

DATA LEAKAGE DETECTION USING CLOUD COMPUTING

M.Sai Charan Reddy¹ T. Venkata Satya Yaswanth² T. Gopal³ L. Raji⁴ Dr. K.Vijaya⁵

¹²³ Student, Department of Computer Science, R.M.K. Engineering College

⁴ Assistant professor, R.M.K. Engineering College

⁵ Head of Department, Information Technology, R.M.K. Engineering College

-----***-----

Abstract- This paper presents the conception of data leakage, its causes of leakage and different techniques to protect and detect the data leakage. The value of the data is incredible, so it should not be leaked or mishandled. In the world of IT, a huge database is being consumed. This database is shared with several people at a time. But during this distribution of the data, there are huge chances of data exposure, insecurity or alteration. So, to nullify these problems, a data leakage detection system has been proposed. This paper includes a curt idea about data leakage detection and a methodology to detect data leakage persons.

Keywords-Watermarking guilty agent; Explicit data; DLP (Data Leakage Prevention)

1.Introduction

Data leakage is explained as the accidental or unmeant distribution of private or sensitive data to an unlicensed entity. Sensitive data of companies and firms include intellectual property (IP), financial information, patient information, personal credit-card data, and other information hanging on the business and the industry. Besides, in many cases, sensitive data is split among various patrons such as employees working from outside the organizational places (e.g., on laptops), business partners and customers.

In the route of doing business, sometimes sensitive data must be handed over to supposedly trusted third parties. For example, a hospital may give patient records to testers who will formulate new treatments. Identically, a company may have collaboration with other companies that need sharing customer data. Other endeavors may obtain its data processing, so

data must be given to diverse companies. We call the owner of the data the distributor and the purposely trusted third parties as the agents. Our motive is to note when the distributor's sensitive data has been leaked by agents, and if possible to identify the agent that leaked the data.

1.1 How Was Ingress To The Data Gained?

The "How was ingress to the data gained?" attribute extends the "Who created the leak?" attribute.

These attributes are not replaceable, but fairly complementary and the various ways to obtain ingress to sensitive data can be clustered into the following groups.

The categorization by leakage channel is important in order to see how the incidents may be stopped in the future and can be classified as physical or logical.

Physical leakage channel is that physical media (e.g., HDD, laptops, workstations, CD/DVD, USB devices) holding vulnerable information or the document itself was moved away from the organization. This more frequently means that the control over data was lost even before it left the organization.

2. Literature Survey

a) Agent Guilt Model

Assume an agent U_i is guilty if it grants one or more objects to the target. The event that agent U_i is guilty for a given leaked set S denoted by $G_i | S$. Now estimate $Pr \{G_i | S\}$, i.e., the probability that agent G_i is guilty given evidence S .

To compute the $\Pr \{G_i | S\}$, estimate the probability that values in S_{bcean} “guessed” by the target. For instance, say some of the objects in t are emails of individuals. Conduct an experiment and ask a person to find the email of say 50 individuals, the person may only discover say 10, leading to an estimate of 0.2. Assume this estimate as p_t , the probability that object ‘ t ’ can be guessed by the target.

Consider two assumptions regarding the relationship between the various leakage events.

Assumption 1: For all $t, t \in S$ such that $t \neq T$ the provenance of t is independent of the provenance of T .

The term provenance in this assumption statement refers to the source 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.

Assumption 2: An object $t \in S$ can only be obtained by the target in one of two ways.

A single agent U_i leaked t from its own R_i set, or

The target guessed without the help of any of the n agents. To observe the probability that an agent U_i is guilty given a set S , consider the target guessed t_1 with probability p and that agent leaks t_1 to S with probability $1-p$. First calculate the probability that he leaks a single object t to S . To compute this, define the set of agents $V_t = \{U_i | t \in R_i\}$ that have t in their data sets. Then use Assumption 2 and known probability p , We have the, $\Pr \{\text{agent leaked } t \text{ to } S\} = 1 - p$

Assuming that all agents that belong to V_t can leak t to S with equal probability and using Assumption 2 obtain, $\Pr \{U_i \text{ leaked } t \text{ to } S\}$

Given that agent U_i is guilty if he leaks at least one value to S , compute the probability $\Pr \{G_i | S\}$, agent U_i is guilty, $\Pr \{G_i | S\}$

b) Data Allocation Problem

The distributor “intelligently” delivers data to agents to improve the chances of finding a guilty

agent. There are four illustrations of this problem, hanging on the type of data requests made by agents and whether “fake objects” are allowed. An agent makes two types of requests, called sample and explicit. Based on the requests the fakes objects are attached to the data list.

