

Study of text classifier based data mining methods for intrusion detection domain

¹RN Phursule, ²Dr. YP Singh

¹ PhD Research Student, Department Kalinga University, Computer Science and Engineering, Raipur, Chhattisgarh, India

² Research Guide, Department of Computer Science and Engineering, Kalinga University, Raipur, Chhattisgarh, India

Abstract

The intrusion detection do an important part in cyber security. Method of extricating invisible, formerly undisclosed as well as functional data from vast databases is mention to Data Mining. Hence data mining methods assist to recognize models in data set as well as utilize these models for identify destiny obrusion. Data Mining based Intrusion Detection System is combined with Multi-Agent System to improve the performance of the IDS. In the current era, there is ample knowledge in using Internet in social networks (such as instant messaging, video conferencing, etc.), the field of healthcare, various areas related to electronic commerce, banking, and services several other fields. As computer systems based on the network plays an ever more important in current community when they possess overtake aim of our opponent with gangsters. So that, we have to search better option for secure our techniques. Safety of computer system is undermined throughout an intrusion exists. Intrusion can be stated as "a set of actions which target to undermine morality, familiarity else accessibility of facility," e.g., unlawfully acquire online user advantages for strike as well as proceed of technique (i.e., Denial of Service) and so on.

Keywords: cyber security, data mining, IDS, multi-agent

1. Introduction

Data mining is procedure of extricating the data or comprehension regularly as well as smartly from a large quantity of information. Here in procedure of data mining, tactful data can be revealed through undermining the independents principled to isolation. Expanding requirement of Isolation conservation in data mining offers me way to investigate regarding isolation conserving data mining.

Supposing instant promotion in novelty, e.g. web, data accumulating, data making systems, we should offer cautious assumption in the direction of safety protecting data mining. To secure extensive substructure we not only require imagining regarding the garnishing of data in addition of data initiation. With the headway and "extension of data mining" there is a vast degree and need of a region which can fill the need of different areas. Combination of systems from data mining, dialect, data process recovery and visual comprehension made an interdisciplinary field called content mining. Content data mining, alluded to as content mining is a procedure of extricating the data from an unstructured content. With a specific end goal to get high content data, a procedure of example division and patterns is finished. For a proficient content mining framework, the unstructured content is parsed and connected or evacuated some level of etymological component, subsequently making it organized content. A standard content mining methodology will include order of content, content bunching, and extraction of ideas, granular scientific categorizations generation, notion examination, record outline and displaying". As a rule, arrangement is the activity of doling out an article to a classification as per the attributes of the item. In data mining, characterization alludes to the undertaking of breaking down an arrangement of pre-ordered information articles to take in a model (or a capacity) that can be utilized to group a

concealed information object into one of the few predefined classes. An information object alluded to for instance, is portrayed by an arrangement of properties or variables. One of the properties portrays the class that a case has a place with and is in this manner called the class characteristic or class variable. Different qualities are regularly called free or indicator properties (or variables). The arrangement of illustrations used to take in the order model is known as the preparation dataset. Assignments identified with characterization incorporate relapse, which assembles a model from preparing information to foresee numerical values, and bunching, which bunches case to frame classifications. Characterization has a place with the class of directed taking in recognized from unsupervised learning. In managed taking in, the preparation information comprises of sets of information, and craved yields, while in unsupervised realizing there is no priori yield.

An intrusion detection technique supervisor web trading for apprehensive actions as well as alerts the system or network administrator in order to take evasive action. It has a very important position in the network information security and it is considered as the second security gate after firewall. In recent years, intrusion detection method and key technology has become one of research focus in network security field.

Intrusion detection is categorization work, as well as that contains constructing a divinatory technique that can recognize strike occurrences. On the one hand, there are various attributes or characteristics that may consist of false connection; categorization of inconsistent intrusion detection technique is complicated task. Furthermore, various attributes may be inapplicable else unessential. For this cause, feature selection methods techniques can be utilized to obtain void of inapplicable as well as unessential attributes in the absence of decreasing performance. On the other hand, IDS usually runs

day by day in real world. And the instance in IDS datasets is very huge and they take time to do classification or clustering. It will take several days to get classification results from a dataset which has over 1 million instances. And if a dataset has a large number of instances and features, it will take large memory and computation resources to run. Thus, feature selection is very necessary to IDS datasets since they usually involve a huge count of occasion as well as characteristics. Formation of IDS construction is shown in figure 1. IDS is composed by four parts as follows

1. Monitoring object. It is monitored and it can be a host or a network.
2. Data collection and storage. This part collects all data from every event, and converts the data to a proper format to store.
3. Data analysis and management. This is a core part in IDS. It searches suspected actions and generates a signal when it detects an attack. Then, IDS deals with the attack or send a signal to network administrator to handle.
4. Signal. It can be seen as an output of IDS. The output is an automatic response or an alarm to network administrator.

