



# File Clustering using Forensic Analysis System

G. Madan Kumar<sup>1</sup>, Sunil Kumar. V<sup>2</sup>

<sup>1</sup>M.Tech 2<sup>nd</sup> year, Department of CSE, PBR VITS, Kavali, Nellore, A.P, India

<sup>2</sup>Associate Professor, Department of CSE, PBR VITS, Kavali, Nellore, A.P, India

<sup>1</sup>madankumar688@gmail.com; <sup>2</sup> sunil.vemula1981@gmail.com

---

*Abstract- In this paper computer forensic analysis investigation, thousands of files are generally surveyed. In this much of the data in those files consists of formless manuscript, whose investigation by computer examiners is very tough to accomplish. Clustering is the unverified organization of designs that is data items, remarks, or feature vectors into groups (clusters). To find a noble clarification for this automated method of analysis are of great interest. In particular, algorithms such as K-means, K-medics, Single Link, Complete Link and Average Link can simplify the detection of new and valuable information from the documents under investigation. In This paper we are going to present a tactic that applies text clustering algorithms to forensic examination of computers seized in police investigations using multithreading technique for data clustering. Our experiments show that the Average Link and Complete Link algorithms provide the best results for our application domain. If suit-ably initialized, partition algorithms (K-means and K-medoids) can also yield to very good results. Finally, we also present and discuss several practical results that can be useful for researchers and practitioners of forensic computing.*

*Keywords- Forensic computing, text mining, multithreading, K-Means, Clustering*

---

## I. INTRODUCTION

Extremely vast increase in crime relating to Internet and computers has caused a growing need for computer forensics. In document clustering computer forensics identifies evidence when computers are used in the police investigations of crimes. In this particular application domain, it usually involves examining the thousands of files per computer. This activity exceeds the expert's ability of analysis and understanding of data. In general, for computer forensic analysis we need computer forensic tools that can exist in the form of computer software. Such tools have been developed to help computer forensic investigators in a computer investigation. However, because storage media is growing in size, day by day investigators may have difficulty in locating their points of interest from a large pool of data. In addition, the format in which the data is presented may result in misinforming and difficulty for the investigators. As a result, the process of analyzing large volumes of data may consume a very large amount of time. It may happen that data generated by computer forensic tools may be meaningless at times, due to the amount of data that can be stored on a storage medium and the fact that current computer forensic tools are not able to present a visual overview of all the objects (e.g. files) found on the storage medium [1].

Essentially this paper used for the police investigations through forensic data analysis. Clustering algorithms are typically used for examining data analysis, where there is little or no prior as shown in table there are various algorithm with their parameters like distance which has cosine as well as levenshtein distance which is nothing but a string metric for measuring the difference between two sequences. Informally, the Levenshtein distance between two words is the minimum number of single-character edits (i.e. insertions, deletions or substitutions) required to change one word into the other. The application for levenshtein distance is to in approximate string matching; the objective is to find matches for short strings in many longer texts, in situations where a small number of differences is to be expected. Table also gives the initialization of each algorithm [1].

Seized digital devices [1] can provide precious information and evidences about facts. In this large amount of data analysis purpose we use Digital text analysis text mining technique. In this technique to search string is difficult. Solve the problem in using forensic acquisition and early analysis and textual information extraction and text clustering. Supervised learning tools to categorize data on already defined categories for investigate purposes. In computer forensic [4] analysis hundreds of thousands of files are usually examined. Much of the data in those files consists of unstructured text, whose analysis by computer examiners is difficult to be performed. To overcome these problems applies clustering algorithms to forensic analysis of computer seized in police investigations. Clustering includes labels. Examiner identifies easy and also content quick search. Difficult to [3] identifies specific text string. To solve this problem we are using ranking and indexing algorithms. Automatic approaches for clustering labeling. The assignment of labels to clusters may enable the expert examiner to identify the semantic content of each cluster more quickly. Improve the quality of data analysis. Make a automatic procedure for inferring accurate and [2] easily understandable expert-system-like rules from forensic data. Methodology is based in the fuzzy set theory. TO overcome these problem using fuzzy set theory, and it produces the best result comparing k-means and k-mediods. The accuracy of rules inferred was very high and clearly better than the minimum level required to make them usable in a particular string. Complicates reduce communication experts.

