

Honghua Dai · James N.K. Liu  
Evgueni Smirnov *Editors*

# Reliable Knowledge Discovery

 Springer

# Reliable Knowledge Discovery



Honghua Dai • James N.K. Liu • Evgueni Smirnov  
Editors

# Reliable Knowledge Discovery

 Springer

*Editors*

Honghua Dai  
Deakin University  
Burwood, Victoria, Australia

James N.K. Liu  
The Hong Kong Polytechnic University  
Hong Kong

Evgueni Smirnov  
Maastricht University  
Maastricht, The Netherlands

ISBN 978-1-4614-1902-0                      e-ISBN 978-1-4614-1903-7

DOI 10.1007/978-1-4614-1903-7

Springer New York Dordrecht Heidelberg London

Library of Congress Control Number: 2012932410

© Springer Science+Business Media, LLC 2012

All rights reserved. This work may not be translated or copied in whole or in part without the written permission of the publisher (Springer Science+Business Media, LLC, 233 Spring Street, New York, NY 10013, USA), except for brief excerpts in connection with reviews or scholarly analysis. Use in connection with any form of information storage and retrieval, electronic adaptation, computer software, or by similar or dissimilar methodology now known or hereafter developed is forbidden.

The use in this publication of trade names, trademarks, service marks, and similar terms, even if they are not identified as such, is not to be taken as an expression of opinion as to whether or not they are subject to proprietary rights.

Printed on acid-free paper

Springer is part of Springer Science+Business Media ([www.springer.com](http://www.springer.com))

# Preface

## 1 Description

With the rapid development of the data mining and knowledge discovery, a key issue which could significantly affect the real world applications of data mining is the reliability issues of knowledge discovery. It is natural that people will ask if the discovered knowledge is reliable. Why do we trust the discovered knowledge? How much can we trust the discovered knowledge? When it could go wrong. All these questions are very essential to data mining. It is especial crucial to the real world applications.

One of the essential requirements of data mining is validity. This means both the discovery process itself and the discovered knowledge should be valid. Reliability is a necessary but not sufficient condition for validity. Reliability could be viewed as stability, equivalence and consistency in some ways.

This special volume of the book on the reliability issues of Data Mining and Knowledge Discovery will focus on the theory and techniques that can ensure the discovered knowledge is reliable and to identify under which conditions the discovered knowledge is reliable or in which cases the discovery process is robust. In the last 20 years, many data mining algorithms have been developed for the discovery of knowledge from given data bases. However in some cases, the discovery process is not robust or the discovered knowledge is not reliable or even incorrect in certain cases. We could also find that in some cases, the discovered knowledge may not necessary be the real reflection of the data. Why does this happen? What are the major factors that affect the discovery process? How can we make sure that the discovered knowledge is reliable? What are the conditions under which a reliable discovery can be assured? These are some interesting questions to be investigated in this book.

## 2 Scope and Topics of this Book

The topics of this book covers the following:

- The theories on reliable knowledge discovery
- Reliable knowledge discovery methods
- Reliability measurement criteria of knowledge discovery
- Reliability estimation methods
- General reliability issues on knowledge discovery
- Domain specific reliability issues on knowledge discovery
- The criteria that can be used to assess the reliability of discovered knowledge.
- The conditions under which we can confidently say that the discovered knowledge is reliable.
- The techniques which can improve reliability of knowledge discovery
- Practical approaches that can be used to solve reliability problems of data mining systems.
- The theoretical work on data mining reliability
- The practical approaches which can be used to assess if the discovered knowledge is reliable.
- The analysis of the factors that affect data mining reliability
- How reliability can be assessed
- In which condition, the reliability of the discovered knowledge is assured.

## 3 The Theme and Related Resources

The main purpose of this book is to encourage the use of Reliable Knowledge Discovery from Databases (RKDD) in critical-domain applications related to society, science, and technology. The book is intended for practitioners, researchers, and advanced-level students. It can be employed primarily as a reference work and it is a good compliment to the excellent book on reliable prediction Algorithmic learning in a random world by Vladimir Vovk, Alex Gammerman, and Glenn Shafer (New York: Springer, 2005). Extra information sources are the proceedings of the workshops Reliability Issues in Knowledge Discovery held in conjunction with the IEEE International Conferences on Data Mining. Other relevant conferences are the Annual ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD), the International Conference on Machine Learning (ICML), The pacific-Asia Conference on Knowledge Discovery (PAKDD), and the European Conference on Machine Learning and Principles and Practice of Knowledge Discovery in Databases (ECML PKDD). Many AI-related journals regularly publish work in RKDD. Among others it is worth mentioning the Journal of Data Mining and Knowledge Discovery, the Journal of Machine Learning Research, and the Journal of Intelligent Data Analysis.

