

Data Linkage and Leakage Detection in Data Mining Using E-Random and S-Random

Mr. SreedharAmbala , Mr.Mamidala Sagar , Ms.K.Lavanya

^{1,2}Asst.Professor. Siddhartha Institute of Technology and Sciences, Hyderabad.

³Pursuing M.Tech Siddhartha Institute of Technology and Sciences, Hyderabad.

Email: sridhar.Ambala@gmail.com , Email: Mamidala.sagar@gmail.com , Email: lavanya701@gmail.com

Abstract:

A data distributor has given sensitive data to a set of supposedly trusted agents (third parties). Some of the data has leaked and found in an unauthorized place (e.g., on the web or somebody's laptop). The distributor should assess the likelihood of the leaked data came from one or more agents, as opposed to having independently gathered by others. We propose data allocation strategies (across the agents) that improve the probability of identifying leakages and linkage the data. These methods don't rely on alterations of the released data (e.g., watermarks). In some cases we can also inject "realistic but fake" data records to further improve our chances of detecting leakage and identifying the guilty party. The main objective of this paper E-Random algorithm and S-Random algorithm is used by adding additional information (i.e. Fake data) to the original data to detect the abnormal access in database records very effectively.

Keywords: E -Random, S-Random, M -score, guilty agent, Distributer.

I. INTRODUCTION:

Ever since in Data Linkage techniques, there is a dramatic challenge for identifying the fraud detection information illegally over the internet (i.e. license key, Applications, Bank Transaction Details etc..) to different unauthorized users. It makes critical problem for handling confidential and many other

secret information. In existing system, many techniques can be used to for detecting illegal distributing users. In [4] One-Class Clustering Tree (OCCT) method fraud detection is used to detect the abnormal access to database records that might indicate a data misuse. The main disadvantage of this method is that it takes more time to detect data misuse. In proposed system, E-Random algorithm and S-Random algorithm is used by adding additional information (i.e. Fake data) to the original data. This method is used to detect the abnormal access in database records very effectively and it takes less time to detect the illegal file distribution compared to OCCT.

Our aim is to detect when the distributor's sensitive data have been leaked by agents, and if possible to identify the particular agent that leaked the data. Perturbation is most useful technique where the data has modified and made "less sensitive" before being watermarking is used to handle the leakage detection. We announce the need for watermarking database relations to deter their piracy, and identify the unique characteristics of relational data which pose new challenges for watermarking, and provide desirable properties of watermarking system for relational data. A watermark can be applied to any of the database relation having attributes which are such that changes in a few of their values do not affect the applications. Watermarking means a unique code is embedded in each distributed copy If that copy is later discovered

in the hands of an unauthorized party, the leaker can be identified. Furthermore, watermarks can sometimes be destroyed if the data recipient is malicious.

2. LITERATURE SURVEY:

The Paper [1] used to represent a one class clustering tree approach which performs one-to-many record linkage. This method is based on a one class decision tree model which sums up the knowledge of which records to be linked together.

To summarize, this method allows performing one-to-many linkage while the traditional methods followed one-to-one linkage. Then, we have used a one-class approach which results in matching pairs are only required in the training set, as more number of non-matching (negative) pairs will confuse the model and it will lead to a less accurate model. Another advantage of using OCCT model is that the solution can be easily transformed to rules.

The paper [2] proposes the vast amount of information on the World Wide Web, a typical short query of 1-3 words submitted to a search engine usually get a result list of tens of thousands web pages, while only a tiny part of these pages is useful for users. Web data cleansing with key resource selection based on K-means clustering makes it possible to get better retrieval performance with fewer pages indexed.

This paper [3] has been to demonstrate that the technology for building call trees from examples is fairly sturdy. Current industrial systems square measure powerful tools that have achieved noteworthy successes. The groundwork has been fin

ished advances that may allow such tools to deal even with shouting, incomplete information typical of advanced real-world applications. Work is constant at many centers to enhance the performance of the underlying algorithms. 2 samples of modern analysis provide some tips to the directions during which the sector is moving. Where as call trees generated by the higher than systems square measure quick to execute and might be terribly correct, they leave a lot of to be desired as representations of knowledge.