## II. RELATED WORK

The use of clustering has been reported by only few studies in the computer forensics field.[1] Basically, The use of classic algorithm for clustering data is described by most of the studies such as Expectation-Maximization (EM) for unsupervised learning of Gaussian Mixture Models, K-means, Fuzzy C-means (FCM), and Self-Organizing Maps (SOM). These algorithms have well-known properties and are widely used in practice. [4]

An integrated environment for mining e-mails for forensic analysis, using classification and clustering algorithms, was presented in [4]. In a related application domain, e-mails are grouped by using lexical, syntactic, structural, and domain-specific features [6]. Three clustering algorithms (K-means, Bisecting K-means and EM) were used. The problem of clustering e-mails for forensic analysis was also addressed, where a Kernel-based variant of K-means was applied [7]. The obtained results were analyzed subjectively, and the result was concluded that they are interesting and useful from an investigation perspective. More recently, a FCM-based method for mining association rules from forensic data was described [3].

In this paper when we talk about computer forensics there are so many tools, algorithms and methods to do it. So this paper presents those algorithms and methods are going to discuss one by one.

Document Clustering for Forensic Analysis: An Approach for Improving Computer Inspection uses various algorithms and preprocessing technique for giving result as cluster data. Finally in their conclusion they have shown that, the approach presented by them applies document clustering methods to forensic analysis of computers seized in police investigations. Also, they are reported and discussed with several practical results that can be very useful for researchers and practitioners of forensic computing. More specifically, in their experiments the hierarchical algorithms known as Average Link and Complete Link presented the best results. Despite their usually high computational costs, they have shown that those algorithm are particularly suitable for the studied application domain because the dendrograms that they provide offer summarize views of the documents being inspected, thus being helpful tools for forensic examiners that analyze textual documents from seized computers.[1]

A Comparative Study on Unsupervised Feature Selection Methods for Text Clustering describe one of the central problems in text mining and information retrieval area is text clustering. Performance of clustering algorithms will considerably reject for the high dimensionality of feature space and the inherent data sparsity, two techniques are used to deal with this problem: feature extraction and feature selection. Feature selection methods have been successfully applied to text categorization but seldom applied to text clustering due to the unavailability of class label information. Four unsupervised feature selection methods such as DF, TC, TVQ, and a new proposed method TV were introduced in that paper. Experiments were taken to show that feature selection methods can improves efficiency as well as accuracy of text clustering. [5]

Fuzzy Methods for Forensic Data Analysis is again describes a methodology and an automatic procedure for inferring accurate and easily understandable expert-system-like rules from forensic data. In most data analysis environments the methodology and the algorithms used were proven to be easily implementable. By discussing the applicability of different fuzzy methods to improve the effectiveness and the quality of the data analysis phase for crime investigation the fuzzy set theory would get implemented. [3]

In mining write prints from anonymous e-mails for forensic investigation, basically they are collecting e-mails written by multiple anonymous authors and focusing on the problem of mining the writing styles of those e-mails. The general idea is to first cluster the anonymous e-mail by the Stylometric (Stylometry is the application of the study of linguistic style, usually to written language, but it has successfully been applied to music and to fine-art paintings as well) features and then extract the write print, i.e., the unique writing style, from each cluster. [4]

They have mainly focus on lexical and syntactic features of an e-mail as when we talk about lexical features they are used to learn about the preferred use of isolated characters and words of an individual. Following table gives some of the commonly used character-based features, these include frequency of individual alphabets (26 letters of English), total number of upper case letters, capital letters used in the beginning of sentences, average number of characters per word, and average number of characters per sentence. To indicate the preference of an individual for certain special characters or symbols or the preferred choice of selecting certain units the use of such features come in picture. For example most of the people prefer to use „\$“ symbol instead of word „dollar“, „%“ for „percent“, and „#“ instead of writing the word „number“. [4]

