



## Privacy Preserving Using Direct and Indirect Discrimination Rule Method

C.V.Nithya

Dept. of CSE

Vivekanandha College of Technology for Women  
Namakkal, India

A.Jeyasree

Dept. of CSE

Vivekanandha College of Technology for Women  
Namakkal, India

---

**Abstract:** *The Privacy preserving Data mining (PPDM) has become an important issue in recent years because of abundance of sensitive information on the internet. Discrimination is a presuppose privileges where provide to the each separate group for the safety of the data which is stored. Discrimination is two type direct and indirect discrimination. Direct discrimination is based on sensitive information. In direct discrimination is based on non sensitive information. In existing system standard algorithm is used. Sometimes the data should be lost. In that data model sensitive information should be free. It does not successfully handle the indirect discrimination problems which are associated to direct discrimination. In the system using the new techniques to prevent the sensitive information. In term of data quality and discrimination detach for both direct and indirect discrimination.*

**Keyword:** *direct discrimination, indirect discrimination, privileges, sensitive data, non sensitive data*

---

### I. INTRODUCTION

In sociology, discrimination is the prejudicial treatment of an individual based on their membership in a certain collection. It involves denying to members of one group opportunities that are available to other groups. There is a catalog of antidiscrimination acts, which are laws designed to avoid discrimination on the basis of a number of elements (e.g., race, religion, gender, nationality, disability, marital status, and age) in various settings (e.g., employment and training, access to public services station, credit and insurance, etc.). For example, the European Union implements the primary of equal treatment between men and women in the access to and supply of goods and services in or in matters of employment and occupation in. Although there are some laws abutting discrimination, all of them is reactive, not proactive. Technology can add proactively to legislation by contributing discrimination discovery and prevention techniques. Services in the information society allow for automatic and routine collection of large amounts of data. Those data are often used to train association/classification rules in view of making automatic decisions, like loan granting/denial, insurance premium computation, personnel selection, etc. At first place, automatic decisions may give a sense of fairness: classification rules do not guide themselves by personal favorite. However, at a closer look, one realizes that classification rules are actually learned by the system (e.g., loan granting) from the training data. If the training data are essential biased for or against a particular community (e.g., foreigners), the learned model may show a discriminatory prejudice behavior. In other words, the system may infer that just being foreign is a legitimate reason for loan denial. Find such potential biases and eliminating them from the training data without harming their decision-making utility is therefore highly necessary. One must data avoid mining from becoming itself a source of discrimination, because of data mining tasks generating discriminatory models from biased data sets as part of the automated decision making. In, it is demonstrated that data mining can be both a source of discrimination and a means for discovering discrimination. Discrimination is classified two types direct or indirect (also called systematic). Direct discrimination consists of rules or procedures that integrally mention minority or disadvantaged groups based on sensitive discriminatory attributes related to group membership. Indirect discrimination consists of law or process that, while not explicitly mentioning discriminatory.

### II. PRIVACY PRESERVING DATA MINING (PPDM)

Data mining is the process of gathering information about the user specific data, also called knowledge discovery, on internet. The problem with data mining output is that it also discloses some information, which is considered to be private and personal. Effortless access to such personal data causes a peril to individual privacy. Official statistics, Health information, and E-commerce are some key concern for privacy. Privacy preserving data mining technique gives novel way to solve this problem. The main purpose of privacy preserving data mining is to design competent frameworks and algorithms that can extract relevant knowledge from a large amount of data without revealing of any sensitive information. It protects sensitive information by providing sanitized database of original database on the internet or a process is used in such a way that private data and private knowledge remain private even after the mining process. It is PPDM due to which the benefits of data mining be enjoyed, without compromising the privacy of concerned individuals.

PPDM Techniques can be classified over five dimensions. The first dimension is related to distribution of data i.e. Centralized or Distributed. The second dimension refers to the modification of original values of data that are to be released for data mining task. Modification is carried out using perturbation, blocking, aggregation, merging, swapping or sampling or any combination of these. The third dimension is that of data mining algorithms. The data mining algorithm are applied on the transformed data to get useful nuggets of information that were hidden previously. The fourth dimension refers to whether the raw data or aggregated data should be hidden. The fifth and the final dimension refer to the techniques that are used for protecting privacy. Based on these dimensions, different PPDM techniques may be classified into following five categories.

#### A. Anonymization based PPDM

Actually, when quasi identifiers [set of attributes that could potentially identify a record] are linked to publicly available data, identity of individual can be predicted with higher probability. Such attacks are called as linking attacks. Anonymization approach conceal identity or/and sensitive data about record owners using generalization and suppression in anonymized dataset. Replacing a value with less specific but semantically consistent value is called as generalization and suppression involves blocking the values. Such data when released for mining reduces the risk of identification when combined with publicly available data. But, besides, accuracy of the applications on the transformed data is reduced.

