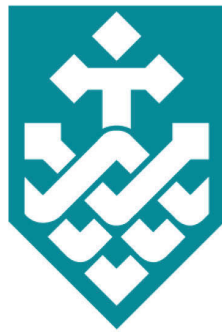# Multinomial Latent Logistic Regression

Zhe Xu

Faculty of Engineering and Information Technology

University of Technology, Sydney

A thesis submitted for the degree of

*Doctor of Philosophy*

November 2016

# 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: Zhe Xu

Date: 21/10/2016

I would like to dedicate this thesis to my loving parents

*Nan Huang* and *Yuan Xu*

# Acknowledgements

I would like to take this good opportunity to appreciate my advisors, colleagues, friends and also my family for their significant help during my doctoral study in University of Technology, Sydney.

First of all, I would like to express my sincere appreciation and deep gratitude to my advisor supervisor **Prof. Dacheng Tao** for his unlimited patience, generous support, and supportive guidance. I'm particularly impressed by his incredible enthusiasm and high standard in academic research, which encourage me to breakthrough my own setting limits and submit papers to the leading journals or conferences in my research field.

I would also like to thank my advisor **Prof. Ya Zhang** from Shanghai Jiao Tong University. She always gives me plenty of freedom to explore and timely constructive suggestions to help me out of difficulties. Without the effect by her, Prof. Xiaokang Yang and Prof. Chengqi Zhang, I would not have the opportunity to study here in UTS as a dual-degree PhD student.

I have been fortunate to work in UTS and Centre for Quantum Computation and Intelligent Systems (QCIS) directed by Prof. Chengqi Zhang. QCIS provides full support for me to attend top conferences including ECCV, ICCV and KDD, where I got the opportunities to learn from many world-famous experts. It's really a pleasure for me to work in thus a great team and around so many brilliant minds in QCIS. Studying in QCIS and also UTS will be a fantastic memory that I will never forget.

I am also deeply indebted to Dr. Jun Zhu who led me into the field of computer vision. He shows me how interesting my research objective

is, which encourages me constantly during the journey of exploration in the following years. Moreover, I also want to give special thanks to my excellent collaborators: Zhibin Hong and Shaoli Huang for their brilliant work and timely support, and also to my dear colleagues and friends I met in QCIS: Dr. Bozhong Liu, Chunyang Liu, Meng Fang, Tongliang Liu, Mingming Gong, Maoying Qiao, Qiang Li, Runxin Wang, Changxing Ding, Zhiguo Long, Caishi Fang, Wei Bian, Shirui Pan, Yong Luo, Xiao Liu, Dianshuang Wu, Weilong Hou, Chang Xu, Chen Gong, Sujuan Hou, Haishuang Wang, Jia Wu, Wei Yang, Qin Zhang, Yali Du, Hao Xiong, Jiankang Deng, Xiao Liu, Peng Hao, Liu Liu, Guodong Long, Jing Jiang, Peng Zhang, Barbara Munday, Prof. Bo Du, Prof. Xianhua Ben, Prof. Wankou Yang, Prof. Shigang Liu, Prof. Xianhua Zeng, for the inspiring discussions, kind support and companionship.

I am also grateful to all the other friends: Guanbo Huang, Wei Xu, Weiyuan Chen, Xiaohang Ren, Zhiyi Tan, Jialin Li, Qing Wang, Weiyuan Chen, for their support and company during both joyful and stressful times.

Finally, I would like to express my deeply felt gratitude to my parents, who never excoriate me when I make mistakes but show me how to do it the right way, who never asks for anything but give me everything they have, who gives me such a wonderful place to grow up. It's you who make me the man who I am now. Thank you, from the bottom of my heart.

# Abstract

We are arriving at the era of big data. The booming of data gives birth to more complicated research objectives, for which it is important to utilize the superior discriminative power brought by explicitly designed feature representations. However, training models based on these features usually requires detailed human annotations, which is being intractable due to the exponential growth of data scale.

A possible solution for this problem is to employ a restricted form of training data, while regarding the others as latent variables and performing latent variable inference during the training process. This solution is termed *weakly supervised learning*, which usually relies on the development of latent variable models. In this dissertation, we propose a novel latent variable model - **multinomial latent logistic regression (MLLR)**, and present a set of applications on utilizing the proposed model on weakly supervised scenarios, which, at the same time, cover multiple practical issues in real-world applications.

We first derive the proposed MLLR in Chapter 3, together with theoretical analysis including the concave and convex property, optimization methods, and the comparison with existing latent variable models on structured outputs. Our key discovery is that by performing "maximization" over latent variables and "averaging" over output labels, MLLR is particularly effective when the latent variables have a large set of possible values or no well-defined graphical structure is existed, and when probabilistic analysis is preferred on the output predictions. Based on it, the following three sections will discuss the application of MLLR in a variety of tasks on weakly supervised learning.

In Chapter 4, we study the application of MLLR on a novel task of architectural style classification. Due to a unique property of this

task that rich inter-class relationships between the recognizing classes make it difficult to describe a building using "hard" assignments of styles, MLLR is believed to be particularly effective due to its ability to produce probabilistic analysis on output predictions in weakly supervised scenarios. Experiments are conducted on a new self-collected dataset, where several interesting discoveries on architectural styles are presented together with the traditional classification task.

In Chapter 5, we study the application of MLLR on an extreme case of weakly supervised learning for fine-grained visual categorization. The core challenge here is that the inter-class variance between subordinate categories is very limited, sometimes even lower than the intra-class variance. On the other hand, due to the non-convex objective function, latent variable models including MLLR are usually very sensitive to the initialization. To conquer these problems, we propose a novel multi-task co-localization strategy to perform warm start for MLLR, which in turn takes advantage of the small inter-class variance between subordinate categories by regarding them as related tasks. Experimental results on several benchmarks demonstrate the effectiveness of the proposed method, achieving comparable results with latest methods with stronger supervision.

In Chapter 6, we aim to further facilitate and scale weakly supervised learning via a novel knowledge transferring strategy, which introduces detailed domain knowledge from sophisticated methods trained on strongly supervised datasets. The proposed strategy is proved to be applicable in a much larger web scale, especially accounting for the ability of performing noise removal with the help of the transferred domain knowledge. A generalized MLLR is proposed to solve this problem using a combination of strongly and weakly supervised training data.

# Contents

# List of Figures

# List of Tables