 2 represents the category of various clustering techniques but the algorithms which can be applied on the textual data are: K-means and K-medoids in Partitioned-based, Agglomerative and BIRCH in Hierarchical-based, DBSCAN in Density-based, Expectation Maximization in Model-based, STING and CLIQUE in Grid-based. Each algorithm has its own usage, we cannot say that any of the algorithm is best. It totally depends on the dataset and our requirements.

TEXT DOCUMENT CLUSTERING PROCEDURE

To cluster the text document, we need to perform a number of steps which are as follows:

2.1 Data Collection: In this step, the dataset is collected from any source to which we want to apply the clustering. This paper is for text clustering, so we have to choose some textual data or textual documents.

2.2 Data preprocessing: The data to be analyzed is preprocessed so as to remove the unwanted data that is not so much useful for making the analysis. The preprocessing of data includes a number of steps which are discussed as below:

2.2.1 Tokenization: At the very first, all the data is converted into the form of tokens, the smallest unit of the document, all the whitespaces are removed from the data.

2.2.2 Stop words Removal: Stop words means to the words which are not so much relevant for the purpose of making the analysis.

2.2.3 Stemming: There are so many words which are originated from the same word having same meaning. The process of stemming converts the words into their root words so that they both are to be considered as the similar words for making the analysis easier.

2.3 Document Representation: The document is represented using some model, as we generally use the vector space model. The vector space model considers each term as a vector. Let us assume the document D_x as a vector of the document, can be represented as n terms, $D_x = (t_1, t_2, \dots, t_n)$. If there are m documents in the dataset containing n number of unique words, is then represented as $m \times n$ matrix, where every document is a vector of q dimensions.

2.4 Similarity Measures: There are various types of similarity measures which are used to find the similarity of two documents or words or vectors. Some of them which are used for clustering the text, are as below:

2.4.1 Cosine measure: The cosine measure is used to compute the cosine angle between the two vectors of the documents. The value of cosine measure lies in between 0 to 1, as increases with more similarity.



2.4.2 Jaccard measure: The Jaccard measure is also used to compute the similarity among the documents by obtaining the common terms and unique terms of both the documents respectively. Same as of the cosine measure, the value of jaccard measure lies in between 0 to 1.

2.4.3 Correlation Coefficient: The correlation coefficient calculates the relativity of the two vectors or the documents. The value it returns lies in between -1 to 1.

2.5 Clustering Techniques: After calculating the similarity, now it is the turn to apply the clustering technique as per our requirement of clusters and according to our dataset. As in the Figure 2, there is a category of different types of text clustering techniques.

LITERATURE REVIEW

Text Clustering is utilized to make the groups without the knowledge of the categories which are utilized in our data. In this field, many looks into have been inspected which are made by various authors which are as follows:

3.1 Partitioning-based Clustering

An algorithm was designed by Xinwu [7], by consolidating the superiorities of K-means and SOM to build the efficiency of the algorithm than the ordinary K-means algorithm which has less stability of the clusters as expressed by the author utilizing the comparison made, F measure is used for making comparison of accuracy of the proposed improved clustering approach with the ordinary k-means, the scope of dispersing of the F-values if there should be an occurrence of proposed approach than the conventional approach. The restriction of k-means algorithm of introductory cluster center was progressed.

Jadon et al. [9], proposed a new approach in which the author has taken distinctive number of keywords from different branches of engineering and isolated them into specific number of domains as the predefined value of k and then flags are taken to incorporate the keywords into various clusters until each of the clusters get wrapped up. The performance is tried by taking the purity, entropy and F-measure contrasted with k-means algorithm and found that the results of new approach demonstrate preferable results in all measures over the k-means algorithm.

Balabantaray et al. [11], compared two clustering algorithms: K-means and K-medoids, for clustering 100 documents of five different categories using the Weka tool and implemented all the steps which are necessary for processing the text documents. The author analyzed that the results of K-means clustering outperforms over the K-medoids clustering in terms of efficiency, by using two distance measures: Euclidean and Manhattan distance.

Arora et al. [16], compared two clustering techniques, k-means and k-medoids [4] on transaction10k dataset of KEEL containing 10000 transactions of the purchased things and analyzed that the cluster head determination and also the complexity of clusters overlying is far great if there should be an occurrence of k-medoids rather than k-means. Also the k-medoids shows superior results in all manners like in terms of execution time, insensitive to unpredictable data, and turn down the noise as in contrast with k-means, which decreases the measure of heterogeneity of data points. In case of k-medoids, by substituting the representative data points with the non-representative data points will upgrade the nature of clusters generated.

