Advanced Information and Knowledge Processing

Springer-Verlag London Ltd.

Shichao Zhang, Chengqi Zhang and
Xindong Wu

# Knowledge Discovery
# in Multiple Databases

With 21 Figures

Springer

Shichao Zhang, PhD, MSc
Chengqi Zhang, PhD, MSc, BSc, DSc
FIT, University of Technology Sydney, Australia

Xindong Wu, PhD, MSc
Department of Computer Science, University of Vermont, USA

*Series Editors*
Xindong Wu
Lakhmi Jain

# Preface

Many organizations have an urgent need of mining their multiple databases inherently distributed in branches (distributed data). In particular, as the Web is rapidly becoming an information flood, individuals and organizations can take into account low-cost information and knowledge on the Internet when making decisions. How to efficiently identify quality knowledge from different data sources has become a significant challenge.

This challenge has attracted a great many researchers including the authors who have developed a local pattern analysis, a new strategy for discovering some kinds of potentially useful patterns that cannot be mined in traditional multi-database mining techniques. Local pattern analysis delivers high-performance pattern discovery from multiple databases. There has been considerable progress made on multi-database mining in such areas as hierarchical meta-learning, collective mining, database classification, and peculiarity discovery. While these techniques continue to be future topics of interest concerning multi-database mining, this book focuses on these interesting issues under the framework of local pattern analysis.

The book is intended for researchers and students in data mining, distributed data analysis, machine learning, and anyone else who is interested in multi-database mining. It is also appropriate for use as a text supplement for broader courses that might also involve knowledge discovery in databases and data mining.

The book consists of ten chapters. Chapter 1 states the multi-database mining problem and its importance. Chapter 2 lays a common foundation for subsequent material. This includes the preliminaries on data mining and multi-database mining, as well as necessary concepts, previous efforts, and applications. Chapter 3 introduces the framework of local pattern analysis. The later chapters are essentially self-contained and may be read selectively, and in any order. Chapters 4, 5, and 6 develop techniques for preprocessing the data in multi-databases. Chapters 7, 8, and 9 presents techniques for identifying interesting patterns from multi-databases based on local pattern analysis. And Chapter 10 presents a summary of the previous chapters and demonstrates some open problems.

Beginners should read Chapters 1 and 2 before selectively reading other chapters. Although the opening problems are very important, techniques in

other chapters may be helpful for experienced readers who want to attack such problems.

Shichao Zhang

Chengqi Zhang

Xindong Wu

# Acknowledgments

# Contents