

# Information Security in Large Amount of Data: Privacy and Data Mining

Prof. Harish Barapatre <sup>[1]</sup>, Ms. Sampada Samel <sup>[2]</sup>

Department of Computer Science and Engineering  
YTIET, Bhivpuri Road  
Mumbai - India

## ABSTRACT

Now a days, In our day to day life, development in data mining becomes very much popular. But, growing popularity and development in data mining technologies brings serious threat to the security of individual's sensitive information. So, To avoid access to one's sensitive information, Privacy Preserving Data Mining (PPDM) is used. In this technique, data gets modified in order to secure one's sensitive information. PPDM technique is mainly focused on how privacy maintained at Data Mining. However, Sensitive information can be retrieved at Data Collection, Data Publishing and Information Delivering processes. In this Paper, we briefly discuss the Privacy Preserving Data Mining with respect to user such as Data Provider, Data Collector, Data Miner and Decision Make. We will discuss privacy concerns related to each user. We also discuss about the Game Theory approaches.

**Keywords:-** Data Mining, Privacy, information, KDD, PPDM, etc

## I. INTRODUCTION

In recent years, Data mining has gained a lot public-ity. Data Mining is the process of discovering interesting patterns and knowledge from large amounts of data.[1]

### 1.1. THE PROCESS OF KDD

To retrieve useful information from data, following steps are followed:

Step 1: Data Processing. It includes basic operations like Data selection, Data cleaning and Data integration. Data selection process retrieves the data relevant to KDD task from database. Data cleaning process is used to remove noise and inconsistent data and Data Integration process combines data from multiple sources.

Step 2: Data Transformation. Data Transformation performs feature selection and feature transformation.

Step 3: Data Mining. Data Mining process extracts the useful data from large amount of data.

Step 4: Pattern evaluation and presentation. It includes basic operations like identifying the truly interesting patterns which represent knowledge and presenting the mined knowledge in an easy-to-understand fashion.

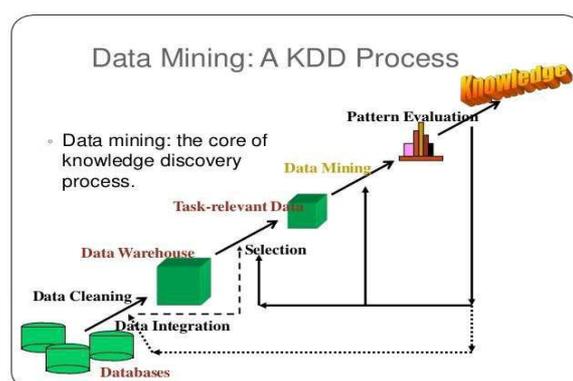


Figure 1. An overview of KDD process

### 1.2. THE PRIVACY CONCERN AND PPDM

In Recent years, Data Mining becomes very valuable in areas such as Web search, scientific discovery, digital libraries, etc. But, In Data Mining, individual privacy may be compromised due to unauthorized access to the personal information. Due to unauthorized access to one's private data, people may suffer. So, To avoid this PPDM technique has gained a great development in recent years. The main objective of PPDM (Privacy Preserving Data Mining)[2] is to safeguard sensitive information from unnecessary disclosure of personal information.

### 1.3. USER ROLE-BASED METHODOLOGY

PPDM technique mainly focus on how to hide sensitive information from data mining operations. But, It is also true that security issues arise in data collecting, data processing also. As shown in Fig.1, we can identify different types of users in data mining scenario.[3]

**Data Provider:** The user who owns some data that are desired by data mining task.

**Data Collector:** The user who collects data from data providers and then publish the data to the data miner.

**Data Miner:** The user who performs data mining tasks on the data.

**Decision Maker:** The user who makes decisions based on the data mining results in order to achieve certain goals.

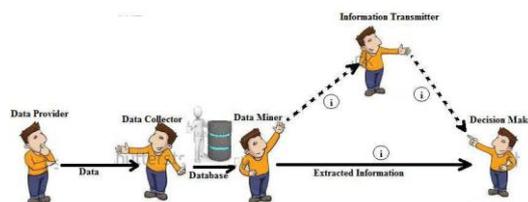
We will discuss about each user's privacy concerns one by one.[4]

**1.3.1. DATA PROVIDER.** Data provider has to ensure that the data provided by him to the data collector must not contain any privacy information. the provider should be able to make his very private data, namely the data containing information that he does not want anyone else to know, inaccessible to the data collector. On the other hand, if the provider has to provide some data to the data collector, he wants to hide his sensitive information as much as possible and get enough compensations for the possible loss in privacy.

**1.3.2. DATA COLLECTOR.** The data received from data provider may contain sensitive information. Directly sending such a data to the data miner, may violate data provider's privacy, hence data modification is required. But major concern of data collector is to guarantee that the modified data contain no sensitive information and still preserve high utility. If data after modification is not useful then collecting data from data provider becomes meaning-less.

**1.3.3. DATA MINER.** The data miner receives data from data collector and performs different mining algorithms to achieve certain goal. Data miner wishes to extract useful information from the data provider. Data miner has major responsibility of protecting sensitive data, while data miner can focus on how to hide sensitive information after mining from untrusted parties.

**1.3.4. DECISION MAKER.** As shown in the fig.2, Decision maker gets information from Data Miner and from Information transmitter. It is likely that the information transmitter may change mining results intentionally or un-intentionally. So, main concern of Decision maker is that whether mining results are credible.



