Artificial Intelligence formulated this projection for compatibility purposes from the original article published at Global Journals. However, this technology is currently in beta. *Therefore, kindly ignore odd layouts, missed formulae, text, tables, or figures.*

A New Approach for Improving Computer Inspections by using Fuzzy Methods for Forensic Data Analysis P. Jyothi¹ ¹ MITS college madanapalli Received: 13 December 2014 Accepted: 5 January 2015 Published: 15 January 2015

7 Abstract

Now a day?s digital world data in computers has great significance and this data is extremely 8 critical in perspective for upcoming position and learn irrespective of different fields. 9 Therefore we the assessment of such data is vital and imperative task. Computer forensic 10 analysis a lot of data there in the digital campaign is study to extract data and computers 11 consist of hundreds of thousands of files which surround shapeless text or data here clustering 12 algorithms is of plays a great interest. Clustering helps to develop analysis of documents 13 under deliberation. This document clustering analysis is extremely useful to analyze the data 14 from seized devices like computers, laptops, hard disks and tablets etc. There are total six 15 algorithms used for clustering of documents like K-means, K-medoids, single link, complete 16 link, Average Link and CSPA. These six algorithms are used to cluster the digital documents. 17 Existing document clustering algorithms are operated in single document at a time. In the 18 proposed approach of these working algorithm applied on multiple documents at a time. Now 19 we using clustering technique named as agglomerative hierarchical clustering which gives 20

²¹ better finer clusters compared to existing techniques.

22

23 Index terms— clustering, forensic analysis, clustering algorithms, hierarchical agglomerative clustering 24 algorithm.

25 1 Introduction

n computer forensic process is impacted by large amount of data. This has roughlydistinct asrestraints that
mergeelement of law and computer sciences to gather and examine information from computer systems. In our
study there are hundreds of files are there in instructed format. For this analysis they have some methods like
machine learning and data mining are of great importance. Clustering algorithms are usually needed to grouping
data in files, where there is practically no prior knowledge about the information [1] [13].

From a more specialized perspective, our datasets comprise of unlabeled objects. In addition, actually expecting 31 that named datasets could be accessible from previous analysis, there is very nearly no hope that the same classes 32 would be still legitimate for the upcoming information, got from different computers also related to different 33 34 examinations. More definitely, it is likely that the new information would come from different locations. In this 35 way, the utilization of clustering algorithms, which are fit for discovering latent patterns from content documents 36 found in seized computers, can improve the analysis performed by the expert examiner. The methods of rational 37 clustering algorithm objects within a substantial group are more like one another than they are two objects belongs to alternative group [1]. Along those data partition has been actuated from data. The export examiner 38 may concentrate on interesting on delegated documents from the obtained set of groups by performing this task 39 of examination of all documents. In a more functional and sensible situations, domain experts are rare and have 40 limited time accessible for performing examinations. Therefore it is sensible to expect that finding a significant 41 document. The examiner could prioritized the investigation of different documents belongs to the cluster interest. 42

6 FRAMEWORK REQUIREMENTS

Clustering algorithm has been mulled over for a considerable length of time and the literature on the subject is huge. Therefore, we decided to demonstrate the capability of the proposed methodology, namely: the partition algorithms K-Means, K-Medoids, the hierarchical single link, complete link, average link and the cluster ensemble algorithm known as CSPA [3]. It is well known that the number of clusters demonstrating parameter of many algorithms and it is generally having an earlier knowledge. However the number of clusters has not been examined in the computer forensics. Really we could not even spot one work that is sensibly close in its application area and that reports the utilization of number of algorithms capable finding the number of clusters [3].

50 **2** II.

51 3 Review of Related Research

In our software development process research is the most important one. In this is based on the time factors, economy and company strength we can determine the developing process. Once the programmer start the work based on experts suggestions and gather related information to different websites based on their work. Before building the system each and developer can maintain the above requirements report.

C. M. Fung et al. [7] have present an agglomerative and divisive hierarchical clustering to group the documents into clusters, cluster in a documents have high similarity to each other, the dissimilar documents into other group. Likewise no labeled documents are provided in clustering. Hence these are known as unsupervised learning. Hence the clusters are analyzed into a tree that facilitates browsing. The parent-child relationship among the nodes in the tree can be viewed as a topic-subtopic relationship.