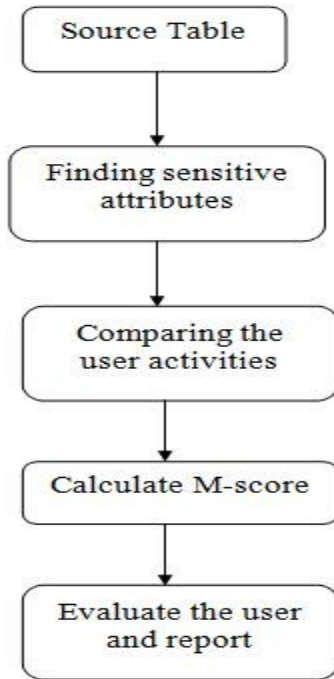
3. SYSTEM ARCHITECTURE:

Fig 1 and Fig 1.1 shows overall flow of our proposed system. A knowledge distributor has given sensitive data to a collection of purportedly sure agents (third parties). A number of the information is leaked and located in AN unauthorized place (e.g., on

the net or somebody's laptop). The distributor should assess the chance that the leaked knowledge came from one or additional agents, as critical having been severally gathered by different suggests that. We have a tendency to propose knowledge allocation ways (across the agents) that improve the chance of distinguishing leakages. These ways don't have confidence alterations of the discharged knowledge (e.g., watermarks). In some cases we are able to additionally inject "realistic however

fake" knowledge records to additional improve our probabilities of police work discharge and distinguishing the guilty one. Our goal is to notice once the distributor's sensitive knowledge has been leaked by

Fig 1.1 system flow diagram agents, and if doable to spot the agent that leaked the information.



The main focus of our project is that the knowledge allocation drawback as however will the distributor “intelligently” provides knowledge to agents so as to enhance the probabilities of police work a guilty agent.

Fake objects are objects generated by the distributor so as to extend the probabilities of police work agents that leak knowledge.

The distributor could also be able to add faux objects to the distributed knowledge so as to enhance his effectiveness in police work guilty agents. Our use of faux objects is impressed by the employment “trace”

4. EXISTING SYSTEM:

- One-Class Clustering Tree (OCCT) method[4] is used for fraud detection in data linkage.
- Datasets are constructed in tree structures.
- Datasets is splitted into number of subsets to

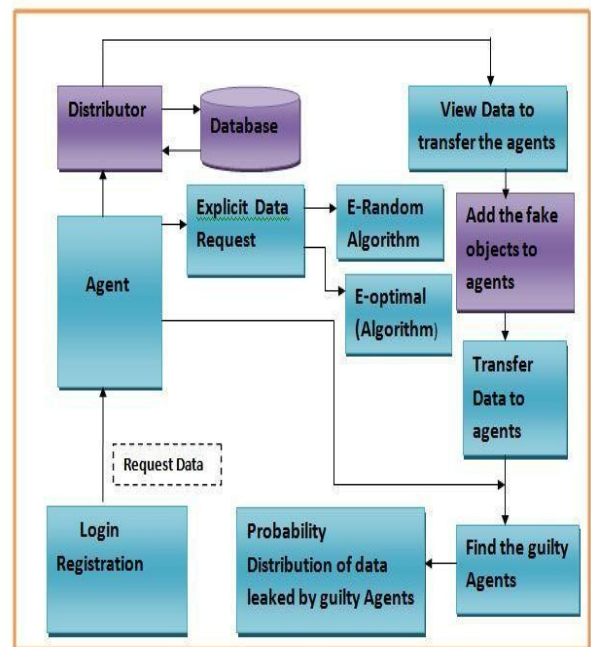
perform comparison using CGJ (Coarse-Grained Jaccard).

4.1 DISADVANTAGES

- Tree is constructed in small amount of nodes.
- It is complex to process large number of nodes.
- It takes more time to compare the data linkage in large number of subsets.

5. PROPOSED SYSTEM:

- In proposed system, we present a new concept, Misuseability Weight, for estimating the risk emanating from data exposed to insiders.
- This concept focuses on assigning a score



that represents the sensitivity level of the

records in mailing lists.

data exposed to the user and by that predicts the ability of the user to maliciously exploit this data.

- E-Random algorithm and S-Random algorithm is used by adding additional

information (i.e. Fake data) to the original data and randomly allocating fake objects to the agent.

5.1 ADVANTAGES

- It is more flexible to process large number of nodes.
- It takes less time to compare the data linkage in large number of subsets.
- It is easy to detect and calculate the abnormal access in the database.