**Figure 2. A simple illustration of the application scenario with data mining at the core.**

## II. DATA PROVIDER

### 2.1. CONCERNS OF DATA PROVIDER

As shown in fig.2, There are two data providers. Data provider which provides data to the data collector. and Data collector which provides data to the Data miner. We will consider here ordinary Data provider which owns some data which is desired by Data mining process. If the data provider provides his sensitive data to the Data collector, then his privacy might get compromised due to un-expected data breach or exposure of sensitive information.

### 2.2. APPROACHES TO PRIVACY PROTECTION

**2.2.1. LIMIT THE ACCESS.** Data provider can provide his data to the Data collector in active way or passive way.

**Active Way:** Data provider voluntarily opts in a survey initiated by the Data collector or fill in some registration form to create an account in a website.

**Passive Way:** Data collector collects the Data provider's data by the Provider's routine activities. Data collector retrieves the data by recording provider's routine activities unaware of Data Provider.

Data provider can avoid tracking his routine activities by emptying cache, deleting cookies, clearing usage records of applications, etc. Current security tools that are developed for internet environment to protect provider's data can be categorized into three types:

1. Anti-tracking extensions: Data collector can retrieve user's sensitive data by tracking his/her routine activities. To avoid this unauthorized access to the provider's data, provider can use anti-tracking tools such as Disconnect, Do Not Track Me, Ghostery, etc.
2. Advertisement and script blockers: By adding browser extensions such as AdblockPlus, NoScript, FlashBlock, etc. user can block advertisements on the sites and kill scripts and widgets that send user's sensitive data to unknown third party.

3. Encryption tools: A user can utilise encryption tools such as MailCloak, TorChat to encrypt mails to make sure that a private communication between two parties cannot be intercepted by third parties.

In addition to all these tools, user can use anti-virus, anti-malware tools to protect data. Using such a tools user can limit the access of his/her sensitive data to third parties.

### 2.2.2. TRADE PRIVACY FOR BENEFIT.

In some cases, provider needs to make tradeoff between the loss of privacy and the benefits brought by participating in Data mining. Consider shopping website. If website tracks user's routine activities and find out user interested products, then it will be beneficial for user also. User can fill better shopping experience in this case. Now suppose user has to enter information about salary on shopping website, then the website can show the interested item in user's budget. so, disclosure of sensitive information such as salary is more beneficial as it reduces the searching time of the user.

### 2.2.3. PROVIDE FALSE DATA.

Data providers takes efforts to hide sensitive information from data collector. Data collector takes efforts to hide sensitive information from Data miner. But, In today's internet age, internet users cannot completely stop the unwanted access to user's personal information. So, instead of trying to limit the access, the data provider can provide false information to untrustworthy Data collectors. Following methods can be used to falsify the data.

1. sockpuppets: A sockpuppet is a false online identity. By using multiple sockpuppets, the data produced by one individual's activities will be deemed as data belonging to different individuals. Data collector do not have enough knowledge to relate different sockpuppets with one individual. So, user's true activities are unknown to others and user's sensitive information cannot be easily retrieved.
2. Clone identity: This technique can protect user's privacy by creating fake identity. This clone identity automatically makes some actions which are totally different from user's actions. So if the third party tries to retrieve user's data, then it will get data from clone identity which is completely different. By this way, user's sensitive data is inaccessible to the unwanted user.
3. MaskMe: By adding MaskMe browser extension user can hide his/her sensitive data. Wherever user perform online transaction, user has to enter his sensitive information such as email id, Bank details, etc. Using this extensions, many aliases are created. so, user's sensitive data can be secured.

## 2.3. CONCLUSION

Data Provider can secure his data by three ways: he can limit the access to his online activities or data by using

anti-tracking extensions, advertisement and script block-ers or by using encryption tools to encrypt emails between two private parties. Data Provider can also demand for high price to disclose his private data with others. Nowadays, whatever you try, but the hackers can get your secure information. So, Data provider can provide false data to misguide such a hackers. Using sockpuppet, Data provider can make different sockpuppets. Data provider can use MaskMe to mask his sensitive information.

## III. DATA COLLECTOR

### 3.1. CONCERNS OF DATA COLLECTOR

As shown in Fig.2, Data collector collects the data from Data provider and provide this data to Data Miner. Data collected from Data provider may contain sensitive information of the individual. If such a data is directly send to the Data Miner, then individual's sensitive information disclosed to the unwanted third parties or Data miner. So, before sending data to the Data miner, Data collector has to check whether data contains sensitive information or not. If so then Data collector has to encrypt the data collected from Data provider and then send it to the Data miner.

Data collector has to modify the data before releasing it to the Data miner. But, After using modification techniques, there will be loss in data utility. So the main concern of Data miner is that the data must retained utility after the modification. Otherwise collecting data is waste process.

The data modification process adopted by Data collector with the goal of preserving privacy and utility simultaneously is called as Privacy Preserving Data Publishing (PPDP)[5].

### 3.2. APPROACHES TO PRIVACY PROTECTION

3.2.1. BASICS OF PPDP. The original data is in the form of table with multiple records. Each record consists of four types of attributes.

1. Identifier (ID): Attributes that uniquely identifies user on cloud
2. Quasi-identifier (QID): Attributes that linked with the external data to re identify user.
3. Sensitive Attribute (SA): Attributes that the user wants to hide for privacy.
4. Non-Sensitive Attribute (NSA): Attributes that user don't matter to disclose with anyother.

