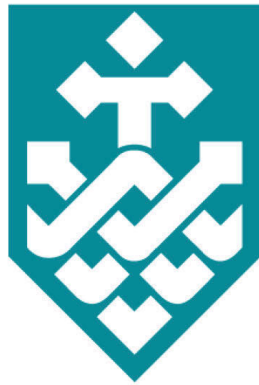# Causal and Causally-inspired Learning



Mingming Gong

Faculty of Engineering and Information Technology

University of Technology Sydney

A thesis submitted for the degree of

*Doctor of Philosophy*

2017

To my loving parents, wife, and son.

# 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.

Mingming Gong
28/03/2017

# Acknowledgements

I would like to sincerely thank everyone who has helped me to finish my doctoral studies.

First of all, I would like to express my sincere appreciation and deep gratitude to my supervisor **Prof. Dacheng Tao**. He has given me trust and freedom to pursue my research interests, and provided constructive suggestions to help me out of difficulties. I can always benefit and learn a lot from various detailed discussions with him, and be excited and energised by his amazing insight, unlimited patience, generous support, and constant encouragement. I feel very lucky to have had him as my supervisor.

I also wish to express my sincere appreciation to my co-supervisor **Dr. Kun Zhang** who hosted my visit to Max Planck Institute for Intelligent Systems. He has been leading me and my research with his passion, patience, optimism, intelligence, and creativity. His research styles and way of thinking have influenced me very deeply.

I also would like to give special thanks to my excellent collaborators: Prof. Bernhard Schölkopf, Prof. Clark Glymour, Dr. Dominik Janzing, Prof. Yuhong Wang, Prof. Changyin Sun, A/Prof. Di Huang, A/Prof. Junliang Xing, A/Prof. Wankou Yang, Dr. Chaohui Wang, Dr. Tongliang Liu, Dr. Chang Xu, Mr. Philipp Geiger, Mr. Ruxing Wang, Mr. Qiang Li, Mr. Shaoli Huang, Mr. Huan Fu, for their brilliant work and timely support. I have also been fortunate to work and have discussions with many other brilliant researchers: Dr. Nannan Wang, Dr. Fei Gao, Prof. Chen Gong, Dr. Yong Luo, Dr. Lianyang Ma, Dr. Weilong Hou, A/Prof. Bo Du, A/Prof. Shigang Liu, A/Prof.

# Abstract

A main goal of statistics and machine learning is to discover statistical dependencies between random variables, and these dependencies will be used to perform predictions on future observations. However, many scientific investigations involve causal predictions, the aim of which is to infer how the data generating system should behave under changing conditions, for example, changes induced by external interventions. To perform causal predictions, we need both statistical dependencies as well as causal structures to determine the behaviour of the system. The standard way to identify causal structures is to use randomized controlled experiments. However, conducting these experiments is usually expensive or even impossible in many scenarios. As a consequence, inferring cause and effect relationships from purely observational data, known as causal discovery or causal learning, has drawn much attention.

Various causal discovery methods have been proposed in the past decades, including constraint-based methods, structural equation models-based methods, and time series-based methods. Among these methods, time series-based methods, e.g., Granger causality, are relatively well-established as the temporal information excludes the case that effects happen before causes. Many of the existing time series-based methods assume that the data are measured at the right frequency; however, in practice the sampling frequency of the data is often lower than the true causal frequency. In this thesis, we consider learning high-resolution causal relationships at the causal frequency from subsampled time series. Existing methods suffer from the identifiability problems: under the Gaussianity assumption of the data, the solutions are generally not unique. We prove that, however, if the noise

terms are non-Gaussian, the underlying model is identifiable from subsampled time series under mild conditions. We then propose an Expectation-Maximization approach and a variational inference approach to recover causal relations from subsampled data.

More recently, researchers began to touch upon implications of causal models for machine learning tasks such as semi-supervised learning and domain adaptation. In this thesis, we develop causally-inspired learning methods for domain adaptation in both multi-source and single-source settings. In particular, we use causal models to represent the relationship between the features and labels, and consider possible situations where different modules of the causal model change with the domain. In each situation, we investigate what knowledge is appropriate to transfer and find the optimal target-domain hypothesis. Furthermore, we propose methods to correct distribution shift in the general situation where the marginal distribution of features and conditional distribution of labels given features both change, under the assumption that labels are causes for features. We provide theoretical analysis and empirical evaluation on both synthetic and real-world data to show the effectiveness of our methods.

# Contents

# List of Figures

# Nomenclature

## Abbreviations

| | |
|---|---|
| EM | expectation maximization |
| DAG | directed acyclic graph |
| AR | autoregressive |
| BN | Bayesian network |
| CBN | causal Bayesian network |
| SEM | structural equation model |
| RCT | randomized controlled experiments |
| PC | Peter-Clark |
| GES | greedy equivalence search |
| DA | domain adaptation |
| SSL | semi-supervised learning |
| MCMC | Markov chain Monte Carlo |
| ICA | independent component analysis |
| LS | location-scale |
| MMD | maximum mean discrepancy |
| IC | invariant components |
| CIC | conditional invariant components |
| CTC | conditional transferable components |
| RKHS | reproducing kernel Hilbert space |