

Instructive of Ooze Information

Sudheer Kumar Kotha¹ and CH.S.V.V.S.N.Murthy²

¹ Jawaharlal Nehru Technical University, Kakinada.

Received: 8 December 2011 Accepted: 2 January 2012 Published: 15 January 2012

Abstract

We study the following problem: A data distributor has given sensitive data to a set of supposedly trusted agents (third parties). Some of the data are leaked and bring into being in an unconstitutional place (e.g., on the web or somebody's laptop). The distributor must evaluate the likelihood that the leaked data came from one or more agents, as opposed to having been independently gathered by other means. We propose data distribution strategies (across the agents) that improve the likelihood of identifying leakages. These methods do not rely on alterations of the released data (e.g., watermarks). In some cases, we can also inject "realistic but replica" data records to further improve our chances of detecting leakage and identifying the guilty party. In the course of doing business, sometimes sensitive data must be handed over to supposedly trusted third parties. For example, a hospital may give patient records to Researchers who will devise new treatments. Similarly, a company may have partnerships with other companies that require sharing customer data. Another enterprise may outsource its data processing, so data must be given to various other companies. There always remains a risk of data getting leaked from the agent. Perturbation is a very valuable technique where the data are modified and made "less sensitive" before being handed to agents. For example, one can add random noise to certain attributes, or one can replace exact values by ranges. But this technique requires modification of data. Leakage detection is handled by watermarking, e.g., a unique code is implanted in each distributed copy. If that copy is later discovered in the hands of an unconstitutional party, the leaker can be identified. But again it requires code modification. Watermarks can sometimes be destroyed if the data recipient is malicious.

Index terms— Allocation strategies, data leakage, data privacy, fake records, leakage model.

1 Introduction

In the course of doing business, sometimes sensitive data must be handed over to supposedly trusted third parties. For example, a hospital may give patient records to researchers who will devise new treatments. Similarly, a company may have partnerships with other companies that require sharing customer data. Another enterprise may outsource its data processing, so data must be given to various other companies. We call the owner of the data the distributor and the supposedly trusted third parties the agents. Our goal is to detect when the distributor's sensitive data has been leaked by agents, and if possible to identify the agent that leaked the data.

We consider applications where the original sensitive data cannot be perturbed. Perturbation is a very useful technique where the data is modified and made "less sensitive" before being handed to agents. For example, one can add random noise to certain attributes, or one can replace exact values by ranges [18]. However, in some cases it is important not to alter the original distributor's data. For example, if an outsourcer is doing our payroll, he must have the exact salary and customer bank account numbers. If medical researchers will be treating patients (as opposed to simply computing statistics), they may need accurate data for the patients.

Traditionally, leakage detection is handled by watermarking, e.g., 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. Watermarks can be very useful in some cases, but again, involve some modification of the original data. Furthermore, watermarks can sometimes be destroyed if the data recipient is malicious. In this paper we study unobtrusive (Not attracting unnecessary attention) techniques for detecting leakage of a set of objects or records. Specifically, we study the following scenario: After giving a set of objects to agents, the distributor discovers some of those same objects in an unauthorized place. (For example, the data may be found on a web site, or may be obtained through a legal discovery process).

At this point the distributor can assess the likelihood that the leaked data came from one or more agents, as opposed to having been independently gathered by other means. If the distributor sees "enough evidence" that an agent leaked data, he may stop doing business with him, or may initiate legal proceedings.

In this paper we develop a model for assessing the "guilt" of agents. We also present algorithms for distributing Objects to agents, in a way that improves our chances of identifying a leaker. Finally, we also consider the option of adding fake" objects to the distributed set. Such objects do not correspond to real entities but appear realistic to the agents. In a sense, the fake objects acts as a type of watermark for the entire set, without modifying any individual members. If it turns out an agent was given one or more fake objects that were leaked, then the distributor can be more confident that agent was guilty.

2 Problem setup and notation

Entities and Agents: A distributor owns a set $T=\{t_1, t_m\}$ of valuable data objects. The distributor wants to share some of the objects with a set of agents U_1, U_2, \dots, U_n , but does not wish the objects be leaked to other third parties. The objects in T could be of any type and size, e.g., they could be tuples in a relation, or relations in a database. An agent U_i receives a subset of objects R_i request or an explicit request:

Sample request $R_i = \text{SAMPLE}(T, m_i)$: Any subset of M_i records from T can be given to U_i .

Explicit request $R_i = \text{EXPLICIT}(T, \text{condi})$: Agent U_i receives all the T objects that satisfy condi .