Data anonymization is a type of information sanitization whose intent is privacy protection. It is the process of either encrypting or removing personally identifiable information from data sets so that identity and sensitive attribute values hidden from adversaries. Record linkage (RL) refers to the task of finding records in a data set that refer to the

same entity across different data sources (e.g., data files, books, websites, databases). In Attribute linkage (AL), the adversary may not precisely identify the record of the target victim, but could infer his/her sensitive values from the published data, based on the sensitive values associated to the group, that the victim belongs to. In Table linkage (TL), the attack seeks to determine the presence or absence of victim's record in the released table. Probabilistic linkage, takes a different approach to the record linkage problem by taking into account a wider range of potential identifiers, computing weights for each identifier based on its estimated ability to correctly identify a match or non-match, and using these weights to calculate probability that two given records refer to the same entity. Different privacy models includes k-anonymity, l-diversity, t-closeness, epsilon-differential privacy.

k-anonymity is used for record linkage.

l-diversity is used for preventing record and attribute linkage.

t-closeness is used for preventing attribute and probabilistic linkage.

epsilon-differential is used for preventing table and probabilistic linkage.

Among all these, k anonymity is widely used. In k-anonymity, attributes are suppressed or generalized until each row is identical with atleast k-1 other rows. Thus it prevents definite database linkages. K-anonymity guarantees that the data released is accurate.

Consider following table which gives idea about k-anonymity. Now consider above table A and table B which denotes Raw table and anonymized table respectively. Using K-anonymous technique, Data collector can hide Identifiers and Quasi-identifier fields from third parties. As shown in fig.B, quasi-identifier fields such as age, sex zipcode are replaced by either special characters or range values or common attribute. So, by using such anonymous table, adversaries are unable to track particular individual. then, the probability that the individual's record being identified by the adversary will not exceed 1/K.

To satisfy privacy model conditions, following operations can be done.

Generalization: Replace some values in the table with with parent value in the taxonomy of an attribute.

Suppression: Replace some values in the table with a special character ("\*"), as shown in the column "Zip-code" in table.B.

Anatomization: De-associates the relationship between two.

Permutation: De-associates the relationship between a quasi-identifier and a numerical sensitive attribute by partitioning a set of data records into groups and shuffling their sensitive values within each group.

Age	Sex	Zipcode	Disease
5	Female	12000	HIV
9	Male	14000	dyspepsia
6	Male	18000	dyspepsia
8	Male	19000	bronchitis
12	Female	21000	HIV
15	Female	22000	cancer
17	Female	26000	pneumonia
19	Male	27000	gastritis
21	Female	33000	flu
24	Female	37000	pneumonia

(a)

Age	Sex	Zipcode	Disease
[1, 10]	People	1****	HIV
[1, 10]	People	1****	dyspepsia
[1, 10]	People	1****	dyspepsia
[1, 10]	People	1****	bronchitis
[11, 20]	People	2****	HIV
[11, 20]	People	2****	cancer
[11, 20]	People	2****	pneumonia
[11, 20]	People	2****	gastritis
[21, 60]	People	3****	flu
[21, 60]	People	3****	pneumonia

(b)

Figure 3. An example of 2-anonymity where QID = Age,Sex,Zipcode.(a)Original Table (b)2-anonymous table

Perturbation: Replace original data values with some synthetic data values.

But, All these privacy model information results into information loss.

3.2.2. PRIVACY PRESERVING PUBLISHING OF SOCIAL NETWORK DATA. Social network data is always represented in the form of graph. where vertex represents an entity and edge represents the relationship between two entities. So, In case of social network PPDP deals with the anonymized graph data. Anonymizing social network data[6] is much more challenging than that of relational data.

it is much more challenging to model background knowledge of adversaries and attacks about social network data than that about relational data. On relational data, it is often assumed that a set of attributes serving as a quasiidentifier is used to associate data from multiple tables, and attacks mainly come from identifying individuals from the quasi-identifier. However, in a social network, many pieces of information can be used to identify individuals, such as labels of vertices and edges, neighborhood graphs, induced subgraphs, and their combinations. It is much more complicated and much more difficult than the relational case.

it is much more challenging to measure the information loss in anonymizing social network data than that in anonymizing relational data. Typically, the information loss in an anonymized table can be measured

using the sum of information loss in individual tuples. Given one tuple in the original table and the corresponding anonymized tuple in the released table, we can calculate the distance between the two tuples to measure the information loss at the tuple level. However, a social network consists of a set of vertices and a set of edges. It is hard to compare two social networks by comparing the vertices and edges individually. Two social networks having the same number of vertices and the same number of edges may have very different network-wide properties such as connectivity, betweenness, and diameter. Thus, there can be many different ways to assess information loss and anonymization quality

it is much more challenging to devise anonymization methods for social network data than for relational data. Divide-and-conquer methods are extensively applied to anonymization of relational data due to the fact that tuples in a relational tables are separable in anonymization. In other words, anonymizing a group of tuples does not affect other tuples in the table. However, anonymizing a social network is much more difficult since changing labels of vertices and edges may affect the neighborhoods of other vertices, and removing or adding vertices and edges may affect other vertices and edges as well as the properties of the network.

3.2.3. ATTACK MODEL. In Anonymized social network data, adversaries often rely on the background knowledge to de-anonymize individuals and learn relationships between de-anonymized individuals.

“Seed and Grow” algorithm[7] invented by Peng et al. is used to identify users from an anonymized social graph, based solely on graph structure. The seed stage plants a small specially designed sub graph into undirected graph before its release. After anonymized graph is released, the attacker locates sub graph in anonymized graph, so, the vertices are readily identified and serves as the initial seeds to be grown. The grow stage is essentially comprised of a structure based vertex matching, which further identifies vertices adjacent to initial seeds. This is self reinforcing process, in which the seeds grow larger as more vertices are identified.

“Structural Attack”[8] is the attack that de-anonymize social graph data. This attack uses cumulative degree of a vertex.

