

AN ABETTOR INVASION ON KERNEL-BASED DATA MINING SYSTEMS

Subhami mohan¹
M.Tech Student
Department of CSE
Lourdes Matha College of
Science and Technology

Renetha J B²
Associate Professor
Department of CSE
Lourdes Matha College of
Science and Technology

Abstract- Kernel is the central core of the computer system that has control over every event that occurs in the system. Data Mining is the process of discovering patterns in large datasets for analyzing and summarizing data into useful information. Threats and privacy issues are of great concern in the privacy preserving work of kernel-based data mining systems. One of the most significant attacks associated with kernel-based data mining are the insider attacks. There is a need of an automated system to assess the risk of data mining applications. This paper, proposes various mechanisms adopted to counter insider attacks. Finally there is a description of the drawbacks that lead to the failure of these adopted methods.

Index Terms— kernel, data mining, insider attack, privacy preserving

I. INTRODUCTION

Data Mining is the process of discovering patterns in large datasets for analysing and summarizing data into useful information. Data mining applications store huge amount of information and hence privacy preserving requires great attention in this area. Insider attack is the common type of data breaching problems that have risen significantly in countries like the United States. These attacks arise from the staff or person working within an organization or company. Hence these attacks are now considered to be the fastest growing type of attacks. Insiders or abettors can provide sensitive information to the outsiders as there are no technical barriers. Several privacy preserving schemes are currently proposed to counter insider attacks. But no prior work can be considered to be robust.

Three types of privacy preserving mechanisms are currently available. One of the major privacy preserving mechanism is the use of Support Vector Machines. Support Vector Machines is a data mining technique used with the kernel trick to map data into a higher dimensional space. Another type of privacy preserving includes the m-privacy technique. In m-privacy technique, an anonymized view of integrated data is published. The third type of mechanism is a declarative approach, called DASAI, which is a declarative approach to identify and prevent insider attacks.

II. EXISTING SYSTEMS

A. Support Vector Machines

Support Vector Machine (SVM) is generally utilized as a part of classification and regression analysis [2]. SVM, a standout amongst the most prevalent fields of extraordinary research was initially recommended by

Vapnik in 1960 for classification. The two key components in the usage of SVM are Mathematical programming and kernel functions. SVM changes the original training data into higher dimensions and searches for a linear optimal separating hyper plane among this dimension. For this purpose SVM additionally utilizes the kernel trick which is a nonlinear mapping technique.

The distributed data can be formulated to frame a global SVM classification model, in which the accommodated multiple parties and their representative part is kept encoded or disclosed to each other. The disclosure of data is prevented and made purely authorized by the so using kernel matrix. The kernel matrix is therefore considered as the central structure in a SVM as it acts as an intermediate profile that can generate the global model without disclosing any local data information. The support vectors in the learned classifier arise a privacy threat of SVM based classifications.

The disclosure of sensitive information about the original owner of the training data may happen while the releasing of SVM classifier as the support vectors are intact instances taken from training data. A hyperbolic tangent kernel based SVM classifier is developed for privacy-preserving aspects. It is an approximation of the original one accounting with the kernel function in the classifier. The increasing number of support vectors denotes the level of accurate privacy preserving for it is another term point to the degree of approximation.

A privacy-preserving SVM classifier is designed with Gaussian kernel function to compensate the violation of individual's privacy through the release of support vectors. Transforming the original decision function is realised by Privacy-preserving. It is determined by support vectors through an infinite series of linear combinations of monomial feature mapped support vectors. This linear combination will destroy the sensitive content of support

vectors, while the decision function can precisely approximate the original one.

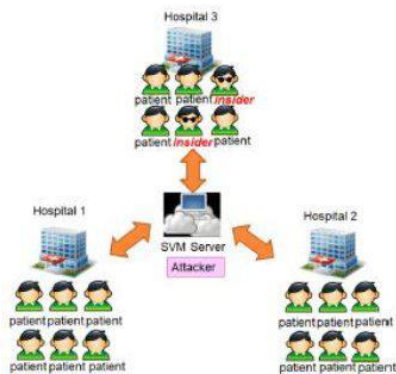


Figure1:Insiders in the hospitals help the outsider attacker to launch attacks.

B. m-privacy

As far as the collaborative data publishing problem is considered, it anonymises the horizontally partitioned data at multiple data providers. The data providers who use their own data records which is a subset of the overall data introduce a new type of attack called m-adversary [3]. These are actually in addition to the external background knowledge to infer the data records contributed by other data providers. The disclosure of sensitive information by an m-adversary is prevented by the guaranteeing m-privacy.

The presenting adaptive ordering techniques and heuristic algorithms exploiting equivalence group monotonicity of privacy constraints check m-privacy. To ensure high utility and m-privacy of anonymized data a provider-aware anonymization algorithm is introduced with adaptive m-privacy checking strategies. As m-privacy ensures efficiency than existing algorithms almost experiments confirmed that this method covers utility standards. A proper privacy fitness score for different privacy constraints is one of the remaining research questions. When data are distributed in a vertical or ad-hoc fashion m-privacy constraints do not support to address and model the data knowledge of its providers.

C. DASAI

DASAI is a logic rule-based static analysis approach for determining the possibility of an attack can take place on a process. If an attack is suppose to be possible, DASAI can perform on the process and ascertain the rogue insiders involved and determine the ways of this attack [4]. When a process is vulnerable to an attack dataflow based process and attack models are considered, and a holistic perspective is used that looks at steps, data, annotations on data and controlling agents to determine. A graph matching-based search problem is made to reduce the effort identifying the attack determination problems.