Fake objects are objects created by the distributor that are not in set T . The objects are planned to look like real objects and are distributed to agents together with the T objects, in order to increase the rates of detecting agents that leak data.

c) Optimization Problem

The distributor’s data assignment to agents has one constraint and one objective. The distributor’s curb is to satisfy agents’ requests, by giving them the number of objects they request or with all handy objects that satisfy their conditions. His goal is to be able to find an agent who leaks any portion of his data.

We consider the curb as strict. The distributor may not deny providing an agent request and may not serve agents with different perturbed versions of the same objects. The fake object distribution as the only possible constraint relaxation. The goal is to maximize the chances of detecting a guilty agent that leaks all his data objects.

The $\Pr \{G_i | S = R_i\}$ or simply $\Pr \{G_i | R_i\}$ is the probability that agent U_a is guilty if the distributor discovers a leaked table S that contains all R_i objects. The difference functions $\Delta (i, j)$ is defined as:

$$\Delta (i, j) = \Pr \{G_i | R_i\} - \Pr \{G_j | R_j\}$$

i. Problem Definition

Let the distributor have data requests from n agents. The distributor wants to give tables

R_1, \dots, R_n to agents, U_1, \dots, U_n respectively, so that

- Distribution satisfies agents’ requests; and
- Maximizes the guilt probability differences $\Delta (i, j)$ for all $i, j = 1, \dots, n$ and $i \neq j$.

Assuming that the sets satisfy the agents' requests, we can express the problem as a multi-criterion

ii. Optimization Problem

Maximize $(\dots, \Delta(i, j), \dots) \quad i \neq j$ (Over R_1, \dots, R_n)

The estimation of objective of the above equation does not rely on the agent's probabilities and therefore minimize the relative overlap among the agents as

Minimize $(\dots, (|R_i \cap R_j|) / R_i, \dots) \quad i \neq j$ (over R_1, \dots, R_n)

This guess is valid if minimizing the relative overlap, $(|R_i \cap R_j|) / R_i$ maximizes $\Delta(i, j)$.

3. Allocation Strategies Algorithm

There are two types of strategies algorithms

a) Explicit data Request

In case of the explicit data request with dupes not allowed, the distributor is not allowed to add fake objects to the distributed data. So Data allocation is fully defined by the agent's data request. In case of the explicit data request with fake allowed, the distributor cannot remove or alter the requests R from the agent. However, the distributor can add a fake object. In algorithm for data allocation for the explicit request, the input to this is a set of request R_1, R_2, \dots, R_n from n agents and different conditions for requests. The e-optimal algorithm finds the agents that are eligible to receive fake objects. Then create one fake object in

iteration and allocate it to the agent selected. The optimal algorithm minimizes every term of the objective summation by adding maximum number b_i of fake objects to every set R_i yielding an optimal solution. Algorithm 1: Allocation for Explicit Data Requests (EF) Input: $R_1, \dots, R_n, \text{cond}_1, \dots, \text{cond}_n, b_1, \dots, b_n, B$ Output: $R_1, \dots, R_n, F_1, \dots, F_n$

Step 1: $R \leftarrow \emptyset$, Agents that can receive fake objects

Step 2: for $i = 1, \dots, n$ do

Step 3: if $b_i > 0$ then

Step 4: $R \leftarrow R \cup \{i\}$

Step 5: $F_i \leftarrow \emptyset$;

Step 6: while $B > 0$ do

Step 7: $i \leftarrow \text{SELECTAGENT}(R, R_1, \dots, R_n)$

Step 8: $f \leftarrow \text{CREATEFAKEOBJECT}(R_i, F_i, \text{cond}_i)$

Step 9: $R_i \leftarrow R_i \cup \{i\}$

Step 10: $F_i \leftarrow F_i \cup \{i\}$

Step 11: $b_i \leftarrow b_i - 1$

Step 12: if $b_i = 0$ then

Step 13: $R \leftarrow R \setminus \{R_i\}$

Step 14: $B \leftarrow B - 1$.

