

Structural Community Detection in Big Graphs

by

DONG WEN

A THESIS SUBMITTED IN PARTIAL FULFILLMENT OF
THE REQUIREMENTS FOR THE DEGREE OF
DOCTOR OF PHILOSOPHY



Centre for Artificial Intelligence (CAI)
Faculty of Engineering and Information Technology (FEIT)
University of Technology Sydney (UTS)

June, 2018

CERTIFICATE OF AUTHORSHIP/ORIGINALITY

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.

Production Note:

Signature removed prior to publication.

Signature of Candidate



ACKNOWLEDGEMENTS

First, I would like to deliver my sincere gratitude to my supervisor Dr. Lu Qin for his continuous support and guidance of my PhD career. He is professional, efficient, patient and diligent. His valuable ideas always guided me and extended my knowledge not only in the database area, but also about the practical research skills. Additionally, Lu is a good mentor and friend for me. He is selfless and glad to share his own experience with me. Thanks to his encouragement, I am always positive and brave when experiencing challenges and even failures. This thesis could not reach its present form without his illuminating instructions.

Secondly, I would like to express my great gratitude to my co-supervisor Prof. Ying Zhang for his guidance and advice, especially for his strong confidence in me and support for my academic career. He gave me many wonderful ideas and inspirations. His guidance significantly extended my knowledge and helped me pursue a correct research direction all the time. He is also visionary, and often encouraged and guided me to optimize the plan for the career development.

Thirdly, I would like to thank Prof. Xuemin Lin and Dr. Lijun Chang for supporting the works in this thesis, as most of the works were conducted in collaboration with them. I thank Prof. Lin for offering an interesting but rigorous

research environment. I learned the characteristics of an excellent researcher from Prof. Lin — passion, preciseness and earnest. I thank Dr. Chang for his brilliant ideas and confidence in me. His wonderful research works always inspired me, and he gave me many accurate and valuable suggestions to improve the quality of the works in this thesis.

I would also like to thank the following people: Prof. Jeffrey Xu Yu, Dr. Wenjie Zhang, Dr. Xin Cao, Dr. Zengfeng Huang, Dr. Ling Chen, Dr. Xiaoyang Wang, Dr. Shiyu Yang, for sharing valuable ideas and experiences. Thanks to Dr. Fan Zhang, Dr. Long Yuan, Dr. Longbin Lai, Dr. Xing Feng, Dr. Xiang Wang, Dr. Jianye Yang, Dr. Fei Bi, Mr. Xubo Wang, Mr. Wei Li, Mr. Haida Zhang, Ms. Chen Zhang, Dr. Qing Bing, Dr. Shenlu Wang, Mr. Yang Yang, Mr. Kai Wang, Mr. You Peng, Mr. Hanchen Wang, Mr. Xuefeng Chen, Mr. Boge Liu, Mr. Jiahui Yang, Dr. Wei Wu, Dr. Qian Zhang, Dr. Zhaofeng Su, Mr. Dian Ouyang, Mr. Wentao Li, Mr. Bohua Yang, Mr. Mingjie Li, Ms. Conggai Li, Mr. Wei Song, for sharing the happiness with me in my PhD career.

Last but not least, I would like to thank my father Mr. Rongxiao Wen and my mother Mrs. Tiefang Liu, for bringing me a wonderful life, and other relatives for their understanding, encouragement and love. Thanks to my wife Mrs. Qing Lin, for her company and support in my PhD study.

ABSTRACT

Community detection in graphs is a fundamental problem widely experienced across industries. Given a graph structure, one popular method to identify communities is classifying the vertices, which is formally named graph clustering. Additionally, community structures are always dense and highly connected in graphs. There are also a large number of research works focusing on mining cohesive subgraphs for community detection. Even though the community detection problem is extensively studied, challenges still exist. With the development of social media, graphs are highly dynamic, and the size of graphs is sharply increasing. The large time and space cost of traditional solutions may hardly be endured in big and dynamic graphs.

In this thesis, we propose an index-based algorithm for the structural graph clustering (SCAN). Based on the proposed index structure, the time expended to compute structural clustering depends only on the result size, not on the size of the original graph. The space complexity of the index is bounded by $O(m)$, where m is the number of edges in the graph. We also propose algorithms and several optimization techniques for maintaining our index in dynamic graphs.

