



# Applying Data Mining Principles in the Extraction of Digital Evidence

**Raburu George<sup>1</sup>, Omollo Richard<sup>2</sup>, Okumu Daniel<sup>3</sup>**

<sup>1,2,3</sup>Department of Computer Science and Software Engineering

Jaramogi Oginga Odinga University of Science and Technology, Kenya

Email: [graburu@hotmail.com](mailto:graburu@hotmail.com)<sup>1</sup>; [comolor@hotmail.com](mailto:comolor@hotmail.com)<sup>2</sup>; [okumudaniel6@gmail.com](mailto:okumudaniel6@gmail.com)<sup>3</sup>

---

*Abstract - Data mining in databases is part of the interdisciplinary field of knowledge discovery used for extracting patterns and relationships amongst the data sets from a data source. This paper focuses on mining data from a database. The study further identifies how data mining techniques can be used to support digital forensic investigation in a digital crime case. Digital forensic models are reviewed, concepts mapped onto experimental cases and tested on a Pentium (R) Core™ 2 Duo CPU, 1.89 GHz with Windows XP Professional OS. The findings supported the benefits of integrating data mining principles in extracting data for digital evidence.*

*Keywords: - Data Mining techniques, Digital forensics, digital Investigation, Data Recovery*

---

## I. INTRODUCTION

Data mining research began in the 1980s and rapidly grew in the 1990s. Specific techniques that have been developed within disciplines such as machine learning, artificial intelligence, and pattern recognition have been successfully employed in data mining. Data mining has been successfully introduced in many different fields with World Wide Web as the most recent, others including field of criminal forensics (Digital forensics). It includes detecting deceptive criminal identities, identifying groups of criminals who are engaging in various illegal activities and many others. Data mining techniques typically aim to produce insight evidence from large volumes of data. Digital forensics is a sophisticated and cutting edge area of breakthrough research. Canvass of digital forensic investigation and application is growing at a rapid rate with mammoth digitization of an information economy. Law enforcement and military organizations have heavy reliance on digital forensic today. As information age is revolutionizing at a speed inconceivable and information being stored in digital form, the need for accurate intellectual interception, timely retrieval, and nearly zero fault processing of digital data is crux of the issue.

The goal of traditional forensic analysis is to provide accurate information derived through the use of proven and well-understood methodologies. Forensic Science applied in courts of law has sought to use commonly applied techniques and tools only after rigorous, repetitive testing and thorough scientific analysis. In 2006 Richard and Roussev discussed the urgent need for new tools and strategies for the rapid turnaround of large forensics targets. Their focus was on the acquisition and analysis of forensic evidence and argued that current forensics tools are inadequate given the increased complexity of cases, increased size of targets, better awareness of the capabilities of digital forensics, and multi-computing scenarios. Data mining is the application of algorithms for extracting patterns from data [1]. These extracted patterns will provide useful knowledge to decision makers. As such, there has been an increasing demand for data mining tools to help

organizations uncover knowledge that can guide in decision making.

For law enforcement agencies it is tedious and time consuming task to find out the user from the large quantity of the acquired digital evidence drive, here we propose a method for finding user ownership information based on attribute analysis of evidence drive. With the growing sizes of databases, law enforcement and intelligence agencies face the challenge of analyzing large volumes of data involved in criminal and terrorist activities [2]. Thus, a suitable scientific method for digital forensics is data mining.

The useful seamless integration of data mining techniques with digital forensic science has been depicted at analysis phase. This helps in boosting up the performance and the reliability of investigations of the subjects. The formal methodology of data mining includes following basic steps [3]:

- (i) Determine the nature and structure of the representation of the data sets.
- (ii) Decide how to quantify the data; compare how well different representations fit the data.
- (iii) Choose an algorithmic process to optimize the scoring function.
- (iv) Decide what principles of data management are required to implement the algorithms efficiently.

Data mining functionalities are used to specify the various types of patterns to be looked for. The first part of this paper will study and review the current know-how in the field of digital forensics. Subsequently, common thread of conventional forensic tools will be investigated. The results of this analysis can be used to create a basic classification of computer forensic tools. This classification serves as a foundation for the identification of inherent limitations of current computer forensic tools and recommendations for the breakthrough improvements through data mining techniques for ushering in state of the art digital forensic tools.