6. SYSTEM IMPLEMENTATION:

6.1. DATA ALLOCATION PROBLEM

A. Fake objects:

Fake objects are objects generated by the distributor in order to increase the chances of detecting agents that leak data. The distributor may be add fake objects to the distributed data in order to improve his effectiveness in detecting guilty agents. Our use of fake objects is inspired by the use of

“trace” records in mailing lists. The idea of perturbing data to detect leakage is not new, e.g., [1]. However, in most cases, individual objects are Perturbed, e.g., by adding random noise to sensitive salaries, or adding fake elements.

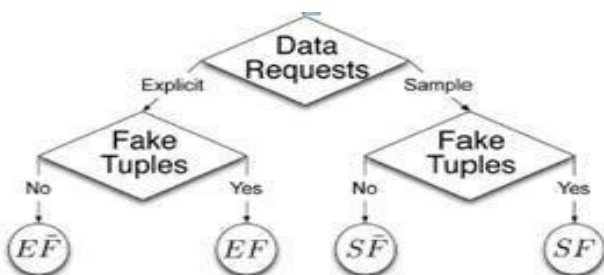


Fig 2. Leakage problem Instances

In some applications, fake objects may cause fewer problems than perturbing real objects. Creation. The creation of fake but real-looking objects is a nontrivial problem whose thorough investigation is beyond the scope of this paper. Here, we model the creation of a fake object for agent U_i as a black box function $CREATEFAKEOBJECT(R_i, F_i, condi)$ that takes as input the set of all objects R_i , the subset of fake objects F_i that U_i has received so far, and $condi$, and returns a new fake object. This function needs $condi$ to produce a valid object that

satisfies U_i 's condition. Set R_i is needed as input so that the created fake object is not only valid but also indistinguishable from other real objects.

Although we do not deal with the implementation of $CREATEFAKEOBJECT()$, we note that there are two main design options. The function can either produce a fake object on demand every time it is called or it can return an appropriate object from a pool of objects created in advance. We are using the following strategies to add the fake object to finding guilty agent.

B. optimization problem

The Optimization Module is the distributor's data allocation to agents has one constraint and one objective. The distributor's constraint is to satisfy agents' requests, by providing them with the number of objects they request or with all available objects that satisfy their conditions. His objective is to be able to detect an agent who leaks any portion of his data. We consider the constraint as strict. The distributor may not deny serving an agent request and may not provide agents with different perturbed versions of the same objects as in [1]. We consider fake object distribution as the only possible constrain

relaxation. Our detection objective is ideal and intractable.

A data leakage is the intentional or unintentional release of secure information to an untrusted environment. Other terms for this phenomenon include unintentional information disclosure, databreach and also data spill. Incidents range from concerted attack by black hats with the backing of organized crime or national governments to careless disposal of used computer equipment or data storage media.

C. Objective Approximation

We can approximate the objective of (2) with (3) that does not depend on agents' guilt probabilities, and therefore, on \mathbf{p} :

$$\text{Maximise } (\dots |R_i \cap R_j| / |R_i|, \dots) \quad i \neq j \quad \dots \dots \dots (3)$$

This approximation is valid if minimizing the relative overlap

$$|R_i \cap R_j| / |R_i| \text{ maximizes } \Delta(i, j).$$

Therefore, we can scalarize the problem objective by assigning the same weights to all vector objectives.

$$\text{Maximize } \sum_{j=1}^n 1/|R_i| |R_i \cap R_j| \quad \dots \dots \dots (4a)$$

(Over $R_1 \dots R_n$)

$$\text{Maximize } \max |R_i \cap R_j| / |R_i| \quad \dots \dots \dots (4b)$$

(Over $R_1 \dots R_n$)

Both scalar optimization problems yield the optimal solution of the problem of (3), if

such solution exists. If there is no global optimal solution, the sum-objective yields

the Pareto optimal solution that allows the distributor solution that guarantees that the distributor will detect the guilty agent with certain confidence in the worst case. Such guarantee may adversely impact the average performance of the distribution.

6.2. ALLOCATION STRATEGIES

The main focus of our project is the data allocation problem as how can the distributor "intelligently" give data to agents in order to improve the chances of detecting a guilty agent.

6.3. ALGORITHM USED

6.3.1 Explicit and Sample Random Algorithm

A. Explicit Data Requests

In problems of class EF, the distributor is not allowed to add fake objects to the distributed data.