Type of data leakage: In order to implement the appropriate protective measures, we must first understand what we are protecting. Based on publicly disclosed Data Leakage breaches, the type of data leaked is broken down as follows: Type of information leaked Percentage Confidential information -15% Intellectual property -4% Customer data -73% Health records -8% Guilty Agents: Suppose that after giving objects to agents, the distributor discovers that a set S leaked. This means that some third party called the target has been caught in possession of S . For example, this target may be displaying S on its web site, or perhaps as part of a legal discovery process, the target turned over S to the distributor.

Since the agents U_1, U_n has some of the data, it is reasonable to suspect them leaking the data. However, the agents can argue that they are innocent, and that the S data was obtained by the target through other means.

For example, say one of the objects in S represents a customer X . Perhaps X is also a customer of some other company, and that company provided the data to the target. Or perhaps X can be reconstructed from various publicly Available sources on the web.

Our goal is to estimate the likelihood that the leaked data came from the agents as opposed to other sources. Intuitively, the more data in S , the harder it is for the agents to argue they did not leak anything. Similarly, the rarer" the objects, the harder it is to argue that the target obtained them through other means. For instance, if one of the S objects was only given to agent U_1 , while the other objects were given to all agents, we may suspect U_1 more. The model we present next captures this intuition.

We say an agent U_i is guilty and if it contributes one or more objects to the target. We denote the event that agent U_i is guilty as G_i and the event that agent U_i is guilty for a given leaked set S as $G_i|S$. Our next step is to estimate $\Pr \{G_i|S\}$, i.e., the probability that agent U_i is guilty given evidence S .

3 III.

4 Related works

The guilt detection approach we present is related to the data provenance problem [3]: (whether it is genuine or not problem) tracing the lineage of S objects implies essentially the detection of the guilty agents. Tutorial [4] provides a good overview on the research conducted in this field. Suggested solutions are domain specific, such as lineage tracing for data warehouses [5], and assume some prior knowledge on the way a data view is created out of data sources.

Our problem formulation with objects and sets is more general and simplifies lineage tracing, since we do not consider any data transformation from R_i sets to S . As far as the data allocation strategies are concerned, our work is mostly relevant to watermarking that is used as a means of establishing original ownership of distributed objects. Watermarks were initially used in images [16], video [8] and audio data [6] whose digital representation includes considerable redundancy. Recently, [1], [17], [10], [7] and other works have also studied marks insertion to relational data. Our approach and watermarking are similar in the sense of providing agents with some kind of receiveridentifying information.

However, by its very nature, a watermark modifies the item being watermarked. If the object to be watermarked cannot be modified then a watermark cannot be inserted. In such cases methods that attach watermarks to the

distributed data are not applicable. Finally, there are also lots of other works on mechanisms that allow only authorized users to access sensitive data through access control policies [9], [2]. Such approaches prevent in some sense data leakage by sharing information only with trusted parties. However, these policies are restrictive and may make it impossible to satisfy agents' requests.

5 IV.

6 Agent guilt model

To compute this $\Pr\{Gi|S\}$, we need an estimate for the probability that values in S can be guessed" by the target. For instance, say some of the objects in S are emails of individuals. We can conduct an experiment and ask a person with approximately the expertise and resources of the target to find the email of say 100 individuals. If this person can find say 90 emails, then we can reasonably guess that the probability of finding one email is 0.9. On the other hand, if the objects in question are bank account numbers, the person may only discover say 20, leading to an estimate of 0.2. Probability pt is analogous to the probabilities used in designing fault-tolerant systems. That is, to estimate how likely it is that a system will be operational throughout a given period, we need the probabilities that individual components will or will not fail. A component failure in our case is the event that the target guesses an object of S . while we use the probability of guessing to identify agents that have leaked information.

The component failure probabilities are estimated based on experiments, just as we propose to estimate the pt 's. Similarly, the component probabilities are usually conservative estimates, rather than exact numbers. For example, say we use a component failure probability that is higher than the actual probability, and we design our system to provide a desired high level of reliability. Then we will know that the actual system will have at least that level of reliability, but possibly higher. In the same way, if we use pt 's that are higher than the true values, we will know that the agents will be guilty with at least the computed probabilities.