Mishra et al. [18] utilized k-means strategy for clustering the documents in view of subjects present in every one. The fundamental supposition was that a document may manage different points. The proposed algorithm, named inter-passage based algorithm for clustering, was connected to group documents portions in view of closeness. After sections were preprocessed, keywords were recognized for each portion utilizing TF-IDF and sentiment conflict scores. All the sections were then expressed for utilizing keywords and a fragment count was computed. At last, k-means was connected to entire portions. The subsequent groups indicated high intra-cluster closeness and low inter-cluster closeness.

A multi-objective approach of text clustering using k-means algorithm was designed by Abualigah et al. [19], which is known for optimization, is compared using seven datasets by testing it with traditional k-means algorithm and found superior results for the proposed one, by taking the F-measure method to check the accuracy of both the algorithms. Al-Anazi et al. [20], compared three partitioning clustering techniques: k-means, k-means fast and k-medoids on capstone projects documents using three similarity measures: cosine, jaccard and correlation coefficient and figure out that k-medoids and k-means algorithms using cosine measure achieve good outcomes because cosine measure doesn't depends on the document length. In partitioning algorithms, we have to predefine the value of k, number of clusters which is one of the restrictions of partitioning algorithm. Also, as amount of k goes up, quality of the cluster gets more accurate and efficient.

Bide et al. [17], proposed an improved version of k-means algorithm, in which divide and conquer strategy is utilized which partitions the dataset into various terms subsequent to applying the feature extraction on the isolated documents and then conquer every one of the documents. By implementing this thought, it makes the

algorithm more efficient and also boost up the performance of the algorithm. The author analyzed that the F1-score of the designed approach is high as of the traditional algorithm and the time consumed by the proposed approach is significantly lesser in contrast with existing one.

To group categorical items, Zhexue Huang et al. [21] designed 2 approaches: k-modes and k-prototypes. K-modes approach utilizes basic coordinating uniqueness measure for managing categorical articles, take over the means of groups with modes, and utilizes a frequency based technique to refresh modes in the clustering procedure to limit the clustering cost function. By using these augmentations, k-modes approach can cluster the categorical information in a manner like k-means approach. This k-prototypes approach concludes the meaning of consolidated difference measure, additionally incorporates the k-means and k-modes approach to take into consideration grouping objects expressed by blended numeric and categorical qualities.

3.2 Hierarchical-based Clustering

Steinbach et al. [3], contrasted k-means algorithm and hierarchical algorithm in which two variants of k-means are utilized i.e. simple k-means [2] and bisecting k-means, are tried with the dataset retrieved from TREC, REUTERS-21578 and WebACE. The author analyzed that the best performance comes when we apply the bisecting k-means algorithm with the datasets as it manages with the uniform size of clusters as opposed to the other algorithms. Three kinds of hierarchical clustering techniques have been contrasted i.e. intra-cluster similarity, centroid similarity and UPGMA algorithms and found the UPGMA as the best algorithm as of the others algorithm. Additionally, the F-measure and entropy of the both k-means and UPGMA with refinements are tested, bisecting k-means shows better results over the others, and simple k-means is superior to UPGMA with refinement.

Abbas et al. [6], compared various algorithms: Expectation Maximization, K-means, Self-Organizing Maps and Hierarchical clustering algorithm by utilizing two dataset taken from the UCI and KDnuggets and conclude that as the k value gets high, the SOM algorithm perform low however higher than the hierarchical algorithm in same case. At the time when the data is huge, then the k-means and EM algorithm works superior to the other algorithms and are conscious to unpredictable data.

3.3 Density-based Clustering

Raviya and Dhinoja [10], compared two clustering techniques i.e. K-means and DBSCAN which includes BD scan and SNN using weka tool by comparing with seven parameters which demonstrate that the K-means algorithm takes less time as of the DBSCAN algorithm. In DBSCAN, the area having high density expresses the cluster existence whereas the are having low density of points expresses the outliers or noise.

3.4 Model-based Clustering

As already discussed [6], various clustering algorithms are compared from which, the Expectation Maximization and K-means works better when there is huge data.

