

# Document Clustering Using K-Means Algorithm

<sup>#1</sup>Akshay Pandav, <sup>#2</sup>Ankush Bharekar, <sup>#3</sup>Chaitanya Shewale, <sup>#4</sup>Pradeep Bhutare, <sup>#5</sup>Prof. Shilpa Bhosale



<sup>1</sup>akshay.pandav@gmail.com  
<sup>2</sup>bharekarankush91@gmail.com  
<sup>3</sup>chaitanyashewale99@gmail.com  
<sup>4</sup>pradeepbhutare@gmail.com

<sup>#1234</sup>Student, Computer Engineering  
<sup>#5</sup>Prof, Computer Engineering

NBSSSOE, Pune, India - 411041.

## ABSTRACT

Data mining brings theories among several field including databases and optimization and data visualizations and has been applied to various real-life applications. In digital computer forensic analysis, hundreds of thousands of documents are saving into the digital formats these documents are usually in the unstructured textual data format and examining of that unstructured data is very difficult to examine by computer examiners. In this paper, we build a system for grouping a document clustering method for grouping whole data into cluster format. We reported an approach that applied to document clustering algorithms to forensic examination of system used in police investigations. In data mining, document Clustering is dependent on K-means and K-medoid algorithm for clustering the documents by utilizing centroid informational point. Here we provide security to our database server.

**Keywords:** Clustering, Text mining

## ARTICLE INFO

### Article History

Received: 24<sup>th</sup> May 2016

Received in revised form :  
24<sup>th</sup> May 2016

Accepted: 27<sup>th</sup> May 2016

**Published online :**

28<sup>th</sup> May 2016

## I. INTRODUCTION

The size of data is increasing with new day. This huge amount of data is impacting over computer forensics. To categorize the data manually required large time and it is hard work. Manual analysis of data is done by human and hence, there are chances of generating errors in analysis. So we had defined computerized analysis of data. For solving these problem by using K-means and K-medoid algorithms. The main aim of clustering is to track information and in the present days, to locate most relevant electronic resources. Clustering is a technique in which, one make cluster of same type of informational objects that are somehow similar in behaviour. Our data may be labelled or not-labelled. Labelled data tracking is easy as compared to unlabelled analysis of data. The concept behind these algorithms is that data objects within a valid clusters are very identical to each other than they are to objects belonging to a different cluster. Thus, once a data distribution has been introduced from data, the examiner might focus on rechecking representative document from the obtained set of clusters. Then, after of preliminary analysis, he may eventually decide to categorize other documents from each cluster. By doing so, one can tackle the tedious task of checking all the documents but, even if so required, it still could be done.



Fig 1. Introduction to Forensic Analysis

From the diagram it is obvious that forensic analysis perform collecting of evidence.

Forensic analysis puts same data into the first phase. It is the selective storage. It involves two steps viz.

1. Theoretical extraction.
2. Theoretical data analysis

Digital forensics were commonly used in cybercrime activity and private analysis of process. It has been similar with law, where proofs were collected to current hypothesis before the courts. As with other sector of forensics this is always part of larger analytic principles. In this type the

gathered proofs are used as a form of intelligence, used for other needs than court laws. As a result, intelligence collections are sometimes held to a non-discipline forensic standard. A general instance would be following not authorised network intrusion identification. A forensic analysis into the know-how of the attack is operated as a damage ceasing act.

Two of will compare within to get an opportunity to find out the subject. Such malicious attempts have impressively operated over mobile lines during the last 40 years, but in the present time are likely to use Internetwork.

## II. RELATED WORK

Document clustering is a strategy in which, the data that is matched is stored together. For improving the performance of finding and then getting back in database, the number of disk accesses is to be reduced. In clustering, due to the objects with matched characteristics are placed in similar class of objects, a single allowance to the disk can get back the whole class of objects. If the clustering occurs in some already assigned algorithmic space, we may make population into subsets with unique properties, and then lower the obstacle space by operation on only a single cluster head from each cluster. Clustering is timely technique of minimizing a large amount of information's to be formatted groups. For relating and calculations reducing complexities, these groups may include of "similar" information and objects. There are two ways to document clustering, simplified in data retracing; they are called as terms and item clustering. Terms clustering is a way in which groups show terms, and this grouping minimizes, unwanted information and getting most of assignment. There are lesser studies presenting the application of grouping or clustering algorithms in the Computer Forensics domain. Importantly, most of the studies describe the use of orthodox algorithms for clustering data. e.g., Expectation-Maximization (EM) for learning of Gaussian Mixture Models, K-means. These algorithms have similar characteristics and are frequently in practice. For example K-means, K-medoids, Single Link, Complete Link and Average Link, can be seen as particular cases of Expectation maximization. The theory on Computer Forensics only presents the use of algorithms that assure that the number of clusters is well defined and static by the user. Aim at relaxing this assumptions, which is generally not in order in real life utilization, a general approach in other sectors includes assumptions and the number of clusters from information. Significantly, it induces variations in information divisions and then corrects them with a comparative time index in demand to assume the improved price for the number of clusters. In this section, we discuss related work on document clustering and clustering algorithm.