For the cohesive subgraph detection, we study both degree-constrained (k -core) and connectivity-constrained (k -VCC) cohesive subgraph metrics. A k -core is a maximal connected subgraph in which each vertex has degree at least k . We study I/O efficient core decomposition following a semi-external model, which only allows vertex information to be loaded in memory. We propose an optimized

I/O efficient algorithm for both core decomposition and core maintenance. In addition, we extend our algorithm to compute the graph degeneracy order, which is an important graph problem that is highly related to core decomposition.

A k -vertex connected component (k -VCC) is a connected subgraph in which the removal of any $k - 1$ vertices will not disconnect the subgraph. A k -VCC has many outstanding structural properties, such as high cohesiveness, high robustness, and subgraph overlapping. The state-of-the-art solution enumerates all k -VCCs following a partition-based framework. It requires high computational cost in connectivity tests. We prove the upper bound of the number of partitions, which implies the polynomial running time of this framework. We propose two effective optimization strategies, namely neighbor sweep and group sweep, to largely reduce the number of local connectivity tests.

We conducted extensive performance studies using several large real-world datasets to show the efficiency and effectiveness of all our approaches.

PUBLICATIONS

- *Dong Wen, Lu Qin, Xuemin Lin, Ying Zhang, and Lijun Chang, Enumerating k -Vertex Connected Components in Large Graphs, in submission.*
- *Dong Wen, Lu Qin, Ying Zhang, Xuemin Lin, and Jeffrey Xu Yu, I/O Efficient Core Graph Decomposition: Application to Degeneracy Ordering, TKDE, revision submitted.*
- *Dong Wen, Lu Qin, Ying Zhang, Lijun Chang, and Xuemin Lin, Efficient Structural Graph Clustering: an Index-Based Approach, to appear in PVLDB 2018.*
- *Lingxi Yue, Dong Wen, Lizhen Cui, Lu Qin, and Yongqing Zheng, K -Connected Cores Computation in Large Dual Networks, to appear in DAS-FAA 2018.*
- *Dong Wen, Lu Qin, Ying Zhang, Xuemin Lin, and Jeffrey Xu Yu, I/O Efficient Core Graph Decomposition at Web Scale, ICDE2016, Best Paper Award.*

TABLE OF CONTENT

CERTIFICATE OF AUTHORSHIP/ORIGINALITY	ii
ACKNOWLEDGEMENTS	iii
ABSTRACT	v
PUBLICATIONS	vii
TABLE OF CONTENT	viii
LIST OF FIGURES	xi
LIST OF TABLES	xiii
Chapter 1 INTRODUCTION	1
1.1 Structural Graph Clustering	2
1.2 Degree-Based Cohesive Subgraph Detection	7
1.3 Connectivity-Based Cohesive Subgraph Detection	11
1.4 Graph Model	15
Chapter 2 LITERATURE REVIEW	18
2.1 Graph Clustering	18
2.2 Cohesive Subgraph Detection	19
2.2.1 Global Cohesiveness	19
2.2.2 Local Degree and Triangulation	20
2.2.3 Connectivity Cohesiveness	21
Chapter 3 STRUCTURAL GRAPH CLUSTERING BASED COMMUNITY DETECTION	23
3.1 Overview	23
3.2 Preliminary	24
3.3 Existing Solutions	27
3.3.1 SCAN	27

3.3.2	pSCAN	28
3.4	Index-Based Algorithms	30
3.4.1	Index Overview	31
3.4.2	Index Construction	37
3.4.3	Query Processing	39
3.5	Index Maintenance	44
3.5.1	Basic Solutions	44
3.5.2	Improved Algorithms	48
3.6	Performance Studies	52
3.6.1	Performance of Query Processing	53
3.6.2	Performance of Index Construction	56
3.6.3	Performance of Index Maintenance	59
3.7	Chapter Summary	61
Chapter 4 DEGREE-CONSTRAINED COMMUNITY DETECTION		63
4.1	Chapter Overview	63
4.2	Preliminary	64
4.3	Existing Solutions	66
4.4	I/O Efficient Core Decomposition	69
4.4.1	Basic Semi-external Algorithm	69
4.4.2	Optimal Vertex Computation	73
4.5	I/O Efficient Core Maintenance	77
4.5.1	Edge Deletion	77
4.5.2	Edge Insertion	79
4.5.3	Optimization for Edge Insertion	82
4.6	I/O Efficient Degeneracy Ordering	87
4.6.1	Degeneracy Order Computation	89
4.6.2	Degeneracy Order Maintenance	91
4.7	Performance Studies	96
4.7.1	Core Decomposition	98
4.7.2	Core Maintenance	100
4.7.3	Degeneracy Order Computation	101
4.7.4	Degeneracy Order Maintenance	102
4.7.5	Scalability Testing	103
4.8	Chapter Summary	107
Chapter 5 CONNECTIVITY-CONSTRAINED COMMUNITY DETECTION		108
5.1	Overview	108
5.2	Preliminary	109