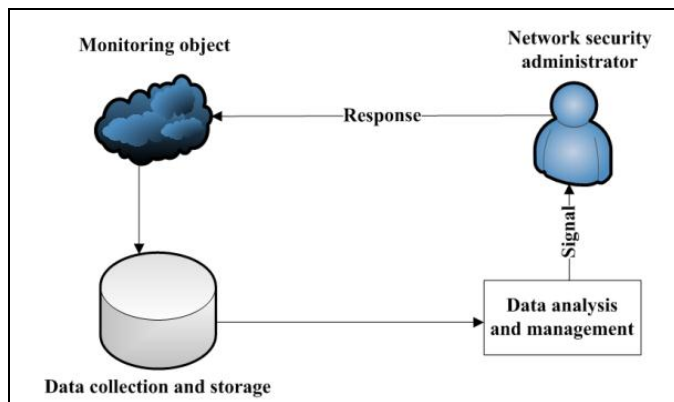


Fig 1: IDS structure

Important quantity of investigate has been coordinated to establish insightful intrusion detection methods, that assist accomplish well network security. Appropriated expanding-depend upon C5 decision trees as well as Kernel Miner are two of previous strive for construct intrusion detection strategies. Systems recommended to prosperously being claimed machine learning systems, including SVM (Support Vector Machine), to categorize web transport models which are not suite usual network traffic. Both techniques were providing in the company of 5 discrete classifiers for recognize usual congestion as well as 4 distinct kinds of strike (i.e., DoS, probing, U2R as well as R2L). Unessential as well as inapplicable feature in information have generated a long-term issue in network traffic classification. These features not just fall off procedure of categorization however also restrain classifier from preparing correct resolutions, specifically when surviving accompanied by big data. Despite enlarging consciousness of network security, the occurred results persist helpless of entire preserving internet applications as well as computer networks across prospect from ever-advancing cyber-attack systems e.g. DoS attack as well as computer malware.

2. Study on Methods

Yang, Z., Zhong, S., & Wright, R., N (2015): In [3], Assurance of security has turned into a vital issue in information mining. Specifically, people have turned out to be progressively unwilling to share their information, much of the time bringing about people either declining to share their information or giving inaccurate information. Thusly, such issues in information accumulation can influence the accomplishment of information mining, which depends on adequate measures of precise information with a specific end goal to deliver significant outcomes. Arbitrary irritation and randomized reaction systems can give some level of security in information accumulation; however they have a related cost in precision. Cryptographic security safeguarding information mining strategies give great protection and exactness properties. In any case, keeping in mind the end goal to be productive, those arrangements must be customized to particular mining errands, along these lines misplacing universality. In paper, we suggest constructive cryptographic techniques to online data aggregation in that data from the comprehensive count of communicator is collected namelessly, in the absence of encouragement of a believed stranger. i.e. our responses sanction the exploiter to collect initial data from each responsive, excluding such that digger can't connect a responsive data to responsive. Approving stage of such a response is, to the point that, since it doesn't change the genuine information, its prosperity does not rely upon the hidden data mining issue. We give verifications of the rightness and protection of our answer, and in addition exploratory information that shows its effectiveness. We likewise stretch out our answer for endure certain sorts of malignant conduct of the members.

Vaidya, J., & Clifton, C. (2014): In [4], Data mining is beneath strike from isolation supporter due to a misinterpretation regarding what it surprisingly is as well as a logical affect regarding how it is commonly over. That object presents how technique from safety society can alter data mining for better, supplying entire its profits if yet conserving isolation.

Aggarwal, C., Pei, J., & Zhang, B. (2010): In [5], Privacy preserving data processing has being an significant topic recently because of progress in hardware techniques that have conduct to extensive growth of statistic as well as defensive information. A basic option to conserve isolation is to easily conceal the data in some of defensive area sharpest through a individual. Although, the system is much from sufficient in its capability to conserve antipathetic data mining. Actual information registers are not accidently separated. As conclusion, some areas in register may be associated on the company of each other. While association is satisfactorily huge, it may be feasible to an competitor to estimate some of emotional areas utilizing another scopes. In paper, creator works the issue of isolation conservation across adversative data mining, that is to conceal a nominal set of entrances hence isolation of emotional scopes are sufficient conserved. In another term, still via data mining, an opponent still cannot correctly recuperate the concealed data entrances. A vast interpretation work is gathered on artificial as well as actual information sets to analyze usefulness of resembles.