### A.DOCUMENT CLUSTERING

Extraction and faster data access. Document clustering includes explanation of a extraction. Descriptors were group of words that defines the information within cluster. Document cluster is generally assumed to be centralized activity in web document clustering. Application are available both over and across the internet.

### B.PRE-PROCESSING STEPS

Cease words, sentences, commas before clustering algorithm. It explains deleting of pronouns, prepositions, articles and not related document. It enables/disables snowball on line surfs. Mining of text using pre-existing words satisfies the way of access. Finds out blueprints vector-space. In this type we own effectiveness, proficiency, efficiency, and clustering algorithm. Changes in vector gets a number of attributes that had been used namely, cosine-based, computer-based distance.

### C.CLUSTERING ALGORITHM

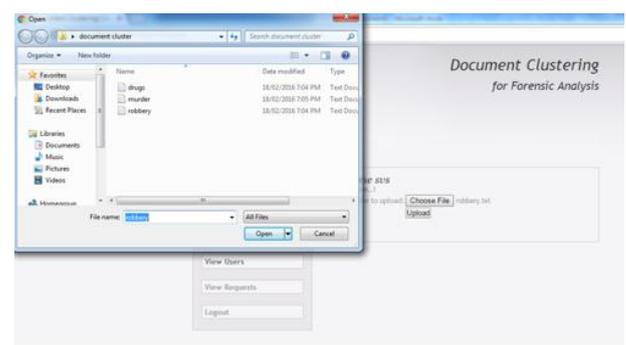
Data mining fields involves Cluster Ensemble Based Algorithm including centroids. This characteristic makes it significantly interesting for utilities in which,

- 1) Centroid cannot be computed, and
- 2) Paths between pairs of objects are available k-MEANS AND k-medias are sensitive to start assuming division algorithms. Every data object is reported by the diagrams subsequently selecting best outcomes. CSPA algorithm specially gets a clustering from a cluster group made by a set of variations in divisions. After forming cluster algorithms to the information giving objects a inter same number is computed. Each member of this matrix shows same pair-wise between objects. The similarity in information objects is not complex then fraction of the clustering solutions in which those two objects remains within the single cluster using the Template.

## III. IMPLEMENTATION OF MODULE

### Module 1:

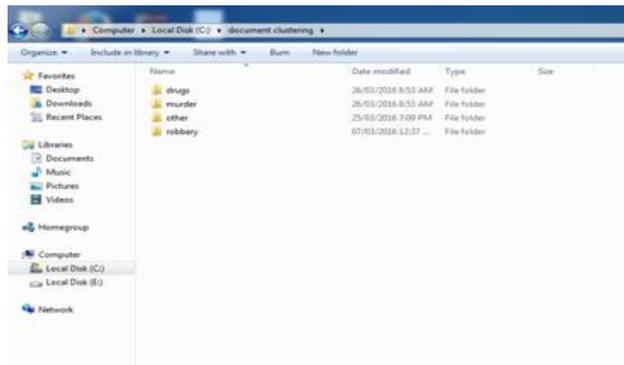
In the module 1 user upload the document. While uploading document partitioned into number of parts then it is forwarded to the server, after uploading, it is processed for making clusters.



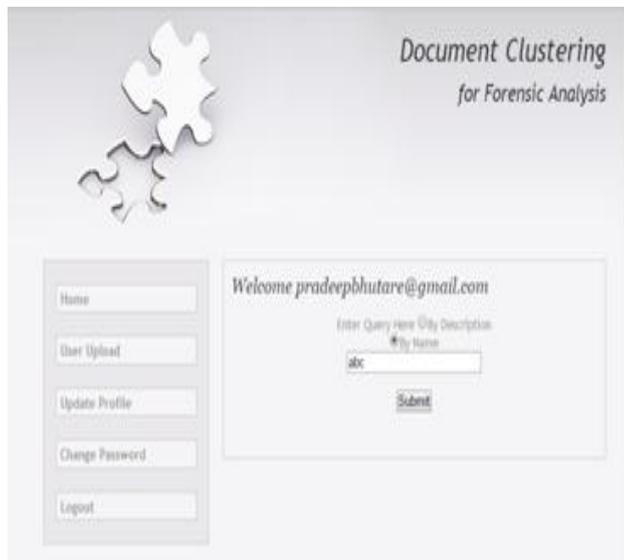
### Module 2:

In Module 2 uploaded file is examined and after the examination Pre-processing is done on this document. In Pre-processing we delete the stop words. These documents will be clustered based on the data in document. Before clustering files are analyzed and its keywords are stored on the repository which will be used for finding for user. We are attaching criminal name with the file as we are using this system for visa processing These are our clusters,  
Drugs  
Murder

Robbery  
Others



Module 3:  
This is our last module here we are finding user files based on the name or based keywords. For finding we are using keywords from repositories which are stored while examining, and also we can find by name of criminal so our system will be useful for visa processing.



**IV. RESULT**

**Test case 1:**

<b>MODULE NAME: CLUSTERING FILE</b>		
<b>ACTION</b>	<b>INPUT</b>	<b>EXPECTED OUTPUT</b>
Get Analyse File	Clustering	Clustering Of File Done
<b>RESULT : SUCCESS</b>		

**Test case 2**

<b>MODULE NAME: NAME ATTACHED</b>		
<b>ACTION</b>	<b>INPUT</b>	<b>EXPECTED OUTPUT</b>
Get Clustering File	Name Attach	Criminal Name attaché successfully
<b>RESULT : SUCCESS</b>		

**Test case 3**

<b>MODULE NAME: USER UPLOAD FILE</b>		
<b>ACTION</b>	<b>INPUT</b>	<b>EXPECTED OUTPUT</b>
Select File	Upload Action	File Upload Successfully
<b>RESULT : SUCCESS</b>		

**Test case 4**

<b>MODULE NAME: FILE ANALYZE</b>		
<b>ACTION</b>	<b>INPUT</b>	<b>EXPECTED OUTPUT</b>
Get Upload File	Analyze Start	File Analysis Done
<b>RESULT : SUCCESS</b>		

**Test case 5**

<b>MODULE USER UPLOAD FILE</b>		
<b>ACTION</b>	<b>INPUT</b>	<b>EXPECTED OUTPUT</b>
Select File	Upload Action	File Not Upload Successfully Due to Some Path Problem
<b>RESULT : SUCCESS</b>		

**V. CONCLUSION**

We have successfully implemented the system of Document clustering by utilization of k-means algorithm, which forms faster searching of unframed data as well as structured data. This system can also be applicable for the

passport identification purpose because we have made available crime verification approach, which helps in getting any human whether he is accused or not..

### REFERENCES

- 1].L. Liu Pu, and W. Han, "XWRAP: An XML-Enabled Wrapper Construction System for Web Information Sources," Proceedings of the 16th IEEE International Conference on Data Engineering (ICDE), San Diego, California, 2000, pp. 611-621.
- 2].Gibson D, Punera K, Tomkins A. The volume and evolution of web page templates. In: Proceedings of WWW'05. New York, NY, USA, 2005: 830-839.
- 3].NoDoSE (Northwestern Document Structure Extractor): An interactive tool for extracting data from semi-structured documents (plain text or HTML pages) [Adelberg, 1998].
- 4].Shine N. Das, Pramod K. Vijayaraghavan, Midhun Mathew, "Eliminating Noisy Information in Web Pages using featured DOM tree", International Journal of Applied Information Systems (IJAIS) – ISSN : 2249-0868.
- 5].S. Mythili, T. Vetrivelvi," Analytics of Noisy Data in Web Documents Using a Dom Tree", International Journal of Advanced Research in Computer Science and Software Engineering.
- 6].Vinayak B. Kadam , Ganesh K. Pakle ,"DEUDS: Data Extraction Using DOM Tree and Selectors", International Journal of Computer Science and Information Technologies, Vol. 5 (2) , 2014, 1403-1410.
- 7].Vivek D. Mohod, Mrs. J. V. Megha,"A Survey on Data Extraction of Web Pages Using Tag Tree Structure", International Journal of Computer Science and Information Technology.
- 8].Laender, 2002b; Ribeiro-Neto,"DEByE (Data Extraction By Example): An interactive data extraction system [ 1999]."
- 9].Appukuti Chandrashekhar,Dr. P. Venkata Subba Readdy, " Html Tag Based Extraction and Tree Merging From Template Page", International Journal of Advance Research in Computer Science and Management Studies.
- 10].Chia-Hui Chang<sup>1</sup>, Mohammed Kayed<sup>2</sup>, Moheb Ramzy Girgis<sup>3</sup> And Khaled Shaalan," Criteria For Evaluating Information Extraction Systems"