“Mutual Friend Attack” is de-anonymized data based on the number of social common friends of two directly connected individuals. As shown in Fig.4, The anonymization mapping  $f$ , is a random, secret mapping. Naive anonymization prevents re-identification when adversary has no information about individual in original graph. In practice the adversary may

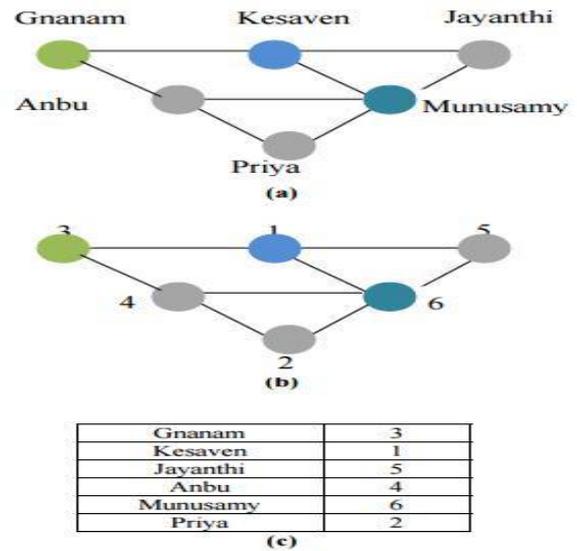


Figure 4. Example of mutual friend attack: (a)Original network; (b)Naive anonymized network. (c)Mapping Function (f)

have access to external information about the entities in the graph and their relationships. This information may be available through a public source beyond the control of the data owner, or may be obtained by the adversary's malicious actions. For example, for the graph in Figure 1, the adversary might know that Munusamy has three or more neighbors, or that Gnanam is connected to at least two nodes, each with degree 2. Such information allows the adversary to reduce the set of candidates in the anonymized graph for each of the targeted individuals. Although an adversary may also have information about the attributes of nodes, the focus of this paper is structural re-identification, where the adversary's information is about graph structure. Re-identification with attribute knowledge has been well studied, as have techniques for resisting it. More importantly, many network analyses are concerned exclusively with structural properties of the graph; therefore safely publishing an unlabeled network is a legitimate goal.

1. PRIVACY MODEL No. of privacy models are proposed for graph data based on classic k-anonymity model.[9]

- (a) k-NMF anonymity: It protects the privacy of relationship from the mutual friend attack.
- (b)  $K^2$ -degree anonymity: It protects information loss due to friendship attack.
- (c) k-structural diversity anonymization (k-SDA): It protects information loss due to degree attack.

3.2.4. PRIVACY-PRESERVING PUBLISHING OF TRAJECTORY DATA. In recent years, LBS(Location

Based Services)[10] becomes very popular. Using these services user can able to find out interesting places near him/her. If he/she wants information about nearest bank. then he/she can use such a services and able to find out nearest bank location. To provide location-based services, commercial entities (e.g. a telecommunication company) and public entities (e.g. a transportation company) collect large amount of individuals' trajectory data, i.e. sequences of consecutive location readings along with time stamps. If the data collector publish such spatio-temporal data to a third party (e.g. a data-mining company), sensitive information about individuals may be disclosed. To realize a privacy-preserving publication, anonymization techniques can be applied to the trajectory data set, so that no sensitive location can be linked to a specific individual.

### 3.3. CONCLUSION

Data Collector receives data from Data provider and sends that data to the Data Miner. Before sending this data to the Data miner, Data collector has to check whether data contains any private information or not. data Collector has to develop different attack models to check whether data contains any private information about data provider

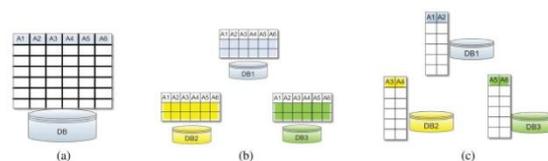
## IV. DATA MINER

### 4.1. CONCERNS OF DATA MINER

Data collector sends the data after modification to the Data miner. Then, the Data miner has to retrieve the important data using different data mining techniques. The primary concern of data miner is how to prevent sensitive information from appearing in the mining results. To perform a privacy preserving data mining, the data miner usually needs to modify the data he got from the data collector. As a result, the decline of data utility is inevitable. Similar to data collector, the data miner also faces the privacy utility tradeoff problem. But in the context of PPDM, quantifications of privacy and utility are closely related to the mining algorithm employed by the data miner.

### 4.2. APPROACHES TO PRIVACY PROTECTION

Privacy preserving data mining approaches are classified into two main categories i.e. Approaches for centralized data mining and Approaches for Distributed data mining. Distributed Data mining again further classified as horizontally partitioned data and vertically partitioned data as shown in Fig.5. Now, most services are using distributed data mining where secure multi-party computation is used. SMS (Secure Multi-party Computation) is a subfield of cryptography. SMC assumes that there are number of participants  $P_1; P_2; P_3; \dots; P_m$  with having private data  $D_1; D_2; D_3; \dots; D_m$  respectively. The participants want to compute the value of public function  $f$ . We can say that SMC protocol is secure if, at the end of computation, par-



**Figure 5. Data Distribution (a)Centralized Data (b)Horizontally Partitioned Data (c)Vertically Partitioned Data.**

icipant can able to view only their own data. So, the main goal of SMC protocol is to find correct data mining results without revealing participants data with others.

#### 4.2.1. PRIVACY-PRESERVING ASSOCIATION RULE MINING.

Association rule mining is a two-step process:

- (1) Finding all frequent itemsets;
- (2) Generating strong association rules from the frequent itemsets.

