

# Scalable Convex Optimization Methods for Semidefinite Programming

Thèse N° 9598

Présentée le 27 août 2019

à la Faculté des sciences et techniques de l'ingénieur  
Laboratoire de systèmes d'information et d'inférence  
Programme doctoral en informatique et communications

pour l'obtention du grade de Docteur ès Sciences

par

**Alp YURTSEVER**

Acceptée sur proposition du jury

Prof. O. N. A. Svensson, président du jury

Prof. V. Cevher, directeur de thèse

Prof. J. A. Tropp, rapporteur

Prof. S. Sra, rapporteur

Prof. M. Jaggi, rapporteur

2019



Devranın deęiřmez dzeni vardır.  
zen, zemeyen yenilir ona.  
Onu okumanın bir dili vardır.  
Bilimler bilimi denilir ona.

— Hseyin Yurtsever (2012)  
*Sayılar Bilimi - Bilinmezlik Deryası*

To my family,  
dedicated to the memory of my father, Hseyin Yurtsever.







# Acknowledgements

Before we go into the technical content, I would like to add a few words to express my appreciation to the people who supported me during my studies as a PhD candidate.

A very special gratitude goes to my PhD advisor Volkan Cevher. I am grateful to him not only for our research discussions and his technical supervision; but also for being an advisor with genuine concern for my welfare, for sharing the excitement of my accomplishments, and for being a heartfelt friend.

I spent three months at MIT as a visiting student in the fall semester of 2018. I am thankful to Suvrit Sra for hosting me during this visit, and to Volkan for arranging this research visit and making it possible in the first place. I really enjoyed my time at MIT and I am looking forward to joining there as a postdoc. I would also like to thank Joel Tropp for his guidance on our collaborated research for more than 3 years, for all our research discussions, and for inviting and hosting me at Caltech multiple times during this collaboration.

I was fortunate to have worked with so many great collaborators during my PhD. I am deeply grateful to all my collaborators. In particular, this dissertation is based on some joint works with Volkan Cevher, Olivier Fercoq, Ya-Ping Hsieh, Francesco Locatello, Suvrit Sra, Quoc Tran-Dinh, Joel Tropp, and Madeleine Udell.

I was honored to have Volkan Cevher, Martin Jaggi, Suvrit Sra, Ola Svensson, and Joel Tropp as my thesis defense committee members. I am grateful for their time and valuable comments. I would also like to thank my colleagues Ahmet Alacaoglu, Ali Kavis, Fatih Sahin, and Maria Vladarean for reading my thesis and providing feedback, and to Paul Rolland for proofreading my French in the abstract.

I am thankful to our lab members and the alumni for the collaborations, research discussions, and the friendly working environment. I am indebted to our administrative assistant Gosia Baltaian for her help on many occasions and for her friendliness.

I would also like to take this opportunity to thank my friends. My PhD years coincided with some unfortunate incidents in my personal life. Without the love and support of my friends, I would never have recovered. My life is more meaningful with the memorable times that we spent together. Finally, I acknowledge the effort of my teachers over the years, which played a critical role in shaping my academic identity and interests.

Last but certainly not least, I would like to thank my parents Çiğdem and Hüseyin, my sister Bengi, and my brother Emre for their love and constant support. I dedicate this dissertation to the memory of my father. Without his encouragement, I would never have started this journey.

*Lausanne, 10 July 2019*

A. Y.





# Abstract

With the ever-growing data sizes along with the increasing complexity of the modern problem formulations, contemporary applications in science and engineering impose heavy computational and storage burdens on the optimization algorithms. As a result, there is a recent trend where heuristic approaches with unverifiable assumptions are overtaking more rigorous, conventional methods at the expense of robustness and reproducibility.

My recent research results show that this trend can be overturned when we jointly exploit dimensionality reduction and adaptivity in optimization at its core. I contend that even the classical convex optimization did not reach yet its limits of scalability.

Many applications in signal processing and machine learning cast a fitting problem from limited data, introducing spatial priors to be able to solve these otherwise ill-posed problems. Data is small, the solution is compact, but the search space is high in dimensions. These problems clearly suffer from the wasteful use of storage resources by the classical optimization algorithms.

This problem is prevalent. Storage is a critical bottleneck that prevents us from solving many important large-scale optimization problems. For example, semidefinite programs (SDP) often have low-rank solutions. However, SDP algorithms require us to store a matrix decision variable in the ambient dimensions.

This dissertation presents a convex optimization paradigm which makes it possible to solve trillion dimensional SDP relaxations to combinatorial decision problems on a regular personal computer. The key idea is to use an optimization procedure that performs compact updates, and to maintain only a small sketch of the decision variable. Once the iterations terminate, we can recover an approximate solution from the information stored in this sketch.

We start by formalizing this idea for a model problem with low-rank solutions. Based on the classical conditional gradient method, we propose the first convex optimization algorithm with optimal storage guarantees for this model problem. Then, we develop and describe the key ingredients to extend our results for a broader set of problems, SDP in particular.

SDP formulations are key in signal processing, machine learning, and other engineering applications. By greatly enhancing the scalability of optimization algorithms for solving SDP formulations, we can potentially unlock new results in important scientific applications.

Key words: Convex optimization, semidefinite programming, low-rank matrix optimization, primal-dual methods, conditional gradient methods, low-rank matrix sketching





## Résumé

Avec la taille toujours croissante des données massives et la complexité des formulations de problèmes modernes, les applications actuelles en sciences et en ingénierie imposent de lourdes charges de calcul et de stockage. En conséquence, il existe une tendance récente où les approches heuristiques avec des hypothèses invérifiables dépassent les méthodes plus rigoureuses et conventionnelles au détriment de la robustesse et de la reproductibilité.

Mes résultats récents montrent que cette tendance peut être renversée si nous exploitons conjointement la réduction de la dimensionnalité et l'adaptivité dans l'optimisation. Je soutiens que même l'optimisation convexe classique n'a pas encore atteint ses limites.

De nombreuses applications dans le traitement du signal et l'apprentissage automatique posent un problème d'ajustage à partir de données limitées, en introduisant des informations préalables spatiales afin de pouvoir résoudre ces problèmes. Les données et la solution sont compactes, mais l'espace de recherche est de grande dimension. Clairement, ces problèmes souffrent de l'utilisation inutile des ressources de stockage par les algorithmes d'optimisation. Ce problème est très répandue. Le stockage est un goulot d'étranglement critique qui nous empêche de résoudre de nombreux problèmes d'optimisation à grande échelle. Par exemple, les programmes semi-définies (SDP) ont souvent des solutions de rang faible. Cependant, les algorithmes SDP nous obligent à stocker une variable de décision matricielle dans les dimensions ambiantes.

Cette thèse présente un paradigme d'optimisation convexe qui permet de résoudre des instances de programmation semi-définie à grande échelle de trillions de dimensions sur un ordinateur personnel ordinaire. L'idée clé est d'utiliser un algorithme d'optimisation qui effectue des mises à jour compactes, et de ne conserver qu'un petit *sketch* de la variable de décision. Une fois la procédure terminée, nous pouvons récupérer une solution approximative à partir des informations stockées dans le sketch.

Nous commençons par la formalisation de cette idée pour un problème de modèle avec des solutions à faible rang. Sur la base de la méthode classique du gradient conditionnel, nous proposons le premier algorithme d'optimisation convexe avec des garanties de stockage optimales pour ce problème de modèle. Ensuite, nous développons et décrivons les ingrédients clés pour étendre nos résultats à d'autres modèles de problèmes, SDP en particulier.

Les formulations SDP jouent un rôle clé dans le traitement du signal, l'apprentissage automatique, et d'autres applications d'ingénierie. En améliorant l'évolutivité des méthodes d'optimisation pour la résolution des formulations SDP, nous pouvons potentiellement obtenir de nouveaux résultats dans des applications scientifiques importantes.

## Résumé

---

Mots clés : Optimisation convexe, programmation semi-définie, optimisation matricielle de rang faible, méthodes primales-duales, méthodes de gradient conditionnel, sketch pour des matrices de rang faibles

# Bibliographic Note

This dissertation is based on the following publications:

- [YTDC15a] Alp Yurtsever, Quoc Tran-Dinh, Volkan Cevher. "A Universal Primal-Dual Convex Optimization Framework." Conference on Neural Information Processing Systems (NeurIPS), 2015.
- [YHC15] Alp Yurtsever, Ya-Ping Hsieh, Volkan Cevher. "Scalable Convex Methods for Phase Retrieval." IEEE International Workshop on Computational Advances in Multi-Sensor Adaptive Processing (CAMSAP), 2015.
- [YUTC17] Alp Yurtsever, Madeleine Udell, Joel Aaron Tropp, Volkan Cevher. "Sketchy Decisions: Convex Low-rank Matrix Optimization with Optimal Storage." International Conference on Artificial Intelligence and Statistics (AISTATS), 2017.
- [TYUC17b] Joel Aaron Tropp, Alp Yurtsever, Madeleine Udell, Volkan Cevher. "Practical Sketching Algorithms for Low-rank Matrix Approximation." SIAM Journal on Matrix Analysis and Applications (SIMAX), 2017.
- [TYUC17a] Joel Aaron Tropp, Alp Yurtsever, Madeleine Udell, Volkan Cevher. "Fixed-rank Approximation of a Positive-Semidefinite Matrix from Streaming Data." Conference on Neural Information Processing Systems (NeurIPS), 2017.
- [YFLC18] Alp Yurtsever, Olivier Fercoq, Francesco Locatello, Volkan Cevher. "A Conditional Gradient Framework for Composite Convex Minimization with Applications to Semidefinite Programming." International Conference on Machine Learning (ICML), 2018.
- [YFC19] Alp Yurtsever, Olivier Fercoq, Volkan Cevher. "A Conditional Gradient-Based Augmented Lagrangian Framework." International Conference on Machine Learning (ICML), 2019.
- [YSC19] Alp Yurtsever, Suvrit Sra, Volkan Cevher. "Conditional Gradient Methods via Stochastic Path-Integrated Differential Estimator." International Conference on Machine Learning (ICML), 2019.
- [TYUC19] Joel Aaron Tropp, Alp Yurtsever, Madeleine Udell, Volkan Cevher. "Streaming Low-Rank Matrix Approximation with an Application to Scientific Simulation." Accepted to SIAM Journal on Scientific Computing (SISC), 2019.

## Bibliographic Note

---

My other publications relevant to this research, but not covered in this dissertation, are:

- [YVC16] Alp Yurtsever, Bang Cong Vu, Volkan Cevher. "Stochastic Three-Composite Convex Minimization." Conference on Neural Information Processing Systems (NeurIPS), 2016.
- [OLY<sup>+</sup>16] Gergely Odor, Yen-Huan Li, Alp Yurtsever, Ya-Ping Hsieh, Quoc Tran-Dinh, Marwa El Halabi, Volkan Cevher. "Frank-Wolfe works for non-Lipschitz Continuous Gradient Objectives: Scalable Poisson Phase Retrieval." IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), 2016.
- [CVY18] Volkan Cevher, Bang Cong Vu, Alp Yurtsever. "Stochastic Forward Douglas-Rachford Splitting Method for Monotone Inclusions." in Large-Scale and Distributed Optimization, Chapter 7. Springer, Cham, (2018)
- [LYC18] Kfir Levy, Alp Yurtsever, Volkan Cevher. "Online Adaptive Methods, Universality and Acceleration." Conference on Neural Information Processing Systems (NeurIPS), 2018.
- [DYC<sup>+</sup>19] Lijun Ding, Alp Yurtsever, Volkan Cevher, Joel Aaron Tropp, Madeleine Udell. "An Optimal-Storage Approach to Semidefinite Programming Using Approximate Complementarity." Manuscript (arXiv:1902.03373), 2019.
- [CVY19] Volkan Cevher, Bang Cong Vu, Alp Yurtsever. "Inertial Three-Operator Splitting Method and Applications." Manuscript (arXiv:1904.12980), 2019.
- [LYFC19] Francesco Locatello, Alp Yurtsever, Olivier Fercoq, Volkan Cevher. "Stochastic Conditional Gradient Method for Composite Convex Minimization." Manuscript (arXiv:1901.10348), 2019.



# Contents

<b>Acknowledgements</b>	<b>v</b>
<b>Abstract (English/Français)</b>	<b>vii</b>
<b>Bibliographic Note</b>	<b>xi</b>
<b>List of figures</b>	<b>xvi</b>
<b>List of tables</b>	<b>xvii</b>
<b>1 Introduction</b>	<b>1</b>
1.1 Problem Formulation . . . . .	2
1.2 Semidefinite Programming . . . . .	3
1.2.1 Relax and Round . . . . .	4
1.3 Arithmetics: Crisis and Opportunity . . . . .	5
1.4 Storage: Crisis and Opportunity . . . . .	6
1.5 State of the Research in the Field . . . . .	7
1.6 Organization and Contributions . . . . .	9
1.7 Notation and Terminology . . . . .	10
<b>2 Storage-Efficient Convex Optimization</b>	<b>13</b>
2.1 Conditional Gradient . . . . .	15
2.1.1 Opportunity . . . . .	15
2.1.2 CGM Iteration . . . . .	16
2.2 ThinCGM . . . . .	17
2.2.1 Thin SVD Updates . . . . .	17
2.2.2 ThinCGM Iteration . . . . .	18
2.2.3 A Numerical Study on the Rank Expansion . . . . .	20
2.3 SketchyCGM . . . . .	21
2.3.1 Randomized Sketch . . . . .	21
2.3.2 SketchyCGM Iteration . . . . .	22
2.3.3 Convergence Results for SketchyCGM . . . . .	24
2.4 Extensions for Semidefinite Programming . . . . .	25
2.5 Numerical Experiments . . . . .	26
2.5.1 Matrix completion . . . . .	26

## Contents

---

2.5.2	Phase Retrieval Problems . . . . .	27
<b>3</b>	<b>Universal Primal-Dual Methods</b>	<b>33</b>
3.1	Preliminaries . . . . .	35
3.2	Universal Primal-Dual Method . . . . .	35
3.2.1	Efficiency Considerations and Hölder Smoothness . . . . .	36
3.2.2	UPD Iteration . . . . .	37
3.2.3	Guarantees . . . . .	38
3.2.4	Application to the SDP template . . . . .	39
3.3	Accelerated Universal Primal-Dual Method . . . . .	41
3.3.1	AUPD Iteration . . . . .	41
3.3.2	Guarantees . . . . .	42
3.3.3	Application to the SDP template . . . . .	42
3.4	Numerical Experiments . . . . .	43
3.4.1	Quantum tomography with Pauli operators . . . . .	43
3.4.2	Matrix completion with MovieLens dataset . . . . .	44
3.5	CGM is Dual Averaging Subgradient Method . . . . .	46
3.6	Appendix: Proofs . . . . .	48
<b>4</b>	<b>CGM for Composite Problems</b>	<b>59</b>
4.1	Preliminaries . . . . .	61
4.1.1	Nesterov Smoothing . . . . .	61
4.1.2	Quadratic Penalty . . . . .	62
4.2	Homotopy CGM . . . . .	63
4.2.1	HCGM Iteration . . . . .	63
4.2.2	Stopping Criterion . . . . .	64
4.2.3	Guarantees . . . . .	64
4.2.4	Convergence with Inexact Oracles . . . . .	65
4.2.5	Application to the SDP template . . . . .	67
4.3	Applications & Related Work . . . . .	68
4.3.1	Smooth Problems . . . . .	68
4.3.2	Regularized Problems . . . . .	68
4.3.3	Non-Smooth Problems . . . . .	69
4.3.4	Minimax Problems . . . . .	69
4.3.5	Problems with Affine Constraints . . . . .	70
4.3.6	Minimization via Splitting . . . . .	71
4.4	Numerical Experiments . . . . .	72
4.4.1	Clustering the MNIST dataset . . . . .	72
4.4.2	Robust PCA . . . . .	73
4.5	Appendix: Proofs . . . . .	75

<b>5</b>	<b>CGM via Augmented Lagrangian</b>	<b>81</b>
5.1	Augmented Lagrangian Penalty . . . . .	82
5.2	Conditional Gradient Augmented Lagrangian Method . . . . .	83
5.2.1	CGAL Iteration . . . . .	83
5.2.2	Guarantees . . . . .	84
5.2.3	Application to the SDP template . . . . .	85
5.3	Numerical Experiments . . . . .	86
5.3.1	Max-cut . . . . .	87
5.3.2	k-means Clustering . . . . .	87
5.3.3	Generalized Eigenvector Problem . . . . .	89
5.4	Appendix: Proofs . . . . .	93
<b>6</b>	<b>Low-Rank Matrix Sketching from Streaming Data</b>	<b>99</b>
6.1	Two Component Sketch . . . . .	101
6.1.1	The Sketch . . . . .	101
6.1.2	The Basic Reconstruction Algorithm . . . . .	102
6.1.3	The Fixed-Rank Reconstruction Algorithm . . . . .	103
6.1.4	Computing a PSD Approximation . . . . .	103
6.2	Three Component Sketch . . . . .	104
6.2.1	The Sketch . . . . .	104
6.2.2	Computing Truncated Low-Rank Approximations . . . . .	105
6.2.3	Analysis of Initial and Truncated Approximations . . . . .	105
6.3	Nyström Sketch for PSD Matrices . . . . .	106
6.3.1	The Sketch . . . . .	106
6.3.2	Fixed-Rank Nyström Approximation . . . . .	107
6.4	Numerical Experiments . . . . .	108
6.4.1	Sketching and Reconstruction Methods in our Experiments . . . . .	108
6.4.2	Experimental Setup . . . . .	109
6.4.3	Classes of Input Matrices . . . . .	109
6.4.4	Comparison of Formulas for General Low-Rank Matrices . . . . .	111
6.4.5	Comparison of Formulas for Low-Rank PSD Matrices . . . . .	111
<b>7</b>	<b>CGM with Stochastic Path-Integrated Differential Estimator</b>	<b>117</b>
7.1	Related Works . . . . .	118
7.1.1	Conditional Gradient Sliding . . . . .	119
7.1.2	Stochastic Path-Integrated Differential Estimator . . . . .	119
7.2	Preliminaries . . . . .	119
7.3	SPIDER Frank-Wolfe . . . . .	121
7.3.1	Convex Finite-Sum . . . . .	122
7.3.2	Convex Expectation Minimization . . . . .	123
7.3.3	Non-convex Finite-Sum . . . . .	123
7.3.4	Non-convex Expectation Minimization . . . . .	124
7.4	SPIDER Conditional Gradient Sliding . . . . .	125

## Contents

---

7.4.1	Convex Finite-Sum . . . . .	125
7.4.2	Convex Expectation Minimization . . . . .	127
7.4.3	Non-convex Finite-Sum . . . . .	127
7.4.4	Non-convex Expectation Minimization . . . . .	128
7.5	Comparison & Discussion . . . . .	129
7.5.1	Convex Optimization Camp . . . . .	129
7.5.2	Non-convex Optimization Camp . . . . .	130
7.5.3	Results from Concurrent Research . . . . .	131
<b>8</b>	<b>Conclusion &amp; Future Directions</b>	<b>151</b>
	<b>Bibliography</b>	<b>166</b>
	<b>Curriculum Vitae</b>	<b>167</b>

# List of Figures

1.1	Hierarchies in convex optimization. . . . .	3
1.2	Graph cut on a toy graph. . . . .	4
1.3	Black-box model for solving an SDP relaxation. . . . .	6
2.1	Rank expansion of CGM estimates. . . . .	20
2.2	Performance of SketchyCGM for solving matrix completion (MovieLens 100K). . . . .	26
2.3	Performance of SketchyCGM for solving matrix completion (MovieLens 10m). . . . .	27
2.4	Memory scaling for five convex optimization algorithms for phase retrieval. . . . .	28
2.5	Gaussian and Poisson phase retrieval under Poisson noise. . . . .	29
2.6	Evolution of the Reconstruction Quality of SketchyCGM for Fourier Ptychography. . . . .	30
2.7	Three algorithms for Fourier ptychographic imaging via phase retrieval. . . . .	31
3.1	Comparison of UPD and CGM for solving quantum tomography problem. . . . .	45
3.2	Comparison of UPD and CGM for solving matrix completion problem. . . . .	45
4.1	HCGM for clustering SDP with preprocessed MNIST dataset. . . . .	72
4.2	HCGM for image inpainting with robust PCA. . . . .	74
4.3	PSNR and SSIM vs iteration counter for formulations with $\ell_1$ and $\ell_2$ loss. . . . .	74
5.1	Empirical performance of various methods for solving max-cut problem. . . . .	88
5.2	Empirical comparison of UPD, HCGM and CGAL for max-cut problem. . . . .	88
5.3	Comparison of CGAL and HCGM for clustering SDP. . . . .	88
5.4	Comparison of CGAL, HCGM and UPD for generalized eigenvector SDP. . . . .	89
5.5	Comparison of projection-free methods for max-cut SDP [Part 1]. . . . .	91
5.6	Comparison of the projection-free methods for max-cut SDP [Part 2]. . . . .	92
6.1	Spectra of input matrices for numerical experiments. . . . .	110
6.2	Comparison of fixed-rank reconstruction formulas: Synthetic examples. . . . .	112
6.3	Comparison of fixed-rank reconstruction formulas: Real data examples. . . . .	113
6.4	Comparison of fixed-rank PSD reconstruction formulas: Synthetic examples. . . . .	114
6.5	Comparison of fixed-rank PSD reconstruction formulas: Real data examples. . . . .	115
6.6	Numerical stability in implementing Nyström approximation formulas. . . . .	116





## List of Tables

7.1	Comparison of conditional gradient methods for stochastic optimization. . . .	129
-----	---	-----





# 1 Introduction

We can recognize almost every problem in science and engineering as an optimization problem. Whether it is a finance problem where we want to maximize the profit based on a financial model, a problem in physics where the aim is to identify a minimum energy state under the constraints imposed by the laws of physics, or a machine learning problem for designing an autonomous self-driving vehicle, we can model it as a mathematical optimization problem.

Convex optimization is an important subclass of mathematical optimization. One major benefit of formulating an application as a convex optimization problem is its robustness. Convex optimization provides a unified modeling framework for solving problems across diverse disciplines. To find an approximate solution, convex optimization algorithms rely on favorable geometric structures and functional properties present in the convex optimization templates, avoiding further assumptions on the problem data as much as possible.

Unfortunately, with the ever-growing data sizes along with the increasing complexity of the modern applications, contemporary problems in science and engineering impose heavy computational and storage burdens on the optimization algorithms. As a result, there is a recent trend where heuristic approaches with unverifiable assumptions on the problem data are overtaking more rigorous, conventional methods at the expense of robustness and reproducibility.

My recent research results show that this trend can be overturned, when we jointly exploit dimensionality reduction and adaptivity in optimization at its core. I contend that even the classical convex optimization did not reach yet its limits of scalability. In particular, this dissertation presents a convex optimization paradigm which makes it possible to solve some trillion dimensional semidefinite programming problems on a regular laptop.

The design of an end-product algorithm with a more detailed cost analysis, along with a well-tuned, stable and robust implementation and numerical demonstrations in cutting-edge applications are deferred to a follow-up paper.

## 1.1 Problem Formulation

The importance of convex optimization in machine learning has increased dramatically in the last decade. Indeed, a large class of learning formulations can be addressed by the following constrained convex minimization template:

$$\underset{\mathbf{x}}{\text{minimize}} \quad f(\mathbf{x}) \quad \text{subject to} \quad \mathcal{A}\mathbf{x} \in \mathcal{K}, \quad \mathbf{x} \in \mathcal{X}. \quad (1.1)$$

where  $\mathcal{X} \subset \mathbb{R}^p$  is a convex and compact (nonempty, bounded, closed) set, and its 0-dimensional faces (*i.e.*, its vertices) are called *atoms*.  $f: \mathcal{X} \rightarrow \mathbb{R} \cup \{+\infty\}$  is a proper closed convex function,  $\mathcal{A}: \mathbb{R}^p \rightarrow \mathbb{R}^d$  is a linear map, and  $\mathcal{K} \subseteq \mathbb{R}^d$  is a convex set.

This template covers in particular semidefinite programming (SDP) formulations. Let us focus on the following SDP template:

$$\begin{aligned} &\underset{\mathbf{X}}{\text{minimize}} \quad \langle \mathbf{C}, \mathbf{X} \rangle \\ &\text{subject to} \quad \langle \mathbf{A}_i, \mathbf{X} \rangle = b_i \quad \text{for } i = 1, \dots, d \\ &\quad \mathbf{X} \in \mathbb{S}_+^n \quad \text{and} \quad \text{Tr } \mathbf{X} = \alpha \end{aligned} \quad (1.2)$$

Here,  $\mathbb{S}_+^n$  denotes the cone of  $n \times n$  positive semidefinite matrices.  $\mathbf{C}, \mathbf{A}_1, \mathbf{A}_2, \dots, \mathbf{A}_d \in \mathbb{S}^n$  (*i.e.*,  $n \times n$  dimensional and symmetric), and  $b_1, b_2, \dots, b_d \in \mathbb{R}$ . This is the standard SDP template except the  $\text{Tr } \mathbf{X} = \alpha$  constraint that we add to ensure the boundedness of the problem domain. This additional constraint is not restrictive, since we can compute (or approximate) the trace of a solution a priori, either from the problem formulation or from the data in many applications.

We will use the following standard compact notation for the affine constraints

$$\begin{aligned} \mathcal{A}: \mathbb{S}_+^n &\rightarrow \mathbb{R}^d \quad \text{where} \quad \mathcal{A}\mathbf{X} = \begin{bmatrix} \langle \mathbf{A}_1, \mathbf{X} \rangle & \dots & \langle \mathbf{A}_d, \mathbf{X} \rangle \end{bmatrix}; \\ \mathcal{A}^*: \mathbb{R}^d &\rightarrow \mathbb{S}_+^n \quad \text{where} \quad \mathcal{A}^* \mathbf{z} = \sum_{i=1}^d z_i \mathbf{A}_i. \end{aligned} \quad (1.3)$$

Then, we can reformulate the model problem (1.2) as

$$\underset{\mathbf{X}}{\text{minimize}} \quad \langle \mathbf{C}, \mathbf{X} \rangle \quad \text{subject to} \quad \mathcal{A}\mathbf{X} = \mathbf{b}, \quad \underbrace{\mathbf{X} \in \mathbb{S}_+^n, \quad \text{Tr } \mathbf{X} = \alpha}_{\mathbf{X} \in \mathcal{X}}, \quad (1.4)$$

where the vector  $\mathbf{b} := (b_1, \dots, b_d) \in \mathbb{R}^d$  is the constraint (measurement) vector. Clearly (1.4) is an instance of (1.1). We restrict ourselves to equality constraint  $\mathcal{A}\mathbf{X} = \mathbf{b}$  for notational simplicity. All results that we present in this dissertation can be extended for the affine inclusion constraint  $\mathcal{A}\mathbf{X} \in \mathcal{K}$ . Similarly, extension for  $\text{Tr } \mathbf{X} \leq \alpha$  constraint is straightforward.

In most applications, the cost matrix  $\mathbf{C}$  and the constraint matrices  $\mathbf{A}_i$  have low intrinsic dimensionality, *e.g.*, they are low rank, or sparse. Hence, storing and applying these matrices do not put a substantial cost on our computational budget.

## 1.2 Semidefinite Programming

There is a common misconception between convexity and tractability of an optimization problem. In full generality, convex optimization problems are not solvable in polynomial time, unless  $P = NP$ . There are examples of NP-hard combinatorial decision problems that can be formulated as a convex optimization problem over the copositive cone<sup>1</sup>. As a consequence, it is not possible to design a generic solver for the complete class of convex optimization problems, and we need to define (or identify) a problem subclass before designing (resp. applying) a convex optimization algorithm.

This identification process of an appropriate algorithmic template requires field expertise. In order to circumvent this need, disciplined convex programming introduces a set of conventions to follow when formulating optimization problems. This yields a natural hierarchy of programming templates in convex optimization. Semidefinite programming (SDP) sits somewhere on top of this hierarchy. It covers most of the other well-known templates, from linear programming to second-order cone programming, as a special case.

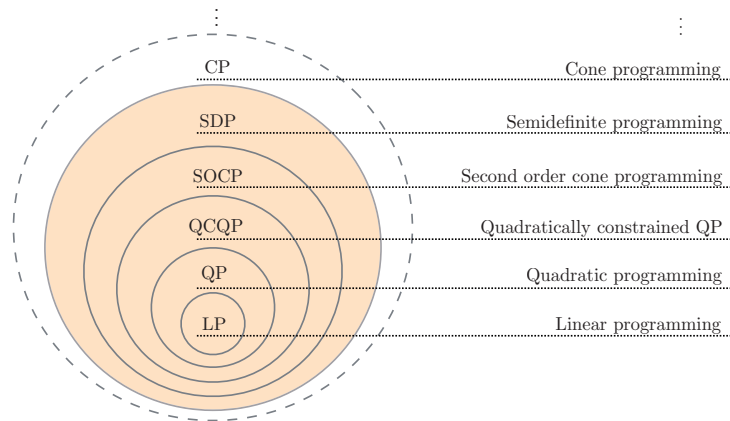


Figure 1.1 – Disciplined convex programming yield a natural hierarchy in convex optimization. SDP sits somewhere on top of this hierarchy.

Semidefinite programming relaxations are powerful for obtaining approximate solutions to countless problems in computer science [KN12, Lov03], machine learning [MNS15, SSGB07], and signal processing [Sin11, SS11, CKS15]. Many recent studies have revealed the deep and surprising phenomenon that, for many of these problems, one cannot obtain better approximations beyond the semidefinite relaxation. A famous example (see [Rag08]) in computer science is that, assuming a strengthened  $P \neq NP$  conjecture, *called the Unique Games Conjecture*, it is NP-hard to beat semidefinite relaxation for a large class of combinatorial optimization problems.

<sup>1</sup>The definition of copositivity can be traced back to [Mot52]. For a recent survey, see [HUS10].

### 1.2.1 Relax and Round

Let us derive a basic SDP relaxation for a combinatorial decision problem, maximum cut problem (max-cut), in order to present some common structures in these formulations. Similar ideas arise in many other SDP applications.

Given an undirected graph  $G = (V, E)$  with vertices  $\mathcal{V} = \{v_1, v_2, \dots, v_n\}$  and a set of weights  $c : E \rightarrow \mathbb{R}_+$  associated to each edge on the graph, the aim is to define cut that partitions the vertices into two distinct sets, such that the sum of the weights that correspond to the edges that we cross is maximized.

We can mathematically formulate this problem as a combinatorial optimization problem:

$$\underset{x \in \mathbb{R}^n}{\text{maximize}} \quad \frac{1}{2} \sum_{\{i,j\} \in E} c_{i,j}(1 - x_i x_j) \quad \text{subject to} \quad x_i \in \{-1, +1\}, \quad i = 1, 2, \dots, n. \quad (1.5)$$

This problem is NP-hard, and the difficulty arises from the quadratic terms in the objective value  $x_i x_j$ , and the integer valued constraint  $x_i \in \{-1, +1\}$ .

To deal with this difficulty, we can *lift* the problem into the matrix space: Define the matrix variable  $\mathbf{X} = \mathbf{x}\mathbf{x}^*$ . Then, we can equivalently formulate (1.5) as a rank constrained matrix optimization problem:

$$\underset{\mathbf{X}}{\text{maximize}} \quad \frac{1}{2} \sum_{\{i,j\} \in E} c_{i,j}(1 - X_{ij}) \quad \text{subject to} \quad \text{diag } \mathbf{X} = \mathbf{1} \quad \mathbf{X} \in \mathbb{S}_+^n, \quad \text{rank } \mathbf{X} = 1. \quad (1.6)$$

where the linear map  $\text{diag} : \mathbb{S}_+^n \rightarrow \mathbb{R}^n$  extracts the diagonal of a matrix. This problem is still NP-hard, but we translated the difficulty into the non-convex rank constraint.

Nest step is to *relax* by dropping the rank constraint:

$$\underset{\mathbf{X}}{\text{maximize}} \quad \frac{1}{2} \sum_{\{i,j\} \in E} c_{i,j}(1 - X_{ij}) \quad \text{subject to} \quad \text{diag } \mathbf{X} = \mathbf{1} \quad \mathbf{X} \in \mathbb{S}_+^n. \quad (1.7)$$

This is the SDP relaxation of the combinatorial max-cut problem [GW95]. We can add  $\text{Tr } \mathbf{X} = n$  constraint to (1.7), which is already implied by  $\text{diag } \mathbf{X} = \mathbf{1}$ .

We can solve (1.7), in polynomial time, and get an approximate solution  $\mathbf{X}_\star$  to (1.7). However,  $\mathbf{X}_\star$  is not rank-one in general. Hence, in order to provide an approximate solution to the original problem (1.5), we need a *rounding* step. Goemans and Williamson [GW95] describe a randomized rounding procedure which produces approximately optimal cuts with theoretical guarantees.

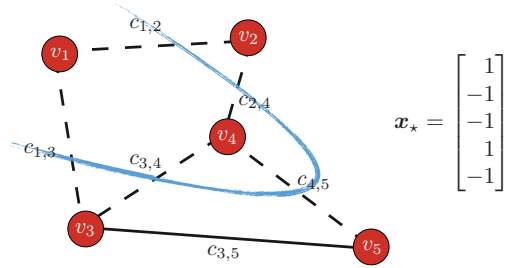


Figure 1.2 – Graph cut on a toy graph.

### 1.3 Arithmetics: Crisis and Opportunity

Within the last two decades, developments in machine learning applications have led to revolutionary changes in our expectations from optimization in many ways. The stupendous amount of data involved in the problem descriptions forced us to consider improving scalability of the optimization algorithms as the primary goal.

Given a fixed amount of computational resources, and a fixed target approximation accuracy, the scalability (to problem size) of an optimization procedure for a specific task can be theoretically bounded. In other words, in order to scale our algorithms to very large problems beyond some theoretical bound, we have to sacrifice from the approximation accuracy. It is for this reason that we typically consider the problem dimensions and the accuracy as variables when we derive iteration complexity of an optimization algorithm, and suppress all other data-dependent constants.

Fortunately, this compromise is theoretically well-justified in machine learning. In many machine learning applications, one minimizes a loss function that corresponds to a statistical estimator of the true model (to which we do not have access). This incurs modeling error to the overall system. As a consequence, solving the optimization problem to very high accuracy, orders beyond the estimation error, can only provide an insignificant amount of extra information about the overall system. Indeed, high accuracy solutions to these problems are not even desirable because it can cause overfitting. Accordingly, at first, first-order methods, and more recently their stochastic extensions became increasingly popular despite their slow sublinear convergence rates.

We can argue that this phenomenon is not specific to machine learning applications. Consider an SDP relaxation to a combinatorial decision problem, such as the max-cut example. We relax the original NP-hard problem as an SDP formulation by dropping the non-convex rank constraint. We approximately solve the SDP problem using an iterative optimization procedure. Then, we execute a rounding step to produce an approximate solution to the original problem. We can view this whole procedure, which consists of *relaxation*, *optimization*, and *rounding*, as a black-box model for solving the original combinatorial optimization problem, see Figure 1.3.

The goals of optimization and modeling are tightly connected. Relaxation imposes a modeling error to the overall system. Hence, solving the optimization step to very high accuracy, beyond the relaxation accuracy, do not have a significant effect on the final estimation quality. We demonstrate this phenomenon numerically for the SDP formulation of a clustering problem in Section 4.4.1. These arguments exclude a small class of examples for which the empirical study suggests that the relaxation is exact.



The goals of **optimization** and **modeling** are tightly connected:

$$\underbrace{\|X^t - X^\natural\|}_{\text{estimation accuracy}} \leq \underbrace{\|X^t - X^*\|}_{\text{optimization accuracy}} + \underbrace{\|X^* - X^\natural\|}_{\text{relaxation accuracy}}$$

$X^\natural$ : solution of the true model (lifted,  $X^\natural = x^\natural x^{\natural*}$ )  
 $X^*$ : solution obtained by SDP relaxation  
 $X^t$ : estimation obtained at iteration  $t$

Figure 1.3 – Black-box model for solving an SDP relaxation: Relax, Optimize and Round.

## 1.4 Storage: Crisis and Opportunity

The majority of the semidefinite relaxation formulations relies on *lifting* the decision variable to a higher dimensional search-space. Recall the max-cut example, the original formulation (1.5) corresponds to an  $n$ -dimensional combinatorial problem, whereas the SDP relaxation (1.7) lies in an  $(n \times n)$ -dimensional space.

Storage cost is the critical bottleneck that prevents us from solving many important problems in large-scale. This becomes evident even in a simple max-cut example, considering the size of modern social or transportation network graphs, which can easily reach to millions of nodes. Even just the decision variable requires storing a huge, millions by millions matrix!

Decision variable is a core internal variable for optimization. Any reasonable optimization algorithm must output an approximate solution to the problem. Consequently, we cannot completely discard the decision variable from our optimization procedure, and the only hope that remains is to exploit the structures of a solution for storing it in intrinsic dimensions.

Solutions to SDP problems are often low-rank, or close to low-rank. Consider an SDP relaxation to a combinatorial optimization problem, such as the max-cut example. Even though we drop the rank constraint, roughly speaking, a solution to the relaxation is informative only if it is close to the solution set of the original problem. Accordingly, the solution (to SDP) typically exhibits a fast decay in singular value spectrum, even if it is not exactly low-rank.

Moreover, many SDP problems have low-rank solutions due to the widely known Pataki-Barvinok bound, see [Bar95] and [Pat98]. Under some mild technical assumptions, SDP problem (1.4) has a solution with rank  $r \leq r^{PB}$  for some  $r^{PB} = \Theta(\sqrt{d})$ . In other words, SDP problems provably have low-rank solutions when  $d \ll n^2$ . Recall that  $d = n$  in max-cut.

Unfortunately the storage issue remains, even if the solution has low intrinsic dimensionality. Virtually all convex optimization methods require storing the  $n \times n$  dimensional matrix decision variable in the ambient dimensions, simply because the intermediate estimates might be incompressible. In fact, this fundamental burden has been seen as a natural limit for scalability, and the researchers switched their focus to non-convex frameworks to attain scalable solutions at the expense of robustness. In 2015, while proposing the non-convex Wirtinger-Flow method for solving phase retrieval problems, Candès et al. [CLS15b] explain the reason of using a non-convex template as follows:

*“For certain random models, some recent SDP relaxations such as PhaseLift [CESV13] are known to provide exact solutions (up to global phase) to the generalized phase retrieval problem using a near minimal number of sampling vectors [CLS15a, CSV12]. While in principle SDP based relaxations offer tractable solutions, they become computationally prohibitive as the dimension of the signal increases. Indeed, for a large number of unknowns in the tens of thousands, say, the memory requirements are far out of reach of desktop computers so that these SDP relaxations are de facto impractical.”*

Our recent research results show that this trend can be overturned. In Chapter 2, we describe a storage-optimal convex optimization paradigm for solving large-scale matrix optimization problems with low-rank solutions. The key idea is to maintain only a sketch of the matrix decision variable, and to use an algorithm that is compatible with these sketches. To demonstrate the scalability of our framework, we exhibit numerical solutions to a phase-retrieval imaging problem with  $25'600$  pixel images containing human blood cells from a working Fourier ptychography system [HCO<sup>+</sup>15]. Our approach can scale to problems of millions of dimensions on a regular personal computer.

## 1.5 State of the Research in the Field

The first reliable algorithms for solving semidefinite programs were based on interior-point methods, introduced independently by Nesterov & Nemirovski [NN89, NN94] and Alizadeh [Ali91, Ali93]. Unfortunately, interior point methods are not effective for solving large-scale problems. This is simply because they typically rely on a second order optimization algorithm as an oracle to query at each iteration. In other words, each iteration of an interior point method itself is a difficult optimization problem that does not scale well.

The success in solving SDP problems led researchers to investigate new applications. In return, invention of new applications put the existing optimization frameworks under more and more pressure to accommodate larger-scale problems. As a result of the increased popularity of SDP formulations in machine learning and signal processing applications, researchers started to search for more scalable procedures for solving SDP problems. Accordingly, first-order methods are proposed as an alternative to interior point methods for large-scale applications. The majority of these algorithms are direct adaptations of some classical methods in the

## Chapter 1. Introduction

---

literature for an SDP template, such as the proximal gradient approach and its accelerated variants [Roc70, BT09] or the alternating direction method of multipliers [BPC<sup>+</sup>11].

A substantial challenge for the first-order methods is the computational cost of the projection oracle. Recall that projection onto positive semidefinite cone might require a full eigendecomposition, which imposes  $\mathcal{O}(n^3)$  arithmetic cost per-iteration. Coupled with the slow, sublinear convergence rate of first-order methods, this creates an undesirable computational burden. The conditional gradient method (CGM, also known as the Frank-Wolfe algorithm) is extremely powerful in this setup, since it avoids projection steps by leveraging the so-called linear minimization oracle that can be solved efficiently using iterative linear algebra routines such as the power method or Lanczos algorithm [FW56, Haz08, Jag13].

Hazan [Haz08] proposed using CGM for solving smooth problems over positive semidefinite cone. However, CGM is not as flexible as the proximal first-order methods, and no extant practical variant of CGM exists in the literature for solving the SDP template (1.4). To this end, [Haz08, GH11, GH16] suggest solving a sequence of feasibility problems with the least squares loss and applying binary search to the objective value. Renegar [Ren14] develops a subgradient method for linear and semidefinite programming with similar per-iteration cost as CGM. None of these methods, however, addresses the storage bottleneck.

Despite the advancements in reducing the arithmetic costs, our progress in reducing the storage footprint is minimal in the convex realm. Almost all convex optimization methods require storing the  $n \times n$  dimensional matrix decision variable, which eliminates them in large-scale.

This promoted a new trend of non-convex algorithms that apply to non-convex reformulations of the SDP template. The majority of these heuristics are based on the factorization idea of Burer & Monteiro [BM03], where the matrix variable is factorized as  $\mathbf{X} = \mathbf{U}\mathbf{U}^*$ , such that  $\mathbf{U} \in \mathbb{R}^{n \times r}$  for some splitting rank  $r \leq n$ :

$$\text{minimize } \langle \mathbf{C}, \mathbf{U}\mathbf{U}^* \rangle \quad \text{subject to } \mathcal{A}(\mathbf{U}\mathbf{U}^*) = \mathbf{b}. \quad (1.8)$$

One can apply a wide range of optimization methods to find a stationary point of this problem with respect to the factor  $\mathbf{U}$ . See [BM03, JNS13, Bou15, BKS16, CLS15b] for some examples. The main advantage of the non-convex heuristics is the reduced size of the decision variable. Some methods also benefit from locally linear convergence rates. On the other hand, these algorithms are provable only under stringent, unverifiable and (sometimes) unrealistic statistical assumptions on the problem data.

Very recently, we have developed a new storage-optimal algorithm for solving SDP in standard form, based on the approximate complementarity principle [DYC<sup>+</sup>19]. The main idea is to solve the low-dimensional dual SDP approximately, and then to recover the primal SDP by solving a compressed SDP problem to the eigenspace with the small eigenvalues of the dual slack matrix. We do not cover this work in this dissertation.



## 1.6 Organization and Contributions

This dissertation presents a novel convex optimization paradigm to design scalable algorithms for solving the general SDP template (1.4). In particular, the contributions of each chapter are as follows:

- In Chapter 2, we propose a novel paradigm for designing optimization algorithms for large-scale convex problems with structured solutions. It specifically focuses on a fundamental class of convex matrix optimization problems with low-rank solutions. The two key ideas are (i) to maintain only a sketch of the decision variable and (ii) to use algorithmic templates that are compatible with these sketches. It develops and analyzes a new sketch-based algorithm for a model problem from this class. This algorithm, SketchyCGM, modifies a standard convex optimization scheme, the conditional gradient method. In contrast to nonconvex heuristics, the guarantees for SketchyCGM do not rely on statistical models for the problem data. Numerical evidence establishes the benefits of SketchyCGM.

In order to adapt SketchyCGM for solving the SDP template (1.4), we need to design compatible algorithms for (1.4). To this end, Chapters 2 and 3 develop new optimization methods for (1.1), with similar characteristics as the conditional gradient method.

- In Chapter 3, we describe a new primal-dual optimization method for the prototypical constrained convex optimization template (1.1). In contrast to existing primal-dual algorithms, our framework avoids the proximity operator of the objective function altogether. We instead leverage computationally cheaper, Fenchel-type operators, which are the main workhorses of the generalized conditional gradient type methods. Unlike standard CGM, our algorithm can also handle the affine constraints  $\mathcal{A}\mathbf{x} \in \mathcal{K}$ . Our algorithms are *universal* in the sense that they can automatically adapt to the unknown Hölder smoothness orders within the template.
- In Chapter 4, we introduce a generalized conditional gradient method, HomotopyCGM, for the constrained convex optimization template (1.1). Our approach combines smoothing (quadratic penalty) and homotopy techniques under the CGM framework, and provably achieves  $\mathcal{O}(1/\sqrt{t})$  convergence rate, where  $t$  denotes the iteration counter. We demonstrate that the same rate holds if the linear subproblems are solved approximately with additive or multiplicative error.
- In Chapter 5, We extend our HomotopyCGM algorithm from quadratic penalty to an augmented Lagrangian framework. The resulting algorithm, CGAL, retains the strong theoretical guarantees of HomotopyCGM while exhibiting significantly superior empirical performance. Numerical evidence demonstrates the benefits of CGAL in various SDP applications.

SketchyCGM adopts a sketching model as a natural mechanism for reducing the dimensionality of the decision variable. Consequently, designing practical sketching algorithms for low-rank matrix approximation is a fundamental aspect of our research.

## Chapter 1. Introduction

---

- In Chapter 6, we develop a suite of new algorithms for fixed-rank approximation from a sketch. These methods can preserve structural properties of the input matrix, such as positive-semidefiniteness. Each method is accompanied by an informative error bound. The algorithms are simple, accurate, numerically stable, and provably correct. We implement our algorithms as a MATLAB Toolbox open for public access. Computer experiments show that the proposed methods dominate alternative techniques across a wide range of examples.

Many SDP algorithms require the problem data to be readily available. However, most problems in engineering and machine learning contains some degree of randomness. Therefore, an interesting future research direction is the stochastic extension of our findings.

- In Chapter 7, we present our preliminary study for the stochastic extension, where we aim to identify and present the stochastic CGM variants with the best theoretical guarantees. We study expectation minimization and finite-sum settings, both with and without the convexity assumption. We propose a class of novel variance-reduced stochastic CGM algorithms, and improve the best known rates in the literature for various problem settings.
- Finally in Chapter 8, we summarize the main contributions of this dissertation and describe some future research directions.

### 1.7 Notation and Terminology

For notational simplicity, we work on the  $\mathbb{R}^p / \mathbb{R}^{n \times n}$  and  $\mathbb{R}^d$  spaces with the Euclidean norms. We write  $\|\cdot\|$  for the Euclidean norm,  $\|\cdot\|_F$  for the Frobenius norm, and  $\|\cdot\|_{S_1}$  for the Schatten 1-norm (aka the *trace norm* or the *nuclear norm*). Depending on context,  $\langle \cdot, \cdot \rangle$  refers to the Euclidean or Frobenius inner product. The symbol  $*$  denotes the conjugate transpose of a vector or matrix, as well as the adjoint of a linear map. The dagger  $^\dagger$  refers to the pseudoinverse. The symbol  $[X]_r$  stands for a best rank- $r$  Frobenius-norm approximation of the matrix  $X$ . The symbol  $\succ$  denotes the semidefinite order. The function  $\text{dist}(\mathbf{z}; \mathcal{K})$  returns the minimum Euclidean distance from  $\mathbf{z}$  to a set  $\mathcal{K}$ , i.e.,  $\text{dist}(\mathbf{z}; \mathcal{K}) := \min_{\mathbf{u} \in \mathcal{K}} \|\mathbf{z} - \mathbf{u}\|$ . Similarly, the function  $\text{dist}(\mathbf{X}; \mathcal{X})$  returns the minimum Frobenius-norm distance from  $\mathbf{X}$  to a set  $\mathcal{X}$ . For a convex function  $f$ , we use  $\nabla f$  both for its subgradient and gradient. We denote the diameter of  $\mathcal{X}$  by  $D_{\mathcal{X}} = \max_{\mathbf{x}_1, \mathbf{x}_2 \in \mathcal{X}} \|\mathbf{x}_1 - \mathbf{x}_2\|$ . We use the computer science interpretation of the order notation  $\mathcal{O}, \tilde{\mathcal{O}}, \Omega, \Theta$ .

Next, we recall some basic notions from convex analysis that will be used especially in Chapters 3 and 4 to design new optimization methods for solving (1.1).

**Subdifferential.** The subdifferential of a function  $f$  at a point  $\mathbf{x} \in \mathbb{R}^p$  is defined as

$$\partial f(\mathbf{x}) = \left\{ \mathbf{u} \in \mathbb{R}^p : f(\mathbf{v}) - f(\mathbf{x}) \geq \langle \mathbf{v} - \mathbf{x}, \mathbf{u} \rangle, \forall \mathbf{v} \in \mathbb{R}^p \right\}. \quad (1.9)$$

Elements of the subdifferential set are called subgradients, and we denote an arbitrary subgradient of  $f$  at  $\mathbf{x}$  by  $\nabla f(\mathbf{x})$ . If  $\partial f(\mathbf{x})$  is a singleton, then  $f$  is a differentiable function, and we call  $\nabla f(\mathbf{x})$  as the gradient of  $f$  at  $\mathbf{x}$ .

