

Chapter 5.21

Heuristics in Medical Data Mining

Susan E. George

University of South Australia, Australia

HISTORICAL PERSPECTIVE

Deriving—or discovering—information from data has come to be known as data mining. Within health care, the knowledge from medical mining has been used in tasks as diverse as patient diagnosis (Brameier et al., 2000; Mani et al., 1999; Cao et al., 1998; Henson et al., 1996), inventory stock control (Bansal et al., 2000), and intelligent interfaces for patient record systems (George et al., 2000). It has also been a tool of medical discovery itself (Steven et al., 1996). Yet, it remains true that medicine is one of the last areas of society to be “automated,” with a relatively recent increase in the volume of electronic data, many paper-based clinical record systems in use, a lack of standardisation (for example, among coding schemes), and still some reluctance among health-care providers to use computer technology. Nevertheless, the rapidly increasing volume of electronic medical data is perhaps one of the domain’s current distinguishing characteristics, as one of the last components of society to be “automated.”

Data mining presents many challenges, as “knowledge” is automatically extracted from data

sets, especially when data are complex in nature, with many hundreds of variables and relationships among those variables that vary in time, space, or both, often with a measure of uncertainty, as is common within medicine. Cios and Moore (2001) identified a number of unique features of medical data mining, including the use of imaging and need for visualisation techniques, the large amounts of unstructured nature of free text within records, data ownership and the distributed nature of data, the legal implications for medical providers, the privacy and security concerns of patients requiring anonymous data used, where possible, together with the difficulty in making a mathematical characterisation of the domain.

Strictly speaking, many ventures within medical data mining are better described as exercises in “machine learning,” where the main issues are, for example, discovering the complexity of relationships among data items, or making predictions in light of uncertainty, rather than “data mining,” in large, possibly distributed, volumes of data that are also highly complex. Large data sets mean not only increased algorithmic complexity but also often the need to employ special-purpose methods

to isolate trends and extract “knowledge” from data. However, medical data frequently provide just such a combination of vast (often distributed) complex data sets.

Heuristic methods are one way in which the vastness, complexity, and uncertainty of data may be addressed in the mining process. A heuristic is something that aids discovery of a solution. Artificial intelligence (AI) popularised the heuristic as something that captures, in a computational way, the knowledge that people use to solve everyday problems. AI has a classic graph search algorithm known as A* (Hart et al., 1968), which is a heuristic search (under the right conditions). Increasingly, heuristics refer to techniques that are inspired by nature, biology, and physics. The genetic search algorithm (Holland, 1975) may be regarded as a heuristic technique. More recent population-based approaches have been demonstrated in the Memetic Algorithm (Moscato, 1989), and specific modifications of such heuristic methods in a medical mining context can be noted (Brameier et al., 2000).

Aside from the complexity of data with which the medical domain is faced, there are some additional challenges. Data security, accuracy, and privacy are issues within many domains, not just the medical (Walhstrom et al., 2000). Also, while ethical responsibility is an issue in other contexts, it is faced by the medical world in a unique way, especially when heuristic methods are employed. One of the biggest ethical issues concerns what is done with the knowledge derived combined with a “forward-looking responsibility” (Johnson et al., 1995). Forward-looking responsibility is accountable for high-quality products and methods and requires appropriate evaluation of results and justification of conclusions.

George (2002) first identified and proposed a set of guidelines for heuristic data mining within medical domains. The proposed guidelines relate to the evaluation and justification of data-mining results (so important when heuristic “aids to discovery” are utilised that “may” benefit a solution)

and extend to both where and how the conclusions may be utilised and where heuristic techniques are relevant in this field. The remainder of this article summarises some heuristic data-mining applications in medicine and clarifies those proposed guidelines.

BACKGROUND

First, we will explain some of the heuristic methods that have been employed in medical data mining, examining a range of application areas. We broadly categorise applications as clinical, administrative, and research, according to whether they are used (or potentially used) in a clinical context, are infrastructure related, or are exploratory, in essence. We also note that with the exception of some medical imaging applications and mining of electronic medical records, the databases are small.

There is a wide variety of automated systems that have been designed for diagnosis—systems that detect a problem, classify it, and monitor change. Brameier and Banzhaf (2000) described the application of linear genetic programming to several diagnosis problems in medicine, including tests for cancer, diabetes, heart conditions, and thyroid conditions. Their focus was upon an efficient algorithm that operates with a range of complex data sets, providing a population-based heuristic method that is based upon biological principles. Their heuristic method is based on an inspiration from nature about how “introns” (denoting DNA segments with information removed before proteins are synthesised) are used in generating new strings. They suggest that introns may help to reduce the number of destructive recombinations between chromosomes by protecting the advantageous building blocks from being destroyed by crossover. Massive efficiency improvements in the algorithm are reported.

An interesting administrative application of data mining in a medical context comes in the

4 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/heuristics-medical-data-mining/7780

Related Content

Marketing Data Mining

Victor S.Y. Lo (2005). *Encyclopedia of Data Warehousing and Mining* (pp. 698-704).

www.irma-international.org/chapter/marketing-data-mining/10687

Ontology-Based Construction of Grid Data Mining Workflows

Peter Brezany, Ivan Janciak and A. Min Tjoa (2008). *Data Warehousing and Mining: Concepts, Methodologies, Tools, and Applications* (pp. 913-941).

www.irma-international.org/chapter/ontology-based-construction-grid-data/7680

Data Driven vs. Metric Driven Data Warehouse Design

John M. Artz (2005). *Encyclopedia of Data Warehousing and Mining* (pp. 223-227).

www.irma-international.org/chapter/data-driven-metric-driven-data/10597

Biomedical Data Mining Using RBF Neural Networks

Feng Chu and Lipo Wang (2008). *Data Warehousing and Mining: Concepts, Methodologies, Tools, and Applications* (pp. 1706-1713).

www.irma-international.org/chapter/biomedical-data-mining-using-rbf/7726

Benchmarking Data Mining Algorithms

Balaji Rajagopalan and Ravi Krovi (2002). *Data Warehousing and Web Engineering* (pp. 77-99).

www.irma-international.org/chapter/benchmarking-data-mining-algorithms/7862