The purpose of privacy preserving is to discover accurate patterns to achieve specific task without precise access to the original data. The algorithm of association rule mining is to mine the association rule based on the given minimal support and minimal confidence. Therefore, the most direct method to hide association rule is to reduce the support or confidence of the association rule below the minimal support of minimal confidence. With regard to association rule mining, the proposed methodology that is effective at hiding sensitive rules is implemented mainly by depressing the support and confidence.

various kinds of approaches have been proposed to perform association rule hiding. These approaches are classified as five categories.

**Heuristic Distortion Approaches:** selects appropriate data sets for data modification.

**Heuristic blocking approaches:** reduces the degree of support and confidence of the sensitive association rules by replacing certain attributes of some data item with specific symbol.

**Probabilistic distortion approaches:** distorts the data through random numbers generated from predefined probability distortion function.

**Exact database distortion approaches:** formulates the solution of the hiding problem as a constraint satisfaction problem (CSP), and apply linear programming approaches to its solution.

**Reconstruction-based approaches:** generates a database from the scratch that is compatible with a given set of non-sensitive association rules.

#### 4.2.2. PRIVACY-PRESERVING CLASSIFICATION.

Data classification is a two step process.

1. Step1 : Learning Step: algorithm generate classification model.
2. Step2 : Classification: Develops different classification models such as Decision Tree, Bayesian Model, Support Vector Machine (SMC), etc.

**1. DECISION TREE**

A decision tree[12] is defined as a predictive modeling technique from the field of machine learning and statistics that builds a simple tree-like structure to model the underlying pattern of data. Decision tree is one of the popular methods to classify is able to handle both categorical and numerical data and perform classification with minimal computation. Decision trees are often easier to understand and compute. Decision tree is a classifier which is a directed tree with a node having no incoming edges called root. All the nodes except root have exactly one incoming edge. Each non-leaf node called internal node or splitting node contains a decision and most appropriate target value assigned to one class is represented by leaf node. Decision tree classifier is able to divide a complex process into number of simpler processes. The complex decision is subdivided into simpler decision on the basis of splitting of complex process into simple processes. It divides complete data set into smaller subsets. Information gain, gain ratio, gini index are three basic splitting criteria to select attribute as a splitting point. Decision trees can be built from historical data they are often used for explanatory analysis as well as a form of supervised learning. The algorithm is designed in such a way that it works on all the data that is available and as perfect as possible. According to Breiman et al. the tree complexity has a crucial effect on its accuracy performance. The tree complexity is explicitly controlled by the pruning method employed and the stopping criteria used. Usually, the tree complexity is measured by one of the following metrics:

- The total number of nodes
- Total number of leaves
- Tree depth
- Number of attributes used

Decision tree induction is closely related to rule induction. Each path from the root of a decision tree to one of its leaves can be transformed into a rule simply by conjoining the tests along the path to form the antecedent part, and taking the leaf's class prediction as the class value. The resulting rule set can then be simplified to improve its accuracy and comprehensibility to a human user.

**2. NAIVE BAYESIAN CLASSIFICATION**

The Naive bayesian classifier[13] is a simple but efficient baseline classifier. This classifier used for text

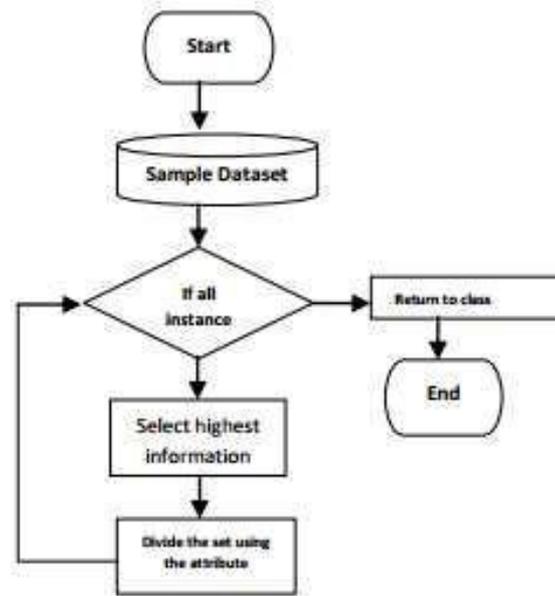


Figure 6. Flowchart for Decision Tree Based Classification.

classification. Naive bayesian is based on a bayesian formulation of the classification problem which uses the simplifying assumption of attribute independence. It is simple to compute and computation calculates good results. Thus, preliminary evaluation is carried out using the Naive Bayesian classifier to serve both as a baseline and to decide whether more sophisticated solutions are required. The problem of secure distributed classification is an important one. The goal is to have a simple, efficient, easy to compute and privacy-preserving classifier. The ideal would be for all parties to decide on a model. Jointly select/discover the appropriate parameters for the model and then use the model locally as and when necessary. We discuss the specifics in the context of the Naive Bayesian classifier later. In this, data is assumed to be horizontally partitioned. This means that many parties collect the same set of information about different entities. Parties want to improve classification accuracy as much as possible by leveraging other parties data. They do not want to reveal their own instances or the instance to be classified. Thus, what we have is a collaboration for their own advantage. One way to solve this is to decide on a model. The model parameters are generated jointly from the local data. Classification is performed individually without involving the other parties. Thus, the parties decide on sharing the model, but not the training set nor the instance to be classified. This is quite realistic. For example, consider banks which decide to leverage all data to identify fraudulent credit card usage, or insurance companies which jointly try to identify high-risk customers. In this paper, we use / extend several existing

cryptographic techniques to create a privacy preserving Naive Bayesian Classifier for horizontally partitioned data