## II. LITERATURE REVIEW

### 2.1 Principles of Data Mining

This explains and explores the techniques of data mining, that is to say, classification, association rule mining and clustering.

- Classification: mining patterns that can classify future data into known classes.
- Association rule mining: mining any rule of the form  $X \rightarrow Y$ , where  $X$  and  $Y$  are sets of data items.
- Clustering: identifying a set of similarity groups in the data

Other techniques are discussed later in this research paper.

### 2.2 Digital Forensic Investigation Process Model

The first framework developed at the DFRWS broke digital forensics down into a seven step process. They are identification, preservation, collection, examination, analysis, presentation, and lastly decision [4]. It was this time also that data mining was classified as a key area of research under the analysis step illustrates in figure1.

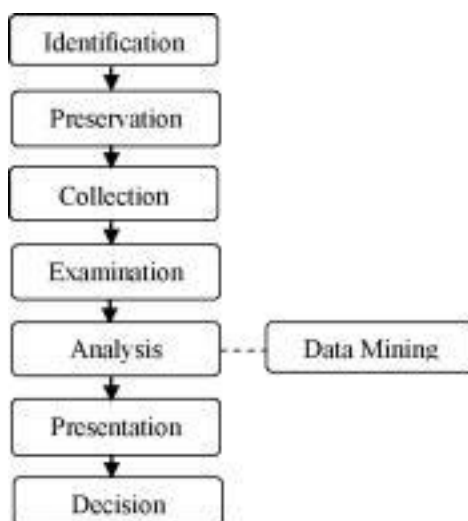


Fig. 1 Data Mining In Digital Forensic Investigation Process

### **2.2.1 Identification**

This recognizes an incident from indicators and determines its type.

### **2.2.2 Preservation**

This phase includes packaging, transportation and storage. Appropriate procedures should be followed and documented to ensure that the electronic evidence collected is not altered or destroyed. All potential sources of evidence should be identified and labeled properly before packing.

### **2.2.3 Collection**

Evidence collection of the digital or mobile devices is an important step and required a proper procedure or guideline to make them work. This can be categorized into two categories: Volatile Evidence Collection, and Non-Volatile Evidence Collection

### **2.2.4 Examination**

This phase involves examining the contents of the collected evidence by forensic expert and extracting information, which is critical for proving the case. Appropriate number of evidence backups must be created before proceeding to examination. This phase aims at making the evidence visible, while explaining its originality and significance. Huge volumes of data collected during the volatile and non-volatile collection phases need to be converted into a manageable size and form for future analysis. Data filtering, validation, pattern matching and searching for particular keywords with regard to the nature of the crime or suspicious incident, recovering relevant ASCII as well as non-ASCII data etc. are some of the major steps performed during this phase. Personal organizer information data like address book, appointments, calendar, scheduler etc, text messages, voice messages, documents and emails are some of the common sources of evidence, which are to be examined in detail. Finding evidence for system tampering, data hiding or deleting utilities, unauthorized system modifications etc. should also be performed. Detecting and recovering hidden or obscured information is a major tedious task involved. Data should be searched thoroughly for recovering passwords, finding unusual hidden files or directories, file extension and signature mismatches etc.

### **2.2.5 Analysis**

This step is more of a technical review conducted by the investigative team on the basis of the results of the examination of the evidence. Identifying relationships between fragments of data, analyzing hidden data, determining the significance of the information obtained from the examination phase, reconstructing the event data, based on the extracted data and arriving at proper conclusions etc. are some of the activities to be performed at this stage.

### **2.2.6 Presentation**

This phase includes packaging, transportation and storage. Appropriate procedures should be followed and documented to ensure that the electronic evidence collected is not altered or destroyed. All potential sources of evidence should be identified and labeled properly before packing. Use of ordinary plastic bags may cause static electricity. Hence antistatic packaging of evidence is essential. The device and accessories should be put in an envelope and sealed before placing it in the evidence bag.

### **2.2.7 Decision**

The final stage in the model is the decision phase. This involves reviewing all the steps in the investigation and identifying areas of improvement. As part of the decision phase, the results and their subsequent interpretation can be used for further refining the gathering, examination and analysis of evidence in future investigations. In many cases, much iteration of examination and analysis phases are required to get the total picture of an incident or crime. This information will also help to establish better policies and procedures in place in future.

