

INTRODUCTION

1.1. OVERVIEW:

The objective of data mining is to generalize across populations, rather than reveal information about individuals. The hitch is that data mining works by evaluating individual data that is subject to privacy concerns. Thus, the true problem is not data mining, but the way data mining is done. However, the concern among privacy advocates is well founded, as bringing data together to support data mining makes misuse easier. Much of this information has already been collected, however it is held by various organizations. Separation of control and individual safeguards prevent correlation of this information, providing acceptable privacy in practice. However, this separation also makes it difficult to use the information for purposes that would benefit society, such as identifying criminal activity. Proposals to share information across agencies, most recently to combat terrorism, would eliminate the safeguards imposed by separation of the information.

1.2. MOTIVATION:

The title "Privacy Preservation in Collaborative Data Mining as Goal Oriented Attack Model" is justified by the implementation of Homomorphic Encryption to secure the mined data between the user and the server.

The goal of data mining is to extract or mine knowledge from large amounts of data. However, details often collected by several different sites. Privacy, legal and commercial concerns restrict centralized access to this data. Theoretical results from the area of secure multi party computation in cryptography prove that assuming the existence of trapdoor permutations; one may provide secure protocols for any two party's computation as well as for any multiparty computation with honest majority. However, the general methods are far too inefficient and impractical for computing complex

functions on inputs consisting of large sets of data. What remains open is come up with a set of techniques to achieve this efficiently within a quantifiable security framework. The distributed data model considered is the heterogeneous databases scenario with different features of the same set of data being collected by different sites.

This thesis argues that it is indeed possible to have efficient and practical techniques for useful privacy-preserving mining of knowledge from large amounts of data. The dissertation presents several privacy preserving data mining algorithms operating over vertically partitioned data. This set of underlying techniques solving independent sub-problems are also presented. Together, these enable these secure "mining" of knowledge.

In today's information age, data collection is ubiquitous, and every transaction is recorded somewhere. The resulting datasets can consist of terabytes or even petabytes of data, so efficiency and scalability is the primary consideration of most data mining algorithms. Data mining technology has emerged as a means of identifying patterns and trends from large quantities of data. Most tools operate by gathering all data into a central site, then running an algorithm against that data. However, privacy concerns can prevent building a centralized warehouse and data may be distributed among several custodians, none of which are allowed to transfer their data to another site. The problem is that computing association rules. The goal is to produce association rules that hold globally while limiting the information shared about each site. Previous work in privacy preserving data mining has addressed two issues. In one, the aim is preserving customer privacy by distorting the data values. The idea is that the distorted data does not reveal private information and thus is safe to use for mining. The key result is that the distorted data, and information on the distribution of the random data used to distort the data, can be used to generate an approximation to the original data values.

Recent advances in data collection, data dissemination and related technologies have inaugurated a new era of research where existing data mining algorithms should be reconsidered from the point of view of privacy preservation. The need for privacy is sometimes due to law (e.g., for medical databases) or can be motivated by business interests. However, there are situations where the sharing of data can lead to mutual