## 4 An Overview of the Book

This book presents the recent advances in the emerging field of **Reliable Knowledge Discovery from Data (RKDD)**. In this field the knowledge is considered as reliable in the sense that its generalization performance can be set in advance. Hence, RKDD has a potential for a broad spectrum of applications, especially in critical domains like medicine, finance, military etc. The main material presented in the book is based on three consequent workshops Reliability Issues in Knowledge Discovery held in conjunction with the IEEE International Conferences on Data Mining (ICDM) in 2006, 2008, and 2010, respectively. In addition we provided an opportunity to authors to publish the results of their newest research related to RKDD.

This book is organized in seventeen chapters divided into four parts.

### **Part I includes three chapters on Reliability Estimation.**

Chapter 1 provides an overview of typicalness and transductive reliability estimation frameworks. The overview is employed for introducing an approach for accessing reliability of individual classifications called joint confidence machine. Chapter 1 describes an approach that compensates the weaknesses of typicalness-based confidence estimation and transductive reliability estimation by integrating them into a joint confidence machine. It provides better interpretation of the performance of any classifiers. Experimental results performed with different machine learning algorithms in several problem domains show that there is no reduction of discrimination performance and is more suitable for applications with risk-sensitive problems with strict confidence limits.

Chapter 2 introduces new approaches to estimating and correcting individual predictions in the context of stream mining. It investigates the online reliability estimation of individual predictions. It proposes different strategies and explores techniques based on local variance and local bias, of local sensitivity analysis and online bagging of predictors. Comparison results on benchmark data are given to demonstrate the improvement of prediction accuracy.

Chapter 3 deals with the problem of quantifying the reliability in the context of neural networks. It elaborates on new approaches to estimation of confidence and prediction intervals for polynomial neural networks.



## **Part II includes seven chapters on Reliable Knowledge Discovery Methods.**

Chapter 4 investigates outliers in regression targeting robust diagnostic regression. The chapter discusses both robust regression and regression diagnostics, presents several contemporary methods through numerical examples in linear regression.

Chapter 5 presents a conventional view on the definition of reliability; points out the three major categories of factors that affect the reliability of knowledge discovery, examined the impact of model complexity, weak links, varying sample sizes and the ability of different learners to the reliability of graphical model discovery, proposed reliable graph discovery approaches.

Chapter 6 provides a generalization of version spaces for reliable classification implemented using support vector machines.

Chapter 7 presents a unified generative model ONM which characterizes the life cycle of a ticket. The model uses maximum likelihood estimation to capture reliable ticket transfer profiles which can reflect how the information contained in a ticket is used by human experts to make reliable ticket routing decisions.

Chapter 8 applies the methods of aggregation functions for the reliable web based knowledge discovery from network traffic data.

Chapter 9 gives two new versions of SVM for the regression study of features in the problem domain. It provides means for feature selection and weighting based on the correlation analysis to give better and reliable result.

Chapter 10 describes in detail an application of transductive confidence machines for reliable handwriting recognition. It introduces a TCM framework which can enhance classifiers to reduce the computational costs and memory consumption required for updating the non-conformity scores in the offline learning setup of TCMs. Results are found to have outperformed previous methods on both relatively easy data and on difficult test samples.

## **PART III includes four Chapters on Reliability Analysis.**

Chapter 11 addresses the problem of reliable feature selection. It introduces a generic-feature-selection measure together with a new search approach for globally optimal feature-subset selection. It discusses the reliability in the feature-selection process of a real pattern-recognition system, provides formal measurements and allows consistent search for relevant features in order to attain global optimal solution.

Chapter 12 provides three detailed case studies to show how the reliability of an induced classifier can be influenced. The case study results reveal the impact of data-oriented factors to the reliability of the discovered knowledge.

Chapter 13 analyzes recently-introduced instance-based penalization techniques capable of providing more accurate predictions.

Chapter 14 investigates subsequence frequency measurement and its impact on the reliability of knowledge discovery in single sequences.