**Lipschitz continuity.** We say that a function  $g : \mathbb{R}^d \rightarrow \mathbb{R}$  is  $L$ -Lipschitz continuous if it satisfies

$$|g(\mathbf{x}_1) - g(\mathbf{x}_2)| \leq L \|\mathbf{x}_1 - \mathbf{x}_2\|, \quad \forall \mathbf{x}_1, \mathbf{x}_2 \in \mathbb{R}^d.$$

**Smoothness.** A differentiable function  $f : \mathcal{X} \rightarrow \mathbb{R}$  is  $L_f$ -smooth if its gradient  $\nabla f$  is  $L_f$ -Lipschitz continuous:

$$\|\nabla f(\mathbf{x}_1) - \nabla f(\mathbf{x}_2)\| \leq L_f \|\mathbf{x}_1 - \mathbf{x}_2\|, \quad \forall \mathbf{x}_1, \mathbf{x}_2 \in \mathcal{X}.$$

**Proximal operator.** The proximal operator of a function  $g : \mathbb{R}^d \rightarrow \mathbb{R}$  is defined as follows

$$\text{prox}_g(\mathbf{z}) = \arg \min_{\mathbf{y} \in \mathbb{R}^d} \left\{ g(\mathbf{y}) + \frac{1}{2} \|\mathbf{y} - \mathbf{z}\|^2 \right\}. \quad (1.10)$$

Roughly speaking, the proximal operator is tractable when the computation of (1.10) is cheap. If  $g$  is the indicator function of a nonempty, closed convex subset  $\mathcal{K}$ , its proximity operator is the projection operator on  $\mathcal{K}$ .

**Linear minimization oracle.** As the problem dimensions become increasingly larger, the proximal tractability assumption can be restrictive. This fact increased the popularity of the conditional gradient methods (CGM, *aka* Frank-Wolfe algorithms), which instead leverage the following linear minimization oracle:

$$\text{lmo}_{\mathcal{X}}(\mathbf{x}) = \arg \min_{\mathbf{u} \in \mathcal{X}} \langle \mathbf{u}, \mathbf{x} \rangle. \quad (1.11)$$

**Lagrange saddle point.** Introduce a slack variable  $\mathbf{v} \in \mathcal{K}$ , and define the Lagrangian of (a reformulation of) the problem (1.1) by

$$\mathcal{L}(\mathbf{x}, \mathbf{v}, \mathbf{y}) := f(\mathbf{x}) + \langle \mathbf{y}, \mathcal{A}\mathbf{x} - \mathbf{v} \rangle. \quad (1.12)$$

For the SDP template (1.4), we define the Lagrangian as  $\mathcal{L}(\mathbf{X}, \mathbf{y}) = f(\mathbf{X}) + \langle \mathbf{y}, \mathcal{A}\mathbf{X} - \mathbf{b} \rangle$ .

Then, we can formulate the primal and dual problems as follows:

$$\underbrace{\max_{\mathbf{y} \in \mathbb{R}^d} \min_{\substack{\mathbf{x} \in \mathcal{X} \\ \mathbf{v} \in \mathcal{K}}} \mathcal{L}(\mathbf{x}, \mathbf{y})}_{\text{dual}} \leq \underbrace{\min_{\substack{\mathbf{x} \in \mathcal{X} \\ \mathbf{v} \in \mathcal{K}}} \max_{\mathbf{y} \in \mathbb{R}^d} \mathcal{L}(\mathbf{x}, \mathbf{y})}_{\text{primal}}. \quad (1.13)$$

We assume that the strong duality holds, *i.e.*, the relation above holds with equality.

## Chapter 1. Introduction

---

Slater's condition is a sufficient condition for strong duality. By Slater's condition, we mean

$$\text{relint}(\mathcal{X} \times \mathcal{K}) \cap \{(\mathbf{x}, \mathbf{v}) \in \mathbb{R}^p \times \mathbb{R}^d : \mathcal{A}\mathbf{x} = \mathbf{v}\} \neq \emptyset, \quad (1.14)$$

where  $\text{relint}$  stands for the relative interior.

**Solution set.** We denote an exact solution of (1.1) by  $\mathbf{x}_\star$  (*resp.*, solution of (1.4) by  $\mathbf{X}_\star$ ), and the set of all solutions by  $\mathcal{X}_\star$ .

Similarly, we denote a solution of the dual problem by  $\mathbf{y}_\star$ , and the solution set by  $\mathcal{Y}_\star$ .

Throughout this dissertation, we assume that the solution sets  $\mathcal{X}_\star$  and  $\mathcal{Y}_\star$  are non-empty.

**$\varepsilon$ -solution.** Our goal is to approximately solve (1.1) (*resp.* (1.4)) to obtain  $\mathbf{x}_\varepsilon$  in the following sense: Given an accuracy level  $\varepsilon > 0$ , a point  $\mathbf{x}_\varepsilon \in \mathcal{X}$  is said to be an  $\varepsilon$ -solution of (1.1) if

$$f(\mathbf{x}_\varepsilon) - f_\star \leq \varepsilon, \text{ and } \text{dist}(\mathcal{A}\mathbf{x}_\varepsilon, \mathcal{K}) \leq \varepsilon. \quad (1.15)$$

Here, we call  $|f(\mathbf{x}_\varepsilon) - f_\star|$  the primal objective residual and  $\text{dist}(\mathcal{A}\mathbf{x}_\varepsilon, \mathcal{K})$  the feasibility gap. We use the same  $\varepsilon$  for the objective residual and the feasibility gap for notational simplicity. Distinct choices can be handled simply by scaling  $f$ .

Remark that  $f(\mathbf{x}_\varepsilon) - f_\star$  can take negative values, but the following inequality characterizes the maximum super-optimality in terms of the feasibility gap:

$$f(\mathbf{x}_\varepsilon) - f_\star \geq -\|\mathbf{y}_\star\| \cdot \text{dist}(\mathcal{A}\mathbf{x}_\varepsilon, \mathcal{K}) \quad \text{for any dual solution } \mathbf{y}_\star. \quad (1.16)$$

This is an algorithm-independent bound.

*Proof.* From the Lagrange saddle point theory, we know that the following bound holds  $\forall \mathbf{x} \in \mathcal{X}$  and  $\forall \mathbf{v} \in \mathcal{K}$ :

$$f_\star \leq \mathcal{L}(\mathbf{x}, \mathbf{v}, \mathbf{y}_\star) = f(\mathbf{x}) + \langle \mathbf{y}_\star, \mathcal{A}\mathbf{x} - \mathbf{v} \rangle \leq f(\mathbf{x}) + \|\mathbf{y}_\star\| \|\mathcal{A}\mathbf{x} - \mathbf{v}\| \quad (1.17)$$

where the last inequality holds due to the Cauchy-Schwarz inequality. By setting  $\mathbf{x} = \mathbf{x}_\varepsilon \in \mathcal{X}$  in this inequality and optimizing with respect to  $\mathbf{v} \in \mathcal{K}$ , we get

$$f(\mathbf{x}_\varepsilon) - f_\star \geq -\min_{\mathbf{v} \in \mathcal{K}} \|\mathbf{y}_\star\| \|\mathcal{A}\mathbf{x}_\varepsilon - \mathbf{v}\| = -\|\mathbf{y}_\star\| \text{dist}(\mathcal{A}\mathbf{x}_\varepsilon, \mathcal{K}). \quad (1.18)$$

□

## 2 Storage-Efficient Convex Optimization

This chapter introduces a storage optimal convex optimization paradigm. As a model example, we consider a fundamental class of convex matrix optimization problems with low-rank solutions. We argue that the main obstacle that prevents us from solving these problems at scale is the storage cost. As a proof of concept, we exhibit the first provably correct algorithm, SketchyCGM, for these problems, with optimal storage.

This chapter is based on two distinct joint works, with Ya-Ping Hsieh and Volkan Cevher [YHC15], and with Madeleine Udell, Joel A. Tropp and Volkan Cevher [YUTC17].

### Introduction

We consider a fundamental class of matrix optimization problems:

$$\underset{X \in \mathbb{R}^{m \times n}}{\text{minimize}} \quad f(\mathcal{A}X) \quad \text{subject to} \quad \|X\|_{S_1} \leq \alpha. \quad (2.1)$$

- ▷  $f$  is a smooth convex loss function;
- ▷  $\mathcal{A}: \mathbb{R}^{m \times n} \rightarrow \mathbb{R}^d$  is a given linear map;
- ▷  $\alpha \in \mathbb{R}_+$  is a given model parameter.

Nuclear norm constraint is a powerful proxy that promotes low-rank solutions [Faz02]. When  $\alpha$  is well-tuned, we expect solutions to problem (2.1) to be (approximately) low-rank. Validating template (2.1) for low-rank matrix optimization problems is beyond the scope of our work. Our focus is to overcome computational challenges for solving (2.1) in large-scale.

Problem (2.1) is smooth minimization within a nuclear norm-ball of radius  $\alpha$ , which is considerably easier to solve than the SDP template (1.4), especially in terms of the arithmetic cost. Nevertheless, this template suffers from the same storage burden, which is evident especially in the data-limited setting where  $d \ll mn$ . This is a natural assumption since  $d$  typically reflects the amount of data, which is typically smaller than the ambient dimensions of the problem. An example of this phenomenon is matrix completion, where the aim is to recover the complete matrix using only a small sample set of its entries.

## Chapter 2. Storage-Efficient Convex Optimization

---

Let us formalize the challenges that we face in storing the matrix variable.

- We need only  $d$  units of memory space to store measurements  $\mathcal{A}X$ , or the problem data. We also assume that  $\mathcal{A}$  has an intrinsic structure that makes it easy to store and access. If we think about the matrix completion example, storing indices that are associated with the given sample set is enough to implement  $\mathcal{A}$ .
- Suppose that a solution to problem (2.1) is (approximately) low-rank, denoting the rank by  $r \ll \min(m, n)$ . Then, we can store this solution (resp. a rank- $r$  approximation of this solution) using only  $r(m + n)$  units of memory space in factorized form.
- Almost all classical algorithms in the convex optimization literature for solving (2.1) requires initializing and storing a matrix decision variable in the ambient dimensions, which requires a storage unit of size  $\Theta(mn)$ , much larger than the storage we need for storing the data and a solution!

Our goal is to find a solution to problem (2.1) with rigorous approximation guarantees, but without ever initializing a matrix decision variable in the ambient dimensions.

### Contributions

In this chapter, we design two variants of conditional gradient method (CGM) with reduced memory footprint. We can summarize our contributions as follows:

- We introduce ThinCGM. ThinCGM operates much like CGM, but it stores a singular value decomposition (SVD) of the matrix variable. While ThinCGM is not theoretically guaranteed to alleviate the expansion of the memory usage, it often keeps a low memory footprint in comparison with the classical CGM. The procedure is exactly the same as CGM, except the way that matrix variable is stored in the memory. Therefore, ThinCGM should be viewed as a careful implementation of CGM rather than a new algorithm.
- We initiate the discussion about exploring the limits of storage efficiency for convex optimization algorithms. We formalize the concept of storage-optimality, and we introduce a new paradigm to design storage-optimal convex optimization methods.
- We propose SketchyCGM. SketchyCGM is the first storage-optimal convex optimization method for problem (2.1). SketchyCGM maintains a small randomized linear image (a *sketch*) of the matrix variable. After the optimization procedure is terminated, this sketch is used to extract a provably good low-rank approximation of the solution.
- We present computational evidence that SketchyCGM can solve large-scale instances of (2.1) that are not accessible to other convex optimization algorithms in the literature.

Most importantly, our results demonstrate that storage alone is not a reason to drop convexity.

## 2.1 Conditional Gradient

Conditional gradient method (CGM) plays a crucial role in designing convex optimization algorithms for solving large-scale low-rank matrix optimization problems.

CGM is proposed for the first time in the seminal work of Frank and Wolfe [FW56] for solving smooth convex optimization problem on a polytope. Following the machine learning revolution in optimization, this classical method underwent a significant increase in popularity, especially for matrix factorization problems. It is extended for the smooth convex minimization over the simplex by Clarkson [Cla10], for the spectrahedron by Hazan [Haz08], and finally for an arbitrary compact convex set by Jaggi [Jag13].

**Intuition.** In essence, at each iteration, CGM solves a linear program formed by the best linear underestimate of the smooth loss function at the current iterate  $X_t$ . Then, a solution to the original problem is obtained as a convex combination of the solutions to these linear programs.

One can also view CGM as an instance of the dual averaging subgradient method of Nesterov [Nes09]. See Section 3.5 for the details on this relation.

### 2.1.1 Opportunity

CGM does not need a projection oracle. This is the characteristic feature that makes CGM an irreplaceable algorithm for sparse signal recovery. The reason is simple: The solution to a linear program over a compact set, is provably an extreme point of this set. Since CGM constructs the solution as a convex combination of these extreme points, it is automatically feasible, leaving no reason to perform a projection step.

Why is this a big deal? Because the efficiency of implementing the projection step becomes questionable in large-scale problems. For the low-rank matrix optimization problem, for instance, projection corresponds to a soft-thresholding on the singular values. Hence, it requires computing a full SVD first. On the other hand, a linear minimization over the same nuclear norm-ball requires only the top singular vectors, which is arguably an easier task. In addition, CGM is robust against noise in the linear minimization step. Consequently, we can implement it with an approximate singular vector, which can be obtained efficiently using iterative linear algebra routines such as power method.

An auxiliary benefit is the interpretability of the primal averaging sequence of CGM. CGM forms the decision variable as a convex combination of the extreme points of the problem domain. The extreme points, naturally, have low intrinsic dimensionality. For example, in our model problem (2.1), extreme points are rank-one. Therefore, CGM constructs a solution to (2.1) as the sum of a streaming rank-one matrices. This model of streaming atomic updates is key to tackle the storage burden.

### 2.1.2 CGM Iteration

We describe the application of CGM for solving (2.1). Choose a feasible initial estimate:

$$\mathbf{X}_0 \in \mathbb{R}^{m \times n} \quad \text{where} \quad \|\mathbf{X}_0\|_{S_1} \leq \alpha. \quad (2.2)$$

As described, at each iteration  $t = 0, 1, 2, \dots$ , CGM minimizes a linear approximation:

$$\begin{aligned} \mathbf{H}_t &= \arg \min_{\mathbf{H}} \left\{ f(\mathcal{A}\mathbf{X}_t) + \langle \mathbf{H} - \mathbf{X}_t, \mathcal{A}^*(\nabla f(\mathcal{A}\mathbf{X}_t)) \rangle : \|\mathbf{H}\|_{S_1} \leq \alpha \right\} \\ &= \arg \min_{\mathbf{H}} \left\{ \langle \mathbf{H}, \mathcal{A}^*(\nabla f(\mathcal{A}\mathbf{X}_t)) \rangle : \|\mathbf{H}\|_{S_1} \leq \alpha \right\}. \end{aligned} \quad (2.3)$$

We can implement (2.3) using the following operations:

$$\mathbf{H}_t = -\alpha \mathbf{u}_t \mathbf{v}_t^* \quad \text{where} \quad (\mathbf{u}_t, \mathbf{v}_t) = \text{MaxSingVec}(\mathcal{A}^*(\nabla f(\mathcal{A}\mathbf{X}_t))). \quad (2.4)$$

MaxSingVec returns the left and right singular vectors that correspond to the maximum singular value. Note that we can efficiently implement MaxSingVec using power-method or Lanczos algorithm.

Finally, we update the decision variable:

$$\mathbf{X}_{t+1} = (1 - \eta_t) \mathbf{X}_t + \eta_t \mathbf{H}_t \quad \text{where} \quad \eta_t = 2/(t+2). \quad (2.5)$$

#### Stopping Rule

CGM comes with a natural and easily computable duality gap certificate that can be used to implement a simple stopping criterion.

From convexity of  $f$ , we know

$$\begin{aligned} f^* &\geq f(\mathcal{A}\mathbf{X}_t) + \langle \mathbf{X}_* - \mathbf{X}_t, \mathcal{A}^*(\nabla f(\mathcal{A}\mathbf{X}_t)) \rangle \\ &\geq \min_{\mathbf{H}} \left\{ f(\mathcal{A}\mathbf{X}_t) + \langle \mathbf{H} - \mathbf{X}_t, \mathcal{A}^*(\nabla f(\mathcal{A}\mathbf{X}_t)) \rangle : \|\mathbf{H}\|_{S_1} \leq \alpha \right\}. \end{aligned} \quad (2.6)$$

Hence, once we compute  $\mathbf{H}_t$  in (2.4), we can use it to compute the following surrogate:

$$\langle \mathcal{A}\mathbf{X}_t - \mathcal{A}\mathbf{H}_t, \nabla f(\mathcal{A}\mathbf{X}_t) \rangle \geq f(\mathcal{A}\mathbf{X}_t) - f^*. \quad (2.7)$$

As a consequence, given a target suboptimality parameter  $\varepsilon > 0$ , we can terminate CGM when

$$\langle \mathcal{A}\mathbf{X}_t - \mathcal{A}\mathbf{H}_t, \nabla f(\mathcal{A}\mathbf{X}_t) \rangle \leq \varepsilon. \quad (2.8)$$

This ensures that  $f(\mathcal{A}\mathbf{X}_t) - f^* \leq \varepsilon$ .



## 2.2 ThinCGM

The goal is to solve a matrix optimization problem. We are not allowed to form and store a matrix explicitly in ambient dimensions. We assume that the solution is low-rank so that it can be stored efficiently. What is the most straightforward action? Trying to store the matrix variable in a factored form, for example as an SVD.

This idea, however, should be taken with care. None of the algorithmic steps should require the explicit form of the matrix variable internally. In particular, we need an algorithmic procedure to compute SVD of the new estimate after an additive rank-one update (2.5).

### 2.2.1 Thin SVD Updates

In this section, we overview an algorithmic procedure that computes SVD of  $\mathbf{X}_{t+1}$  from SVD of  $\mathbf{X}_t$ , without generating  $\mathbf{X}_t$  in the ambient dimensions. This procedure is a direct application of Brand's ThinSVD updates in [Bra06].

Let  $\mathbf{X}$  have rank  $\hat{r}$ , and  $(\mathbf{U}, \mathbf{\Sigma}, \mathbf{V}) \in \mathbb{R}^{m \times \hat{r}} \times \mathbb{R}^{\hat{r} \times \hat{r}} \times \mathbb{R}^{n \times \hat{r}}$  be its SVD, so that

$$\mathbf{X} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^*. \quad (2.9)$$

We want to compute SVD of the next estimate defined by the update rule (2.5). Let us define the internal variables

$$\mathbf{m} = \mathbf{U}^* \mathbf{u}, \quad \mathbf{p} = \mathbf{u} - \mathbf{U}\mathbf{m}, \quad \hat{\mathbf{p}} = \frac{\mathbf{p}}{\|\mathbf{p}\|}, \quad (2.10)$$

$$\mathbf{n} = \mathbf{V}^* \mathbf{v}, \quad \mathbf{q} = \mathbf{v} - \mathbf{V}\mathbf{n}, \quad \hat{\mathbf{q}} = \frac{\mathbf{q}}{\|\mathbf{q}\|}. \quad (2.11)$$

Then, we compute  $(\hat{\mathbf{U}}, \hat{\mathbf{\Sigma}}, \hat{\mathbf{V}})$  by solving a small  $(\hat{r} + 1) \times (\hat{r} + 1)$  dimensional SVD problem

$$\hat{\mathbf{U}}\hat{\mathbf{\Sigma}}\hat{\mathbf{V}}^* = (1 - \eta) \begin{bmatrix} \mathbf{\Sigma} & \mathbf{0} \\ \mathbf{0} & 0 \end{bmatrix} - \eta \alpha \begin{bmatrix} \mathbf{m} \\ \|\mathbf{p}\| \end{bmatrix} \begin{bmatrix} \mathbf{n} \\ \|\mathbf{q}\| \end{bmatrix}^* \quad (2.12)$$

Now we can perform the CGM update (2.5) as follows:

$$\mathbf{U} \leftarrow \begin{bmatrix} \mathbf{U}\hat{\mathbf{U}} & \hat{\mathbf{p}}\hat{\mathbf{U}} \end{bmatrix} \quad \text{and} \quad \mathbf{V} \leftarrow \begin{bmatrix} \mathbf{V}\hat{\mathbf{V}} & \hat{\mathbf{q}}\hat{\mathbf{V}} \end{bmatrix} \quad \text{and} \quad \mathbf{\Sigma} \leftarrow \hat{\mathbf{\Sigma}}, \quad (2.13)$$

without ever generating  $\mathbf{X}_t$  in the ambient space.

**Remark.** This scheme can also be generalized for additive modifications of any fixed rank. See [Bra06] for more details.

### 2.2.2 ThinCGM Iteration

Choose a low-rank initial estimate directly in the form of an SVD. We consider  $\mathbf{X}_0 = \mathbf{0}$ . Then, perform the initialization:

$$\mathbf{U}_0 = \mathbf{0} \in \mathbb{R}^m, \quad \mathbf{V}_0 = \mathbf{0} \in \mathbb{R}^n, \quad \Sigma_0 = \mathbf{0} \in \mathbb{R}, \quad \text{and} \quad \mathbf{y}_0 = \mathcal{A}\mathbf{U}\Sigma\mathbf{V}^* = \mathbf{0} \in \mathbb{R}^d. \quad (2.14)$$

Optionally, fix an  $\varepsilon > 0$  for the stopping condition, in order to prescribe a target accuracy.

At each iteration  $t = 0, 1, 2, \dots$ , evaluate the linear minimization following the recipe

$$\mathbf{h}_t = \mathcal{A}(-\alpha \mathbf{u}_t \mathbf{v}_t^*) \quad \text{where} \quad (\mathbf{u}_t, \mathbf{v}_t) = \text{MaxSingVec}(\mathcal{A}^*(\nabla f(\mathbf{y}_t))). \quad (2.15)$$

Set the learning rate  $\eta_t = 2/(t+2)$ . Update the "dual iterate"  $\mathbf{y}_t$ , and also the factors of the primal variable using the procedure described in Section 2.2.1:

$$\begin{aligned} \mathbf{y}_{t+1} &= (1 - \eta_t) \mathbf{y}_t + \eta_t \mathbf{h}_t; \\ (\mathbf{U}_{t+1}, \Sigma_{t+1}, \mathbf{V}_{t+1}) &= \text{SVDUPDATE}(\mathbf{U}_t, \Sigma_t, \mathbf{V}_t, \mathbf{u}_t, \alpha, \mathbf{v}_t, \eta_t). \end{aligned} \quad (2.16)$$

If a stopping criteria is defined, continue until the condition is satisfied:

$$\langle \mathbf{y}_t - \mathbf{h}_t, \nabla f(\mathbf{y}_t) \rangle \leq \varepsilon. \quad (2.17)$$

See Algorithm 2.1 for a complete pseudocode.

An important detail here is that we can evaluate  $\text{MaxSingVec}$  in (2.15) without forming  $\mathcal{A}^*(\nabla f(\mathbf{y}_t))$  in ambient dimensions! To emphasize how, let us consider the compact notation of  $\mathcal{A}$ , in the same spirit as (1.3):

$$\begin{aligned} \mathcal{A} : \mathbb{R}^{m \times n} &\rightarrow \mathbb{R}^d \quad \text{where} \quad \mathcal{A}\mathbf{X} = \begin{bmatrix} \langle \mathbf{A}_1, \mathbf{X} \rangle & \dots & \langle \mathbf{A}_d, \mathbf{X} \rangle \end{bmatrix}; \\ \mathcal{A}^* : \mathbb{R}^d &\rightarrow \mathbb{R}^{m \times n} \quad \text{where} \quad \mathcal{A}^* \mathbf{z} = \sum_{i=1}^d z_i \mathbf{A}_i. \end{aligned} \quad (2.18)$$

By assumption,  $\mathcal{A}$  has low intrinsic complexity, which translates into (2.18) as each coefficient matrix  $\mathbf{A}_i \in \mathbb{R}^{m \times n}$  has only a few degrees of freedom. For the matrix completion problem, for instance, each  $\mathbf{A}_i$  is one-sparse.

To evaluate  $\text{MaxSingVec}$ , we just need to run an iterative procedure such as power method, Lanczos algorithm or randomized SVD [HMT11]. These procedures require only an oracle to compute matrix-vector multiplications involving the matrix  $\mathcal{A}^*(\nabla f(\mathbf{y}))$ :

$$\begin{aligned} \mathcal{A}^*(\nabla f(\mathbf{y})) : \mathbb{R}^n &\rightarrow \mathbb{R}^m \quad \text{where} \quad \mathcal{A}^*(\nabla f(\mathbf{y})) \mathbf{u} = \sum_{i=1}^d z_i \mathbf{A}_i \mathbf{u} \\ (\mathcal{A}^*(\nabla f(\mathbf{y})))^* : \mathbb{R}^m &\rightarrow \mathbb{R}^n \quad \text{where} \quad (\mathcal{A}^*(\nabla f(\mathbf{y})))^* \mathbf{v} = \sum_{i=1}^d z_i \mathbf{A}_i^* \mathbf{v} \end{aligned} \quad (2.19)$$

Notice that these oracles work in  $\mathbb{R}^m$  and  $\mathbb{R}^n$  dimensions.

**Algorithm 2.1** ThinCGM for model problem (2.1)**Require:** Data for (2.1); (optional) suboptimality  $\varepsilon$ 


---

```

1: function THINCGM
2:    $(\mathbf{U}, \mathbf{\Sigma}, \mathbf{V}) \leftarrow (\mathbf{0}, \mathbf{0}, \mathbf{0})$  and  $\mathbf{y} \leftarrow \mathbf{0}$ 
3:   for  $t \leftarrow 0, 1, 2, 3, \dots$  do
4:      $(\mathbf{u}, \mathbf{v}) \leftarrow \text{MaxSingVec}(\mathcal{A}^*(\nabla f(\mathbf{y})))$ 
5:      $\mathbf{h} \leftarrow \mathcal{A}(-\alpha \mathbf{u} \mathbf{v}^*)$ 
6:     if  $\langle \mathbf{y} - \mathbf{h}, \nabla f(\mathbf{y}) \rangle \leq \varepsilon$  then break for
7:     end if
8:      $\eta \leftarrow 2/(t+2)$ 
9:      $\mathbf{y} \leftarrow (1-\eta)\mathbf{y} + \eta\mathbf{h}$ 
10:     $(\mathbf{U}, \mathbf{\Sigma}, \mathbf{V}) \leftarrow \text{SVDUPDATE}(\mathbf{U}, \mathbf{\Sigma}, \mathbf{V}, \mathbf{u}, \alpha, \mathbf{v}, \eta)$ 
11:  end for
12:  return  $(\mathbf{U}, \mathbf{\Sigma}, \mathbf{V})$ 
13: end function

```

---

```

14: function SVDUPDATE( $\mathbf{U}, \mathbf{\Sigma}, \mathbf{V}, \mathbf{u}, \alpha, \mathbf{v}, \eta$ )
15:    $\mathbf{m} \leftarrow \mathbf{U}^* \mathbf{u}$  and  $\mathbf{n} \leftarrow \mathbf{V}^* \mathbf{v}$ 
16:    $\mathbf{p} \leftarrow \mathbf{u} - \mathbf{U} \mathbf{m}$  and  $\mathbf{q} \leftarrow \mathbf{v} - \mathbf{V} \mathbf{n}$ 
17:    $\hat{\mathbf{p}} \leftarrow \mathbf{p} / \|\mathbf{p}\|$  and  $\hat{\mathbf{q}} \leftarrow \mathbf{q} / \|\mathbf{q}\|$ 
18:    $(\hat{\mathbf{U}}, \hat{\mathbf{\Sigma}}, \hat{\mathbf{V}}) \leftarrow \text{svd} \left( (1-\eta) \begin{bmatrix} \mathbf{\Sigma} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix} - \eta \alpha \begin{bmatrix} \mathbf{m} \\ \|\mathbf{p}\| \end{bmatrix} \begin{bmatrix} \mathbf{n} \\ \|\mathbf{q}\| \end{bmatrix}^* \right)$ 
19:    $\mathbf{U} \leftarrow [\mathbf{U} \hat{\mathbf{U}} \quad \hat{\mathbf{p}} \hat{\mathbf{U}}]$  and  $\mathbf{V} \leftarrow [\mathbf{V} \hat{\mathbf{V}} \quad \hat{\mathbf{q}} \hat{\mathbf{V}}]$  and  $\mathbf{\Sigma} \leftarrow \hat{\mathbf{\Sigma}}$ 
20:   return  $(\mathbf{U}, \mathbf{\Sigma}, \mathbf{V})$ 
21: end function

```

---

Suppose that the CGM iteration (2.2)–(2.5) generates the sequence  $(\mathbf{X}_t : t = 0, 1, 2, \dots)$ . It is easy to verify that the ThinCGM iteration generates the sequences  $(\mathbf{U}_t, \mathbf{\Sigma}_t, \mathbf{V}_t : t = 0, 1, 2, \dots)$  such that  $\mathbf{U}_t \mathbf{\Sigma}_t \mathbf{V}_t^* = \mathbf{X}_t$ . In other words, ThinCGM is a careful implementation of the standard CGM, rather than a new method. Consequently, ThinCGM maintains the same convergence guarantees as the standard method.

While ThinCGM is not theoretically guaranteed to alleviate the expansion of the rank, it often keeps a low memory footprint in comparison with the classical CGM. The weak point of ThinCGM is that the rank of  $\mathbf{X}_t$  can increase with  $t$ . There is no bound on the rank, beyond the fact that it is clearly less than  $t$  (recall that each step brings a rank-one update). CGM exhibits a slow, sublinear convergence rate, which means that we need to run it for many iterations before reaching an approximate solution. As a result, the peak rank of the iterates is often much larger than the rank of the solution.

**Remark.** After we submitted [YHC15] for review, a discussion on the idea of using thin SVD updates to implement CGM-type methods appeared also in a concurrent work by Freund et al., see Section 3.6 in [FGM17].

### 2.2.3 A Numerical Study on the Rank Expansion

ThinCGM is a CGM implementation for (2.1), which attempts to maintain a low-rank representation of the decision variable. The intermediate estimates, however, are not guaranteed to be low-rank, even when the sequence starts from and converges towards low-rank points.

Consider the matrix completion problem that arises in machine learning applications such as collaborative filtering [SRJ04]. Suppose that we are trying to construct an unknown, low-rank matrix  $\mathbf{X}_\dagger \in \mathbb{R}^{m \times n}$  from a small sample set  $E$  of its entries:

$$b_{ij} = (\mathbf{X}_\dagger)_{ij} + \xi_{ij} \quad \text{for } (i, j) \in E, \quad (2.20)$$

where  $\xi_{ij}$  models the noise. Let us consider a Gaussian noise model, under which we can formulate the problem as an instance of (2.1) with the following setting:

- ▷ Each coefficient matrix (2.18) associated with  $\mathcal{A}$  is a distinct, one-sparse  $(m \times n)$  matrix.
- ▷  $\mathbf{b} := \mathcal{A}\mathbf{X}_\dagger + \boldsymbol{\xi} \in \mathbb{R}^d$  is a vectorized representation of samples (2.20).
- ▷  $f(\mathbf{z}) := \frac{1}{2} \|\mathbf{z} - \mathbf{b}\|^2$  corresponds to negative log-likelihood estimator of the Gaussian model.

We consider MovieLens 100K dataset [HK16], which consists of 100'000 ratings (1 to 5) from 943 users in 1682 movies. We use a very basic preprocessing step, which only removes the movies that are not rated by any user, and the users that have not provided any ratings.

Figure 2.1 shows the expansion of the rank. On the left panel, we draw the evolution of  $\varepsilon$ -rank of the estimate. By  $\varepsilon$ -rank, we mean the number of singular values that exceed  $\varepsilon\sigma_1$ , where  $\sigma_1$  is the largest singular value. On the right panel, we plot the singular singular value spectrum of the ground truth, obtained by solving the optimization problem to high accuracy. It is clear from this picture that the rank of the intermediate iterates can be much higher than the rank of the solution.

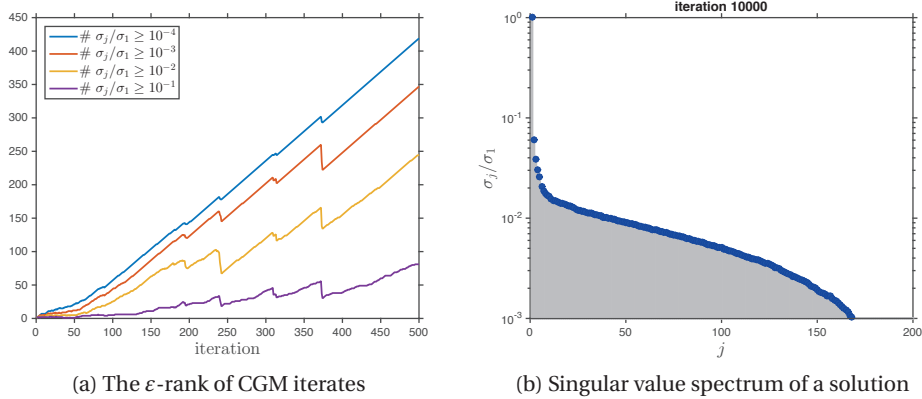


Figure 2.1 – Rank expansion of CGM estimates. ThinCGM is applied to the matrix completion problem with MovieLens 100K dataset.

## 2.3 SketchyCGM

Our numerical study on matrix completion problem demonstrates that the rank increases steadily, even when the solution has an approximately low-rank representation with a fast decay in the singular value spectrum.

There also appears an auxiliary question in implementing ThinCGM. Due to numerical error, strictly speaking, each rank-one update increases the numerical rank. Then, we need to answer the question which singular values should we consider as noise and truncate. Without going into these details, we present a new approach with provable guarantees on the storage.

Notice that ThinCGM does not use variables  $(\mathbf{U}, \mathbf{\Sigma}, \mathbf{V})$  at any step, except for updating the variables itself. In other words, these variables are maintained only to output a solution in the end. Instead,  $\mathbf{y}_t$  contains all information about  $\mathbf{X}_t$  that we need to drive CGM iterations. This observation is the important piece to create a CGM variant with optimal storage guarantees.

We propose a new approach, SketchyCGM, which completely discards the exact storage of the decision variable  $\mathbf{X}$ , or its factors. Instead, we maintain a randomized linear image (a sketch) of  $\mathbf{X}$  which contains enough information to reconstruct a low-rank approximation of the solution once we converge.

Since the iterations are driven by the low-dimensional *dual* variable  $\mathbf{y}$ , SketchyCGM follows exactly the same solution path as the standard CGM. It does not, however, store matrix decision variable, and reconstructs it from a randomized sketch at the end of the optimization procedure. This reconstruction is provably accurate, if the solution attained by CGM is (approximately) low-rank.

Let us summarize a randomized sketching method that we introduced in [TYUC17b]. We provide a more complete overview of this approach in Chapter 6.

### 2.3.1 Randomized Sketch

Draw and fix two independent standard normal matrices  $\mathbf{\Omega}$  and  $\mathbf{\Psi}$  where

$$\mathbf{\Omega} \in \mathbb{R}^{n \times k} \quad \text{with} \quad k = 2r + 1 \quad \text{and} \quad \mathbf{\Psi} \in \mathbb{R}^{\ell \times m} \quad \text{with} \quad \ell = 4r + 3. \quad (2.21)$$

The sketch consists of two matrices  $\mathbf{Y}$  and  $\mathbf{W}$  that capture the range and co-range of  $\mathbf{X}$ :

$$\mathbf{Y} = \mathbf{X}\mathbf{\Omega} \in \mathbb{R}^{m \times k} \quad \text{and} \quad \mathbf{W} = \mathbf{\Psi}\mathbf{X} \in \mathbb{R}^{\ell \times n}. \quad (2.22)$$

We can efficiently update the sketch  $(\mathbf{Y}, \mathbf{W})$  to reflect a rank-one update (2.5).

## Chapter 2. Storage-Efficient Convex Optimization

---

The following procedure provides a rank- $r$  approximation  $\hat{\mathbf{X}}$  of matrix  $\mathbf{X}$  from its sketch:

$$\mathbf{Y} = \mathbf{Q}\mathbf{R} \quad \text{and} \quad \mathbf{B} = (\mathbf{\Psi}\mathbf{Q})^\dagger \mathbf{W} \quad \text{and} \quad \hat{\mathbf{X}} = \mathbf{Q}[\mathbf{B}]_r, \quad (2.23)$$

where  $\mathbf{Q}$  and  $\mathbf{R}$  are the orthogonal and triangular factors. This procedure yields a rank- $r$  approximation  $\hat{\mathbf{X}}$  that satisfies

$$\mathbb{E} \|\mathbf{X} - \hat{\mathbf{X}}\|_F \leq 3\sqrt{2} \|\mathbf{X} - [\mathbf{X}]_r\|_F. \quad (2.24)$$

Similar bounds hold with high probability. See Theorem 6.2 or Theorem 4.3 in [TYUC17b] for more details.

**Remark.** The recommended sketch size parameters  $(k, \ell)$  in (2.21) are chosen to balance storage against reconstruction quality.

### 2.3.2 SketchyCGM Iteration

Let us introduce SketchyCGM, the first storage-optimal convex optimization method that produces a low-rank approximation of a solution to (2.1) with rigorous guarantees.

Fix the target rank  $r$ . Optionally, also set a suboptimality parameter  $\varepsilon$ . Independently draw and fix standard normal test matrices  $\mathbf{\Omega} \in \mathbb{R}^{n \times k}$  and  $\mathbf{\Psi} \in \mathbb{R}^{\ell \times m}$  as in (2.21).

Initialize the iterate and the sketches:

$$\mathbf{y}_0 = \mathbf{0} \in \mathbb{R}^d \quad \text{and} \quad \mathbf{Y}_0 = \mathbf{0} \in \mathbb{R}^{m \times k} \quad \text{and} \quad \mathbf{W}_0 = \mathbf{0} \in \mathbb{R}^{\ell \times n}. \quad (2.25)$$

At each iteration  $t = 0, 1, 2, \dots$ , resolve a solution to the linear minimization subproblem:

$$\mathbf{h}_t = \mathcal{A}(-\alpha \mathbf{u}_t \mathbf{v}_t^*) \quad \text{where} \quad (\mathbf{u}_t, \mathbf{v}_t) = \text{MaxSingVec}(\mathcal{A}^*(\nabla f(\mathbf{y}_t))). \quad (2.26)$$

Set the learning rate  $\eta_t = 2/(t+2)$ , and update the variables:

$$\begin{aligned} \mathbf{y}_{t+1} &= (1 - \eta_t) \mathbf{y}_t + \eta_t \mathbf{h}_t; \\ \mathbf{Y}_{t+1} &= (1 - \eta_t) \mathbf{Y}_t + \eta_t (-\alpha \mathbf{u}_t \mathbf{v}_t^*) \mathbf{\Omega}; \\ \mathbf{W}_{t+1} &= (1 - \eta_t) \mathbf{W}_t + \eta_t \mathbf{\Psi} (-\alpha \mathbf{u}_t \mathbf{v}_t^*). \end{aligned} \quad (2.27)$$

If a stopping criterion is defined, continue this loop until the condition is satisfied:

$$\langle \mathbf{y}_t - \mathbf{h}_t, \nabla f(\mathbf{y}_t) \rangle \leq \varepsilon. \quad (2.28)$$

Otherwise, continue until the maximum number of iterations is reached.

Once the iterations terminate, form a rank- $r$  approximate solution  $\hat{\mathbf{X}}_t$  using the procedure (2.23).

See Algorithm 2.2 for a complete pseudocode.

**Algorithm 2.2** SketchyCGM for model problem (2.1)**Require:** Data for (2.1); suboptimality  $\varepsilon$ ; target rank  $r$ 

```

1: function SKETCHYCGM
2:   SKETCH.INIT( $m, n, r$ )
3:    $\mathbf{y} \leftarrow \mathbf{0}$ 
4:   for  $t \leftarrow 0, 1, 2, 3, \dots$  do
5:      $(\mathbf{u}, \mathbf{v}) \leftarrow \text{MaxSingVec}(\mathcal{A}^*(\nabla f(\mathbf{y})))$ 
6:      $\mathbf{h} \leftarrow \mathcal{A}(-\alpha \mathbf{u} \mathbf{v}^*)$ 
7:     if  $\langle \mathbf{y} - \mathbf{h}, \nabla f(\mathbf{y}) \rangle \leq \varepsilon$  then break for
8:     end if
9:      $\eta \leftarrow 2/(t+2)$ 
10:     $\mathbf{y} \leftarrow (1-\eta)\mathbf{y} + \eta\mathbf{h}$ 
11:    SKETCH.CGMUPDATE( $-\alpha \mathbf{u}, \mathbf{v}, \eta$ )
12:  end for
13:   $(\mathbf{U}, \mathbf{\Sigma}, \mathbf{V}) \leftarrow \text{SKETCH.RECONSTRUCT}()$ 
14:  return  $(\mathbf{U}, \mathbf{\Sigma}, \mathbf{V})$ 
15: end function

```

---

**Methods for SKETCH object**


---

```

16: function SKETCH.INIT( $m, n, r$ )
17:    $k \leftarrow 2r+1$  and  $\ell \leftarrow 4r+3$ 
18:    $\mathbf{\Omega} \leftarrow \text{randn}(n, k)$  and  $\mathbf{\Psi} \leftarrow \text{randn}(\ell, m)$ 
19:    $\mathbf{Y} \leftarrow \text{zeros}(m, k)$  and  $\mathbf{W} \leftarrow \text{zeros}(\ell, n)$ 
20: end function
21: function SKETCH.CGMUPDATE( $\mathbf{u}, \mathbf{v}, \eta$ )
22:    $\mathbf{Y} \leftarrow (1-\eta)\mathbf{Y} + \eta \mathbf{u}(\mathbf{v}^* \mathbf{\Omega})$ 
23:    $\mathbf{W} \leftarrow (1-\eta)\mathbf{W} + \eta(\mathbf{\Psi} \mathbf{u}) \mathbf{v}^*$ 
24: end function
25: function SKETCH.RECONSTRUCT()
26:    $\mathbf{Q} \leftarrow \text{orth}(\mathbf{Y})$ 
27:    $\mathbf{B} \leftarrow (\mathbf{\Psi} \mathbf{Q}) \setminus \mathbf{W}$ 
28:    $(\mathbf{U}, \mathbf{\Sigma}, \mathbf{V}) \leftarrow \text{svds}(\mathbf{B}, r)$ 
29:   return  $(\mathbf{Q} \mathbf{U}, \mathbf{\Sigma}, \mathbf{V})$ 
30: end function

```

---

The cost of storing a dual variable  $\mathbf{y}_t$  is  $\Theta(d)$ , and the test matrices  $(\mathbf{\Omega}, \mathbf{\Psi})$  and the sketches is  $(\mathbf{Y}, \mathbf{W})$  is  $\Theta(r(m+n))$ . The MaxSingVec subroutine can be implemented with a working storage of size  $\Theta(d+m+n)$ . Overall, total storage cost of SketchyCGM is  $\Theta(d+r(m+n))$ .

**Definition.** We call the amount of storage that we need to specify the problem data and to write an approximate solution as the optimal storage cost. For example, the optimal storage cost of the matrix completion problem in Section 2.2.3 is  $\Theta(d+r(m+n))$ , where  $\Theta(d)$  is required to store  $\mathbf{b}$  and  $\Theta(r(m+n))$  to write a rank- $r$  approximate solution. If the total storage cost of an algorithm does not exceed the optimal storage cost, we say that the algorithm is storage-optimal. On this regard, SketchyCGM is a storage-optimal method!

### 2.3.3 Convergence Results for SketchyCGM

Suppose that the CGM iterations yield a sequence  $(\mathbf{X}_t : t = 0, 1, \dots)$  which is converging towards a low-rank point. Then, the SketchyCGM iterations maintain the sketch of  $\mathbf{X}_t$ , from which we can recover an accurate approximation of the low-rank solution. The following theorem formalizes this intuition.

**Theorem 2.1.** *Suppose that the sequence  $(\mathbf{X}_t : t = 0, 1, \dots)$  constructed by CGM for solving (2.1) is converging to a matrix  $\mathbf{X}_{\text{cgm}}$ . Let  $\hat{\mathbf{X}}_t$  be the rank- $r$  approximation of  $\mathbf{X}_t$  produced by SketchyCGM. Then,*

$$\lim_{t \rightarrow \infty} \mathbb{E} \|\hat{\mathbf{X}}_t - \mathbf{X}_{\text{cgm}}\|_F \leq 3\sqrt{2} \|\mathbf{X}_{\text{cgm}} - [\mathbf{X}_{\text{cgm}}]_r\|_F. \quad (2.29)$$

*In particular, if  $\text{rank}(\mathbf{X}_{\text{cgm}}) \leq r$ , then*

$$\mathbb{E} \|\hat{\mathbf{X}}_t - \mathbf{X}_{\text{cgm}}\|_F \rightarrow 0. \quad (2.30)$$

**Corollary 2.1.** *Suppose that all solutions to problem (2.1) are rank- $r$  or less. Then SketchyCGM produces a sequence of approximations that approaches to the solution set  $\mathcal{X}_\star$ :*

$$\mathbb{E} \text{dist}_F(\hat{\mathbf{X}}_t, \mathcal{X}_\star) \rightarrow 0. \quad (2.31)$$

Recall that the rank of the intermediate estimates can get much higher than the solution rank. Obviously, the information stored in the sketches might not contain enough information to recover an accurate approximation of these intermediate estimates. As a result, a convergence rate guarantee for the general problem template is not achievable.

Nevertheless, we identified a setting where it is possible to prove a bound on the convergence rate. Naturally, this requires stronger assumptions. In particular, we assume that the unique solution is stable in the sense that the objective value increases as we get far from the solution.

**Theorem 2.2.** *Fix  $\kappa > 0$  and  $\nu > 0$ . Suppose the (unique) solution  $\mathbf{X}_\star$  of (2.1) has  $\text{rank}(\mathbf{X}_\star) \leq r$ , and that it is stable in the sense that*

$$f(\mathcal{A}\mathbf{X}) - f(\mathcal{A}\mathbf{X}_\star) \geq \kappa \|\mathbf{X} - \mathbf{X}_\star\|_F^\nu \quad (2.32)$$

*for all feasible  $\mathbf{X}$ . Then we have the error bound*

$$\mathbb{E} \|\hat{\mathbf{X}}_t - \mathbf{X}_\star\|_F \leq 6 \left( \frac{2C\kappa^{-1}}{t+2} \right)^{1/\nu} \text{ for } t = 0, 1, 2, \dots \quad (2.33)$$

*where  $C$  is the curvature constant of the problem (2.1). See Eqn. (3) in [Jag13] for the definition of the curvature constant. We can also replace  $C$  by  $L_f D^2$ , where  $L_f$  is the smoothness constant of  $f$  and  $D$  is the domain diameter.*

We refer to the supplements of [YUTC17] for the proofs.



## 2.4 Extensions for Semidefinite Programming

The principal idea of SketchyCGM can be adopted to design storage-optimal convex optimization algorithms for other structured convex programs. In particular, this dissertation focuses on the SDP template (1.4). We describe the key ingredients to generalize SketchyCGM for solving a broad set of SDP formulations.

First, notice the difference of the domain constraint between the two formulations. The specific domain constraint in CGM affects the linear minimization subproblem. Consider the positive semidefinite (PSD) cone constraint with bounded trace,  $\mathcal{X} = \{\mathbf{X} : \mathbf{X} \in \mathbb{S}_+^n, \text{Tr } \mathbf{X} = \alpha\}$ . In order to adapt ThinCGM and SketchyCGM to handle problems in this domain, we replace computation (2.26) with

$$\mathbf{h}_t = \mathcal{A}(-\alpha \mathbf{u}_t \mathbf{v}_t^*) \quad \text{where} \quad (\lambda_t, \mathbf{u}_t) = \text{MinEig}(\mathcal{A}^*(\nabla f(\mathbf{y}_t))). \quad (2.34)$$

MinEig returns the minimum eigenvalue  $\lambda_t$  and an associated eigenvector  $\mathbf{u}_t$  of a conjugate symmetric matrix. Similarly, we can also handle  $\mathcal{X} = \{\mathbf{X} : \mathbf{X} \in \mathbb{S}_+^n, \text{Tr } \mathbf{X} \leq \alpha\}$  by using

$$\mathbf{h}_t = \begin{cases} \mathcal{A}(\alpha \mathbf{u}_t \mathbf{u}_t^*), & \lambda_t \leq 0 \\ \mathbf{0}, & \text{otherwise,} \end{cases} \quad \text{where} \quad (\lambda_t, \mathbf{u}_t) = \text{MinEig}(\mathcal{A}^*(\nabla f(\mathbf{y}_t))). \quad (2.35)$$

We also make small changes in the matrix sketch, in order to preserve PSD structure in the reconstruction algorithm, as described in our Algorithm 9 in [TYUC17b]. Alternatively, we can also develop more efficient sketching methods specifically for PSD matrices, see Section 6.3 and Algorithm 3 [TYUC17a].