## **2.3 Data Mining in Digital Forensics**

Data mining can be defined as the analysis of large observational data sets to find un-suspected relationships and to summarize the data in novel ways that are both understandable and useful to the owner [5].

Data mining can be relevant when the data sets of interest is large; if they were not then it might be feasible to

manually explore the data and make a decision. There are varying scales of data sets that may be considered to be large, but this paper focuses on data and files on a single hard drive which may easily number to millions of files at a time. Once validated on a single hard drive, the goal is scale the data mining effort to much larger data sets across multiple hard drives and possibly networks. Secondly as defined, a goal of data mining is to find unsuspected or unknown relationships within data. Obviously there is no reason to report or to repackage already known relationships through data mining. We found the relationships among files or data, specifically the ascription or ownership of the data. From a forensics perspective such correlations or relationships discovered may be used to tie criminals, terrorists, or people of interest together. Lastly, the results of data mining must be understandable and useful.

### **2.3.1 Role of Data Mining in Digital Forensic**

Data mining and soft computing has several applications in digital forensics. These include identifying correlations in forensic data (association), discovering and sorting forensic data into groups based on similarity (classification), locating groups of latent facts/clustering, and discovering patterns in data that may lead to useful predictions/forecasting[6]. While this technique is ideal for association, classification, clustering and forecasting, it is also particularly useful for visualization. [7]. Visualization enables digital investigators to locate vital information that is of interest rapidly and efficiently. In addition, it can guides digital investigators towards the best next step in their search so that digital evidence recovery is carried out in a more efficient and effective manner [8]. In 2003, the Artificial Intelligence Lab at the University of Arizona, presented an overview of case studies done with relation to their COPLINK project. The project's specific interest was how information overload hindered the effective analysis of criminal and terrorist activities by Law enforcement and national security personnel.

Their work proposed the use of data mining to aid in solving these issues. In their report they define data mining in the context of crime and intelligence analysis to include entity extraction, clustering techniques, deviation detection, classification, and lastly string comparators. Four case studies in the report showed how data mining was useful in extracting entity information from police narrative reports, detecting criminal identity deceptions, authorship analysis in cybercrime, and lastly criminal network analysis. Today, COPLINK is software that has been successfully deployed in the field, and works by consolidating, sharing, and identifying the information from online databases and criminal records [9]. Work done by Hewlett Packard in 2005 applied data mining to solve their problem of finding similar files in large document repositories [10]. The end analysis yielded clusters of related files and was further enhanced by applying a graph bipartite partitioning algorithm [11]. In 2006, Galloway and Simoff experimented with a case study redefining an approach to network data mining. In their work, they defined network data mining as identifying emergent Networks between large sets of individual data items. [12]. Shatz, Mohay, and Clark in 2006 explored a correlation method for establishing provenance of time stamped data for use as digital evidence. This work has a deep and relevant impact on digital forensics research as it reiterated the complexity issues of dealing with timestamps because of clock skew, drift, offsets, and possible human tampering. [13]. In 2006 as well, research done by Abraham explored event data mining to develop problems for computer forensic investigation purposes. He analyzed computer data in search of discovering owner or usage profiles based upon sequences of events which may occur on a system. He categorized an owner profile with four different attributes: subject, object, action, and time stamp [14]. In 2007, Beebe and Clark in their work proposed pre-retrieval and post-retrieval clustering of digital forensics text string search results. Though their work is focused on text mining, the data clustering algorithms used have shown success in efficiency and improving information retrieval efforts [15].

There exists a formal methodology for data mining which includes these basic steps:

- (i) Determine the nature and structure of the representation of the data sets to be used.
- (ii) Decide how to quantify the data; compare how well different representations fit the data.
- (iii) Choose an algorithmic process to optimize the scoring function. Decide what principles of data management are required to implement the algorithms efficiently [16]
- (iv) Determining the causal effect from relationships obtained which is relevant to an investigator

### **2.4 Applications of Data Mining in Digital Forensics**

Data mining techniques are designed for large volumes of data. Hence, they are able to support digital investigations. While such techniques have been employed in other fields, their application in digital forensics

is still relatively unexplored.