Now when we talked about syntactic features, they are also called as style markers which consist of all purpose function words such as „though“, „where“, „your“, punctuation such as „!“ and „:“, parts-of-speech tags and hyphenation etc. as shown in table. [4]

**Table 1** Lexical and Syntactic Features

<b>LEXICAL AND SYNTACTIC FEATURES</b>	
<b>Features type</b>	<b>Features</b>
Lexical: character- based	1. Character count (N) 2. Ratio of digits to N 3. Ratio of letters to N 4. Ratio of uppercase letters to N 5. Ratio of spaces to N 6. Ratio of tabs to N 7. Occurrences of alphabets (A-Z) (26 features) 8. Occurrences of special characters: < > % j { } [ ] / \ @ # w p _ * \$ ^ & O (21 features)
Lexical: word-based	9. Token count(T) 10. Average sentence length in terms of characters 11. Average token length 12. Ratio of characters in words to N 13. Ratio of short words (1e3 characters) to T 14. Ratio of word length frequency distribution to T (20 features) 15. Ratio of types to T 16. Vocabulary richness (Yule's K measure) 17. Hapax legomena 18. Hapax dislegomena
Syntactic features	19. Occurrences of punctuations, . ? ! : ; " " (8 features) 20. Occurrences of function words (303 features)

### III. K-MEANS ALGORITHM IMPLEMENTATION

K-means algorithm is one of the simplest unsupervised learning algorithms that partition feature vectors into k clusters so that the within group sum of squares is minimized. K-means clustering is a method of vector quantization originally from signal processing that is popular for cluster analysis in data [9].

Mining from the fig K-Means follows a simple way to classify a given dataset and looks like Steps:

1. Place randomly initial group centroids into the 2<sup>nd</sup> space.
2. Assign each object to the group that has the closest centroid.
3. Recalculate the positions of the centroids.
4. Finally if the positions of the centroids did not change go to the next step else go to the step2.
5. End.