Recall that CGM does not apply for problems with affine constraints. Consequently, we cannot directly use SketchyCGM for solving SDP in standard form. One quick remedy is to use the standard reduction approach of the SDP template to a sequence of feasibility problems by performing a binary search over the optimal value of  $\langle \mathbf{C}, \mathbf{X} \rangle$ , see Section 1.2 in [Gar12]. In large-scale instances, however, the binary search over the model parameters imposes additional computational cost.

Instead, we design new practical optimization algorithms that are compatible with sketching for solving SDP template (1.4). We describe these methods in the following chapters.

We say SketchyCGM is storage-optimal because the working storage is  $\mathcal{O}(d + r(m + n))$ , which does not exceed the storage required to store the data and a solution. We conclude this section by highlighting a recent result from Waldspurger and Waters [WW18], which demonstrates that the non-convex Burer-Monteiro factorization idea (1.8) cannot lead to storage-optimal algorithms for the standard form SDP template. In particular, they show that the factorization rank must be set in the order Pataki-Barvinok rank even for some problems with a unique rank-one solution, in order to avoid spurious stationary points in the problem formulation. Accordingly, convexity seems to be the key to achieve storage-optimality!

## 2.5 Numerical Experiments

We present our numerical experiments on matrix completion and phase retrieval problems, to demonstrate flexibility, scalability and robustness of our algorithms.

### 2.5.1 Matrix completion

Instate the matrix completion problem description from Section 2.2.3. SketchyCGM is flexible, i.e., it can solve any instance of (2.1). In Section 2.2.3, we considered a Gaussian noise model, which yields a quadratic loss function. In this section, we also consider Gauss–Laplace and Bernoulli noise models, which correspond to the huber and logistic loss, respectively.

With logistic loss, we consider a binary classification model. Only for this experiment, we binarize the data using the formula  $\mathbf{b} \leftarrow \text{sign}(\mathbf{b} - 3.5)$  as part of the preprocessing. For other experiments, with quadratic and huber loss, we keep the ratings intact, with no preprocessing.

For the MovieLens 100K dataset, we use the default ub train and test partition of the data. For each loss function, we sweep  $\alpha$  from 3'000 to 10'000 in steps of 500. We select the value of  $\alpha$  that provides the best test error after 10'000 iterations of CGM. A similar procedure applies to the MovieLens 10M dataset with the rb partition. This time, we sweep  $\alpha$  from 50'000 to 250'000 in steps of 25'000.

We compute the test error as follows: At the end of each iteration, we execute sketch reconstruction procedure to obtain a rank- $r$  approximate solution. Then, we evaluate the loss function for each datapoint in the test set, and compute the average loss of the test partition. We report this average value as the test error.

Figures 2.2 and 2.3 present the evolution of test error for SketchyCGM as a function of the iteration counter, for various values of the rank parameter  $r$  in SketchyCGM. For the 100K dataset, rank  $r = 50$  yields test error similar with the CGM solution. For the 10M dataset, rank  $r = 200$  yields equivalent performance.

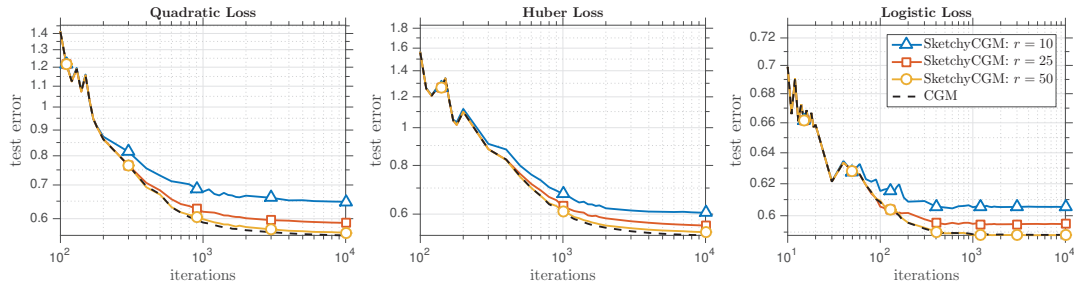


Figure 2.2 – *MovieLens 100K Dataset*. Convergence of test error for SketchyCGM with three loss functions. The parameter  $\alpha$  is chosen by cross-validation on the final CGM solution:  $\alpha = 7,000$  for quadratic loss, 7,500 for Huber loss, and 4,500 for logistic loss.

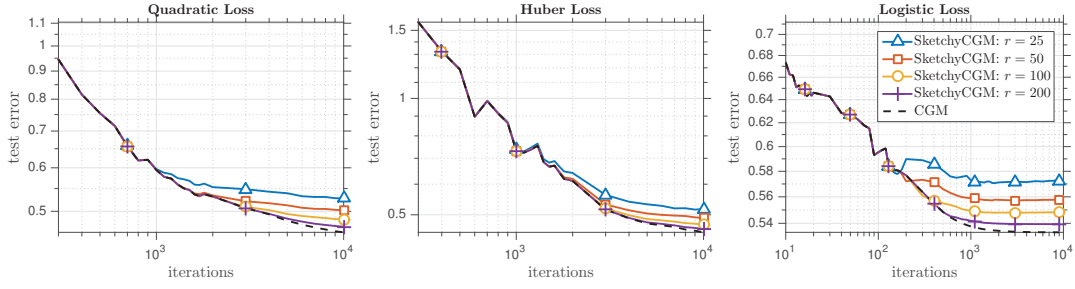


Figure 2.3 – *MovieLens 10M Dataset*. Convergence of test error for SketchyCGM with three loss functions. The parameter  $\alpha$  is chosen by cross-validation on the final CGM solution:  $\alpha = 150,000$  for quadratic and logistic losses, and 175,000 for Huber loss.

### 2.5.2 Phase Retrieval Problems

The problem of retrieving the phase of a signal from its intensity-only measurements has regained significant attention recently [CSV12, SEC<sup>+</sup>15, CZH<sup>+</sup>17]. Formally, phase retrieval problem aims to recover a signal  $\mathbf{x}^\natural \in \mathbb{C}^n$  from  $d$  noisy quadratic measurements

$$b_i = |\langle \mathbf{a}_i, \mathbf{x}_i^\natural \rangle|^2 + \xi_i \quad \text{for } i = 1, \dots, d. \quad (2.36)$$

where  $\mathbf{a}_i \in \mathbb{C}^n$  are the given measurement vectors, and  $\xi_i$  model the unknown noise. This problem arises in many applications, including X-ray crystallography, astronomical imaging and microscopy imaging.

As an estimation problem, the nonlinear observation model (2.36) poses significant difficulties, since the standard maximum likelihood estimators yield non-convex optimization problems. Convex relaxation is useful in this setting. Let us recast the quadratic measurements (2.36) as

$$\begin{aligned} b_i &= |\langle \mathbf{a}_i, \mathbf{x}_i^\natural \rangle|^2 + \xi_i \\ &= \mathbf{a}_i^* \mathbf{x}_i^\natural \mathbf{x}_i^{\natural*} \mathbf{a}_i + \xi_i \\ &= \text{Tr}(\mathbf{a}_i^* \mathbf{x}_i^\natural \mathbf{x}_i^{\natural*} \mathbf{a}_i) + \xi_i \\ &= \text{Tr}(\mathbf{a}_i \mathbf{a}_i^* \mathbf{x}_i^\natural \mathbf{x}_i^{\natural*}) + \xi_i \\ &= \langle \mathbf{A}_i, \mathbf{X}_i \rangle + \xi_i \quad \text{where} \quad \mathbf{A}_i = \mathbf{a}_i \mathbf{a}_i^*, \quad \text{and} \quad \mathbf{X}_i = \mathbf{x}_i^\natural \mathbf{x}_i^{\natural*} \in \mathbb{S}_+^n \end{aligned} \quad (2.37)$$

This leads to the following linear observation model in the *lifted* dimensions:

$$\mathbf{b} = \mathcal{A} \mathbf{X}_\natural + \boldsymbol{\xi}, \quad \text{where} \quad \mathbf{X}_\natural = \mathbf{x}_\natural \mathbf{x}_\natural^* \in \mathbb{S}_+^n \quad (2.38)$$

As a result, data fidelity costs can be measured using convex functions, albeit in terms of matrix variables with rank constraint. We relax the problem to a convex formulation by dropping the rank constraint and obtain an SDP problem.

### Scalability Test with Synthetic Phase Retrieval with Gaussian Noise

Let us first investigate the scalability of our approach on a synthetic phase retrieval setup. We consider a measurement model based on random coded diffraction patterns with octonary modulation; see [CLS15b].

We construct a synthetic data, here are the details. We draw a vector  $\mathbf{x}_1 \in \mathbb{C}^n$  from the standard normal distribution. Then we acquire  $d = 10n$  intensity measurements (2.36), with *iid* Gaussian noise such that the SNR is 20 dB. We choose  $\alpha = d^{-1} \sum_{i=1}^d b_i$ , which gives an accurate estimate. See Section II of [YHC15] for more details on this particular choice of  $\alpha$ .

We compare five convex optimization algorithms: the classic proximal gradient method (PGM) [Roc76]; the Auslender–Teboulle (AT) accelerated method [AT06]; the classic CGM algorithm [Jag13]; the PSD variants of ThinCGM and SketchyCGM (with rank parameter  $r = 1$ ).

We solve the following convex optimization template

$$\underset{\mathbf{X} \in \mathbb{C}^{n \times n}}{\text{minimize}} \quad f(\mathcal{A}\mathbf{X}) \quad \text{subject to} \quad \mathbf{X} \in \mathbb{S}_+^n, \quad \text{Tr} \mathbf{X} \leq \alpha, \quad (2.39)$$

with quadratic loss ( $f(\mathcal{A}\mathbf{X}) := \frac{1}{2} \|\mathcal{A}\mathbf{X} - \mathbf{b}\|^2$ ).

Figure 2.4(A) displays storage costs for each algorithm as the signal length  $n$  increases. We approximate memory usage by reporting the total workspace allocated by MATLAB for the algorithm. PGM, AT, and CGM have static allocations, but they use a matrix variable of size  $n^2$ . ThinCGM attempts to maintain a low-rank approximation of the decision variable, but the rank increases steadily, so the algorithm fails after  $n = 10^5$ . In contrast, SketchyCGM has a static memory allocation of  $\Theta(n)$ . It already offers the best memory footprint for  $n = 10^2$ , and it still works when  $n = 10^6$ .

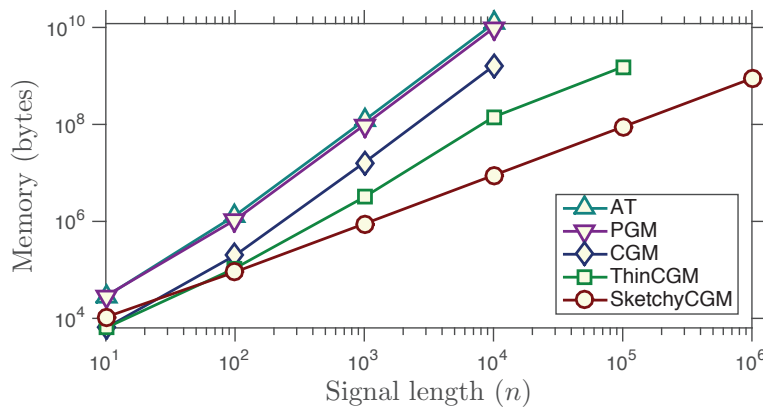


Figure 2.4 – *Memory Scaling*. Memory scaling for five convex optimization algorithms applied to a synthetic instance of the convex phase retrieval problem.

### Flexibility Test with Synthetic Phase Retrieval with Poisson Noise

Poisson noise model can capture the physics of an imaging system better than a Gaussian noise model [SEC<sup>+</sup>15, OLY<sup>+</sup>16, CLDM18]. In this section, we demonstrate the flexibility of SketchyCGM by solving (2.39) with Poisson loss function:  $f(\mathbf{z}) := \sum_i z_i - b_i \log(z_i)$ .

Remark that Poisson loss is not a smooth function, hence the standard CGM does not directly apply. In [OLY<sup>+</sup>16], however, we show that CGM can be adapted to solve this problem with minor modifications on the step-size and initialization: We initialize the algorithm with the dual vector  $\mathbf{y}_0 = d^{-1/2} \mathbf{1}$  and set the learning rate  $\eta_t = 2/(3+t)$ .

We consider a similar setup to the previous section with Gaussian noise model, but we use real images rather than random signals so that we can visualize the effects of choosing the appropriate noise model. We choose a small image  $\mathbf{x}_0 \in \mathbb{C}^n$  with  $n = 240 \times 320 = 76'800$  pixels. We acquire  $d = 20n$  measurements of the form (2.36) using the same coded diffraction model as in Gaussian noise model. Each  $\xi_i \in \mathbb{R}^d$  is drawn *iid* from a Poisson distribution. We choose the mean of this distribution so that the SNR is 20 dB. Once again, we set  $\alpha$  to the average of the measurements  $b_i$ .

We set the rank parameter  $r = 1$ , and run SketchyCGM for 100 iterations. To recover an approximate solution of the imaging problem, we take the top singular vector of the solution.

Figure 2.5 displays the recoveries of obtained. We compare against the recovery obtained with the Gaussian noise model. As expected, the Poisson formulation performs better.

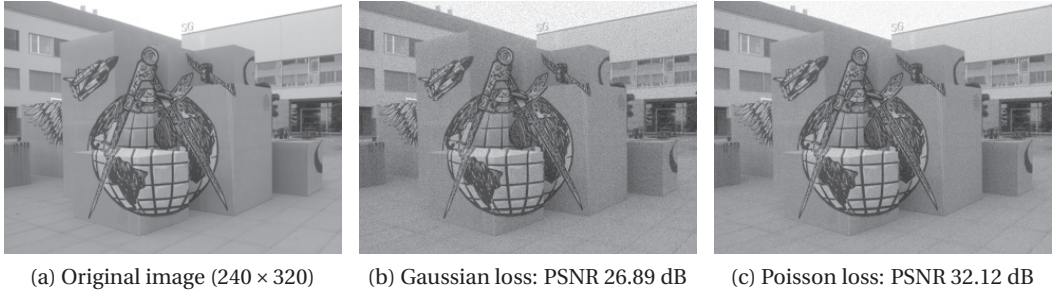


Figure 2.5 – *Gaussian and Poisson phase retrieval under Poisson noise*. See Section 2.5.2 for details.

### Fourier Ptychography

Lifting yields an obvious curse of dimensionality, expanding the search space from  $n$  to  $n^2$  dimensions. As a result, SDP relaxations have been thought to be impractical for phase retrieval problems of a realistic size.

In this section, we demonstrate an application of SketchyCGM to a phase retrieval problem from a real-world microscopy imaging application. The authors of [HCO<sup>+</sup>15] provided measurements of a slide containing human blood cells from a working Fourier ptychography microscopy system. See [HCO<sup>+</sup>15] for the technical details of this imaging modality.

The measurement vectors  $\mathbf{a}_i$  are obtained from windowed discrete Fourier transforms. The data consists of 29 image samples of  $(80 \times 80)$  pixels from different illuminations. Hence, the number of measurements  $d = 29 \cdot 80 \cdot 80 = 185'600$ . The final image resolution is  $(160 \times 160)$ , which we treat as a vector of size  $25'600$ . After lifting, this yields a complex valued  $(25'600 \times 25'600)$  dimensional PSD matrix problem. We consider the phase retrieval problem with Gaussian noise model (*i.e.*, with quadratic loss function), and we choose  $\alpha = 1'400$ .

We run SketchyCGM with rank parameter  $r = 1$  for  $10'000$  iterations. We use the standard basic rounding step, which takes the leading eigenvector of the solution. We compare against two non-convex approaches. The Wirtinger Flow method [CLS15b] that we implemented with the recommended parameters, and the Burer–Monteiro method [BM03] (reconstruction provided by the authors of [HCO<sup>+</sup>15]).

First row of Figure 2.7 displays the phase of the reconstruction obtained as a result of this procedure. Roughly, the phase indicates the thickness of the sample at a given location. Second and third rows show the phase gradient to provide an alternative view. SketchyCGM produces a superior quality image than the non-convex heuristics. These results are consistent with Figure 4 in [HCO<sup>+</sup>15], which indicates 5–10 dB improvement of convex optimization over non-convex heuristics.

Figure 2.6 displays snapshots of the SketchyCGM iterates as the algorithm proceeds. We see that SketchyCGM already achieve a diagnostic quality image even after  $1'000$  iterations.

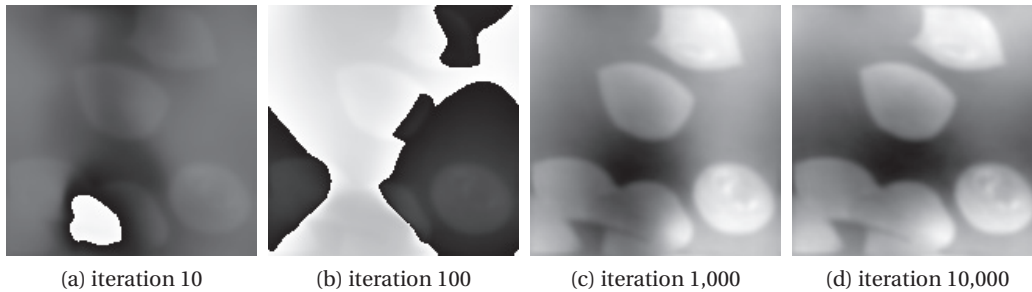


Figure 2.6 – *Evolution of the Reconstruction Quality of SketchyCGM for Fourier Ptychography.* Hue indicates the complex phase of the signal.



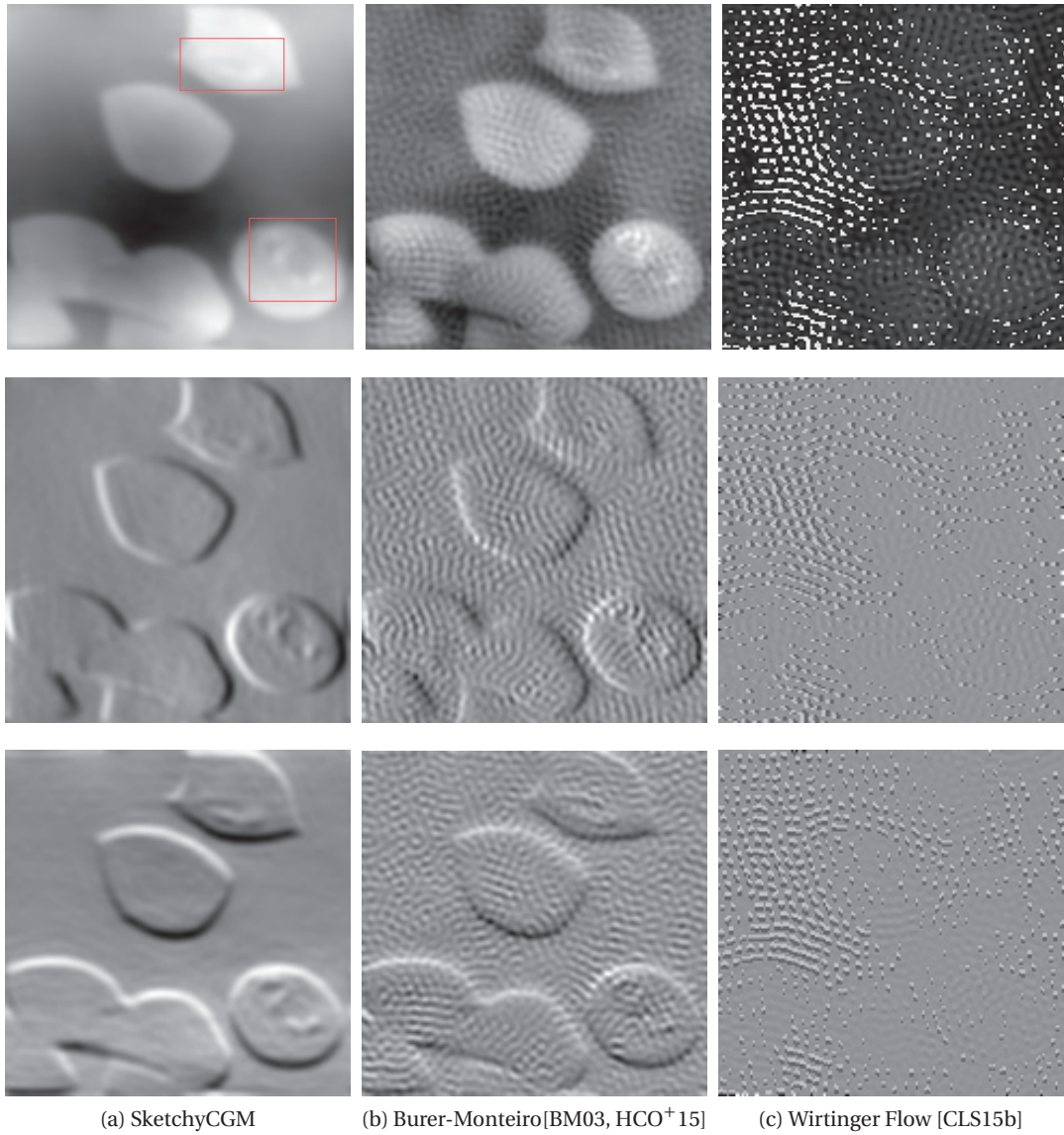


Figure 2.7 – *Three algorithms for Fourier ptychographic imaging via phase retrieval.* [top] Brightness indicates the complex phase of a pixel. Only relative differences in brightness are meaningful. Red boxes mark malaria parasites in blood cells. [middle] The horizontal differences of the phase maps presented. [bottom] The vertical differences of the phase maps.





## 3 Universal Primal-Dual Methods

In Chapter 2, we described a storage-optimal convex optimization paradigm for solving low-rank matrix optimization problems. Two major efficiencies of the conditional gradient method (CGM) play a crucial role in our paradigm:

- CGM executes rank one updates on the decision variable  $\mathbf{X}$ ,
- Even if we discard the decision variable  $\mathbf{X}$  from the workspace, CGM iterations can be driven by a low-dimensional “dual” variable.

In order to extend our framework for solving the SDP formulation (1.4), we need to develop algorithms for the model problem, which maintain these two characteristic features.

In this chapter, we describe our primal-dual (sub)gradient methods with these two features. Our algorithms are based on the Nesterov’s universal gradient method [Nes15] applied to the non-smooth dual formulation of the model problem (1.1). We equip our algorithms with a new primal averaging mechanism to construct the primal decision variable through structured updates, which exhibits rank one updates when we use our algorithms for solving instances from the SDP formulation (1.4).

The content in this chapter is based on the joint work with Quoc Tran-Dinh and Volkan Cevher [YTDC15a].

### Introduction

Processing the affine constraint ( $\mathcal{A}\mathbf{x} \in \mathcal{K}$ ) in (1.4) require significant computational effort in large-scale setting. The majority of scalable numerical algorithms for constrained optimization are primal-dual-type, including decomposition, augmented Lagrangian and alternating direction methods: See [CP08, BPC<sup>+</sup>11, GEB13, ST14, TDC14, LM15] and the references therein.

Primal-dual methods are a natural choice for designing storage-optimal methods. This is simply because we can use the low-dimensional dual variable to drive iterations in the dual space. Nevertheless, the majority of the primal-dual methods rely on a proximal-type (or

### Chapter 3. Universal Primal-Dual Methods

---

projection-type) oracle to recover the primal variable from the dual methods. As the problem dimensions become increasingly larger, the proximal tractability assumption can be restrictive. More importantly, the proximal operator requires full dimensional computations in the primal problem space, which eliminates the storage benefits.

#### Contributions

This chapter describes our universal primal-dual convex optimization algorithms.

- Our algorithmic approach features a gradient method and its accelerated variant that operates on the dual formulation of (1.1). Given the dual method, we then use a new averaging scheme to construct primal-iterates for the problem template (1.1). This averaging scheme results in streaming rank one updates on the primal decision variable when solving our model SDP template (1.4). Consequently, our algorithms are amenable to thin SVD updates for increased storage efficiency, or to sketching for storage optimality.
- Our algorithms trade-off the computational difficulty in the primal subproblems with the difficulty of optimization in the dual subproblems by avoiding the proximal mappings and using a subgradient approach instead. The dual subgradient mapping requires evaluation of a Fenchel-type oracle that we call as the *sharp operator*. Sharp operator can be viewed as a generalization of the linear minimization oracle (lmo) of CGM. In particular, when applied to the SDP template (1.4), our sharp operator can be cast as an instance of the lmo.
- Unlike CGM, our approach can handle non-smooth objectives and affine constraints in (1.4) efficiently. Due to the special structure of the dual problem, we can always assume that it has some Hölder smoothness (also called as weak smoothness) of order  $\nu \in [0, 1]$ . As a result, we tailor the universal gradient algorithms of Nesterov [Nes15] for solving the dual problem. Although our algorithms do not require the knowledge of the smoothness parameters a priori, they achieve the optimal worst-case complexity guarantees in the sense of first-order black-box models [NY83], with respect to the smoothness order of the *dual problem*.
- We extend the convergence analysis in Nesterov’s universal gradient methods [Nes15] to the primal-dual setting. We provide a rigorous analysis of the convergence rate and the worst-case complexity of these algorithms.
- We present an accelerated variant of our algorithm based on the combination of our analysis with the FISTA approach [BT09]. This variant requires fewer evaluations of projection onto  $\mathcal{K}$  than a direct extension of the fast universal gradient method of Nesterov [Nes15].
- We present numerical experiments in matrix completion and quantum tomography applications to demonstrate the flexibility and scalability of the proposed approach.

### 3.1 Preliminaries

Let us first recall some basic notions from the primal-dual optimization.

Introduce a slack variable  $\mathbf{v} = \mathcal{A}\mathbf{x} \in \mathcal{K}$ , and form the Lagrangian:

$$\mathcal{L}(\mathbf{x}, \mathbf{v}, \mathbf{y}) := f(\mathbf{x}) + \langle \mathbf{y}, \mathcal{A}\mathbf{x} - \mathbf{v} \rangle. \quad (3.1)$$

Minimize the Lagrangian over the constraint set to obtain the dual function  $\phi$ :

$$\begin{aligned} \phi(\mathbf{y}) &:= \min_{\mathbf{x} \in \mathcal{X}} \min_{\mathbf{v} \in \mathcal{K}} \mathcal{L}(\mathbf{x}, \mathbf{v}, \mathbf{y}) = \min_{\mathbf{x} \in \mathcal{X}} \min_{\mathbf{v} \in \mathcal{K}} f(\mathbf{x}) + \langle \mathbf{y}, \mathcal{A}\mathbf{x} - \mathbf{v} \rangle \\ &= \min_{\mathbf{x} \in \mathcal{X}} f(\mathbf{x}) + \langle \mathcal{A}^* \mathbf{y}, \mathbf{x} \rangle - \max_{\mathbf{v} \in \mathcal{K}} \langle \mathbf{y}, \mathbf{v} \rangle. \end{aligned} \quad (3.2)$$

Now we can write the dual problem as

$$\underset{\mathbf{y} \in \mathbb{R}^d}{\text{maximize}} \quad \phi(\mathbf{y}). \quad (3.3)$$

This is an unconstrained concave maximization problem.

To follow the common notational convention, we define the negation of the dual function as  $g(\mathbf{y}) := -\phi(\mathbf{y})$  and work with a convex minimization formulation. In the rest of this chapter, we will call  $g$  as the dual function relaxing the mathematical terminology, but recall that it is in fact the negation of it. We decompose  $g$  as

$$\begin{aligned} g(\mathbf{y}) &= g_x(\mathbf{y}) + g_v(\mathbf{y}) \quad \text{where} \quad g_x(\mathbf{y}) = -\min_{\mathbf{x} \in \mathcal{X}} [f(\mathbf{x}) + \langle \mathcal{A}^* \mathbf{y}, \mathbf{x} \rangle] \\ &\quad \text{and} \quad g_v(\mathbf{y}) = \max_{\mathbf{v} \in \mathcal{K}} \langle \mathbf{y}, \mathbf{v} \rangle. \end{aligned} \quad (3.4)$$

Now we can write an equivalent (negated) dual problem as

$$\underset{\mathbf{y} \in \mathbb{R}^d}{\text{minimize}} \quad g_x(\mathbf{y}) + g_v(\mathbf{y}). \quad (3.5)$$

This is a composite convex minimization formulation. Since  $g_x$  and  $g_v$  are generally non-smooth, FISTA and its proximal-based analysis [BT09] are not directly applicable.

### 3.2 Universal Primal-Dual Method

In this section, we introduce our universal primal-dual (sub)gradient method, UPD, for solving (1.1). UPD applies Nesterov's universal gradient method from [Nes15] for solving the dual problem (3.5). To recover the primal variable, we introduce an ergodic averaging mechanism.

Let us first define a Fenchel-type oracle, that we call as the *sharp operator* in the rest of this dissertation. The sharp operator plays an important role in our algorithmic design. Given an

### Chapter 3. Universal Primal-Dual Methods

---

input variable  $\mathbf{x} \in \mathbb{R}^p$ , the sharp operator returns  $\mathbf{h} \in \mathbb{R}^p$  such that

$$\mathbf{h} = \text{sharp}_{f, \mathcal{X}}(\mathbf{x}) \quad \text{where} \quad \mathbf{h} \in \arg \min_{\mathbf{u} \in \mathcal{X}} f(\mathbf{u}) + \langle \mathbf{x}, \mathbf{u} \rangle. \quad (3.6)$$

We can view the sharp operator as a generalization of the lmo of CGM, as it becomes a linear minimization subproblem over the domain  $\mathcal{X}$  when  $f$  is a linear function.

Using the sharp operator, we can write a subgradient of  $g_x$ , defined in (3.5), as follows:

$$\nabla g_x(\mathbf{y}) = -\mathcal{A}\mathbf{h} \quad \text{where} \quad \mathbf{h} = \text{sharp}_{f, \mathcal{X}}(\mathcal{A}^* \mathbf{y}). \quad (3.7)$$

#### 3.2.1 Efficiency Considerations and Hölder Smoothness

The difficulty of solving an unconstrained optimization problem depends on the problem geometry, *i.e.*, the functional properties of the objective function, such as the smoothness and the strong convexity.

For instance, when the objective function is smooth, we can use accelerated first-order methods with  $\mathcal{O}(1/\sqrt{\varepsilon})$  iteration complexity, to achieve an  $\varepsilon$ -suboptimal solution. On the other hand, iteration complexity for solving non-smooth problems is incomparably worse, on the order of  $\mathcal{O}(1/\varepsilon^2)$ .

The classical optimization literature tackles smooth and non-smooth (with bounded subgradient) problems with distinct algorithms, specifically designed to get the optimal guarantees for the target problem. However, designing a single algorithm which works for both cases, without any prior information about the smoothness, is desirable for two practical reasons:

1. Smoothness of the problem might be unknown a priori, or it might be expensive to evaluate the smoothness parameters.
2. Locally, the problem might possess more favorable structures than the global geometry. We might expect a better empirical performance from an adaptive method in this case.

Motivated by its practical significance, there is a recent attempt to design *universal* gradient methods. Universal, in this context, means that the method can adapt to the unknown smoothness level and smoothness parameters of the problem, as defined by Nesterov in [Nes15].

Nesterov's universal gradient methods are the first to achieve this goal [Nes15]. His methods are based on an inexact line-search strategy with an additional slackness term. More recently, some line-search-free universal methods are introduced [Lev17]. These methods typically rely on an online-to-offline conversion of the AdaGrad method [DHS11]. We have also contributed to this line of research, by introducing the first line-search-free accelerated method with learning rate adaptation [LYC18].

**Hölder smoothness.** To close the gap between the complexity bounds for smooth and non-smooth problems, Nesterov's universal gradient methods also consider the Hölder smoothness as the intermediate class.

Let  $\nabla g_x$  be a subgradient of  $g_x$ . We define the Hölder smoothness constant of order  $\nu \geq 0$  as

$$M_\nu := \max_{\substack{\mathbf{y}_1, \mathbf{y}_2 \in \mathbb{R}^d \\ \mathbf{y}_1 \neq \mathbf{y}_2}} \frac{\|\nabla g_x(\mathbf{y}_1) - \nabla g_x(\mathbf{y}_2)\|}{\|\mathbf{y}_1 - \mathbf{y}_2\|^\nu} \quad (3.8)$$

Note that the parameter  $M_\nu$  explicitly depends on  $\nu$ . We are interested in the case  $\nu \in [0, 1]$ , and *especially the two extremal cases*:  $\nu = 1$  corresponds to the smooth problems with Lipschitz continuous gradients, and  $\nu = 0$  corresponds to the non-smooth problems with bounded subgradients.

**Assumption.** We assume that there exists a finite smoothness constant  $M$  defined as

$$M := \min_{0 \leq \nu \leq 1} M_\nu(g) < +\infty. \quad (3.9)$$

This is a reasonable assumption. We explain this claim with the following examples:

◦ If  $\mathcal{X}$  is bounded, then  $\nabla g_x(\cdot)$  is also bounded. Indeed, we have

$$\begin{aligned} \|\nabla g_x(\mathbf{y})\| &= \|\mathcal{A}^* \mathbf{h}\| \leq D_{\mathcal{X}} \|\mathcal{A}\| \quad \text{since} \quad \mathbf{h} = \text{sharp}_{f, \mathcal{X}}(\mathcal{A}^* \mathbf{y}) \\ &\in \arg \min_{\mathbf{u} \in \mathcal{X}} f(\mathbf{u}) + \langle \mathcal{A}^* \mathbf{y}, \mathbf{u} \rangle. \end{aligned} \quad (3.10)$$

Hence, we can choose  $\nu = 0$  and  $M = D_{\mathcal{X}} \|\mathcal{A}\|$ . Our SDP template (1.4) fits into this case.

◦ If  $f$  is uniformly convex with convexity parameter  $\mu_f > 0$  and degree  $q \geq 2$ , *i.e.*,

$$\langle \nabla f(\mathbf{x}_1) - \nabla f(\mathbf{x}_2), \mathbf{x}_1 - \mathbf{x}_2 \rangle \geq \mu_f \|\mathbf{x}_1 - \mathbf{x}_2\|^q \quad \forall \mathbf{x}_1, \mathbf{x}_2 \in \mathbb{R}^p, \quad (3.11)$$

then  $g_x$  satisfies (3.8) with  $\nu = \frac{1}{q-1}$  and  $M = (\mu_f^{-1} \|\mathcal{A}\|^2)^{\frac{1}{q-1}}$ , as shown in [Nes15].

◦ In particular, if  $f$  is  $\mu_f$ -strongly convex (*i.e.*, uniformly convex with degree  $q = 2$ ), then  $\nu = 1$  and  $g_x$  is smooth with parameter  $M = \mu_f^{-1} \|\mathcal{A}\|^2$ .

### 3.2.2 UPD Iteration

Set the desired accuracy  $\varepsilon > 0$ , and an initial step-size  $\sigma_{-1}$ . Initialize the primal and dual variable, and the algorithm parameters

$$\mathbf{x}_0 \in \mathbb{R}^p \quad \text{and} \quad \mathbf{y}_0 \in \mathbb{R}^d \quad \text{and} \quad S_{-1} = 0. \quad (3.12)$$

### Chapter 3. Universal Primal-Dual Methods

---

At each iteration  $t = 0, 1, 2, \dots$ , compute a subgradient using the sharp operator:

$$\nabla g_x(\mathbf{y}_t) = -\mathcal{A}\mathbf{h}_t \quad \text{where} \quad \mathbf{h}_t = \text{sharp}_{f, \mathcal{X}}(\mathcal{A}^* \mathbf{y}_t). \quad (3.13)$$

Then, evaluate the partial dual objective

$$g_x(\mathbf{y}_t) = -f(\mathbf{h}_t) - \langle \mathbf{y}_t, \mathcal{A}\mathbf{h}_t \rangle \quad (3.14)$$

UPD executes a line-search to find an appropriate dual step-size as follows:

For  $j = -1, 0, 1, \dots$ , compute a trial point

$$\begin{aligned} \hat{\mathbf{y}} &= \text{prox}_{\hat{\sigma} g_v}(\mathbf{y}_t - \hat{\sigma} \nabla g_x(\mathbf{y}_t)) \\ &= \mathbf{y}_t - \hat{\sigma} \nabla g_x(\mathbf{y}_t) - \hat{\sigma} \text{proj}_{\mathcal{K}}(\hat{\sigma}^{-1}(\mathbf{y}_t - \hat{\sigma} \nabla g_x(\mathbf{y}_t))) \quad \text{where} \quad \hat{\sigma} = 2^j \sigma_{t-1}. \end{aligned} \quad (3.15)$$

The second line follows from the Moreau decomposition, since  $g_v$  is the support function of  $\mathcal{K}$ . We continue to the line-search procedure until the following condition is satisfied:

$$g_x(\hat{\mathbf{y}}) \leq g_x(\mathbf{y}_t) + \langle \nabla g_x(\mathbf{y}_t), \hat{\mathbf{y}} - \mathbf{y}_t \rangle + \frac{1}{2\hat{\sigma}} \|\hat{\mathbf{y}} - \mathbf{y}_t\|^2 + \frac{\varepsilon}{2}. \quad (3.16)$$

This is the usual descent lemma for smooth convex minimization, except the additional slackness term  $\varepsilon/2$ . Once the condition is satisfied, we set  $\mathbf{y}_{t+1} = \hat{\mathbf{y}}$  and  $\sigma_t = \hat{\sigma}$ .

It only remains to update the primal variable. First, update the accumulated sum of the old step-sizes,  $S_t$ . Then, compute the primal step-size  $\eta_t$  as follows:

$$\eta_t = \sigma_t / S_t \quad \text{where} \quad S_t = S_{t-1} + \sigma_t. \quad (3.17)$$

Finally, perform CGM-like update on the primal variable

$$\mathbf{x}_{t+1} = (1 - \eta_t) \mathbf{x}_t + \eta_t \mathbf{h}_t. \quad (3.18)$$

#### 3.2.3 Guarantees

**Theorem 3.1.** *Suppose that we apply UPD for solving the constrained minimization problem (1.1). Assume that strong duality holds. Further assume that the objective function  $f$  has Hölder continuous (sub)gradients of some order  $\nu \in [0, 1]$ , i.e., the assumption (3.9) holds. Then, the sequence  $(\mathbf{x}_t : t = 0, 1, \dots)$  generated by UPD satisfies*

$$f(\mathbf{x}_t) - f_\star \geq -\|\mathbf{y}_\star\| \cdot \text{dist}(\mathcal{A}\mathbf{x}_t, \mathcal{K}) \quad (3.19)$$

$$f(\mathbf{x}_t) - f_\star \leq \frac{\hat{M}_\varepsilon \|\mathbf{y}_0\|^2}{t+1} + \frac{\varepsilon}{2} \quad (3.20)$$

$$\text{dist}(\mathcal{A}\mathbf{x}_t, \mathcal{K}) \leq \frac{4\hat{M}_\varepsilon}{t+1} \|\mathbf{y}_0 - \mathbf{y}_\star\| + \sqrt{\frac{2\varepsilon\hat{M}_\varepsilon}{t+1}} \quad (3.21)$$

where  $\mathbf{y}_\star$  denotes any solution of the dual problem, and

$$\hat{M}_\varepsilon := \left[ \frac{1-\nu}{1+\nu} \frac{1}{\varepsilon} \right]^{\frac{1-\nu}{1+\nu}} M_\nu^{\frac{2}{1+\nu}}, \quad \text{for a given target accuracy } \varepsilon > 0. \quad (3.22)$$

**Corollary 3.1** (Iteration complexity). Instate the assumptions of Theorem 3.1. Then, UPD provably achieves an  $\varepsilon$ -solution of (1.1) after at most  $t_\varepsilon$  iterations, where

$$t_\varepsilon = \mathcal{O} \left( \left( \frac{M_\nu}{\varepsilon} \right)^{\frac{2}{1+\nu}} \right). \quad (3.23)$$

This matches the iteration complexity of the non-accelerated first-order methods, which is optimal for  $\nu = 0$ , but not for  $\nu \in (0, 1]$ . See [NY83].

**Corollary 3.2** (Oracle complexity). Instate the assumptions of Theorem 3.1. Then, the total number of oracle queries  $N_t$  after  $t$  iterations of UPD is bounded as:

$$N_t \leq 2(t+1) + 1 + \log_2(\sigma_{-1}) + \inf_{\nu \in [0,1]} \left\{ \frac{1-\nu}{1+\nu} \log_2 \left( \frac{1-\nu}{1+\nu} \frac{1}{\varepsilon} \right) + \frac{2}{1+\nu} \log_2(M_\nu) \right\}. \quad (3.24)$$

As a consequence,  $N_t \approx 2t$  in the long run (as  $t \rightarrow \infty$ ), *i.e.*, UPD requires approximately 2 oracle queries per iteration on the average.

### 3.2.4 Application to the SDP template

We can use UPD for solving instances from the SDP template (1.4), by setting

$$f(\mathbf{X}) = \langle \mathbf{C}, \mathbf{X} \rangle \quad \text{and} \quad \mathcal{X} = \{\mathbf{X} : \mathbf{X} \in \mathbb{S}_+^n, \text{Tr } \mathbf{X} = \alpha\} \quad \text{and} \quad \mathcal{K} = \{\mathbf{b}\}. \quad (3.25)$$

Then, we can formulate the sharp operator as

$$\mathbf{H}_t = \text{lmo}_{\mathcal{X}}(\mathcal{A}^* \mathbf{y}_t) \in \arg \min_{\mathbf{X} \in \mathcal{X}} \langle \mathbf{C} + \mathcal{A}^* \mathbf{y}_t, \mathbf{X} \rangle. \quad (3.26)$$

This is the same form as an lmo of CGM. Consequently, we can implement it using

$$\mathbf{H}_t = \alpha \mathbf{u}_t \mathbf{u}_t^* \quad \text{where} \quad (\mathbf{u}_t, \lambda_t) = \text{MinEigVec}(\mathbf{C} + \mathcal{A}^* \mathbf{y}_t). \quad (3.27)$$

We can form the partial dual objective and the subgradient as

$$\nabla g_x(\mathbf{y}_t) = -\mathcal{A} \mathbf{H}_t \quad \text{where} \quad g_x(\mathbf{y}_t) = -\alpha \lambda_t. \quad (3.28)$$

Since  $\text{proj}_{\mathcal{K}}(\cdot) = \mathbf{b}$  independently of its input, we can cast the dual update rule (3.15) as

$$\hat{\mathbf{y}} = \mathbf{y}_t + \hat{\sigma}(\mathcal{A} \mathbf{H}_t - \mathbf{b}) \quad \text{where} \quad \hat{\sigma} = 2^j \sigma_{t-1}. \quad (3.29)$$

Algorithm 3.1 shows the complete pseudocode of UPD for this template.

### Chapter 3. Universal Primal-Dual Methods

---

**Remark.** After the first iteration, we can reuse  $(\hat{\mathbf{u}}, \hat{\lambda})$  computed in the line-search procedure of the previous iteration to implement (3.27), in order to avoid unnecessary repetitive calls of the sharp operator. As a consequence, UPD requires  $j_t + 2$  queries of the sharp operator at iterations  $t$ , also counting  $j = -1$  and 0.

---

**Algorithm 3.1** UPD for the standard SDP formulation (1.4)

---

**Require:** Data for (1.4); target accuracy  $\varepsilon > 0$ ; initial dual variable  $\mathbf{y} \in \mathbb{R}^d$

```
1: function UPD
2:    $\mathbf{X} \leftarrow \mathbf{0}$  and  $S \leftarrow 0$ 
3:   for  $t \leftarrow 0, 1, 2, 3, \dots$  do
4:      $(\mathbf{u}, \lambda) \leftarrow \text{MinEigVec}(\mathbf{C} + \mathcal{A}^* \mathbf{y})$ 
5:      $\mathbf{H} \leftarrow \alpha \mathbf{u} \mathbf{u}^*$  and  $g_x \leftarrow -\alpha \lambda$ 
6:     for  $j \leftarrow -1, 0, 1, \dots$  do
7:        $\hat{\sigma} \leftarrow 2^{-j} \sigma$ 
8:        $\hat{\mathbf{y}} = \mathbf{y} + \hat{\sigma}(\mathcal{A}\mathbf{H} - \mathbf{b})$ 
9:        $(\hat{\mathbf{u}}, \hat{\lambda}) \leftarrow \text{MinEigVec}(\mathbf{C} + \mathcal{A}^* \hat{\mathbf{y}})$ 
10:       $\hat{g}_x \leftarrow -\alpha \hat{\lambda}$ 
11:      if  $\hat{g}_x \leq g_x + \langle \nabla g_x(\mathbf{y}_t), \hat{\mathbf{y}} - \mathbf{y}_t \rangle + \frac{1}{2\hat{\sigma}} \|\hat{\mathbf{y}} - \mathbf{y}_t\|^2 + \frac{\varepsilon}{2}$  then break for
12:      end if
13:    end for
14:     $\mathbf{y} \leftarrow \hat{\mathbf{y}}$  and  $\sigma \leftarrow \hat{\sigma}$ 
15:     $S \leftarrow S + \sigma$ 
16:     $\eta \leftarrow \sigma / S$ 
17:     $\mathbf{X} \leftarrow (1 - \eta)\mathbf{X} + \eta\mathbf{H}$ 
18:  end for
19:  return  $\mathbf{X}$ 
20: end function
```

---



### 3.3 Accelerated Universal Primal-Dual Method

We now develop AUPD, an accelerated variant of UPD. For the accelerated scheme, we study an alternative to the fast universal gradient method of Nesterov [Nes15]. We design our accelerated variant based on FISTA scheme [BT09], in order to reduce the number of queries to the  $\text{proj}_{\mathcal{K}}$  oracle per iteration.

#### 3.3.1 AUPD Iteration

Set the desired accuracy  $\varepsilon > 0$ , and an initial step-size  $\sigma_{-1}$ . Initialize the primal and the dual variables, as well as the algorithm parameters:

$$\mathbf{x}_0 \in \mathbb{R}^p, \quad \mathbf{y}_0 = \mathbf{z}_0 \in \mathbb{R}^d, \quad \tau_0 = 1, \quad \text{and} \quad S_{-1} = 0. \quad (3.30)$$

At each iteration  $t = 0, 1, 2, \dots$ , compute a (sub)gradient using the sharp operator:

$$\nabla g_x(\mathbf{y}_t) = -\mathcal{A}\mathbf{h}_t \quad \text{where} \quad \mathbf{h}_t = \text{sharp}_{f, \mathcal{X}}(\mathcal{A}^* \mathbf{y}_t). \quad (3.31)$$

Then, compute the partial dual objective  $g_x(\mathbf{y}_t)$  using (3.14).

AUPD executes a line-search procedure to find an appropriate dual step-size as follows:

For  $j = 0, 1, 2, \dots$ , compute a trial point

$$\begin{aligned} \hat{\mathbf{z}} &= \text{prox}_{\hat{\sigma} g_v}(\mathbf{y}_t - \hat{\sigma} \nabla g_x(\mathbf{y}_t)) \\ &= \mathbf{y}_t - \hat{\sigma} \nabla g_x(\mathbf{y}_t) - \hat{\sigma} \text{proj}_{\mathcal{K}}(\hat{\sigma}^{-1}(\mathbf{y}_t - \hat{\sigma} \nabla g_x(\mathbf{y}_t))) \quad \text{where} \quad \hat{\sigma} = 2^{-j} \sigma_t \end{aligned} \quad (3.32)$$

until the following line-search condition holds:

$$g_x(\hat{\mathbf{z}}) \leq g_x(\mathbf{y}_t) + \langle \nabla g_x(\mathbf{y}_t), \hat{\mathbf{z}} - \mathbf{y}_t \rangle + \frac{1}{2\hat{\sigma}} \|\hat{\mathbf{z}} - \mathbf{y}_t\|^2 + \frac{\varepsilon}{2\tau_t}. \quad (3.33)$$

Once the condition is satisfied, set  $\mathbf{z}_{t+1} = \hat{\mathbf{z}}$  and  $\sigma_t = \hat{\sigma}$ .

Update the accumulated sum  $S_t$  of old step-sizes, then compute the primal step-size  $\eta_t$  as

$$\eta_t = \tau_t \sigma_t / S_t \quad \text{where} \quad S_t = S_{t-1} + \tau_t \sigma_t. \quad (3.34)$$

Update the acceleration parameter  $\tau_t$ , and perform the acceleration step

$$\mathbf{y}_{t+1} = \mathbf{z}_{t+1} + \frac{\tau_t - 1}{\tau_{t+1}}(\mathbf{z}_{t+1} - \mathbf{z}_t) \quad \text{where} \quad \tau_{t+1} = \frac{1}{2}(1 + \sqrt{1 + 4\tau_t^2}). \quad (3.35)$$

Finally, update the primal-variable

$$\mathbf{x}_{t+1} = (1 - \eta_t)\mathbf{x}_t + \eta_t \mathbf{h}_t. \quad (3.36)$$

### Chapter 3. Universal Primal-Dual Methods

**Remark.** We use a decaying  $\frac{\varepsilon}{2t_t}$  slackness term in the line-search condition (3.33) of AUPD in order to prevent error accumulation, in contrast to the constant  $\frac{\varepsilon}{2}$  slackness in UPD.