B. Fei et al. [3] have discusses the application of a self-organizing map (SOM) is to support decision making by computer forensic investigators and assist them in conducting data analysis in a more efficient manner and also SOM produces patterns similarity in data sets. Author explores great ability to interpret, explore data generated by computer forensic tools.

Alexander Strehlet al. [2] hasintroduces three effective and efficient techniques to obtaining high quality combiners. In first combiner induce Partitioning and re-clustering of objects is based on similar measure, second is based on hyper-graph and third is based on collapse group of clusters into meta-clusters which participate to find individual object to the combined clustering. By using the three approaches to provide the low computational costs and feasible to use a supra-consensus function against the objective function and provide the best results.

L. F. Nassif et al. [5] have present an approach that applies document clustering algorithms to forensic analysis of computers seized in police examinations. Author represents experimentation with six well-known clustering calculations (K-Means, K-medoids, Single Link, Complete Link, Average Link, and CSPA) applied to five certifiable datasets acquired from computers seized in true examinations. Investigations have been performed with different combination of parameters, resulting in 16 different instantiations of algorithms. Moreover, two relative legitimacy records were utilized to consequently appraise the quantity of clusters. In the event that suitably introduced, partition algorithms (Kmeans and K-medoids) can likewise respect great results.

Ying Zhao et al. [8] have use high quality clustering algorithms play an essential part in giving intuitive
navigation and browsing mechanism by sorting large amount of data into a little number of meaningful clusters.
Specifically clustering algorithm, hierarchical clustering that assemble meaningful hierarchies out of large amount
of accumulations. This all concentrates on document clustering algorithm that manufactures such various leveled
solutions.

⁸² 4 a) Hierarchical Agglomerative Algorithm

Hierarchical clustering algorithms are either topdown or bottom up. Bottom-up algorithms treat each one record
as a single cluster at the beginning and after that progressively merge (or agglomerate) sets of clusters until all
groups have been merged into a single group that contains all documents. ??ottom

⁸⁶ 5 i. Cluster Formation

By passing up from the bottom to top process of clusters the dendrograms used to reproduce the historical environment of consolidations that brought about the represented clustering. For example, we see that the two documents permittedAddressBook_Deletecopy.html and AddressBook_New_Action.html were combining in Figure (1) and that the last combine added Visitors.html to a cluster comprising of the other 27 documents. ii. Finding Similarity The Hierarchical Agglomerative clustering is usually defined as a dendrograms as

⁹² illustrate in Fig (1). Each combination is signified by a horizontal line. That line represents similarity between

 $_{93}$ two documents, where documents are viewed as single clusters. We call this similarity the combination similarity

of the combined documents. For example, the combination similarity of the documents Schedule_New.html and

ToDos_Index.html consisting of Figure (??) is ?0.19. This defines the cosine similarity of clusters as 1.0.

97 6 Framework Requirements

For finding the meaningful data from the dataset, researchers have used data mining techniques, in which of clustering is one of the popular techniques. Let DS will taking as our dataset represented as $DS = \{d \ 1 \ d \ 2, d \ 1 \ d \ 2\}$? d n }; 1? I ? n, where n is the number documents in a dataset DS. In our propose system basically there are
 three important steps which are as follows 1) Preprocessing 2) Cluster Formulation 3) Forensic analysis

102 7 Preprocessing

In preprocessing step there are three steps such as a) fetch a file contents, b) Stemming, c) Stop word Removal. 103 These 3 steps are used to remove the noise and inconsistent data. In first step fetch the dataset and perform the 104 second operation with the help of porter stemming. In this stemming is based on the idea that the suffixes in the 105 English language are mostly made up of a combination of smaller and simpler suffixes. If the words end with ed, 106 ing, ly etc that words are removed. This step is a linear step stemmer [16]. In this last step is remove the stop 107 words with the help of Stop token filter.[17] Stop words in a document like to, I, has, the, be, or etc. stop words 108 are the foremost frequent words with in the English language. Stop words blot your index while not providing 109 any additional worth. 110

¹¹¹ 8 Global Journal of C omp uter S cience and T echnology

Volume XV Issue V Version I Year 2015 () C frequent words with in the English language. Stop words blot your index while not providing any additional worth. At that point, we received a customary statistical methodology for text mining, in which documents are meant in a vector space model. In this each one model, each one document is denoted by vector containing the frequencies of events of words. To process the distance between reports, two measures have been utilized specifically: cosine-based separation and hierarchical agglomerative clustering. After these steps our data will be relevant.

