

## Data Leakage Detection

Devayani A. Sirbhate

Computer Engineering, Sipna C.O.E.T., S.G.B. University, Amravati

### Abstract

*This paper contains concept of data leakage, cause of data leakage, and different techniques to protect and detect the data leakage. The value of data is incredible, so it should not be leaked or altered. In the field of IT huge database is being used. This database is shared with multiple people at a time. But during this sharing of the data, there are huge chances of data vulnerability, leakage or alteration. So, to prevent these problems, a data leakage detection system has been invent. This paper includes brief idea about data leakage detection and a methodology to detect the data leakage persons.*

**Keywords-** FakeObject, Guilty agent, data distributor, watermarking, data leakage.

### 1. Introduction

The unauthorized transfer of classified information from a computer or data center to the outside world. Data leakage is simply defined as the accidental or intentional distribution of private or sensitive data to an unauthorized entity. Sensitive data contain confidential document like financial information, patient information, personal credit card data and other information depending upon the business and the industry. Furthermore, in many cases, sensitive data shared among various stakeholders such as employees working from outside the organizational premises, business partners and customers. This increases the risk of confidential information falling into unauthorized hands. Furthermore, in many cases, sensitive data is shared.

### 2. Existing System

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. E.g. 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.

### 3. Proposed System

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 develop unobtrusive techniques for detecting leakage of a set of objects or records.

In this section 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.

### 4. Methodology

#### 4.1 Entities And Agents

Distributor owns set of data objects  $T = \{t_1, t_2, \dots, t_n\}$ . Distributor has to share some of the objects with set of agents  $U_1, U_2 \dots U_n$ , but does not wish the object 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  $R_i$  of objects  $T$ , determined either by a sample request or an explicit request.

- **Evaluation of Explicit Data Request Algorithms**

In this request the agent will send the request with appropriate condition. Agent gives the input as request with input as well as the condition for the request. After processing on the data he will get the new data by adding fake object using watermarking technique. Explicit request  $R_i = \text{EXPLICIT}(T, \text{condi})$ : Agent  $U_i$  receives all  $T$  objects that satisfy  $\text{condi}$ .

- **Evaluation of Sample Data Request Algorithms**

In this request agent request does not have condition. The agent sends the request without condition as per his query he will get the data by adding fake object using watermarking technique. Sample request =  $\text{SAMPLE}(T, m_i)$ : Any subset of  $m_i$  records from  $T$  can be given to  $U_i$ .

## 4.2 Guilty Agents

Suppose that after giving objects to agents, the distributor discovers that a set  $S$  belongs to  $T$  has leaked. This means that the third party, called the target, has been caught in possession of  $S$ . We say an agent  $U_i$  is guilty and if it contributes one or more object to 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  $S$ .

## 5. User Model

To compute this  $\Pr\{G_i|S\}$ , we need an estimate for the probability that values in  $S$  can be “guessed” by the target. We call this estimate  $pt$ , the probability that object  $t$  can be guessed by the target. Probability  $P_t$  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  $S$ . The component failure is used to compute the overall system reliability, 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  $P_t$ '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  $P_t$ 's that are higher than the true values, we will know that the agents will be guilty with at least the computed probabilities.

## 6. Data Allocation Problem

The main focus of our paper is the data allocation problem: how can the distributor “intelligently” give data to agents in order to improve the chances of detecting a guilty agent?

### 6.1 Fake Objects

The idea of perturbing data to detect leakage is not new. 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. For example, say the distributed data objects are medical records and the agents are hospitals. In this case, even small modifications to the records of actual patients may be undesirable. However, the addition of some fake medical records may be acceptable, since no patient matches these records, and hence no one will ever be treated based on fake records. Our use of fake objects is inspired by the use of “trace” records in mailing lists. In this case, company A sells to company B a mailing list to be used once (e.g., to send advertisements). Company A adds trace records that contain addresses owned by company A. Thus, each time company B uses the purchased mailing list, A receives copies of the mailing. These records are a type of fake objects that help identify improper use of data. The distributor creates and adds fake objects to the data that he distributes to agents. We let  $F_i \subseteq R_i$  be the subset of fake objects that agent  $U_i$  receives. Fake objects must be created carefully so that agents cannot distinguish them from real objects.

### 6.2.1 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.

## 8. Conclusion

In a perfect world, there would be no need to hand over sensitive data to agents that may unknowingly or maliciously leak it. And even if we had to hand over sensitive data, in a perfect world, we could watermark each object so that we could trace its origins with absolute certainty. However, in many cases, we must indeed work with agents that may not be 100 percent trusted, and we may not be certain if a leaked object came from an agent or from some other source, since certain data cannot admit watermarks.

In spite of these difficulties, we have shown that it is possible to assess the likelihood that an agent is responsible for a leak, based on the overlap of his data with the leaked data and the data of other agents, and based on the probability that objects can be “guessed” by other means. Our model is relatively simple, but we believe that it captures the essential trade-offs. 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 where there is large overlap in the data that agents must receive.

## 9. Acknowledgement

I am very grateful to my **Prof. V.K.Shandilya** who helped us and supported a lot. And I also want to thank our head of the department **Dr.R.B.Gavande** . Finally, I pay sincere thanks to all those who indirectly and directly helped me towards the successful completion of the paper.

## 10. References

- [1] P. Papadimitriou and H. Garcia-Molina, “Data leakage detection,” IEEE Transactions on Knowledge and Data Engineering, pages 51-63, volume 23, 2011.
- [2] 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.
- [3] Sandip A. Kale, Prof. S.V. Kulkarni/ IOSR Journal of Computer Engineering (IOSRJCE) ISSN:2278-0661 Volume 1, Issue 6 (July-Aug 2012), PP 32-35 www.iosrjournals.org / page Data Leakage Detection : A Survey.
- [4] Technical Report TR-BGU-2409-2010 24 Sept. 2010 1 A Survey of Data Leakage Detection and

Prevention Solutions P.P (1 -5, 24-25) A. Shabtai, a. Gershman, M. Kopeetsky, y. Elovici Deutsche Telekom Laboratories at Ben-Gurion University, Israel.

[5] Mr.V.Malsoru, Naresh Bollam/ International Journal of Engineering Research and Applications (IJERA)ISSN: 2248 -9622 www.ijera.com Vol. 1, Issue 3, pp.1088-1091 1088 | P a g e REVIEW ON DATALEAKAGE DETECTION.

[6] IEEE TRANSACTIONS ON KNOWLEDGE AND DATA ENGINEERING, VOL. 22, NO. 3, MARCH2011 Data Leakage Detection Panagiotis Papadimitriou, Member, IEEE, Hector Garcia-Molina,Member, IEEE P.P (2,4-5)

[7] Data Leakage: Affordable Data Leakage Risk Management by Joseph A. Rivela Senior Security Consultant P.P (4-6)

[8] Data Leakage Prevention: A news letter for IT Professionals Issue 5 P.P (1-3)

[9] The Who, What, When & Why of Data Leakage Prevention/Protection Presented by: Archie Alimagno

California Department of Insurance P.P (2-7)

[10] An ISACA White Paper Data Leak Prevention P.P (3-7)

[11] Mr.V.Malsoru, Naresh Bollam/ International Journal of Engineering Research and Applications (IJERA)

ISSN: 2248-9622 www.ijera.com Vol. 1, Issue 3, pp.1088-1091 1088 | P a g e REVIEW ON DATA LEAKAGE DETECTION