Xiong, L., Chitti, S., & Liu, L. (2007): In [8], Advances in

distributed service-oriented computing and Internet technology have formed a strong technology push for outsourcing and information sharing. There is an increasing need for organizations to share their data across organization boundaries both within the country and with countries that may have lesser privacy and security standards. Ideally, Authors wish to share certain statistical data and extract the knowledge from the private databases without revealing any additional information of each individual database apart from the aggregate result that is permitted. In this article, Authors describe two scenarios for outsourcing data aggregation services and present a set of decentralized peer-to-peer protocols for supporting data sharing across multiple private databases while minimizing the data disclosure among individual parties. Authors basic protocols include a set of novel probabilistic computation mechanisms for important primitive data aggregation operations across multiple private databases such as max, min, and top k selection. Authors provide an analytical study of our basic protocols in terms of precision, efficiency, and privacy characteristics. Authors advanced protocols implement an efficient algorithm for performing KNN classification across multiple private databases. Authors provide a set of experiments to evaluate the proposed protocols in terms of their correctness, efficiency, and privacy characteristics.

Zhu, M., & Liu, L. (2004): Randomization is an financial as well as proficient resemble to privacy preserving data mining (PPDM) in [9]. In subsequent of guarantee the accomplishment of data mining as well as preservation of separate isolation, ideal organization program require to be hired. Paper reveals the fabrication of ideal organization programs for privacy preserving future approximation. Writer suggests a normal substructure for organization utilizing integrated replicas. Effect of organization on data mining is expressed through interpretation humiliation as well as interchange able at a reduction, if privacy as well as privacy reduction is expressed through intermission-based emphasis. 2 distinct kinds of issues are stated to recognize ideal

organization for PPDM. Explanatory examples as well as reflection solutions are announced.

Fayyad, U., Piatetsky-Shapiro, G., with Smyth P (2006): Data mining as well as knowledge discovery in databases has been stimulating a remarkable quantity of investigation, trade, as well as media attentiveness of late in [12]. What is total agitation regarding? This object supplies an review of this appearing area, simplifying how data mining as well as knowledge discovery in databases are connected both to one another as well as to connected areas, e.g. machine learning, statistics, as well as databases. The object states specific real-world applications, particular data-mining methods, provocations included in real-world applications of knowledge discovery, as well as recent with destiny investigation towards in area.

Lin, T. Y. (2004): In [15], One of the organic outcome of extensive behavior of Internet is a requirement for larger safety. End-technique association across various webs presents an compulsory exposure. This declaration is depending on safety defect which occur in network applications as well as network facilities. Increasing in percentage to always enlarging utilization of internet is evolution of safety instruments. Specifically, variation/intrusion detection techniques are kinds of safety instruments become established. One region of network safety which is yet undersized is utilization of brilliant instruments.

This work develops a soft computing technique to identify irregular nature making contributions that include being a novel technique for anomaly detection, having real-time applicability, resistance to scalability, minimal computational complexity, adaptability, and robustness

3. Comparative analysis

Techniques calculated above are resembled in concept of profits, drawbacks, systems as well as correction presentation. Table 1 displays relative work between these systems as well as fig 2 is displaying accuracy differentiation.

Table 1: Relative Study of Picture Fabrication Detection Systems

Paper Title	Key Techniques and Methods	Advantages	Disadvantages
Mining prodigious insufficient data sets through ideal reformation	Key Intuitions, Covariance Matrix	effective conceptual representations	Very amenable for reconstruction.
Anonymity-preserving data collection	Random perturbation, Cryptographic privacy-preserving data mining	resolution to allow definite types of malicious nature of the contributors	Face worker to gather the real data from every respondent
Privacy-preserving data mining: Why, how, as well as when.	privacy, data mining	Method to protect individual info, release aggregate info.	Exact aggregate info may leak individual info.
On Privacy Preservation across Antipathetic Data Mining	Privacy preservation, data mining, association rules.	does not change very much for distinct values of user-specified maintain	efficient from a computational perspective as well as need a few seconds across various experimental Settings.
Conserving information privacy in outsourcing data aggregation services.	Privacy, Confidentiality, Classification, Outsourcing	Extricate comprehension from personal databases	Does not consider data sharing in global scale.
Optimal randomization for privacy preserving data mining	Consecutive model mining; Anonymity; Randomization; Secure multiparty computation;	Computed as well as differentiated in concepts of precision, recall as well as accuracy ratio.	It is complicated technique as well as processing duration is not computed.
Behavior Comprehension Space-Based integration for Copy-Move Fabrication Detection	Multi scale BKS, Random Forest, SVM.	Almost affectionate region in data mining as well as for conserving privacy in extricating progression	Ideal organization for PPDM