So, the data allocation is fully defined by the agents' data requests. Therefore, there is nothing to optimize. In EF problems, The distributor cannot remove or alter the R_1 or R_2 data to decrease the overlap $R_1 \cap R_2$. However, say that the distributor can create one fake object ($B=1$) and both agents can receive

one fake object ($b_1= b_2= 1$). In this case, the distributor can add one fake object to either R_1 or R_2 to increase the corresponding denominator of the summation term. Assume that the distributor creates a fake object f and he gives it to agent R_1 . Agent U_1 has now $R_1= \{t_1, t_2, f\}$ and $F_1=\{f\}$ and the value of the sum-objective decreases to $1/3+1/1=1.33<1.5$.

Algorithm 1. Allocation for Explicit Data Requests (EF)

Input: $R_1 \dots R_n, cond_1; \dots; cond_n, b_1, \dots, b_n, B$
 Output: $R_1 \dots R_n, F_1 \dots F_n$
 1: $R \leftarrow \phi$ Agents that can receive fake objects
 2: for $i=1 \dots n$ do

```

3: if  $b_i > 0$  then
4:  $R \leftarrow R \cup \{i\}$ 
5:  $F_i \leftarrow \phi$ 
6: while  $B > 0$  do
7:  $i \leftarrow \text{SELECTAGENT}(R, R_1, \dots, R_n)$ 
8:  $f \leftarrow \text{CREATEFAKEOBJECT}(R_i, F_i, \text{condi})$ 
9:  $R_i \leftarrow R_i \cup \{f\}$ 
10:  $F_i \leftarrow F_i \cup \{f\}$ 
11:  $b_i \leftarrow b_i - 1$ 
12: if  $b_i = 0$  then
13:  $R \leftarrow R - \{R_i\}$ 
14:  $B \leftarrow B - 1$ 

```

Algorithm 2. Agent Selection for e-random

```

1: function SELECTAGENT (R, R1; . . . , Rn)
2:  $i \leftarrow$  select at random an agent from R
3: return i

```

In lines 1-5, Algorithm 1 finds agents that are eligible to receiving fake objects in $O(n)$ time. Then, in the main loop in lines 6-14, the algorithm creates one fake object in every iteration and allocates it to random agent. The main loop takes $O(B)$ time. Hence, the running time of the algorithm is

$O(n + B)$. If $B \geq \sum_{i=1}^n b_i$, the algorithm minimizes every term of the objective summation by adding the maximum number b_i of fake objects to every set R_i , yielding the optimal solution. Otherwise, if $B < \sum_{i=1}^n b_i$ (as in our example where $B=1 < b_1+b_2=2$), the algorithm just selects at random the agents that are provided with fake objects. We return back to our example and see how the objective would change if the distributor adds fake object f to R_2 instead of R_1 . In this case, the sum-objective would be $1/2 + 1/2 = 1 < 1.33$. The reason why we got a greater improvement is that the addition of a fake object to R_2 has greater impact on the corresponding summation terms, since

$$1/|R_1| - 1/|R_1| + 1 = 1/6 < 1/|R_2| - 1/|R_2| + 1 = 1/2.$$

The left-hand side of the inequality corresponds to the objective improvement after the addition of a fake

object to R_1 and the right-hand side to R_2 .

Algorithm 3. Agent Selection for e-optimal

```

1: function SELECTAGENT (R, R1; . . . , Rn)
2:  $i \leftarrow \text{argmax} (1/|R_i| - 1/|R_i| + 1) \sum |R_i' \cap R_j|$ 
3: return i

```

Algorithm 3 makes a greedy choice by selecting the agent that will yield the greatest improvement in the sum objective

the cost of this greedy choice is $O(n^2)$ in every iteration. The overall running time of e-optimal is $O(n + n^2B) = O(n^2B)$. Theorem 2 shows that this greedy approach finds an optimal distribution with respect to both optimization objectives defined in (4).

Theorem 2. Algorithm 3 e-optimal yields an object allocation that minimizes both sum and max-objective in problem instances of class EF.

