

Chapter 8.4

Pattern Discovery as Event Association

Andrew K. C. Wong
University of Waterloo, Canada

Yang Wang
Pattern Discovery Technology, Canada

Gary C. L. Li
University of Waterloo, Canada

INTRODUCTION

A basic task of machine learning and data mining is to automatically uncover patterns that reflect regularities in a data set. When dealing with a large database, especially when domain knowledge is not available or very weak, this can be a challenging task. The purpose of pattern discovery is to find non-random relations among events from data sets. For example, the “exclusive OR” (XOR) problem concerns 3 binary variables, A, B and $C=A \otimes B$, i.e. C is true when either A or B, but not both, is true. Suppose not knowing that it is the XOR problem, we would like to check whether or not the occurrence of the compound event $[A=T, B=T, C=F]$ is just a random happening. If we could

estimate its frequency of occurrences under the random assumption, then we know that it is not random if the observed frequency deviates significantly from that assumption. We refer to such a compound event as an event association pattern, or simply a pattern, if its frequency of occurrences significantly deviates from the default random assumption in the statistical sense. For instance, suppose that an XOR database contains 1000 samples and each primary event (e.g. $[A=T]$) occurs 500 times. The expected frequency of occurrences of the compound event $[A=T, B=T, C=F]$ under the independence assumption is $0.5 \times 0.5 \times 0.5 \times 1000 = 125$. Suppose that its observed frequency is 250, we would like to see whether or not the difference between the observed and expected frequencies (i.e. $250 - 125$) is significant enough to indicate that the compound event is not a random happening.

DOI: 10.4018/978-1-60960-818-7.ch8.4

In statistics, to test the correlation between random variables, contingency table with chi-squared statistic (Mills, 1955) is widely used. Instead of investigating variable correlations, pattern discovery shifts the traditional correlation analysis in statistics at the variable level to association analysis at the event level, offering an effective method to detect statistical association among events.

In the early 90's, this approach was established for second order event associations (Chan & Wong, 1990). A higher order pattern discovery algorithm was devised in the mid 90's for discrete-valued data sets (Wong & Yang, 1997). In our methods, patterns inherent in data are defined as statistically significant associations of two or more primary events of different attributes if they pass a statistical test for deviation significance based on residual analysis. The discovered high order patterns can then be used for classification (Wang & Wong, 2003). With continuous data, events are defined as Borel sets and the pattern discovery process is formulated as an optimization problem which recursively partitions the sample space for the best set of significant events (patterns) in the form of high dimension intervals from which probability density can be estimated by Gaussian kernel fit (Chau & Wong, 1999). Classification can then be achieved using Bayesian classifiers. For data with a mixture of discrete and continuous data (Wong & Yang, 2003), the latter is categorized based on a global optimization discretization algorithm (Liu, Wong & Yang, 2004). As demonstrated in numerous real-world and commercial applications (Yang, 2002), pattern discovery is an ideal tool to uncover subtle and useful patterns in a database.

In pattern discovery, three open problems are addressed. The first concerns learning where noise and uncertainty are present. In our method, noise is taken as inconsistent samples against statistically significant patterns. Missing attribute values are also considered as noise. Using a standard statistical hypothesis testing to confirm statistical patterns from the candidates, this method is a less ad hoc

approach to discover patterns than most of its contemporaries. The second problem concerns the detection of polythetic patterns without relying on exhaustive search. Efficient systems for detecting monothetic patterns between two attributes exist (e.g. Chan & Wong, 1990). However, for detecting polythetic patterns, an exhaustive search is required (Han, 2001). In many problem domains, polythetic assessments of feature combinations (or higher order relationship detection) are imperative for robust learning. Our method resolves this problem by directly constructing polythetic concepts while screening out non-informative pattern candidates, using statistics-based heuristics in the discovery process. The third problem concerns the representation of the detected patterns. Traditionally, if-then rules and graphs, including networks and trees, are the most popular ones. However, they have shortcomings when dealing with multilevel and multiple order patterns due to the non-exhaustive and unpredictable hierarchical nature of the inherent patterns. We adopt attributed hypergraph (AHG) (Wang & Wong, 1996) as the representation of the detected patterns. It is a data structure general enough to encode information at many levels of abstraction, yet simple enough to quantify the information content of its organized structure. It is able to encode both the qualitative and the quantitative characteristics and relations inherent in the data set.

BACKGROUND

In the ordinary sense, "discovering regularities" from a system or a data set implies partitioning the observed instances into classes based on similarity. Michalski and Stepp (1983) pointed out that the traditional distance-based statistical clustering techniques make no distinction among relevant, less relevant and irrelevant attributes nor do they render conceptual description of the clusters with human input. They proposed CLUSTER/2 as a conceptual clustering algorithm in a noise-free

7 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/pattern-discovery-event-association/56234

Related Content

Innovations in Education: Integrating Explainable AI Into Educational Intelligence

Shugufta Fatima, C. Kishor Kumar Reddy, Akshita Sunerahand Srinath Doss (2025). *Internet of Behavior-Based Computational Intelligence for Smart Education Systems* (pp. 19-52).

www.irma-international.org/chapter/innovations-in-education/358973

A Survey of Different Approaches for the Class Imbalance Problem in Software Defect Prediction

Abdul Waheed Darand Sheikh Umar Farooq (2022). *International Journal of Software Science and Computational Intelligence* (pp. 1-26).

www.irma-international.org/article/a-survey-of-different-approaches-for-the-class-imbalance-problem-in-software-defect-prediction/301268

The Optimal Path Finding Algorithm Based on Reinforcement Learning

Ganesh Khekare, Pushpneel Verma, Urvashi Dhanre, Seema Rautand Shahrukh Sheikh (2020). *International Journal of Software Science and Computational Intelligence* (pp. 1-18).

www.irma-international.org/article/the-optimal-path-finding-algorithm-based-on-reinforcement-learning/262585

On Cognitive Models of Causal Inferences and Causation Networks

Yingxu Wang (2013). *Advances in Abstract Intelligence and Soft Computing* (pp. 103-113).

www.irma-international.org/chapter/cognitive-models-causal-inferences-causation/72776

Machine Learning-Based Categorization of COVID-19 Patients

Tanvi Arora (2022). *Applications of Computational Science in Artificial Intelligence* (pp. 214-233).

www.irma-international.org/chapter/machine-learning-based-categorization-of-covid-19-patients/302068