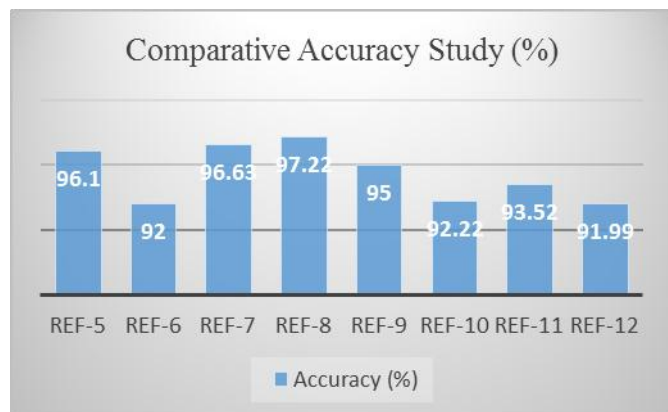


Fig 2: Accuracy Analysis of Studied Methods

4. Conclusion and Future Work

Data mining techniques are widely used because of their capability to drastically improve the performance and usability of intrusion detection systems. Different data mining techniques like classification, clustering and association rule mining are very helpful in analyzing the network data. Since large amount of network traffic needs to be collected for intrusion detection, clustering is more suitable than classification in the domain of intrusion detection as it does not require labeled data set thereby reducing manual efforts. Data mining techniques can detect known as well as unknown attacks. Data mining technology helps to understand normal behavior inside the data and use this knowledge for detecting unknown intrusions. Text classification is the primary requirement of text retrieval systems, which retrieve texts in response to a user query, and text understanding systems, which transform text in some way such as producing summaries, answering questions or extracting data.

5. References

- Han Jiawei, Kamber M, Data Mining. Concepts and Techniques, Beijing: China Machine Press. 2006, 1-40.
- Aggarwal C, Parthasarathy S. Mining massively incomplete data sets by conceptual reconstruction. Proceedings of the 7th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining. 2010, 227-232.
- Yang Z, Zhong S, Wright RN. Anonymity-preserving data collection. Study presented at the Proceeding of the eleventh ACM SIGKDD international conference on Knowledge discovery in data mining, Chicago, Illinois, USA, 2015.
- Vaidya J, Clifton C. Privacy-preserving data mining: Why, how, and when. IEEE Security & Privacy. 2014; 2(6):19-27.
- Aggarwal C, Pei J, Zhang B. On Privacy Preservation against Adversarial Data Mining. Proceedings of the 12th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining. 2010, 510-516.
- Vaidya J, Clifton C, Zhu M. Privacy Preserving Data Mining. New York: Springer, 2012.
- Xiong L, Chitti S, Liu L. Preserving data privacy in outsourcing data aggregation services. ACM Transactions on Internet Technology TOIT. 2007; 7(3).
- Zhu M, Liu L. Optimal randomization for privacy preserving data mining. Proceedings of the 2004 ACM SIGKDD International Conference on Knowledge Discovery and Data Mining. 2004, 761-766.
- Introduction to Data Mining and Knowledge Discovery, Third Edition ISBN: 1-892095-02- 5, Two Crows Corporation, 10500 Falls Road, Potomac, MD 20854 (U.S.A.), 2009.
- Dunham MH, Sridhar S. Data Mining: Introductory and Advanced Topics, Pearson Education, New Delhi, ISBN: 81-7758-785-4, 1st Edition, 2010
- Fayyad U, Piatetsky-Shapiro G, Smyth P. From Data Mining to Knowledge Discovery in Databases, AI Magazine, American Association for Artificial Intelligence, 2006.
- Morgenstern M. Security and Inference in Multilevel Database and Knowledge Base Systems,” Proceedings of the ACM SIGMOD Conference, San Francisco, CA, 2007.
- Database Security IX Status and Prospects Edited by D. L. Spooner, S. A. Demurjian and J. E. Dobson ISBN 0 412 729202. 1996, 391-399.
- Lin TY. Anomaly Detection -- A Soft Computing Approach”, Proceedings in the ACM SIGSAC New Security Paradigm Workshop, This study reappeared in the Proceedings of National Computer Security Center Conference under the title Fuzzy Patterns in data, 2004, 1994, 44-53.
- Dileep Kumar Singh, Vishnu Swaroop. Review and Analysis of Data Security in Data Mining, IRACST - International Journal of Computer Science and Information Technology & Security (IJCSITS), ISSN: 2249-9555. 2012; 2(4).