B. Sample Data Requests:

With sample data requests, each agent U_i may receive any T subset out of $(|T|)$ different ones. Hence, there are $i=1(|T|)$ different object allocations. In every allocation, the distributor can permute T objects and keep the same chances of guilty agent detection. The reason is that the guilt

probability depends only on which agents have received the leaked objects and not on the identity of the leaked objects. Therefore, from the distributor's perspective, different

allocations. The distributor's problem is to

pick one out so that he optimizes his objective. We formulate the problem as a nonconvex QIP that is NP-hard.

Note that the distributor can increase the number of possible allocations by adding fake objects (and increasing $|T|$) but the

problem is essentially the same. So, in the rest of this section, we will only deal with problems of class SF, but our algorithms are applicable to SF problems as well.

C. Random: An object allocation that satisfies requests and ignores the distributor's objective is to

We present S-random in two parts: Algorithm 4 is a general allocation algorithm that is used by other algorithms in this section. In line 6 of Algorithm 4, there is a call to function SELECTOBJECT () whose implementation differentiates algorithms that rely on Algorithm 4. Algorithm 5 shows function SELECTOBJECT () for s-random. Algorithm 4. Allocation for Sample Data Requests (SF)

Input: $m_1, \dots, m_n, |T|$. Assuming $m_i < |T|$

Output: R_1, \dots, R_n

1: $a \leftarrow 0|T|$. $a[k]$: number of agents who have received object tk

2: $R_1 \leftarrow \emptyset, \dots, R_n \leftarrow \emptyset$

3: $\text{remaining} \leftarrow \sum_{i=1}^n m_i$

4: while $\text{remaining} > 0$ do

 5: for all $i = 1 \dots n$: $|R_i| < m_i$ do

 6: $k \leftarrow \text{SELECTOBJECT}(i, R_i)$ May also use additional parameters

 7: $R_i \leftarrow R_i \cup \{tk\}$

 8: $a[k] \leftarrow a[k] + 1$

 9: $\text{remaining} \leftarrow \text{remaining} - 1$

Algorithm 5. Object Selection for s-random

1: function SELECTOBJECT (i, R_i)

2: $k \leftarrow$ select at random an element from set $\{K' / tk \in R_i\}$

3: return k

7. CONCLUSION:

In this paper We have shown that it's attainable to assess the likelihood that associate agent is chargeable for a leak, based on the overlap of his knowledge with the leaked knowledge and therefore the data of different agents, and supported the likelihood that objects are often "guessed" by give each agent U_i a randomly selected subset of T of size

m_i . We denote this algorithm by S-random and we use it as our baseline. different suggests that. We are proposing our increased approach for detection the guilty agents. during this technique we have a tendency to use the linguistics inference model that represents the likelihood of possible colluding attacks from any agents to the different knowledge allocation methods. we present a new concept, Misuseability Weight, for estimating the risk emanating from data exposed to insiders. In this we are calculating Misuseability weight and probability value based on how much data misused or leakage data by agents. so that distributor can easily estimate misused data and guilty agent.

In future the extension of our allocation methods will handle agent requests in an internet fashion (the presented methods assume that there's a hard and fast set of agents with requests proverbial in advance) are often implemented.

REFERENCES

- [1] Sunandhini.S and Suguna.M," Improved One-to-Many Record Lin kage using One Class Clustering Tree", (IJCA)(0975- 8887) ,International Conference on Simulations in Computing Nexus, ICSCN-2014.
- [2] YiqunLIU,MinZHANG, Canhui WANG, Shaoping MA,"Learning-based Web Data Cleansing for Information Retrieval" , Journal of Computational Information Systems 1:2 (2005) 203-213,june 200
- [3] J.R. QUINLAN,"Induction of Decision Trees", Machine Learning 1: 81-106, 1986
- [4] Ma'ayan Dror,Asaf Shabtai,Lior Rokach and Yuval Elovici, "OOCT:A One class clustering Tree for implementing One-to- many Data lin kage" IEEE,vol.26,ZNo.3, March 2014.
- [5] B.Mungamuru and H.Garcia Molina,"privacy,preservation and Performance: The 3 p's of Distributed Data Management," technical report, Stanford univ.,2008
- [6] R. Agrawal and J. Kiernan, "Watermarking Relational Databases,"Proc. 28th Int'l Conf. Very Large Data Bases (VLDB '02), VLDB Endowment, pp. 155-166, 2002.
- [7] L. Sweeney, "Achieving K-Anonymity Privacy Protection Using Generalization and Suppression," <http://en.scientificcommons.org/43196131>, 2002.