Some of the data mining techniques applied in digital forensics are as follows:

#### 2.4.1 Association rules

It has been employed to profile user behavior and identify irregularities in log files such irregularities can assist in locating evidence that might be crucial to a digital investigation [17]. In digital forensic, association rule mining can play an important role as to extract the login information of a user from log files of a computer system. By generating the rule sets, with the help of behavioral profiles of the user, forensic expert can detect the behavioral anomalies of users.

#### 2.4.2 Outlier analysis

It has been utilized to locate potential evidence in files and directories that have been hidden or that are different from their surrounding files and directories [18]. Outlier analysis applied in digital forensic is to locate hidden files, directory structure of the files, and the characteristics of each file within a directory are compared to detect potential outliers. This approach is similar to that used when locating hidden directories where the characteristics of directories at the same level are compared.

#### 2.4.3 Support vector machines (SVM)

The SVM have been utilized in several research areas in the field of digital forensics. A support vector machine (SVM) is an algorithm for classification that seeks categorized data based on certain fundamental features of the data [19]. In one instance, a support vector machine was applied to determine the gender of the author of an e-mail based on the gender-preferential language used by the author. In another instance, a support vector machine was applied to determine the authorship of an email. Based on the content of the e-mail, each e-mail was classified according to its likely author. Image mining is one of the many activities undertaken during a digital investigation. A support vector machine can also be used to recognize certain patches or areas of an image [20]. Other instances where support vector machines have been utilized include image retrieval and in executing search queries on images containing suspicious objects [10].

#### 2.4.4 Discriminant analysis

Discriminant analysis has been employed in digital forensic to determine whether contraband images, such as child pornography, were intentionally downloaded or downloaded without the consent of the user [6]. Often, individuals prosecuted for crimes based on digital evidence claim that a Trojan horse or virus installed on their computer system was responsible. In this instance, discriminant analysis provided a mechanism for event reconstruction and enabled digital investigators to counter the Trojan defence by examining the characteristics of the data.

#### 2.4.5 Bayesian networks

It has been used to automate digital investigations. Bayesian networks are based on Baye's theorem of posterior probability [21]. A Bayesian network is a directed acyclic graph which models probabilistic relationships among a set of random variables. The aim was to gather information about likely attacks, actions performed by attackers, the most vulnerable software systems and the investigation techniques that should be used.

### 2.5 Digital Forensic Investigation

A framework, for seamless communication, between the technical members of the digital forensic investigation team and the non-technical members of the judicial team, is very necessary. Defining a generic model for digital forensic investigation, sometimes pose a problem taking into account the varied devices available today. This framework is logical in its outline, scientific in its approach though it is to be adapted to comply with all the legal requirements of the country where the incident has occurred. We are proposing an efficient model for both as economical and time factor. The architecture of the model is illustrates in figure 2.

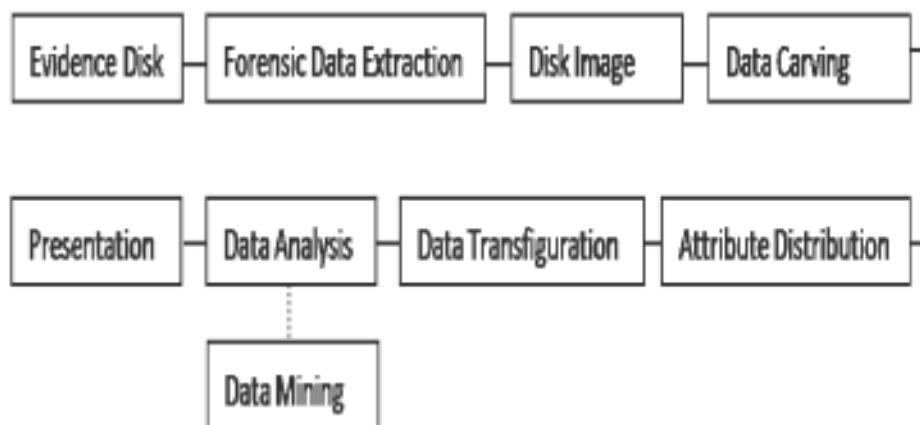


Fig. 2 Architecture of the Model [22]