Umale and Nilav [15], gave the overview of K-means and Expectation maximization techniques for document clustering for forensic analysis, by which a forensic examiner can find the similar words in a cluster using both techniques. The Expectation Maximization technique is the variant of K-means clustering that belongs to model-based clustering, having two steps: expectation which assigns the data objects and maximization which finds the maximum parameters that maximizes the weight or likelihood of the probabilistic algorithm.

3.5 Grid-based Clustering

Mann and Kaur [22], compared two clustering techniques i.e. Grid Density Based Clustering which is hybrid of Grid based and Density based and second one is K-means, from which the author analyzed that the grid based density algorithm is better than the K-means algorithm as it works with outlier detection, no specification of initial clusters as of K-means which requires initial seed for clustering, grid density based can be used for spatial database and high dimensional database and no restriction on the type of data. Also, the time complexity of grid density base clustering is superior to other algorithms.

3.6 Other Text Clustering Algorithms

An evolutionary approach of document clustering was proposed by Akter et al. [12], which depends on Genetic Algorithm. The author has taken the dataset from REUTERS-21578 and isolate the dataset into a few groups and then implemented the genetic algorithm to each partition independently and then implemented the second step on genetic algorithm on the consequences of first step. Cosine measure is used to calculate the similarity of two documents or vectors. The final result demonstrates that the proposed approach indicates better outcomes as opposed to the existing one. Gupta et al. [13], compared two different clustering algorithms, one is partitioning based i.e. k-means clustering algorithm and other is hierarchical i.e. BIRCH clustering algorithm and analyzed that the BIRCH algorithm shows better cluster results for a expansive dataset however it is less effective in time and memory as it takes additional time than the k-means clustering.

A new algorithm was proposed by Kadhim et al. [14], in which TF-IDF and SVD are consolidated for the reduction of large measurements. The author had taken the BBC sport and BBC news datasets to make the

examinations. In the paper, documents are preprocessed by utilizing various steps and afterward concept of term weighting is used to weight the terms for computing the relevance of the term for the document, after which dimensionality reduction is performed to reduce the dimensions so as to make the performance higher for which SVD is used and for clustering, k-means clustering algorithm is used to generate the clusters. The outcomes of proposed algorithm demonstrate the superior performance of the system and boost the approval degree between the documents. It is seen that as the dimensions get diminished by using the designed algorithm, the accuracy of the clustering get increased.

Figure 3 represents the distribution of various text clustering techniques applied in the literature.

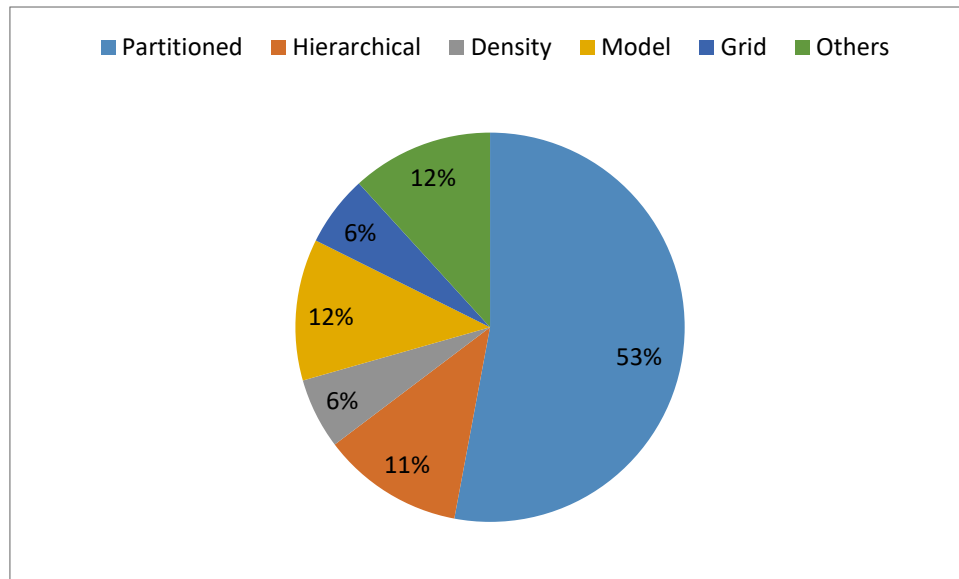


Figure 3: Distribution of various Text Clustering Techniques

From Figure 3, it is evident that the partitioned-based text clustering has been used by most of the researchers and Grid-based and Density-based is least used by the researchers.