#### 3.3.2 Guarantees

**Theorem 3.2.** *Instate the assumptions of Theorem 3.1. Then, the primal sequence  $(\mathbf{x}_t : t = 0, 1, \dots)$  generated by AUPD satisfies*

$$f(\mathbf{x}_t) - f_\star \geq -\|\mathbf{y}_\star\| \cdot \text{dist}(\mathcal{A}\mathbf{x}_t, \mathcal{K}) \quad (3.37)$$

$$f(\mathbf{x}_t) - f_\star \leq \frac{4\hat{M}_\varepsilon \|\mathbf{y}_0\|^2}{(t+2)^{\frac{1+3\nu}{1+\nu}}} + \frac{\varepsilon}{2}, \quad (3.38)$$

$$\text{dist}(\mathcal{A}\mathbf{x}_t, \mathcal{K}) \leq \frac{16\hat{M}_\varepsilon}{(t+2)^{\frac{1+3\nu}{1+\nu}}} \|\mathbf{y}_0 - \mathbf{y}_\star\| + \sqrt{\frac{8\varepsilon\hat{M}_\varepsilon}{(t+2)^{\frac{1+3\nu}{1+\nu}}}} \quad (3.39)$$

where  $\mathbf{y}_\star$  denotes any solution to the dual problem, and  $\hat{M}_\varepsilon$  is defined as in (3.22).

**Corollary 3.3** (Iteration complexity). *Instate the assumptions of Theorem 3.1. Then, AUPD provably gets an  $\varepsilon$ -solution of (1.1) after at most  $t_\varepsilon$  iterations, where*

$$t_\varepsilon = \mathcal{O}\left(\left(\frac{M_\nu}{\varepsilon}\right)^{\frac{2}{1+3\nu}}\right). \quad (3.40)$$

This is the optimal iteration complexity for any smoothness order  $\nu \in [0, 1]$ , in the sense of first-order black box models [NY83], with respect to the smoothness of the dual problem.

**Corollary 3.4** (Oracle complexity). *Instate the assumptions of Theorem 3.1. Then, the total number of oracle queries  $N_t$  after  $t$  iterations of AUPD is bounded as:*

$$N_t \leq 2(t+1) + 1 + \log_2(\sigma_{-1}) + \inf_{\nu \in [0,1]} \left\{ \frac{1-\nu}{1+\nu} \log_2\left(\frac{1-\nu}{1+\nu} \frac{t+1}{\varepsilon}\right) + \frac{2}{1+\nu} \log_2(M_\nu) \right\}. \quad (3.41)$$

Roughly speaking, AUPD requires approximately 2 oracle queries per iteration on average on the long-run.

#### 3.3.3 Application to the SDP template

Algorithm 3.2 shows the complete pseudocode of AUPD for the standard SDP formulation (1.4). We omit the details for brevity, which are similar to the details of UPD described in Section 3.2.4.

**Remark.** Unlike UPD, we cannot reuse the output of the sharp operator on the step 5 of Algorithm 3.2, because the input of `MinEigVec` are different in steps 5 and 10. As a result, AUPD requires  $j_t + 2$  queries of the sharp operator at iterations  $t$ . That is, one for the step 6, and  $j_t + 1$  in the line-search procedure also counting  $j = 0$ .

---

**Algorithm 3.2** AUPD for the standard SDP formulation (1.4)
 

---

**Require:** Data for (1.4); target accuracy  $\varepsilon > 0$ ; initial dual variable  $\mathbf{y} \in \mathbb{R}^d$

```

1: function AUPD
2:    $\mathbf{X} \leftarrow \mathbf{0}$  and  $\mathbf{z} \leftarrow \mathbf{y}$ 
3:    $\tau \leftarrow 1$  and  $S \leftarrow 0$ 
4:   for  $t \leftarrow 0, 1, 2, 3, \dots$  do
5:      $(\mathbf{u}, \lambda) \leftarrow \text{MinEigVec}(\mathbf{C} + \mathcal{A}^* \mathbf{y})$ 
6:      $\mathbf{H} \leftarrow \alpha \mathbf{u} \mathbf{u}^*$  and  $g_x \leftarrow -\alpha \lambda$ 
7:     for  $j \leftarrow 0, 1, \dots$  do
8:        $\hat{\sigma} \leftarrow 2^{-j} \sigma$ 
9:        $\hat{\mathbf{z}} = \mathbf{y} + \hat{\sigma}(\mathcal{A} \mathbf{H} - \mathbf{b})$ 
10:       $(\hat{\mathbf{u}}, \hat{\lambda}) \leftarrow \text{MinEigVec}(\mathbf{C} + \mathcal{A}^* \hat{\mathbf{z}})$ 
11:       $\hat{g}_x \leftarrow -\alpha \hat{\lambda}$ 
12:      if  $\hat{g}_x \leq g_x + \langle \nabla g_x(\mathbf{y}_t), \hat{\mathbf{z}} - \mathbf{y}_t \rangle + \frac{1}{2\hat{\sigma}} \|\hat{\mathbf{z}} - \mathbf{y}_t\|^2 + \frac{\varepsilon}{2\tau}$  then break for
13:      end if
14:    end for
15:     $\sigma \leftarrow \hat{\sigma}$ 
16:     $\tau^+ \leftarrow 0.5(1 + \sqrt{1 + 4\tau^2})$ 
17:     $\mathbf{y} \leftarrow \hat{\mathbf{z}} + \frac{\tau-1}{\tau^+} (\hat{\mathbf{z}} - \mathbf{z})$ 
18:     $\mathbf{z} \leftarrow \hat{\mathbf{z}}$ 
19:     $S \leftarrow S + \tau \sigma$ 
20:     $\eta \leftarrow \tau \sigma / S$ 
21:     $\mathbf{X} \leftarrow (1 - \eta) \mathbf{X} + \eta \mathbf{H}$ 
22:     $\tau \leftarrow \tau^+$ 
23:  end for
24:  return  $\mathbf{X}$ 
25: end function
    
```

---

## 3.4 Numerical Experiments

In this section, we illustrate the scalability and flexibility of our primal-dual framework, by providing numerical experiments in quantum tomography and matrix completion problems.

### 3.4.1 Quantum tomography with Pauli operators

We consider the quantum tomography problem which aims to extract information from a physical quantum system. A  $q$ -qubit quantum system is mathematically characterized by its density matrix, which is a complex  $n \times n$  positive semidefinite Hermitian matrix  $\mathbf{X}_q$ , where  $n = 2^q$ .

We can provably deduce the quantum state by performing compressive linear measurements based on the Pauli operators [GLF<sup>+</sup>10]:  $\mathbf{b} = \mathcal{A} \mathbf{X} \in \mathbb{C}^d$ . While the size of the density matrix grows exponentially in  $q$ , a significantly fewer compressive measurements, in the order of

### Chapter 3. Universal Primal-Dual Methods

---

$d = \mathcal{O}(n \log n)$ , suffices to recover a pure state  $q$ -qubit density matrix by solving the following convex optimization problem:

$$\underset{\mathbf{X} \in \mathbb{C}^{n \times n}}{\text{minimize}} \quad \frac{1}{2} \|\mathcal{A}\mathbf{X} - \mathbf{b}\|^2 \quad \text{subject to} \quad \mathbf{X} \in \mathbb{S}_+^n, \quad \text{Tr } \mathbf{X} = 1 \quad (3.42)$$

Since the objective function is smooth, we can also apply CGM to this problem. Since the formulation is tuning-free, this problem provides an ideal scalability test for comparing our framework against CGM. We remain within the noiseless setting, so that we can verify the performance of the algorithms with respect to the ground-truth in large-scale. The behavior of the algorithms are similar under polarization and additive Gaussian noise.

Here are the details of our test setup. We generate a random pure quantum state (a rank-one  $\mathbf{X}_{\text{q}} = \mathbf{x}_{\text{q}} \mathbf{x}_{\text{q}}^*$ ), and we take  $d = 2n \log n$  random Pauli measurements. For  $q = 14$  qubits system, this corresponds to a  $268'435'456$  dimensional problem with  $n = 317'983$  measurements.

To apply our algorithms, we recast (3.42) into (1.1) by introducing a slack variable  $\mathbf{u} = \mathcal{A}\mathbf{X} - \mathbf{b}$ :

$$\underset{\mathbf{u}, \mathbf{X}}{\text{minimize}} \quad \frac{1}{2} \|\mathbf{u}\|^2 \quad \text{subject to} \quad \mathcal{A}\mathbf{X} - \mathbf{u} = \mathbf{b}, \quad \mathbf{X} \in \mathbb{S}_+^n, \quad \text{Tr } \mathbf{X} = 1 \quad (3.43)$$

We set  $\varepsilon = 2 \times 10^{-4}$  for our methods and we set a wall-time of  $2 \times 10^4$  seconds for all methods. We performed the experiments in MATLAB, using a computational resource of 16 CPUs of 2.40 GHz and 512 GB memory space.

Computing the sharp operator requires one query of MinEigVec, that we implement by using the MATLAB's built-in `eigs` function.

Figure 3.1 illustrates iteration and time complexities. UPD, with an average of 1.978 line-search steps per iteration, exhibits a similar empirical performance to CGM with standard step-size. AUPD, with an average of 1.057 line-search steps, exhibits a superior performance.

The line-search variant of CGM improves over the standard approach. We also implement a similar line-search strategy in the primal-step for our methods. In this heuristic, we choose  $\eta_t$  greedily such that it minimizes the objective function. We observe a significant improvement in the practical performance of our methods with this line-search variant. Note however, the improvement due to line-search is typically less significant in more realistic noisy problem setups, both for CGM and our methods.

#### 3.4.2 Matrix completion with MovieLens dataset

To demonstrate the flexibility of our algorithms, we consider the popular matrix completion problem. In matrix completion, we seek to estimate a low-rank matrix  $\mathbf{X} \in \mathbb{R}^{m \times n}$  from its subsampled entries  $\mathbf{b} \in \mathbb{R}^d$ , where  $\mathcal{A} : \mathbb{R}^{m \times n} \rightarrow \mathbb{R}^d$  samples  $d$  entries from the matrix that it is applied to.

### 3.4. Numerical Experiments

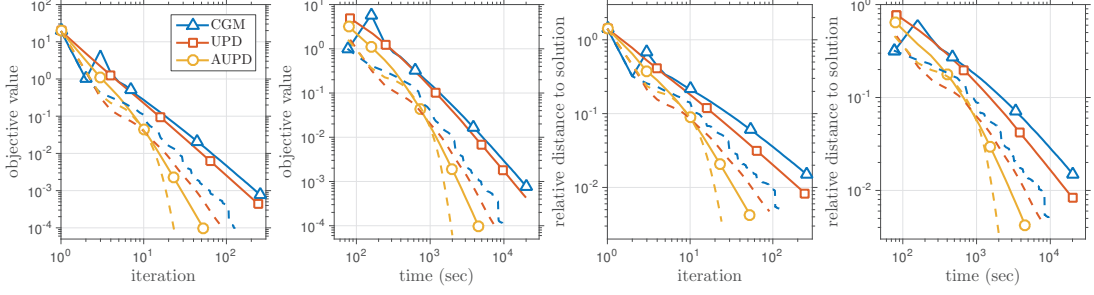


Figure 3.1 – The convergence behavior of algorithms for the  $q = 14$  qubits quantum tomography problem. [Solid lines] correspond to the theoretical averaging scheme. [Dashed lines] correspond to the line-search variants (in the primal step).

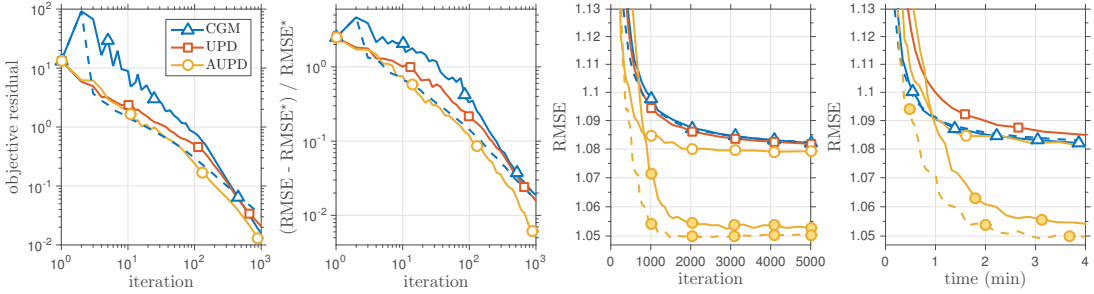


Figure 3.2 – The performance of the algorithms for the matrix completion problems. [Dashed lines] correspond to the line-search variants (in primal steps). [Empty] and the [Filled] markers correspond to the formulation (3.44) and (3.45) respectively.

Convex formulations involving the nuclear norm have been shown to be effective in estimating low-rank matrices from limited number of measurements [CR12]. For instance, we can solve the following least-squares formulation with nuclear norm constraint:

$$\underset{X \in \mathbb{R}^{m \times n}}{\text{minimize}} \quad \frac{1}{2} \|\mathcal{A}X - \mathbf{b}\|^2 \quad \text{subject to} \quad \|X\|_{S_1} \leq \alpha. \quad (3.44)$$

We can solve (3.44) with CGM. In this formulation,  $\alpha$  is a tuning parameter.

Unlike CGM, we can apply our algorithms to the following parameter-free formulation:

$$\underset{X \in \mathbb{R}^{m \times n}}{\text{minimize}} \quad \frac{1}{d} \|X\|_{S_1}^2 \quad \text{subject to} \quad \mathcal{A}X = \mathbf{b} \quad (3.45)$$

While the nonsmooth objective of (3.45) lets us drop the tuning parameter, it clearly burdens the computational efficiency of the convex optimization algorithms.

We solve the matrix completion problem with MovieLens 100K dataset. We use our algorithms for solving (3.44) and (3.45). In comparison, we use CGM for solving (3.44). Recall that CGM does not apply to (3.45).

### Chapter 3. Universal Primal-Dual Methods

---

For this experiment, we do not perform any pre-processing, and we use the default ub test and training data partition. We set the target accuracy  $\varepsilon = 10^{-3}$ , and we choose the tuning parameter  $\alpha = 9975/2$  as in [JS10].

The sharp operator requires the computation of the top singular vectors. We use `lansvd` function from PROPACK<sup>1</sup> to implement the sharp operator.

The first two plots in Figure 3.2 draw the performance of the algorithms for solving (3.44). Our metrics are the normalized objective residual and the root mean squared error (test RMSE). Other two plots in Figure 3.2 compare the performance of the formulations (3.44) and (3.45) which are represented by empty and filled markers respectively. The dashed line for AUPD corresponds to the line-search heuristic, where we use greedy  $\eta_t$  that minimizes the feasibility gap (which is equivalent to the training RMSE).

### 3.5 CGM is Dual Averaging Subgradient Method

This section reviews a simple relation between CGM and the dual averaging subgradient method of Nesterov [Nes09]. We presented this relation in the appendix of [YTDC15b], which is an initial version of [YTDC15a]. For a more detailed study of the duality between the subgradient methods and CGM, we refer to [Bac15].

Consider the CGM template

$$\underset{\mathbf{x}}{\text{minimize}} \quad f(\mathbf{x}) \quad \text{subject to} \quad \mathbf{x} \in \mathcal{X}, \quad (3.46)$$

with the smooth and convex objective function  $f: \mathbb{R}^p \rightarrow \mathbb{R}$ , and convex a compact domain  $\mathcal{X}$ .

Let us reformulate this problem by introducing a slack variable  $\mathbf{u} = \mathbf{x}$ :

$$\underset{\mathbf{u}, \mathbf{x}}{\text{minimize}} \quad f(\mathbf{u}) \quad \text{subject to} \quad \mathbf{u} = \mathbf{x}, \quad \mathbf{x} \in \mathcal{X}. \quad (3.47)$$

The dual averaging subgradient method is perhaps one of the most efficient variants in the class of subgradient methods for solving (3.47). In order to apply this method, we first derive the dual function:

$$\begin{aligned} \phi(\mathbf{y}) &:= \min_{\mathbf{u} \in \mathbb{R}^p} \{f(\mathbf{u}) - \langle \mathbf{u}, \mathbf{y} \rangle\} + \min_{\mathbf{x} \in \mathcal{X}} \langle \mathbf{y}, \mathbf{x} \rangle \\ &= -f^*(\mathbf{y}) + \min_{\mathbf{x} \in \mathcal{X}} \langle \mathbf{y}, \mathbf{x} \rangle \end{aligned} \quad (3.48)$$

where  $f^*$  is the Fenchel conjugate of  $f$ .

---

<sup>1</sup>R.M.Larsen. PROPACK Software for large and sparse SVD calculations. Available from "<http://sun.stanford.edu/~rmunk/PROPACK/>"

### 3.5. CGM is Dual Averaging Subgradient Method

We can characterize the dual averaging subgradient method for solving the formulation (3.48) with the following two-step iterative update scheme:

$$\begin{aligned} \mathbf{x}_{t+1} &= \mathbf{x}_t + \eta_t \nabla \phi(\mathbf{y}_t) \\ \mathbf{y}_{t+1} &= \arg \max_{\mathbf{y}} \{ \langle \mathbf{y}, \mathbf{x}_{t+1} \rangle - \beta_t \psi(\mathbf{y}) \} \end{aligned} \quad (3.49)$$

Here,  $\beta_t$  is a non-increasing sequence of positive numbers, and  $\psi : \mathbb{R}^p \rightarrow \mathbb{R}$  is a strongly convex function (also called as proximity function) to be chosen.

Let us choose  $\beta_t = 1$  for all  $t$ , and  $\psi = f^*$ . Recall that  $f$  is smooth, hence  $f^*$  is strongly convex due to the Baillon-Haddad theorem [BC11].

Then, from the identity of the second step,  $\mathbf{y}_{t+1} = \arg \max_{\mathbf{y}} \{ \langle \mathbf{y}, \mathbf{x}_{t+1} \rangle - \beta_t f^*(\mathbf{y}) \}$ , we have

$$\mathbf{x}_{t+1} \in \partial f^*(\mathbf{y}_{t+1}), \quad \text{or equivalently,} \quad \mathbf{y}_{t+1} = \nabla f(\mathbf{x}_{t+1}). \quad (3.50)$$

This equivalence follows from the standard properties of the Fenchel duality.

$\nabla \phi(\mathbf{y}_t)$  in (3.49) denotes any subgradient of  $\phi$  at  $\mathbf{y}_t$ . We can compute a subgradient by

$$\nabla \phi(\mathbf{y}_t) = \mathbf{h}_t - \mathbf{u}_t, \quad \text{where} \quad \mathbf{h}_t := \text{lmo}_{\mathcal{X}}(\mathbf{y}_t), \quad \forall \mathbf{u}_t \in \partial f^*(\mathbf{y}_t). \quad (3.51)$$

In particular, we can choose  $\mathbf{u}_t := \mathbf{x}_t$ , since  $\mathbf{x}_t \in \partial f^*(\mathbf{y}_t)$ , as shown in (3.50).

Considering these choices, we can formulate an instance of (3.52) as follows:

$$\begin{aligned} \mathbf{x}_{t+1} &= \mathbf{x}_t + \eta_t (\mathbf{h}_t - \mathbf{x}_t) \quad \text{where} \quad \mathbf{h}_t = \text{lmo}_{\mathcal{X}}(\mathbf{y}_{t+1}) \\ \mathbf{y}_{t+1} &= \nabla f(\mathbf{x}_{t+1}) \end{aligned} \quad (3.52)$$

This is nothing but a simple rearrangement of the standard CGM for solving (3.46).

### 3.6 Appendix: Proofs

The following Lemma that we borrow from [Nes15] provides the key properties for constructing universal gradient algorithms. We refer to [Nes15] for the proof of this lemma.

**Lemma 3.3** (Lemma 2 of [Nes15]). Let  $g_x$  be Hölder smooth of order  $\nu$  for some  $\nu \in [0, 1]$ . Then, for any  $\delta > 0$  and

$$M \geq \hat{M}_\delta := \left[ \frac{1-\nu}{1+\nu} \frac{1}{\delta} \right]^{\frac{1-\nu}{1+\nu}} M_\nu^{\frac{2}{1+\nu}}, \quad (3.53)$$

the following statement holds:

$$g_x(\mathbf{z}) \leq \underbrace{g_x(\mathbf{y}) + \langle \nabla g_x(\mathbf{y}), \mathbf{z} - \mathbf{y} \rangle}_{Q_M(\mathbf{z}, \mathbf{y})} + \frac{M}{2} \|\mathbf{z} - \mathbf{y}\|^2 + \frac{\delta}{2}, \quad \forall \mathbf{y}, \mathbf{z} \in \mathbb{R}^d. \quad (3.54)$$

This lemma introduces an approximate quadratic upper bound for  $g_x$  with the additional slackness term. The bound depends on the choice of the slackness parameter  $\delta$  and the smoothness parameter  $\nu$ . If we know in advance that  $\nu = 1$ , then we can set  $M$  to the smoothness constant, and we can drop the slackness term  $\delta/2$ .

UPD and AUPD are based on the proximal (sub)gradient step (3.15), which comes with the following estimation guarantee.

**Lemma 3.4.** Suppose that  $\mathbf{y}^+$ , which is defined by

$$\mathbf{y}^+ := \text{prox}_{\sigma g_\nu}(\mathbf{y} - \sigma \nabla g_x(\mathbf{y})) \quad (3.55)$$

for some  $\sigma > 0$  satisfies

$$g_x(\mathbf{y}^+) \leq g_x(\mathbf{y}) + \langle \nabla g_x(\mathbf{y}), \mathbf{y}^+ - \mathbf{y} \rangle + \frac{1}{2\sigma} \|\mathbf{y}^+ - \mathbf{y}\|^2 + \frac{\delta}{2} \quad \text{for some } \delta \in \mathbb{R}. \quad (3.56)$$

Then, the following inequality holds  $\forall \mathbf{z} \in \mathbb{R}^d$ :

$$g(\mathbf{y}^+) \leq g_x(\mathbf{y}) + \langle \nabla g_x(\mathbf{y}), \mathbf{z} - \mathbf{y} \rangle + g_\nu(\mathbf{z}) + \frac{1}{2\sigma} (\|\mathbf{z} - \mathbf{y}\|^2 - \|\mathbf{z} - \mathbf{y}^+\|^2) + \frac{\delta}{2}. \quad (3.57)$$

*Proof.* Start with the optimality condition of (3.55)

$$0 \in \nabla g_x(\mathbf{y}) + \frac{1}{\sigma}(\mathbf{y}^+ - \mathbf{y}) + \partial g_\nu(\mathbf{y}^+) \implies \mathbf{y} - \mathbf{y}^+ = \sigma(\nabla g_x(\mathbf{y}) + \nabla g_\nu(\mathbf{y}^+)), \quad (3.58)$$

where we denote a subgradient of  $g_\nu$  at  $\mathbf{y}^+$  by  $\nabla g_\nu(\mathbf{y}^+) \in \partial g_\nu(\mathbf{y}^+)$ .



The rest of the proof is as follows:

$$\begin{aligned} & \frac{1}{2\sigma} (\|\mathbf{z} - \mathbf{y}^+\|^2 - \|\mathbf{z} - \mathbf{y}\|^2) \\ &= \frac{1}{\sigma} \langle \mathbf{y} - \mathbf{y}^+, \mathbf{z} - \mathbf{y}^+ \rangle - \frac{1}{2\sigma} \|\mathbf{y} - \mathbf{y}^+\|^2 \end{aligned} \quad (3.59)$$

$$= \langle \nabla g_x(\mathbf{y}) + \nabla g_v(\mathbf{y}^+), \mathbf{z} - \mathbf{y}^+ \rangle - \frac{1}{2\sigma} \|\mathbf{y} - \mathbf{y}^+\|^2 \quad (3.60)$$

$$= -\langle \nabla g_x(\mathbf{y}), \mathbf{y}^+ - \mathbf{y} \rangle + \frac{1}{2\sigma} \|\mathbf{y}^+ - \mathbf{y}\|^2 + \langle \nabla g_v(\mathbf{y}^+), \mathbf{z} - \mathbf{y}^+ \rangle + \langle \nabla g_x(\mathbf{y}), \mathbf{z} - \mathbf{y} \rangle \quad (3.61)$$

$$\leq g_x(\mathbf{y}) - g_x(\mathbf{y}^+) + \frac{\delta}{2} + \langle \nabla g_v(\mathbf{y}^+), \mathbf{z} - \mathbf{y}^+ \rangle + \langle \nabla g_x(\mathbf{y}), \mathbf{z} - \mathbf{y} \rangle \quad (3.62)$$

$$\leq g_x(\mathbf{y}) - g_x(\mathbf{y}^+) + \frac{\delta}{2} + g_v(\mathbf{z}) - g_v(\mathbf{y}^+) + \langle \nabla g_x(\mathbf{y}), \mathbf{z} - \mathbf{y} \rangle \quad (3.63)$$

$$= g_x(\mathbf{y}) + \langle \nabla g_x(\mathbf{y}), \mathbf{z} - \mathbf{y} \rangle + g_v(\mathbf{z}) + \frac{\delta}{2} - g(\mathbf{y}^+) \quad (3.64)$$

(3.60) follows from (3.58), (3.62) uses the line-search condition (3.56), and (3.63) is based on the convexity of  $g_v$ .  $\square$

Clearly, (3.56) holds if  $\sigma \leq 1/\hat{M}_\delta$  due to Lemma 3.3, where  $\hat{M}_\delta$  is defined by (3.53). If  $v$  and  $M_v$  are known, we can directly set  $\sigma = 1/\hat{M}_\delta$ , as it provably satisfies the line-search condition (3.56). In general, however, we do not know  $v$  and  $M_v$  a priori. In this case, we determine the step-size  $\sigma$  by using a line-search procedure.

The following lemma guarantees that the line-search terminates after a finite number of trials.

**Lemma 3.5.** The line-search procedure in UPD terminates after at most

$$j_t = \left\lfloor \log_2(\hat{M}_\varepsilon \sigma_{t-1}) \right\rfloor + 1 \text{ iterations.} \quad (3.65)$$

Similarly, the line-search procedure in AUPD terminates after at most

$$j_t = \left\lfloor \log_2\left(\frac{t+1}{\varepsilon}\right) + \log_2\left(M_v^{\frac{2}{1+v}} \sigma_{t-1}\right) \right\rfloor + 1 \text{ iterations.} \quad (3.66)$$

*Proof.* First recall that  $M_v$  is finite due to the Hölder smoothness assumption.

In UPD, we choose the slackness parameter  $\delta = \varepsilon$ , where  $\varepsilon$  is the target accuracy. For a given  $\varepsilon > 0$ ,  $\hat{M}_\varepsilon$  is also finite. Moreover, our line-search condition is provably satisfied when  $\sigma \leq 1/\hat{M}_\varepsilon$ . Since  $\hat{\sigma} = 2^{-j} \sigma_{t-1}$ , the line-search procedure is terminated after at most  $j = \lfloor \log_2(\hat{M}_\varepsilon \sigma_{t-1}) \rfloor + 1$  trials.

In AUPD, we choose the slackness parameter  $\delta = \varepsilon/\tau_t$ . Remark that

$$\tau_{t+1} = \frac{1}{2} \left( 1 + \sqrt{1 + 4\tau_t^2} \right) \leq \frac{1}{2} (1 + (1 + 2\tau_t)) = \tau_t + 1, \quad \text{and} \quad \tau_0 = 1. \quad (3.67)$$

### Chapter 3. Universal Primal-Dual Methods

Hence, by induction, we can show  $\tau_t \leq t + 1$ . Using this bound, we get

$$\hat{M}_{\varepsilon/\tau_t} = \left[ \frac{1-\nu}{1+\nu} \frac{\tau_t}{\varepsilon} \right]^{\frac{1-\nu}{1+\nu}} M_v^{\frac{2}{1+\nu}} \leq \left[ \frac{\tau_t}{\varepsilon} \right]^{\frac{1-\nu}{1+\nu}} M_v^{\frac{2}{1+\nu}} \leq \left[ \frac{t+1}{\varepsilon} \right]^{\frac{1-\nu}{1+\nu}} M_v^{\frac{2}{1+\nu}}. \quad (3.68)$$

The line-search condition is provably satisfied if  $\sigma \leq 1/\hat{M}_{\varepsilon/\tau_t}$ . As a result, our line-search procedure is terminated after at most  $j = \lfloor \log_2(\frac{t+1}{\varepsilon}) + \log_2(M_v^{\frac{2}{1+\nu}} \sigma_{t-1}) \rfloor + 1$  trials. We can show this simply by combining (3.68) with the fact that  $\hat{\sigma} = 2^{-j} \sigma_{t-1}$ .  $\square$

### Convergence analysis of UPD

We first provide the convergence guarantee of the dual function in Theorem 3.6. Then, we prove the convergence rate and the iteration complexity guarantees of UPD.

**Theorem 3.6.** *Let  $(\mathbf{y}_t : t = 0, 1, 2, \dots)$  be a sequence generated by UPD. Then,*

$$g(\bar{\mathbf{y}}_t) - g(\mathbf{y}) \leq \bar{g}_t - g(\mathbf{y}) \leq \frac{\hat{M}_\varepsilon}{t+1} \|\mathbf{y}_0 - \mathbf{y}\|^2 + \frac{\varepsilon}{2} \quad \forall \mathbf{y} \in \mathbb{R}^d, \quad (3.69)$$

where the two averaging sequences  $(\bar{\mathbf{y}}_t : t = 0, 1, 2, \dots)$  and  $(\bar{g}_t : t = 0, 1, 2, \dots)$  are defined as

$$\bar{\mathbf{y}}_t := \frac{1}{S_t} \sum_{i=0}^t \sigma_i \mathbf{y}_{i+1} \quad \text{and} \quad \bar{g}_t := \frac{1}{S_t} \sum_{i=0}^t \sigma_i g(\mathbf{y}_{i+1}). \quad (3.70)$$

In particular, by choosing  $\mathbf{y} = \mathbf{y}_\star$  in (3.69), we get

$$g(\bar{\mathbf{y}}_t) - g(\mathbf{y}_\star) \leq \frac{\hat{M}_\varepsilon}{t+1} \|\mathbf{y}_0 - \mathbf{y}_\star\|^2 + \frac{\varepsilon}{2}. \quad (3.71)$$

*Proof.* The following inequality follows directly from Lemma 3.4:

$$g(\mathbf{y}_{i+1}) \leq g(\mathbf{y}) + \frac{\varepsilon}{2} + \frac{1}{2\sigma_i} [\|\mathbf{y} - \mathbf{y}_i\|^2 - \|\mathbf{y} - \mathbf{y}_{i+1}\|^2], \quad \forall \mathbf{y} \in \mathbb{R}^d. \quad (3.72)$$

Taking the weighted sum of this inequality over  $i$ , we get

$$\bar{g}_t \leq g(\mathbf{y}) + \frac{\varepsilon}{2} + \frac{1}{2S_t} [\|\mathbf{y} - \mathbf{y}_0\|^2 - \|\mathbf{y} - \mathbf{y}_{t+1}\|^2], \quad \forall \mathbf{y} \in \mathbb{R}^d. \quad (3.73)$$

Lemma 3.3 proves that the line-search condition (3.56) is satisfied if  $\sigma \leq 1/\hat{M}_\varepsilon$ , where  $\hat{M}_\varepsilon$  defined as in (3.22). Considering the algorithmic design of our line-search procedure, this ensures that  $\sigma_i \geq 1/(2\hat{M}_\varepsilon)$ . As a consequence, we have

$$S_t = \sum_{i=0}^t \sigma_i \geq \sum_{i=0}^t \frac{1}{2\hat{M}_\varepsilon} = \frac{t+1}{2\hat{M}_\varepsilon}. \quad (3.74)$$

Substituting (3.74) into (3.73), we obtain (3.69). Since  $g$  is convex,  $g(\bar{\mathbf{y}}_t) \leq \bar{g}_t$ .  $\square$

**Proof of Theorem 3.1**

We use the following identities to relate the convergence in the dual problem to the primal:

$$\begin{aligned} g_x(\mathbf{y}_i) &= -f(\mathbf{h}_i) - \langle \mathbf{y}_i, \mathcal{A}\mathbf{h}_i \rangle \\ \nabla g_x(\mathbf{y}_i) &= -\mathbf{h}_i \\ g(\mathbf{y}_i) &\geq g_\star := g(\mathbf{y}_\star) = -f_\star, \quad \text{for } i = 0, 1, 2, \dots \end{aligned} \quad (3.75)$$

Substituting these expressions into Lemma 3.4, we get the following estimate

$$f(\mathbf{h}_i) - f_\star \leq -\langle \mathcal{A}\mathbf{h}_i, \mathbf{y} \rangle + g_\nu(\mathbf{y}) + \frac{\varepsilon}{2} + \frac{1}{2\sigma_i} (\|\mathbf{y} - \mathbf{y}_i\|^2 - \|\mathbf{y} - \mathbf{y}_{i+1}\|^2), \quad \forall \mathbf{y} \in \mathbb{R}^d. \quad (3.76)$$

Taking the weighted sum of this inequality over  $i$ , and using the convexity of  $f$ , we get

$$f(\mathbf{x}_t) - f_\star \leq -\langle \mathcal{A}\mathbf{x}_t, \mathbf{y} \rangle + g_\nu(\mathbf{y}) + \frac{\varepsilon}{2} + \frac{1}{2S_t} (\|\mathbf{y} - \mathbf{y}_0\|^2 - \|\mathbf{y} - \mathbf{y}_{k+1}\|^2), \quad \forall \mathbf{y} \in \mathbb{R}^d. \quad (3.77)$$

Choosing  $\mathbf{y} = \mathbf{0}$  (note that  $g_\nu(\mathbf{0}) = 0$ ), we get (3.20)

$$f(\mathbf{x}_t) - f_\star \leq \frac{\varepsilon}{2} + \frac{\|\mathbf{y}_0\|^2}{2S_t} \leq \frac{\varepsilon}{2} + \frac{\hat{M}_\varepsilon \|\mathbf{y}_0\|^2}{t+1}. \quad (3.78)$$

(3.19) follows from the Lagrange saddle point formulation:

$$f_\star \leq \mathcal{L}(\mathbf{x}, \mathbf{r}, \mathbf{y}_\star) = f(\mathbf{x}) + \langle \mathbf{y}_\star, \mathcal{A}\mathbf{x} - \mathbf{r} \rangle \leq f(\mathbf{x}) + \|\mathbf{y}_\star\| \|\mathcal{A}\mathbf{x} - \mathbf{r}\|, \quad \forall \mathbf{r} \in \mathcal{K} \text{ and } \forall \mathbf{x} \in \mathcal{X}. \quad (3.79)$$

The last inequality follows from the Cauchy-Schwarz inequality. We get (3.19) by setting  $\mathbf{x} = \mathbf{x}_t$ .

Next, we prove (3.38). Start from the following saddle point formulation:

$$f_\star \leq \mathcal{L}(\mathbf{x}, \mathbf{r}, \mathbf{y}_\star) = f(\mathbf{x}) + \langle \mathbf{y}_\star, \mathcal{A}\mathbf{x} - \mathbf{b} - \mathbf{r} \rangle, \quad \forall \mathbf{r} \in \mathcal{K} \text{ and } \forall \mathbf{x} \in \mathcal{X}. \quad (3.80)$$

Set  $\mathbf{x} = \mathbf{x}_t$  in this inequality, and substitute it into (3.77):

$$\min_{\mathbf{r} \in \mathcal{K}} \left\{ \langle \mathcal{A}\mathbf{x}_t - \mathbf{r}, \mathbf{y} - \mathbf{y}_\star \rangle - \frac{1}{2S_t} [\|\mathbf{y} - \mathbf{y}_0\|^2 - \|\mathbf{y} - \mathbf{y}_{t+1}\|^2] \right\} \leq \frac{\varepsilon}{2}, \quad \forall \mathbf{y} \in \mathbb{R}^d \quad (3.81)$$

$$\implies \min_{\mathbf{r} \in \mathcal{K}} \left\{ \langle \mathcal{A}\mathbf{x}_t - \mathbf{r}, \mathbf{y} - \mathbf{y}_\star \rangle - \frac{1}{2S_t} \|\mathbf{y} - \mathbf{y}_0\|^2 \right\} \leq \frac{\varepsilon}{2}, \quad \forall \mathbf{y} \in \mathbb{R}^d \quad (3.82)$$

$$\implies \max_{\mathbf{y} \in \mathbb{R}^d} \min_{\mathbf{r} \in \mathcal{K}} \left\{ \langle \mathcal{A}\mathbf{x}_t - \mathbf{r}, \mathbf{y} - \mathbf{y}_\star \rangle - \frac{1}{2S_t} \|\mathbf{y} - \mathbf{y}_0\|^2 \right\} \leq \frac{\varepsilon}{2} \quad (3.83)$$

$$\iff \min_{\mathbf{r} \in \mathcal{K}} \max_{\mathbf{y} \in \mathbb{R}^d} \left\{ \langle \mathcal{A}\mathbf{x}_t - \mathbf{r}, \mathbf{y} - \mathbf{y}_\star \rangle - \frac{1}{2S_t} \|\mathbf{y} - \mathbf{y}_0\|^2 \right\} \leq \frac{\varepsilon}{2} \quad (3.84)$$

$$\iff \min_{\mathbf{r} \in \mathcal{K}} \left\{ \langle \mathcal{A}\mathbf{x}_t - \mathbf{r}, \mathbf{y}_0 - \mathbf{y}_\star \rangle + \frac{S_t}{2} \|\mathcal{A}\mathbf{x}_t - \mathbf{r}\|^2 \right\} \leq \frac{\varepsilon}{2}, \quad (3.85)$$

(3.81) follows the definition of  $g_\nu$ , (3.82) removes a positive term on the left-hand-side, (3.83)

### Chapter 3. Universal Primal-Dual Methods

---

maximizes over the free variable  $\mathbf{y}$ , (3.84) holds due to the Sion's minimax theorem [Sio58], finally (3.85) analytically solves the inner maximization subproblem.

Denote by

$$\bar{\mathbf{r}} \in \arg \min_{\mathbf{r} \in \mathcal{K}} \left\{ \langle \mathcal{A}\mathbf{x}_t - \mathbf{r}, \mathbf{y}_0 - \mathbf{y}_\star \rangle + \frac{S_t}{2} \|\mathcal{A}\mathbf{x}_t - \mathbf{r}\|^2 \right\} \quad (3.86)$$

Then, we can rewrite (3.85) as

$$\langle \mathcal{A}\mathbf{x}_t - \bar{\mathbf{r}}, \mathbf{y}_0 - \mathbf{y}_\star \rangle + \frac{S_t}{2} \|\mathcal{A}\mathbf{x}_t - \bar{\mathbf{r}}\|^2 \leq \frac{\varepsilon}{2}. \quad (3.87)$$

Using Cauchy-Schwarz inequality, this implies

$$-\|\mathcal{A}\mathbf{x}_t - \bar{\mathbf{r}}\| \cdot \|\mathbf{y}_0 - \mathbf{y}_\star\| + \frac{S_t}{2} \|\mathcal{A}\mathbf{x}_t - \bar{\mathbf{r}}\|^2 \leq \frac{\varepsilon}{2}. \quad (3.88)$$

Solving this second order inequality for  $\|\mathcal{A}\mathbf{x}_t - \bar{\mathbf{r}}\|$ , we get

$$\begin{aligned} \|\mathcal{A}\mathbf{x}_t - \bar{\mathbf{r}}\| &\leq \frac{1}{S_t} \left( \|\mathbf{y}_0 - \mathbf{y}_\star\| + \sqrt{\|\mathbf{y}_0 - \mathbf{y}_\star\|^2 + \varepsilon S_t} \right) \\ &\leq \frac{1}{S_t} \left( 2\|\mathbf{y}_0 - \mathbf{y}_\star\| + \sqrt{\varepsilon S_t} \right). \end{aligned} \quad (3.89)$$

By definition of the distance function,

$$\text{dist}(\mathcal{A}\mathbf{x}_t, \mathcal{K}) \leq \|\mathcal{A}\mathbf{x}_t - \bar{\mathbf{r}}\|, \quad \forall \mathbf{r} \in \mathcal{K}. \quad (3.90)$$

Finally by substituting  $S_t \geq \frac{t+1}{2\hat{M}_\varepsilon}$ , we get (3.38).

#### Proof of Corollary 3.1

To provably get an  $\varepsilon$ -solution, we need to guarantee both  $\text{dist}(\mathcal{A}\mathbf{x}_t, \mathcal{K}) \leq \varepsilon$  and  $f(\mathbf{x}_t) - f_\star \leq \varepsilon$ .

If we combine these conditions with the bounds from Theorem 3.1, we obtain the following conditions on  $t_\varepsilon$ :

$$\frac{4\hat{M}_\varepsilon}{t_\varepsilon + 1} \|\mathbf{y}_0 - \mathbf{y}_\star\| + \sqrt{\frac{2\varepsilon\hat{M}_\varepsilon}{t_\varepsilon + 1}} \leq \varepsilon, \quad \text{and} \quad \frac{\hat{M}_\varepsilon \|\mathbf{y}_0\|^2}{t_\varepsilon + 1} + \frac{\varepsilon}{2} \leq \varepsilon. \quad (3.91)$$

For the objective residual, we can get the bound on  $t_\varepsilon$  by a simple rearrangement

$$t_\varepsilon + 1 \geq \frac{2\hat{M}_\varepsilon \|\mathbf{y}_0\|^2}{\varepsilon}. \quad (3.92)$$

For the feasibility gap, we need to solve a second order inequality with respect to  $\frac{1}{\sqrt{t_\varepsilon+1}}$

$$\frac{1}{\sqrt{t_\varepsilon+1}} \leq \frac{-\sqrt{2\varepsilon\hat{M}_\varepsilon} + \sqrt{2\varepsilon\hat{M}_\varepsilon + 16\varepsilon\hat{M}_\varepsilon\|\mathbf{y}_0 - \mathbf{y}_\star\|}}{8\hat{M}_\varepsilon\|\mathbf{y}_0 - \mathbf{y}_\star\|}. \quad (3.93)$$

Rearranging and simplifying, we get

$$t_\varepsilon + 1 \geq \frac{(8\hat{M}_\varepsilon\|\mathbf{y}_0 - \mathbf{y}_\star\|)^2}{4\varepsilon\hat{M}_\varepsilon + 16\varepsilon\hat{M}_\varepsilon\|\mathbf{y}_0 - \mathbf{y}_\star\| - \sqrt{16\varepsilon^2\hat{M}_\varepsilon^2 + 128\varepsilon^2\hat{M}_\varepsilon\|\mathbf{y}_0 - \mathbf{y}_\star\|}} \quad (3.94)$$

$$= \frac{(8\|\mathbf{y}_0 - \mathbf{y}_\star\|)^2}{4 + 16\|\mathbf{y}_0 - \mathbf{y}_\star\| - \sqrt{16 + 128\|\mathbf{y}_0 - \mathbf{y}_\star\|}} \frac{\hat{M}_\varepsilon}{\varepsilon} \quad (3.95)$$

Both (3.92) and (3.95) imply that  $t_\varepsilon = \mathcal{O}\left(\frac{\hat{M}_\varepsilon}{\varepsilon}\right)$ . Using the definition (3.22) of  $\hat{M}_\varepsilon$ , we can get the following order representation in terms of  $\varepsilon$ :

$$t_\varepsilon = \mathcal{O}\left(\inf_{v \in [0,1]} \left(\frac{1-v}{1+v}\right)^{\frac{1-v}{1+v}} \left(\frac{M_v}{\varepsilon}\right)^{\frac{2}{1+v}}\right) = \mathcal{O}\left(\inf_{v \in [0,1]} \left(\frac{M_v}{\varepsilon}\right)^{\frac{2}{1+v}}\right). \quad (3.96)$$

### Proof of Corollary 3.2

Similar to the analysis in [Nes15], we estimate the total number of oracle quires of UPD.

The total number of oracle quires up to the iteration  $t$  is given by the formula

$$N_t = \sum_{i=0}^t (j_i + 2). \quad (3.97)$$

We designed the line-search procedure such that  $\sigma_i = 2^{j_i} \sigma_{i-1}$ . Consequently,

$$\begin{aligned} N_t &= \sum_{i=0}^t (\log_2(\sigma_i / \sigma_{i-1}) + 2) \\ &= \sum_{i=0}^t (\log_2(\sigma_i) - \log_2(\sigma_{i-1}) + 2) \\ &= 2(t+1) + \log_2(\sigma_{-1}) - \log_2(\sigma_t). \end{aligned} \quad (3.98)$$

Recall that  $\sigma_t \geq 1/(2\hat{M}_\varepsilon)$ , hence we get

$$\begin{aligned} N_t &\leq 2(t+1) + 1 + \log_2(\sigma_{-1}) + \log_2(\hat{M}_\varepsilon) \\ &= 2(t+1) + 1 + \log_2(\sigma_{-1}) + \inf_{v \in [0,1]} \left\{ \frac{1-v}{1+v} \log_2\left(\frac{1-v}{1+v} \frac{1}{\varepsilon}\right) + \frac{2}{1+v} \log_2(M_v) \right\}. \end{aligned} \quad (3.99)$$

The second line simply follows from the definition (3.22) of  $\hat{M}_\varepsilon$ .

### Convergence analysis of AUPD

Consider the following dual update scheme of  $\mathbf{z}_{t+1}$  and  $\mathbf{w}_{t+1}$  from  $\mathbf{y}_t$  and  $\mathbf{w}_t$ :

$$\begin{cases} \mathbf{y}_t &= (1 - \tau_t)\mathbf{z}_t + \bar{\tau}_t\mathbf{w}_t \\ \mathbf{z}_{t+1} &= \text{prox}_{\sigma_t g_v}(\mathbf{y}_t - \sigma_t \nabla g_x(\mathbf{y}_t)) \\ \mathbf{w}_{t+1} &= \mathbf{w}_t - \frac{1}{\bar{\tau}_t}(\mathbf{y}_t - \mathbf{z}_{t+1}), \end{cases} \quad (3.100)$$

where  $\mathbf{w}_0 = \mathbf{z}_0$ ,  $\bar{\tau}_0 = 1$  and

$$\bar{\tau}_t^2 = \bar{\tau}_{t-1}^2(1 - \bar{\tau}_t). \quad (3.101)$$

The parameter  $\sigma_t$  is determined based on the following line-search condition:

$$g(\mathbf{z}_{t+1}) \leq g(\mathbf{y}_t) + \langle \nabla g(\mathbf{y}_t), \mathbf{z}_{t+1} - \mathbf{y}_t \rangle + \frac{1}{2\sigma_t} \|\mathbf{z}_{t+1} - \mathbf{y}_t\|^2 + \frac{\varepsilon \bar{\tau}_t}{2}, \quad (3.102)$$

with  $\sigma_t \leq \sigma_{t-1}$  for  $t \geq 0$ .

Next, we show that (3.100) is equivalent to AUPD update scheme in the dual.

**Lemma 3.7.** The scheme (3.100) can be restated as follows:

$$\begin{cases} \mathbf{z}_{t+1} &= \text{prox}_{\sigma_t g_v}(\mathbf{y}_t - \sigma_t \nabla g_x(\mathbf{y}_t)) \\ \tau_{t+1} &= \frac{1}{2}(1 + \sqrt{1 + 4\bar{\tau}_t^2}) \\ \mathbf{y}_{t+1} &= \mathbf{z}_{t+1} + \frac{\tau_t - 1}{\tau_{t+1}}(\mathbf{z}_{t+1} - \mathbf{z}_t), \end{cases} \quad (3.103)$$

where  $\mathbf{y}_0 = \mathbf{z}_0$  and  $\tau_0 = 1$ , and  $\sigma_t$  is determined based on the line-search.

This dual scheme is in the form of FISTA [BT09], except for the line-search procedure.

*Proof.* Set  $\tau_t = \bar{\tau}_t^{-1}$ , then  $t_0 = \tau_0^{-1} = 1$ . From (3.100), we have

$$\mathbf{w}_t - \mathbf{w}_{t+1} = \frac{1}{\bar{\tau}_t}(\mathbf{y}_t - \mathbf{z}_{t+1}) = \tau_t(\mathbf{y}_t - \mathbf{z}_{t+1}). \quad (3.104)$$

We also have  $\mathbf{y}_t = (1 - \bar{\tau}_t)\mathbf{z}_t + \bar{\tau}_t\mathbf{w}_t$ , which leads to

$$\mathbf{w}_t = \frac{1}{\bar{\tau}_t}(\mathbf{y}_t - (1 - \bar{\tau}_t)\mathbf{z}_t) = \tau_t(\mathbf{y}_t - (1 - \tau_t^{-1})\mathbf{z}_t). \quad (3.105)$$

Combining these expressions, we get

$$\begin{aligned} \tau_t(\mathbf{y}_t - \mathbf{z}_{t+1}) &= \mathbf{w}_t - \mathbf{w}_{t+1} \\ &= \tau_t(\mathbf{y}_t - (1 - \tau_t^{-1})\mathbf{z}_t) - \tau_{t+1}(\mathbf{y}_{t+1} - (1 - \tau_{t+1}^{-1})\mathbf{z}_{t+1}). \end{aligned} \quad (3.106)$$

Then, we simplify as follows:

$$\begin{aligned}\tau_{t+1}\mathbf{y}_{t+1} &= \tau_t\mathbf{z}_{t+1} + \tau_{t+1}(1 - \tau_{t+1}^{-1})\mathbf{z}_{t+1} - \tau_t(1 - \tau_t^{-1})\mathbf{z}_t \\ &= (\tau_t + \tau_{t+1} - 1)\mathbf{z}_{t+1} - (\tau_t - 1)\mathbf{z}_t.\end{aligned}\quad (3.107)$$

Hence  $\mathbf{y}_{t+1} = \mathbf{z}_{t+1} + \frac{\tau_t - 1}{\tau_{t+1}}(\mathbf{z}_{t+1} - \mathbf{z}_t)$ , which is the third step of (3.103).

