

# Preface

知者行之始，行者知之成  
—王陽明

*“The best theory is inspired by practice.  
The best practice is inspired by theory.”*  
— Donald Knuth

Optimization is central to machine learning (ML), which in turn forms the foundation of artificial intelligence (AI). From training deep neural networks to fine-tuning Large Language Models, almost every advancement in AI relies on solving some form of optimization problem. While classical methods based on empirical risk minimization (ERM) have powered much of early progress in ML, they are no longer sufficient to address the growing complexity of today’s AI challenges. This book aims to bridge that gap by offering a systematic treatment of the emerging optimization paradigm known as **compositional optimization** and its applications in modern AI. Many critical optimization problems in ML now exhibit intricate compositional structures as  $f(g)$  or  $\sum_{i=1}^n f_i(g_i)$  that go beyond traditional frameworks, where both  $f$  and  $g$  are non-linear functions and potentially non-convex, extending beyond the scope of traditional optimization paradigms. However, most existing texts remain focused on classical stochastic optimization and ERM, overlooking the depth and diversity of these newer challenges.

## Motivation of writing the book

Optimization once held a central spotlight at leading ML venues such as NeurIPS and ICML. In recent years, however, the field has seen an influx of new topics in AI, capturing the interest of students and early-career researchers. While attention has increasingly shifted toward foundation models and AGI, the importance and impact of optimization remain as vital as ever.

As someone working at the intersection of optimization and machine learning, I feel a dual responsibility. **First**, to bring cutting-edge optimization techniques to the

---

broader ML/AI community. When I speak with researchers in ML/AI and mention my focus on optimization for machine learning, I am often met with questions like, “*What problems are you working on?*” or “*Are these theories truly useful, given that they rely on assumptions that may not be easily verified in practice?*” Some even remarked that optimization’s only practical contribution to AI is the Adam algorithm. This reflects a common misconception that optimization in ML is limited to training algorithms like SGD or Adam, which is far from the truth. **Second**, I feel a responsibility to encourage researchers in mathematical optimization to engage more deeply with the challenges of modern AI. Many researchers in traditional optimization are eager to contribute, but the rapid pace of AI along with the constant influx of new models and terminology can make it difficult to identify core problems where optimization insights are most needed. Working at this intersection gives me a unique perspective: recognizing fundamental challenges in modern AI, such as the training of large foundation models, and abstracting them into rigorous mathematical frameworks where optimization methods can offer meaningful solutions. I hope this book contributes to bridging the gap between the AI and optimization communities and inspires new collaborations across these fields.

At first glance, the focus on compositional optimization in this book may seem narrow, but it is deeply connected to fundamental learning and optimization principles including discriminative learning and robust optimization, and has broad applicability across ML and modern AI, which will be shown in this book. In particular, this book introduces a new family of risk functions termed X-risks, in which the loss function of each data involves comparison with many others. We formulate empirical X-risk minimization as finite-sum coupled compositional optimization (FCCO) - a new family of compositional optimization. After five years of intensive research on this subject, we have explored different aspects of FCCO, from upper bounds to lower bounds, from smooth objectives to non-smooth objectives, from convex problems to non-convex problems, and from theoretical complexity analysis to applications in training large foundation models. While significant progress has been made, many open questions remain. Nevertheless, we believe it is time to share this advanced body of knowledge with the broader community in the form of a comprehensive book.

## Structure of the book

This book is crafted to engage both theory-oriented and practice-driven audiences. It presents rigorous theoretical analysis with deep insights, complemented by practical implementation tips, Github code repositories, and empirical evidence—effectively bridging the gap between theory and application. It is intended for graduate students, applied researchers, and anyone interested in the intersection of optimization and machine learning. The readers are assumed to have some basic knowledge in ML. The materials in this book have been used in my graduate-level course on stochastic optimization for ML.

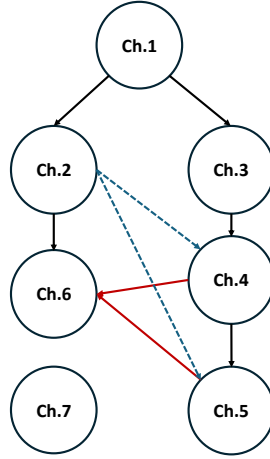


Fig. 0.1: Structure of the Book Chapters. Dashed lines indicate motivation. The red solid lines indicate application. Other solid lines indicate dependency.