Table 1 describes the summary of techniques used in Text Mining, by different authors and the year in which they made the research.

Table 1: Summary of Techniques used in Text Mining

Year	Author	Techniques	Critic
2000	Steinbach et al. [3]	K-means Bisecting K-means Hierarchical	The bisecting K-means works better among the others because it produces the clusters of uniform size however K-means produce clusters of widely different sizes
2008	Abbas et al. [6]	K-means Hierarchical Self Organizing Maps Expectation – Maximization	The author compared various algorithms: k-means, hierarchical, self-organizing maps and expectation maximization clustering algorithm and conclude that as the k value gets high, the SOM algorithm perform low however higher than the hierarchical algorithm in same case. At the time when the data is huge, then the k-means and EM algorithm works superior to the other algorithms and are conscious to unpredictable data.
2008	Xinwu [7]	K-means Self Organizing Model	The author proposed an approach by consolidating the advantages of K-means and SOM to increase the efficiency of algorithm

			than the ordinary K-means algorithm which has less stability of the clusters. The restriction of k-means algorithm of introductory cluster center was progressed
2013	Jadon et al. [9]	K-means	The author proposed a new approach in which different number of keywords from different branches of engineering are taken and separated them into specific number of domains as the predefined value of k by comparing it with the K-means algorithm. The results of new approach demonstrate better results in all measures over the k-means algorithm
2013	Raviya and Dhinoja [10]	K-means DBSCAN	The author compared two clustering techniques i.e. K-means and DBSCAN using weka tool by using seven parameters and found the results which demonstrate that the K-means algorithm takes less time as of the DBSCAN algorithm
2013	Balabantaray et al. [11]	K-means K-medoids	The author compared two clustering algorithms: K-means and K-medoids, for clustering 100 documents of five different categories using the Weka tool and analyzed that the results of K-means clustering outperforms over the K-medoids clustering in terms of efficiency.
2013	Akter et al. [12]	Genetic Algorithm	The author proposed an algorithm which depends on Genetic Algorithm. The data is partitioned and then implemented the genetic algorithm to each partition independently and then implemented the second phase on genetic algorithm on the consequences of first phase. The final result demonstrate that the proposed approach indicates better outcomes as opposed to the existing one
2014	Gupta et al. [13]	K-means BIRCH	The author compared two different clustering algorithms, one is partitioning based i.e. k-means clustering algorithm and other is hierarchical i.e. BIRCH clustering algorithm and analyzed that the BIRCH algorithm shows better cluster results for large dataset however it is less effective in time and memory as it takes additional time than the k-means clustering.
2015	Bide et al. [17]	K-means	The author compared two clustering techniques, k-means and k-medoids and analyzed that the cluster head determination and also the complexity of clusters overlying is far great if there should be an occurrence of k-medoids rather than k-means. Also the k-medoids shows superior results in all manners like in terms of execution time, insensitive to unpredictable data, and turn down the noise as in contrast with k-means, which decreases the measure of heterogeneity of data points.

2016	Al-Anazi et al. [20]	K-means K-means fast K-medoids	The author compared three partitioning clustering techniques: k-means, k-means fast and k-medoids on capstone projects documents using three similarity measures: cosine, jaccard and correlation coefficient and figure out that k-medoids and k-means using cosine measure achieves good outcomes. In partitioning algorithms, the restriction of the partitioning algorithm is predetermination of value of k.
------	----------------------	--------------------------------------	---

CONCLUSION

For clustering the text documents, we can use different algorithms as per our requirement of the output clusters. The normal clustering is different from the text document clustering because we have to convert the documents into a form which can be used for making the analysis. We discussed the process of text document clustering with the researches made for different types of clustering techniques, each technique has its own advantages and restrictions, like for huge dataset k-means performs better than the others, DBSCAN performs superior to other techniques in case of outlier detection.