Next, from the condition (3.101), we have  $\tau_{t+1}^2 - \tau_{t+1} - \tau_t^2 = 0$ . Hence,  $\tau_{t+1} = \frac{1}{2}(1 + \sqrt{1 + 4\tau_t^2})$ , which is exactly the second step of (3.103).  $\square$

### Proof of Theorem 3.2

Start from Lemma 3.4:

$$\begin{aligned}g(\mathbf{z}_{t+1}) &\leq (g_x(\mathbf{y}_t) + \langle \nabla g_x(\mathbf{y}_t), \mathbf{y} - \mathbf{y}_t \rangle + g_v(\mathbf{z})) + \frac{\varepsilon}{2\tau_t} \\ &\quad + \frac{1}{\sigma_t} \langle \mathbf{y}_t - \mathbf{z}_{t+1}, \mathbf{y}_t - \mathbf{y} \rangle - \frac{1}{2\sigma_t} \|\mathbf{z}_{t+1} - \mathbf{y}_t\|^2, \quad \forall \mathbf{y} \in \mathbb{R}^d\end{aligned}\quad (3.108)$$

$$\leq g(\mathbf{z}) + \frac{\varepsilon}{2\tau_t} + \frac{1}{\sigma_t} \langle \mathbf{y}_t - \mathbf{z}_{t+1}, \mathbf{y}_t - \mathbf{y} \rangle - \frac{1}{2\sigma_t} \|\mathbf{z}_{t+1} - \mathbf{y}_t\|^2, \quad \forall \mathbf{y} \in \mathbb{R}^d. \quad (3.109)$$

Next, subtract  $g_\star$  from (3.108) to get

$$\begin{aligned}g(\mathbf{z}_{t+1}) - g_\star &\leq (g_x(\mathbf{y}_t) + \langle \nabla g_x(\mathbf{y}_t), \mathbf{y} - \mathbf{y}_t \rangle + g_v(\mathbf{y}) - g_\star) + \frac{\varepsilon}{2\tau_t} \\ &\quad + \frac{1}{\sigma_t} \langle \mathbf{y}_t - \mathbf{z}_{t+1}, \mathbf{y}_t - \mathbf{y} \rangle - \frac{1}{2\sigma_t} \|\mathbf{z}_{t+1} - \mathbf{y}_t\|^2, \quad \forall \mathbf{y} \in \mathbb{R}^d.\end{aligned}\quad (3.110)$$

Set  $\mathbf{y} = \mathbf{z}_t$  in (3.109), and subtract  $g_\star$  from the both sides:

$$g(\mathbf{z}_{t+1}) - g_\star \leq g(\mathbf{z}_t) - g_\star + \frac{\varepsilon}{2\tau_t} - \frac{1}{2\sigma_t} \|\mathbf{z}_{t+1} - \mathbf{y}_t\|^2 + \frac{1}{\sigma_t} \langle \mathbf{y}_t - \mathbf{z}_{t+1}, \mathbf{y}_t - \mathbf{z}_t \rangle. \quad (3.111)$$

Multiply (3.110) by  $\tau_t^{-1}$ , (3.111) by  $(1 - \tau_t^{-1})$ , and sum them up

$$\begin{aligned}g(\mathbf{z}_{t+1}) - g_\star &\leq (1 - \tau_t^{-1})(g(\mathbf{z}_t) - g_\star) + \frac{1}{2\sigma_t\tau_t^2} (\|\mathbf{w}_t - \mathbf{y}\|^2 - \|\mathbf{w}_{t+1} - \mathbf{y}\|^2) + \frac{\varepsilon}{2\tau_t} \\ &\quad + \frac{1}{\tau_t} (g_x(\mathbf{y}_t) + \langle \nabla g_x(\mathbf{y}_t), \mathbf{y} - \mathbf{y}_t \rangle + g_v(\mathbf{y}) - g_\star), \quad \forall \mathbf{y} \in \mathbb{R}^d.\end{aligned}\quad (3.112)$$

Multiply (3.112) by  $\sigma_t\tau_t^2$ , and sum over  $t$  to get

$$\begin{aligned}\sum_{i=0}^t \sigma_i \tau_i^2 (g(\mathbf{z}_{i+1}) - g_\star) &\leq \sum_{i=0}^t \left[ \sigma_i \tau_i^2 (1 - \tau_i^{-1})(g(\mathbf{z}_i) - g_\star) + \frac{1}{2} (\|\mathbf{w}_i - \mathbf{y}\|^2 - \|\mathbf{w}_{i+1} - \mathbf{y}\|^2) \right. \\ &\quad \left. + \frac{\sigma_i \tau_i \varepsilon}{2} + \sigma_i \tau_i (g(\mathbf{y}_i) + \langle \nabla g_x(\mathbf{y}_i), \mathbf{y} - \mathbf{y}_i \rangle + g_v(\mathbf{y}) - g_\star) \right], \quad \forall \mathbf{y} \in \mathbb{R}^d.\end{aligned}\quad (3.113)$$

### Chapter 3. Universal Primal-Dual Methods

---

By design of the update rules and the line-search procedure (since  $\sigma_t \leq \sigma_{t-1}$ ), we have

$$\tau_0 = 1 \quad \text{and} \quad \sigma_t \tau_t^2 (1 - \tau_t^{-1}) \leq \sigma_{t-1} \tau_{t-1}^2 \quad \text{for } t = 1, 2, \dots \quad (3.114)$$

Substitute (3.114) into (3.113) to get

$$\begin{aligned} \sum_{i=0}^t \sigma_i \tau_i^2 (g(\mathbf{z}_{i+1}) - g_\star) &\leq \frac{1}{2} (\|\mathbf{w}_0 - \mathbf{y}\|^2 - \|\mathbf{w}_{t+1} - \mathbf{y}\|^2) + \sum_{i=1}^t [\sigma_{i-1} \tau_{i-1}^2 (g(\mathbf{z}_i) - g_\star)] \\ &\quad + \sum_{i=0}^t \left[ \sigma_i \tau_i \left( g_x(\mathbf{y}_i) + \langle \nabla g_x(\mathbf{y}_i), \mathbf{y} - \mathbf{y}_i \rangle + g_v(\mathbf{y}) - g_\star + \frac{\varepsilon}{2} \right) \right], \quad \forall \mathbf{y} \in \mathbb{R}^d. \end{aligned} \quad (3.115)$$

This implies

$$\begin{aligned} \frac{\sigma_t \tau_t^2}{S_t} (g(\mathbf{z}_{t+1}) - g_\star) &\leq \frac{1}{S_t} \sum_{i=0}^t \left[ \sigma_i \tau_i \left( g_x(\mathbf{y}_i) + \langle \nabla g_x(\mathbf{y}_i), \mathbf{y} - \mathbf{y}_i \rangle + g_v(\mathbf{y}) - g_\star \right) \right] \\ &\quad + \frac{1}{2S_t} [\|\mathbf{w}_0 - \mathbf{y}\|^2 - \|\mathbf{w}_{t+1} - \mathbf{y}\|^2] + \frac{\varepsilon}{2}, \quad \forall \mathbf{y} \in \mathbb{R}^d. \end{aligned} \quad (3.116)$$

Note that  $g(\mathbf{z}_{t+1}) - g_\star \geq 0$ , hence by rearranging we get

$$-\frac{1}{S_t} \sum_{i=0}^t \sigma_i \tau_i \left( g_x(\mathbf{y}_i) + \langle \nabla g_x(\mathbf{y}_i), \mathbf{y} - \mathbf{y}_i \rangle + g_v(\mathbf{y}) - g_\star \right) \leq \frac{1}{2S_t} (\|\mathbf{w}_0 - \mathbf{y}\|^2 - \|\mathbf{w}_{t+1} - \mathbf{y}\|^2) + \frac{\varepsilon}{2}. \quad (3.117)$$

Now, we use the following identities to map this estimate for the primal variable

$$\begin{aligned} g(\mathbf{y}_i) &= -f(\mathbf{h}_i) - \langle \mathbf{y}_i, \mathcal{A}\mathbf{h}_i \rangle \\ \nabla g_x(\mathbf{y}_i) &= -\mathcal{A}\mathbf{h}_i, \\ g_\star &= -f_\star. \end{aligned} \quad (3.118)$$

Then, using Jensen's inequality and convexity of  $f$  we get

$$f(\mathbf{x}_t) - f_\star \leq -\langle \mathcal{A}\mathbf{x}_t, \mathbf{y} \rangle + g_v(\mathbf{y}) + \frac{\varepsilon}{2} + \frac{1}{2S_t} (\|\mathbf{w}_0 - \mathbf{y}\|^2 - \|\mathbf{w}_{t+1} - \mathbf{y}\|^2), \quad \forall \mathbf{y} \in \mathbb{R}^d \quad (3.119)$$

$$\leq \frac{\varepsilon}{2} + \frac{\|\mathbf{y}_0\|^2}{2\hat{S}_k}, \quad (3.120)$$

where (3.120) follows by setting  $\mathbf{y} = \mathbf{0}$  in (3.119).

We can reformulate the update rule of  $\tau_t$  in (3.103) as

$$\tau_t = \tau_t^2 - \tau_{t-1}^2, \quad \text{for } t = 1, 2, \dots \quad (3.121)$$

Also note that

$$\frac{t+2}{2} \leq \tau_t < t+2, \quad \text{for } t = 0, 1, 2, \dots \quad (3.122)$$



Using (3.121), (3.122), and the fact that

$$\sigma_0 \geq \sigma_i \geq \sigma_t \geq \frac{1}{2\hat{M}_\varepsilon/\tau_t} = \frac{1}{2\tau_t^{\frac{1-v}{1+v}}\hat{M}_\varepsilon} \geq \frac{1}{2(t+2)^{\frac{1-v}{1+v}}\hat{M}_\varepsilon}, \quad (3.123)$$

one can show

$$\begin{aligned} S_t = \sum_{i=0}^t \sigma_i \tau_i &\geq \sum_{i=0}^t \frac{\tau_i}{2\hat{M}_\varepsilon/\tau_i} = \frac{1}{2\hat{M}_\varepsilon/\tau_t} \left[ 1 + \sum_{i=1}^t (\tau_i^2 - \tau_{i-1}^2) \right] \\ &\geq \frac{\tau_t^2}{2(t+2)^{\frac{1-v}{1+v}}\hat{M}_\varepsilon} \geq \frac{(t+2)^{\frac{1+3v}{1+v}}}{8\hat{M}_\varepsilon}. \end{aligned} \quad (3.124)$$

We get (3.38) by substituting (3.124) into (3.120).

Inequality (3.37) follows the saddle point formulation (3.79) by setting  $\mathbf{x} = \mathbf{x}_t$ .

Finally, the proof for (3.39) follows similarly to the proof in Theorem 3.1. One can show

$$\text{dist}(\mathcal{A}\mathbf{x}_t, \mathcal{K}) \leq \frac{2\|\mathbf{y}_0 - \mathbf{y}_\star\|}{S_t} + \sqrt{\frac{\varepsilon}{S_t}}, \quad (3.125)$$

and we complete the proof by substituting (3.124) into this estimate.

### Proof of Corollary 3.3

Using the bounds from Theorem 3.2, we obtain the following conditions on  $t_\varepsilon$ :

$$\frac{16\hat{M}_\varepsilon}{(t_\varepsilon + 2)^{\frac{1+3v}{1+v}}} \|\mathbf{y}_0 - \mathbf{y}_\star\| + \sqrt{\frac{8\varepsilon\hat{M}_\varepsilon}{(t_\varepsilon + 2)^{\frac{1+3v}{1+v}}}} \leq \varepsilon, \quad \text{and} \quad \frac{4\hat{M}_\varepsilon\|\mathbf{y}_0\|^2}{(t_\varepsilon + 2)^{\frac{1+3v}{1+v}}} + \frac{\varepsilon}{2} \leq \varepsilon. \quad (3.126)$$

For the objective residual, we can get the bound on  $t_\varepsilon$  by a simple rearrangement

$$t_\varepsilon + 2 \geq \left( \frac{8\hat{M}_\varepsilon\|\mathbf{y}_0\|^2}{\varepsilon} \right)^{\frac{1+v}{1+3v}}. \quad (3.127)$$

For the feasibility gap, we solve a second order inequality with respect to  $\frac{1}{\sqrt{(t_\varepsilon + 2)^{\frac{1+3v}{1+v}}}}$

$$\frac{1}{\sqrt{(t_\varepsilon + 2)^{\frac{1+3v}{1+v}}}} \leq \frac{-\sqrt{8\varepsilon\hat{M}_\varepsilon} + \sqrt{8\varepsilon\hat{M}_\varepsilon + 64\varepsilon\hat{M}_\varepsilon\|\mathbf{y}_0 - \mathbf{y}_\star\|}}{32\hat{M}_\varepsilon\|\mathbf{y}_0 - \mathbf{y}_\star\|}. \quad (3.128)$$

Rearranging and simplifying, we get

$$t_\varepsilon + 2 \geq \left( \frac{(2^5\|\mathbf{y}_0 - \mathbf{y}_\star\|)^2}{2^4 + 2^6\|\mathbf{y}_0 - \mathbf{y}_\star\| - \sqrt{2^8 + 2^{11}\|\mathbf{y}_0 - \mathbf{y}_\star\|}} \frac{\hat{M}_\varepsilon}{\varepsilon} \right)^{\frac{1+v}{1+3v}} \quad (3.129)$$

### Chapter 3. Universal Primal-Dual Methods

---

Using the definition (3.22) of  $\hat{M}_\varepsilon$ , together with the bounds (3.127) and (3.129), we can show

$$t_\varepsilon = \mathcal{O} \left( \inf_{v \in [0,1]} \left( \frac{1-v}{1+v} \right)^{\frac{1-v}{1+3v}} \left( \frac{M_v}{\varepsilon} \right)^{\frac{2}{1+3v}} \right) = \mathcal{O} \left( \inf_{v \in [0,1]} \left( \frac{M_v}{\varepsilon} \right)^{\frac{2}{1+3v}} \right). \quad (3.130)$$

#### Proof of Corollary 3.4

Similar to the proof of Corollary 3.2, the total number of oracle quires up to the iteration  $t$  is

$$N_t = \sum_{i=0}^t (j_i + 2) = 2(t+1) + \log_2(\sigma_{-1}) - \log_2(\sigma_t). \quad (3.131)$$

Recall that  $\sigma_t \geq 1/(2\hat{M}_{\varepsilon/\tau_t})$ . Using the definition of  $\hat{M}_{\varepsilon/\tau_t}$ , and (3.122), we get

$$\begin{aligned} N_t &\leq 2(t+1) + 1 + \log_2(\sigma_{-1}) + \inf_{v \in [0,1]} \left\{ \frac{1-v}{1+v} \log_2 \left( \frac{1-v}{1+v} \frac{\tau_t}{\varepsilon} \right) + \frac{2}{1+v} \log_2(M_v) \right\} \\ &\leq 2(t+1) + 1 + \log_2(\sigma_{-1}) + \inf_{v \in [0,1]} \left\{ \frac{1-v}{1+v} \log_2 \left( \frac{1-v}{1+v} \frac{t+2}{\varepsilon} \right) + \frac{2}{1+v} \log_2(M_v) \right\}. \end{aligned} \quad (3.132)$$

## 4 CGM for Composite Problems

CGM features significantly reduced computational costs (in comparison with the projected gradient methods) in key machine learning applications. As we discussed in Chapter 2, it also exhibits structured updates on the decision variable which are amenable for further storage benefits via sketching. Unfortunately, standard CGM cannot handle affine inclusion constraint  $\mathcal{A}\mathbf{x} \in \mathcal{K}$  efficiently, which restricts its applicability to the key machine learning applications as well as the SDP formulations.

To this end, we describe a generalized CGM approach for a composite convex minimization template with broad applications. In particular, the proposed method applies to our model problems (1.1) and (1.4).

This chapter is based on the joint work with Francesco Locatello, Olivier Fercoq and Volkan Cevher [YFLC18].

### Introduction

In this chapter, we consider the following composite convex minimization formulation:

$$\underset{\mathbf{x} \in \mathcal{X}}{\text{minimize}} \quad F(\mathbf{x}) := f(\mathbf{x}) + g(\mathcal{A}\mathbf{x}), \quad (4.1)$$

Again,  $f : \mathcal{X} \rightarrow \mathbb{R}$  is a smooth proper closed and convex function, and  $g : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{+\infty\}$  is a proper closed convex function which is possibly non-smooth.  $\mathcal{A} : \mathbb{R}^p \rightarrow \mathbb{R}^d$  is a given linear map. This template covers affine inclusion constraint  $\mathcal{A}\mathbf{x} \in \mathcal{K}$  as a special case, since we can choose  $g$  as the indicator function of  $\mathcal{K}$ .

By using the powerful operator-splitting framework, many problems that belong to the template (4.1) can be solved efficiently. These methods make use of the gradient of the smooth part  $f$ , along with the prox-operator of the non-smooth part  $g$  and the projection onto the domain  $\mathcal{X}$ .

## Chapter 4. CGM for Composite Problems

---

CGM-type algorithms (also called as projection-free methods) are desirable when the projection onto the problem domain is a computational bottleneck. A classical example is the positive semi-definite cone, for which the projection requires computing a full eigendecomposition.

Unfortunately, the standard CGM cannot handle the non-smooth term in (4.1). When the non-smooth part is an indicator function, one could take the linear minimization oracle (lmo) with respect to the intersection between  $\mathcal{X}$  and the set represented by  $g(\mathcal{A}\mathbf{x})$ . Unfortunately, even the lmo itself can be a difficult optimization problem in this case.

To this end, we propose a novel CGM framework for solving the composite problem (4.1) with rigorous convergence guarantees. Our approach retains the simplicity of projection-free methods, but it allows us to disentangle the complexity of the lmo from the non-smooth component  $g(\mathcal{A}\mathbf{x})$ .

Our method combines the ideas of smoothing [Nes05] and homotopy under the CGM framework. Lan [Lan14] proposed a similar approach for non-smooth problems, which is also extended for the conditional gradient sliding framework in their follow-up papers [LZ16] and [LPZZ17]. Their analysis, however, is restricted with a Lipschitz continuity assumption on the non-smooth part. As a consequence, their methods and analyses do not apply to the problems with affine constraints. See Sections 4.3.5 and 4.3.6.

### Contributions

Our main contributions in this chapter are as follows:

- We introduce a simple, easy to implement CGM framework for solving the instances from composite minimization template (4.1).
- When the non-smooth part is a Lipschitz-continuous regularizer, we prove that the proposed method, HCGM, achieves the optimal  $\mathcal{O}(1/\sqrt{t})$  convergence rate. When the non-smooth term is an indicator function, *i.e.*, when the problem has affine constraints, we prove that the proposed method gets  $\mathcal{O}(1/\sqrt{t})$  convergence rate, both in the objective residual and the feasibility gap.
- We analyze the convergence guarantees of the proposed method for the case where the linear minimization oracle is noisy, with additive and/or multiplicative errors. We show that our approach is robust.
- We present key instances of our framework, including the non-smooth minimization, minimization with affine inclusion constraints, and convex splitting. We provide a summary of the related work for each case.
- We present empirical evidence supporting our findings.

## 4.1 Preliminaries

Our method is based on the idea of combining smoothing and homotopy under the CGM framework. Hence, we first review the smoothing approach of Nesterov [Nes05], which plays a crucial role in our algorithmic design.

### 4.1.1 Nesterov Smoothing

Nesterov smoothing exploits an important class of non-smooth functions  $\psi(\mathbf{x})$  that can be written in the following *max-form*:

$$\psi(\mathbf{x}) = \max_{\mathbf{u} \in \mathcal{U}} \{ \langle \mathcal{A}\mathbf{x}, \mathbf{u} \rangle - \hat{\phi}(\mathbf{u}) \}, \quad (4.2)$$

for some convex and compact set  $\mathcal{U} \subset \mathbb{R}^d$  and a convex function  $\hat{\phi} : \mathcal{U} \rightarrow \mathbb{R}$ .

Let us consider a prox-function on  $\mathcal{U}$  (i.e., a strongly convex continuous function), and denote it by  $\delta(\mathbf{u})$ . In this work, we use the Euclidean prox-function  $\delta(\mathbf{u}) = \frac{1}{2} \|\mathbf{u}\|^2$ .

Smooth approximation  $\psi_\beta(\mathbf{x})$  with smoothness parameter  $\beta > 0$  is defined as

$$\psi_\beta(\mathbf{x}) = \max_{\mathbf{u} \in \mathcal{U}} \{ \langle \mathcal{A}\mathbf{x}, \mathbf{u} \rangle - \hat{\phi}(\mathbf{u}) - \beta \delta(\mathbf{u}) \}. \quad (4.3)$$

Then,  $\psi_\beta$  is well defined, differentiable, convex and smooth. Moreover, it uniformly approximates  $\psi$ , in the sense that it satisfies the following envelop property:

$$\psi_\beta(\mathbf{x}) \leq \psi(\mathbf{x}) \leq \psi_\beta(\mathbf{x}) + \beta D_{\mathcal{U}} \quad (\forall \mathbf{x} \in \mathcal{X}), \quad (4.4)$$

where  $D_{\mathcal{U}} = \max_{\mathbf{u} \in \mathcal{U}} \delta(\mathbf{u})$ . See Theorem 1 in [Nes05] for the proof and more details.

If we assume  $g(\mathcal{A} \cdot)$  is Lipschitz continuous, then we can write it in the max form by choosing  $\psi(\mathbf{x}) = g(\mathcal{A}\mathbf{x})$  and  $\hat{\phi}(\mathbf{u}) = g^*(\mathbf{u})$ . Here,  $g^*$  is the Fenchel conjugate of  $g$ , defined as

$$g^*(\mathbf{u}) = \max_{\mathbf{z}} \{ \langle \mathbf{u}, \mathbf{z} \rangle - g(\mathbf{z}) \}. \quad (4.5)$$

Since  $g$  is convex and lower semicontinuous, Fenchel duality holds, and we have

$$g(\mathcal{A}\mathbf{x}) = g^{**}(\mathcal{A}\mathbf{x}). \quad (4.6)$$

Moreover, Lipschitz continuity assumption on  $g$  imposes the boundedness of the dual domain. See Lemma 5 in [DFTJ16] for a proof of this well-known result in the convex analysis.

### 4.1.2 Quadratic Penalty

Quadratic penalty approach is an effective proxy for handling the affine constraint  $\mathcal{A}\mathbf{x} \in \mathcal{K}$ . It works by replacing the constraint with a penalty function which favors the feasibility of the solution. We consider the squared Euclidean distance as the penalty function:  $\frac{1}{2\beta} \text{dist}^2(\mathcal{A}\mathbf{x}, \mathcal{K})$ , where  $\beta > 0$  is the so-called penalty parameter. Surprisingly, quadratic penalty approach is structurally equivalent to a *de facto* instance of the Nesterov smoothing.

Let us start by writing the Fenchel conjugate of the indicator function  $\iota_{\mathcal{K}}(\cdot)$ ,

$$\iota_{\mathcal{K}}^*(\mathbf{z}) = \max_{\mathbf{v} \in \mathcal{K}} \langle \mathbf{v}, \mathbf{z} \rangle, \quad \text{where} \quad \iota_{\mathcal{K}}(\mathbf{y}) = \begin{cases} 0 & \text{if } \mathbf{y} \in \mathcal{K} \\ +\infty & \text{otherwise.} \end{cases} \quad (4.7)$$

Then, we can write the affine constraint in the max form by choosing  $\hat{\phi}(\mathbf{z}) = \iota_{\mathcal{K}}^*(\mathbf{z})$ , and using the following relation:

$$\begin{aligned} \iota_{\mathcal{K}}(\mathcal{A}\mathbf{x}) &= \max_{\mathbf{z}} \left\{ \min_{\mathbf{v} \in \mathcal{K}} \langle \mathcal{A}\mathbf{x} - \mathbf{v}, \mathbf{z} \rangle \right\} \\ &= \max_{\mathbf{z}} \left\{ \langle \mathcal{A}\mathbf{x}, \mathbf{z} \rangle - \iota_{\mathcal{K}}^*(\mathbf{z}) \right\}. \end{aligned} \quad (4.8)$$

Choosing the standard Euclidean prox-function, we get the following “smooth approximation”:

$$\iota_{\mathcal{K}\beta}(\mathcal{A}\mathbf{x}) = \max_{\mathbf{z}} \left\{ \min_{\mathbf{v} \in \mathcal{K}} \langle \mathcal{A}\mathbf{x} - \mathbf{v}, \mathbf{z} \rangle - \frac{\beta}{2} \|\mathbf{z}\|^2 \right\} \quad (4.9)$$

$$= \min_{\mathbf{v} \in \mathcal{K}} \max_{\mathbf{z}} \left\{ \langle \mathcal{A}\mathbf{x} - \mathbf{v}, \mathbf{z} \rangle - \frac{\beta}{2} \|\mathbf{z}\|^2 \right\} \quad (4.10)$$

$$= \frac{1}{2\beta} \text{dist}^2(\mathcal{A}\mathbf{x}, \mathcal{K}), \quad (4.11)$$

where the inversion of min and max holds due to the Sion’s minimax theorem [Sio58].

In summary, we can obtain the quadratic penalty form by applying Nesterov smoothing procedure to the indicator of an affine constraint. However, quadratic penalty does not serve as a uniform approximation, because the dual domain is unbounded and the envelop property (4.4) does not hold. Consequently, the common analysis techniques for smoothing does not apply for quadratic penalty methods.

Nevertheless, we exploit this structural similarity to design algorithms that work for both cases; composite problems with smoothing-friendly non-smooth regularizers and problems with affine constraints.

## 4.2 Homotopy CGM

Consider the following auxiliary optimization problem:

$$\underset{\mathbf{x} \in \mathcal{X}}{\text{minimize}} \quad F_\beta(\mathbf{x}) := f(\mathbf{x}) + g_\beta(\mathcal{A}\mathbf{x}) \quad (4.12)$$

where  $g_\beta$  is the smooth approximation of  $g$ :

$$g_\beta(\mathcal{A}\mathbf{x}) = \max_{\mathbf{u} \in \mathbb{R}^d} \left\{ \langle \mathcal{A}\mathbf{x}, \mathbf{u} \rangle - g^*(\mathbf{u}) - \frac{\beta}{2} \|\mathbf{u}\|^2 \right\}. \quad (4.13)$$

Note that  $F_\beta$  is a smooth and convex function.

Solutions of (4.12) and (4.1) do not coincide for any particular  $\beta$ . Nevertheless, as  $\beta \rightarrow 0$ , any sequence of solutions of (4.12) converges to a solution of (4.1).

Based on these observations, we design HCGM as follows: At iteration  $t$ , we take a conditional gradient step with respect to the smooth approximation  $F_{\beta_t}$ . Then, we choose the next penalty parameter  $\beta_{t+1}$  slightly smaller than  $\beta_t$  to guide the iterates toward the solution of the original problem (4.1).

### 4.2.1 HCGM Iteration

Let  $\beta_0 > 0$  be an initial smoothing parameter. Begin with an arbitrary choice  $\mathbf{x}_0 \in \mathcal{X}$  for the decision variable. At each iteration  $t = 0, 1, 2, \dots$ , we compute the gradient  $\nabla F_{\beta_t}(\mathbf{x}_t)$ ;

$$\nabla F_{\beta_t}(\mathbf{x}_t) = \nabla f(\mathbf{x}_t) + \nabla g_{\beta_t}(\mathcal{A}\mathbf{x}_t), \quad \text{where} \quad \beta_t = \frac{\beta_0}{\sqrt{t+2}}. \quad (4.14)$$

The argument of the maximization subproblem (4.13) can be written as  $\text{prox}_{\beta^{-1}g^*}(\beta^{-1}\mathcal{A}\mathbf{x})$ . Hence, we can compute the gradient of  $g_{\beta_t}$  by using

$$\begin{aligned} \nabla g_{\beta_t}(\mathcal{A}\mathbf{x}_t) &= \mathcal{A}^* \text{prox}_{\beta^{-1}g^*}(\beta^{-1}\mathcal{A}\mathbf{x}_t) \\ &= \beta^{-1} \mathcal{A}^* \left( \mathcal{A}\mathbf{x}_t - \text{prox}_{\beta g}(\mathcal{A}\mathbf{x}_t) \right), \end{aligned} \quad (4.15)$$

where the second line follows from the well-known Moreau decomposition.

Next, we evaluate the linear minimization oracle

$$\mathbf{h}_t = \text{lmo}_{\mathcal{X}}(\nabla F_{\beta_t}(\mathbf{x}_t)). \quad (4.16)$$

Finally, we update the matrix variable as follows:

$$\mathbf{x}_{t+1} = (1 - \eta_t)\mathbf{x}_t + \eta_t \mathbf{h}_t, \quad \text{where} \quad \eta_t = 2/(t+2). \quad (4.17)$$

### 4.2.2 Stopping Criterion

We can design a stopping criterion for HCGM. If  $g : \mathbb{R}^d \rightarrow \mathbb{R}$  is a Lipschitz continuous function, then we can use the following condition to stop our algorithm:

$$\langle \nabla F_{\beta_t}(\mathbf{x}_t), \mathbf{x}_t - \mathbf{h}_t \rangle \leq \varepsilon, \quad \text{which ensures} \quad F(\mathbf{x}_t) - F_\star \leq \varepsilon. \quad (4.18)$$

If  $g$  is an indicator function, then we can use the following criterion to guarantee an  $\varepsilon$ -solution in the sense (1.15):

$$\langle \nabla F_{\beta_t}(\mathbf{x}_t), \mathbf{x}_t - \mathbf{h}_t \rangle \leq \varepsilon \quad \text{and} \quad \|\mathcal{A}\mathbf{x}_t - \text{proj}_{\mathcal{K}}(\mathcal{A}\mathbf{x}_t)\| \leq \varepsilon. \quad (4.19)$$

We can also choose different target accuracies for the objective residual and the feasibility gap.

### 4.2.3 Guarantees

**Theorem 4.1.** *Sequence  $(\mathbf{x}_t : t = 1, 2, \dots)$  generated by HCGM satisfies*

$$F_{\beta_{t-1}}(\mathbf{x}_t) - f_\star \leq 2D_{\mathcal{X}} \left( \frac{L_f}{t+1} + \frac{\|\mathcal{A}\|^2}{\beta_0 \sqrt{t+1}} \right). \quad (4.20)$$

Theorem 4.1 does not directly guarantee the convergence of  $\mathbf{x}_t$  to a solution of (1.1), since the bound is on smoothed gap  $F_{\beta_{t-1}}(\mathbf{x}_t) - f_\star$ . To relate  $F_{\beta_{t-1}}$  back to  $F$ , one usually assumes Lipschitz continuity. This well known perspective leads us to Theorem 4.2, which is a direct extension of Theorem 4 in [Lan14] for composite functions.

**Theorem 4.2.** *Assume that  $g : \mathbb{R}^d \rightarrow \mathbb{R}$  is  $L_g$ -Lipschitz continuous. Then, the sequence  $(\mathbf{x}_t : t = 1, 2, \dots)$  generated by HCGM satisfies the following convergence bound:*

$$F(\mathbf{x}_t) - F_\star \leq 2D_{\mathcal{X}}^2 \left( \frac{L_f}{t+1} + \frac{\|\mathcal{A}\|^2}{\beta_0 \sqrt{t+1}} \right) + \frac{\beta_0 L_g^2}{2\sqrt{t+1}}. \quad (4.21)$$

Furthermore, if the constants  $D_{\mathcal{X}}$ ,  $\|\mathcal{A}\|$  and  $L_g$  are known or easy to approximate, we can choose  $\beta_0 = 2D_{\mathcal{X}}\|\mathcal{A}\|/L_g$  to get the following convergence rate:

$$F(\mathbf{x}_t) - F_\star \leq \frac{2D_{\mathcal{X}}^2 L_f}{t+1} + \frac{2D_{\mathcal{X}}\|\mathcal{A}\|L_g}{\sqrt{t+1}}. \quad (4.22)$$

Lipschitz continuity assumption in Theorem 4.2 leaves many important applications out (see Sections 4.3.5 and 4.3.6). In Theorem 4.3, we take a step further and characterize the convergence when the non-smooth part is an indicator function.



**Theorem 4.3.** *Sequence  $(\mathbf{x}_t : t = 1, 2, \dots)$  generated by HCGM satisfies:*

$$f(\mathbf{x}_t) - f_\star \geq -\|\mathbf{y}_\star\| \cdot \text{dist}(\mathcal{A}\mathbf{x}_t, \mathcal{K}) \quad (4.23)$$

$$f(\mathbf{x}_t) - f_\star \leq 2D_{\mathcal{X}} \left( \frac{L_f}{t+1} + \frac{\|\mathcal{A}\|^2}{\beta_0 \sqrt{t+1}} \right) \quad (4.24)$$

$$\text{dist}(\mathcal{A}\mathbf{x}_t, \mathcal{K}) \leq \frac{2}{\sqrt{t+1}} \left( \beta_0 \|\mathbf{y}_\star\| + D_{\mathcal{X}} \sqrt{\frac{\beta_0 L_f}{\sqrt{t+1}} + \|\mathcal{A}\|^2} \right) \quad (4.25)$$

**Remark.**  $\mathbf{x}_t$  is not guaranteed to be a feasible point, since the condition  $\mathcal{A}\mathbf{x}_t \in \mathcal{K}$  is not guaranteed. Nevertheless, it converges towards the feasible set, and the sequence  $(\mathbf{x}_t : t = 0, 1, 2, \dots)$  is asymptotically feasible.

**Remark.** Similar to the classical CGM, we can consider a line-search variant of HCGM (which replaces the step size by  $\eta_t = \min_{\eta \in [0,1]} F_{\beta_t}((1-\eta)\mathbf{x}_t + \eta\mathbf{h}_t)$ ), and a fully corrective variant (which replaces the last step by  $\mathbf{x}_{t+1} = \arg\min_{\mathbf{x} \in \text{conv}(\mathbf{h}_1, \dots, \mathbf{h}_t)} F_{\beta_t}(\mathbf{x})$ ). All results that we present in this chapter remain valid for these two design variants.

#### 4.2.4 Convergence with Inexact Oracles

Finding an exact solution of lmo can be expensive in practice. On the other hand, approximate solutions can be much more efficient. For instance, exact computation of MinEigVec that we use in the SDP examples is costly, while we can use Lanczos method for an approximate computation.

Different notions of inexact lmo are already explored for CGM and greedy optimization frameworks, see [LJSP13, LKT17, LTR17]. We revisit the notion of additive and multiplicative errors which we adapt here for our setting.

##### Inexact Oracle with Additive Error

At iteration  $t$ , for the given gradient direction  $\nabla F_{\beta_t}(\mathbf{x}_t)$ , we assume that the approximate lmo returns an element  $\tilde{\mathbf{h}}_t \in \mathcal{X}$  such that:

$$\langle \nabla F_{\beta_t}(\mathbf{x}_t), \tilde{\mathbf{h}}_t \rangle \leq \langle \nabla F_{\beta_t}(\mathbf{x}_t), \mathbf{h}_t \rangle + \delta \frac{\eta_t}{2} D_{\mathcal{X}}^2 \left( L_f + \frac{\|\mathcal{A}\|^2}{\beta_t} \right) \quad \text{for some } \delta > 0, \quad (4.26)$$

where  $\mathbf{h}_t$  denotes the exact solution of the lmo. As for the classical CGM, we require the accuracy of lmo to increase [Jag13] as the algorithm progresses.

Replacing the exact lmo with the approximate oracles of the form (4.26) in HCGM, we get the following convergence guarantees:

## Chapter 4. CGM for Composite Problems

**Theorem 4.4.** Sequence  $(\mathbf{x}_t : t = 1, 2, \dots)$  generated by HCGM with approximate lmo (4.26) satisfies:

$$F_{\beta_{t-1}}(\mathbf{x}_t) - f_\star \leq 2D_{\mathcal{X}}^2 \left( \frac{L_f}{t+1} + \frac{\|\mathcal{A}\|^2}{\beta_0 \sqrt{t+1}} \right) (1 + \delta). \quad (4.27)$$

**Theorem 4.5.** Assume that  $g$  is  $L_g$ -Lipschitz continuous. Then, the sequence  $\mathbf{x}_t$  generated by HCGM with approximate lmo (4.26) satisfies:

$$F(\mathbf{x}_t) - F_\star \leq 2D_{\mathcal{X}}^2 \left( \frac{L_f}{t+1} + \frac{\|\mathcal{A}\|^2}{\beta_0 \sqrt{t+1}} \right) (1 + \delta) + \frac{\beta_0 L_g^2}{2\sqrt{t+1}}. \quad (4.28)$$

**Theorem 4.6.** Sequence  $(\mathbf{x}_t : t = 1, 2, \dots)$  generated by HCGM with approximate lmo (4.26) satisfies:

$$f(\mathbf{x}_t) - f_\star \geq -\|\mathbf{y}_\star\| \cdot \text{dist}(\mathcal{A}\mathbf{x}_t, \mathcal{K}) \quad (4.29)$$

$$f(\mathbf{x}_t) - f_\star \leq 2D_{\mathcal{X}}^2 \left( \frac{L_f}{t+1} + \frac{\|\mathcal{A}\|^2}{\beta_0 \sqrt{t+1}} \right) (1 + \delta) \quad (4.30)$$

$$\text{dist}(\mathcal{A}\mathbf{x}_t, \mathcal{K}) \leq \frac{2}{\sqrt{t+1}} \left( \beta_0 \|\mathbf{y}_\star\| + D_{\mathcal{X}} \sqrt{\left( \frac{\beta_0 L_f}{\sqrt{t+1}} + \|\mathcal{A}\|^2 \right) (1 + \delta)} \right). \quad (4.31)$$

### Inexact Oracle with Multiplicative Error

We consider the multiplicative inexact oracle. Given a gradient direction  $\nabla F_{\beta_t}(\mathbf{x}_t)$ , we assume that the approximate lmo returns an element  $\tilde{\mathbf{h}}_t \in \mathcal{X}$  such that:

$$\langle \nabla F_{\beta_t}(\mathbf{x}_t), \tilde{\mathbf{h}}_t - \mathbf{x}_t \rangle \leq \delta \cdot \langle \nabla F_{\beta_t}(\mathbf{x}_t), \mathbf{h}_t - \mathbf{x}_t \rangle \quad \text{for some } \delta \in (0, 1], \quad (4.32)$$

where  $\mathbf{h}_t$  denotes the exact solution of the lmo. Replacing the exact lmo with the approximate oracles of the form (4.32) in HCGM, we get the following convergence guarantees:

**Theorem 4.7.** Sequence  $(\mathbf{x}_t : t = 1, 2, \dots)$  generated by HCGM with approximate lmo (4.32), and modifying  $\eta_t = \frac{2}{\delta t + 2}$  and  $\beta_t = \frac{\beta_0}{\sqrt{\delta(t+1)+1}}$  satisfies:

$$F_{\beta_{t-1}}(\mathbf{x}_{t+1}) - f_\star \leq \frac{2}{\delta} \left( \frac{D_{\mathcal{X}}^2 L_f + \delta \mathcal{E}}{\delta(t+1)+2} + \frac{D_{\mathcal{X}}^2 \|\mathcal{A}\|^2}{\beta_0 \sqrt{\delta(t+1)+2}} \right) \quad \text{where } \mathcal{E} = F(\mathbf{x}_0) - f_\star. \quad (4.33)$$

**Theorem 4.8.** Assume that  $g$  is  $L_g$ -Lipschitz continuous. Then, the sequence  $(\mathbf{x}_t : t = 1, 2, \dots)$  generated by HCGM with approximate lmo (4.32), and modifying  $\eta_t = \frac{2}{\delta t + 2}$  and  $\beta_t = \frac{\beta_0}{\sqrt{\delta(t+1)+1}}$  satisfies:

$$F(\mathbf{x}_t) - F_\star \leq \frac{2}{\delta} \left( \frac{D_{\mathcal{X}}^2 L_f + \delta \mathcal{E}}{\delta(t+1)+1} + \frac{D_{\mathcal{X}}^2 \|\mathcal{A}\|^2}{\beta_0 \sqrt{\delta(t+1)+1}} \right) + \frac{\beta_0 L_g^2}{2\sqrt{\delta(t+1)+1}}, \quad \text{where } \mathcal{E} = F(\mathbf{x}_0) - F_\star.$$

**Theorem 4.9.** Sequence  $(\mathbf{x}_t : t = 1, 2, \dots)$  generated by HCGM with approximate lmo (4.32), and modifying  $\eta_t = \frac{2}{\delta t + 2}$  and  $\beta_t = \frac{\beta_0}{\sqrt{\delta(t+1)+1}}$  satisfies:

$$f(\mathbf{x}_t) - f_\star \geq -\|\mathbf{y}_\star\| \cdot \text{dist}(\mathcal{A}\mathbf{x}_t, \mathcal{K}) \quad (4.34)$$

$$f(\mathbf{x}_t) - f_\star \leq \frac{2}{\delta} \left( \frac{D_{\mathcal{X}}^2 L_f + \delta \mathcal{E}}{\delta(t+1)+1} + \frac{D_{\mathcal{X}}^2 \|\mathcal{A}\|^2}{\beta_0 \sqrt{\delta(t+1)+1}} \right) \quad (4.35)$$

$$\text{dist}(\mathcal{A}\mathbf{x}_t, \mathcal{K}) \leq \frac{2\beta_0}{\sqrt{\delta(t+1)+1}} \left( \|\mathbf{y}_\star\| + \sqrt{\frac{D_{\mathcal{X}}^2 C_0 + \delta \mathcal{E}}{\beta_0 \delta}} \right) \quad (4.36)$$

where  $\mathcal{E} = F(\mathbf{x}_0) - f_\star$  and  $C_0 = L_f + \|\mathcal{A}\|^2 / \beta_0$ .

#### 4.2.5 Application to the SDP template

We can use HCGM for solving instances from the SDP template (1.4), by choosing

$$f(\mathbf{X}) = \langle \mathbf{C}, \mathbf{X} \rangle, \quad \text{and} \quad \mathcal{X} = \{\mathbf{X} : \mathbf{X} \in \mathbb{S}_+^n, \text{Tr } \mathbf{X} = \alpha\}, \quad \text{and} \quad \mathcal{K} = \{\mathbf{b}\}. \quad (4.37)$$

We can compute the gradient of the smooth-objective  $F_{\beta_t}$  at  $\mathbf{X}_t$  by

$$\nabla F_{\beta_t}(\mathbf{X}_t) = \mathbf{C} + \beta_t^{-1} \mathcal{A}^*(\mathcal{A}\mathbf{X} - \mathbf{b}) \quad \text{where} \quad \beta_t = \frac{\beta_0}{\sqrt{t+2}}. \quad (4.38)$$

Then, we can evaluate  $\text{lmo}_{\mathcal{X}}$  by using MinEigVec which returns the minimum eigenvalue eigenvector pair:

$$\mathbf{H}_t = \alpha \mathbf{u}_t \mathbf{u}_t^* \quad \text{where} \quad (\mathbf{u}_t, \lambda_t) = \text{MinEigVec}(\mathbf{C} + \beta_t^{-1} \mathcal{A}^*(\mathcal{A}\mathbf{X} - \mathbf{b})) \quad (4.39)$$

Algorithm 4.1 shows the complete pseudocode of HCGM for this template.

---

**Algorithm 4.1** HCGM for the standard SDP formulation (1.4)

---

**Require:** initial state  $\mathbf{X}$ ; initial penalty parameter  $\beta_0$ ; (optional) target accuracy  $\varepsilon$

```

1: function HCGM
2:   for  $t \leftarrow 0, 1, 2, \dots$  do
3:      $\beta \leftarrow \beta_0 / \sqrt{t+2}$ 
4:      $(\mathbf{u}, \lambda) \leftarrow \text{MinEigVec}(\mathbf{C} + \beta^{-1} \mathcal{A}^*(\mathcal{A}\mathbf{X} - \mathbf{b}))$ 
5:      $\mathbf{H} \leftarrow \alpha \mathbf{u} \mathbf{u}^*$ 
6:     if  $\langle \mathbf{X} - \mathbf{H}, \mathbf{C} + \beta^{-1} \mathcal{A}^*(\mathcal{A}\mathbf{X} - \mathbf{b}) \rangle \leq \varepsilon$  and  $\|\mathcal{A}\mathbf{X} - \mathbf{b}\| \leq \varepsilon$  then break for
7:     end if
8:      $\eta \leftarrow 2/(t+2)$ 
9:      $\mathbf{X} \leftarrow (1-\eta)\mathbf{X} + \eta\mathbf{H}$ 
10:  end for
11:  return  $\mathbf{X}$ 
12: end function

```

---

### 4.3 Applications & Related Work

CGM is proposed for the first time in the seminal work of Frank and Wolfe [FW56] for solving smooth convex optimization on a polytope. It is then progressively generalized for more general settings in [LP66, DH78, Dun79, Dun80]. Nevertheless, with the introduction of the fast gradient methods with  $\mathcal{O}(1/t^2)$  rate by Nesterov [Nes83], the development of CGM-type methods entered into a stagnation period.

The recent developments in machine learning applications with vast data brought the scalability of the first order methods under scrutiny. As a result, there has been a renewed interest in CGM in the last decade. We compare our framework with the recent developments of CGM literature in different camps of problem templates below.

#### 4.3.1 Smooth Problems

CGM is extended for the smooth convex minimization over a simplex by Clarkson [Cla10], for a spectrahedron by Hazan [Haz08], and for an arbitrary compact convex set by Jaggi [Jag13].

When applied to smooth problems, HCGM is equivalent to the classical CGM, and it recovers the known optimal  $\mathcal{O}(1/t)$  convergence rate. We refer to [Jag13] for a review on the applications of the smooth template.

We need to mention that Nesterov [Nes17] relaxes the smoothness assumption, by showing that CGM also converges for weakly-smooth objectives, *i.e.*, objectives with Hölder continuous gradients of order  $\nu \in (0, 1]$ .

#### 4.3.2 Regularized Problems

CGM for composite problems is considered by Nesterov [Nes17] and Xu [Xu17]. A similar but slightly different template, where  $\mathcal{X}$  and  $g$  are assumed to be a closed convex cone and a norm respectively, is also studied by Harchaoui et al. [HJN15]. However, these works are based on the resolvents of a modified oracle instead of the standard lmo,

$$\arg \min_{\mathbf{x} \in \mathcal{X}} \langle \mathbf{x}, \mathbf{v} \rangle + g(\mathcal{A}\mathbf{x}), \quad (4.40)$$

which can be expensive to compute unless  $\mathcal{X} \equiv \mathbb{R}^p$ , or  $g = 0$ .

HCGM applies to the problem template (4.1) by leveraging the prox of the regularizer and the lmo of the domain separately. This allows us to consider additional sparsity, group sparsity and structured sparsity promoting regularizations, elastic-net regularization, total variation regularization and many others under the CGM framework.

The semi-proximal mirror-prox method [HH15] is also based on the smoothing idea, but the motivation of this method is fundamentally different than ours. This method considers the regularizers for which the prox is difficult to compute, but can be approximated via CGM.

### 4.3.3 Non-Smooth Problems

Template (4.1) covers the non-smooth convex minimization template as a special case:

$$\underset{\mathbf{x} \in \mathcal{X}}{\text{minimize}} \quad g(\mathcal{A}\mathbf{x}). \quad (4.41)$$

Unfortunately, the classical CGM provably cannot handle the non-smooth minimization, as shown by the following counter-example from [Nes17].

*Example.* Let  $\mathcal{X}$  be the unit Euclidean norm ball in  $\mathbb{R}^2$ , and  $g(\mathbf{x}) = \max\{x_{(1)}, x_{(2)}\}$ . Clearly,  $\mathbf{x}_\star = [-\frac{1}{\sqrt{2}}, -\frac{1}{\sqrt{2}}]^*$ . Choose an initial point  $\mathbf{x}_0 \neq \mathbf{x}_\star$ . We can use an oracle that returns a subgradient  $\nabla f(\mathbf{x}) \in [\frac{1}{0}, \frac{0}{1}]$  at any point  $\mathbf{x} \in \mathcal{X}$ . Therefore, lmo returns  $[\frac{-1}{0}]$  or  $[\frac{0}{-1}]$  at each iteration, and  $\mathbf{x}_t$  belongs to the convex hull of  $\{\mathbf{x}_0, [\frac{-1}{0}], [\frac{0}{-1}]\}$  which does not contain the solution.