The book is organized as follows. Chapter 1 reviews the fundamentals of convex optimization essential for the material presented in this book. Chapter 2 introduces advanced learning methods that go beyond traditional ERM framework so as to motivate compositional optimization. Chapter 3 presents classical stochastic optimization algorithms and their complexity analysis in both convex and non-convex settings. Chapter 4 delves into stochastic compositional optimization (SCO) problems with algorithms and complexity analysis. Chapter 5 explores algorithms and analysis for solving FCCO problems. Chapter 6 presents applications of SCO and FCCO in supervised and self-supervised learning for training predictive models, generative models, and representation models. Chapter 5 and 6 are largely devoted to the original research conducted by the author and his team. The dependencies and flow among the chapters are illustrated in Figure 0.1. Practitioners may focus on Chapter 2 and Chapter 6. For theory-oriented audiences who are interested in ML applications, I strongly recommend reading Chapter 2 and Chapter 6 as well.

### Acknowledgments

This book would not have been possible without the dedication and contributions of my students. I would like to thank my former and current Ph.D. students, visiting students and postdoc who contributed to both the theoretical and empirical results presented in the book. In particular, I acknowledge significant theoretical contributions from Bokun Wang, Quanqi Hu, Zhishuai Guo, Wei Jiang, Yan Yan, Qi Qi, Ming Yang, Xingyu Chen, Yao Yao, Yi Xu, Mingrui Liu and Linli Zhou, and significant empirical contributions from Zhuoning Yuan, Gang Li, Xiyuan Wei, Dixian Zhu, Siqi Guo, Zihao Qiu, Vicente Balmaseda and Anant Mehta. Special thanks go to

---

Bokun Wang for his help on preparation of the lecture notes for my course with the initial version of proofs of many methods covered in this book. I thank Ning Ning for proofreading some chapters.

I am grateful to my academic collaborators Qihang Lin, Yiming Ying, Lijun Zhang, Tuo Zhao, Yunwen Lei, Shuiwang Ji, Nitesh Chawla, Zhaosong Lu, Jiebo Luo, Xiaodong Wu, My T. Thai, Milan Sonka, Zhengzhong Tu, Tomer Galanti, Yin-bing Liang, Hongchang Gao, Bang An, Ilgee Hong, Guanghui Wang, Limei Wang, Youzhi Luo, Haiyang Yu, and Zhao Xu, and industrial collaborators Rong Jin, Wotao Yin, Denny Zhou, Wei Liu, Xiaoyu Wang, Ming Lin, Liangliang Cao, Xuanhui Wang, Yuexin Wu, and Xianzhi Du. I would like to thank my long-term collaborator Qihang Lin. We have worked together on the application of FCCO to constrained optimization featured in the book. Special thanks to Guanghui Lan, Chih-Jen Lin and Stephen Wright for their insightful discussions on subjects covered in this work. They have inspired me to solve some of the fundamental optimization problems covered in this book. I am especially thankful to My T. Thai for encouraging me to publish this book.

I owe a great debt of gratitude to my PhD advisor, Dr. Rong Jin, who introduced me to the world of optimization and taught me the value of focus.

I am deeply grateful to my department head, Scott Schaefer, as well as to all my colleagues in the Department of Computer Science and Engineering, for fostering such a positive and collaborative atmosphere. I also like to thank my former colleagues at the University of Iowa.

Finally, I am grateful for support from the National Science Foundation for my research under my career award #1844403, the RI core grant #2246756, and the FAI grant #2246757.

College Station, TX, USA,  
January, 2026

Tianbao Yang

# Contents

<b>1</b>	<b>Basics: Convex Optimization</b>	1
1.1	Notations and Definitions	3
1.2	Verification of Convexity	6
1.3	Fenchel Conjugate	8
1.4	Convex Optimization	9
1.4.1	Local Minima and Global Minima	10
1.4.2	Optimality Conditions	10
1.4.3	Karush–Kuhn–Tucker (KKT) Conditions	11
1.5	Basic Lemmas	16
1.6	History and Notes	21
<b>2</b>	<b>Introduction: Advanced Machine Learning</b>	23
2.1	Empirical Risk Minimization	25
2.1.1	Discriminative Label Prediction	25
2.1.2	Discriminative Loss Functions	26
2.1.3	Need of Optimization Algorithms	29
2.1.4	Generalization Analysis	30
2.2	Robust Optimization	31
2.2.1	Distributionally Robust Optimization	31
2.2.2	Optimized Certainty Equivalent	35
2.2.3	Group Distributionally Robust Optimization	38
2.3	Empirical X-risk Minimization	39
2.3.1	AUC Losses	40
2.3.2	Average Precision Loss	44
2.3.3	Partial AUC Losses	46
2.3.4	Ranking Losses	50
2.3.5	Contrastive Losses	52
2.4	Discriminative Data Prediction	53
2.4.1	A Discriminative Probabilistic Modeling Approach	54
2.4.2	A Robust Optimization Approach	59
2.5	History and Notes	62