There are $T=\{t1,t2,t3\}; R1=\{t1,t2\}; R2=\{t2,t3\}; S=\{t1,t2,t3\}$

In this case, all three of the distributor's objects have been leaked and appear in S . Let us first consider how the target may have obtained object $t1$, which was given to both agents. The target either guessed $t1$ or one of $U1$ or $U2$ leaked it. We know that the probability of the former event is p , so assuming that probability that each of the two agents leaked $t1$ is the same we have the following cases:

? The target guessed $t1$ with probability p ; ? Agent $U1$ leaked $t1$ to S with probability $(1 - p)/2$; ? Agent $U2$ leaked $t1$ to S with probability $(1 - p)/2$; Similarly, we find that agent $U1$ leaked $t2$ to S with Probability $1 - p$ since he is the only agent that has $t2$. Given these values, the probability that agent $U1$ is not Guilty, namely that $U1$ did not leak either object is: $(1 - (1 - p)/2) - (1 - (1 - p))$; (1)

And the probability that $U1$ is guilty is: $1 - \Pr\{G1\}$ Note that if did not hold, our analysis would be more complex because we would need to consider joint events, e.g., the target guesses $t1$ and at the same time one or two agents leak the value. In our simplified analysis we say that an agent is not guilty when the object can be guessed, regardless of whether the agent leaked the value. Since we are "not counting" instances when an agent leaks information, the Simplified analysis yields conservative values (smaller Probabilities).

To simplify the formulas that we present in the rest of the paper, we assume that all T objects have the same pt , which we call p . Our equations can be easily generalized to diverse pt 's though they become cumbersome to display. Next, we make two assumptions regarding the relationship among the various leakage events. The first assumption simply states that an agent's decision to leak an object is not related to other objects. In [14] we study a scenario where the actions for different objects are related, and we study how our results are impacted by the different independence assumptions.

7 The term

provenance" in this assumption statement refers to the sources of a value t that appears in the leaked set. The Source can be any of the agents who have t in their sets or the target itself (guessing). To simply our formulas, the following assumption states that join events have a negligible probability. As we argue in the example below, this assumption gives us more conservative estimates for the guilt of agents, which is consistent with our goals.

A single agent leaked t from its own set. The target guessed (or obtained through other means) t without the help of any of the n agents.

In other words, for all target guesses t and the events that agents leaks objects t are disjoint. Before we present the general formula for computing the probability $\Pr\{Gi|S\}$ that an agent Ui is guilty, we provide a simple example. Assume that the distributor set T , the agent sets R 's and the target set S are: $T = \{t1, t2, t3\}$, $R1 = \{t1, t2\}$, $R2 = \{t1, t3\}$, $S = \{t1, t2, t3\}$.

In this case, all three of the distributor's objects have been leaked and appear in S . Let us first consider how the target may have obtained object $t1$, which was given to both agents. From Assumption 2, the target either guessed $t1$ or one of $U1$ or $U2$ leaked it. We know that the probability of the former event is p , so assuming that probability that each of the two agents leaked $t1$ is the same we have the following cases:

The target guessed $t1$ with probability p ; Agent $U1$ leaked $t1$ to S with probability $(1 - p)/2$; Agent $U2$ leaked $t1$ to S with probability $(1 - p)/2$;

Similarly, we find that agent U_1 leaked t_2 to S with probability $1 - p$ since he is the only agent that has t_2 . Given values, the probability U_1 is guilty is:

(2) Note that if Assumption 2 did not hold, our analysis would be more complex because we would need to consider joint events, e.g., the target guesses t_1 and at the same time one or two agents leak the value. Since we are not counting instances when an agent leaks information, the simplified analysis yields conservative values (smaller probabilities). In the general case (with our assumptions), to find the probability that an agent U_i is guilty given a set S , first we compute the probability that he leaks a single object t to S . To compute this we define the set of agents Assumption 1: For all $t, t' \in S$ such that $t \neq t'$ the provenance of t is independent of the provenance of t' .

8 "

Assumption 2: An Object $t \in S$ can only be obtained by the target in one of the two ways as follows: $\exists U_i \in R$ the event that the $t \in S$, U_i ($i=1, \dots, n$) $\neq \emptyset$ that have t in their data sets. Then using Assumption 2 and known probability p , we have: $\Pr \{\text{some agent leaked } t \text{ to } S\} = 1 - p^n$.

(Assuming that all agents that belong to V_t can leak t to S with equal probability and using Assumption 2 we obtain: otherwise (4) Given that agent U_i is guilty if he leaks at least one value to S , with Assumption 1 and Equation ?? we can compute the probability $\Pr \{G_i | S\}$ that agent U_i is guilty:

(5) Fake Objects: The distributor may be able to add fake objects to the distributed data in order to improve his effectiveness in detecting guilty agents. However, fake objects may impact the correctness of what agents do, so they may not always be allowable. Perturbed, e.g., by adding random noise to sensitive salaries, or adding a watermark to an image. In our case, we are perturbing the set of distributor objects by adding fake elements. In some applications, fake objects may cause fewer problems than perturbing real objects.