Our framework escapes such issues by leveraging the prox of the objective function  $g$ . In the example above,  $\text{prox}_g$  corresponds to the projection onto the simplex. Often times the cost of  $\text{prox}_g$  is negligible in comparison to the cost of lmo $_{\mathcal{X}}$ .

Assume that  $g : \mathbb{R}^d \rightarrow \mathbb{R}$  is  $L_g$ -Lipschitz continuous. As a consequence of Theorem 4.2, HCGM for solving (4.41) by choosing  $\beta_0 = 2D_{\mathcal{X}}\|\mathcal{A}\|/L_g$  satisfies

$$g(\mathcal{A}\mathbf{x}_t) - g_\star \leq \frac{2D_{\mathcal{X}}\|\mathcal{A}\|L_g}{\sqrt{t+1}}. \quad (4.42)$$

We recover the method proposed by Lan [Lan14] in this specific setting. He also proves that this rate is optimal for any algorithm that approximates the solution of (4.41) as a convex combination of lmo outputs.

We extend the analysis in this setting for inexact oracles. In contrast to the smooth case, where the additive error should decrease by  $\mathcal{O}(1/t)$  rate, definition (4.26) implies that we can preserve the convergence rate in the non-smooth case as long as the additive error is  $\mathcal{O}(1/\sqrt{t})$ .

### 4.3.4 Minimax Problems

We consider the minimax problems of the following form:

$$\min_{\mathbf{x} \in \mathcal{X}} \max_{\mathbf{y} \in \mathcal{Y}} \mathcal{L}(\mathcal{A}\mathbf{x}, \mathbf{y}) \quad (4.43)$$

where  $\mathcal{L}$  is a smooth convex-concave function, *i.e.*,  $\mathcal{L}(\cdot, \mathbf{y})$  is convex  $\forall \mathbf{y} \in \mathcal{Y}$  and  $\mathcal{L}(\mathcal{A}\mathbf{x}, \cdot)$  is concave  $\forall \mathbf{x} \in \mathcal{X}$ . Note that this formulation is a special instance of (4.41) with  $g(\mathcal{A}\mathbf{x}) = \max_{\mathbf{y} \in \mathcal{Y}} \mathcal{L}(\mathcal{A}\mathbf{x}, \mathbf{y})$ . Consequently, we can apply HCGM if  $\text{prox}_g$  is tractable.

## Chapter 4. CGM for Composite Problems

---

When  $\mathcal{Y}$  admits an efficient projection oracle,  $\text{prox}_g$  is also efficient for bilinear saddle point problems  $\mathcal{L}(\mathcal{A}\mathbf{x}, \mathbf{y}) = \langle \mathcal{A}\mathbf{x}, \mathbf{y} \rangle$ . By Moreau decomposition, we have

$$\text{prox}_g(\mathcal{A}\mathbf{x}_t) = \mathcal{A}\mathbf{x}_t - \text{proj}_{\mathcal{Y}}(\mathcal{A}\mathbf{x}_t), \quad (4.44)$$

hence  $\nabla F_{\beta_t}(\mathbf{x}_t)$  takes the form

$$\nabla F_{\beta_t}(\mathbf{x}_t) = \nabla f(\mathbf{x}_t) + \frac{1}{\beta_t} \mathcal{A}^* \text{proj}_{\mathcal{Y}}(\mathcal{A}\mathbf{x}_t). \quad (4.45)$$

Gidel et al. [GJLJ17] propose a CGM variant for the smooth convex-concave saddle point problems. This method processes both  $\mathbf{x}$  and  $\mathbf{y}$  via the lmo, and hence it also requires  $\mathcal{Y}$  to be bounded. Our method, on the other hand, is more suitable when  $\text{proj}_{\mathcal{Y}}$  is easy.

Bilinear saddle point problem covers the maximum margin estimation of structured output models [TLJJ06] and minimax games [VNM44]. In particular, it also covers an important semidefinite programming formulation [GH16], where  $\mathcal{X}$  is a spectrahedron and  $\mathcal{Y}$  is the simplex. Our framework fits perfectly here since the projection onto the simplex can be computed efficiently. We defer the extension of our framework with the entropy Bregman smoothing for future.

Note that the CGM is applicable also for the variational inequality problems beyond (4.1), see [Ham84], [JN16] and [CJN17].

### 4.3.5 Problems with Affine Constraints

Our algorithms apply to smooth convex minimization problems with affine constraints over a convex compact set:

$$\underset{\mathbf{x} \in \mathcal{X}}{\text{minimize}} \quad f(\mathbf{x}) \quad \text{subject to} \quad \mathcal{A}\mathbf{x} \in \mathcal{K}, \quad (4.46)$$

by setting  $g(\mathcal{A}\mathbf{x})$  in (4.1) as indicator function of set  $\mathcal{K}$ .

Since the prox operator of the indicator function is the projection,  $\nabla F_{\beta_t}(\mathbf{x}_t)$  becomes

$$\nabla F_{\beta_t}(\mathbf{x}_t) = \nabla f(\mathbf{x}_t) + \frac{1}{\beta_t} \mathcal{A}^* (\mathcal{A}\mathbf{x}_t - \text{proj}_{\mathcal{K}}(\mathcal{A}\mathbf{x}_t)). \quad (4.47)$$

We implicitly assume that  $\text{proj}_{\mathcal{K}}$  is tractable. We can use a splitting scheme when it is computationally more advantageous to use  $\text{lmo}_{\mathcal{K}}$  instead. See Section 4.3.6 for the details.

This template covers the standard semidefinite programming in particular. Applications include clustering [PW07], optimal power-flow [LL12], sparse PCA [dGJL07], kernel learning [LCG<sup>+</sup>04], blind deconvolution [ARR14], community detection [BBV16], *etc.* Besides machine learning applications, this formulation has a crucial role in the convex relaxation of combinatorial problems.

A significant example is the problems over the doubly nonnegative cone (*i.e.*, the intersection of the positive semidefinite cone and the positive orthant) with a bounded trace norm [YM10]. Note that the lmo over this domain can be costly since the lmo can require full dimensional updates [HJL96, LTRJ17]. Our framework can handle these problems ensuring the positive semidefiniteness by  $\text{lmo}_{\mathcal{K}}$ , and can still ensure the convergence to the first orthant via  $\text{proj}_{\mathcal{K}}$ .

#### 4.3.6 Minimization via Splitting

We can take advantage of splitting since we can handle affine constraints. This lets us to disentangle the complexity of the constraints. Consider the following optimization template:

$$\begin{aligned} & \underset{\mathbf{x} \in \mathcal{X}_1 \cap \mathcal{X}_2}{\text{minimize}} && f(\mathbf{x}) + g(\mathcal{A}_1 \mathbf{x}) \\ & \text{subject to} && \mathcal{A}_2 \mathbf{x} \in \mathcal{K}, \quad \mathcal{A}_3 \mathbf{x} \in \mathcal{S} \end{aligned} \tag{4.48}$$

where  $\mathcal{X}_1, \mathcal{X}_2 \subset \mathbb{R}^p$  are two convex compact sets,  $\mathcal{A}_1, \mathcal{A}_2$  and  $\mathcal{A}_3$  are known linear maps.

Suppose that

- $\text{lmo}_{\mathcal{X}_1}$  and  $\text{lmo}_{\mathcal{X}_2}$  are easy to compute, but not  $\text{lmo}_{\mathcal{X}_1 \cap \mathcal{X}_2}$
- $\text{prox}_g$  is easy to compute
- $\mathcal{K}$  is a simple convex set and  $\text{proj}_{\mathcal{K}}$  is efficient
- $\mathcal{S}$  is a convex compact set with an efficient lmo.

We can reformulate this problem introducing slack variables  $\xi \in \mathcal{X}_2$  and  $\zeta \in \mathcal{S}$  as follows:

$$\begin{aligned} & \underset{\substack{\mathbf{x} \in \mathcal{X}_1, \xi \in \mathcal{X}_2 \\ \zeta \in \mathcal{S}}}{\text{minimize}} && f(\mathbf{x}) + g(\mathcal{A}_1 \mathbf{x}) \\ & \text{subject to} && \mathcal{A}_2 \mathbf{x} \in \mathcal{K}, \quad \mathcal{A}_3 \mathbf{x} = \zeta, \quad \mathbf{x} = \xi. \end{aligned} \tag{4.49}$$

This formulation is in the form of (4.1) with respect to the variable  $(\mathbf{x}, \xi, \zeta) \in \mathcal{X}_1 \times \mathcal{X}_2 \times \mathcal{S}$ . It is easy to verify that HCGM leverages  $\text{lmo}_{\mathcal{X}_1}$ ,  $\text{lmo}_{\mathcal{X}_2}$ ,  $\text{lmo}_{\mathcal{S}}$ ,  $\text{prox}_g$  and  $\text{proj}_{\mathcal{K}}$  separately. This approach can be generalized for an arbitrary finite number of non-smooth terms in a straightforward way.

## 4.4 Numerical Experiments

This section presents numerical experiments supporting our theoretical findings in clustering and robust PCA examples.

### 4.4.1 Clustering the MNIST dataset

We consider the model-free  $k$ -means clustering SDP formulation of Peng and Wei [PW07]:

$$\text{minimize } \langle \mathbf{C}, \mathbf{X} \rangle \quad \text{subject to } \underbrace{\mathbf{X}\mathbf{1} = \mathbf{1}, \mathbf{X} \geq 0}_{\mathcal{AX} \in \mathcal{K}}, \underbrace{\mathbf{X} \in \mathbb{S}_+^n, \text{Tr } \mathbf{X} = k}_{\mathbf{X} \in \mathcal{X}}. \quad (4.50)$$

where  $\mathbf{C} \in \mathbb{R}^{n \times n}$  is the Euclidean distance matrix.

We use the setup described in [MVW17] and made available online by the authors. We can briefly describe this setup as follows: First, meaningful features from MNIST dataset<sup>1</sup> which consists of  $28 \times 28$  grayscale images that can be stacked as  $784 \times 1$  vectors, are extracted using a one-layer neural network. This gives us a weight matrix  $\mathbf{W} \in \mathbb{R}^{784 \times 10}$  and a bias vector  $\mathbf{b} \in \mathbb{R}^{10}$ . Then, the trained neural network is applied to the first 1000 elements of the test set, which gives the probability vectors for these 1000 test points, where each entry represents the probability of being each digit.

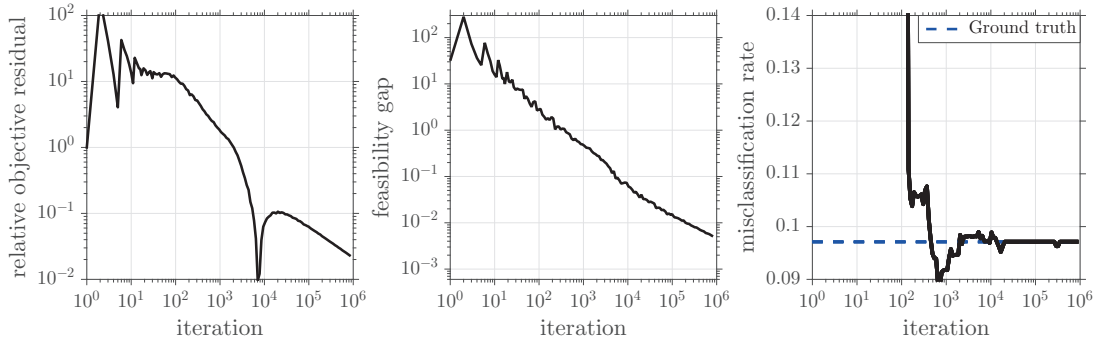


Figure 4.1 – Clustering MNIST: Empirical performance of our method. Blue dashed line on the misclassification plot represents the value reported in [MVW17].

[MVW17] runs a relax-and-round algorithm which solves (4.50) by SDPNAL+ [YST15] followed by a rounding scheme (see Section 5 of [MVW17] for details), and compares the results against MATLAB's built-in  $k$ -means++ implementation. Relax-and-round method is reported to achieve a misclassification rate of 0.0971. This rate matches with the all-time best rate for  $k$ -means++ after 100 different runs with random initializations.

For this experiment, we solve (4.50) by using HCGM. Then, we cluster data using the same rounding scheme as [MVW17]. We initialize our method from  $\mathbf{0}$  and choose  $\beta_0 = 1$ . We implement the lmo using the built-in MATLAB function `eigs` with tolerance  $10^{-9}$ .

<sup>1</sup>LeCun & Cortes. MNIST handwritten digit database. Available at "<http://yann.lecun.com/exdb/mnist/>"



We present the results of this experiment in Figure 4.1. We observe empirical  $\mathcal{O}(1/\sqrt{t})$  rate both in the objective residual and the feasibility gap. Surprisingly, the method attains the best misclassification rate around 1000 iterations, achieving 0.0914. This improves the value reported in [MVW17] by 5.8%.

This example supports the claims that a low-to-medium accuracy solution is enough for many important SDP formulations, where we need a rounding step. In this example, the low-accuracy solution generalizes as well as the optimal point, even better.

#### 4.4.2 Robust PCA

Suppose that we are given a large matrix that can be decomposed as the summation of a low-rank and a sparse (in some representation) matrix. Robust PCA aims to recover these components accurately, and it has many applications in machine learning and data science, such as collaborative filtering, system identification, genotype imputation, etc. Here, we focus on an image decomposition problem so that we can visualize the decomposition error results.

Our setting is similar to the setup described in [ZS18]. We consider a scaled grayscale photograph with pattern from [LMWY13], and we assume that we only have access to an occluded image. Moreover, the image is contaminated by salt and pepper noise of density 1/10. We seek to approximate the original from this noisy image.

This is essentially a matrix completion problem, and most of the scalable techniques rely on the Gaussian noise model. Note however the corresponding least-squares formulation is a good model against outliers:

$$\underset{\|X\|_{s_1} \leq \rho}{\text{minimize}} \quad \frac{1}{2} \|\mathcal{A}X - \mathbf{b}\|^2 \quad \text{subject to} \quad 0 \leq X \leq 1, \quad (4.51)$$

where  $\mathcal{A} : \mathbb{R}^{p \times p} \rightarrow \mathbb{R}^d$  is the sampling operator.

Our framework also covers the following least absolute deviations formulation :

$$\underset{\|X\|_{s_1} \leq \rho}{\text{minimize}} \quad \|\mathcal{A}X - \mathbf{b}\|_1 \quad \text{subject to} \quad 0 \leq X \leq 1, \quad (4.52)$$

We solve both formulations with HCGM, starting from all zero matrix, running 1000 iterations, and assuming that we know the true nuclear norm of the original image. We choose  $\beta_0 = 1$  in both cases.

This experiment demonstrates the implications of the flexibility of our framework in a simple machine learning setup. We compile the results in Figure 4.2, where the non-smooth formulation recovers a better approximation with 5dB higher peak signal to noise ratio (PSNR) and 0.27 higher structural similarity index (SSIM). Evaluation of PSNR and SSIM vs iteration counter are shown in Figure 4.3.

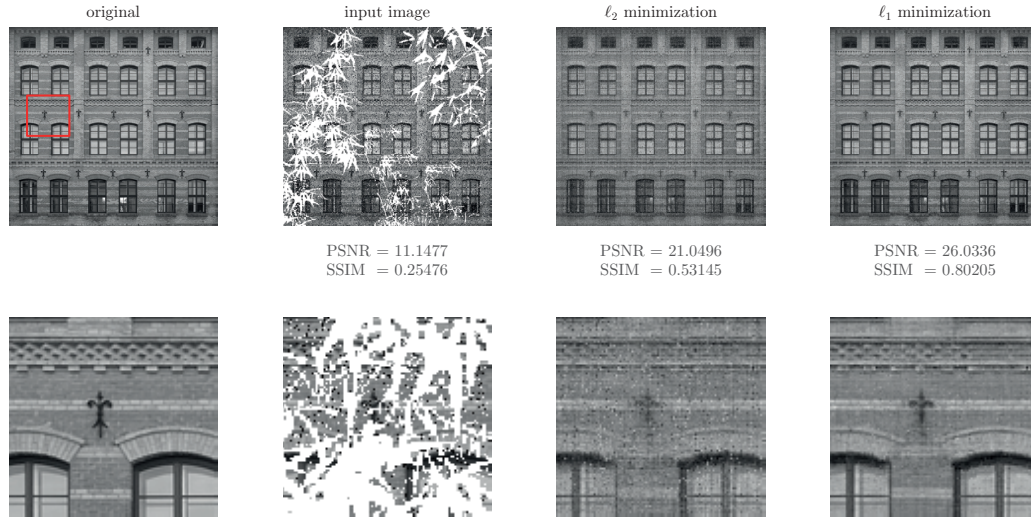


Figure 4.2 – Image inpainting from noisy test image ( $493 \times 517$ ): Robust PCA recovers a better approximation with 5dB higher PSNR.

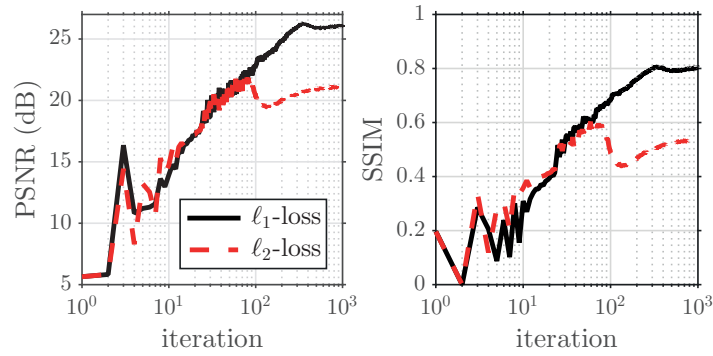


Figure 4.3 – PSNR and SSIM vs iteration counter for formulations with  $\ell_1$  and  $\ell_2$  loss.

## 4.5 Appendix: Proofs

### Preliminaries

The following properties of smoothing are key to derive the convergence rate of our algorithms.

Let  $g : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{+\infty\}$  be a proper, closed and convex function, and denote its smooth approximation by

$$g_\beta(\mathbf{z}) = \max_{\mathbf{y} \in \mathbb{R}^d} \langle \mathbf{z}, \mathbf{y} \rangle - g^*(\mathbf{y}) - \frac{\beta}{2} \|\mathbf{y}\|^2 \quad (4.53)$$

where  $g^*$  represents the Fenchel conjugate of  $g$  and  $\beta > 0$  is the smoothing parameter. Then,  $g_\beta$  is convex and  $\frac{1}{\beta}$ -smooth. Let us denote the unique maximizer of this concave problem by

$$\mathbf{y}_\beta^*(\mathbf{z}) = \arg \max_{\mathbf{y} \in \mathbb{R}^d} \langle \mathbf{z}, \mathbf{y} \rangle - g^*(\mathbf{y}) - \frac{\beta}{2} \|\mathbf{y}\|^2 \quad (4.54)$$

$$= \arg \min_{\mathbf{y} \in \mathbb{R}^d} \frac{1}{\beta} g^*(\mathbf{y}) - \frac{1}{\beta} \langle \mathbf{z}, \mathbf{y} \rangle + \frac{1}{2} \|\mathbf{y}\|^2 + \frac{1}{2} \left\| \frac{1}{\beta} \mathbf{z} \right\|^2 \quad (4.55)$$

$$= \arg \min_{\mathbf{y} \in \mathbb{R}^d} \frac{1}{\beta} g^*(\mathbf{y}) + \frac{1}{2} \left\| \mathbf{y} - \frac{1}{\beta} \mathbf{z} \right\|^2 \quad (4.56)$$

$$= \text{prox}_{\beta^{-1}g^*}(\beta^{-1}\mathbf{z}) = \frac{1}{\beta} (\mathbf{z} - \text{prox}_{\beta g}(\mathbf{z})) \quad (4.57)$$

where the last equality is known as the Moreau decomposition. Then, the followings hold for  $\forall \mathbf{z}_1, \mathbf{z}_2 \in \mathbb{R}^d$  and  $\forall \beta, \gamma > 0$

$$g_\beta(\mathbf{z}_1) \geq g_\beta(\mathbf{z}_2) + \langle \nabla g_\beta(\mathbf{z}_2), \mathbf{z}_1 - \mathbf{z}_2 \rangle + \frac{\beta}{2} \|\mathbf{y}_\beta^*(\mathbf{z}_2) - \mathbf{y}_\beta^*(\mathbf{z}_1)\|^2 \quad (4.58)$$

$$g(\mathbf{z}_1) \geq g_\beta(\mathbf{z}_2) + \langle \nabla g_\beta(\mathbf{z}_2), \mathbf{z}_1 - \mathbf{z}_2 \rangle + \frac{\beta}{2} \|\mathbf{y}_\beta^*(\mathbf{z}_2)\|^2 \quad (4.59)$$

$$g_\beta(\mathbf{z}_1) \leq g_\gamma(\mathbf{z}_1) + \frac{\gamma - \beta}{2} \|\mathbf{y}_\beta^*(\mathbf{z}_1)\|^2 \quad (4.60)$$

Proofs can be found in Lemma 10 from [TDFC18].

Suppose that  $g$  is  $L_g$ -Lipschitz continuous. Then, for any  $\beta > 0$  and any  $\mathbf{z} \in \mathbb{R}^d$ , the following bound holds:

$$g_\beta(\mathbf{z}) \leq g(\mathbf{z}) \leq g_\beta(\mathbf{z}) + \frac{\beta}{2} L_g^2 \quad (4.61)$$

Proof follows from equation (2.7) in [Nes05] with a remark on the duality between Lipschitzness and bounded support (see Lemma 5 in [DFTJ16] for a statement of this well-known duality).

### Convergence analysis

We skip the proofs of Theorems 4.1 to 4.3 since we can get these results as a special case by setting  $\delta = 0$  in Theorems 4.4 to 4.6.

#### Proof of Theorem 4.4

First, we use smoothness of  $F_{\beta_t}$  to upper bound the progress. Note that  $F_{\beta_t}$  is  $(L_f + \|\mathcal{A}\|^2/\beta_t)$ -smooth:

$$F_{\beta_t}(\mathbf{x}_{t+1}) \leq F_{\beta_t}(\mathbf{x}_t) + \eta_t \langle \nabla F_{\beta_t}(\mathbf{x}_t), \tilde{\mathbf{h}}_t - \mathbf{x}_t \rangle + \frac{\eta_t^2}{2} \|\tilde{\mathbf{h}}_t - \mathbf{x}_t\|^2 (L_f + \frac{\|\mathcal{A}\|^2}{\beta_t}) \quad (4.62)$$

$$\leq F_{\beta_t}(\mathbf{x}_t) + \eta_t \langle \nabla F_{\beta_t}(\mathbf{x}_t), \tilde{\mathbf{h}}_t - \mathbf{x}_t \rangle + \frac{\eta_t^2}{2} D_{\mathcal{X}}^2 (L_f + \frac{\|\mathcal{A}\|^2}{\beta_t}), \quad (4.63)$$

where  $\tilde{\mathbf{h}}_t$  denotes the atom selected by the inexact lmo, and the second inequality follows since  $\tilde{\mathbf{h}}_t \in \mathcal{X}$ .

By definition of inexact oracle (4.26), we have

$$\langle \nabla F_{\beta_t}(\mathbf{x}_t), \tilde{\mathbf{h}}_t - \mathbf{x}_t \rangle \leq \langle \nabla F_{\beta_t}(\mathbf{x}_t), \mathbf{h}_t - \mathbf{x}_t \rangle + \delta \frac{\eta_t}{2} D_{\mathcal{X}}^2 (L_f + \frac{\|\mathcal{A}\|^2}{\beta_t}) \quad (4.64)$$

$$\leq \langle \nabla F_{\beta_t}(\mathbf{x}_t), \mathbf{x}_\star - \mathbf{x}_t \rangle + \delta \frac{\eta_t}{2} D_{\mathcal{X}}^2 (L_f + \frac{\|\mathcal{A}\|^2}{\beta_t}) \quad (4.65)$$

$$= \langle \nabla f(\mathbf{x}_t), \mathbf{x}_\star - \mathbf{x}_t \rangle + \langle \mathcal{A}^* \nabla g_{\beta_t}(\mathcal{A}\mathbf{x}_t), \mathbf{x}_\star - \mathbf{x}_t \rangle + \delta \frac{\eta_t}{2} D_{\mathcal{X}}^2 (L_f + \frac{\|\mathcal{A}\|^2}{\beta_t}), \quad (4.66)$$

where the second line follows since  $\mathbf{h}_t$  is a solution of  $\min_{\mathbf{x} \in \mathcal{X}} \langle \nabla F_{\beta_t}(\mathbf{x}_t), \mathbf{x} \rangle$ .

Now, convexity of  $f$  ensures  $\langle \nabla f(\mathbf{x}_t), \mathbf{x}_\star - \mathbf{x}_t \rangle \leq f(\mathbf{x}_\star) - f(\mathbf{x}_t)$ . Using property (4.59), we have

$$\langle \mathcal{A}^* \nabla g_{\beta_t}(\mathcal{A}\mathbf{x}_t), \mathbf{x}_\star - \mathbf{x}_t \rangle = \langle \nabla g_{\beta_t}(\mathcal{A}\mathbf{x}_t), \mathcal{A}\mathbf{x}_\star - \mathcal{A}\mathbf{x}_t \rangle \quad (4.67)$$

$$\leq g(\mathcal{A}\mathbf{x}_\star) - g_{\beta_t}(\mathcal{A}\mathbf{x}_t) - \frac{\beta_t}{2} \|\mathbf{y}_{\beta_t}^*(\mathcal{A}\mathbf{x}_t)\|^2. \quad (4.68)$$

Putting these altogether, we get the following bound

$$\begin{aligned} F_{\beta_t}(\mathbf{x}_{t+1}) &\leq F_{\beta_t}(\mathbf{x}_t) + \eta_t \left( f(\mathbf{x}_\star) - f(\mathbf{x}_t) + g(\mathcal{A}\mathbf{x}_\star) - g_{\beta_t}(\mathcal{A}\mathbf{x}_t) - \frac{\beta_t}{2} \|\nabla \mathbf{y}_{\beta_t}^*(\mathcal{A}\mathbf{x}_t)\|^2 \right) \\ &\quad + \frac{\eta_t^2}{2} D_{\mathcal{X}}^2 (L_f + \frac{\|\mathcal{A}\|^2}{\beta_t}) (1 + \delta) \\ &= (1 - \eta_t) F_{\beta_t}(\mathbf{x}_t) + \eta_t F_\star - \frac{\eta_t \beta_t}{2} \|\nabla \mathbf{y}_{\beta_t}^*(\mathcal{A}\mathbf{x}_t)\|^2 + \frac{\eta_t^2}{2} D_{\mathcal{X}}^2 (L_f + \frac{\|\mathcal{A}\|^2}{\beta_t}) (1 + \delta). \end{aligned} \quad (4.69)$$

Now, using (4.60), we get

$$F_{\beta_t}(\mathbf{x}_t) = f(\mathbf{x}_t) + g_{\beta_t}(\mathcal{A}\mathbf{x}_t) \quad (4.70)$$

$$\leq f(\mathbf{x}_t) + g_{\beta_{t-1}}(\mathcal{A}\mathbf{x}_t) + \frac{\beta_{t-1} - \beta_t}{2} \|\mathbf{y}_{\beta_t}^*(\mathcal{A}\mathbf{x}_t)\|^2 \quad (4.71)$$

$$= F_{\beta_{t-1}}(\mathbf{x}_t) + \frac{\beta_{t-1} - \beta_t}{2} \|\mathbf{y}_{\beta_t}^*(\mathcal{A}\mathbf{x}_t)\|^2. \quad (4.72)$$

We combine this with (4.69) and subtract  $F_\star$  from both sides to get

$$\begin{aligned} F_{\beta_t}(\mathbf{x}_{t+1}) - F_\star &\leq (1 - \eta_t)(F_{\beta_{t-1}}(\mathbf{x}_t) - F_\star) + \frac{\eta_t^2}{2} D_{\mathcal{X}}^2(L_f + \frac{\|\mathcal{A}\|^2}{\beta_t})(1 + \delta) \\ &\quad + ((1 - \eta_t)(\beta_{t-1} - \beta_t) - \eta_t \beta_t) \frac{1}{2} \|\mathbf{y}_{\beta_t}^*(\mathcal{A}\mathbf{x}_t)\|^2. \end{aligned} \quad (4.73)$$

Let us choose  $\eta_t$  and  $\beta_t$  in a way to vanish the last term. By choosing  $\eta_t = \frac{2}{t+2}$  and  $\beta_t = \frac{\beta_0}{\sqrt{t+2}}$  for  $t \geq 0$  with some  $\beta_0 > 0$ , we get  $(1 - \eta_t)(\beta_{t-1} - \beta_t) - \eta_t \beta_t < 0$ . Hence, we end up with

$$F_{\beta_t}(\mathbf{x}_{t+1}) - F_\star \leq (1 - \eta_t)(F_{\beta_{t-1}}(\mathbf{x}_t) - F_\star) + \frac{\eta_t^2}{2} D_{\mathcal{X}}^2(L_f + \frac{\|\mathcal{A}\|^2}{\beta_t})(1 + \delta). \quad (4.74)$$

By recursively applying this inequality, we get

$$\begin{aligned} F_{\beta_t}(\mathbf{x}_{t+1}) - F_\star &\leq \prod_{j=0}^t (1 - \eta_j) (F_{\beta_{j-1}}(\mathbf{x}_t) - F_\star) + \frac{1}{2} D_{\mathcal{X}}^2(1 + \delta) \sum_{\ell=0}^t \eta_\ell^2 \prod_{j=\ell}^t (1 - \eta_j) (L_f + \frac{\|\mathcal{A}\|^2}{\beta_\ell}) \end{aligned} \quad (4.75)$$

$$\leq \prod_{j=0}^t (1 - \eta_j) (F_{\beta_{j-1}}(\mathbf{x}_t) - F_\star) + \frac{1}{2} D_{\mathcal{X}}^2(L_f + \frac{\|\mathcal{A}\|^2}{\beta_t})(1 + \delta) \sum_{\ell=0}^t \eta_\ell^2 \prod_{j=\ell}^t (1 - \eta_j) \quad (4.76)$$

$$= \frac{1}{2} D_{\mathcal{X}}^2(L_f + \frac{\|\mathcal{A}\|^2}{\beta_t})(1 + \delta) \sum_{\ell=0}^t \eta_\ell^2 \prod_{j=\ell}^t (1 - \eta_j), \quad (4.77)$$

where the second inequality follows since  $\beta_t \leq \beta_j$  for any positive integer  $j \leq t$ , and the last line since  $\eta_0 = 1$ .

Now, we use the following relation

$$\sum_{\ell=0}^t \eta_\ell^2 \prod_{j=\ell}^t (1 - \eta_j) = \sum_{\ell=0}^t \frac{4}{(\ell+2)^2} \prod_{j=\ell}^t \frac{j}{j+2} = \sum_{\ell=0}^t \frac{4}{(\ell+2)^2} \frac{\ell(\ell+1)}{(t+1)(t+2)} \leq \frac{4}{t+2}, \quad (4.78)$$

which yields the first result of Theorem 4.4 as

$$F_{\beta_t}(\mathbf{x}_{t+1}) - F_\star \leq \frac{2}{t+2} D_{\mathcal{X}}^2(L_f + \frac{\|\mathcal{A}\|^2}{\beta_t})(1 + \delta) = 2D_{\mathcal{X}}^2(\frac{L_f}{t+2} + \frac{\|\mathcal{A}\|^2}{\beta_0 \sqrt{t+2}})(1 + \delta). \quad (4.79)$$

**Proof of Theorem 4.5**

Now, we further assume that  $g : \mathbb{R}^d \rightarrow \mathbb{R} \cup \{+\infty\}$  is  $L_g$ -Lipschitz continuous. From (4.61), we get

$$g(\mathcal{A}\mathbf{x}_{t+1}) \leq g_{\beta_t}(\mathcal{A}\mathbf{x}_{t+1}) + \frac{\beta_t L_g^2}{2} = g_{\beta_t}(\mathcal{A}\mathbf{x}_{t+1}) + \frac{\beta_0 L_g^2}{2\sqrt{t+2}}. \quad (4.80)$$

We complete the proof by adding  $f(\mathbf{x}_{t+1}) - F_\star$  to both sides:

$$F(\mathbf{x}_{t+1}) - F_\star \leq F_{\beta_t}(\mathbf{x}_{t+1}) - F_\star + \frac{\beta_0 L_g^2}{2\sqrt{t+2}}. \quad (4.81)$$

**Proof of Theorem 4.6**

From the Lagrange saddle point theory, we know that the following bound holds  $\forall \mathbf{x} \in \mathcal{X}$  and  $\forall \mathbf{v} \in \mathcal{K}$ :

$$f_\star \leq \mathcal{L}(\mathbf{x}, \mathbf{v}, \mathbf{y}_\star) = f(\mathbf{x}) + \langle \mathbf{y}_\star, \mathcal{A}\mathbf{x} - \mathbf{v} \rangle \leq f(\mathbf{x}) + \|\mathbf{y}_\star\| \|\mathcal{A}\mathbf{x} - \mathbf{v}\|, \quad (4.82)$$

Since  $\mathbf{x}_{t+1} \in \mathcal{X}$ , we get

$$f(\mathbf{x}_{t+1}) - f_\star \geq -\min_{\mathbf{v} \in \mathcal{K}} \|\mathbf{y}_\star\| \|\mathcal{A}\mathbf{x}_{t+1} - \mathbf{v}\| = -\|\mathbf{y}_\star\| \text{dist}(\mathcal{A}\mathbf{x}_{t+1}, \mathcal{K}). \quad (4.83)$$

This proves the first bound in Theorem 4.6.

The second bound directly follows by Theorem 4.4 as

$$f(\mathbf{x}_{t+1}) - f_\star \leq f(\mathbf{x}_{t+1}) - f_\star + \frac{1}{2\beta_t} \text{dist}^2(\mathcal{A}\mathbf{x}_{t+1}, \mathcal{K}) := F_{\beta_t}(\mathbf{x}_{t+1}) - f_\star \quad (4.84)$$

$$\leq 2D_{\mathcal{X}}^2 \left( \frac{L_f}{t+2} + \frac{\|\mathcal{A}\|^2}{\beta_0 \sqrt{t+2}} \right) (1 + \delta). \quad (4.85)$$

Denote by  $\phi_t := 2D_{\mathcal{X}}^2 \left( \frac{L_f}{t} + \frac{\|\mathcal{A}\|^2}{\beta_0 \sqrt{t}} \right) (1 + \delta)$ . Then, we combine this with (4.83), and we get

$$-\|\mathbf{y}_\star\| \text{dist}(\mathcal{A}\mathbf{x}_{t+1}, \mathcal{K}) + \frac{1}{2\beta_t} \text{dist}^2(\mathcal{A}\mathbf{x}_{t+1}, \mathcal{K}) \leq \phi_{t+2}. \quad (4.86)$$

This is a second order inequality in terms of  $\text{dist}(\mathcal{A}\mathbf{x}_t, \mathcal{K})$ . Solving this inequality, we get

$$\text{dist}(\mathcal{A}\mathbf{x}_{t+1}, \mathcal{K}) \leq \beta_t \left( \|\mathbf{y}_\star\| + \sqrt{\|\mathbf{y}_\star\|^2 + 2\frac{\phi_{t+2}}{\beta_t}} \right) \quad (4.87)$$

$$\leq \frac{2}{\sqrt{t+2}} \left( \beta_0 \|\mathbf{y}_\star\| + D_{\mathcal{X}} \sqrt{\left( \frac{L_f \beta_0}{\sqrt{t+2}} + \|\mathcal{A}\|^2 \right) (1 + \delta)} \right). \quad (4.88)$$

**Proof of Theorem 4.7**

Let us define the multiplicative error  $\delta$  of the LMO:

$$\langle \nabla F_{\beta_t}(\mathbf{x}_t), \tilde{\mathbf{h}}_t - \mathbf{x}_t \rangle \leq \delta \langle \nabla F_{\beta_t}(\mathbf{x}_t), \mathbf{h}_t - \mathbf{x}_t \rangle \quad (4.89)$$

For the proof we assume that  $\mathbf{x}_0$  is feasible. First, we use the smoothness of  $F_{\beta_t}$  to upper bound the progress. Note that  $F_{\beta_t}$  is  $(L_f + \|\mathcal{A}\|^2/\beta_t)$ -smooth.

$$F_{\beta_t}(\mathbf{x}_{t+1}) \leq F_{\beta_t}(\mathbf{x}_t) + \eta_t \langle \nabla F_{\beta_t}(\mathbf{x}_t), \tilde{\mathbf{h}}_t - \mathbf{x}_t \rangle + \frac{\eta_t^2}{2} \|\tilde{\mathbf{h}}_t - \mathbf{x}_t\|^2 (L_f + \frac{\|\mathcal{A}\|^2}{\beta_t}) \quad (4.90)$$

$$\leq F_{\beta_t}(\mathbf{x}_t) + \eta_t \langle \nabla F_{\beta_t}(\mathbf{x}_t), \tilde{\mathbf{h}}_t - \mathbf{x}_t \rangle + \frac{\eta_t^2}{2} D_{\mathcal{X}}^2 (L_f + \frac{\|\mathcal{A}\|^2}{\beta_t}), \quad (4.91)$$

where  $\tilde{\mathbf{h}}_t$  denotes the atom selected by the inexact linear minimization oracle, and the second inequality follows since  $\tilde{\mathbf{h}}_t \in \mathcal{X}$ . By definition of inexact oracle (4.89), we have

$$\langle \nabla F_{\beta_t}(\mathbf{x}_t), \tilde{\mathbf{h}}_t - \mathbf{x}_t \rangle \leq \delta \langle \nabla F_{\beta_t}(\mathbf{x}_t), \mathbf{h}_t - \mathbf{x}_t \rangle \quad (4.92)$$

$$\leq \delta \langle \nabla F_{\beta_t}(\mathbf{x}_t), \mathbf{x}_\star - \mathbf{x}_t \rangle \quad (4.93)$$

$$= \delta \langle \nabla f(\mathbf{x}_t), \mathbf{x}_\star - \mathbf{x}_t \rangle + \delta \langle \mathcal{A}^* \nabla g_{\beta_t}(\mathcal{A}\mathbf{x}_t), \mathbf{x}_\star - \mathbf{x}_t \rangle, \quad (4.94)$$

where the second line follows since  $\mathbf{h}_t$  is a solution of  $\min_{\mathbf{x} \in \mathcal{X}} \langle \nabla F_{\beta_t}(\mathbf{x}_t), \mathbf{x} \rangle$ .

Now, convexity of  $f$  ensures  $\langle \nabla f(\mathbf{x}_t), \mathbf{x}_\star - \mathbf{x}_t \rangle \leq f(\mathbf{x}_\star) - f(\mathbf{x}_t)$ . Using property (4.59), we have

$$\langle \mathcal{A}^* \nabla g_{\beta_t}(\mathcal{A}\mathbf{x}_t), \mathbf{x}_\star - \mathbf{x}_t \rangle = \langle \nabla g_{\beta_t}(\mathcal{A}\mathbf{x}_t), \mathcal{A}\mathbf{x}_\star - \mathcal{A}\mathbf{x}_t \rangle \quad (4.95)$$

$$\leq g(\mathcal{A}\mathbf{x}_\star) - g_{\beta_t}(\mathcal{A}\mathbf{x}_t) - \frac{\beta_t}{2} \|\mathbf{y}_{\beta_t}^*(\mathcal{A}\mathbf{x}_t)\|^2. \quad (4.96)$$

Putting these altogether, we get the following bound

$$\begin{aligned} F_{\beta_t}(\mathbf{x}_{t+1}) &\leq F_{\beta_t}(\mathbf{x}_t) + \eta_t \delta \left( f(\mathbf{x}_\star) - f(\mathbf{x}_t) + g(\mathcal{A}\mathbf{x}_\star) - g_{\beta_t}(\mathcal{A}\mathbf{x}_t) - \frac{\beta_t}{2} \|\nabla \mathbf{y}_{\beta_t}^*(\mathcal{A}\mathbf{x}_t)\|^2 \right) \\ &\quad + \frac{\eta_t^2}{2} D_{\mathcal{X}}^2 (L_f + \frac{\|\mathcal{A}\|^2}{\beta_t}) \\ &= (1 - \delta\eta_t) F_{\beta_t}(\mathbf{x}_t) + \delta\eta_t F_\star - \frac{\delta\eta_t\beta_t}{2} \|\nabla \mathbf{y}_{\beta_t}^*(\mathcal{A}\mathbf{x}_t)\|^2 + \frac{\eta_t^2}{2} D_{\mathcal{X}}^2 (L_f + \frac{\|\mathcal{A}\|^2}{\beta_t}). \end{aligned} \quad (4.97)$$

Now, using (4.60), we get

$$F_{\beta_t}(\mathbf{x}_t) = f(\mathbf{x}_t) + g_{\beta_t}(\mathcal{A}\mathbf{x}_t) \quad (4.98)$$

$$\leq f(\mathbf{x}_t) + g_{\beta_{t-1}}(\mathcal{A}\mathbf{x}_t) + \frac{\beta_{t-1} - \beta_t}{2} \|\mathbf{y}_{\beta_t}^*(\mathcal{A}\mathbf{x}_t)\|^2 \quad (4.99)$$

$$= F_{\beta_{t-1}}(\mathbf{x}_t) + \frac{\beta_{t-1} - \beta_t}{2} \|\mathbf{y}_{\beta_t}^*(\mathcal{A}\mathbf{x}_t)\|^2. \quad (4.100)$$

## Chapter 4. CGM for Composite Problems

We combine this with (4.97) and subtract  $F_\star$  from both sides to get

$$F_{\beta_t}(\mathbf{x}_{t+1}) - F_\star \leq (1 - \delta\eta_t)(F_{\beta_{t-1}}(\mathbf{x}_t) - F_\star) + \frac{\eta_t^2}{2} D_{\mathcal{X}}^2(L_f + \frac{\|\mathcal{A}\|^2}{\beta_t}) \quad (4.101)$$

$$+ ((1 - \delta\eta_t)(\beta_{t-1} - \beta_t) - \delta\eta_t\beta_t) \frac{1}{2} \|\nabla \mathbf{y}_{\beta_t}^*(\mathcal{A}\mathbf{x}_t)\|^2. \quad (4.102)$$

By choosing  $\eta_t = \frac{2}{\delta t + 2}$  and  $\beta_t = \frac{\beta_0}{\sqrt{\delta(t+1)+1}}$  for some  $\beta_0 > 0$ , we get  $(1 - \delta\eta_t)(\beta_{t-1} - \beta_t) - \delta\eta_t\beta_t < 0$  for any  $t \geq 0$  (with the convention  $\beta_{-1} = \beta_0$ , hence we end up with

$$F_{\beta_t}(\mathbf{x}_{t+1}) - F_\star \leq (1 - \delta\eta_t)(F_{\beta_{t-1}}(\mathbf{x}_t) - F_\star) + \frac{\eta_t^2}{2} D_{\mathcal{X}}^2(L_f + \frac{\|\mathcal{A}\|^2}{\beta_t}). \quad (4.103)$$

Let us call for simplicity  $C := D_{\mathcal{X}}^2(L_f + \frac{\|\mathcal{A}\|^2}{\beta_t})$ , and  $\mathcal{E}_{t+1} := F_{\beta_t}(\mathbf{x}_{t+1}) - F_\star$ . Therefore, we have

$$\mathcal{E}_{t+1} \leq (1 - \delta\eta_t)\mathcal{E}_t + \frac{\eta_t^2}{2} C \quad (4.104)$$

We now show by induction that:

$$E_t \leq 2 \frac{\frac{1}{\delta} C + E_1}{\delta t + 2} \quad (4.105)$$

The base case  $t = 0$  is trivial as  $C > 0$ . Call for simplicity  $K := \delta t + 2$ . Note that  $K \geq 2$ . Under this notation we can write  $\eta_t = \frac{2}{\delta t + 2} = \frac{2}{K}$ . For the induction step, we add a positive term ( $\mathcal{E}_0$  is positive as  $\mathbf{x}_0$  is assumed feasible) to (4.104) and use the induction hypothesis:

$$E_{t+1} \leq (1 - \delta\eta_t)E_t + \frac{\eta_t^2}{2} C + 2\delta \frac{E_1}{K^2} \quad (4.106)$$

$$\leq (1 - \delta \frac{2}{K})E_t + \frac{2}{K^2} C + 2\delta \frac{E_1}{K^2} \quad (4.107)$$

$$\leq (1 - \delta \frac{2}{K}) 2 \frac{\frac{1}{\delta} C + E_1}{K} + \frac{2}{K^2} C + 2\delta \frac{E_1}{K^2} \quad (4.108)$$

$$= (1 - \delta \frac{2}{K}) 2 \frac{\frac{1}{\delta} C + E_1}{K} + 2\delta \left( \frac{\frac{1}{\delta} C}{K^2} + \frac{E_1}{K^2} \right) \quad (4.109)$$

$$= 2 \frac{\frac{1}{\delta} C + E_1}{K} \left( 1 - \delta \frac{2}{K} + \frac{\delta}{K} \right) \quad (4.110)$$

$$= 2 \frac{\frac{1}{\delta} C + E_1}{K} \left( 1 - \frac{\delta}{K} \right) \leq 2 \frac{\frac{1}{\delta} C + E_1}{K + \delta} \quad (4.111)$$

noting that  $K + \delta = \delta(t + 1) + 2$  concludes the proof.

Proof of Theorems 4.8 and 4.9 follows similarly to the proofs of Theorems 4.5 and 4.6.



## 5 CGM via Augmented Lagrangian

In the previous chapter, we introduced a CGM extension with quadratic penalty (HCGM) for solving problems from template (1.4). HCGM is easy to implement with simple and interpretable steps. Unfortunately, many times, it performs poor in practice despite its strong theoretical guarantees.

To this end, we extend HCGM from quadratic penalty to an augmented Lagrangian framework in this chapter. We call the new variant as conditional gradient augmented Lagrangian framework (CGAL). CGAL retains similar guarantees as HCGM but it performs better in practice. We demonstrate the empirical superiority of CGAL against HCGM, as well as some other projection-free algorithms in the literature, in various numerical experiments.

This chapter is based on the joint work with Olivier Fercoq and Volkan Cevher [YFC19].

### Introduction

As we have shown in Theorem 4.6, convergence rate of CGM is bounded above by  $\mathcal{O}(1/\sqrt{t})$ , both in objective residual and feasibility gap. However, in various numerical experiments, we observed that the method exhibits this “worst-case” rate in practice even when we are solving an arguably simple problem instance. Unfortunately, this undesirable behavior strains the practical impact of the method for solving problems at scale. In this chapter, we introduce CGAL, an extension HCGM from quadratic penalty to an augmented Lagrangian framework, in the spirit of [Ber76]. To address this issue, we introduce an extension of HCGM, from quadratic penalty to an augmented Lagrangian framework in the spirit of [Ber76].

### Contributions

- We introduce a new algorithm for solving the problem template (1.1), based on a combination of CGM iterations and an augmented Lagrangian formulation. We show that the proposed method attains  $\mathcal{O}(1/\sqrt{t})$  convergence rate both in objective residual and feasibility gap.

- We identify an implementable step-size rule for the dual updates thanks to the simplicity of our analysis. The cost of evaluating the proposed step-size is negligible. It retains the desirable theoretical convergence rates of HCGM while significantly enhancing the practical performance.
- We demonstrate the empirical superiority of CGAL in various SDP examples, in comparison with HCGM and some other CGM-type methods from the literature.

## 5.1 Augmented Lagrangian Penalty

First, we present a short overview of the augmented Lagrangian penalty and describe its relationship with the Nesterov smoothing and quadratic penalty techniques.