#### B. Randomized Response based PPDM

Randomized response is statistical technique to solve a survey problem. In Randomized response, the data is jumbled in such a way that the central place cannot tell with probabilities better than a pre-defined threshold, whether the data from a customer contains truthful information or false information. The information received from each individual user is snarled and if the number of users is significantly large, the aggregate information of these users can be predictable with good amount of accuracy. The data collection process in randomization method is carried out using two steps. During first step, the data providers randomize their data and transmit the randomized data to the data receiver. In second step, the data receiver reconstructs the original distribution of the data by employing a distribution reconstruction algorithm. The randomization method can be implemented at data collection time. It does not need a trusted server to contain all the original records in order to perform the anonymization process. The weakness of a randomization response based PPDM technique is that it treats all the records equal irrespective of their local density. This leads to a problem where the outlier records become more prone to adversarial attacks than to records in more dense regions in the data. One solution to this is to unnecessarily adding noise to all the records in the data. But, it reduces the utility of the data for mining purposes as the reconstructed distribution may not yield results in consistency of the purpose of data mining.

#### C. Condensation approach based PPDM

Condensation approach constructs constrained clusters in dataset and then generates fake data from the statistics of these clusters. It is called as condensation because of its approach of using condensed statistics of the clusters to generate fake data. It constructs groups of non-homogeneous size from the data, such that it is guaranteed that each record lies in a group whose size is at least equal to its secrecy level. Subsequently, fake data is generated from each group so as to create a synthetic data set with the same aggregate distribution as the original data. This approach can be effectively used for the problem of classification. The use of fake data provides an additional layer of protection, as it becomes difficult to perform adversarial attacks on synthetic data. Moreover, the aggregate behavior of the data is preserved, making it useful for a variety of data mining problems. This approach helps in better privacy preservation as compared to other techniques as it uses fake data rather than modified data. Moreover, it works even without redesigning data mining algorithms since the fake data has the same format as that of the original data. It is very effective in case of data stream problems where the data is highly dynamic. At the same time, data mining results get affected as large amount of information is lost because of the condensation of a larger number of records into a single statistical group entity.

#### D. Cryptography based PPDM

Cryptographic techniques are ideally meant for such scenarios where multiple parties collaborate to compute results or share non sensitive mining results and thereby avoiding disclosure of sensitive information. Cryptographic techniques find its utility in such scenarios because of two reasons: First, it offers a well defined model for privacy that includes methods for proving and quantifying it. Second a vast set of cryptographic algorithms and constructs to implement privacy preserving data mining algorithms are available in this domain. Although cryptographic techniques ensure that the transformed data is exact and secure but this approach fails to deliver when more than a few parties are involved.

#### E. Perturbation based PPDM

Perturbation has a long back history, being used in statistical disclosure control as it has an inherent property of simplicity, efficiency and ability to preserve statistical information. In perturbation the original values are replaced with some synthetic data values so that the statistical information computed from the perturbed data does not differ from the statistical information computed from the original data to a larger extent. The perturbed data records do not correspond to real-world record owners, so the attacker cannot perform the sensitive linkages or recover sensitive information from the published data. In perturbation approach, records released is synthetic i.e. it does not correspond to real world entities represented by the original data. Therefore the individual records in the perturbed data are meaningless to the human recipient as only statistical properties of the records are preserved. Perturbation can be done by using additive noise or data swapping or synthetic data generation. Since the perturbation method does not reconstruct the original values but only the distributions, new algorithms are to be developed for mining of the data.

### III. LITERATURE SURVEY

#### 1) Classification with no discrimination by preferential sampling

We can remove the sensitive data instead of relabeling it. The new solution to the CND problem by introducing a sampling scheme for making the discrimination free instead of relabeling the data set. The algorithm is used in this paper is classification algorithm. The goal of classification is to accurately predict the target class for each case in the data. Predicts categorical labels and classify the data based on the training set and the values in a classifying attribute and uses it in classifying new data. The techniques used in this paper is Pre-processing, Preferential sampling, Over sampling, Uniform sampling. In preprocessing is If there is much irrelevant and redundant information present or noisy and unreliable data, then knowledge discovery during the training phase is more complicated. Data preparation and filtering steps can considerable amount of processing period. Data pre-processing includes cleaning, normalization, transformation, characteristic extraction selection. In Preferential sampling arises when the process that determines the data location and the process being modeled are stochastically dependent. In the over sampling is the process of sampling a signal with a sampling frequency significantly higher than the twice the band width or highest frequency of the signal being sampled. Over sampling helps avoid aliasing, improves resolution and reduces noise. The equation is used  $f_s = 2b$  Where  $f_s$  is the sampling frequency and  $b$  is the bandwidth or highest frequency of the signal. Nyquist rate is then  $2b$ . The Uniform sampling defined as each Data objects probability is uniform. In this paper disadvantage is Discrimination were removed in ethical and legal region