### 3. SUPPORT VECTOR MACHINE

Support vector Machine (SVM)[14] is one of the most developed classification methodology in data mining. It provides properties such as the margin maximization and nonlinear classification via kernel tricks and has proven to be effective in many real world applications. Privacy preserving SVM classification solution, PP-SVM which constructs the global SVM classification model from the data distributed from a multiple parties. The data may be partitioned horizontally, vertically or in an arbitrary manner between the parties. The data of each party is kept private, while the final model is constructed at an independent site. This independent site then performs classification of new instances. of sense in many different contexts. For example, consider a clearing house for a consortium of banks. The different banks collect data of their customers. The features collected such as age, gender, balance, average monthly income, etc. are the same for all bank. Thus, the data is horizontally distributed. the clearing house is an independent entity, unrelated to any of the banks. The classification model is constructed at the clearing house while preserving the privacy of the individual data from each of the banks. When a bank has a new instance it wants to classify, it goes through a secure protocol with the clearing house to classify just this instance. The clearing house learns nothing. This would allow all of the banks to leverage the global data without compromising on privacy at all.

### 4.3. CONCLUSION

Data Miner has to retrieve the important data using different data mining techniques. so, The primary concern of data miner is how to prevent sensitive information from appearing in the mining results. To perform a privacy preserving data mining, the data miner usually needs to modify the data he got from the data collector. As a result, the decline of data utility is inevitable. Similar to data collector, the data miner also faces the privacy utility tradeoff problem. By using different algorithm techniques such as Decision tree, Support Vector Machine, Naive Bayesian Techniques Data collector modifies data and sends it to the Decision maker.

## V. DECISION MAKER

### 5.1. CONCERNS OF DECISION MAKER

The final goal of data mining process is to provide useful information to the decision maker, so that the decision maker can choose a result which is better way to achieve his objective. As we can see, Data provider sends data

to Data collector, Data collector sends data to the Data miner and finally Data miner sends data to the Decision Maker. So, we can say that Decision maker is less responsible for the data security. The data mining results provided by the data miner are of high importance to the decision maker. If the results are disclosed to someone else, e.g. a competing company, the decision maker may suffer a loss. That is to say, from the perspective of decision maker, the data mining results are sensitive information. On the other hand, if the decision maker does not get the data mining results directly from the data miner, but from someone else which we called information transmitter, the decision maker should be skeptical about the credibility of the results, in case that the results have been distorted. Therefore, the privacy concerns of the decision maker are twofold: how to prevent unwanted disclosure of sensitive mining results, and how to evaluate the credibility of the received mining results.

### 5.2. APPROACHES TO PRIVACY PROTECTION

#### 5.2.1. DATA PROVENANCE

Usually, Decision maker receives data from data miner but, in some cases if the decision maker does not get the data mining results directly from the data miner i.e. receives data from other sources, then he wants to know how the results are delivered to him and what kind of modifications are applied to the results, so that he can decide whether the results are trusted or not. This is why "provenance" is needed. The term provenance [15] originally refers to the custody of the data. In computer science, data provenance refers to the information that helps determine the derivation history of the data, starting from the original source. Two kinds of information can be found in the provenance of the data: the ancestral data from which current data evolved, and the transformations applied to ancestral data that helped to produce current data. With such information, people can better understand the data and judge the credibility of the data. data provenance has been extensively studied in the fields of databases and work flows. Several surveys are now available. The following five aspects are used to capture the characteristics of a provenance system:

1. Application of provenance. Provenance systems may be applied in many fields to support a number of uses, such as estimate data quality and data reliability, trace the audit trail of data, repeat the derivation of data, etc.
2. Subject of provenance. Provenance information can be collected about different sources and at various levels of detail.
3. Representation of provenance. There are mainly two types of methods to represent provenance information, one is annotation and the other is inversion. The annotation uses metadata. Using inversion method, derivations are inverted to find out inputs to the derivations.

4. Provenance storage. Provenance is tightly coupled with the data it describes and located in the same data storage system or even be embedded within the data. Alternatively, provenance can be stored separately with other metadata or simply by itself.
5. Provenance dissemination. A provenance system can use different ways to provide the provenance information, such as providing a derivation graph that users can browse and inspect.

### 5.2.2. WEB INFORMATION CREDIBILITY.

Be-cause of the lack of publishing barriers, the low cost of dissemination, and the lax control of quality, credibility of web information[16] has become a serious issue. Tudjman et al. identify the following five criteria that can be employed by Internet users to differentiate false information from the truth:

1. Authority: the real author of false information is usually unclear.
2. Accuracy: false information does not contain accurate data or approved facts.
3. Objectivity: false information is often prejudicial.
4. Currency: for false information, the data about its source, time and place of its origin is incomplete, out of date, or missing.
5. Coverage: false information usually contains no effective links to other information online.

### 5.3. CONCLUSION

the privacy concerns of Decision Maker are twofold: how to prevent unwanted disclosure of sensitive mining results, and how to evaluate the credibility of the received mining results.

## VI. GAME THEORY IN DATA PRIVACY

### 6.1. GAME THEORY PRELIMINARIES

Till now, we have discussed about the privacy concerns for each users and control measures taken by each user. Now, we will focus on the iterations among different users.

