



An Efficient Privacy Preserving Approach Using Id3 Decision Tree Learning Algorithm

T. Satya Narayana Murthy
Asst. Professor, M.Tech, CSE
Vignan University, Vadlamudi,
Guntur (D) India

P. Jeevitha Lakshmi
Student, M.Tech, CSE
Vignan University, Vadlamudi,
Guntur (D) India

Abstract— *Privacy-preserving is an important issue in the areas of data mining and security. The aim of privacy preserving data mining is to develop algorithms to modify the original dataset so that the privacy of confidential information remains preserved and as such, no confidential information could be revealed as a result of applying data mining tasks. In existing system they introduced a new privacy preserving approach via data set complementation which confirms the utility of training data sets for decision tree learning. This approach converts the original data sets, T_S , into some unreal data sets such that any original data set is not constructable if an unauthorized party were to steal some portion of unrealized datasets. Meanwhile, there remains only a low probability of random matching of any original data set to the stolen data sets, T_L . This work covers the application of new privacy preserving approach with the ID3 decision tree learning algorithm. The problem in existing system is insufficient storage mechanism and this ID3 only can be implemented for discrete-valued attributes. To support continuous-valued attributes c5.0 algorithm is used.*

Index Terms—*Classification, Data mining, Machine learning, Privacy protection and Security.*

I. INTRODUCTION

The problem of privacy-preserving in data mining has become more important in recent year because the ability to store personal data about users is increased, and the increasing knowledge about the data mining algorithms to control this information. There are number of techniques such as randomization and k-anonymity have been suggested in order to perform privacy-preserving data mining. Also this problem has been discussed in many communities such as the statistical disclosure control community, database community and the cryptography community. In some cases, the different communities have explored parallel lines of work which are quite similar. This paper will try to explore different topics from the perspective of different communities, and will try to give a fused idea of the work in different communities. Preserving privacy is more important for machine learning and data mining, but the measures designed to protect private information sometimes result in a degradation and reduced utility of the training samples.

This work introduces an approach that can be applied to decision-tree learning, without concurrent loss of accuracy. It describes a privacy preservation approach for the collected data samples in cases when information of the sample database has been partially lost. This approach converts the original datasets into a group of unreal datasets [1], in which the original data cannot be reconstructed without the entire group of unreal datasets if some portion of the unreal datasets is stolen. This approach does not suitable when sample datasets have low frequency or low variance in the distribution of all samples. However, this problem can be resolved through an alternative implementation of the approach introduced later in this work, by using some extra storage. The key directions in the field of privacy-preserving data mining [4], [5], [15] are as follows:

A. Privacy-Preserving Data Publishing

These techniques tend to study different transformation methods associated with privacy. These approaches include methods such as randomization [17], k-anonymity, and l-diversity. A related issue is how the perturbed data can be used along with classical data mining methods such as association rule mining [3], [13]. Other related problems include that of determining privacy preserving methods to keep the underlying data useful or the problem of studying the various privacy definitions, and how they compare in terms of effectiveness in different states.

B. Modifying the results of Data Mining Applications to preserve privacy

In many cases, the results of data mining applications such as association rule [13] or classification rule mining can compromise the privacy of the data. This has generate a field of privacy in which the results of data mining algorithms such as association rule mining are modified in order to preserve the privacy of the data. A classic example of such techniques are association rule hiding methods [3], in which some of the association rules are suppressed in order to preserve privacy. Likewise many techniques are available to modify the results of the data mining applications.

C. Cryptographic techniques for Distributed Privacy

In many cases, the data may be distributed across many sites, and the owners of the data across these different sites may wish to compute a common function. In those cases, a variety of cryptographic protocols [2], [14] may be used to communicate among various sites, so that secure function computation is possible without revealing the sensitive information.

II. RESEARCH BACKGROUND AND OBJECTIVE

Even when databases of samples with sensitive information are protected securely, partial information of the databases can be lost through procedural mistakes or privacy attacks which can be from anywhere within a network. This work focuses on analyzing privacy preservation following the loss of some training datasets from the whole sample database used for decision-tree learning [11], [20]. On this basis, we make the following assumptions for the scope of this work: first, as is the norm in data collection processes, a large number of sample datasets have been collected to achieve significant data mining results that cover the whole research target. Second, the number of datasets lost constitutes a small portion of the entire sample database. Third, for decision-tree data mining, no attribute is designed for distinctive values, because such values negatively affect decision classification.

