

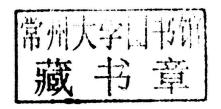
Giovanni Ponti

Advances in Mining Complex Data: Modeling and Clustering



Giovanni Ponti

Advances in Mining Complex Data: Modeling and Clustering



LAP LAMBERT Academic Publishing

Impressum / Imprint

Bibliografische Information der Deutschen Nationalbibliothek: Die Deutsche Nationalbibliothek verzeichnet diese Publikation in der Deutschen Nationalbibliografie; detaillierte bibliografische Daten sind im Internet über http://dnb.d-nb.de abrufbar.

Alle in diesem Buch genannten Marken und Produktnamen unterliegen warenzeichen-, marken- oder patentrechtlichem Schutz bzw. sind Warenzeichen oder eingetragene Warenzeichen der jeweiligen Inhaber. Die Wiedergabe von Marken, Produktnamen, Gebrauchsnamen, Handelsnamen, Warenbezeichnungen u.s.w. in diesem Werk berechtigt auch ohne besondere Kennzeichnung nicht zu der Annahme, dass solche Namen im Sinne der Warenzeichen- und Markenschutzgesetzgebung als frei zu betrachten wären und daher von jedermann benutzt werden dürften.

Bibliographic information published by the Deutsche Nationalbibliothek: The Deutsche Nationalbibliothek lists this publication in the Deutsche Nationalbibliografie; detailed bibliographic data are available in the Internet at http://dnb.d-nb.de.

Any brand names and product names mentioned in this book are subject to trademark, brand or patent protection and are trademarks or registered trademarks of their respective holders. The use of brand names, product names, common names, trade names, product descriptions etc. even without a particular marking in this works is in no way to be construed to mean that such names may be regarded as unrestricted in respect of trademark and brand protection legislation and could thus be used by anyone.

Coverbild / Cover image: www.ingimage.com

Verlag / Publisher:
LAP LAMBERT Academic Publishing
ist ein Imprint der / is a trademark of
AV Akademikerverlag GmbH & Co. KG
Heinrich-Böcking-Str. 6-8, 66121 Saarbrücken, Deutschland / Germany
Email: info@lap-publishing.com

Herstellung: siehe letzte Seite /

Printed at: see last page ISBN: 978-3-659-30522-1

Zugl. / Approved by: Arcavacata di Rende (CS), University of Calabria, DEIS Dept., 2010

Copyright © 2013 AV Akademikerverlag GmbH & Co. KG Alle Rechte vorbehalten. / All rights reserved. Saarbrücken 2013

Giovanni Ponti

Advances in Mining Complex Data: Modeling and Clustering

Acknowledgments

It is a honor for me to be part of a research group where people work with devotion and passion in a friendly atmosphere. This environment has provided me with incitements and support during my studies.

The first acknowledgment goes to my advisor Prof. Sergio Greco. He was the first who believed in my capabilities and introduced me in the scientific research world. His suggestions were very precious and his intuitions opened my mind to many interesting research opportunities.

A special thanks to Andrea Tagarelli, who led me carefully to and through my research activities. His critical observations and his constructive remarks contributed to the growth of my research capabilities. He also revealed to be a friend who provided encouragements in difficulties. His presence is invaluable, a reference point in work and in life.

I also acknowledge Prof. Luigi Palopoli, who directed the Ph.D. course. He was always interested in the evolution of my activities and his suggestions were valuable to enrich my studies.

The majority of the topics shown in this thesis are result of research activities conducted with my room mate, Francesco Gullo. Such results could not be obtained without his contribution. We worked together sharing opinions and providing precious remarks each other. Apart from work, we shared many moments of friendship. For these reasons, he deserves a great thanks.

There were many people at DEIS who stood by me during these years. All of them gave me an important contribution to improve my knowledge with helpful feedback, discussions, and complementary points of view. In particular, I would thank Bettina Fazzinga, Sergio Flesca, Filippo Furfaro, Massimiliano Mazzeo, Cristian Molinaro, Francesco Parisi, Andrea Pugliese, Francesca Spezzano, and Ester Zumpano for work and joy moments spent together. I also thank Mario Cannataro, Pierangelo Veltri, and Giuseppe Tradigo from UNICZ, who provided me a valuable support in several issues of my research experience. I acknowledge Giovanni

Costabile and Gabriele Gigliotti for their help in bureaucratic affairs.

Apart from the work environment, I would thank people with whom I spent my leisure time. For this reason, a great thanks goes to all my friends of "Post Cresima". I get known them since the youth and we grew up together, sharing lots of joyful and spiritual moments.

An invaluable thanks goes to my parents Luigi and Mena, as well as my sister Barbara. They always supported me in pursuing my objectives. A whole life would not be enough to express all my gratitude to them.

Last but not least, my deepest thanks is for Mari, my better half. She is the most important person in my life, and her love gave me the strength to overcome difficulties I faced during these years. Also, discussions and moments spent together were very important for my overall growth. My life is better since we met.