### 2.5.1 Forensic Data Extraction

In this section, we describe the forensic data extraction process. We conducted an empirical study using selected digital forensic tools that are predominantly used in practice. Since each utility does some specific functionality, collection of such tools were necessary to perform a comprehensive set of functionalities. Hence, the following forensic utilities / tools were adopted to conduct the experimental investigation in this research work:

In an investigation, the analysis phase is the one that most relies on the investigator's skills and experience. To analyze the data collected about a case, an investigator wants to understand and know where the suspect might have hidden the data and in what formats and what application he might have used. Some patterns in the data are important; once found and fully examined; they can lead to more evidence. In order to achieve a successful analysis, many tools are adopted to aid investigators analyze the collected data. In this research we are using *Encase forensics* Guidance Software Inc. [23] for extracting disk image from the suspect hard drive. This tool displays the files of a storage media and allows the user to navigate through the files similar to traditional file explorers. However, they provide additional features that are useful in forensics context such as displaying file headers and opening compressed files. Some of these tools provide more contextual analysis features such as queries and a time-line view of the files. However, the investigator is responsible for manually performing the analysis and gathering knowledge from the extracted data.

## III. METHODOLOGY AND DISCUSSIONS

The following six steps involved for experimental result that fulfill our goal.

### Step 1: Setup of a test computer

For the experimental investigation of the effectiveness of the above tools, we created test data on a Pentium (R) Core (TM) 2 Due CPU, 1.89 GHz, 0.99 of RAM, 160 GB hard drive, write blocker device with Windows XP professional.

### Step 2: Forensically Image the Drive using Encase forensics version 6.2

Here we create an image and analyze the registry of the suspect hard drive using the Encase program. The Encase forensics is a simple but concise tool. It saves an image of a hard disk in one file or in segments that may be later on reconstructed. The data is extracted from the original device, taking care that there is no process that writes on to the digital device under investigation. It calculates MD5 hash values and confirms the integrity of the data before closing the files. The result is an Encase raw image that we have to save in system or another fresh hard drive of at least equal capacity. Now the raw image created by Encase program is used for analysis and examination purpose.

### Step 3: Data extraction/carving

Given the working Encase image we will now extract whole imaged data in readable form with the Encase copy folder program. We recovered all files on to system and/or fresh hard drive of the forensic computer. The time frame for the actual data recovery depends on the duration and frequency of usage of the hard drive.

### Step 4: Attribute Distribution

In figure 4 the attributes or the metadata extracted from each file off the hard drive is also listed. For

our experiment, the attribute id is purely a sequential generated number starting with the number one, associated with each instance of data. The partition attribute details which partition of the hard drive the file belongs to; typically there are 1 primary and two secondary partitions and its data type is numeric. The file size details the size of the file as a numeric data type. The *mtime*, *ctime*, *atime*, and *dtime* attributes represent the modified time, created time, access time, and deleted time of the file in year, month, day, hour, minutes, and second's format.