The main elements of Game theory[17] are players, actions, payoffs and information. Players perform some actions at particular times in the game. As per the performed actions, player will receive payoffs. This payoff is depends on the both player's own actions and other player's actions. Information sets represents a player's knowledge about the values of different variables in the game. A player's strategy is a rule that suggests him which action to choose at what time of the game. A strategy profile is an ordered set consisting of one strategy for each of the players in the game. An equilibrium is a strategy profile consisting of a best strategy for each of the players in

the game. The most important equilibrium concept for the majority of games is Nash equilibrium. A strategy profile is a Nash equilibrium if no player has incentive to deviate from his strategy, given that other players do not deviate.

In many fields such as, Computer Science, Economics, Politics, etc game theory has been successfully applied.

### 6.2. PRIVATE DATA COLLECTION AND PUBLICATION

Data collector collects data from Data provider. So, in order to get benefit, Collector may need to negotiate with the Data provider about the "price" of the sensitive data and the level of privacy protection. Adl build a sequential game model to analyze private data collection process. In this model, an end user who wants to buy a data from the data collector, makes an offer to the collector at the beginning of the game. If the data collector accepts offer, then he announces some incentives to data provider's data. Before selling the collected data to the end user, the data collector applies anonymization technique to the data and sends anonymized data in order to protect the privacy of data providers at certain level. On the price and incentives offered by the collector, Data provider decides whether to send data or not.

### 6.3. PRIVACY PRESERVING DISTRIBUTED DATA MINING

#### 6.3.1. SMC-BASED PRIVACY PRESERVING DISTRIBUTED DATA MINING

. Secure Multi-party Computation (SMC) is used in privacy preserving data mining[18]. In this computation, a set of distrustful parties, each with a private input, jointly computes a function over their inputs. SMC protocol is established to ensure that each party will get only computation result and their own data will remain unaffected and secure. However, during the execution of the protocol, a party may take one of the following actions in order to get more benefits:

Semi-honest Adversary: one follows the established protocol and follows the rules of that protocol and correctly performs the computation but attempts to analyze others' private inputs to get benefit.

Malicious Adversary: one arbitrarily deviates from the established protocol which leads to the failure of computation so, quality of whole computation result declined.

Collusion: one party colludes with several other parties to expose the private input of another party who doesn't participate in the collusion.

Nanvati and Jinwala model the secret sharing in PPDARM as a repeated game, where a Nash equilibrium is achieved when all parties send their shares and attain a non-collusive

behavior. Based on the game model, they develop punishment policies which aim at getting the maximum possible participants involved in the game so that they can get maximum utilities.

### 6.3.2. RECOMMENDER SYSTEM

. Many online sites asked the user to recommend the product of that company to the others using FIVE star to re-view that product. The recommendation system predict user's preference by analyzing the item ratings provided by users. So, user can protect his private preference by providing false data i.e. by providing false ratings to the product. however, this may cause decline of the quality of the recommendation. Halkidi et al. employ game theory to address the trade-off between privacy and recommendation. In the proposed game model, users are treated as a players, and the rating data provided to the recommender server are seen as user's strategies. It has been shown that the Nash equilibrium strategy for each user is to declare false rating only for one item, the one that is highly ranked in his private profile and less correlated with items. To find the equilibrium strategy, data exchange between users and the recommender server is modeled as an iterative process. At each iteration, by using the ratings provided by other users at previous iteration, each user computes a rating vector that can maximize the preservation of his privacy, with respect to a constraint of the recommendation quality. Then the user declare this rating vector to the recommender server. After several iterations, the process converges to a Nash equilibrium.

### 6.3.3. LINEAR REGRESSION AS A NON-COOPERATIVE GAME

. [19]Data collector receives the data from different Data providers. It may be possible that the data provider may add some noise to protect his data, which affects the accuracy of the model. The interactions among individuals are modeled as a non-cooperative game, where each individual selects the variance level of the noise to minimize his cost. The cost relates to both the privacy loss incurred by the release of data and the accuracy of the estimated linear regression model. It is shown that under appropriate assumptions on privacy and estimation costs.

## 6.4. DATA ANONYMIZATION

Chakravarthy propose a k-anonymity method which utilizes coalitional game theory to achieve a proper privacy level, given the threshold for information loss. This method consider each tuple in the data table as a player in the game theory, and calculates the payoff to each player according to a concept hierarchy tree (CHT) of quasi-identifiers. The equivalent class in the anonymous table is formed by establishing a coalition among different tuples based on their payoffs. Given the affordable information loss, this method can automatically find the most feasible value of k, while traditional methods need to fix up the value of k before the starting anonymization process.

## 6.5. ASSUMPTIONS OF THE GAME MODEL

Most of the proposed approaches adopt the following research paradigm.[20]

Define elements of the game such as Players and their actions and payoffs.

Determine whether game is static or dynamic, with complete information or incomplete information.

Solve the game to find equilibrium.

Analyze the equilibriums to find implications.

Apart from these assumptions, we have to make few assumptions while developing the game model. But, Too many assumptions will reduce game model's accuracy. also, An improper assumptions may comprise game model's accuracy.

## VII. NON-TECHNICAL SOLUTIONS TO PRIVACY PROTECTION

Apart from the technical solutions to privacy protection, there are non-technical solutions. Legislation on privacy protection has been a prime concern with respect to people. Many countries have established laws to avoid such attacks on individual's private information to regulate acts involving personal information. Despite the many laws and regulations, now a days, the definition of the right to privacy and the boundary of "legitimate" practice on personal data are still vague.

Besides laws and regulation, industry conventions are very important. Agreement between different organizations on how personal data should be received, placed and modified, can help to build a privacy safe environment for data mining applications. Many countries develop different laws and regulations to avoid such a attack on individual's personal information. If someone tries to do so, he/she will have to go through the legal procedure.

