

Summarizing Data with Representative Patterns



Chunyang Liu

Faculty of Engineering and Information Technology
University of Technology, Sydney

A thesis submitted for the degree of
Doctor of Philosophy

March 2016

To my parents and my wife.

Certificate of Original Authorship

I certify that the work in this thesis has not previously been submitted for a degree nor has it been submitted as part of requirements for a degree except as fully acknowledged within the text.

I also certify that the thesis has been written by me. Any help that I have received in my research work and the preparation of the thesis itself has been acknowledged. In addition, I certify that all information sources and literature used are indicated in the thesis.

Student: Chunyang Liu

Date: 05/03/2016

Acknowledgements

I benefited and learned a lot from my advisors, my friends, my colleagues and my family during PhD study in University of Technology, Sydney. I would like to take this good opportunity to appreciate their significant helps to me.

First of all, I would like to express my appreciation and gratitude to my academic supervisors Dr. Ling Chen and Prof. Chengqi Zhang. I benefited significantly from various discussions and communications with them. They always give me advices, encouragement, and sufficient freedom to think and explore. I also want to thank them for kindly offering me invaluable suggestions and experienced instructions on my research career and life. I cannot image how this thesis can be accomplished without their high scientific criterion, endless patience, generous support, and constant guidance.

Next, I wish to thank other researchers who have gave me helpful guidance and encouragement during my PhD study and on conferences: Prof. Jian Pei, A/Prof. Ivor W. Tsang, Dr. Lu Qin, and Dr. Yi Yang. From discussions with them I learned not only much knowledge but also their attitude for doing high quality research.

I have been fortunate to work in a center gathering the most brilliant researchers and best friends in the past four years: Guodong Long, Jing Jiang, Wei Bian, Tianyi Zhou, Meng Fang, Bozhong Liu, Zhe Xu, Mingming Gong, Ruxin Wang, Zhibin Hong, Sujuan Hou, Wei Wu, Haishuai Wang, Xueping Peng, Alan Wang, Zhenxing Qin, Shirui Pan, Jia Wu, Lianhua Chi, Maoying Qiao, Tongliang Liu, Weiwei Liu, Anjin Liu, Junfu Yin, Jinjiu Li, Changxing Ding, Nannan Wang, Naiyang Guan, Li Wan, Shengzheng

Wang, Xianhua Zeng, Zhijing Xu, Bo Du, Hongshu Chen, Ting Guo, Chao Ma, Shaoli Huang, Qiang Li, Dianshuang Wu, Zhiguo Long, Weilong Hou, Yi Ji, Ming Xie and many others. I enjoyed the invaluable friendships with them, their kindly support and accompany are always my source of strength and courage in both research and daily life. I am also grateful to my dearest friends Jiawang Liu, Weipeng Zhang, Mingjian Gao, Junfeng Ye, and Jia Chen since my college, and Sen Li and Song Li since my high school. They are the ones who have given me support during both joyful and stressful times, to whom I will always be thankful.

Finally, it is my greatest honor to thank my family: my parents and my wife. They always believe in me, encourage me, give me invaluable suggestions, and fully support all my decisions. No words could possibly express my deepest gratitude for their endless love, self-sacrifice and unwavering help. To them I dedicate this dissertation.

Abstract

The advance of technology makes data acquisition and storage become unprecedentedly convenient. It contributes to the rapid growth of not only the volume but also the veracity and variety of data in recent years, which poses new challenges to the data mining area. For example, uncertain data mining emerges due to its capability to model the inherent veracity of data; spatial data mining attracts much research attention as the widespread of location-based services and wearable devices. As a fundamental topic of data mining, how to effectively and efficiently summarize data in this situation still remains to be explored.

This thesis studied the problem of summarizing data with representative patterns. The objective is to find a set of patterns, which is much more concise but still contains rich information of the original data, and may provide valuable insights for further analysis of data. In the light of this idea, we formally formulate the problem and provide effective and efficient solutions in various scenarios.

We study the problem of summarizing probabilistic frequent patterns over uncertain data. Probabilistic frequent pattern mining over uncertain data has received much research attention due to the wide applicabilities of uncertain data. It suffers from the problem of generating an exponential number of result patterns, which hinders the analysis of patterns and calls for the need to find a small number of representative patterns to approximate all other patterns. We formally formulate the problem of *probabilistic representative frequent pattern (P-RFP) mining*, which aims to find the minimal set of patterns with sufficiently high probability to represent all other patterns. The bottleneck turns out to be checking whether a pattern can probabilistically represent another, which involves the computation of a joint probability of the supports of two patterns. We propose a novel dynamic programming-based approach to address the problem and devise effective optimization strategies to improve the computation efficiency.