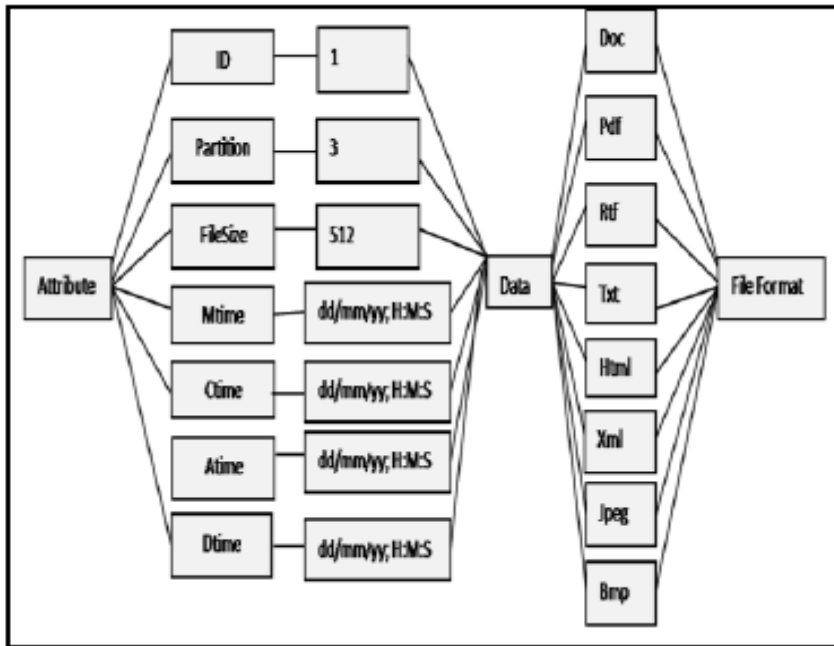


Fig. 3 Attribute distribution format

**Step 5: Data transfiguration**

Use any spreadsheet which is available on your laptop/personal computer. The role played by the spreadsheet can also be achieved by running the data conversion process at the database level too. The major data transformations are conversion of the data into any standard format (comma separated value used here), generation of the parent directories and extraction of the file extensions [24].

**Step 6: Data Analysis**

With the creation of our attribute distribution file, we now have the file metadata that has been extracted from our test drive. As is common in data mining, before running tests on data instances, it was necessary to clean and prepare our data for use into the WEKA workbench. We can also test our metadata statistically using Bartlett's test of sphericity and Kaiser-Meyer-Olkin (KMO) measure of sampling adequacy. After the data preparation was done, WEKA now could be used to run its suite of algorithms on the test data. The complete dataset, consisting of the file tree, the file attributes the timestamps, file size and the deleted flag, is loaded into Weka, open source software, for analyses. On an average, the hard drive, used in this study, recovered around 10000 instances of data (excluding operating system files).

The CSV files we have created with free tool available reflect all information that we mentioned in the attribute distribution format (Figure 3).

The occurrence of all different file formats (including deleted) are listed below:

- (i) The occurrences of all Document files are more than 275.
- (ii) The occurrences of all PDF files are more than 1200.
- (iii) The occurrences of all TXT and RTF files are more than 3000.
- (iv) The occurrences of all other files such as BMP, JPEG, XML, HTML etc are more than 3000.

After categorized all data files, the occurrence of document, pdf and text files are more than any other file format.

File typ <sup>e</sup>	%	File Type	%	File Type	%
DOC	7%	RTF	5%	JPEG	5%
PDF	20%	XML	2%	BMP	7%
XT	25%	HTML	4%	Other	20%

The decision tree output of the algorithm C4.5, suggests the usage pattern of the hard drive. It is evident from Table I as the disk used to store a lot of text and pdf files.

#### IV. CONCLUSIONS

In conclusion, based upon the experimental results, the distribution of different types of files has been demonstrated. It is possible to identify a specific user group hard drive with the help of attribute classification of the data retrieve from the digital evidence. The occurrence of text files in the hard drive are more than other files, and the metadata of the document/pdf files are reflecting that the data contains in the hard drive are preferably belongs to any academic person such as research scholar and the contents of the files can be strongly used to identify the behavior and the area of the person which he/she belongs. Within the experiment, ownership was defined by parsing out user profiles from the windows file system directory structure.

Given a different heuristic, it would be interesting to apply this technique in future to other file systems other than that of Windows and to compare results. This research utilized a basic set of metadata from the files found on each hard drive; for example, file-size, partition, and file format.

