# Cohesive Subgraph Mining on Social Networks

*by*

**Fan Zhang**

B.E. ZHEJIANG UNIVERSITY, 2014

the Centre for Artificial Intelligence (CAI)

the Faculty of Engineering and Information Technology (FEIT)

the University of Technology Sydney (UTS)

August, 2017

# 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.

Signature of Candidate

**UNIVERSITY OF TECHNOLOGY SYDNEY**

# ACKNOWLEDGEMENTS

First and foremost, I would like to deliver my sincere gratitude to my supervisor Prof. Ying Zhang for his continuous support of my PhD study and research, especially for his professionalism, patience, passion and diligence. His guidance extends my knowledge in computer science, improves my capacity in scientific research and elevates my love in exploring the fields which are significant, undiscovered and challenging. Besides the character of supervisor, Ying has also been a friend and mentor. Thanks to his confidence in me, I never lost hope when experiencing failures and was always positive during my PhD study. Without his consistent and illuminating instruction, this thesis could not have reached its present form.

Secondly, I would like to express my great gratitude to my co-supervisor Dr. Lu Qin for his constant encouragement and guidance, especially for his brilliant ideas and inspirations. His work efficiency gave me the hope to conduct good research without abandoning too many of other interests, which prevented me to be negative in the early stage of my PhD study. Lu always has confidence in solving research problems, regardless of their complexities, which encourages me to keep thinking and challenging myself.

# PUBLICATIONS

- **Fan Zhang**, Ying Zhang, Lu Qin, Wenjie Zhang, Xuemin Lin. When Engagement Meets Similarity: Efficient (k,r)-Core Computation on Social Networks. PVLDB 2017. (Chapter 3)

- **Fan Zhang**, Wenjie Zhang, Ying Zhang, Lu Qin, Xuemin Lin. OLAK: An Efficient Algorithm to Prevent Unraveling in Social Networks. PVLDB 2017. (Chapter 4)

- **Fan Zhang**, Ying Zhang, Lu Qin, Wenjie Zhang, Xuemin Lin. Efficiently Reinforcing Social Networks over User Engagement and Tie Strength. Under submission. (Chapter 4)

- **Fan Zhang**, Ying Zhang, Lu Qin, Wenjie Zhang, Xuemin Lin. Finding Critical Users for Social Network Engagement: The Collapsed k-Core Problem. AAAI 2017. (Chapter 5)

# TABLE OF CONTENT

# LIST OF FIGURES

# LIST OF TABLES

# ABSTRACT

Graphs are widely used to represent the abundant information in social networks for discovering promising communities, reinforcing network stability, and finding critical users, to name a few. Cohesive subgraph mining, as one of the most fundamental problems in graphs, gains increasing popularity in social network study for its effectiveness. In this thesis, some basic social components are considered in cohesive subgraphs to better accommodate various real-life applications.

Firstly, we investigate the problem of $(k,r)$-core which intends to find cohesive subgraphs on social networks considering both user engagement and similarity. Efficient algorithms are proposed to enumerate all *maximal* $(k,r)$-cores and find the *maximum* $(k,r)$-core, where both problems are shown to be NP-hard. Effective pruning techniques and search orders substantially reduce the search space of two algorithms. A novel upper bound enhances performance of the maximum $(k,r)$-core computation. Comprehensive experiments on real-life data demonstrate that the algorithms efficiently find interesting communities.

Secondly, we study the problem of the anchored $k$-core, which was introduced by Bhawalkar and Kleinberg *et* al. in the context of user engagement in social networks. The problem has been shown to be NP-hard and inapproximable. We propose an efficient algorithm, namely `OLAK`, as the first to solve the problem on general graphs. An *onion layer* structure is designed together with efficient

candidates exploration, early termination and pruning techniques to significantly simplify computation and greatly reduce the search space.

Besides considering user engagement, we further explore the unraveling phenomenon with tie strength, which leads us to the model of $k$-truss. We then investigate the anchored $k$-truss problem which is also NP-hard and propose an *edge onion layer* structure based algorithm, namely AKT. Efficient candidate exploration and pruning techniques are designed based on the *edge onion layers*. Comprehensive experiments on real-life graphs for the above two problems demonstrate the effectiveness and efficiency of our proposed methods.

Finally, we study the leave of critical users, which may greatly break network engagement. Accordingly, we propose the collapsed $k$-core problem to find the vertices whose leave can lead to the smallest $k$-core. We prove the problem is NP-hard. Then, an efficient algorithm is proposed, which significantly reduces the number of candidate vertices to speed up computation. Comprehensive experiments on real-life social networks demonstrate effectiveness of the model and efficiency of the proposed techniques.