Preface

In the last years, there has been a great production of data that come from different application contexts. However, although technological progress provides several facilities to digitally encode any type of event, it is important to define a suitable representation model which underlies the main characteristics of the data. This aspect is particularly relevant in fields and contexts where data to be archived can not be represented in a fix structured scheme, or that can not be described by simple numerical values. We hereinafter refer to these data with the term *complex data*.

Although it is important define ad-hoc representation models for complex data, it is also crucial to have analysis systems and data exploration techniques. Analysts and system users need new instruments that support them in the extraction of patterns and relations hidden in the data. The entire process that aims to extract useful information and knowledge starting from raw data takes the name of *Knowledge Discovery in Databases* (KDD). It starts from raw data and consists in a set of specific phases that are able to transform and manage data to produce models and knowledge. There have been many knowledge extraction techniques for traditional structured data, but they are not suitable to handle complex data.

Investigating and solving representation problems for complex data and defining proper algorithms and techniques to extract models, patterns and new information from such data in an effective and efficient way are the main challenges which this thesis aims to face. In particular, two main aspects related to complex data management have been investigated, that are the way in which complex data can be modeled (i.e., data modeling), and the way in which homogeneous groups within complex data can be identified (i.e., data clustering). The application contexts that have been objective of such studies are time series data, uncertain data, text data, and biomedical data.

It is possible to illustrate research contributions of this thesis by dividing them into four main parts, each of which concerns with one specific area and data type:

Time Series — A time series representation model has been developed, which is conceived to support accurate and fast similarity detection. This model is called *Derivative time series Segment Approximation (DSA)*, as it achieves a concise yet feature-rich time series representation by combining the notions of derivative estimation, segmentation and segment approximation.

Uncertain Data — Research in uncertain data mining went into two directions. In a first phase, a new proposal for partitional clustering has been defined by introducing the Uncertain K-medoids (UK-medoids) algorithm. This approach provides a more accurate way to handle uncertain objects in a clustering task, since a cluster representative is an uncertain object itself (and not a deterministic one). In addition, efficiency issue has been addressed by defining a distance function between uncertain objects that can be calculated offline once per dataset.

In a second phase, research activities aimed to investigate issues related to hierarchical clustering of uncertain data. Therefore, an agglomerative centroidbased linkage hierarchical clustering framework for uncertain data (*U-AHC*) has been proposed. The key point lies in equipping such scheme with a more accurate distance measure for uncertain objects. Indeed, it has been resorted to *information theory* field to find a measure able to compare probability distributions of uncertain objects used to model uncertainty.

Text Data — Research results on text data can be summarized in two main contributions. The first one regards clustering of multi-topic documents, and a framework for hard clustering of documents according to their mixtures of topics has been proposed. Documents are assumed to be modeled by a generative process, which provides a mixture of probability mass functions (pmfs) to model the topics that are discussed within any specific document. The framework combines the expressiveness of generative models for document representation with a properly chosen information-theoretic distance measure to group the documents.

The second proposal concerns distributional clustering of XML documents, focusing on a the development of a distributed framework for efficiently clustering XML documents. The distributed environment consists of a peer-to-peer network where each node in the network has access to a portion of the whole document collection and communicates with all the other nodes to perform a

clustering task in a collaborative fashion. The proposed framework is based on modeling and clustering XML documents by structure and content. Indeed, XML documents are transformed into transactional data based on the notion of tree tuple. The framework is based on the well-known paradigm of centroid-based partitional clustering to conceive the distributed, transactional clustering algorithm.

Biomedical Data — Research results on time series and uncertain data have been involved to support effective and efficient biomedical data management. The focus regarded both proteomics and genomics, investigating Mass Spectrometry (MS) data and microarray data. In the specific, a Mass Spectrometry Data Analysis (MaSDA) system has been defined. The key idea consists in exploiting temporal information implicitly contained in MS data and model such data as time series. The major advantages of this solution are the dimensionality and the noise reduction. As regards micrarray data, U-AHC has been employed to perform clustering of microarray data with probe-level uncertainty. A strategy to model probe-level uncertainty has been defined, together with a hierarchical clustering scheme for analyzing such data. This approach performs a gene-based clustering to discover clustering solutions that are well-suited to capture the underlying gene-based patterns of microarray data.

The effectiveness and the efficiency of the proposed techniques in clustering complex data are demonstrated by performing intense and exhaustive experiments, in which such proposals are extensively compared with the main state-of-the-art competitors.