¹¹⁸ 9 Cluster Formulation

This session exhibits the mining of datasets from the preprocessed dataset. For each document the similarity of the concentrated words from the preprocessed step is processed and the top comparability documents are clustered first this. This session depicts the mining of successive item sets from the preprocessed content documents. For each document the recurrence of the concentrated words from the preprocessing step is registered and the top continuous words from each are taking out. From the set of top frequent words, the binary database is framed

124 by getting the unique words.

¹²⁵ 10 a) Hierarchical Agglomerative Clustering Algorithm

Hierarchical agglomerative algorithms treat every one document as a singleton cluster toward the starting and 126 thereafter dynamically consolidation set of clusters until all clusters have been melded into a single cluster that 127 contains all documents. Input: List of Documents D=d 1, d 2? d n Output: Clusters resultC= {c 1, c 2?c 128 n 1. For i=1 to n do 2. For the given list of documents each document is treated as a specified 3. Finding 129 parsers //those are theunique words in documents 4. Suppress non-dictionary words 5. Get unique edges in this 130 documents 6. Initialize clusters a. For n?1 to N b. Applying clustering to the items Constructing histogram 131 //for analyzing clusters h min should be 1.0; h max should be 0.0 A New Approach for Improving Computer 132 Inspections by using Fuzzy Methods for Forensic Data Analysis For T to 1?n-2 For J to 1?n-1 t sim =sim(doc[t], 133 doc[j] {If(h min >t sim) h min =t sim ; If (h max <t sim) h max =t sim ; } 134

¹³⁵ 11 Finding analogousity 3. Forensic Analysis

This will be the last step in proposed method. Here the algorithm process initially provides a topological arrangement between neurons at convergence of documents. Here we can analyze the number of clusters from our selected dataset. At final step this process will calculate the similarity between formed documents with less time compared to other algorithms.

140 **12 III.**

¹⁴¹ 13 Implementation

The implementation process of clusters can done through number of steps that process will needed for the purpose 142 of good cluster similarity between clusters. The follower can do these steps very care full. In this process first 143 collect the documents from local systems then perform preprocessing-In preprocessing step there are three steps 144 145 such as a) fetch a file contents, b) Stemming, c) Stop word Removal. These 3 steps are used to remove the noise 146 and inconsistent data. In first step remove the stop word prepositions, pronouns, irrelevant documents data(147 a, an the etc) ??17] and later on to do stemming on that file which will be removing Portuguese words(ing and edetc) [16] from the upcoming data. At that point, we received a customary statistical methodology for 148 text mining, in which documents are meant in a vector space model. ??19] In this model, each one document is 149 denoted by vector containing the frequencies of events of words. To process the distance between reports, two 150 measures have been utilized specifically: cosine-based separation and hierarchical agglomerative clustering. After 151 these steps our data will be relevant. 152

- 153 14 Global
- 154 15 3.Forensic Analysis
- 155 16 Similarity graph

156 Cluster similarity

157 **17** Cluster Counting

18 Data Analysis Agglomeration Document Clustering Deter mine Clusters

160 A New Approach for Improving Computer Inspections by using Fuzzy Methods for Forensic Data Analysis

The next step of this process will take clusters counting. Here, data will collect from the previous step. 161 Then analyze data for estimating relevant clusters, by using agglomerative algorithm. In this algorithm data is 162 163 analyzes from multiple documents, and then divide similar documents and dissimilar documents. Based on the priority of data high priority documents are saved under one cluster and comparison of first cluster less similarity 164 documents are stored under next cluster based on the algorithmic perspective. Then, the last step algorithm 165 finds the similarity of all cluster consisting documents. In this comparison include clusters containing each and 166 every document is compared and finds similarity between them. This algorithm plays a important role in this 167 process compared to other algorithms. 168

169 **IV**.

170 20 Experimental Results

In this proposed approach experimentation developed by java (JDK 2.0). In the process of running and executing the main file. After executing main file dataset containing documents are loaded that are shown in figure (1). Here we are taking 26 documents[0-25] these all documents are under 8 different areas. Like Address Book [0-2], Members [8][9][10][11][12][13][14] etc those are shown in figure (2). These 8 different partitions are clustered into 4 groups named as C1,C2,C3, C4. In C1 under documents are [0, 1,7,8,9,11,15], C2 under documents are [2,3,4,5,6,13,14,16, ??7, ??8, ??9, ??0, ??1, ??2, ??3, ??4], C3 under documents are [10,12] and C4 under documents are ??25] shown in figure (3) and finding similarity between all documents shown in figure (4).

