Chapter 5 Minimum Database Determination and Preprocessing for Machine Learning

Angel Fernando Kuri-Morales ITAM, Mexico

ABSTRACT

The exploitation of large databases implies the investment of expensive resources both in terms of the storage and processing time. The correct assessment of the data implies that pre-processing steps be taken before its analysis. The transformation of categorical data by adequately encoding every instance of categorical variables is needed. Encoding must be implemented that preserves the actual patterns while avoiding the introduction of non-existing ones. The authors discuss CESAMO, an algorithm which allows us to statistically identify the pattern preserving codes. The resulting database is more economical and may encompass mixed databases. Thus, they obtain an optimal transformed representation that is considerably more compact without impairing its informational content. For the equivalence of the original (FD) and reduced data set (RD), they apply an algorithm that relies on a multivariate regression algorithm (AA). Through the combined application of CESAMO and AA, the equivalent behavior of both FD and RD may be guaranteed with a high degree of statistical certainty.

DOI: 10.4018/978-1-5225-7268-8.ch005

Copyright © 2019, IGI Global. Copying or distributing in print or electronic forms without written permission of IGI Global is prohibited.

INTRODUCTION

Nowadays, commercial enterprises are importantly oriented to continuously improving customer-business (CRM) relationship. With the increasing influence of CRM Systems, such companies dedicate more time and effort to maintain better customer-business relationships. The effort implied in getting to better know the customer involves the accumulation of very large data bases where the largest possible quantity of data regarding the customer is stored.

Data warehouses offer a way to access detailed information about the customer's history, business facts and other aspects of the customer's behavior. The databases constitute the information backbone for any well established company. However, from each step and every new attempted link of the company to its customers the need to store increasing volumes of data arises. Hence databases and data warehouses are always growing up in terms of number of registers and tables which will allow the company to improve the general vision of the customer.

Data warehouses are difficult to characterize when trying to analyze the customers from company's standpoint. This problem is generally approached through the use of data mining techniques (Palpanas, T., 2000; Silva, D. R., 2002; Han, J., Pei, J., & Kamber, M. 2011; Tan, P. N. 2006, Chaudhuri, S., & Dayal, U. (1997). To attempt direct clustering over a data base of several terabytes with millions of registers results in a costly and not always fruitful effort. There have been many attempts to solve this problem. For instance one may use parallel computation, optimization of clustering algorithms, alternative distributed and grid computing and so on. But still the more efficient methods are unwieldy when attacking the clustering problem for databases as considered above. In this work we present a methodology derived from the practical solution of an automated clustering process over large database from a real large sized (over 20 million customers) company. We emphasize the way we used statistical methods to reduce the search space of the problem as well as the treatment given to the customer's information stored in multiple tables of multiple databases.

Because of confidentiality issues the name of the company and the actual final results of the customer characterization are withheld.

CHAPTER OUTLINE

The outline of the chapter is as follows. First, we give an overview of the analysis of large databases; next we give an overview of the methodology we applied. We emphasize the problem of adequately pre-processing non-numerical attributes so that numerical algorithms are applicable, in general. We describe two possible methods

36 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: <u>www.igi-</u> <u>global.com/chapter/minimum-database-determination-and-</u> <u>preprocessing-for-machine-learning/214833</u>

Related Content

A Reengineering Approach for Ensuring Transactional Reliability of Composite Services

Sami Bhiri, Walid Gaalouland Claude Godart (2010). *Web Services Research for Emerging Applications: Discoveries and Trends (pp. 290-316).* www.irma-international.org/chapter/reengineering-approach-ensuring-transactional-reliability/41527

A Hybrid Approach to Big Data Systems Development

Anil K. Aggarwal (2019). Web Services: Concepts, Methodologies, Tools, and Applications (pp. 2271-2288).

www.irma-international.org/chapter/a-hybrid-approach-to-big-data-systems-development/217942

A Predictive and Evolutionary Approach for Cost-Effective and Deadline-Constrained Workflow Scheduling Over Distributed IaaS Clouds

Jiangchuan Chen, Jiajia Jiangand Dan Luo (2019). *International Journal of Web* Services Research (pp. 78-94).

www.irma-international.org/article/a-predictive-and-evolutionary-approach-for-cost-effective-and-deadline-constrained-workflow-scheduling-over-distributed-iaas-clouds/231451

Decentralized Communication for Data Dependency Analysis Among Process Execution Agents

Susan D. Urban, Ziao Liuand Le Gao (2011). *International Journal of Web Services Research (pp. 1-28).*

www.irma-international.org/article/decentralized-communication-data-dependencyanalysis/60164

A Method for Predicting Wikipedia Editors' Editing Interest: Based on a Factor Graph Model

Haisu Zhang, Sheng Zhang, Zhaolin Wu, Liwei Huangand Yutao Ma (2016). *International Journal of Web Services Research (pp. 1-25).*

www.irma-international.org/article/a-method-for-predicting-wikipedia-editors-editinginterest/161800