Giovanni Ponti



Contents

A	ckno	wledgr	ments	i			
P	Preface						
Li	st of	Figur	res	xiv			
Li	st of	Table	s	xvi			
Ι	Tł	ne Ba	sics	1			
1	Intr	oduct	ion	3			
	1.1	Data	heterogeneity and Information Systems	. 3			
	1.2	Know	ledge Discovery in Databases	. 4			
		1.2.1	Data Mining	. 5			
	1.3	Challenges for Complex Data		. 6			
		1.3.1	Trajectories and Time Series	. 7			
		1.3.2	Uncertain data	. 7			
		1.3.3	Text data	. 9			
		1.3.4	Biomedical Data	. 9			
	1.4	Contr	ibutions	. 10			
	1.5	Outlin	ne of the Thesis	. 14			
2	Background						
	2.1	Data	Mining and Clustering	. 17			
		2.1.1	Partitional Clustering	. 19			
		2.1.2	Hierarchical Clustering	. 22			
		2.1.3	Density-based Clustering	. 24			
		2.1.4	Soft Clustering	. 27			

x Contents

		2.1.5	Model-based Clustering	28			
		2.1.6	Evaluation Criteria	28			
	2.2	Time	Series	32			
		2.2.1	Similarity Search	33			
		2.2.2	Dimensionality Reduction	35			
	2.3	Uncer	tain Objects	35			
		2.3.1	Modeling Uncertainty	36			
		2.3.2	Distance Measures	39			
	2.4	2.4 Text Data					
		2.4.1	Text Preprocessing	40			
		2.4.2	Text Modeling	41			
П	т	ime S	Series Data Clustering	47			
			Series Data Grastering	11			
3	Tim	e Seri	es Data Management: State of the Art	49			
	3.1	Simila	rity Measures	49			
	3.2	Dimer	nsionality Reduction	50			
4	А Т	ima S	eries Representation Model for Accurate and Fast Simi-				
4		y Det		- 55			
	4.1		ation and Contributions	55			
	4.2		ative time series Segment Approximation (DSA)				
	1.2	4.2.1	Derivative Estimation	58			
		4.2.2	Segmentation	59			
		4.2.3	Segment Approximation	61			
	4.3	Exper	imental Evaluation	61			
		4.3.1	Settings	62			
		4.3.2	Results	68			
		M DE NEEDS					
5		DSA in Real Case Applications					
	5.1		zing Mass Spectrometry Data	77			
		5.1.1	The Mass Spectrometry Data Analysis (MaSDA) System	78			
		5.1.2	Experimental Evaluation	86			
	5.2		ng Electricity Company Customers	93			
		5.2.1	Low-voltage Electricity Customer Data				
				OF			

Contents xi

		5.2.3	Experimental Evaluation	. 98		
II	I	Uncer	tain Data Clustering	103		
6	Uncertain Data Clustering: State of the Art					
	6.1	Uncert	tain Data Clustering	. 105		
		6.1.1	Partitional Methods	. 105		
		6.1.2	Density-based Methods	. 106		
7	Clu	stering	g of Uncertain Objects via K-medoids	109		
	7.1	Introd	uction	. 109		
	7.2	Uncert	tain Distance	. 110		
	7.3	UK-me	edoids Algorithm	. 112		
	7.4	Experi	imental Evaluation	. 113		
		7.4.1	Settings	. 113		
		7.4.2	${\it Results} \ldots \ldots \ldots \ldots \ldots \ldots \ldots \ldots \ldots \ldots$. 115		
8	Information-Theoretic Hierarchical Clustering of Uncertain Ob-					
	jects					
	8.1	Introduction				
	8.2	Uncertain Prototype				
	8.3	uting Distance between Uncertain Objects	. 124			
		8.3.1	Distance Measure for pdfs	. 124		
		8.3.2	Distance Measure for Uncertain Objects	. 127		
	8.4	4 U-AHC Algorithm				
	8.5	mental Evaluation	. 138			
		8.5.1	Settings	. 138		
		8.5.2	Results	. 140		
	8.6	U-AHO	C for Clustering Microarray Data	. 144		
		8.6.1	Microarray Data and Probe-level Uncertainty	. 144		
		8.6.2	Experimental Evaluation	. 145		
TY	7 1	0 -	Chartering	1.40		
IV	/ J	Jocun	ment Clustering	149		
9	Ger	erative	e Models for Document Representation	151		
	9.1	Introd	uction	. 151		

	9.2	State of the Art	153		
	9.3	A Framework for Topic-based Hard Clustering of Documents	154		
	9.4	Distance Measure for Model-based Documents	155		
	9.5	MuToT-AHC Algorithm for Document Clustering	157		
		9.5.1 Experimental Evaluation	158		
10	Coll	aborative Clustering of XML Documents	165		
	10.1	.1 Introduction			
	10.2 State of the Art				
		10.2.1 XML Document Representation	168		
		10.2.2 Transactional Clustering	169		
	10.3	XML Transactional Representation	169		
10.4 Distributed XML Transactional Clustering			174		
		10.4.1 XML Tree Tuple Item Similarity	174		
		10.4.2 XML Transactional Clustering Algorithm for Collaborative			
		Distributed Environment (CXK-means)	176		
	10.5	Experimental Evaluation	180		
		10.5.1 Settings	180		
		10.5.2 Results	183		
11 Conclusion		clusion	187		
	11.1	Thesis Review	187		
	11.2	Future Works	189		
A	DD	ΓW and DSA Derivative Estimation Models	191		
В	Impact of Preprocessing on Similarity Detection 1				
\mathbf{C}	Dvn	amic Kernel Clustering	199		