



Privacy-Preserving Clustering to Uphold Business Collaboration: A Dimensionality Reduction-Based Transformation Approach

Stanley R. M. Oliveira, Embrapa Informática Agropecuária, Brasil

Osmar R. Zaiane, University of Alberta, Canada

ABSTRACT

While the sharing of data is known to be beneficial in data mining applications and widely acknowledged as advantageous in business, this information sharing can become controversial and thwarted by privacy regulations and other privacy concerns. Data clustering for instance could be more accurate if more information is available, hence the data sharing. Any solution needs to balance the clustering requirements and the privacy issues. Rather than simply hindering data owners from sharing information for data analysis, a solution could be designed to meet privacy requirements and guarantee valid data clustering results. To achieve this dual goal, this article introduces a method for privacy-preserving clustering called dimensionality reduction-based transformation (DRBT). This method relies on the intuition behind random projection to protect the underlying attribute values subjected to cluster analysis. It is shown analytically and empirically that transforming a dataset using DRBT, a data owner can achieve privacy preservation and get accurate clustering with little overhead of communication cost. Such a method presents the following advantages: it is independent of distance-based clustering algorithms, it has a sound mathematical foundation, and it does not require CPU-intensive operations.

Keywords: *business collaboration; dimensionality reduction; privacy-preserving clustering; random projection*

INTRODUCTION

Data clustering is of capital importance in business and it fosters business collaboration, as sharing data for clustering improves the prospects of identifying optimal customer targets, marketing more effectively and understanding customer behaviour. Data clustering

maximizes return on investment supporting business collaboration (Lo, 2002; Berry & Linoff, 1997). Often combining different data sources provides better clustering analysis opportunities. Limiting the clustering on only some attributes of the data confines the correctness of the grouping, while benefiting from additional

attributes could yield more accurate and actionable clusters. For example, it does not suffice to cluster customers based on their purchasing history, but combining purchasing history, vital statistics and other demographic and financial information for clustering purposes can lead to better and more accurate customer behaviour analysis. More often than not, needed data sources are distributed, partitioned and owned by different parties insinuating a requirement for sharing data, often sensitive, between parties. Despite its benefits to support both modern business and social goals, clustering can also, in the absence of adequate safeguards, jeopardize individuals' privacy. The fundamental question addressed in this article is: how can data owners protect personal data shared for cluster analysis and meet their needs to support decision-making or to promote social benefits? To address this problem, data owners must not only meet privacy requirements but also guarantee valid clustering results.

Attaining good clustering may require data sharing between parties and data sharing may jeopardize privacy, a dilemma facing many modern data mining applications. Achieving privacy preservation when sharing data for clustering poses challenges for novel uses of data mining technology. Each application poses a new set of challenges. Let us consider two real-life examples in which the sharing of data poses different constraints:

- Two organizations, an Internet marketing company and an online retail company, have datasets with different attributes for a common set of individuals. These organizations decide to share their data for clustering to find the optimal customer targets so as to maximize return on investments. How can these organizations learn about their clusters using each other's data without learning anything about the attribute values of each other?
- Suppose that a hospital shares some data for research purposes (e.g., to group patients who have a similar disease). The hospital's security administrator may

suppress some identifiers (e.g., name, address, phone number, etc.) from patient records to meet privacy requirements. However, the released data may not be fully protected. A patient record may contain other information that can be linked with other datasets to re-identify individuals or entities (Samarati, 2001; Sweeney, 2002). How can we identify groups of patients with a similar pathology or characteristics without revealing the values of the attributes associated with them?

The above scenarios describe two different problems of privacy-preserving clustering (PPC). We refer to the former as PPC over centralized data and the latter as PPC over vertically partitioned data. To address these scenarios, we introduce a new PPC method called dimensionality reduction-based transformation (DRBT). This method allows data owners to find a trade-off between privacy, accuracy, and communication cost. Communication cost is the cost (typically in size) of the data exchanged between parties in order to achieve secure clustering.

Dimensionality reduction techniques have been studied in the context of pattern recognition (Fukunaga, 1990), information retrieval (Bingham & Mannila, 2001; Faloutsos & Lin, 1995; Jagadish, 1991), and data mining (Fern & Brodley, 2003; Faloutsos & Lin, 1995). To the best of our knowledge, dimensionality reduction has not been used in the context of data privacy in any detail, except in Oliveira & Zaïane (2004b). Although there exists a number of methods for reducing the dimensionality of data, such as feature extraction methods (Kaski, 1999), multidimensional scaling (Young, 1987) and principal component analysis (PCA) (Fukunaga, 1990), this article focuses on random projection, a powerful method for dimensionality reduction. The accuracy obtained after the dimensionality has been reduced, using random projection, is almost as good as the original accuracy (Kaski, 1999; Achlioptas, 2001; Bingham & Mannila, 2001). More formally, when a

22 more pages are available in the full version of this document, which may be purchased using the "Add to Cart" button on the publisher's webpage: www.igi-global.com/article/privacy-preserving-clustering-uphold-business/2459

Related Content

Administering the Semantic Web: Confidentiality, Privacy, and Trust Management

Bhavani Thuraisingham, Natasha Tsybulnik and Ashraf Alam (2007). *International Journal of Information Security and Privacy* (pp. 18-34).

www.irma-international.org/article/administering-semantic-web/2454

Patient Empowerment in IoT for eHealth: How to Deal With Lost Keys

Emmanuel Benoist, Serge Bignens and Alexander Kreutz (2020). *Applied Approach to Privacy and Security for the Internet of Things* (pp. 140-153).

www.irma-international.org/chapter/patient-empowerment-in-iot-for-ehealth/257909

Accurate Classification Models for Distributed Mining of Privately Preserved Data

Sumana M. and Hareesha K.S. (2016). *International Journal of Information Security and Privacy* (pp. 58-73).

www.irma-international.org/article/accurate-classification-models-for-distributed-mining-of-privately-preserved-data/165107

Preserving Privacy in Mining Quantitative Associations Rules

Madhu V. Ahluwalia, Aryya Gangopadhyay and Zhiyuan Chen (2011). *Security and Privacy Assurance in Advancing Technologies: New Developments* (pp. 310-326).

www.irma-international.org/chapter/preserving-privacy-mining-quantitative-associations/49509

The VESP Model: A Conceptual Model of Supply Chain Vulnerability

Arij Lahmar, Habib Chabchoub, François Galasso and Jacques Lamothe (2018). *International Journal of Risk and Contingency Management* (pp. 42-66).

www.irma-international.org/article/the-vesp-model/201074