# Learning with Restricted Data via Gradient Manipulations

**by Jing Li**

Thesis submitted in fulfilment of the requirements for the degree of

**Doctor of Philosophy**

under the supervision of Ivor W. Tsang

# Certificate of Original Authorship

I, Jing Li, 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 own 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.

Signature:

Production Note:
Signature removed prior to publication.

Date:    20 Oct 2022

This thesis is lovingly dedicated to my family. I wish I would be always your pride and joy.

# Acknowledgements

Foremost, I would like to thank my supervisor Prof. Ivor W. Tsang, for his professional guidance and persistent support. He opened up a window for me to have a closer look at the world of *machine learning* which appears complex but fascinating! Without his leading in the passed four-year academic journey, I could not have overcome the obstacles nor found my way to think and do. I am also grateful that Prof. Ivor W. Tsang could always provide appropriate suggestions on my study plan, from each individual project to internship scheduling and thesis preparation. He is such a caring mentor who helps me become confident, mature and independent. What I have learned from him will have a great impact on my future career path.

Thanks also go to my co-supervisor Prof. Yulei Sui who provided a perspective for my work from a security player. His participation makes my work more practical and understandable for people out of the machine learning community. In particular, I appreciate Prof. Yulei Sui's contribution to polishing my papers.

I would like to sincerely thank Dr. Yuangang Pan for his unselfish assistance on my study. His research passion and attitude has encouraged me to keep exploring the interesting topics. I would also thank my best friend Yinghua Yao, who is not only a good flatmate but also reliable co-worker. From Sydney to Singapore, I am fortunate that people around me are helpful. I appreciate Dr. Yan Zhang's invitation last winter and the contribution of Dr. Yueming Lyu in my recent work. I am also grateful for the assistance of Dr. Yaxin Shi, Xiaowei Zhou, Xingrui Yu on my research. My warmest gratitude also goes the Dr Jiangchao Yao, Dr. Xu Chen, Dr. Xiaofeng Xu, Zhuanghua Liu, Jinliang Deng, Feiyang Ye, Bowen Xing, Peiyao Zhao, Yujie Fang, and Cheng Chen. Thanks also to the co-workers out of my group, Shaojun Shi and Tao Zhang.

Many thanks to the University of Technology Sydney for providing me a good studying environment, and Agency of Science, Technology and Research for offering me a rare internship opportunity.

My special thanks to Yume who has accompanied me for the entire PhD journey. Last, Dad and Mom, the sacrifices you have made for me are beyond any description. Thank you, my beloveds. I wish I would always make you proud.

# Abstract

Data has been the fuel that drives modern artificial intelligence. With more and more emerging concerns on data, e.g., data privacy, learning paradigms demand evolving accordingly. In this thesis, I focus on discriminative learning on **restricted data** from the role of learning executors who are in charge of the learning process. That means, my research interest lies in how to design proper learning algorithms while data is considered restricted. Specifically, the following three types of scenarios will be explored.

- *Private data is accessible to learning executors, but the learning process should not expose any data information.* In such a scenario, the learning executors are trusted while any others are restricted from accessing data. Differential Privacy (DP) is a golden principle for this problem which preserves the participation of every data point and thus can defend against strong adversaries. I intend to take a step forward and explore how to ensure a safe learning when data is pairwise labelled, to which DP cannot be directly applied due to the explicit pairwise correlations.

- *Only incomplete data is accessible to learning executors.* Learning on incomplete data is a challenging topic when some data information is restricted from learning executors. For example, not all people would like to answer their demographics in a survey. Suppose a common case where the missing values are from a discrete attribute or the label domain. Inferring them comes to the Semi-Supervised Learning (SSL) problem. I will study how to improve the prediction module for unlabeled data from the aspect of prediction uncertainty.

- *None of data is accessible to learning executors, but some feedbacks are available.* This scenario considers that learning executors cannot access data, which hinders feeding data to the model for the end-to-end back-propagation. However, learning is stilled feasible if some feedbacks from data are provided, e.g., model evaluations. In practice, model tuning tasks are studied in this context, and the tuning efficiency for deep neural networks is particularly investigated.

From top to bottom, one can sense that the restriction on data becomes stricter, which also implies some bigger change should be applied to learning paradigms. Despite the specific requirement in each scenario, I am interested how to deal with them via a general principle. It is known that gradient-based optimization has been popular in machine learning. Given a fixed model structure, the eventually attained model can be attributed to the elaborately designed model gradients (Note that it is not always because of losses). With this insight, I propose to deal with different restricted data by using different meaningful **gradient manipulation** techniques. Concretely, I apply gradient perturbation to compensate for the missing or addition of any interested pairwise data, employ gradient masking to reduce the impact of over-confident unlabeled predictions, adopt gradient estimation to learn from model evaluations. I conclude that although the meaning of restricted data varies across different tasks (which also brings out various challenges), the insight of gradient manipulation constantly offers a good perspective to tackle these problems.

# Table of contents

# List of figures

# List of tables