Algorithm 2 : Agent Selection for e-random

Step 1: function $\text{SELECTAGENT}(R, R_1, \dots, R_n)$

Step 2: i select at random an agent from R

Step 3: return I

Algorithm 3: Agent selection for e-optimal

Step 1: function SELECTAGENT(R;R1; : : : ;Rn)

Step 2: $i \leftarrow \text{argmax}^{(1/R_i - 1/R_i + 1) \sum_j |R_i \cap R_j|}$

Step 3: return i : $R \in R$

Sample Data Request

With sample data requests, each agent U_i may receive any T from a subset out of $(|T|_{n_i})$ different ones. Hence, there are $\pi_{i=1}^n (|T|_{n_i})$ different 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 there are $\pi_{i=1}^n (|T|_{n_i}) / |T|$ different allocations. An object allocation that satisfies requests and ignores the distributor's objective is to give each agent a unique subset of T of size m . The s -max algorithm allocates to an agent the data record that yields the minimum increase of the maximum relative overlap among any pair of agents. The s -max algorithm is as follows.

Algorithm 4: Allocation for Sample Data Requests (SF)

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

Output: R_1, \dots, R_n

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

Step 2: R_1, \dots, R_n ;

Step 3: remaining

Step 4: while remaining > 0 do

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

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

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

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

Step 9: remaining \leftarrow remaining - 1.

Algorithm 5 : Object Selection for s-random

Step 1: function SELECTOBJECT(i, R_i)

Step 2: $k \leftarrow$ select at random an element from set $\{k' \neq tk'\}$

Step 3: return k .

Algorithm 6 : Object Selection for s-overlap

Step 1: function SELECTOBJECT($i; R_i; a$)

Step 2: $K \leftarrow \{k \mid k = \text{argmin } a[k']\}$

Step 3: $k \leftarrow$ select at random an element from set $\{k' \mid k' \in K \wedge tk' \in R_i\}$

Step 4: return k .

Algorithm 7 : Object Selection for s-max

Step1: function SELECTOBJECT($i, R_1, \dots, R_n, m_1, \dots, m_n$)

Step 2: $\text{min_overlap} \leftarrow 1$. The minimum out of the maximum relative overlaps that the allocations of different objects to U_i yield

Step 3: for $k \in \{k' \mid tk' \in R_i\}$ do

Step 4: $\text{max_rel_ov} \leftarrow 0$. The maximum relative overlap between R_i and any set R_j that the allocation of tk to U_i yields

Step 5: for all $j = 1, \dots, n : j \neq i$ and $tk \in R_j$ do

Step 6: $\text{abs_ov} \leftarrow |R_i \cap R_j| + 1$

Step 7: $\text{rel_ov} \leftarrow \text{abs_ov} / \min(m_i, m_j)$

Step 8: $\text{max_rel_ov} \leftarrow \text{MAX}(\text{max_rel_ov}, \text{rel_ov})$

Step 9: if $\max_rel_ov \leq \min_overlap$ then

Step 10: $\min_overlap \leftarrow \max_rel_ov$

Step 11: $ret_k \leftarrow k$

Step 12: return ret_k .

4. Existing System

There are conventional techniques being used and include technical and fundamental analysis. The main issue with these techniques is that they are manual and need laborious work along with experience. Traditionally, escape detection is handled by watermarking, e.g., a novel code is embedded in each distributed copy. If that replicate is later discovered within the hands of AN unauthorized party, the leaker can be identified. Watermarks may be terribly helpful in some cases, but again, involve some modification of the original data. Furthermore, watermarks can sometimes be destroyed if the data recipient is malicious. E.g. A hospital could offer patient records to researchers World Health Organization can devise new treatments. Similarly, an organization could have partnerships with alternative corporations that need sharing client knowledge. Another enterprise may outsource its process, this information should lean to varied different companies. We decide the owner of the info the distributor and therefore the purportedly trusty third parties the agents. The distributor gives the data to the agents. These data will be watermarked. Watermarking is the method of embedding the name or info concerning the corporate. The examples embody the images we've seen within the net. The authors of the pictures are watermarked within it. If anyone tries to repeat the image or knowledge the watermark are going to be present. And thus the data may be unusable by the leakers.

5. Proposed System

We propose data allocation strategies (across the agents) that improve the chance 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 fake" data records to any improve our chances of detecting leakage and identifying the guilty party. We also present an algorithm for distributing the object to an agent.

Our goal is to detect when the distributor's sensitive information has been leaked by agents, and if possible to spot the agent that leaked the information. Perturbation may be a terribly helpful technique wherever the information is modified and made 'less sensitive' before being handed to agents. We develop unobtrusive techniques for detecting leakage of a set of objects or records. In this section, we have a tendency to develop a model for assessing the 'guilt' of agents. We also present algorithms for distributing objects to agents, in an exceedingly means 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 however seem realistic to the agents. In a sense, the faux objects act as a kind of watermark for the entire set, without modifying any individual members. If it turns out an agent was given one or a lot of faux objects that were leaked, then the distributor is a lot of assured that agent was guilty. Today the advancement in technology made the watermarking system a simple technique of data

authorization. There are various software which can remove the watermark from the information and makes the data as original.