REFERENCES

1. Ming-Syan Chen, Jiawei Han and Philip S. Yu, "Data Mining: An overview from a Database perspective", IEEE Transactions on Knowledge and Data Engineering, Vol. 8 No. 6, 1996.
2. John A. Hartigan, "Clustering Algorithms". New York, NY, USA: John Wiley & Sons, Inc.; 99th ed.; 1975.
3. Michael Steinbach, George Karypis and Vipin Kumar, "A Comparison of Document Clustering Techniques", KDD Workshop on Text Mining, 2000.
4. Charles Elkan, "Using the Triangle Inequality to Accelerate k-Means", Machine Learning, Proceedings of the Twentieth International Conference (ICML 2003) p. 147–153.
5. Magnus Rosell KTH CSC, "Introduction to Information Retrieval and Text Clustering", 2006.
6. Osama Abu Abbas, "Comparisons between Data Clustering Algorithms", The International Arab Journal of Information Technology, Vol. 5, No. 3, July 2008.
7. Li Xinwu, "Research on Text Clustering Algorithm Based on K_means and SOM", IITA Workshops 2008, IEEE.
8. Charu C. Aggarwal, ChengXiang Zhai, "A Survey of Text Clustering Algorithms", editors "Mining Text Data", Springer US; 2012, p. 77–128.
9. Chandan Jadon, Ajay Khunteta, "A New Approach of Document Clustering", International Journal of Advanced Research in Computer Science and Software Engineering, Volume 3, Issue 4, April 2013.
10. Kaushik H. Raviya ** Kunjan Dhinoja, An Empirical Comparison of K-Means and DBSCAN Clustering Algorithm, PARIPEX Indian Journal of Research, Vol2, Issue 4, April 2013.
11. Rakesh Chandra Balabantaray, Chandrali Sarma, Monica Jha, Document Clustering using K Means and K-Medoids, International Journal of Knowledge Based Computer System, Vol 1 Issue 1, June 2013
12. Ruksana Akter, Yoojin Chung, "An Evolutionary Approach for Document Clustering", International Conference on Electronic Engineering and Computer Science, Elsevier, 2013, pg. 370-375
13. Mamta Gupta, Anand Rajavat, "Comparison Of Algorithms For Document Clustering", International Conference on Computational Intelligence and Communication Networks, 2014.
14. Ammar Ismael Kadhim, Yu-N Cheah and Nurul Hashimah Ahamed, "Text Document Preprocessing and Dimension Reduction Techniques for Text Document Clustering", International Conference on Artificial Intelligence with Applications in Engineering and Technology 2014, IEEE.
15. Bhagyashree Umale and Nilav M., "Overview of K-means and Expectation Maximization Algorithm for Document Clustering", International Journal of Computer Applications, International Conference on Quality Up-gradation in Engineering, Science and Technology (ICQUEST2014).
16. Preeti Arora, Dr. Deepali, Shipra Varshney, "Analysis of K-Means and K-Medoids Algorithm for Big Data", Elsevier, ICISP2015.



17. Pramod Bide, Rajashree Shedge, "Improved Document Clustering using K-means Algorithm", IEEE International Conference on Electrical, Computer and Communication Technologies (ICECCT), 2015
18. Mishra, R., Saini, K., Bagri, S., "Text document clustering on the basis of inter passage approach by using K-means" In: 2015 International Conference on Computing, Communication Automation (ICCCA). 2015, p. 110–113.
19. Laith Mohammad Abualigah, Ahamad Tajudin Khader, Mohammed Azmi AI-Betar, "Multi-objectives-based text clustering technique using K-mean algorithm", 7th International Conference on Computer Science and Information Technology (CSIT), 2016.
20. Sumayia Al-Anazi, Hind AlMahmoud, Isra Al-Turaiki, "Finding similar documents using different clustering techniques", Symposium on Data Mining Applications (SDMA 2016), ELSEVIER.
21. Zhexue Huang, "Extensions to the k-Means Algorithm for Clustering Large Data Sets with Categorical Values", Data Mining and Knowledge Discovery, vol. 2, no. 3, pp. 283-304, 1998.
22. Amandeep Kaur Mann, Navneet Kaur, "Grid Density Based Clustering Algorithm", International Journal of Advanced Research in Computer Engineering & Technology (IJARCET) Volume 2, Issue 6, June 2013.

CITE AN ARTICLE:

M., Kumari, A. C., Sharma, M., & Sachin, R. (2017). TEXT CLUSTERING TECHNIQUES: A SURVEY. INTERNATIONAL JOURNAL OF ENGINEERING SCIENCES & RESEARCH TECHNOLOGY, 6(5), 248-256. doi:10.5281/zenodo.573535