To enhance the practicability of P-RFP mining, we introduce a novel approximation of the joint probability with both theoretical and empirical proofs. Based on the approximation, we propose an *Approximate P-RFP Mining* (APM) algorithm, which effectively and efficiently compresses the probabilistic frequent pattern set. The error rate of APM is guaranteed to be very small when the database contains hundreds of transactions, which further affirms that APM is a practical solution for summarizing probabilistic frequent patterns.

We address the problem of directly summarizing uncertain transaction database by formulating the problem as *Minimal Probabilistic Tile Cover Mining*, which aims to find a high-quality probabilistic tile set covering an uncertain database with minimal cost. We define the concept of *Probabilistic Price* and *Probabilistic Price Order* to evaluate and compare the quality of tiles, and propose a framework to discover the minimal probabilistic tile cover. The bottleneck is to check whether a tile is better than another according to the *Probabilistic Price Order*, which involves the computation of a joint probability. We prove that it can be decomposed into independent terms and calculated efficiently. Several optimization techniques are devised to further improve the performance.

We analyze the problem of summarizing co-locations mined from spatial databases. Co-location pattern mining finds patterns of spatial features whose instances tend to locate together in geographic space. However, the traditional framework of co-location pattern mining produces an exponential number of patterns because of the downward closure property, which makes it difficult for users to understand, assess or apply the huge number of resulted patterns. To address this issue, we study the problem of mining *representative co-location patterns* (RCP). We first define a covering relationship between two co-location patterns then formally formulate the problem of *Representative Co-location Pattern mining*. To solve the problem of RCP mining, we propose the *RCPFast* algorithm adopting the post-mining framework and the *RCPMS* algorithm pushing pattern summarization into the co-location mining process.

Contents

1	Introduction	1
1.1	Uncertain Data	2
1.2	Spatial Data	3
1.3	Main Contributions and Roadmap	4
1.4	Publications	6
2	Mining Probabilistic Representative Frequent Patterns From Uncertain Data	9
2.1	Introduction	10
2.2	Related Work	12
2.2.1	Frequent pattern mining over uncertain data	12
2.2.2	Frequent pattern summarization	13
2.3	Problem Definition	14
2.4	P-RFP Mining	17
2.4.1	Framework of P-RFP Mining	17
2.4.2	Cover Set Generation	20
2.4.3	Optimization Strategies	22
2.4.4	P-RFP Mining Algorithm	26
2.5	Performance Study	26
2.5.1	Data sets	26
2.5.2	Result analysis	29
2.6	Conclusions	30
3	Summarizing Probabilistic Frequent Patterns: A Fast Approach	33
3.1	Introduction	34
3.2	Related Work	36

3.2.1	Frequent pattern mining over uncertain data	36
3.2.2	Frequent pattern summarization	37
3.3	Background and Preliminary	38
3.4	Approximate P-RFP Mining	41
3.4.1	Framework of APM	42
3.4.2	Cover Set Generation	44
3.4.3	APM Algorithm	46
3.5	Approximation of Joint Support Probability	47
3.5.1	Preparation	47
3.5.2	Proof of Approximation	49
3.6	Performance Study	51
3.6.1	Empirical study of approximation	51
3.6.2	Result analysis	53
3.7	Conclusions	55
4	Summarizing Uncertain Transaction Databases by Probabilistic Tiles	61
4.1	Introduction	62
4.2	Problem Definition	66
4.2.1	Quality of Summarization	69
4.2.2	Parameter Setting	69
4.2.3	NP-Hardness	70
4.3	Algorithm	71
4.3.1	Preliminaries	71
4.3.2	MPTC Mining Framework	72
4.3.3	Generating Candidates	74
4.3.4	Constructing Tiles	75
4.4	Probabilistic Price Order	76
4.5	Optimization Techniques	80
4.5.1	Optimizing Single Transaction Difference	80
4.5.2	Adaptively Computing Cover Quantity	81
4.5.3	Pruning by 3σ Property	82
4.6	Algorithm Analysis	83
4.6.1	Appropriateness of the Greedy Strategy	83
4.6.2	Time Complexity	85