## **PART IV includes three chapters on Reliability Improvement Methods.**

Chapter 15 proposed to use the inexact field learning method and parameter optimized one-class classifiers to improving reliability of unbalanced text mining by reducing performance bias.

Chapter 16 proposes a formal description technique for ontology representation and verification using a high level Petri net approach. It provides the capability of detection and identification of potential anomalies in ontology for the improvement of the discovered knowledge.

Chapter 17 presents an UGDSS framework to provide reliable support for multi-criteria decision making in uncertainty problem domain. It gives the system design and architecture.

## **5 Acknowledgement**

We would like to thank many people that made this book possible. We start with the organizers of the workshops held in conjunction with the IEEE International Conferences on Data Mining (ICDM): Shusaku Tsumoto, Francesco Bonchi, Bettina Berendt, Wei Fan and Wynne Hsu. We express our gratitude to the authors whose contributions can be found in the book. Finally, we thank our colleagues from Springer that made the publication process possible in a short period.

Burwood Victoria (Australia),  
Hong Kong (China),  
Maastricht (The Netherlands),

*Honghua Dai*  
*James Liu*  
*Evgueni Smirnov*  
August 2011



# Contents

## Part I Reliability Estimation

<b>1</b>	<b>Transductive Reliability Estimation for Individual Classifications in Machine Learning and Data Mining</b> .....	<b>3</b>
	Matjaž Kukar	
1.1	Introduction .....	3
1.2	Related work .....	4
1.2.1	Transduction .....	5
1.3	Methods and materials .....	6
1.3.1	Typicalness .....	6
1.3.2	Transductive reliability estimation .....	8
1.3.3	Merging the typicalness and transduction frameworks ...	15
1.3.4	Meta learning and kernel density estimation .....	16
1.3.5	Improving kernel density estimation by transduction principle .....	18
1.3.6	Testing methodology .....	18
1.4	Results .....	20
1.4.1	Experiments on benchmark problems .....	20
1.4.2	Real-life application and practical considerations .....	22
1.5	Discussion .....	23
	References .....	26
<b>2</b>	<b>Estimating Reliability for Assessing and Correcting Individual Streaming Predictions</b> .....	<b>29</b>
	Pedro Pereira Rodrigues, Zoran Bosnić, João Gama, and Igor Kononenko	
2.1	Introduction .....	30
2.2	Background .....	30
2.2.1	Computation and utilization of prediction reliability estimates .....	31

- 2.2.2 Correcting individual regression predictions ..... 32
- 2.2.3 Reliable machine learning from data streams ..... 32
- 2.3 Estimating reliability of individual streaming predictions ..... 34
  - 2.3.1 Preliminaries ..... 34
  - 2.3.2 Reliability estimates for individual streaming predictions ..... 35
  - 2.3.3 Evaluation of reliability estimates ..... 37
  - 2.3.4 Abalone data set ..... 39
  - 2.3.5 Electricity load demand data stream ..... 40
- 2.4 Correcting individual streaming predictions ..... 40
  - 2.4.1 Correcting predictions using the CNK reliability estimate ..... 41
  - 2.4.2 Correcting predictions using the Kalman filter ..... 42
  - 2.4.3 Experimental evaluation ..... 43
  - 2.4.4 Performance of the corrective approaches ..... 44
  - 2.4.5 Statistical comparison of the predictions' accuracy ..... 45
- 2.5 Conclusions ..... 46
- References ..... 48
- 3 Error Bars for Polynomial Neural Networks ..... 51**  
 Nikolay Nikolaev, and Evgeni Smirnov
  - 3.1 Introduction ..... 51
  - 3.2 Genetic Programming of PNN ..... 52
    - 3.2.1 Polynomial Regression ..... 52
    - 3.2.2 Tree-structured PNN ..... 53
    - 3.2.3 Weight Learning ..... 54
    - 3.2.4 Mechanisms of the GP System ..... 54
  - 3.3 Sources of PNN Deviations ..... 56
  - 3.4 Estimating Confidence Intervals ..... 56
    - 3.4.1 Delta Method for Confidence Intervals ..... 57
    - 3.4.2 Residual Bootstrap for Confidence Intervals ..... 59
  - 3.5 Estimating Prediction Intervals ..... 60
    - 3.5.1 Delta Method for Prediction Intervals ..... 61
    - 3.5.2 Training Method for Prediction Bars ..... 62
  - 3.6 Conclusion ..... 65
  - References ..... 65