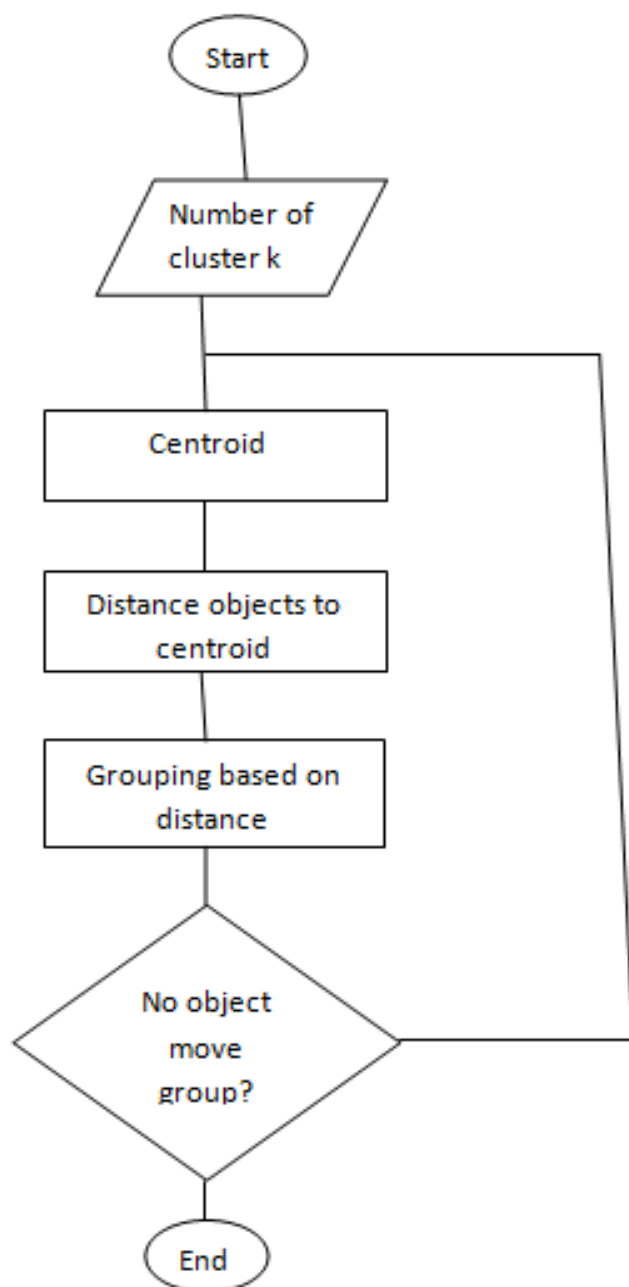


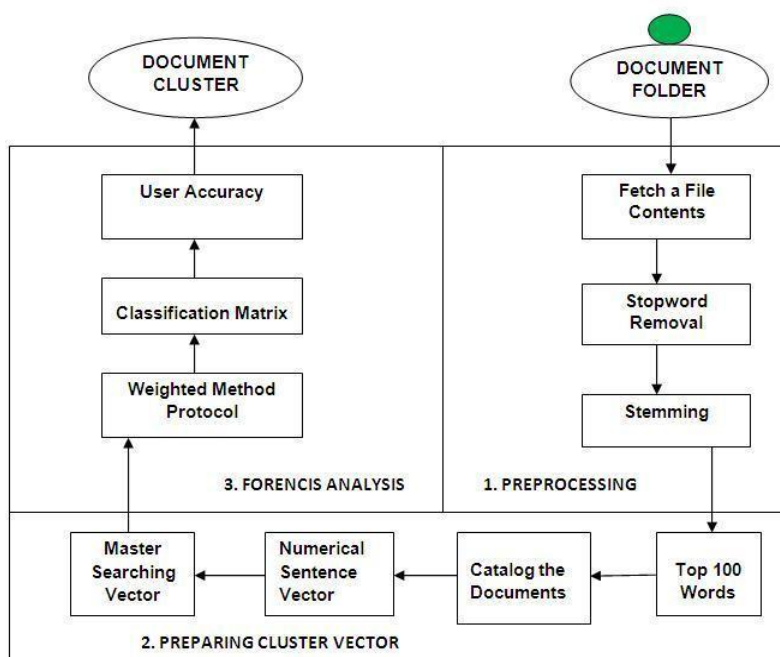
Fig K-Means algorithm

#### IV. PROPOSED FORENSIC ANALYSIS SYSTEM

The proposed forensic analysis system shown in below figure

In our proposed system basically there are three important steps which are as follows

- 1) Preprocessing
- 2) Preparing cluster vector
- 3) Forensic analysis



**Fig:** Architectural diagram of forensic analysis system

- 1) **Preprocessing-** In preprocessing step there are three steps such as a) fetch a file contents, b) stopword removal c) stemming. In all the above steps the basic purpose is to check the file contain and to remove the stop word like a, an ,the etc. and later on to do stemming on that file which will be removing ing and ed words from the given statement.
- 2) **Preparing Cluster Vector-** For preparing the cluster vector one will need to find top 100 words from the file on which preprocessing step is already done. Now from that document or rather way we can say file or data numerical sentences such as the sentence which has numerical word in it that means the sentence which contains date or any kind on number in it.
- 3) **Forensic Analysis-** This will be the last step of proposed method. From the diagram no 1 mention above one can say that for the forensic data analysis classification matrix need to be made with the help of

## V. PERFORMANCE ANALYSIS

### 5.1 Data Set

The data set for forensic analysis will be different number of file in different formant which has information on which data clustering is performed by applying dissimilar algorithm. For the clustering processes this paper makes use of multithreading technique. Later on that data set can be used for police investigation.