The objective of this work is to introduce a new privacy preserving approach to the protection of sample datasets that are utilized for decision-tree data mining. Privacy preservation is applied directly to the samples in storage, so that privacy can be safeguarded even if the data storage were to be threatened by unauthorized parties. Although effective against privacy attacks by any unauthorized party, this approach does not affect the accuracy of data mining results. Moreover, this technique can be applied at any time during the data collection process, so that the protection of privacy can be in effect as early as the first sample is collected.

Privacy preservation converts dataset containing private information into altered or sanitized versions in which private information is hidden from unauthorized parties. Privacy preserving data mining refers to the area of data mining that seeks to safeguard sensitive information from unsanctioned or unsolicited disclosure.

Privacy Preservation Data Mining [1] [2] was introduced to preserve the privacy during mining process to enable conventional data mining technique. Many privacy preservation approaches were developed to protect private information of sample dataset. On the other hand privacy preserving process which hides information may reduce utility of these sanitized dataset. When there utility decreases to a certain level the downgraded information prevents accurate analysis. With the result that primary objective of data mining is compromised. Even when databases of samples with sensitive information are protected securely partial information of database can be lost through procedural mistakes or privacy attacks from anywhere within the network.

III. EXISTING SYSTEM

Iterative Dichotomiser 3 (ID3) selects the test attribute based on the information gain [1], provided by the test outcome. Information gain measures the change of uncertainty level after a classification from an attribute. Fundamentally, this measurement is rooted in information theory. Privacy preservation in data mining activities is of significant importance for many applications. However, the privacy preserving process sometimes reduces the utility of training datasets, which causes inaccurate data mining results. Privacy preservation approaches focus on different areas of a data mining process, and data mining methods also vary. This thesis focuses on privacy protection of the training samples applied for decision tree data mining.

This work presents a new privacy preserving approach via dataset complementation [1], in which the universal set is generated from the original samples. It removes each sample from a set of perturbing datasets. During the privacy preserving process, this set of perturbed datasets is dynamically modified. As the sanitized version of the original samples, these perturbed datasets [1], [11], are stored to enable a modified decision tree data mining method. This method guarantees to provide the same data mining outcomes like the originals, which is mathematically proved and also by a test using one set of sample datasets in this thesis. From the viewpoint of privacy preservation, the original datasets can only be reconstructed in their entirety if someone has all perturbed datasets, which is not supposed to be the case for an unauthorized party.

IV. PROPOSED SYSTEM

FFDs enable de-generalization of the anonymized data and thus lead to privacy breaches. Based on the impact of FFDs [21] to privacy, we distinguish “safe” FFDs that cannot enable any FFD-based attack from the “unsafe” ones that can. The overall proposed work is shown in Fig. 4. The research work presented here uses the C5.0 Algorithm for data mining.

A. Unrealized training set completion

To unrealized the samples, we initialize both set of input sample dataset(T') and perturbing dataset(T_p) as empty sets, i.e. Unrealized training set(T_u) is called. Universal set is generated by using the single instance of all the possible values of the original data set. Consistent with the procedure described above, universal dataset is added as a parameter of the function because reusing pre-computed universal dataset is more efficient than recalculating universal dataset [1]. The recursive function unrealized training-set takes one dataset in input sample dataset in a recursion without any special requirement; it then updates perturbing dataset [11] and set of output training data sets correspondent with the next recursion.