Similar to the quadratic penalty technique, we replace the affine constraint with a "feasibility-promoting" continuous penalty function in the augmented Lagrangian framework. Likewise, this penalty function is parametrized by the penalty parameter  $\beta > 0$ . However, in contrast with the quadratic penalty, we now introduce a "dual variable"  $\mathbf{y} \in \mathbb{R}^d$  (also called as Lagrange multiplier). We can write the augmented Lagrangian penalty function as follows:

$$\min_{\mathbf{v} \in \mathcal{K}} \left\{ \langle \mathbf{y}, \mathcal{A}\mathbf{x} - \mathbf{v} \rangle + \frac{1}{2\beta} \|\mathcal{A}\mathbf{x} - \mathbf{v}\|^2 \right\}. \quad (5.1)$$

We can also view this function as a shifted quadratic penalty, since

$$\arg \min_{\mathbf{x}} f(\mathbf{x}) + \min_{\mathbf{v} \in \mathcal{K}} \left\{ \langle \mathbf{y}, \mathcal{A}\mathbf{x} - \mathbf{v} \rangle + \frac{1}{2\beta} \|\mathcal{A}\mathbf{x} - \mathbf{v}\|^2 \right\} \quad (5.2)$$

$$= \arg \min_{\mathbf{x}} f(\mathbf{x}) + \min_{\mathbf{v} \in \mathcal{K}} \frac{1}{2\beta} \|\mathcal{A}\mathbf{x} - \mathbf{v} + \beta \mathbf{y}\|^2 \quad (5.3)$$

$$= \arg \min_{\mathbf{x}} f(\mathbf{x}) + \frac{1}{2\beta} \text{dist}^2(\mathcal{A}\mathbf{x} + \beta \mathbf{y}, \mathcal{K}). \quad (5.4)$$

Recall our discussion in Section 4.1 on the similarity between the quadratic penalty and Nesterov smoothing techniques. It is not surprising that we can also relate the augmented Lagrangian with Nesterov smoothing. The derivation follows similarly as in Section 4.1.2 for the quadratic penalty, but instead of the simple Euclidean prox-function  $\delta(\mathbf{v}) = \frac{1}{2} \|\mathbf{v}\|^2$ , this time we use a shifted prox-function  $\delta(\mathbf{v}) = \frac{1}{2} \|\mathbf{v} - \mathbf{y}\|^2$ :

$$\iota_{\mathcal{K}\beta}(\mathcal{A}\mathbf{x}) = \max_{\mathbf{z}} \left\{ \min_{\mathbf{v} \in \mathcal{K}} \langle \mathcal{A}\mathbf{x} - \mathbf{v}, \mathbf{z} \rangle - \frac{\beta}{2} \|\mathbf{z} - \mathbf{y}\|^2 \right\} \quad (5.5)$$

$$= \min_{\mathbf{v} \in \mathcal{K}} \max_{\mathbf{z}} \left\{ \langle \mathcal{A}\mathbf{x} - \mathbf{v}, \mathbf{z} \rangle - \frac{\beta}{2} \|\mathbf{z} - \mathbf{y}\|^2 \right\} \quad (5.6)$$

$$= \min_{\mathbf{v} \in \mathcal{K}} \left\{ \langle \mathbf{y}, \mathcal{A}\mathbf{x} - \mathbf{v} \rangle + \frac{1}{2\beta} \|\mathcal{A}\mathbf{x} - \mathbf{v}\|^2 \right\}. \quad (5.7)$$

Sion's minimax theorem [Sio58] justifies the inversion of the order of min and max.

In conclusion, we can say that the augmented Lagrangian formulation is structurally equivalent to a *de facto* instance of the Nesterov smoothing technique, applied to the indicator function of  $\mathcal{K}$ , and with a shifted Euclidean prox-function. Moreover, the dual variable corresponds to the center point of this prox-function.

## 5.2 Conditional Gradient Augmented Lagrangian Method

Consider the augmented Lagrangian of (1.1):

$$\begin{aligned}\mathcal{L}_\beta(\mathbf{x}, \mathbf{y}) &:= f(\mathbf{x}) + \min_{\mathbf{v} \in \mathcal{K}} \left\{ \langle \mathbf{y}, \mathcal{A}\mathbf{x} - \mathbf{v} \rangle + \frac{1}{2\beta} \|\mathcal{A}\mathbf{x} - \mathbf{v}\|^2 \right\} \\ &= f(\mathbf{x}) - \frac{\beta}{2} \|\mathbf{y}\|^2 + \frac{1}{2\beta} \text{dist}^2(\mathcal{A}\mathbf{x} + \beta\mathbf{y}, \mathcal{K}).\end{aligned}\tag{5.8}$$

Clearly  $\mathcal{L}_\beta(\mathbf{x}, \mathbf{y})$  is a smooth and convex function with respect to  $\mathbf{x}$ .

One CGAL iteration is composed of three basic steps:

- ▷ Primal step (conditional gradient step on  $\mathbf{x}$ ),
- ▷ Penalty parameter update (decrement of  $\beta$ ),
- ▷ Dual step (gradient step on  $\mathbf{y}$ ).

### 5.2.1 CGAL Iteration

**Primal step.** CGAL is characterized by a conditional gradient step on the primal variable with respect to  $\mathcal{L}_\beta(\cdot, \mathbf{y})$ . At iteration  $t = 0, 1, \dots$ , we evaluate

$$\nabla_{\mathbf{x}} \mathcal{L}_{\beta_t}(\mathbf{x}_t, \mathbf{y}_t) = \nabla f(\mathbf{x}_t) + \mathcal{A}^* \mathbf{y}_t + \beta_t^{-1} \mathcal{A}^* (\mathcal{A}\mathbf{x}_t - \text{proj}_{\mathcal{K}}(\mathcal{A}\mathbf{x}_t + \beta_t \mathbf{y}_t)).\tag{5.9}$$

Then, we compute the linear minimization oracle

$$\mathbf{h}_t = \text{lmo}_{\mathcal{X}}(\nabla_{\mathbf{x}} \mathcal{L}_{\beta_t}(\mathbf{x}_t, \mathbf{y}_t)).\tag{5.10}$$

Finally, we form the next iterate by a generic conditional gradient step

$$\mathbf{x}_{t+1} = (1 - \eta_t) \mathbf{x}_t + \eta_t \mathbf{h}_t, \quad \text{where } \eta_t = 2/(t+2).\tag{5.11}$$

**Penalty parameter update.** We decrease the penalty parameter in CGAL at a controlled rate:

$$\beta_t = \frac{\beta_0}{\sqrt{t+2}} \quad \text{for some } \beta_0 > 0.\tag{5.12}$$

**Dual step.** After computing  $\mathbf{x}_{t+1}$  and updating the penalty parameter, we update the dual variable as follows:

$$\begin{aligned}\mathbf{y}_{t+1} &= \mathbf{y}_t + \sigma_{t+1} \nabla_{\mathbf{y}} \mathcal{L}_{\beta_{t+1}}(\mathbf{x}_{t+1}, \mathbf{y}_t) \\ &= \mathbf{y}_t + \sigma_{t+1} (\mathcal{A}\mathbf{x}_{t+1} - \text{proj}_{\mathcal{K}}(\mathcal{A}\mathbf{x}_{t+1} + \beta_{t+1}\mathbf{y}_t)).\end{aligned}\tag{5.13}$$

The choice of dual step-size  $\sigma_{t+1}$  has a significant effect on the performance of CGAL. We identify a performant choice as follows: First, define  $\hat{\sigma}_{t+1}$  as

$$\hat{\sigma}_{t+1} := \min\left(\frac{1}{\beta_0}, \frac{(\eta_t \|\mathcal{A}\| D_{\mathcal{X}})^2}{2\beta_{t+1} \|\mathcal{A}\mathbf{x}_{t+1} - \text{proj}_{\mathcal{K}}(\mathcal{A}\mathbf{x}_{t+1} + \beta_{t+1}\mathbf{y}_t)\|^2}\right),\tag{5.14}$$

Then, the recipe for  $\sigma_{t+1}$  is the following:

$$\sigma_{t+1} = \begin{cases} \hat{\sigma}_{t+1} & \text{if } \|\mathbf{y}_t + \hat{\sigma}(\mathcal{A}\mathbf{x}_{t+1} - \text{proj}_{\mathcal{K}}(\mathcal{A}\mathbf{x}_{t+1} + \beta_{t+1}\mathbf{y}_t))\| \leq D_{\mathcal{Y}} \\ 0 & \text{otherwise.} \end{cases}\tag{5.15}$$

Here,  $D_{\mathcal{Y}} > 0$  is the dual search-space parameter to be input by the user.

We emphasize that computing  $\sigma_{t+1}$  does not require an iterative line-search procedure. Instead, we can evaluate it by simple vector operations, hence *the computational cost of finding  $\sigma_{t+1}$  is negligible* in most applications.

### 5.2.2 Guarantees

**Theorem 5.1.** *Sequence  $(\mathbf{x}_t : t = 1, 2, \dots)$  generated by CGAL satisfies:*

$$f(\mathbf{x}_t) - f_{\star} \geq -\|\mathbf{y}_{\star}\| \cdot \text{dist}(\mathcal{A}\mathbf{x}_t, \mathcal{K})\tag{5.16}$$

$$f(\mathbf{x}_t) - f_{\star} \leq 4D_{\mathcal{X}}^2 \left( \frac{L_f}{t+1} + \frac{\|\mathcal{A}\|^2}{\beta_0 \sqrt{t+1}} + \frac{2\|\mathcal{A}\|^2}{\beta_0(t+1)^{1.5}} \right) + \frac{\beta_0 D_{\mathcal{Y}}^2}{2\sqrt{t+1}}\tag{5.17}$$

$$\text{dist}(\mathcal{A}\mathbf{x}_t, \mathcal{K}) \leq \frac{1}{\sqrt{t+1}} \left( \beta_0 (D_{\mathcal{Y}} + 2\|\mathbf{y}_t - \mathbf{y}_{\star}\|) + D_{\mathcal{X}} \sqrt{\frac{8L_f \beta_0}{\sqrt{t+1}} + \left(8 + \frac{1}{t+1}\right) \|\mathcal{A}\|^2} \right).\tag{5.18}$$

Considering these inequalities, we suggest choosing  $D_{\mathcal{Y}}$  proportional to  $D_{\mathcal{X}} \|\mathcal{A}\| / \beta_0$ .

**Remark.** Similar to HCGM, we can extend CGAL for composite problems with Lipschitz continuous regularizers. In this case, dual update of CGAL corresponds to updating the center point of the prox-function. We omit the details.

**Remark.** We omit design variants of CGAL with line-search and fully corrective updates.

### 5.2.3 Application to the SDP template

Similar to HCGM, we can apply CGAL for solving instances of the SDP template (1.4). We can compute the directional derivative of the augmented Lagrangian function by

$$\nabla_X \mathcal{L}_{\beta_t}(\mathbf{X}_t, \mathbf{y}_t) = \mathbf{C} + \mathcal{A}^* \mathbf{y}_t + \beta_t^{-1} \mathcal{A}^* (\mathcal{A} \mathbf{X} - \mathbf{b}). \quad (5.19)$$

Then, we can evaluate the linear minimization oracle by using `MinEigVec`:

$$\mathbf{H}_t = \alpha \mathbf{u}_t \mathbf{u}_t^* \quad \text{where} \quad (\mathbf{u}_t, \beta_t^{-1}) = \text{MinEigVec}(\mathbf{C} + \mathcal{A}^* \mathbf{y}_t + \beta_t^{-1} \mathcal{A}^* (\mathcal{A} \mathbf{X} - \mathbf{b})). \quad (5.20)$$

Finally, since  $\text{proj}_{\mathcal{K}}(\cdot) = \mathbf{b}$  for all inputs, the dual update of CGAL takes the following form:

$$\mathbf{y}_{t+1} = \mathbf{y}_t + \sigma_{t+1} (\mathcal{A} \mathbf{X}_{t+1} - \mathbf{b}), \quad (5.21)$$

where the dual step size is chosen as

$$\sigma_{t+1} = \begin{cases} \hat{\sigma}_{t+1} & \text{if } \|\mathbf{y}_t + \hat{\sigma}(\mathcal{A} \mathbf{X}_{t+1} - \mathbf{b})\| \leq D_Y \\ 0 & \text{otherwise} \end{cases} \quad \text{where} \quad \hat{\sigma}_{t+1} = \min\left(\frac{1}{\beta_0}, \frac{(\eta_t \|\mathcal{A}\|_{D_X})^2}{2\beta_{t+1} \|\mathcal{A} \mathbf{X}_{t+1} - \mathbf{b}\|^2}\right) \quad (5.22)$$

for some  $D_Y > 0$  to be specified by the user. Algorithm 5.1 shows the complete pseudocode.

---

**Algorithm 5.1** CGAL for the standard SDP formulation (1.4)

---

**Require:** Data for (1.4); initial penalty parameter  $\beta_0$ ; dual search-space parameter  $D_Y$

```

1: function CGAL
2:    $\mathbf{X} \leftarrow \mathbf{0}$  and  $\mathbf{y} \leftarrow \mathbf{0}$ 
3:   for  $t \leftarrow 0, 1, 2, \dots$  do
4:      $\beta \leftarrow \beta_0 / \sqrt{t+2}$ 
5:      $(\mathbf{u}, \lambda) \leftarrow \text{MinEigVec}(\mathbf{C} + \mathcal{A}^* \mathbf{y} + \beta^{-1} \mathcal{A}^* (\mathcal{A} \mathbf{X} - \mathbf{b}))$ 
6:      $\mathbf{H} \leftarrow \alpha \mathbf{u} \mathbf{u}^*$ 
7:      $\eta \leftarrow 2/(t+2)$ 
8:      $\mathbf{X} \leftarrow (1-\eta)\mathbf{X} + \eta \mathbf{H}$ 
9:      $\beta_+ \leftarrow \beta_0 / \sqrt{t+3}$ 
10:     $\sigma \leftarrow \text{DUALSTEPSize}(\mathbf{X}, \mathbf{y}, \beta_0, \beta_+, \eta)$ 
11:     $\mathbf{y} \leftarrow (1-\eta)\mathbf{y} + \sigma(\mathcal{A} \mathbf{X} - \mathbf{b})$ 
12:  end for
13:  return  $\mathbf{X}$ 
14: end function

15: function DUALSTEPSize( $\mathbf{X}, \mathbf{y}, \beta_0, \beta_+, \eta$ )
16:    $\sigma \leftarrow \min(\frac{(\eta \|\mathcal{A}\|_{D_X})^2}{2\beta_+ \|\mathcal{A} \mathbf{X} - \mathbf{b}\|^2}, 1/\beta_0)$ 
17:   if  $\|\mathbf{y} + \sigma(\mathcal{A} \mathbf{X} - \mathbf{b})\| > D_Y$  then  $\sigma \leftarrow 0$  ▷ No dual step if  $\|\mathbf{y}\|$  is growing too much
18:   end if
19:   return  $\sigma$ 
20: end function

```

---

### 5.3 Numerical Experiments

In this section, we compare the empirical performance of CGAL against some existing methods in the literature in various SDP examples. Let us first provide a brief overview of the baseline methods that we will use in the comparisons.

#### Baseline methods

- **HCGM** that we introduced in [YFLC18] and described in Chapter 4, is the predecessor of CGAL. It has  $\mathcal{O}(1/\sqrt{t})$  rate guarantee both in the objective residual and feasibility gap.

- **UPD** and **AUPD**, the universal primal-dual methods that we introduced in [YTDC15a] and described in Chapter 3: These methods are based on a primal-dual (sub)gradient formulation. Remark that the main template of UPD and AUPD is different from the template of CGAL. For example, they do not require the smoothness of  $f$ . Instead, they assume Hölder smoothness in the dual formulation. Nevertheless, both methods approaches work for the model SDP formulation (1.4). UPD and AUPD do not directly work with  $\text{lmo}$ . However, their sharp-operator can be cast as an instance of  $\text{lmo}$  for the SDP template (1.4).

These methods adopt an inexact line-search strategy by Nesterov [Nes15], which requires the target accuracy  $\varepsilon$  as an input parameter. Moreover, they are guaranteed to converge only up to  $\varepsilon/2$  accuracy, *i.e.*,  $f(\mathbf{x}) - f_\star \leq \mathcal{O}(1/\sqrt{t}) + \varepsilon/2$ .

- **FW-AL** is a CGM-based augmented Lagrangian framework [GPLJ18]. Similar to CGAL, this method is characterized by one CGM step on  $\mathcal{L}_\beta(\cdot, \mathbf{y}_t)$  followed by one dual gradient ascent step on  $\mathcal{L}_\beta(\mathbf{x}_{t+1}, \cdot)$ . In contrast to CGAL, the penalty parameter of the FW-AL is kept constant throughout the optimization process. Originally, FW-AL is proposed for  $\mathcal{A}\mathbf{x} = \mathbf{0}$  type of constraints, but it can be applied to  $\mathcal{A}\mathbf{x} = \mathbf{b}$  case by using a simple product space technique. The analysis of the FW-AL relies on the error bounds (see Theorem 1 in [GPLJ18] for the conditions, and [BNPS17] for more details about error bounds). Their dual step-size  $\sigma_{t+1}$  depends on the error bound constant  $\alpha$ , as  $\sigma_{t+1} = \frac{2\sigma_0}{t+2}$  with  $\sigma_0 \leq \min\{\frac{2}{\lambda}, \frac{\alpha^2}{2\delta}\}$ . As a results,  $\sigma_0$  comes into algorithm as a tuning parameter, and the method has guarantees only if  $\sigma_0$  is chosen small enough.

- **IAL** is an inexact augmented Lagrangian method, where the Lagrangian subproblems are approximately solved by using CGM up to a prescribed accuracy [LLM19]. With theoretical analysis, the error tolerance for this subproblem is determined as  $\varepsilon_t = \varepsilon_0/t$  for some  $\varepsilon_0 > 0$ . IAL has a double-loop structure, where each outer iteration (or epoch) executes multiple iterations CGM, followed with a dual step.

IAL provably generates an  $\varepsilon$ -solution after  $\mathcal{O}(1/\varepsilon^2)$  outer iterations, by choosing the penalty parameter  $\beta$  appropriately (proportional to  $\sqrt{\varepsilon}$ ). This method, however, requires multiple calls of  $\text{lmo}$  at each iteration. The number of  $\text{lmo}$  calls to achieve the prescribed accuracy in the subproblem is bounded by  $\mathcal{O}(1/\varepsilon_t)$  (see Theorem 2.2 in [LLM19]). Hence, the overall  $\text{lmo}$  complexity of this method becomes  $\mathcal{O}(1/\varepsilon^4)$ .

### 5.3.1 Max-cut

Maximum cut is an NP-Hard combinatorial problem from computer science. Denoting the graph Laplacian matrix by  $\mathbf{L}$ , an SDP relaxation of this problem can be formulated as [GW95]:

$$\text{minimize } -\frac{1}{4}\langle \mathbf{L}, \mathbf{X} \rangle \quad \text{subject to} \quad \text{diag}(\mathbf{X}) = \mathbf{1}, \quad \text{Tr } \mathbf{X} = n, \quad \mathbf{X} \in \mathbb{S}_+^n. \quad (5.23)$$

Tuning all baseline methods requires substantial computational effort, especially because some of these methods have multiple tuning parameters. To this end, we first consider a small scale max-cut instance where we compare all baseline methods. In this setup we use GD97\_b dataset<sup>1</sup>, which corresponds to a  $47 \times 47$  dimensional problem. In Figure 5.1, we present the performance of each method with the best parameter choice obtained after an extensive search. We also provide the performance of each algorithm at all trials at the end of this chapter in Figures 5.5 and 5.6.

Next, we consider a medium scale experiment, where we compare CGAL, HCGM and UPD for max-cut with G1 ( $800 \times 800$ ) and G40 ( $2000 \times 2000$ ) datasets<sup>2</sup>. We compile the results in Figure 5.2. Observe that HCGM converges with  $\mathcal{O}(1/\sqrt{t})$  (which is the worst case bound) while CGAL achieves a faster rate.

### 5.3.2 k-means Clustering

Consider the SDP formulation of the model-free k-means clustering problem [PW07]:

$$\text{minimize } \langle \mathbf{C}, \mathbf{X} \rangle \quad \text{subject to} \quad \underbrace{\mathbf{X}\mathbf{1} = \mathbf{1}, \quad \mathbf{X} \geq 0}_{\mathcal{AX} \in \mathcal{K}}, \quad \underbrace{\mathbf{X} \in \mathbb{S}_+^n, \quad \text{Tr } \mathbf{X} = k}_{\mathbf{X} \in \mathcal{X}}. \quad (5.24)$$

$\mathbf{C}$  is the Euclidean distance matrix and  $k$  is the number of clusters. We denote the vector of ones by  $\mathbf{1}$ , hence  $\mathbf{X}\mathbf{1} = \mathbf{1}$  and  $\mathbf{X} \geq 0$  together implies that each row of  $\mathbf{X}$  is on the unit simplex. The same applies for columns of  $\mathbf{X}$  due to symmetry.

We use the problem setup from [MVW17] that we have also considered in Section 4.4.1. This setup contains a  $1000 \times 1000$  dimensional dataset generated by sampling and preprocessing the MNIST dataset<sup>3</sup> using a one-layer neural network. Further details on this setup and the dataset can be found in [MVW17]. In Figure 5.3, we observe once again that CGAL outperforms HCGM, achieving  $\mathcal{O}(1/t)$  empirical convergence rate.

<sup>1</sup>V. Batagelj and A. Mrvar. Pajek datasets. Available at "<http://vlado.fmf.uni-lj.si/pub/networks/data/>"

<sup>2</sup>Y. Ye. Gset random graphs. Available at "<https://www.cise.ufl.edu/research/sparse/matrices/gset/>"

<sup>3</sup>LeCun & Cortes. MNIST handwritten digit database. Available at "<http://yann.lecun.com/exdb/mnist/>"

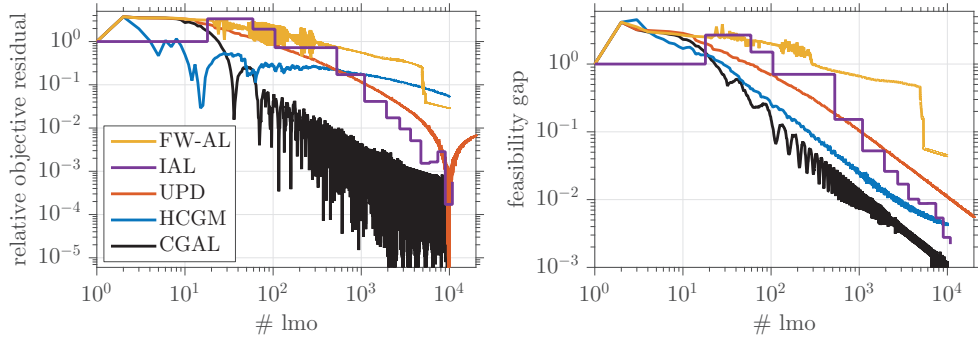


Figure 5.1 – Empirical performance of various methods for solving max-cut problem.

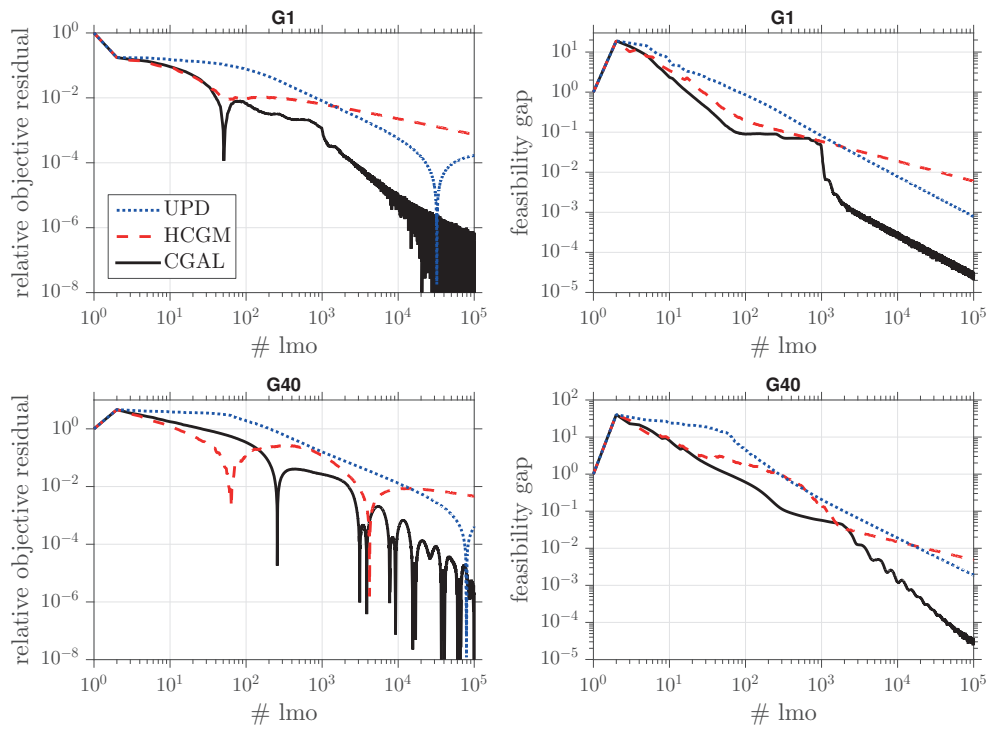


Figure 5.2 – Empirical comparison of UPD, HCGM and CGAL for max-cut problem.

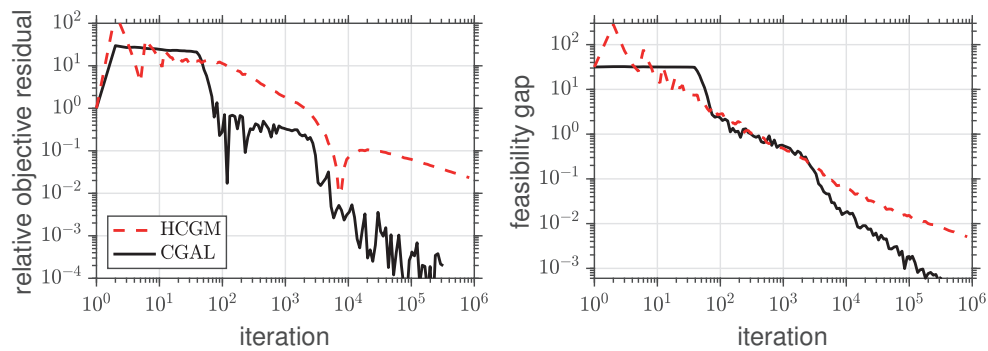


Figure 5.3 – Comparison of CGAL and HCGM for clustering SDP.



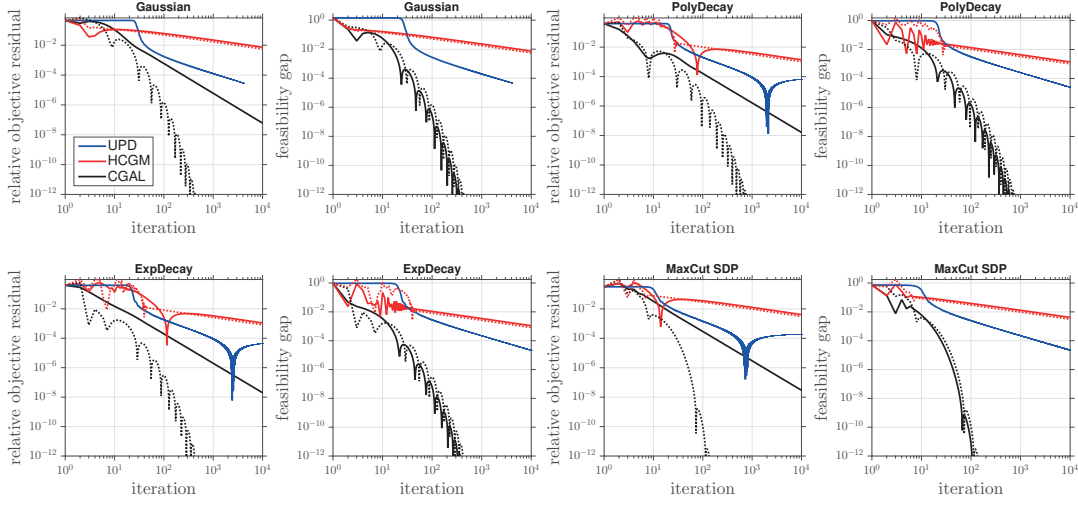


Figure 5.4 – Empirical comparison of CGAL, HCGM and UPD for solving generalized eigenvector problem with 4 different synthetic setups. Dotted lines present objective residual and feasibility gap of the atoms chosen by linear minimization oracle ( $\mathbf{h}_t$ ).

### 5.3.3 Generalized Eigenvector Problem

Consider the SDP relaxation of the generalized eigenvector problem [BVB18]:

$$\text{maximize } \langle \mathbf{C}, \mathbf{X} \rangle \quad \text{subject to } \langle \mathbf{D}, \mathbf{X} \rangle = 1, \quad \mathbf{X} \in \mathbb{S}_+^n, \quad \text{Tr } \mathbf{X} \leq \alpha. \quad (5.25)$$

where  $\mathbf{C}$  and  $\mathbf{D}$  are symmetric matrices and  $\alpha > 0$  is a model parameter. We synthetically generate  $\mathbf{D}$  with *iid* Gaussian entries and consider 4 different cases for  $\mathbf{C}$ :

- Gaussian -  $\mathbf{C}$  generated by taking symmetric part of  $10^3 \times 10^3$  *iid* Gaussian matrix
- PolyDecay -  $\mathbf{C}$  generated by randomly rotating  $\text{diag}(1^{-i}, 2^{-i}, \dots, 1000^{-i})$  ( $i = 1$ )
- ExpDecay -  $\mathbf{C}$  generated by randomly rotating  $\text{diag}(10^{-i}, 10^{-2i}, \dots, 10^{-1000i})$  ( $i = 0.025$ )
- MaxCut SDP -  $\mathbf{C}$  is a solution of the max-cut SDP with G40 dataset ( $2000 \times 2000$ )

This problem highlights an important observation that partially explains the reason why CGAL outperforms HCGM. First, note that this problem has a rank-1 solution by design if we set  $\alpha$  right. In this scenario, if the problem formulation is well-conditioned, we could expect lmo to pick this solution (or some close points). CGAL updates the dual variable (which corresponds to the center point of a quadratic penalty) and adapts better to the geometry. In Figure 5.4, we provide an empirical evidence of this adaptation: Dotted lines correspond to extreme points chosen by lmo. These points ( $\mathbf{h}_t$ ) converge (quickly with linear rates) to a solution for CGAL. On the other hand, we do not observe the same behavior for HCGM, because the condition number grows too much and lmo fails to pick good atoms. We omit lmo outputs of UPD in the figure because it does not converge to a solution.

### Observations & Concluding Remarks

**CGAL** outperformed the alternative methods in all of our experiments. **HCGM** performed with the guaranteed "worst-case" rate in our experiments.

**UPD** and **AUPD** requires the target suboptimality level  $\varepsilon$  as an input parameter, and these methods have convergence guarantees only up to this target accuracy. Indeed, we can observe the saturation effect in the objective residual in our numerical experiments.

In our experiments, we observed similar performance from UPD and AUPD when solving max-cut problems. Remark that both methods have similar guarantees for this problem, since the dual subgradient is Hölder continuous of order  $\nu = 0$ . Surprisingly, AUPD did not exhibit the saturation in practice. We can explain this as follows: Both UPD and AUPD are based on the inexact line-search procedure, but UPD lets a constant  $\frac{\varepsilon}{2}$  slackness term in the line-search condition, while AUPD has a more stringent slackness term  $\approx \frac{\varepsilon}{2t}$ . This is required from the theoretical perspective, in order to prevent error accumulation due to acceleration. However, the actual error accumulation might be much less than the worst-case scenario. As a result, by decreasing the slackness term we prevent the saturation.

Remark that, we can also design a non-saturating variant of UPD, by decreasing the slackness term at a controlled rate. This rate can be characterized from the theory. Note however, as we decrease it, we would also need more accurate evaluations of the sharp operator. In large-scale problems, where we use inexact oracles, this causes stability problems in the line-search procedure. Under noise, our line-search condition might become ill-defined, and the method can get stuck in the line-search procedure.

**FW-AL** iterations are very similar to CGAL, but the empirical behavior of the algorithm is different. The analysis of FW-AL relies on the error bound parameter, which is proved to be positive (assuming that Slater's condition holds). However, this parameter is typically unknown a priori, and we argue that it can be arbitrarily small. Hence, it might be difficult to tune in practice FW-AL for some problems.

**IAL** performed better in practice, in comparison with the  $\mathcal{O}(1/\varepsilon^4)$  lmo complexity bound. Nevertheless, the method is arguably difficult to tune, as it involves multiple tuning parameters. IAL has a double-loop structure, and only the outer iterates are informative. This results in stair-like plots.

### 5.3. Numerical Experiments

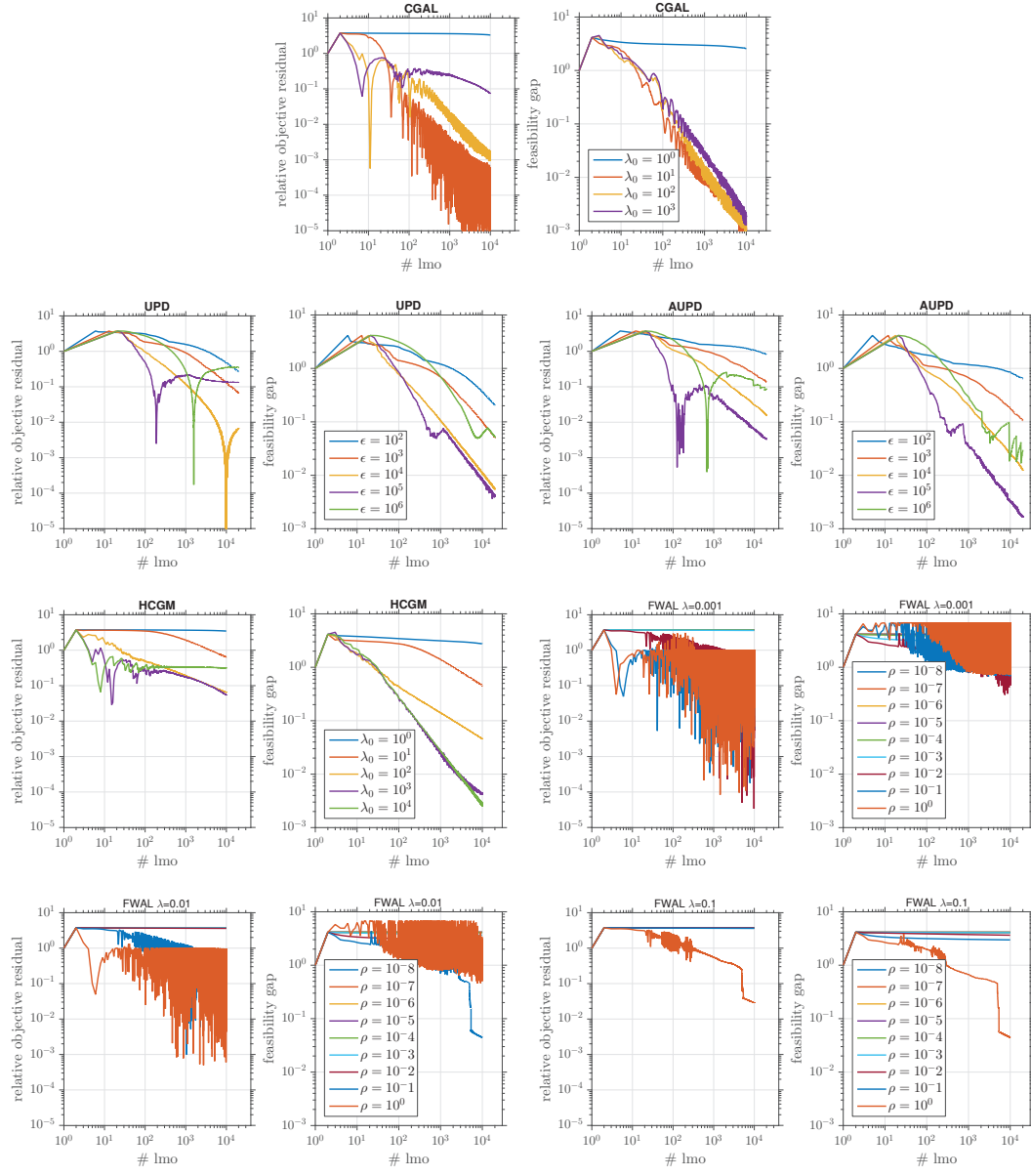


Figure 5.5 – Comparison of projection-free methods for max-cut SDP [Part 1].

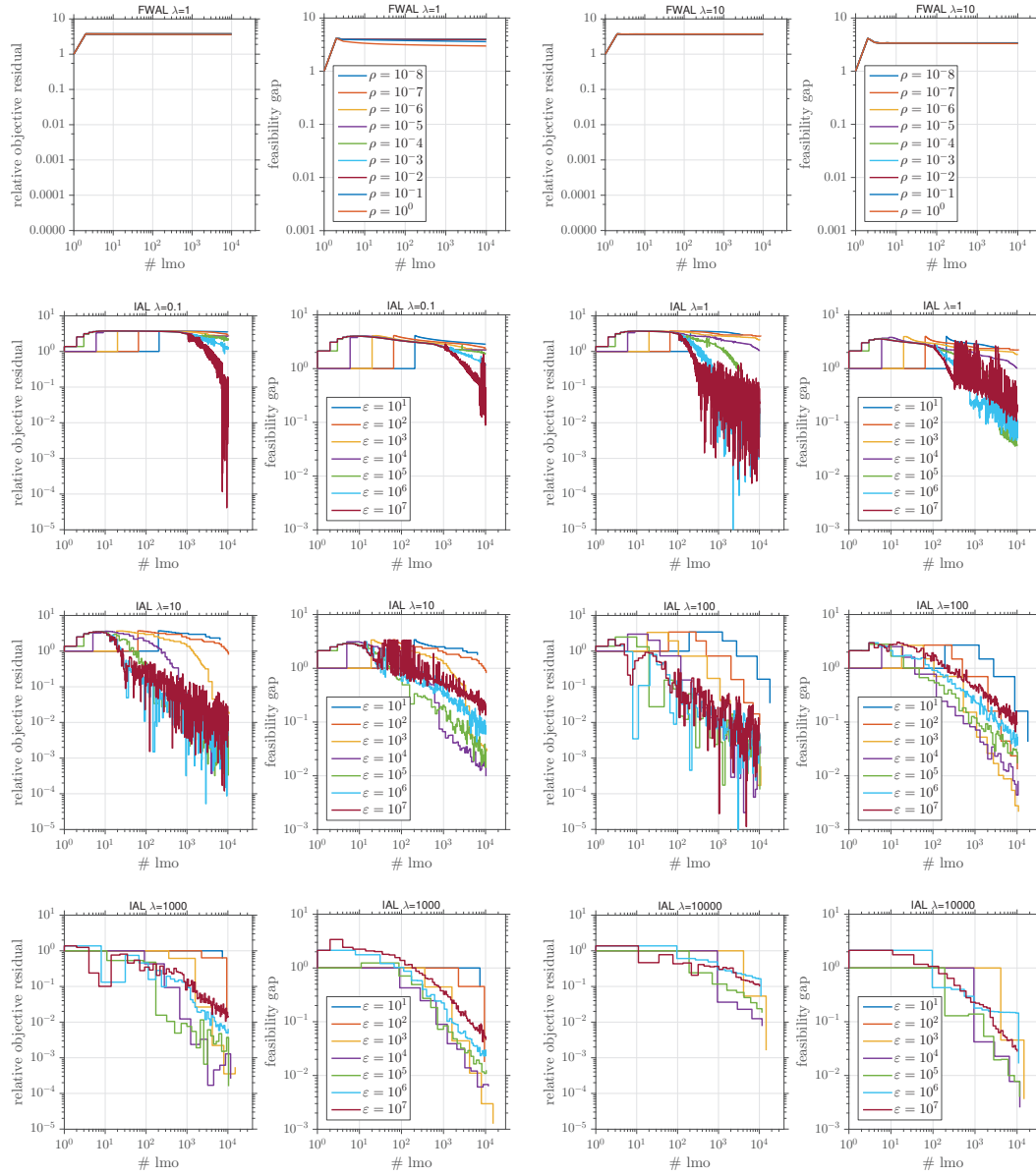


Figure 5.6 – Comparison of the projection-free methods for max-cut SDP [Part 2].

## 5.4 Appendix: Proofs

We formulate the augmented Lagrangian function with three variables, including the slack variable  $\mathbf{v} = \mathcal{A}\mathbf{x} \in \mathcal{K}$ , as follows:

$$\mathcal{L}_\beta(\mathbf{x}, \mathbf{v}, \mathbf{y}) := f(\mathbf{x}) + \langle \mathbf{y}, \mathcal{A}\mathbf{x} - \mathbf{v} \rangle + \frac{1}{2\beta} \|\mathcal{A}\mathbf{x} - \mathbf{v}\|^2, \quad (5.26)$$

where  $\mathbf{y} \in \mathbb{R}^d$  is the Lagrange multiplier and  $\beta > 0$  is the penalty parameter.

We can write the directional derivatives of the augmented Lagrangian function as

$$\nabla_{\mathbf{x}} \mathcal{L}_\beta(\mathbf{x}, \mathbf{v}, \mathbf{y}) = \nabla f(\mathbf{x}) + \mathcal{A}^* \mathbf{y} + \beta^{-1} \mathcal{A}^* (\mathcal{A}\mathbf{x} - \mathbf{v}) \quad (5.27)$$

$$\nabla_{\mathbf{v}} \mathcal{L}_\beta(\mathbf{x}, \mathbf{v}, \mathbf{y}) = -\mathbf{y} - \beta^{-1} (\mathcal{A}\mathbf{x} - \mathbf{v}) \quad (5.28)$$

$$\nabla_{\mathbf{y}} \mathcal{L}_\beta(\mathbf{x}, \mathbf{v}, \mathbf{y}) = \mathcal{A}\mathbf{x} - \mathbf{v}. \quad (5.29)$$

Clearly,  $L_{\beta_t}$  is  $\bar{L}_t$ -smooth, where we denote by  $\bar{L}_t := (L_f + \frac{\|\mathcal{A}\|^2}{\beta_t})$ . Now, let us define

$$\mathbf{v}_t := \text{proj}_{\mathcal{K}}(\mathcal{A}\mathbf{x}_t + \beta_t \mathbf{y}_t) = \arg \min_{\mathbf{v} \in \mathcal{K}} \mathcal{L}_{\beta_t}(\mathbf{x}_t, \mathbf{v}, \mathbf{y}_t). \quad (5.30)$$

Then, using the Taylor expansion, we get the following estimate:

$$\begin{aligned} \mathcal{L}_{\beta_{t+1}}(\mathbf{x}_{t+1}, \mathbf{v}_t, \mathbf{y}_t) &\leq \mathcal{L}_{\beta_{t+1}}(\mathbf{x}_t, \mathbf{v}_t, \mathbf{y}_t) + \langle \nabla_{\mathbf{x}} \mathcal{L}_{\beta_{t+1}}(\mathbf{x}_t, \mathbf{v}_t, \mathbf{y}_t), \mathbf{x}_{t+1} - \mathbf{x}_t \rangle + \frac{\bar{L}_{t+1}}{2} \|\mathbf{x}_{t+1} - \mathbf{x}_t\|^2 \\ &= \mathcal{L}_{\beta_{t+1}}(\mathbf{x}_t, \mathbf{v}_t, \mathbf{y}_t) + \eta_t \langle \nabla_{\mathbf{x}} \mathcal{L}_{\beta_{t+1}}(\mathbf{x}_t, \mathbf{v}_t, \mathbf{y}_t), \mathbf{h}_t - \mathbf{x}_t \rangle + \eta_t^2 \frac{\bar{L}_{t+1}}{2} \|\mathbf{h}_t - \mathbf{x}_t\|^2 \end{aligned} \quad (5.31)$$

$$\begin{aligned} &\leq \mathcal{L}_{\beta_t}(\mathbf{x}_t, \mathbf{v}_t, \mathbf{y}_t) + \eta_t \langle \nabla_{\mathbf{x}} \mathcal{L}_{\beta_t}(\mathbf{x}_t, \mathbf{v}_t, \mathbf{y}_t), \mathbf{h}_t - \mathbf{x}_t \rangle + \eta_t^2 \frac{\bar{L}_{t+1}}{2} D_{\mathcal{X}}^2 \\ &\quad + \frac{1}{2} (\beta_{t+1}^{-1} - \beta_t^{-1}) \|\mathcal{A}\mathbf{x}_t - \mathbf{v}_t\|^2 + \eta_t (\beta_{t+1}^{-1} - \beta_t^{-1}) \langle \mathcal{A}\mathbf{x}_t - \mathbf{v}_t, \mathcal{A}\mathbf{h}_t - \mathcal{A}\mathbf{x}_t \rangle. \end{aligned} \quad (5.32)$$

Now, we can bound the inner product term as follows:

$$\begin{aligned} &\langle \nabla_{\mathbf{x}} \mathcal{L}_{\beta_t}(\mathbf{x}_t, \mathbf{v}_t, \mathbf{y}_t), \mathbf{h}_t - \mathbf{x}_t \rangle \\ &\leq \langle \nabla_{\mathbf{x}} \mathcal{L}_{\beta_t}(\mathbf{x}_t, \mathbf{v}_t, \mathbf{y}_t), \mathbf{x}_\star - \mathbf{x}_t \rangle \end{aligned} \quad (5.33)$$

$$= \langle \nabla f(\mathbf{x}_t), \mathbf{x}_\star - \mathbf{x}_t \rangle + \langle \mathbf{y}_t + \beta_t^{-1} (\mathcal{A}\mathbf{x}_t - \mathbf{v}_t), \mathcal{A}\mathbf{x}_\star - \mathcal{A}\mathbf{x}_t + \mathbf{v}_t - \mathbf{v}_t \rangle \quad (5.34)$$

$$\begin{aligned} &= \langle \nabla f(\mathbf{x}_t), \mathbf{x}_\star - \mathbf{x}_t \rangle - \beta_t^{-1} \|\mathcal{A}\mathbf{x}_t - \mathbf{v}_t\|^2 - \langle \mathbf{y}_t, \mathcal{A}\mathbf{x}_t - \mathbf{v}_t \rangle \\ &\quad + \langle \mathbf{y}_t + \beta_t^{-1} (\mathcal{A}\mathbf{x}_t - \mathbf{v}_t), \mathcal{A}\mathbf{x}_\star - \mathbf{v}_t \rangle \end{aligned} \quad (5.35)$$

$$\begin{aligned} &\leq f_\star - f(\mathbf{x}_t) - \beta_t^{-1} \|\mathcal{A}\mathbf{x}_t - \mathbf{v}_t\|^2 - \langle \mathbf{y}_t, \mathcal{A}\mathbf{x}_t - \mathbf{v}_t \rangle \\ &\quad + \langle \mathbf{y}_t + \beta_t^{-1} (\mathcal{A}\mathbf{x}_t - \mathbf{v}_t), \mathcal{A}\mathbf{x}_\star - \mathbf{v}_t \rangle \end{aligned} \quad (5.36)$$

$$= \mathcal{L}_\star - \mathcal{L}_{\beta_t}(\mathbf{x}_t, \mathbf{v}_t, \mathbf{y}_t) - \frac{1}{2\beta_t} \|\mathcal{A}\mathbf{x}_t - \mathbf{v}_t\|^2 + \langle \mathbf{y}_t + \beta_t^{-1} (\mathcal{A}\mathbf{x}_t - \mathbf{v}_t), \mathcal{A}\mathbf{x}_\star - \mathbf{v}_t \rangle \quad (5.37)$$

$$\leq \mathcal{L}_\star - \mathcal{L}_{\beta_t}(\mathbf{x}_t, \mathbf{v}_t, \mathbf{y}_t) - \frac{1}{2\beta_t} \|\mathcal{A}\mathbf{x}_t - \mathbf{v}_t\|^2. \quad (5.38)$$

## Chapter 5. CGM via Augmented Lagrangian

The first inequality holds because  $\mathbf{h}_t$  is the solution of the linear minimization oracle. The second inequality simply follows from the convexity of  $f$ . Finally, the last inequality comes from the optimality condition. We also use the strong duality assumption to ensure  $f_\star = \mathcal{L}_\star$ .

By definition of  $\mathbf{v}_t$  in (5.30), the following holds:

$$\langle \mathbf{y}_t + \beta_t^{-1}(\mathcal{A}\mathbf{x}_t - \mathbf{v}_t), \mathbf{v} - \mathbf{v}_t \rangle = -\langle \nabla_{\mathbf{v}} \mathcal{L}_{\beta_t}(\mathbf{x}_t, \mathbf{v}_t, \mathbf{y}_t), \mathbf{v} - \mathbf{v}_t \rangle \leq 0, \quad \forall \mathbf{v} \in \mathcal{K}, \quad (5.39)$$

and in particular for  $\mathbf{v} = \mathcal{A}\mathbf{x}_\star \in \mathcal{K}$ .