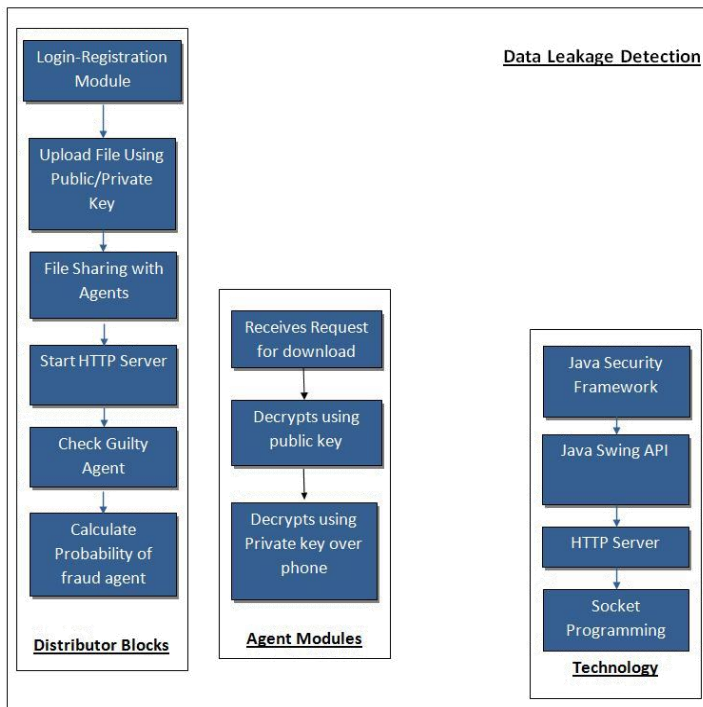


Fig-1: Data Leakage Detection

6. Conclusion

In an excellent world, there would be no ought to hand over sensitive data to agents that may unknowingly or maliciously leak it. And even if we had to hand over sensitive knowledge, in an exceedingly excellent world we tend to may watermark every object so we tend to may trace its origins with absolute certainty. However, in many cases, we must so work with agents which will not be 100% trusted. In spite of these difficulties, we have shown it is possible to assess the chance that AN agent is responsible for a leak, based on the overlap of his data with the leaked knowledge and therefore the knowledge of alternative agents, and supported the chance that objects are often ‘guessed’ by other means. Our model is relatively

simple, however, we tend to believe it captures the essential tradeoffs. The algorithms we have presented implement a variety of data distribution strategies that can improve the distributor’s chances of identifying a leaker. We have shown that distributing objects judiciously can make a significant difference in identifying guilty agents, especially in cases wherever there's giant overlap within the data that agents must receive. It includes the investigation of agent guilt

models that capture leakage scenarios that are not studied in this paper.

References

- [1] Data Leakage Detection, an IEEE paper by Panagiotis Papadimitriou, Member, IEEE, Hector Garcia-Molina, Member, IEEE NOV-2010.
- [2] Watermarking relational databases. In VLDB '02: Proceedings of the 28th international conference on Very Large Data Bases, By R. Agrawal and J. Kiernan, pages 155–166. VLDB Endowment, 2002.
- [3] An algebra for composing access control policies, By P. Bonatti, S. D. C. di Vimercati and P. Samarati, ACM Trans. Inf. Syst. Secur., 5(1):1–35, 2002.
- [4] P. Buneman, S. Khanna, and W. C. Tan. Why and where: A characterization of data provenance. In J. V. den Bussche and V. Vianu, editors, Database Theory - ICDT 2001, 8th International Conference, London, UK, January 4-6, 2001, Proceedings, volume 1973 of Lecture Notes in Computer Science, pages 316–330. Springer, 2001.
- [5] P. Buneman and W.-C. Tan. Provenance in databases. In SIGMOD '07: Proceedings of the 2007 ACM SIGMOD international conference on

Management of data, pages 1171–1173, New York,

NY, USA, 2007. ACM.

[6] Lineage tracing for general data warehouse transformations, By Y. Cui and J. Widom, In The VLDB Journal, pages 471–480, 2001.

[7] Digital music distribution and audio watermarking,

by S. Czerwinski, R. Fromm, and T. Hodes.

[8] Data Leakage Detection by Rajesh Kumar