Therefore, it is obvious that the unrealized training set process can be executed at any point during the sample collection process. The Fig. 1 shows the original table, Fig. 2 shows the perturb table and the Fig. 3 shows the unrealized table.

| workclass | race | sex |
|-----------|-------|--------|
| Private | Black | Female |
| Private | White | Female |
| State gov | White | Male |
| Local gov | Black | Female |
| State gov | white | Male |

Fig 1 – Original Table (T^o)

| workclass | race | sex |
|-----------|-------|--------|
| Private | White | Female |
| Private | Black | Male |
| State gov | Black | Female |
| Local gov | Black | Male |
| State gov | White | Male |
| Local gov | White | Female |

Fig 2 – Perturb Table (T^p)

| workclass | race | sex |
|-----------|-------|--------|
| Private | Black | Male |
| Private | White | Male |
| State gov | White | Male |
| Local gov | White | Female |
| State gov | Black | Male |

Fig 3 –unrealized Table (T^u)

B. Decision tree algorithm C5.0

The research work presented here considers the C5.0 Algorithm for data mining. The enhancement and the optimization of the C4.5 emerge as algorithm C5.0, which exhibits the better performance as compared to the other existing mining algorithms. C5.0 algorithm to build either a decision tree or a rule set. In C5.0 model the sample is split based on the field that provides the maximum information gain.

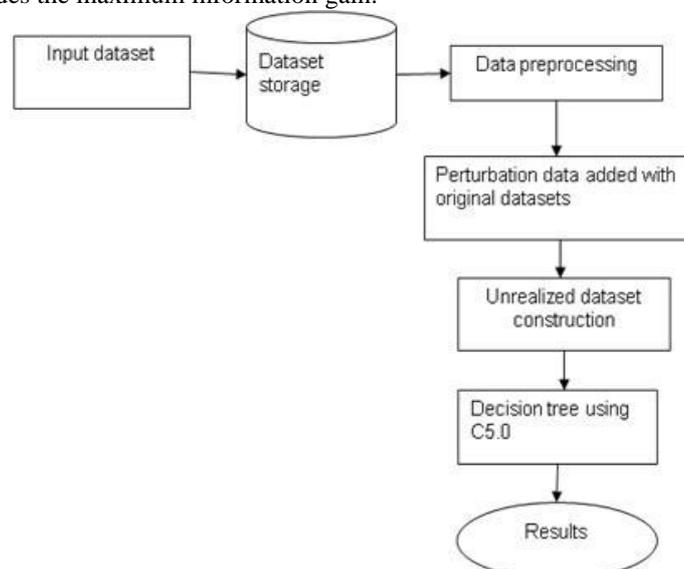


Fig 4- Overall process diagram

Again each subsample defined by the first split is split based on a different field, and this process repeats until the subsamples cannot be split anymore. Finally, the lowest-level splits are evaluated again, and those that do not provide significantly to the value of the model are removed or pruned. C5.0 can produce two varieties of models. The decision tree generated here is a straightforward description of the splits found by the algorithm. Each terminal (or "leaf") node describes a particular subset of the training data and each case in this training data belongs to exactly one terminal node in the tree.

The C5.0 algorithm needs to calculate the entropy and the information gain of the attributes in the table. For calculating the entropy, the equation (1) is used, where p_i is the probability of the class c_i in the table.

$$\text{Entropy}(D) = -\left(\sum_{i=1}^m p_i * \log_2 p_i\right) \quad (1)$$

The Information gain is calculated using the entropy value, the equation (2) is used.

$$\text{Gain}(D) = \text{Entropy}(D) - \sum_{i=1}^n \frac{|D_i|}{|D|} * \text{Entropy}(D_i) \quad (2)$$

The information gain is calculated for all the attributes in the table, to find out the attribute with maximum information. Based on the attribute with maximum gain, the sample has been split.

V. CONCLUSION

This paper covers the new approach for privacy preservation that can be applied to the C5.0 algorithm which support both discrete and continuous value attributes. It optimizes the storage size of the unrealized data. The future work should concerns data with fully functional dependencies [21].