### 5.2 Result Set

The result set produced by this system will be number of clusters formed by applying algorithm on given information.

## VI. CONCLUSION

By doing the survey on computer forensic analysis it can be concluded that clustering on data is not an easy step. There is huge data to be cluster in compute forensic so to overcome this problem, this paper presented an approach that applies document clustering methods to forensic analysis of computers seized in police investigations. Again by using multithreading technique there will be document clustering for forensic data which will be useful for police investigations. It reduces the work of data examiner. It helps to police departments, because the terrorist missing the evidence of device. It searches and examines gives the knowledge about attacks. So it is very helpful to prevent attacks.

## REFERENCES

- [1] L. F. Nassif and E. R. Hruschka, "Document clustering for forensic computing: An approach for improving computer inspection," in Proc. Tenth Int. Conf. Machine Learning and Applications (ICMLA), 2011, vol. 1, pp. 265–268, IEEE Press.
- [2] L. Liu, J. Kang, J. Yu, and Z. Wang, "A comparative study on unsupervised feature selection methods for text clustering," in Proc. IEEE Int. Conf. Natural Language Processing and Knowledge Engineering, 2005, pp. 597–601.
- [3] K. Stoffel, P. Cotofrei, and D. Han, "Fuzzy methods for forensic data analysis," in Proc. IEEE Int. Conf. Soft Computing and Pattern Recognition, 2010, pp. 23–28.
- [4] R. Hadjidj, M. Debbabi, H. Lounis, F. Iqbal, A. Szporer, and D. Benredjem, "Towards an integrated e-mail forensic analysis framework," Digital Investigation, Elsevier, vol. 5, no. 3–4, pp. 124–137, 2009.
- [5] B. K. L. Fei, J. H. P. Eloff, H. S. Venter, and M. S. Oliver, "Exploring forensic data with self-organizing maps," in Proc. IFIP Int. Conf. Digital Forensics, 2005, pp. 113–123.
- [6] F. Iqbal, H. Binsalleeh, B. C. M. Fung, and M. Debbabi, "Mining writeprints from anonymous e-mails for forensic investigation," Digital Investigation, Elsevier, vol. 7, no. 1–2, pp. 56–64, 2010.
- [7] S. Decherchi, S. Tacconi, J. Redi, A. Leoncini, F. Sangiacomo, and R. Zunino, "Text clustering for digital forensics analysis," Computat. Intell. Security Inf. Syst., vol. 63, pp. 29–36, 2009.
- [8] Stoffel .K, Cotofrei .P, and Han.D, "Fuzzy methods for forensic analysis", Proceedings of the International Conference soft computing and pattern Recognition, pp. 23-28, 2010.
- [9] Girolami .M, "Mercer Kernel Based Clustering in featurespace" IEEE Transaction on neural networks, Vol.13, pp. 2780-2784, 2002.

## SHORT BIOGRAPHY



**Mr. G. Madan Kumar** received the **B.Tech** Degree in Computer Science and Engineering from Jawaharlal Nehru Technological University, Anantapur, in **2012**. He currently pursuing **M.Tech (CSE) in Dept of Computer Science and Engineering** in PBR VITS Engg College, kavali, Nellore, under JNTUA University, Anantapur.



**Mr. Vemula.V. Sunil Kumar** has received his B.Tech in Electrical Communication Engineering and M.Tech degree in Computer science from JNTU, Hyderabad in 2002 and 2008 respectively. He is dedicated to teaching field from the last 11 years and he has 1 year industrial experience in BEL at Hyderabad. He has guided 12 P.G Students and 25 U.G students. His research areas included CN- MANETs, Neural Networks, and Image processing, embedded systems. At present he is working as Associate professor in PBR Visvodaya Institute of Technology & Science, Kavali, Andhra Pradesh, India.