Chapter 68 Knowledge Discovery From Massive Data Streams

Sushil Kumar Narang SAS Institute of IT and Research, India

> Sushil Kumar IIT Roorkee, India

Vishal Verma MLN College, India

ABSTRACT

T.S. Eliot once wrote some beautiful poetic lines including one "Where is the knowledge we have lost in information?". Can't say that T.S. Eliot could have anticipated today's scenario which is emerging from his poetic lines. Data in present scenario is a profuse resource in many circumstances and is piling-up and many technical leaders are finding themselves drowning in data. Through this big stream of data there is a vast flood of information coming out and seemingly crossing manageable boundaries. As Information is a necessary channel for educing and constructing knowledge, one can assume the importance of generating new and comprehensive knowledge discovery tools and techniques for digging this overflowing sea of information to create explicit knowledge. This chapter describes traditional as well as modern research techniques towards knowledge discovery from massive data streams. These techniques have been effectively applied not exclusively to completely structured but also to semi-structured and unstructured data. At the same time Semantic Web technologies in today's perspective require many of them to deal with all sorts of raw data.

1. INTRODUCTION TO KNOWLEDGE DISCOVERY

Knowledge Discovery (KD) is a concept which involves the developments of strategies and procedures for making sense out of massive data. In recent years, data have become increasingly available in substantial amounts (petabytes or zettabytes). It has numerous sources including automation of business activities (trading, mobile communication, airline reservation, or credit card usage), online activities

DOI: 10.4018/978-1-5225-5191-1.ch068

(social media, social networking), scientific activities (experiments, simulations, and environmental sensors), biological databases (DNA/RNA/protein structures, gene expression profiles) etc. In addition, new application scenarios like weather forecasting, artificial intelligence, earth observation satellites and so forth produce terabytes of data every day. Clearly, the massive size of data ruled out any manual approach of analyzing (make sense of) collected data. If this massive data will have to be understood at all, it must be analyzed by the use of computers. Although, there are statistical procedures available for data analysis and interpretation, but this explosive growth of data requires new intelligent techniques which can astutely transform the useful data into knowledge. Knowledge discovery is the significant process of digging out meaningful patterns from huge data using automated (or semi-automated) computational tools and techniques (Devedzic, 2002; Piatetsky-Shapiro, 1996). The goals of knowledge discovery are usually identified by business domain. For instance

- Marketing agencies make use of knowledge discovery frameworks to find patterns in the way customers purchase retail items. Once they find that many individuals purchase item A along with item B, they can easily make an appropriate and potentially successful business or marketing announcement.
- Airline companies make use of knowledge discovery systems to find patterns in which their passengers fly (routes, return flights, frequency of flying to a specific destination and so forth). Based on the patterns discovered, they can give promotional offers to frequent travelers, thus attract more customers to the company.
- Banks make use of knowledge discovery frameworks to explore the database of their credits and loans. Based on the patterns discovered, they can more successfully predict the risk of approving loan to their clients, thus increasing the quality of their business decisions.

Of course, most of these goals were well existing even before knowledge discovery was conceptualized. They have been achieved by human expertise, numerical modeling and on the basis of database OLAP (online analytical processing). However, in Knowledge Discovery, these goals are achieved by applying automated (or semi-automated) computational tools and techniques to the huge amount of stored data. Subsequently, knowledge discovery has turned out to be strategically important for big business units, government institutions, and research organizations. However, successfully producing knowledge from massive data sets is very challenging. A lot of research is going on in the area of knowledge discovery to establish stable and well-defined standards which are well understood throughout the community. These standards need to be formalized to make the process more translucent and repeatable.

2. KNOWLEDGE DISCOVERY PROCESS

The Knowledge Discovery Process (KDP) is a long process that not only limited to actual data analysis but includes methods for data collection and preparation; data reduction and projection; data analysis and interpretation, and finally evaluation and action based on the discovered knowledge. Since the 1990s, several KDP models have been proposed (Kurgan & Musilek, 2006; Mariscal, Marban, & Fernandez, 2010). The early models were proposed by academic researchers but quickly followed by some peculiar industry research. The first basic model of the knowledge discovery process was proposed by Fayyad,

25 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/chapter/knowledge-discovery-from-massive-datastreams/198612

Related Content

Investigations On Some Aspects of Reliability of Content Based Routing SOAP based Windows Communication Foundation Services

Subhash Medhi, Abhijit Boraand Tulshi Bezboruah (2017). International Journal of Information Retrieval Research (pp. 17-31).

www.irma-international.org/article/investigations-on-some-aspects-of-reliability-of-content-based-routing-soap-basedwindows-communication-foundation-services/165377

Bates' Berrypicking Model (1989, 2002, 2005)

Linda L. Lillardand YooJin Ha (2018). *Information Retrieval and Management: Concepts, Methodologies, Tools, and Applications (pp. 65-76).* www.irma-international.org/chapter/bates-berrypicking-model-1989-2002-2005/198545

Linked Data, Towards Realizing the Web of Data: An Overview

Leila Zemmouchi-Ghomari (2018). Information Retrieval and Management: Concepts, Methodologies, Tools, and Applications (pp. 292-312). www.irma-international.org/chapter/linked-data-towards-realizing-the-web-of-data/198555

Chatman's Theories of Information Behavior (1996, 1999, 2000)

Abdullahi I. Musa (2015). Information Seeking Behavior and Technology Adoption: Theories and Trends (pp. 136-148).

www.irma-international.org/chapter/chatmans-theories-of-information-behavior-1996-1999-2000/127128

Effective Management of Data Centers Resources for Load Balancing in Cloud Computing

Pradeep Kumar Tiwariand Sandeep Joshi (2018). International Journal of Information Retrieval Research (pp. 40-56).

www.irma-international.org/article/effective-management-of-data-centers-resources-for-load-balancing-in-cloudcomputing/198964