## REFERENCES

- [1] John Galloway, Simeon J. Simoff, "Network data mining: methods and techniques for discovering deep linkage between attributes", In APCCM '06: Proceedings of the 3rd Asia- Pacific conference on Conceptual modelling, pages 21–32. Australian Computer Society, Inc., Darlinghurst, Australia, Australia, 2006. ISBN 1-920-68235-X.
- [2] Guidance Software Inc. Encase Forensics. <http://www.guidancesoftware.com>
- [3] Ms. Smita M. Nirkhii, Dr.R.V.Dharaskar, Director, Dr.V.M.Thakre, Data Mining: A Prospective Approach For Digital Forensics. International Journal of Data Mining & Knowledge Management Process (IJDKP) Pg 44, Nov. 2012
- [4] Padhraic Smyth, David Hand, Mannila Heikki 2001. Principles of Data Mining. The MIT Press, PG
- [5] Simson L. Garfinkel(2006),"Forensic feature extraction and cross-drive analysis", Digital Identification forensics", SIGMOD Rec., 30(4):55–64, ISSN 0163-5808.
- [6] Andrew Clark Bradley Schatz, George Mohay. A correlation method for establishing Provenance of timestamps in digital evidence", *6th Annual Digital Forensic Research Workshop, In Digital Investigation, volume 3, supplement 1, pages 98–107. 2006*
- [7] Brian D. Carrier, Eugene H. Spafford. Automated Digital Evidence Target Definition Using Outlier Analysis and Existing Evidence, *Digital Forensic Research Workshop (DFRWS). 2005.*
- [8] Hsinchun Chen, Wingyan Chung, Yi Qin, Michael Chau, Jennifer Jie Xu, Gang Wang, Rong Zheng, and Homa Atabakhsh," Crime data mining: an overview and case studies", *Proceedings of the 2003 annual national conference on Digital government research, pages 1–5. Digital Government Research Center. 2003.*
- [9] Carney, M. and Rogers, M. The Trojan Made Me Do It: A First Step in Statistical Based Computer Forensics Event Reconstruction. *International Journal of Digital Evidence, 2(4). 1-11. 2004.*
- [10] Fayyad, U., G. Piatetsky-Shapiro and P. Smyth. *International Journal of Computer Applications 2012 by IJCA Journal Pg 1*
- [11] Chen, Y., J.R. Miller, J.A. Francis, G.L. Russell, and F. Aires. Observed and modeled relationships among Arctic climate variables. *J. Geophys. Volume. 108. 2003.*
- [12] Usama M. Fayyad, Gregory Piatetsky-Shapiro, Ramasamy Uthurusamy. (2003),"Summary from the kdd-03 panel:



- data mining: the next 10 years. *SIGKDD Explor. Newsl.*, 5(2): 191–196,. ISSN 1931-0145. Pg. 27. 2003.
- [13] Data Mining Concepts and Techniques, 2ed by Jiawei Han, Kamber M Morgan. Kaufmann Publishers. 2005.
- [14] Jan Guynes, Clark Nicole, Lang Beebe (2007),” Digital forensics text string searching: Improving
- [15] O. de Vel, A. Anderson, M. Corney, and G. Mohay (2001),”Mining e-mail content for author Information retrieval effectiveness by thematically clustering search results”, In 6th Annual Digital Investigation, 3(Supplement-1):71–81 ISBN 1-59593-135-X.
- [16] George Forman, Kave Eshghi, and Stephane Chiochetti” Finding similar files in large. 2005
- [17] Brown, Ross A., Pham, Binh L., & De Vel, Olivier Y. A Grammar for the Specification of Forensic Image Mining Searches. In Lovell, Brian, Campbell, Duncan, & Fookes, Clinton (Eds.) *Eighth Australian and New Zealand Intelligent Information Systems Conference, December, Sydney, Australia. 2003.*
- [18] DFRWS,” A road map for digital forensic research”, DTR - T001-01 FINAL - DFRWS Technical Report, 1(1), August 2001. <http://dfrws.org/2001/dfrws-rm-final>. PDF.
- [19] Brown, Ross A. & Pham, Binh L. Image Mining and Retrieval Using Hierarchical Support Vector Machines. In Chen, Yi-Ping (Ed.) *11th International Conference on Multi-Media Modeling, Jan, Melbourne, Australia 15505502, IEEE. 2005.*
- [20] Joachims T. 2002. Optimizing search engines using click through data. In Proceedings of the ACM Conference on Knowledge Discovery and Data Mining (KDD).
- [21] F. Pernkopf. Detection of Surface Defects on Raw Steel Blocks Using Bayesian Network Classifiers. *Pattern Analysis and Applications, Vol. 7, No. 3, 333-34. 2004.*
- [22] Agrawal, R., Imielinski, T. & Swami A. Mining association rules between sets of items in large databases. *Proceedings of the ACM SIGMOD International Conference on Management of Data, 207 - 216. 1993*
- [23] Veena H Bhat, Member, IAENG, Prasanth G Rao, Abhilash V. R., P. Deepa Shenoy, Venugopal K. R. and L. M. Patnaik. A Data Mining Approach for Data Generation and Analysis for Digital Forensic Application. *IACSIT, Vol.2, No.3, ISSN: 1793-8236. 2010.*
- [24] Conference on Knowledge discovery in data mining, pages 394–400, ACM, New York, NY, USA, Document repositories.”, In *KDD '05: Proceeding of the eleventh ACM SIGKDD international Forensic Research Workshop, volume 4, pp 49–54.*