**Part II Reliable Knowledge Discovery Methods**

- 4 Robust-Diagnostic Regression: A Prelude for Inducing Reliable Knowledge from Regression ..... 69**  
 Abdul Awal Md. Nurunnabi, and Honghua Dai
  - 4.1 Introduction ..... 70
  - 4.2 Background of Reliable Knowledge Discovery ..... 71
  - 4.3 Linear Regression, OLS and Outliers ..... 72

4.4	Robustness and Robust Regression	73
4.4.1	Least Median of Squares Regression	75
4.4.2	Least Trimmed Squares Regression	75
4.4.3	Reweighted Least Squares Regression	76
4.4.4	Robust M (GM)- Estimator	76
4.4.5	Example	77
4.5	Regression Diagnostics	79
4.5.1	Examples	85
4.6	Concluding Remarks and Future Research Issues	89
	References	90
<b>5</b>	<b>Reliable Graph Discovery</b>	<b>93</b>
	Honghua Dai	
5.1	Introduction	93
5.2	Reliability of Graph Discovery	95
5.3	Factors That Affect Reliability of Graph Discovery	96
5.4	The Impact of Sample Size and Link Strength	97
5.5	Testing Strategy	98
5.6	Experimental Results and Analysis	100
5.6.1	Sample Size and Model Complexity	100
5.7	Conclusions	105
	References	105
<b>6</b>	<b>Combining Version Spaces and Support Vector Machines for Reliable Classification</b>	<b>109</b>
	Evgueni Smirnov, Georgi Nalbantov, and Ida Sprinkhuizen-Kuyper	
6.1	Introduction	109
6.2	Task of Reliable Classification	110
6.3	Version Spaces	110
6.3.1	Definition and Classification Rule	111
6.3.2	Analysis of Version-Space Classification	112
6.3.3	Volume-Extension Approach	113
6.4	Support Vector Machines	114
6.5	Version Space Support Vector Machines	115
6.5.1	Hypothesis Space	115
6.5.2	Definition of Version Space Support Vector Machines	118
6.5.3	Classification Algorithm	118
6.6	The Volume-Extension Approach for VSSVMs	119
6.7	Experiments	119
6.8	Comparison with Relevant Work	123
6.8.1	Bayesian Framework	123
6.8.2	Typicalness Framework	124
6.9	Conclusion	125
	References	125

<b>7</b>	<b>Reliable Ticket Routing in Expert Networks</b> .....	127
	Gengxin Miao, Louise E. Moser, Xifeng Yan, Shu Tao, Yi Chen, and Nikos Anerousis	
7.1	Introduction .....	128
7.2	Related Work .....	129
7.3	Preliminaries .....	131
7.4	Generative Models .....	133
	7.4.1 Resolution Model (RM) .....	133
	7.4.2 Transfer Model (TM) .....	134
	7.4.3 Optimized Network Model (ONM) .....	134
7.5	Ticket Routing .....	137
	7.5.1 Ranked Resolver .....	137
	7.5.2 Greedy Transfer .....	138
	7.5.3 Holistic Routing .....	139
7.6	Experimental Results .....	141
	7.6.1 Data Sets .....	141
	7.6.2 Model Effectiveness .....	142
	7.6.3 Routing Effectiveness .....	143
	7.6.4 Robustness .....	144
7.7	Conclusions and Future Work .....	144
	References .....	145
<b>8</b>	<b>Reliable Aggregation on Network Traffic for Web Based Knowledge Discovery</b> .....	149
	Shui Yu, Simon James, Yonghong Tian, and Wanchun Dou	
8.1	Introduction .....	150
8.2	The Reliability of Network Traffic Information .....	151
8.3	Aggregation Functions .....	151
8.4	Information Theoretical Notions of Distance .....	153
8.5	Performance Comparison for Information Distances .....	155
8.6	Summary .....	157
	References .....	158
<b>9</b>	<b>Sensitivity and Generalization of SVM with Weighted and Reduced Features</b> .....	161
	Yan-xing Hu, James N.K.Liu, and Li-wei Jia	
9.1	Introduction .....	161
9.2	Background .....	163
	9.2.1 The Classical SVM Regression Problem .....	163
	9.2.2 Rough Set SVM Regression .....	165
	9.2.3 Grey Correlation Based Feature Weighted SVM Regression .....	167
9.3	Experimental Results and Analysis .....	171
	9.3.1 Data Collection .....	171
	9.3.2 Data Pre-processing .....	172