---

<b>3</b>	<b>Classic: Stochastic Optimization</b>	67
3.1	Stochastic Gradient Descent	69
3.1.1	Smooth Convex Functions	70
3.1.2	Non-smooth Convex Functions	73
3.1.3	Smooth Non-Convex Functions	75
3.1.4	Non-smooth Weakly Convex Functions	77
3.2	Stochastic Proximal Gradient Descent	82
3.2.1	Convex Functions	84
3.2.2	Strongly Convex Functions	86
3.3	Stochastic Coordinate Descent	91
3.4	Stochastic Mirror Descent	96
3.4.1	Non-smooth Composite Problems	99
3.4.2	Non-smooth Problems	101
3.5	Adaptive Gradient Method (AdaGrad)	102
3.6	Stochastic Gradient Descent Ascent	107
3.7	Stochastic Optimistic Mirror Prox	112
3.8	History and Notes	118
<b>4</b>	<b>Foundations: Stochastic Compositional Optimization</b>	123
4.1	Stochastic Compositional Optimization	125
4.2	Stochastic Compositional Gradient Descent	126
4.2.1	Convergence Analysis	127
4.2.2	An Improved Complexity with Smooth Inner Function	131
4.2.3	A Straightforward Approach with a Large Mini-batch	137
4.3	Stochastic Compositional Momentum Methods	138
4.3.1	Moving-Average Gradient Estimator	138
4.3.2	STORM Estimators	147
4.4	Non-smooth (Non-convex) Regularized Problems	154
4.5	Structured Optimization with Compositional Gradient	160
4.5.1	Non-convex Min-Max Optimization	161
4.5.2	Non-convex Min-Min Optimization	166
4.5.3	Non-convex Bilevel Optimization	171
4.6	History and Notes	183
<b>5</b>	<b>Advances: Finite-sum Coupled Compositional Optimization</b>	187
5.1	Finite-sum Coupled Compositional Optimization	189
5.2	Smooth Functions	190
5.2.1	The SOX Algorithm	191
5.2.2	Multi-block Single-Probe Variance Reduction	199
5.3	Non-Smooth Weakly Convex Functions	208
5.3.1	SONX for Non-smooth Inner Functions	210
5.3.2	SONEX for Non-smooth Outer functions	217
5.4	Convex inner and outer functions	222
5.4.1	The ALEXR Algorithm	224
5.4.2	Technical Lemmas	226

5.4.3	Strongly convex objectives	237
5.4.4	Convex objectives with non-smooth outer functions	242
5.4.5	Double-loop ALEXR for weakly convex inner functions	247
5.4.6	Lower Bounds	249
5.5	Stochastic Optimization of Compositional OCE	255
5.5.1	A Basic Algorithm	256
5.5.2	A Geometry-aware Algorithm for Entropic Risk	264
5.6	History and Notes	294
<b>6</b>	<b>Applications: Learning Predictive, Generative and Representation Models</b>	<b>299</b>
6.1	Stochastic Optimization Framework	301
6.1.1	Milestones of Stochastic Optimization	303
6.1.2	Limitations of Existing Optimization Framework	306
6.2	DRO and Group DRO	307
6.2.1	DRO for Imbalanced Classification	307
6.2.2	GDRO for Addressing Spurious Correlation	313
6.3	Extreme Multi-class Classification	315
6.4	Stochastic AUC and NDCG Maximization	318
6.4.1	Stochastic AUC Maximization	319
6.4.2	Stochastic AP Maximization	323
6.4.3	Stochastic Partial AUC Maximization	325
6.4.4	Stochastic NDCG Maximization	331
6.4.5	The LibAUC Library	334
6.5	Discriminative Pretraining of Representation Models	338
6.5.1	Mini-batch Contrastive Losses	338
6.5.2	Contrastive Learning without Large Mini-Batches	341
6.5.3	Contrastive Learning with Learnable Temperatures	344
6.6	Discriminative Fine-tuning of Large Language Models	350
6.6.1	Pipeline of LLM Training	350
6.6.2	DFT for fine-tuning Large Language Models	356
6.6.3	DisCO for Reinforcing Large Reasoning Models	361
6.7	Constrained Learning	367
6.7.1	A General Penalty-based Approach via FCCO	368
6.7.2	Continual Learning with Zero-forgetting Constraints	375
6.7.3	Constrained Learning with Fairness Constraints	379
6.8	Learning Data Compositional Networks	381
6.8.1	Large-scale Graph Neural Networks	381
6.8.2	Multi-instance Learning with Attention	384
6.9	DRRHO Risk Minimization	387
6.10	History and Notes	391
<b>7</b>	<b>Afterword</b>	<b>395</b>
	References	397