4.7	Performance Study	87
4.7.1	Experiments on Synthetic Datasets	87
4.7.2	Experiments on Real World Datasets	93
4.8	Related Work	99
4.8.1	Transaction data summarization	100
4.8.2	Frequent pattern summarization	101
4.9	Conclusions	101
5	RCP Mining: Towards the Summarization of Spatial Co-location Patterns	103
5.1	Introduction	104
5.2	Preliminary	108
5.2.1	Co-location Patterns	108
5.2.2	Co-location Distance Measure	109
5.2.3	Problem Statement	111
5.3	The <i>RCPFast</i> Algorithm	112
5.4	The <i>RCPMS</i> Algorithm	116
5.4.1	Optimization Strategy	118
5.4.2	Approximation Strategy	122
5.4.3	The <i>gen_cover_set()</i> Function	124
5.5	Experimental Study	125
5.5.1	Experiments on Synthetic Data	125
5.5.2	Experiments on Real Data	129
5.6	Related Work	131
5.7	Conclusions	132
6	Conclusion	141
	References	143

List of Tables

1.1	Example transaction databases.	2
2.1	An uncertain database with attribute uncertainty	15
2.2	An example of possible worlds	15
2.3	Probability of different situations in the j th transaction	22
3.1	An example of attribute uncertainty.	40
3.2	An example of possible worlds.	40
3.3	All possible situations of X_1 and X_2 in t_i	48
3.4	Characteristics of Datasets.	53
4.1	Examples of transaction databases.	62
4.2	An example of possible worlds.	67
4.3	Parameters used in experiments.	89
4.4	Probabilistic F1-score w.r.t. the number of transactions n	91
4.5	Probabilistic F1-score w.r.t. the number of items m	92
4.6	Probabilistic F1-score w.r.t. the number of ground truth tiles k	93
4.7	Probabilistic F1-score w.r.t. the probability of noise ϵ	94
4.8	Characteristics of datasets.	94
5.1	A set of prevalent co-location patterns.	107
5.2	Parameters used in synthetic data generation	133

List of Figures

2.1	The Number of P-RFP on Retail	27
2.2	The Number of P-RFP on IIP	27
2.3	Runtime on Retail	28
2.4	Runtime on IIP	28
2.5	Effect of Optimization	29
3.1	Empirical proof of approximation.	52
3.2	The Number of P-RFP on IIP-500.	56
3.3	The Number of P-RFP on Retail-500.	56
3.4	Log Runtime on IIP-500.	57
3.5	Log Runtime on Retail-500.	57
3.6	The Number of P-RFP on IIP.	58
3.7	The Number of P-RFP on Retail.	58
3.8	The Number of P-RFP on Chess.	59
3.9	Runtime on IIP.	59
3.10	Runtime on Retail.	60
3.11	Runtime on Chess.	60
4.1	Performance w.r.t. λ on Equip.	96
4.2	Performance w.r.t. <i>minsup</i> on Equip.	97
4.3	Performance w.r.t. λ on IIP.	97
4.4	Performance w.r.t. <i>minsup</i> on IIP.	98
4.5	Number of Tiles on Equip and IIP.	98
4.6	Performance of optimization techniques.	99
4.7	Number of candidates	99
4.8	Cost and Time w.r.t. u	100

5.1	A motivating example of a spatial data set. Each symbol represents an event corresponding to a spatial feature, and each edge connecting two events represents a neighborhood relationship.	106
5.2	An example illustrating <i>RCPFast</i> algorithm.	114
5.3	An illustration of the optimization strategy based on Theorem 5.3. Each F_i is a co-location pattern of size i . Different types of lines represent different ways of obtaining the co-location distance. Suppose $d_1 + d_2 + d_3 \leq \varepsilon$ and $d_1 + d_2 + d_3 + d_4 > \varepsilon$	121
5.4	Examples of the approximation strategy.	123
5.5	Compression rate tests on synthetic data sets.	134
5.6	Framework comparison on synthetic data sets.	135
5.7	Performance tests with <i>minpi</i> and ε on synthetic data sets.	136
5.8	Co-location distance computation analysis on synthetic data sets.	137
5.9	Compression rate differences between <i>RCPMS</i> and <i>RCPFast</i> on synthetic data sets.	138
5.10	Compression rate tests on EPA and POI data sets.	138
5.11	Performance on EPA and POI data sets.	139