#### 2) Three naive bayes approaches for discrimination free classification

In this method naive bayes is modify for discrimination classification. Discrimination laws do not allow the use of these rules of attributes such as gender, religion. Using decision rules that base their decision on these attributes in classifier. The approaches are used in this paper Naivesbayesmodel, Latent variable model, and Modifiednaivesbayes. The naivesbayes model is a bayes classifier is a simple possibility classifier based on applying bayes theorem with strong statistical independence assumption. Depending on precise nature of the probability model, naivebayes classifiers can be trained very efficiently in supervised learning. A latent variable model is a numerical model that relates a set of variables to set of latent variables. The responses on the indicators or manifest variables are the results of an individual's position on the latent variables. The modified naivebayes is Modify the probability distribution  $p(s/c)$  of the sensitive attribute values  $s$  given the class values.

#### 3) Fast algorithm for mining association rules

Fast algorithm is an efficient algorithm used to avoid the discrimination in data mining. In this paper algorithm apriori, aprioritid, AIS algorithm, apriorihybrid algorithm .The apriori algorithm is The large item sets of the previous pass were extended to get the new candidate item .pruning was done using the fact that any subsection of repeated item set should be frequent. In the aprioritid is related to the apriori algorithm and uses apriori function to determine the candidate sets. The difference for determining the support the database is not used after the first pass. In the AIS algorithm .In the AIS algorithm involves two concepts are extension of an item set, determining what should be in the candidate item set .The apriori hybrid algorithm is Uses apriori in the early passes and later shifts to aprioritid .In this paper disadvantages is An extra cost is sustained when shifting from apriori to aprioritid

#### 4) Discrimination prevention in data mining for intrusion and crime detection

In this paper techniques is used the anti discrimination techniques. Antidiscrimination law refers to the law on the right of people to be treated equally. In the political participation people must be dealt with on equal basis in any case of sex, age, race, nationality. The approaches are used preprocessing, post processing. The preprocessing is data preprocessing is the important process in the data mining. In there is much irrelevant and redundant information present or noisy and unreliable data, and then knowledge discovery during the training phase is more difficult. The analyzing data that has not been carefully screened for such problems can produce misleading results. The post processing is data mining is the process of sorting through large amounts of data and picking our relevant information .Data mining in relation to enterprise resource planning is the statistical and logical analysis of large sets of transaction data. The algorithm used in this paper is not efficient this is main drawback of this paper.

#### 5. Visual Data Mining for Higher-level Patterns: Discrimination- Aware Data Mining and Beyond

In this paper, we propose a visualization approach that can on the one hand be applied to any (classification or association) rules, but that is appropriate to bringing out characteristic of mined patterns that are especially important in discrimination-aware and privacy aware data mining. We define new interestingness proceeding for items and rules and show various ways in which these can help in highlighting information in communicating settings. We conclude by arguing how this approach can lead to a new generation of feedback and awareness tools. The need to inspect mining results carefully for such meta-level relationships between features and outcomes becomes even stronger when specific data, rules and other patterns become the object of scrutiny: The flipside of data mining is that it may make relationships visible that various stakeholders do not wish to become explicit, and that the patterns it finds may suggest actions that various stakeholders do not wish to be taken. Such concerns may lead to a new approach to keep and/or treat these data as private.

#### 6. A Survey of association rule hiding methods for privacy preserving

In this paper, we present taxonomy and a survey of recent approaches that have been applied to PPDM. Association rule hiding refers to the process of modifying the original databases in such a way that certain sensitive association rule disappear without seriously affecting the data and the non-sensitive rules. Finally, we conclude our study by heuristic based algorithms and enumerating interesting future directions in this research body. In recent years, a number of techniques have been proposed for modifying or transforming the data in such a way so as to preserve privacy.

The Association rule hiding approach is one of the widely used approach. Association rule hiding algorithms can be divided into three distinct classes, namely border-based approaches, exact approaches and heuristic approaches.

**1) Border-based Approaches:** These approaches consider the task of sensitive rule hiding through modification of the original borders in the lattice of the frequent and the infrequent patterns in the dataset. In these schemes, the positive and the negative borders in the lattice of all item sets are computed first and then focus on preserving the quality of the computed borders during the hiding process. The quality of the borders directly affects the quality of the sanitized database that is produced, which can be maintained by greedily selecting those modifications that lead to minimal side-effects.

**2) Exact Approaches:** Exact approaches are typically capable of providing superior solutions but at a high computational cost. They achieve this by formulating the sanitization process as a constraint satisfaction problem and by solving it using an integer/linear programming solver. Thus, the sanitization of the dataset is performed as an atomic operation which avoids the local minima.