REFERENCES

- [1] Pui K.Fong and JensH.Weber-Jahnke, "Privacy preserving Decision tree Learning Using Unrealized Data sets", IEEE Trans. Knowledge and Data Eng., vol.24 No. 2, Feb 2012.
- [2] S.Ajmani, R.Morris, and B.Liskov, "A Trusted Third-Party Computation Service," Technical Report MIT-LCS-TR-847, MIT 2001.
- [3] S.L.Wang and A.Jafari, "Hiding Sensitive Predictive Association Rules," Proc.IEEE Int'l Conf. Systems, Man and Cybernetics, pp.164-169, 2005.
- [4] R.Agrawal and R.Srikant, "Privacy Preserving Data Mining," Proc. ACM SIGMOD Conf. Management of Data (SIGMOD'00), pp.439-450, May 2000.
- [5] Q.Ma and P.Deng, "Secure Multi-Party Protocols for Privacy Preserving Data Mining," Proc. Third Int'l Conf. Wireless Algorithms, Systems, and Applications (WASA'08), pp.526-537, 2008.
- [6] J.Gitanjali, J.Indumathi, N.C.Iyengar, and N.Sriman, "A Pristine Clean Cabalistic Foruity Strategize Based approach for Incre-mental Data Stream Privacy Preserving Data Mining," Proc. IEEE Second Int'l Advance Computing Conf. (IACC), pp.410-415,2010.
- [7] N.Lomas,"Data on 84,000 United Kingdom Prisoners Lost," Retrieved Sept.12, 2008, http://news.cnet.com/8301-1009_3-10024550-83.html, Aug.2008.
- [8] BBC News Brown Apologizes for Records Loss. Retrieved Sept. 12, 2008, http://news.bbc.co.uk/2/hi/uk_news/politics/7104945.stm, Nov.2007.
- [9] D.Kaplan, Hackers Steal 22,000 Social Security Numbers from Univ. of Missouri Database, Retrieved Sept.2008, <http://www.scmaga-zineus.com/Hackers-steal-22000-Social-Security-numbers-from-Univ.-of-Missouri-database/article/34964/> May2007.
- [10] D.Goodin,"Hackers In filtrate TD Ameritrade client Database, "Retrieved Sept.2008, http://www.channelregister.co.uk/2007/09/15/ameritrade_database_burgled/, Sept.2007.
- [11] Liu, M.Kantarcioglu, and B.Thuraisingham, "Privacy Preserving Decision Tree Mining from Perturbed Data, "Proc.42nd Hawaii Int'l Conf.System Sciences (HICSS'09), 2009.
- [12] Y.Zhu, L.Huang, W.Yang ,D.Li, Y.Luo, and F.Dong, "Three New Approaches to Privacy-Preserving Add to Multiply Protocol and Its Application,"Proc.Second Int'l Workshop Knowledge Discovery and Data Mining, (WKDD'09), pp.554-558,2009.
- [13] J.Vaidya and C.Clifton, "Privacy Preserving Association Rule Mining in Vertically Partitioned Data," Proc Eighth ACM SIGKDD Int'l Conf. Knowledge Discovery and data Mining (KDD'02), pp.23-26, July 2002.
- [14] M.Shaneck and Y.Kim, "Efficient Cryptographic Primitives for Private Data Mining," System Sciences (HICSS), pp.1-9, 2010.
- [15] C.Aggarwal and P.Yu, Privacy-Preserving Data Mining:, Models and Algorithms. Springer, 2008.
- [16] L.Sweeney, "k-Anonymity: A Model for Protecting Privacy,"Int'l J.Uncertainty, Fuzziness and Systems, vol.10, pp.557-570, May 2002.
- [17] J.Dowd, S.Xu, and W.Zhang, "Privacy-preserving Decision Tree Mining Based on Random Substations," Proc. Int'l Conf. Emerging Trends in Information and Comm. Security (ETRICS'06), pp.145-159, 2006.
- [18] Bu, L.Lakshmanan, R.Ng, and G.Ramesh, "Preservation of Patterns and Input-Output Privacy," Proc. IEEE 23rd Int'l Conf Data Eng., pp.696-705, Apr.2007.
- [19] S.Russell and N.Peter, Artificial Intelligence. A Modern Approach2/ E. Prentice-Hall, 2002.
- [20] P.K.Fong, "Privacy Preservation for Training Data Sets in Database: Application to Decision Tree Learning, master's thesis, Dept.of Computer Science, Univ.of Victoria, 2008.
- [21] Hui Wang and Ruilin Lui, "Privacy Preserving Publishing Micro data with Full Functional Dependencies", Data & Knowledge Engineering 70 (2011) 249-268.