## VIII. CONCLUSION

In this paper, We discuss about security concerns and privacy preserving techniques of each user such as Data Provider, Data Collector, Data Miner and Decision Maker.

For Data Provider, can secure his data by three ways: he can limit the access to his online activities or data by using anti-tracking extensions, advertisement and script blockers or by using encryption tools to encrypt emails between two private parties. Data Provider can also demand for high price to disclose his private data with others. Nowadays, whatever you try, but the hackers can get your secure information. So, Data provider can provide false data to misguide such a hackers. Using sockpuppet, Data provider can make different sockpuppets. Data provider can use MaskMe to mask his sensitive information.

For Data Collector, He receives data from Data provider and sends that data to the Data Miner. Before sending this data to the Data miner, Data collector has to check whether data contains any private information or not. data Collector has to develop different attack models to check whether data contains any private information about data provider

For Data Miner, He has to retrieve the important data using different data mining techniques. so, The primary concern of data miner is how to prevent sensitive information from appearing in the mining results. To perform a privacy preserving data mining, the data miner usually needs to modify the data he got from the data collector. As a result, the decline of data utility is inevitable. Similar to data collector, the data miner also faces the privacy utility trade-off problem. By using different algorithm techniques such as Decision tree, Support Vector Machine, Naive Bayesian Techniques Data collector modifies data and sends it to the Decision maker.

For Decision Maker, the privacy concerns are twofold: how to prevent unwanted disclosure of sensitive mining results, and how to evaluate the credibility of the received mining results.

## REFERENCES

- [1] Jiawei Han and Micheline Kamber, Data Mining: Concepts and Techniques (2nd edition).
- [2] Rakesh Agrawal and Ramkrishnan Srikant, "Privacy-Preserving Data Mining" - 'Privacy Preserving Methods'.
- [3] Abinaya.B, Dr.Amitha .T, "Identifying Users from an Anonymous Social Network" in (IJETCSE) Vol-12, Issue-4, Feb-2015.
- [4] D.Ganesh, Dr. S.K.Mahendran, "Security in Data mining using User role based Methodology" [Page 83] in IJETCCT Vol-1, Issue-2, Nov, 2014.
- [5] BENJAMIN C. M. FUNG, KE WANG, RUI CHEN and PHILIP S. YU, "Privacy-Preserving Data Publishing: A Survey of Recent Developments" in ACM Computing Surveys, Vol. 42, No. 4, Article 14, Publication date: June 2010.
- [6] Bin Zhou, Jian Pei, WoShun Luk, "A Brief Survey on Anonymization Techniques for Privacy Preserving Publishing of Social Network Data".
- [7] Wei Peng, Feng Li, Xukai Zou and Jie Wu, "Seed and Grow: An Attack Against Anonymized Social Networks" - Page 3.
- [8] V.Gnanasekar, S.Jayanthi, "Privacy Preservation of Social Network Data against Structural Attack using K-Auto restructure" - Figure 1: (a) Social Network
- [9] (G), (b) Nave Anonymized Network(G), (c) Mapping function (f).
- [10] LATANYA SWEENEY, "k-ANONYMITY: A MODEL FOR PROTECTING PRIVACY" in International Journal on Uncertainty, Fuzziness and Knowledge-based Systems, 10 (5), 2002; 557-570.
- [11] Giorgos Poulis, Spiros Skiadopoulos, Grigoris Loukides, Aris Gkoulalas-Divanis, "Select-Organize-Anonymize: A framework for trajectory data anonymization".
- [12] Mohammad Azam Chhipa, Prof. Lalit Gehlod, "Survey on Association Rule Hiding Approaches" in IJARCSSE Vol-5, Issue-12, Dec, 2015.
- [13] Mr. Tanupriya Choudhury, Mrs. Vasudha Vashisht, Prof. Vivek Kumar, Mr. Himanshu Srivastava, "Data Mining using Decision Tree for Sales Force Optimization", page [517] in IJARCSSE Vol-3, Issue-4, April 2013.
- [14] Jaideep Vaidya, Anirban Basu, Basit Shafiq, Yuan Hong, "Differentially Private Naive Bayes Classification".
- [15] Robert Burbidge, Bernard Buxton, "An Introduction to Support Vector Machines for Data Mining".
- [16] Boris Glavic, Javed Siddique, Periklis Andritsos, Renee J. Miller, "Provenance for Data Mining".
- [17] Miriam J. Metzger, "Making Sense of Credibility on the Web: Models for Evaluating Online Information and Recommendations for Future Research" in
- [18] JOURNAL OF THE AMERICAN SOCIETY FOR INFORMATION SCIENCE AND TECHNOLOGY, 58(13):20782091, 2007.
- [19] Rosa Karimi Adl, Mina Askari, Ken Barker, and Reihaneh Safavi-Naini, "Privacy Consensus in Anonymization Systems Via Game Theory".
- [20] Jisha Jose Panackal and Dr Anitha S Pillai, "Privacy Preserving Data Mining: An Extensive Survey" in Proc. of Int. Conf. on Multimedia Processing, Communication and Info. Tech., MPCIT.

- [21] Stratis Ioannidis and Patrick Loiseau, "Linear Regression as a Non-Cooperative Game".
- [22] MALVINA TEMA, "ASSUMPTIONS OF THE GAME MODEL" in International Relations Quarterly, Vol-5, No.1.(Spring 2014/1).
- [23] LEI XU, CHUNXIAO JIANG, JIAN WANG, JIAN YUAN AND YONG REN,"Information Security in Big Data: Privacy and Data Mining"
- [24] R K Adl, M Askari, K Barker, and R S-Naini, "Privacy Consensus in Anonymization Systems Via Game Theory".