**3) Heuristic Approaches:** These approaches involve efficient, fast algorithms that selectively sanitize a set of transactions from the database to hide the sensitive knowledge. Due to their efficiency and scalability, the heuristic approaches have been the focus of attention for the vast majority of researchers in the knowledge hiding field. Heuristic Based Approaches can be divided into two groups based on data modification techniques:

A. Data-Distortion : It is based on data perturbation or data transformation, and in particular, the procedure is to change a selected set of 1-values to 0-values (delete items) or 0-values to 1- values (add items) if we consider the transaction database as a two-dimensional matrix. It is aimed to reduce the support or confidence of the sensitive rules below the user pre-defined security threshold.

B. Data-Blocking: It is another data modification approach for association rule hiding. Instead of making data distorted (part of data is altered to false), blocking approach is implemented by replacing certain data items with a question mark “?”. The introduction of this special unknown value brings uncertainty to the data, making the support and confidence of association rule become too uncertain intervals respectively.

Before we conclude our study we have provided an overview of heuristic approaches. These approaches involve efficient, fast algorithms to hide the sensitive knowledge. Due to their efficiency and scalability, the heuristic approaches have been the focus of attention for the vast majority of researchers in the knowledge hiding field. Different algorithms can be designed and developed by taking ideas from existing algorithms and compare their efficiency using metrics.

#### IV. PROPOSED SYSTEM

Our Proposed data transformation methods rule protection and rule generalization are based on measures for both direct and indirect discrimination and can deal with several discriminatory items. We present a unified approach to direct and indirect discrimination anticipation, with finalized algorithms and all possible data transformation methods based on rule protection and or rule generalization that could be applied for direct or indirect discrimination prevention. We propose new utility measures to evaluate the different proposed discrimination prevention methods in terms of data quality and discrimination removal for both direct and indirect discrimination. Direct and indirect discrimination discovery includes identifying discriminatory rules and redlining rules. Using the above transformation methods effectively to identify the categories and remove direct and indirect discrimination method. Finally, discrimination free data models can be produced from the transformed data set without seriously damaging data quality. Discrimination prevention methods in terms of data quality and discrimination removal for both direct and indirect discrimination. The proposed techniques are quite successful in both goals of removing discrimination and preserving data quality.

#### V. CONCLUSION

Along with privacy, discrimination is a very important issue when considering the legal and ethical aspects of data mining .It is more than noticeable that most people do not want to be discriminated because of their gender, religion nationality, age, and so on exceptionally, when those attributes are used for making decisions about them like giving them a job, loan, insurance, etc. The purpose of this paper was to develop a new preprocessing discrimination prevention methodology including different data transformation methods that can prevent direct discrimination indirect discrimination or both of them at them at the same time. To attain this objective, the first steps are measure discrimination and identify categories and groups of individuals that have been direct and indirect discriminated in decision making process. The second step is to transform data in proper way to remove all those discrimination biases at last, discrimination free data models can be produced from the transformed data set without seriously damaging data quality.

#### REFERENCES

- [1] F. Kamiran and T. Calders, “Classification with no Discrimination by Preferential Sampling,” Proc. 19<sup>th</sup> Machine Learning Conf. Belgium and The Netherlands, 2010.
- [2] T. Calders and S. Verwer, “Three Naive Bayes Approaches for Discrimination-Free Classification,” Data Mining and Knowledge Discovery, vol. 21, no. 2, pp. 277-292, 2010.
- [3] R. Agrawal and R. Srikant, “Fast Algorithms for Mining Association Rules in Large Databases,” Proc. 20<sup>th</sup> Int’l Conf. Very Large Data Bases, pp. 487-499, 1994.

- [4] S. Hajian, J. Domingo-Ferrer, and A. Martı́nez-Balleste', "Discrimination Prevention in Data Mining for Intrusion and Crime Detection," Proc. IEEE Symp. Computational Intelligence in Cyber Security (CICS '11), pp. 47-54, 2011.
- [5] D.Pedreschi, S.Ruggeri and F.Turini, "Discrimination Aware Data Mining," Proc. 14<sup>th</sup> ACM Int'l Conf. Knowledge Discovery and Data Mining (KDD '08), pp. 560-568, 2008.
- [6] D.Pedreschi, S. Ruggeri, and F. Turini, "Measuring Discrimination in Socially-Sensitive Decision Records," Proc. Ninth SIAM Data Mining Conf. (SDM '09), pp. 581-592, 2009.
- [7] V.S. Verykios, A.K. Elmagarmid, E. Bertino, Y. Saygin, E. Dasseni, "Association Rule Hiding," IEEE Transactions on Knowledge and Data Engineering, vol. 16, Apr. 2004, pp, 434-447.