



Chapter VII

Advances in Classification of Sequence Data

Pradeep Kumar, University of Hyderabad, Gachibowli, India*

P. Radha Krishna, University of Hyderabad, Gachibowli, India

Raju S. Bapi, University of Hyderabad, Gachibowli, India

T. M. Padmaja, University of Hyderabad, Gachibowli, India

Abstract

In recent years, advanced information systems have enabled a collection of increasingly large amounts of data that are sequential in nature. To analyze huge amounts of sequential data, the interdisciplinary field of knowledge discovery in databases (KDD) is very useful. The most important step within the process of KDD is data mining, which is concerned with the extraction of the valid patterns. Recent research focus in data mining includes stream data mining, sequence data mining, Web mining, text mining, visual mining, multimedia mining, and multi-relational data mining. Sequence data may be discrete or continuous in nature. Most of the research on discrete sequence data concentrated on the discovery of frequently occurring patterns. However, comparatively less amount of work has been carried out in the area of discrete sequence data classification. In this chapter, data taxonomy is introduced with a review of the state of art for sequence data classification. The usefulness of embedding partial subsequence information extracted using sliding window technique into traditional classifier like kNN has been demonstrated. kNN

has been tested with various vector based distance/similarity metrics. Further, with the use of S^3M similarity metric, the full subsequence information embedded in the data sequences is extracted. The experimental data taken is DARPA '98 IDS benchmark dataset collected from UCIML dataset repository. The chapter closes by pointing out various application areas of sequence data and also the open issues in sequence data classification problem.

Introduction

In recent years, the amount of data that is collected by advanced information systems has increased tremendously. The resulting data volume is too large to be examined manually and even the methods for automatic data analysis based on classical statistics and machine learning often face problems when processing large, dynamic data sets consisting of complex data. The interdisciplinary field of knowledge discovery in databases (KDD) helps in analyzing these large volumes of data. KDD employs methods at the cross-post of machine learning, statistics, and database systems. Fayyad, Shapiro, and Smyth (1996) define KDD as follows:

Knowledge discovery in databases is the non-trivial process of identifying valid, novel, potentially useful, and ultimately understandable patterns in data.

According to this definition, data is a set of facts that is somehow accessible in electronic form. The term patterns indicate models and regularities, which can be observed within the data. Patterns have to be valid (i.e., they should be true on new data with some degree of certainty). A novel pattern is not previously known or trivially true. The potential usefulness of patterns refers to the possibility that they lead to an action providing a benefit. A pattern is understandable if it is interpretable by a human user. KDD is a process comprising several steps that are perhaps repeated over several iterations.

The core step of KDD is called *data mining*. Data mining can be defined as an activity that extracts new and nontrivial information contained in large databases. The goal is to discover hidden patterns, unexpected trends or other non-obvious relationships in the data using a combination of techniques such as, machine learning, statistics, and databases.

A huge amount of data is collected every day in the form of sequences. These sequential data are valuable sources of information not only to search for a particular value or event at a specific time, but also to analyze the frequency of certain events or sets of events related by particular temporal/sequential relationship. Examples

30 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/advances-classification-sequence-data/7517

Related Content

A Dynamic and Semantically-Aware Technique for Document Clustering in Biomedical Literature

Min Song, Xiaohua Hu, Illhoi Yoo and Eric Koppel (2009). *International Journal of Data Warehousing and Mining* (pp. 44-57).

www.irma-international.org/article/dynamic-semantically-aware-technique-document/37404

Empirical Investigation of Decision Tree Ensembles for Monitoring Cardiac Complications of Diabetes

Andrei V. Kelarev, Jemal Abawajy, Andrew Stranieri and Herbert F. Jelinek (2013). *International Journal of Data Warehousing and Mining* (pp. 1-18).

www.irma-international.org/article/empirical-investigation-of-decision-tree-ensembles-for-monitoring-cardiac-complications-of-diabetes/105117

A Hyper-Heuristic for Descriptive Rule Induction

Tho Hoan Pham and Tu Bao Ho (2007). *International Journal of Data Warehousing and Mining* (pp. 54-66).

www.irma-international.org/article/hyper-heuristic-descriptive-rule-induction/1778

User Segmentation Based on Twitter Data Using Fuzzy Clustering

Basar Öztaysi and Sezi Çevik Onar (2013). *Data Mining in Dynamic Social Networks and Fuzzy Systems* (pp. 316-333).

www.irma-international.org/chapter/user-segmentation-based-twitter-data/77533

Deep Learning Approach Towards Unstructured Text Data Utilization: Development, Opportunities, and Challenges

Shikha Jain, Shubham Jain and Ajit Kumar Jain (2021). *New Opportunities for Sentiment Analysis and Information Processing* (pp. 29-49).

www.irma-international.org/chapter/deep-learning-approach-towards-unstructured-text-data-utilization/286903