TABLE OF CONTENT

5.3	Algorithm Framework	109
5.4	Basic Solution	112
5.4.1	Find Vertex Cut	112
5.4.2	Algorithm Analysis	116
5.5	Search Reduction	119
5.5.1	Neighbor Sweep	120
5.5.2	Group Sweep	127
5.5.3	The Overall Algorithm	131
5.6	Performance Studies	134
5.6.1	Performance Studies on Real-World Graphs	136
5.6.2	Evaluating Optimization Techniques	138
5.6.3	Scalability Testing	141
5.7	Chapter Summary	142
Chapter 6 EPILOGUE		143
BIBLIOGRAPHY		145

LIST OF FIGURES

1.1	Clusters, hubs and outliers under $\epsilon = 0.7, \mu = 4$	3
1.2	Cohesive subgraphs in graph G .	12
3.1	Clusters, hubs and outliers under $\epsilon = 0.7, \mu = 4$	26
3.2	Clusters under different μ and ϵ	29
3.3	Neighbor-order for each vertex in graph G	36
3.4	Core-order for each μ in graph G	37
3.5	A running example for $\epsilon = 0.7, \mu = 4$	42
3.6	Query time for different ϵ ($\mu = 5$)	54
3.7	Query time for different μ ($\epsilon = 0.6$)	55
3.8	Query time on different datasets	56
3.9	Index size for different datasets	56
3.10	Time cost for index construction	57
3.11	Index construction (vary $ V $)	58
3.12	Index construction (vary $ E $)	58
3.13	Time cost for edge insertion	59
3.14	Time cost for edge removal	60
3.15	Index maintenance (vary $ V $)	61
3.16	Index maintenance (vary $ E $)	62
4.1	A sample graph G and its core decomposition	65
4.2	Number of vertices whose core numbers are changed	73
4.3	A degeneracy order of vertices in Fig. 4.3	88
4.4	A Level-Index for graph G in Fig. 4.1	92
4.5	Core decomposition on different datasets	99
4.6	Core maintenance on different datasets	101
4.7	Time cost and I/Os of computing degeneracy order	102
4.8	Time cost and I/Os of degeneracy order maintenance	103
4.9	Scalability of core decomposition	104
4.10	Scalability of core maintenance	105
4.11	Scalability of degeneracy ordering	106
4.12	Scalability of degeneracy order maintenance	107

LIST OF FIGURES

5.1	An example of overlapped graph partition.	112
5.2	The sparse certificate of given graph G with $k = 3$	116
5.3	Strong side-vertex and vertex deposit when $k = 3$	122
5.4	Increasing deposit with neighbor and group sweep	127
5.5	Performance on different datasets	136
5.6	Against basic algorithm (vary k)	137
5.7	Scalability testing	141

LIST OF TABLES

1.1	Notations	16
3.1	Network statistics	53
4.1	Illustration of SemiCore	72
4.2	Illustration of SemiCore [*]	77
4.3	Illustration of SemiDelete [*] (delete (v_0, v_1))	78
4.4	ILLUSTRATION OF SemiInsert (INSERT (v_4, v_6))	81
4.5	ILLUSTRATION OF SemiInsert [*] (INSERT (v_4, v_6))	86
4.6	Illustration of DInsert [*] (insert (v_0, v_6))	96
4.7	Network statistics (1K = 10^3)	98
5.1	Network statistics	135
5.2	Evaluating pruning rules	139

LIST OF TABLES
