UNIVERSITY OF TECHNOLOGY SYDNEY

Faculty of Engineering and Information Technology

# Local Information and Structures in Analysis and Modelling of Complex Networks

by

**Mingshan Jia**

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

**Doctor of Philosophy**

Sydney, Australia

July 2022

# Certificate of Authorship/Originality

I, Mingshan Jia, declare that this thesis, is submitted in fulfilment of the requirements for the award of Doctor of Philosophy, in the Faculty of Engineering and Information Technology at the University of Technology Sydney.

This thesis is wholly my work unless otherwise referenced or acknowledged. In addition, I certify that all information sources and literature used are indicated in the thesis.

This document has not been submitted for qualifications at any other academic institution.

Mingshan Jia

Production Note:
Signature: Signature removed
prior to publication.

Date: July 2022

# ABSTRACT

## Local Information and Structures in Analysis and Modelling of Complex Networks

by

Mingshan Jia

Abstracting entities and their interactions as nodes and links, networks are a general representation for modelling and studying complex systems. Modelling relational structures of the underlying data, rather than only a set of isolated entities, allows us to build more accurate models for various types of domain data, such as social relationships, molecular interactions, program executions, and many more. Despite being powerful and ubiquitous, networks are also difficult to process, mainly due to their complex topological structures. Therefore, the study of network structure, especially local structure, has been the core theme of studying complex networks. This dissertation aims to provide new understandings of how local structure information is extracted and utilised in studying different types of complex networks.

The dissertation includes three original works in the direction of local structure and information on top of a comprehensive survey. In the review, we propose new taxonomies for graph structures that bring together the notions of centrality measures, motifs, and other local-level metrics. For theoretical understanding, we propose new metrics to quantify the formation of 3-node and 4-node subgraphs and develop new motif patterns that are distinctive features in both network- and node-level analysis. For methodological approaches, we propose the framework to effectively encode edge attributes into the typed-edge graphlet degree vector, for both sociocentric and egocentric networks. Moreover, for practical applications, the proposed metrics and approaches are applied in many different types of complex net-

works and case studies. They are not only proven to be effective in multiple learning and analytical tasks but also lead to new insights and interesting discoveries.

Dissertation directed by Professor Katarzyna Musial-Gabrys and Professor Bogdan Gabrys

School of Computer Science, Data Science Institute, Complex Adaptive Systems Lab, UTS

# Acknowledgements

Although pursuing a PhD is not an easy path to take, looking back on the past three years and nine months, it might be one of the best decisions I have ever made. I am grateful that I did not give up in some very tough situations. It eventually led me to meet my current supervisors and other wonderful people, who have advised me, supported me, worked with me and motivated me.

First and foremost, I am eternally grateful to my supervisors, — Professor Katarzyna Musial-Gabrys and Professor Bogdan Gabrys, for all the support and guidance they have provided throughout my PhD study. They gave me the opportunity to continue my PhD at my most difficult time, and they have supported me, encouraged me and trusted me ever since. They introduced me to the world of network science and provided me with all kinds of advice, from idea selection and experiment design to paper writing and rebuttal. They not only helped me improve in academics, but also recommended me to participate in collaborative research projects with other universities, and gave me multiple opportunities to participate in teaching activities. Without them, I could not have achieved such an all-round growth. Thank you, Professor Katarzyna Musial-Gabrys and Professor Bogdan Gabrys.

Moreover, I would like to thank all the people who have advised me, worked with me and helped me along my PhD journey. I want to thank my co-authors, Maité Van Alboom, Liesbet Goubert and Piet Bracke. Thanks to you and my supervisors, we have conducted an exciting cross-disciplinary study about chronic pain. I want to thank Professor Wei Liu and Dr Yi Zhang for being on my candidature assessment panel. I also thank Professor Pasquale De Meo for the interesting and insightful discussions on multiple research topics.

Next, I want to express my gratitude to my colleagues and friends, who have

# List of Publications

**Journal Papers**

J-1. **M. Jia**, B. Gabrys and K. Musial, "Directed closure coefficient and its patterns," in Plos one 16.6 (2021): e0253822.

J-2. **M. Jia**, B. Gabrys and K. Musial, "Measuring Quadrangle Formation in Complex Networks," in IEEE Transactions on Network Science and Engineering, vol. 9, no. 2, pp. 538-551, 1 March-April 2022.

**Conference Papers**

C-1. **M. Jia**, B. Gabrys and K. Musial, "Closure Coefficient in Complex Directed Networks." International Conference on Complex Networks and Their Applications. Springer, Cham, 2020.

C-2. **M. Jia**, M. Van Alboom, L. Goubert, P. Bracke, B. Gabrys, K. Musial. "Analysing Ego-Networks via Typed-Edge Graphlets: A Case Study of Chronic Pain Patients." International Conference on Complex Networks and Their Applications. Springer, Cham, 2021.

C-3. **M. Jia**, M. Van Alboom, L. Goubert, P. Bracke, B. Gabrys, K. Musial. "Analysing Egocentric Networks via Local Structure and Centrality Measures: A Study on Chronic Pain Patients." 2022 International Conference on Information Networking (ICOIN). IEEE, 2022.

# Contents

# List of Figures