$_{178}$ 21 Conclusion

We use an approach for clustering documents which can become an ideal application forensic analysis of computers. There are several practical results based on our work which are extremely useful. In our work, the algorithm known as hierarchical clustering algorithm that yields the best results. In spite of this algorithm we find the number of similar and dissimilar documents in our input documents and also finding the similarity between the documents. ^{1 2}

 $^{^{1}}$ © 2015 Global Journals Inc. (US)

 $^{^{2}}$ © 2015 Global Journals Inc. (US) 1



Figure 1: Figure 1:



Figure 2: CFigure 1 :



Figure 3: GlobalCFigure (1)CFigure (3)



Figure 4:



Figure 5:

21 CONCLUSION

- [Jain and Dubes ()] Algorithms for Clustering Data, A K Jain , R C Dubes . 1988. Englewood Cliffs, NJ: Prentice Hall.
- [Bishop ()] C M Bishop . Pattern Recognition and Machine Learning, (New York) 2006. Springer-Verlag.
- [Strehl and Ghosh ()] 'Cluster ensembles: A knowledge reuse framework for combining multiple partitions'. A
 Strehl , J Ghosh . J. Mach. Learning Res 2002. 3 p. .
- [Mirkin ()] Clustering for Data Mining: A Data Recovery Approach, B Mirkin . 2005. London, U.K.: Chapman
 & Hall.
- 191 [Hubert and Arabie ()] 'Comparing partitions'. L Hubert, P Arabie . J. Classification 1985. 2 p. .
- [Nassif and Hruschka ()] 'Document clustering for forensic computing: An approach for improving computer
 inspection'. L F Nassif, E R Hruschka. Proc. Tenth Int. Conf. Machine Learning and Applications (ICMLA),
- 194 (Tenth Int. Conf. Machine Learning and Applications (ICMLA)) 2011.
- [Loewenstein et al. ()] 'Effcient algorithms for exact hierarchical clustering of huge datasets: Tackling the entire protein space'. Y Loewenstein , E Portugaly , M Fromer , M Linial . *Bioinformatics* 2008. 24 (13) p. .
- [Zhao and Karypis ()] 'Evaluation of hierarchical clustering algorithms for document datasets'. Y Zhao , G
 Karypis . Proc. CIKM, (CIKM) 2002. p. .
- 199 [Everitt et al. ()] B S Everitt, S Landau, M Leese. Cluster Analysis, (London, U. K) 2001. Arnold.
- [Fei et al. ()] 'Exploring forensic data with self-organizing maps'. B K Fei , J H P Eloff , H S Venter , M S Oliver
 Proc. IFIP Int. Conf. Digital Forensics, (IFIP Int. Conf. Digital Forensics) 2005. p. .
- 202 [Stoffel et al. ()] 'Fuzzy methods for forensic data analysis'. K Stoffel , P Cotofrei , D Han . Proc. IEEE Int.
- 203 Conf. Soft Computing and Pattern Recognition, (IEEE Int. Conf. Soft Computing and Pattern Recognition)
 204 2010. p. .
- [Zhao et al. ()] 'Hierarchical clustering algorithms for document datasets'. Y Zhao , G Karypis , U M Fay Yad .
 Data Min. Knowl. Discov 2005. 10 (2) p. .
- [Kishida ()] 'High-speed rough clustering for very large document collections'. K Kishida . 10.1002/asi.2131. J.
 Amer. Soc. Inf. Sci 2010. 61 p. .
- [Haykin ()] Neural Networks: A Comprehensive Foundation, S Haykin . 1998. Englewood Cliffs, NJ: Prentice Hall.
- [Decherchi et al. ()] 'Text clustering for digital forensics analysis'. S Decherchi , S Tacconi , J Redi , A Leoncini
 F Sangiacomo , R Zunino . Computat. Intell. Security Inf. Syst 2009. 63 p. .
- 213 [Waegel ()] Daniel Waegel . The porter Stemmer". CISC889/Fall, 2011.