V .

9 Future enhancements

In this paper we are using multiple agents, for the purpose of security at the same time we are creating database for separate user, so through this we are strictly find out who is leaked information in internet. ¹

¹© 2012 Global Journals Inc. (US)

188 [Sweeney] *Achieving K-Anonymity Privacy Protection Using Generalization and Suppression*, L Sweeney .
189 <http://en.scientificcommons>

190 [Bonatti et al. ()] ‘An Algebra for Composing Access Control Policies’. P Bonatti , S D C Di Vimercati , P
191 Samarati . *ACM Trans. Information and System Security* 2002. 5 (1) p. .

192 [Guo et al. ()] *An Improved Algorithm to Watermark Numeric Relational Data*, F Guo , J Wang , Z Zhang , X
193 Ye , D Li . 2006. Springer. p. . (Information Security Applications)

194 [Murty ()] ‘Counting the Integer Solutions of a Linear Equation with Unit Coefficients’. V N Murty . *Math.*
195 *Magazine* 1981. 54 (2) p. .

196 [Papadimitriou and Garcia-Molina ()] *Data Leakage Detection*, P Papadimitriou , H Garcia-Molina . 2008.
197 Stanford Univ. (technical report)

198 [Czerwinski et al. ()] *Digital Music Distribution and Audio Watermarking*, S Czerwinski , R Fromm , T Hodes .
199 <http://www.scientificcommons.org/43025658> 2007.

200 [Li et al. (2005)] ‘Fingerprinting Relational Databases: Schemes and Specialties’. Y Li , V Swarup , S Jajodia .
201 *IEEE Trans. Dependable and Secure Computing* Jan.-Mar. 2005. 2 (1) p. .

202 [Jajodia et al. ()] ‘Flexible Support for Multiple Access Control Policies’. S Jajodia , P Samarati , M L Sapino ,
203 V S Subrahmanian . *ACM Trans. Database Systems* 2001. 26 (2) p. .

204 [Cui and Widom ()] ‘Lineage Tracing for General Data Warehouse Transformations’. Y Cui , J Widom . *The*
205 *VLDB J* 2003. 12 p. .

206 [Mungamuru and Garcia-Molina ()] *Privacy, Preservation and Performance: The 3 P’s of Distributed Data*
207 *Management*, B Mungamuru , H Garcia-Molina . 2008. Stanford Univ. (technical report)

208 [Buneman and Tan ()] ‘Provenance in Databases’. P Buneman , W.-C Tan . *Proc. ACM SIGMOD*, (ACM
209 SIGMOD) 2007. p. .

210 [Pardalos and Vavasis ()] ‘Quadratic Programming with One Negative Eigen value Is NP-Hard’. P M Pardalos ,
211 S A Vavasis . *J. Global Optimization* 1991. 1 (1) p. .

212 [Sion et al. ()] ‘Rights Protection for Relational Data’. R Sion , M Atallah , S Prabhakar . *Proc. ACM SIGMOD*,
213 (ACM SIGMOD) 2003. p. .

214 [Nabar et al. ()] ‘Towards Robustness in Query Auditing’. S U Nabar , B Marthi , K Kenthapadi , N Mishra , R
215 Motwani . *Proc. 32nd Int’l Conf. Very Large Data Bases (VLDB ’06), VLDB Endowment*, (32nd Int’l Conf.
216 Very Large Data Bases (VLDB ’06), VLDB Endowment) 2006. p. .

217 [Ruanaidh et al. ()] ‘Watermarking Digital Images for Copyright Protection’. J J K O Ruanaidh , W J Dowling
218 , F M Boland . *IEE Proc. Vision, Signal and Image Processing*, 1996. 143 p. .

219 [Hartung and Girod ()] ‘Watermarking of Uncompressed and Compressed Video’. F Hartung , B Girod . *Signal*
220 *Processing* 1998. 66 (3) p. .

221 [Agrawal and Kiernan ()] ‘Watermarking Relational Databases’. R Agrawal , J Kiernan . *Proc. 28th Int’l Conf.*
222 *Very Large Data Bases (VLDB ’02), VLDB Endowment*, (28th Int’l Conf. Very Large Data Bases (VLDB
223 ’02), VLDB Endowment) 2002. p. .

224 [Buneman et al. (2001)] ‘Why and Where: A Characterization of Data Provenance’. P Buneman , S Khanna ,
225 W C Tan . *Proc. Eighth Int’l Conf. Database Theory (ICDT ’01)*, J V Bussche, V Vianu (ed.) (Eighth Int’l
226 Conf. Database Theory (ICDT ’01) Jan. 2001. p. .