A declarative programming paradigm is used to enumerate automatically the possible ways of matching an attack graph against a process graph through a concept of a valid mapping encoded as logic rules. Every mapping mode gives rise to a possible avenue of attack. DASAI is a logic rule-based approach and very amenable to addition of new constraints to change the definition of an attack mapping. DASAI automatically and opportunistically searches and exploits improvement opportunities in the process, once attack possibilities are determined. Its searching starting from the mostly attacked steps to the lesser attacked ones, to make the process robust against the attack in all possible ways.

Usually generation of automatic attack models is not possible on DASAI from a given process model. as DASAI do not conform to a desired property of a process so current literature describes the use of model checking to automatically generate attack models. The attack model structure is converted into a graphical data-flow based attack once the model checker identifies an attack model that is successful against a process, There DASAI can be used to check other processes and possible ways of successful attack.

III.PROPOSED METHODOLOGY

Each organization maintains a database which stores the corresponding dataset related to the organization. This dataset would be in non-understandable format. So the dataset is converted into an understandable format and subjected through a series of steps to preserve privacy from abettors.

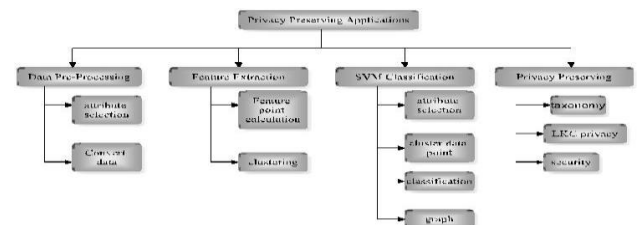


Figure 2: Overview of proposed System

A. DATA PRE-PROCESSING

In this process, the dataset is collected and manipulated. The validity of the datum is checked and verified. As the dataset is in some specific format, the data must be converted accordingly. Therefore, attributes are selected from the obtained data to make it in an understandable form.

B. FEATURE EXTRACTION AND CLUSTERING

REFERENCES

Feature extraction is performed on the available dataset. It builds derived values (features) from the dataset using some functional mapping keeping as much information as possible. This method intended to be informative and non-redundant, facilitating the subsequent learning and generalization steps, and in some cases leading to better human interpretations. Feature extraction is related to dimensionality reduction.

Clustering is the process of grouping of a particular set of data based on their characteristics into clusters [6]. These clusters are then aggregated according to their similarities. Also, Clustering partitions the data most suitably for the desired information analysis. The result of clustering will be displayed in a table format.

C. CLASSIFICATION

The resulting clusters are subjected to classification. Here based on a particular cluster, the data is classified and plotted as a graph. A classifier is used for this purpose which may either be an algorithm that implements classification or a mathematical function that maps input data to a particular category. Each graph shows the dependence of data between clusters within each attribute which is selected.

D. PRIVACY PRESERVING

The current privacy preserving techniques are classified based on a particular taxonomy [5]. Taxonomy relates to a second stage of classification based on a tree structure. Privacy is ensured on the basis of an LKC algorithm. Finally, the entire result is subjected to a security check which interprets the chances of data leakage.

IV. EXPERIMENTS AND DISCUSSION

In this experiment, the proposed insider attack scheme and the methods for preserving privacy are examined. Also, the selected datasets include some medical and bank credit datasets as they usually have stronger privacy concerns.

V. CONCLUSION

This paper presented an efficient method to identify data leakage through abettor invasions on kernel based data mining systems. This paper also explains the drawbacks of existing systems and proposes a method to encounter the attack by estimating the number of insiders. Also, the proposed method does not make any complication regarding the SVM classification.

- [1] Peter Shaojui Wang Feipei Lai (Senior Member, IEEE), Hsu-Chun Hsiao, and Ja-Ling Wu, (Fellow, IEEE), "Insider Collusion Attack on Privacy-Preserving Kernel-Based Data Mining Systems", Digital Object Identifier 10.1109/ACCESS.2016.252561019
- [2] Lei Xu, Chunxiao Jiang, (Member, IEEE), Jian Wang, (Member, IEEE), Jian Yuan, (Member IEEE), and Yong Ren, (Member, IEEE), "Information Security in Big Data: Privacy and Data Mining", IEEE Access, vol. 2, pp. 1149_1176, Oct. 2014.
- [3] A. Sarkar, S. Köhler, S. Riddle, B. Ludaescher, and M Bishop, "Insider attack identification and prevention using a declarative approach," in Proc. IEEE Secur. Privacy Workshops (SPW), May 2014, pp. 265_276.
- [4] S. Goryczka, L. Xiong, and B. C. M. Fung, "m-privacy for collaborative data publishing," IEEE Trans. Knowl. Data Eng., vol. 26, no. 10, pp. 2520_2533, Oct. 2014
- [5] J. Vaidya, H. Yu, and X. Jiang, "Privacy-preserving SVM classification," Knowl. Inf. Syst., vol. 14, no. 2, pp. 161178, Feb. 2008
- [6] D. Vizár and S. Vaudenay, "Cryptanalysis of chosen symmetric homomorphic schemes," Studia Sci. Math. Hungarica, vol. 52, no. 2, pp. 288306, 2015