Substituting (5.38) back into (5.32), we get

$$\begin{aligned} \mathcal{L}_{\beta_{t+1}}(\mathbf{x}_{t+1}, \mathbf{v}_t, \mathbf{y}_t) &\leq (1 - \eta_t) \mathcal{L}_{\beta_t}(\mathbf{x}_t, \mathbf{v}_t, \mathbf{y}_t) + \eta_t \mathcal{L}_\star + (\beta_{t+1}^{-1} - (1 + \eta_t)\beta_t^{-1}) \frac{1}{2} \|\mathcal{A}\mathbf{x}_t - \mathbf{v}_t\|^2 \\ &\quad + \eta_t^2 \frac{\bar{L}_{t+1}}{2} D_{\mathcal{X}}^2 + \eta_t (\beta_{t+1}^{-1} - \beta_t^{-1}) \langle \mathcal{A}\mathbf{x}_t - \mathbf{v}_t, \mathcal{A}\mathbf{h}_t - \mathcal{A}\mathbf{x}_t \rangle. \end{aligned} \quad (5.40)$$

We define  $\bar{\mathbf{v}}_{t+1}$  as

$$\bar{\mathbf{v}}_{t+1} := \text{proj}_{\mathcal{K}}(\mathcal{A}\mathbf{x}_{t+1} + \beta_{t+1}^{-1}\mathbf{y}_t) = \arg\min_{\mathbf{v} \in \mathcal{K}} \mathcal{L}_{\beta_{t+1}}(\mathbf{x}_{t+1}, \mathbf{v}, \mathbf{y}_t). \quad (5.41)$$

Hence, we obtain

$$\mathcal{L}_{\beta_{t+1}}(\mathbf{x}_{t+1}, \bar{\mathbf{v}}_{t+1}, \mathbf{y}_t) \leq \mathcal{L}_{\beta_{t+1}}(\mathbf{x}_{t+1}, \mathbf{v}_t, \mathbf{y}_t) \quad (5.42)$$

$$\begin{aligned} &\leq (1 - \eta_t) \mathcal{L}_{\beta_t}(\mathbf{x}_t, \mathbf{v}_t, \mathbf{y}_t) + \eta_t \mathcal{L}_\star + (\beta_{t+1}^{-1} - (1 + \eta_t)\beta_t^{-1}) \frac{1}{2} \|\mathcal{A}\mathbf{x}_t - \mathbf{v}_t\|^2 \\ &\quad + \eta_t^2 \frac{\bar{L}_{t+1}}{2} D_{\mathcal{X}}^2 + \eta_t (\beta_{t+1}^{-1} - \beta_t^{-1}) \langle \mathcal{A}\mathbf{x}_t - \mathbf{v}_t, \mathcal{A}\mathbf{h}_t - \mathcal{A}\mathbf{x}_t \rangle \end{aligned} \quad (5.43)$$

$$\begin{aligned} &\leq (1 - \eta_t) \mathcal{L}_{\beta_t}(\mathbf{x}_t, \bar{\mathbf{v}}_t, \mathbf{y}_t) + \eta_t \mathcal{L}_\star + (\beta_{t+1}^{-1} - (1 + \eta_t)\beta_t^{-1}) \frac{1}{2} \|\mathcal{A}\mathbf{x}_t - \mathbf{v}_t\|^2 \\ &\quad + \eta_t^2 \frac{\bar{L}_{t+1}}{2} D_{\mathcal{X}}^2 + \eta_t (\beta_{t+1}^{-1} - \beta_t^{-1}) \|\mathcal{A}\mathbf{x}_t - \mathbf{v}_t\| \|\mathcal{A}\| D_{\mathcal{X}}, \end{aligned} \quad (5.44)$$

where the last line follows from the Cauchy-Schwarz inequality.

Now, we use the following bound:

$$\begin{aligned} &(\beta_{t+1}^{-1} - (1 + \eta_t)\beta_t^{-1}) \frac{1}{2} \|\mathcal{A}\mathbf{x}_t - \mathbf{v}_t\|^2 + \eta_t (\beta_{t+1}^{-1} - \beta_t^{-1}) \|\mathcal{A}\mathbf{x}_t - \mathbf{v}_t\| \|\mathcal{A}\| D_{\mathcal{X}} \\ &\leq \frac{\eta_t^2 (\beta_{t+1}^{-1} - \beta_t^{-1})^2}{2((1 + \eta_t)\beta_t^{-1} - \beta_{t+1}^{-1})} \|\mathcal{A}\|^2 D_{\mathcal{X}}^2, \end{aligned} \quad (5.45)$$

which is obtained by considering the maximum value of the second order polynomial in terms of  $\|\mathcal{A}\mathbf{x}_t - \mathbf{v}_t\|$ .

Remark that,

$$\begin{aligned} \frac{\eta_t^2(\beta_{t+1}^{-1} - \beta_t^{-1})^2}{2((1 + \eta_t)\beta_t^{-1} - \beta_{t+1}^{-1})} &\leq \beta_0^{-1} \frac{\frac{2}{(t+2)^2} (\sqrt{t+3} - \sqrt{t+2})^2}{(1 + \frac{2}{t+2})\sqrt{t+2} - \sqrt{t+3}} \leq \beta_0^{-1} \frac{\frac{1}{2(t+2)^3}}{(1 + \frac{2}{t+2})\sqrt{t+2} - \sqrt{t+3}} \\ &\leq \beta_0^{-1} \frac{\frac{1}{2(t+2)^{2.5}}}{(1 + \frac{2}{t+2})(t+2) - (t+2)} \leq \frac{1}{4\beta_0(t+2)^{2.5}}. \end{aligned} \quad (5.46)$$

Combining all these bounds and subtracting  $\mathcal{L}_\star$  from both sides, we end up with

$$\mathcal{L}_{\beta_{t+1}}(\mathbf{x}_{t+1}, \bar{\mathbf{v}}_{t+1}, \mathbf{y}_t) - \mathcal{L}_\star \leq (1 - \eta_t)(\mathcal{L}_{\beta_t}(\mathbf{x}_t, \mathbf{v}_t, \mathbf{y}_t) - \mathcal{L}_\star) + \eta_t^2 \frac{\bar{L}_{t+1}}{2} D_{\mathcal{X}}^2 + \frac{\|\mathcal{A}\|^2 D_{\mathcal{X}}^2}{4\beta_0(t+2)^{2.5}}. \quad (5.47)$$

In order to obtain a recurrence relation, we use

$$\mathcal{L}_{\beta_{t+1}}(\mathbf{x}_{t+1}, \mathbf{v}_{t+1}, \mathbf{y}_{t+1}) \leq \mathcal{L}_{\beta_{t+1}}(\mathbf{x}_{t+1}, \bar{\mathbf{v}}_{t+1}, \mathbf{y}_{t+1}) \quad (5.48)$$

$$= \mathcal{L}_{\beta_{t+1}}(\mathbf{x}_{t+1}, \bar{\mathbf{v}}_{t+1}, \mathbf{y}_t) + \langle \mathbf{y}_{t+1} - \mathbf{y}_t, \mathcal{A}\mathbf{x}_{t+1} - \bar{\mathbf{v}}_{t+1} \rangle \quad (5.49)$$

$$= \mathcal{L}_{\beta_{t+1}}(\mathbf{x}_{t+1}, \bar{\mathbf{v}}_{t+1}, \mathbf{y}_t) + \sigma_{t+1} \|\mathcal{A}\mathbf{x}_{t+1} - \bar{\mathbf{v}}_{t+1}\|^2. \quad (5.50)$$

We choose  $\sigma_{t+1} \geq 0$  such that

$$\sigma_{t+1} \leq 1/\beta_0 \quad \& \quad \sigma_{t+1} \|\mathcal{A}\mathbf{x}_{t+1} - \bar{\mathbf{v}}_{t+1}\|^2 \leq \eta_t^2 \frac{\bar{L}_{t+1}}{2} D_{\mathcal{X}}^2 \quad \& \quad \|\mathbf{y}_{t+1}\| \leq D_{\mathcal{Y}}. \quad (5.51)$$

Note that  $\sigma_{t+1}$  is well defined, in the sense that there always exists  $\sigma_{t+1} \geq 0$  which satisfies all conditions. This is because  $\sigma_{t+1} = 0$  simultaneously satisfies them in a trivial way.

Therefore, we finally obtain

$$\mathcal{L}_{\beta_{t+1}}(\mathbf{x}_{t+1}, \mathbf{v}_{t+1}, \mathbf{y}_{t+1}) - \mathcal{L}_\star \leq (1 - \eta_t)(\mathcal{L}_{\beta_t}(\mathbf{x}_t, \mathbf{v}_t, \mathbf{y}_t) - \mathcal{L}_\star) + \eta_t^2 \bar{L}_{t+1} D_{\mathcal{X}}^2 + \frac{\|\mathcal{A}\|^2 D_{\mathcal{X}}^2}{4\beta_0(t+2)^{2.5}}. \quad (5.52)$$

Applying this inequality recursively, we get

$$\begin{aligned} &\mathcal{L}_{\beta_{t+1}}(\mathbf{x}_{t+1}, \mathbf{v}_{t+1}, \mathbf{y}_{t+1}) - \mathcal{L}_\star \\ &\leq \prod_{j=0}^t (1 - \eta_j) (\mathcal{L}_{\beta_0}(\mathbf{x}_0, \mathbf{v}_0, \mathbf{y}_0) - \mathcal{L}_\star) + \sum_{\ell=0}^t \left( \eta_\ell^2 \bar{L}_{\ell+1} D_{\mathcal{X}}^2 + \frac{\|\mathcal{A}\|^2 D_{\mathcal{X}}^2}{4\beta_0(\ell+2)^{2.5}} \right) \prod_{j=\ell}^t (1 - \eta_j) \end{aligned} \quad (5.53)$$

$$= \sum_{\ell=0}^t \left( \eta_\ell^2 \bar{L}_{\ell+1} D_{\mathcal{X}}^2 + \frac{\|\mathcal{A}\|^2 D_{\mathcal{X}}^2}{4\beta_0(\ell+2)^{2.5}} \right) \prod_{j=\ell}^t (1 - \eta_j) \quad (5.54)$$

$$\leq \sum_{\ell=0}^t \left( \eta_\ell^2 \bar{L}_{t+1} D_{\mathcal{X}}^2 + \frac{\|\mathcal{A}\|^2 D_{\mathcal{X}}^2}{4\beta_0(\ell+2)^{2.5}} \right) \prod_{j=\ell}^t (1 - \eta_j) \quad (5.55)$$

where (5.54) follows since  $\eta_0 = 1$ , and (5.55) holds because  $\bar{L}_{t+1} \geq \bar{L}_{\ell+1}$  for  $\ell = 1, 2, \dots, t$ .

By using

$$\sum_{\ell=0}^t \eta_\ell^2 \prod_{j=\ell}^t (1 - \eta_j) = \sum_{\ell=0}^t \frac{4}{(\ell+2)^2} \prod_{j=\ell}^t \frac{j}{j+2} = \sum_{\ell=0}^t \frac{4}{(\ell+2)^2} \frac{(\ell+1)\ell}{(t+1)(t+2)} \leq \frac{4}{t+2}, \quad \text{and} \quad (5.56)$$

$$\sum_{\ell=0}^t \frac{1}{(\ell+2)^{2.5}} \prod_{j=\ell}^t \frac{j}{j+2} = \sum_{\ell=0}^t \frac{1}{(\ell+2)^{2.5}} \frac{(\ell+1)\ell}{(t+1)(t+2)} \leq \frac{2}{(t+2)^{1.5}}, \quad (5.57)$$

we get the following bound on the augmented Lagrangian:

$$\mathcal{L}_{\beta_{t+1}}(\mathbf{x}_{t+1}, \mathbf{v}_{t+1}, \mathbf{y}_{t+1}) - \mathcal{L}^\star \leq \frac{4}{t+2} D_{\mathcal{X}}^2 \bar{L}_{t+1} + \frac{\|\mathcal{A}\|^2 D_{\mathcal{X}}^2}{2\beta_0(t+2)^{1.5}} \quad (5.58)$$

$$= D_{\mathcal{X}}^2 \left( \frac{4L_f}{t+2} + \frac{4\|\mathcal{A}\|^2/\beta_0}{\sqrt{t+2}} + \frac{\|\mathcal{A}\|^2/\beta_0}{2(t+2)^{1.5}} \right). \quad (5.59)$$

Next, we translate this bound to the convergence guarantees in the objective residual and feasibility gap.

**Convergence in objective residual.** We start from the definition of augmented Lagrangian:

$$\begin{aligned} f(\mathbf{x}_{t+1}) &= \mathcal{L}_{\beta_{t+1}}(\mathbf{x}_{t+1}, \mathbf{v}_{t+1}, \mathbf{y}_{t+1}) + \frac{\beta_{t+1}}{2} \|\mathbf{y}_{t+1}\|^2 - \frac{1}{2\beta_{t+1}} \text{dist}^2(\mathcal{A}\mathbf{x}_{t+1} + \beta\mathbf{y}_{t+1}, \mathcal{K}) \\ &\leq \mathcal{L}_{\beta_{t+1}}(\mathbf{x}_{t+1}, \mathbf{v}_{t+1}, \mathbf{y}_{t+1}) + \frac{\beta_{t+1} D_{\mathcal{Y}}^2}{2}. \end{aligned} \quad (5.60)$$

Then, we subtract  $f_\star$  from both sides (recall that  $f_\star = \mathcal{L}_\star$  due to strong duality)

$$\begin{aligned} f(\mathbf{x}_{t+1}) - f_\star &\leq \mathcal{L}_{\beta_{t+1}}(\mathbf{x}_{t+1}, \mathbf{v}_{t+1}, \mathbf{y}_{t+1}) - \mathcal{L}^\star + \frac{\beta_{t+1} D_{\mathcal{Y}}^2}{2} \\ &\leq D_{\mathcal{X}}^2 \left( \frac{4L_f}{t+2} + \frac{4\|\mathcal{A}\|^2/\beta_0}{\sqrt{t+2}} + \frac{\|\mathcal{A}\|^2/\beta_0}{2(t+2)^{1.5}} \right) + \frac{D_{\mathcal{Y}}^2 \beta_0}{2\sqrt{t+2}}. \end{aligned} \quad (5.61)$$

This completes the proof of the upper bound.

Now, by using the classical Lagrange saddle point theory, we get

$$f_\star \leq \mathcal{L}(\mathbf{x}, \mathbf{v}, \mathbf{y}_\star) = f(\mathbf{x}) + \langle \mathbf{y}_\star, \mathcal{A}\mathbf{x} - \mathbf{v} \rangle \leq f(\mathbf{x}) + \|\mathbf{y}_\star\| \|\mathcal{A}\mathbf{x} - \mathbf{v}\|, \quad \forall (\mathbf{x}, \mathbf{v}) \in \mathcal{X} \times \mathcal{K}. \quad (5.62)$$

In particular, by choosing  $\mathbf{x} = \mathbf{x}_{t+1}$  and  $\mathbf{v} = \text{proj}_{\mathcal{K}}(\mathcal{A}\mathbf{x}_{t+1})$ , and rearranging, we obtain

$$f(\mathbf{x}_{t+1}) - f_\star \geq -\|\mathbf{y}_\star\| \cdot \text{dist}(\mathcal{A}\mathbf{x}_{t+1}, \mathcal{K}), \quad (5.63)$$

which is our lower bound on the objective residual.



**Convergence in feasibility gap.** Denote by  $\phi_t := D_{\mathcal{X}}^2 \left( \frac{4L_f}{t} + \frac{4\|\mathcal{A}\|^2/\beta_0}{\sqrt{t}} + \frac{\|\mathcal{A}\|^2/\beta_0}{2t^{1.5}} \right)$ . We combine (5.58) and (5.62) with  $\mathbf{x} = \mathbf{x}_{t+1}$  and  $\mathbf{v} = \mathbf{v}_{t+1}$ , and use Cauchy-Schwartz inequality:

$$\begin{aligned} \langle \mathbf{y}_{t+1} - \mathbf{y}_\star, \mathcal{A}\mathbf{x}_{t+1} - \mathbf{v}_{t+1} \rangle + \frac{1}{2\beta_t} \|\mathcal{A}\mathbf{x}_{t+1} - \mathbf{v}_{t+1}\|^2 &\leq \phi_{t+2} \\ \implies -\|\mathbf{y}_{t+1} - \mathbf{y}_\star\| \|\mathcal{A}\mathbf{x}_{t+1} - \mathbf{v}_{t+1}\| + \frac{1}{2\beta_t} \|\mathcal{A}\mathbf{x}_{t+1} - \mathbf{v}_{t+1}\|^2 &\leq \phi_{t+2} \end{aligned} \quad (5.64)$$

This is a second order inequality with respect to  $\|\mathcal{A}\mathbf{x}_{t+1} - \mathbf{v}_{t+1}\|$ . Solving this inequality we get

$$\|\mathcal{A}\mathbf{x}_{t+1} - \mathbf{v}_{t+1}\| \leq \beta_t \left( \|\mathbf{y}_{t+1} - \mathbf{y}_\star\| + \sqrt{\|\mathbf{y}_{t+1} - \mathbf{y}_\star\|^2 + \frac{2\phi_{t+2}}{\beta_t}} \right) \quad (5.65)$$

$$\leq \beta_t \left( 2\|\mathbf{y}_{t+1} - \mathbf{y}_\star\| + \sqrt{\frac{2\phi_{t+2}}{\beta_t}} \right) \quad (5.66)$$

$$\leq \frac{1}{\sqrt{t+2}} \left( 2\beta_0 \|\mathbf{y}_{t+1} - \mathbf{y}_\star\| + D_{\mathcal{X}} \sqrt{\frac{8L_f\beta_0}{\sqrt{t+2}}} + \left(8 + \frac{1}{t+2}\right) \|\mathcal{A}\|^2 \right). \quad (5.67)$$

Finally, the bound on the feasibility gap follows by

$$\begin{aligned} \text{dist}(\mathcal{A}\mathbf{x}_{t+1}, \mathcal{K}) &= \|\mathcal{A}\mathbf{x}_{t+1} - \dot{\mathbf{v}}_{t+1}\| \\ &= \|\mathcal{A}\mathbf{x}_{t+1} - \mathbf{v}_{t+1} + \mathbf{v}_{t+1} - \dot{\mathbf{v}}_{t+1}\| \\ &\leq \|\mathcal{A}\mathbf{x}_{t+1} - \mathbf{v}_{t+1}\| + \|\mathbf{v}_{t+1} - \dot{\mathbf{v}}_{t+1}\| \\ &\leq \|\mathcal{A}\mathbf{x}_{t+1} - \mathbf{v}_{t+1}\| + \|\mathcal{A}\mathbf{x}_{t+1} - \mathcal{A}\mathbf{x}_{t+1} + \beta_{t+1}\mathbf{y}_{t+1}\| \\ &\leq \|\mathcal{A}\mathbf{x}_{t+1} - \mathbf{v}_{t+1}\| + D_{\mathcal{Y}}\beta_{t+1}, \end{aligned} \quad (5.68)$$

where  $\dot{\mathbf{v}}_{t+1} = \text{proj}_{\mathcal{K}}(\mathcal{A}\mathbf{x}_{t+1})$ , and the fourth line follows from the non-expansiveness.



## 6 Low-Rank Matrix Sketching from Streaming Data

Our storage-optimal convex optimization paradigm adopts a sketching model as a mechanism for reducing the cost of storing the decision variable. Consequently, designing practical sketching algorithms for low-rank matrix approximation is a fundamental aspect of our research.

This chapter is based on the three joint works with Joel A. Tropp, Madeleine Udell, and Volkan Cevher; [TYUC17a], [TYUC17b], and [TYUC19]. Some details are taken from the preliminary versions [TYUC17c] and [TYUC18].

This section mainly consists of some excerpts quoted verbatim with minor modifications from these works, with the aim of providing a brief overview of our main results. These manuscripts are written by Joel A. Tropp, and he is the first author and the primary contact for these works. The author of this dissertation contributed to this research line primarily on the implementation of software and the execution of the computer experiments. Our motivation for studying low-rank matrix approximation from streaming data comes from the specific optimization application [YUTC17].

### Introduction

In Chapter 2, we introduced SketchyCGM, the first provable algorithm which produces a rank- $r$  approximate solution to a class of SDP problems with low-rank solutions using optimal storage. CGM forms a sequence of approximate low-rank solutions via the iteration

$$\mathbf{X}_1 \leftarrow \mathbf{0} \quad \text{and} \quad \mathbf{X}_{t+1} \leftarrow (1 - \eta_t)\mathbf{X}_t + \eta_t \mathbf{u}_t \mathbf{u}_t^*, \quad \text{for some } \eta_t \in [0, 1]. \quad (6.1)$$

SketchyCGM only maintains a random low-dimensional linear image of the evolving solution  $\mathbf{X}_t$ . When the iteration terminates, it computes a rank- $r$  approximation from this linear image.

This chapter presents a framework for computing structured low-rank approximations of a matrix from a *sketch*, a randomized low-dimensional linear image of the matrix. Suppose that  $\mathbf{X} \in \mathbb{R}^{m \times n}$  (or  $\mathbb{C}^{m \times n}$ ) is an arbitrary matrix. Let  $r$  be a target rank parameter where

$r \ll \min\{m, n\}$ . The computational problem is to produce a low-rank approximation  $\hat{\mathbf{X}}$  of  $\mathbf{X}$  whose error is comparable to a best rank- $r$  approximation:

$$\|\mathbf{X} - \hat{\mathbf{X}}\|_F \approx \min_{\text{rank}(\mathbf{B}) \leq r} \|\mathbf{X} - \mathbf{B}\|_F. \quad (6.2)$$

### Contributions

Our aim is to construct a random sketching map that acquires enough information to solve the low-rank matrix approximation problem (6.2). We must also design algorithms for computing a reconstruction  $\hat{\mathbf{X}}$  from the sketch without direct access to the matrix  $\mathbf{X}$  itself. We can summarize the main features of our approach and the contributions as follows:

- We design sketching maps that minimizes the size of the sketch (and the test matrices) to achieve approximations that satisfy (6.2).
- We propose novel reconstruction algorithms that are *reliable*, *numerically stable*, and *efficient* in terms of storage and computation.
- We present a systematic mechanism to maintain structural properties of the input matrix in our approximations, such as symmetry or positive-semi-definiteness.
- We present explicit *a priori* error bounds for the algorithms so that we can determine a sketch size that suffices to achieve a specific approximation goal.
- Our precise error bounds with respect to the tail energy of the input matrix offers concrete guidance on the best sketch-size parameters. (details excluded in this dissertation, see Section 4.5 in [TYUC17b] and Section 5.4 in [TYUC19])
- We document numerical experiments on real and synthetic data to demonstrate that our methods dominate existing techniques. Our codes are available online<sup>1</sup> as a MATLAB toolbox.
- Our toolkit includes an *a posteriori* error estimator for validating the quality of the approximation. This estimator also provides a principled mechanism for selecting the precise rank of the final approximation. (details excluded in this dissertation, see Section 6 in [TYUC19])
- We demonstrate an application of our approach as a tool to compress large-scale scientific data from several scientific simulations and measurement processes. (details excluded in this dissertation, see Section 7 in [TYUC19])

This chapter overviews only the key features of our results. We omit all proofs and technical details in this short summary. We refer to the original papers for more details.

---

<sup>1</sup>Available at "<https://github.com/alpyurtsever/SKETCH>"

### Preliminaries

The methods apply for the real field ( $\mathbb{F} = \mathbb{R}$ ) and for the complex field ( $\mathbb{F} = \mathbb{C}$ ). We introduce a parameter  $\alpha$  that reflect the field over which we are working:

$$\alpha := \alpha(\mathbb{F}) := \begin{cases} 1, & \mathbb{F} = \mathbb{R} \\ 0, & \mathbb{F} = \mathbb{C} \end{cases} \quad (6.3)$$

Let us also introduce a function that we use to simplify many of our error bounds:

$$\omega(u, v) := \frac{u}{v - u - \alpha} \quad \text{for integers that satisfy } v > u + \alpha > \alpha. \quad (6.4)$$

Observe that the function  $\omega(u, \cdot)$  is decreasing, with range  $(0, u]$ .

**Standard normal matrix.** A matrix  $\mathbf{G} \in \mathbb{R}^{m \times n}$  has the real standard normal distribution if the entries form an independent family of standard normal random variables (i.e., Gaussian with mean zero and variance one). A matrix  $\mathbf{G} \in \mathbb{C}^{m \times n}$  has the complex standard normal distribution if it has the form  $\mathbf{G} = \mathbf{G}_1 + i\mathbf{G}_2$  where  $\mathbf{G}_1$  and  $\mathbf{G}_2$  are independent, real standard normal matrices. Standard normal matrices are also known as Gaussian matrices.

Throughout, we present the analysis of the case where the linear dimension reduction maps (test matrices) are statistically independent and follow the standard normal distribution. We describe other potential distributions for the test matrices in Section 3.9 of [TYUC17b] and Section 3 of [TYUC19].

## 6.1 Two Component Sketch

This section presents a brief summary of our approach introduced in [TYUC17b].

### 6.1.1 The Sketch

To form the sketch of the target matrix, we independently draw and fix two test matrices:

$$\mathbf{\Omega} \in \mathbb{F}^{n \times k} \quad \text{and} \quad \mathbf{\Psi} \in \mathbb{F}^{\ell \times m}. \quad (6.5)$$

The sketch of the target matrix  $\mathbf{X} \in \mathbb{F}^{m \times n}$  consists of two matrices:

$$\mathbf{Y} := \mathbf{X}\mathbf{\Omega} \in \mathbb{F}^{m \times k} \quad \text{and} \quad \mathbf{W} := \mathbf{\Psi}\mathbf{X} \in \mathbb{F}^{\ell \times n}. \quad (6.6)$$

The matrix  $\mathbf{Y}$  collects information about the action of  $\mathbf{X}$ , while the matrix  $\mathbf{W}$  collects information about the action of  $\mathbf{X}^*$ .

### 6.1.2 The Basic Reconstruction Algorithm

The first step in the reconstruction is to compute an orthobasis  $\mathbf{Q}$  for the range of  $\mathbf{Y}$  by means of an orthogonal–triangular factorization:

$$\mathbf{Y} =: \mathbf{Q}\mathbf{R} \quad \text{where} \quad \mathbf{Q} \in \mathbb{F}^{m \times k}. \quad (6.7)$$

The second step uses the co-range sketch  $\mathbf{W}$  to form the matrix

$$\mathbf{B} := (\Psi\mathbf{Q})^\dagger \mathbf{W} \in \mathbb{F}^{k \times n}. \quad (6.8)$$

Then, we report the rank- $k$  approximation

$$\hat{\mathbf{X}} := \mathbf{Q}\mathbf{B} \in \mathbb{F}^{m \times n} \quad \text{where} \quad \mathbf{Q} \in \mathbb{F}^{m \times k} \quad \text{and} \quad \mathbf{B} \in \mathbb{F}^{k \times n}. \quad (6.9)$$

**Intuition.** To motivate the algorithm, we recall a familiar heuristic from randomized linear algebra (see Section 1 in [HMT11]), which states

$$\mathbf{X} \approx \mathbf{Q}\mathbf{Q}^* \mathbf{X}. \quad (6.10)$$

Although we would like to form the rank- $k$  approximation  $\mathbf{Q}(\mathbf{Q}^* \mathbf{X})$ , we cannot compute the factor  $\mathbf{Q}^* \mathbf{X}$  without revisiting the target matrix  $\mathbf{X}$ . Instead, we exploit the information in the co-range sketch  $\mathbf{W} = \Psi\mathbf{X}$ . Notice that

$$\mathbf{W} = \Psi(\mathbf{Q}\mathbf{Q}^* \mathbf{X}) + \Psi(\mathbf{X} - \mathbf{Q}\mathbf{Q}^* \mathbf{X}) \approx (\Psi\mathbf{Q})(\mathbf{Q}^* \mathbf{X}). \quad (6.11)$$

The heuristic (6.10) justifies dropping the second term. Multiplying on the left by the pseudoinverse  $(\Psi\mathbf{Q})^\dagger$ , we arrive at the relation

$$\mathbf{B} = (\Psi\mathbf{Q})^\dagger \mathbf{W} \approx \mathbf{Q}^* \mathbf{X}. \quad (6.12)$$

These considerations suggest that

$$\hat{\mathbf{X}} = \mathbf{Q}\mathbf{B} \approx \mathbf{Q}\mathbf{Q}^* \mathbf{X} \approx \mathbf{X}. \quad (6.13)$$

This explanation is inspired by the discussion in Section 5.5 in [HMT11].

**Theorem 6.1.** *Assume the sketch size parameters satisfy  $k > r + \alpha$  and  $\ell > k + \alpha$ . Draw random test matrices  $\mathbf{\Omega} \in \mathbb{F}^{n \times k}$  and  $\Psi \in \mathbb{F}^{\ell \times m}$  independently from the standard normal distribution. Then the rank- $k$  approximation  $\hat{\mathbf{X}}$  obtained from formula (6.9) satisfies*

$$\mathbb{E} \|\mathbf{X} - \hat{\mathbf{X}}\|_{\mathbb{F}}^2 \leq (1 + \omega(r, k)) \cdot (1 + \omega(k, \ell)) \cdot \|\mathbf{X} - \llbracket \mathbf{X} \rrbracket_r\|_{\mathbb{F}}^2. \quad (6.14)$$

In particular, the selection  $k = 2r + \alpha$  and  $\ell = 2k + \alpha$  yields the error bound

$$\left(\mathbb{E} \|\mathbf{X} - \hat{\mathbf{X}}\|_{\text{F}}^2\right)^{1/2} \leq 2 \cdot \|\mathbf{X} - \llbracket \mathbf{X} \rrbracket_r\|_{\text{F}}. \quad (6.15)$$

We recommend this parameter choice because it achieves a good tradeoff between the size of the sketch and the accuracy of the reconstruction.

Similar bounds hold with high probability.<sup>2</sup>

**Remark.** We can extend the reconstruction error bound in Theorem 6.1, reformulating the bound with respect to tail energy instead of the optimal rank- $r$  error  $\|\mathbf{X} - \llbracket \mathbf{X} \rrbracket_r\|_{\text{F}}$ . This provides a good description of the actual behavior. As a consequence, we can offer a concrete guidance on the sketch size parameters based on the spectral decay of the input matrix. See Theorem 4.3 and Section 4.5 in [TYUC17b] for more details.

### 6.1.3 The Fixed-Rank Reconstruction Algorithm

Suppose that we wish to compute a rank- $r$  approximation of the target matrix  $\mathbf{X} \in \mathbb{F}^{m \times n}$ . First, we form an initial rank- $k$  approximation  $\hat{\mathbf{X}} := \mathbf{Q}\mathbf{B}$  using the procedure (6.9). Then we obtain a rank- $r$  approximation  $\llbracket \hat{\mathbf{X}} \rrbracket_r$  of the target matrix by replacing  $\hat{\mathbf{X}}$  with its best rank- $r$  approximation in Frobenius norm:

$$\llbracket \hat{\mathbf{X}} \rrbracket_r = \llbracket \mathbf{Q}\mathbf{B} \rrbracket_r = (\mathbf{Q}\mathbf{U}) \llbracket \Sigma \rrbracket_r \mathbf{V}^* = \mathbf{Q} \llbracket \mathbf{B} \rrbracket_r, \quad \text{where } \mathbf{B} = \mathbf{U}\Sigma\mathbf{V}^* \text{ is an SVD of } \mathbf{B}. \quad (6.16)$$

**Theorem 6.2.** Assume the sketch size parameters satisfy  $k > r + \alpha$  and  $\ell > k + \alpha$ . Draw random test matrices  $\mathbf{\Omega} \in \mathbb{F}^{n \times k}$  and  $\mathbf{\Psi} \in \mathbb{F}^{\ell \times m}$  independently from the standard normal distribution. Then the rank- $r$  approximation  $\llbracket \hat{\mathbf{X}} \rrbracket_r$  obtained from the formula eq. (6.16) satisfies

$$\mathbb{E} \|\mathbf{X} - \llbracket \hat{\mathbf{X}} \rrbracket_r\|_{\text{F}} \leq \sqrt{1 + \omega(r, k)} \cdot (1 + 2\sqrt{\omega(k, \ell)}) \cdot \|\mathbf{X} - \llbracket \mathbf{X} \rrbracket_r\|_{\text{F}}. \quad (6.17)$$

In particular, the selection  $k = 2r + \alpha$  and  $\ell = 2k + \alpha$  yields the error bound

$$\mathbb{E} \|\mathbf{X} - \llbracket \hat{\mathbf{X}} \rrbracket_r\|_{\text{F}} \leq 3\sqrt{2} \cdot \|\mathbf{X} - \llbracket \mathbf{X} \rrbracket_r\|_{\text{F}}. \quad (6.18)$$

Similar results hold with high probability.

### 6.1.4 Computing a PSD Approximation

Suppose that the target matrix  $\mathbf{X}$  belongs to the closed, convex set  $\mathcal{X} \subset \mathbb{F}^{m \times n}$ . Let  $\hat{\mathbf{X}}_{\text{in}} \in \mathbb{F}^{m \times n}$  be an initial approximation of  $\mathbf{X}$ . We can produce a new approximation  $\text{proj}_{\mathcal{X}}(\hat{\mathbf{X}}_{\text{in}})$  by projecting the initial approximation onto the constraint set. This procedure always improves the approximation quality. This result is well known in convex analysis.

<sup>2</sup> The expectation bounds in this and the following results also describe the typical behavior because of measure concentration effects. One can develop high-probability bounds using the methods from Section 10.3 in [HMT11].

We often encounter the problem of approximating a PSD (or conjugate symmetric) matrix. Projection onto PSD cone can be efficiently performed on the factors, avoiding forming large matrices. See Section 5 in [TYUC17b] for the details.

Last of all, consider the situation where we need to approximate a structured matrix by a structured matrix with fixed rank. The most direct approach is to replace an initial structured approximation by the nearest structured rank- $r$  matrix. Projecting a PSD (respectively, conjugate symmetric) matrix onto the set of fixed-rank matrices preserves PSD (the conjugate symmetric property). The set of matrices with a fixed rank is not convex, however, so the previous analysis does not apply. To this end, we provide the following general bound for fixed-rank approximation with structure. This bound is typically somewhat loose.

**Theorem 6.3.** *Let  $\mathbf{X} \in \mathbb{F}^{m \times n}$  be a target matrix, and let  $\hat{\mathbf{X}}_{\text{in}} \in \mathbb{F}^{m \times n}$  be an approximation. For any rank parameter  $r$ ,*

$$\|\mathbf{X} - \llbracket \hat{\mathbf{X}}_{\text{in}} \rrbracket_r\|_F \leq \|\mathbf{X} - \llbracket \mathbf{X} \rrbracket_r\|_F + 2\|\mathbf{X} - \hat{\mathbf{X}}_{\text{in}}\|_F. \quad (6.19)$$

## 6.2 Three Component Sketch

This section presents the key features of our new sketching approach with three components [TYUC19]. The new algorithm is a hybrid of the methods from [Upa16] and our previous approach with two components. The new algorithm performs better than its predecessors.

### 6.2.1 The Sketch

Independently, draw and fix four randomized linear dimension reduction maps:

$$\begin{aligned} \mathbf{\Upsilon} &\in \mathbb{F}^{k \times m} \quad \text{and} \quad \mathbf{\Omega} \in \mathbb{F}^{k \times n}, \\ \mathbf{\Phi} &\in \mathbb{F}^{s \times m} \quad \text{and} \quad \mathbf{\Psi} \in \mathbb{F}^{s \times n}, \end{aligned} \quad (6.20)$$

such that  $r_0 \leq k \leq s \leq \min\{m, n\}$ , where  $r_0$  is the target rank.

The sketch itself consists of three matrices:

$$\mathbf{W} := \mathbf{\Upsilon} \mathbf{X} \in \mathbb{F}^{k \times n}, \quad \mathbf{Y} := \mathbf{X} \mathbf{\Omega}^* \in \mathbb{F}^{m \times k}, \quad \text{and} \quad \mathbf{Z} := \mathbf{\Phi} \mathbf{X} \mathbf{\Psi}^* \in \mathbb{F}^{s \times s}. \quad (6.21)$$

The first two matrices ( $\mathbf{W}, \mathbf{Y}$ ) capture the co-range and the range of  $\mathbf{X}$ . The core sketch ( $\mathbf{Z}$ ) contains fresh information that improves our estimates of the singular values and singular vectors of  $\mathbf{X}$ ; it is responsible for the superior performance of the new method.



### 6.2.2 Computing Truncated Low-Rank Approximations

The first two components  $(W, Y)$  of the sketch are used to estimate the co-range and the range of the matrix  $X$ . Compute thin QR factorizations:

$$\begin{aligned} W^* &=: PR_1 \quad \text{where } P \in \mathbb{F}^{n \times k}; \\ Y &=: QR_2 \quad \text{where } Q \in \mathbb{F}^{m \times k}. \end{aligned} \quad (6.22)$$

The third sketch  $Z$  is used to compute the core approximation  $B$ , which describes how  $X$  acts between  $\text{range}(P)$  and  $\text{range}(Q)$ :

$$B := (\Phi Q)^\dagger Z ((\Psi P)^\dagger)^* \in \mathbb{F}^{k \times k}. \quad (6.23)$$

This step is implemented by solving a family of least-squares problems. Next, form a rank- $k$  approximation  $\hat{X}$  of the input matrix  $X$  via

$$\hat{X} := QBP^* \quad (6.24)$$

We refer to  $\hat{X}$  as the “initial” approximation. The initial approximation can contain spurious information. To produce an approximation that is fully reliable, we must truncate the rank of the initial approximation (6.24). For a truncation parameter  $r$ , we construct a rank- $r$  approximation by replacing  $\hat{X}$  with its best rank- $r$  approximation in Frobenius norm:

$$[\hat{X}]_r = Q[B]_r P^*. \quad (6.25)$$

We refer to  $[\hat{X}]_r$  as a “truncated” approximation.

**Remark.** We can form other structured approximations of  $X$  by projecting  $\hat{X}$  onto a set of structured matrices as in Section 6.1.4. We omit the details.

### 6.2.3 Analysis of Initial and Truncated Approximations

**Theorem 6.4** (Initial Approximation: Error Bound). *Let  $X \in \mathbb{F}^{m \times n}$  be an arbitrary input matrix. Assume the sketch size parameters satisfy  $s \geq 2k + \alpha$ . Draw independent Gaussian dimension reduction maps  $(Y, \Omega, \Phi, \Psi)$ , as in (6.20). Extract a sketch (6.21) of the input matrix. Then the rank- $k$  approximation  $\hat{X}$ , constructed in (6.24), satisfies the error bound*

$$\mathbb{E} \|X - \hat{X}\|_F^2 \leq (1 + \omega(k, s)) \cdot (1 + 2\omega(r, k)) \cdot \|X - [X]_r\|_F^2. \quad (6.26)$$

In particular,  $k = 4r_0 + \alpha$  and  $s = 2k + \alpha$  yields

$$\mathbb{E} \|X - \hat{X}\|_F^2 \leq \frac{10}{3} \cdot \|X - [X]_r\|_F^2. \quad (6.27)$$

Similar bounds hold with high probability.

**Remark.** In Theorem 6.4, we have imposed the condition  $s \geq 2k + \alpha$  because theoretical analysis and empirical work both suggest that the restriction is useful in practice. The approximation (6.24) only requires that  $k \leq s$ .

**Remark.** We can extend the reconstruction error bound in Theorem 6.4, reformulating the bound with respect to the tail energy, in order to provide a concrete guidance on the sketch size parameters based on the spectral decay of the input matrix. See Theorem 5.1 and Section 5.4 in [TYUC19] for more details.

**Corollary 6.1** (Truncated Approximation: Error Bound). Instate the assumptions of Theorem 6.4. Then the rank- $r$  approximation  $[\hat{X}]_r$  satisfies the error bound

$$\mathbb{E} \|X - [\hat{X}]_r\|_F \leq \left(1 + 2\sqrt{(1 + \omega(k, s)) \cdot (1 + 2\omega(r, k))}\right) \cdot \|X - [X]_r\|_F \quad (6.28)$$

Similar results hold with high probability.

This statement is an immediate consequence of Theorem 6.4 and Theorem 6.3.

### 6.3 Nyström Sketch for PSD Matrices

This section presents a brief summary of our Nystöm Sketch for fixed-rank approximation of a positive-semidefinite (PSD) matrix in [TYUC17a].

#### 6.3.1 The Sketch

Suppose that  $X \in \mathbb{F}^{n \times n}$  is a PSD matrix. Fix a sketch size parameter  $k$  in the range  $r \leq k \leq n$ . Independent from  $X$ , we draw and fix a random test matrix

$$\Omega \in \mathbb{F}^{n \times k}. \quad (6.29)$$

The sketch of the matrix  $X$  takes the form

$$Y = X\Omega \in \mathbb{F}^{n \times k}. \quad (6.30)$$

To find a good approximation, we must set the sketch size  $k$  larger than  $r$ . But storage costs and computation also increase with  $k$ . One of our main contributions is to clarify the role of  $k$ .

**Intuition.** The Nyström method is a general technique for low-rank PSD matrix approximation. Various instantiations appear in the papers [WS00, Pla05, FBCM04, DM05, HMT11, Git11, CD13, Git13, GM16, LLS<sup>+</sup>17].

Given the test matrix  $\mathbf{\Omega}$  and the sketch  $\mathbf{Y} = \mathbf{X}\mathbf{\Omega}$ , the Nyström method constructs a rank- $k$  PSD approximation of the PSD matrix  $\mathbf{X}$  via

$$\hat{\mathbf{X}}^{\text{nys}} = \mathbf{Y}(\mathbf{\Omega}^* \mathbf{Y})^\dagger \mathbf{Y}^*. \quad (6.31)$$

In most work on the Nyström method, the test matrix  $\mathbf{\Omega}$  depends adaptively on  $\mathbf{X}$ , so these approaches are not valid in the streaming setting. Gittens's framework [Git11, Git13, GM16] covers the streaming case.

**Fixed-Rank Nyström Approximation: Prior Art.** To construct a Nyström approximation with exact rank  $r$  from a sketch of size  $k$ , the standard approach is to truncate the center

$$\hat{\mathbf{X}}_r^{\text{nysfix}} = \mathbf{Y}(\llbracket \mathbf{\Omega}^* \mathbf{Y} \rrbracket_r)^\dagger \mathbf{Y}^*. \quad (6.32)$$

The truncated Nyström approximation (6.32) appears in the many papers, including [Pla05, DM05, CD13, GM13]. We have found that the truncation method (6.32) performs poorly in the present setting. This observation motivated us to search for more effective techniques.

**Fixed-Rank Nyström Approximation: Proposal.** Our purpose is to develop, analyze, and evaluate a new approach for fixed-rank approximation of a PSD matrix under the streaming model. We propose a more intuitive rank- $r$  approximation:

$$\hat{\mathbf{X}}_r = \llbracket \hat{\mathbf{X}}^{\text{nys}} \rrbracket_r. \quad (6.33)$$

That is, we report a best rank- $r$  approximation of the full Nyström approximation (6.31).

### 6.3.2 Fixed-Rank Nyström Approximation

We outline a numerically stable and very accurate implementation of (6.33), based on an idea from [Tyg14, LLS<sup>+</sup>17]. In Section 6.4, we show that a poor implementation may produce an approximation with 100% error!

Fix a small parameter  $\nu > 0$ . Instead of approximating the PSD matrix  $\mathbf{X}$  directly, we approximate the shifted matrix  $\mathbf{X}_\nu = \mathbf{X} + \nu \mathbf{I}$  and then remove the shift. Here are the steps:

1. Construct the shifted sketch  $\mathbf{Y}_\nu = \mathbf{Y} + \nu \mathbf{\Omega}$ .
2. Form the matrix  $\mathbf{B} = \mathbf{\Omega}^* \mathbf{Y}_\nu$ .
3. Compute a Cholesky decomposition  $\mathbf{B} = \mathbf{C} \mathbf{C}^*$ .
4. Compute  $\mathbf{E} = \mathbf{Y}_\nu \mathbf{C}^{-1}$  by back-substitution.
5. Compute the (thin) singular value decomposition  $\mathbf{E} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^*$ .
6. Form  $\hat{\mathbf{X}}_r = \mathbf{U} \llbracket \mathbf{\Sigma}^2 - \nu \mathbf{I} \rrbracket_r \mathbf{U}^*$ .

Related, but distinct, methods were proposed by Williams & Seeger [WS00] and analyzed in Gittens's thesis [Git13].

**Theorem 6.5.** Assume  $1 \leq r < k \leq n$ . Let  $\mathbf{X} \in \mathbb{F}^{n \times n}$  be a PSD matrix. Draw a test matrix  $\mathbf{\Omega} \in \mathbb{F}^{n \times k}$  from the Gaussian distribution, and form the sketch  $\mathbf{Y} = \mathbf{X}\mathbf{\Omega}$ . Then the approximation  $\hat{\mathbf{X}}_r$  given by (6.31) and (6.33) satisfies

$$\mathbb{E} \|\mathbf{X} - \hat{\mathbf{X}}_r\|_{S_1} \leq (1 + \omega(r, k)) \cdot \|\mathbf{X} - \llbracket \mathbf{X} \rrbracket_r\|_{S_1}; \quad (6.34)$$

$$\mathbb{E} \|\mathbf{X} - \hat{\mathbf{X}}_r\|_{S_\infty} \leq \|\mathbf{X} - \llbracket \mathbf{X} \rrbracket_r\|_{S_\infty} + \omega(r, k) \cdot \|\mathbf{X} - \llbracket \mathbf{X} \rrbracket_r\|_{S_1}. \quad (6.35)$$

Similar results hold with high probability.

## 6.4 Numerical Experiments

This section presents computer experiments that are designed to evaluate the performance of the proposed sketching algorithms for low-rank matrix approximation. We include comparisons with alternative methods from the literature.

### 6.4.1 Sketching and Reconstruction Methods in our Experiments

We use the following algorithms in our numerical experiments:

- The [HMT11] method in Section 5.5, Remark 5.4 of [HMT11] is a simplification of the method from Section 5.2 of Woolfe et al. [WLRT08], and they perform similarly. There are two sketches, and the sketch size depends on one parameter  $k$ . The total storage cost  $T = k(m + n)$ .
- The [TYUC17b] method is our two component sketch, implemented as described in Algorithm 7 of [TYUC17b]. It involves two sketches, controlled by two parameters  $k, \ell$ . The total storage cost  $T = km + \ell n$ .
- The [Upa16] method in Section 3.3 of [Upa16] simplifies a complicated approach from Theorem 12 in Boutsidis et al. [BWZ16]. This algorithm involves three sketches, controlled by two parameters  $k, s$ . The total storage cost  $T = k(m + n) + s^2$ .
- The [TYUC2019] method in is our three component sketching approach, implemented as described in Algorithm 4.4 of [TYUC19]. It uses three sketches, controlled by two parameters  $k, s$ . The total storage cost  $T = k(m + n) + s^2$ .

For our tests on fixed-rank reconstruction of a PSD in Section 6.4.5, we use:

- The standard Nystrom approximation (6.32), and the proposed approach (6.33), (implemented as described in Algorithm 3 of [TYUC17a]). These methods involve one sketch, controlled by a single parameters  $k$ . The total storage cost  $T = kn$ .
- The [TYUC17b] method is our two component sketch, implemented for PSD matrices as described in Algorithm 9 of [TYUC17b]. It involves two sketches, controlled by two parameters  $k, \ell$ . The total storage cost  $T = (k + \ell)n$ .

### 6.4.2 Experimental Setup

Fix an input matrix  $\mathbf{X} \in \mathbb{F}^{n \times n}$  and a truncation rank  $r$ . Select sketch size parameters. For each trial, draw dimension reduction maps and form the sketch of the input matrix. Compute a rank- $r$  approximation  $\hat{\mathbf{X}}_{\text{out}}$  using a specified reconstruction algorithm. The approximation error is calculated relative to the best rank- $r$  approximation error in Schatten  $p$ -norm:

$$S_p \text{ relative error} = \frac{\|\mathbf{X} - \hat{\mathbf{X}}_{\text{out}}\|_p}{\|\mathbf{X} - [\![\mathbf{X}]\!]_r\|_p} - 1. \quad (6.36)$$

We perform 20 independent trials and report the average error.

### 6.4.3 Classes of Input Matrices

We consider several different types of synthetic and real input matrices. See Figure 6.1 for a plot of the spectra of these input matrices.

**Synthetic Examples.** We work over the complex field  $\mathbb{C}$ . The matrix dimensions  $m = n = 10^3$ , and we introduce an effective rank parameter  $R$ . We compute an approximation with truncation rank  $r = 10$ .

◦ **Low-rank + noise:** Let  $\xi \geq 0$  be a signal-to-noise parameter. These matrices take the form

$$\mathbf{X} = \text{diag}(\underbrace{1, \dots, 1}_R, 0, \dots, 0) + \xi n^{-1} \mathbf{C} \in \mathbb{C}^{n \times n}, \quad (6.37)$$

where  $\mathbf{C} = \mathbf{G}\mathbf{G}^*$  for a standard normal matrix  $\mathbf{G} \in \mathbb{F}^{n \times n}$ .

LowRankLowNoise ( $\xi = 10^{-4}$ ), LowRankMedNoise ( $\xi = 10^{-2}$ ), LowRankHiNoise ( $\xi = 10^{-1}$ )

◦ **Polynomial decay:** For a decay parameter  $p > 0$ , consider matrices

$$\mathbf{X} = \text{diag}(\underbrace{1, \dots, 1}_R, 2^{-p}, 3^{-p}, \dots, (n - R + 1)^{-p}) \in \mathbb{C}^{n \times n}. \quad (6.38)$$

PolyDecaySlow ( $p = 0.5$ ), PolyDecayMed ( $p = 1$ ), PolyDecayFast ( $p = 2$ ).

◦ **Exponential decay:** For a decay parameter  $q > 0$ , consider matrices

$$\mathbf{X} = \text{diag}(\underbrace{1, \dots, 1}_R, 10^{-q}, 10^{-2q}, \dots, 10^{-(n-R)q}) \in \mathbb{C}^{n \times n}. \quad (6.39)$$

ExpDecaySlow ( $q = 0.01$ ), ExpDecayMed ( $q = 0.1$ ), ExpDecayFast ( $q = 0.5$ )

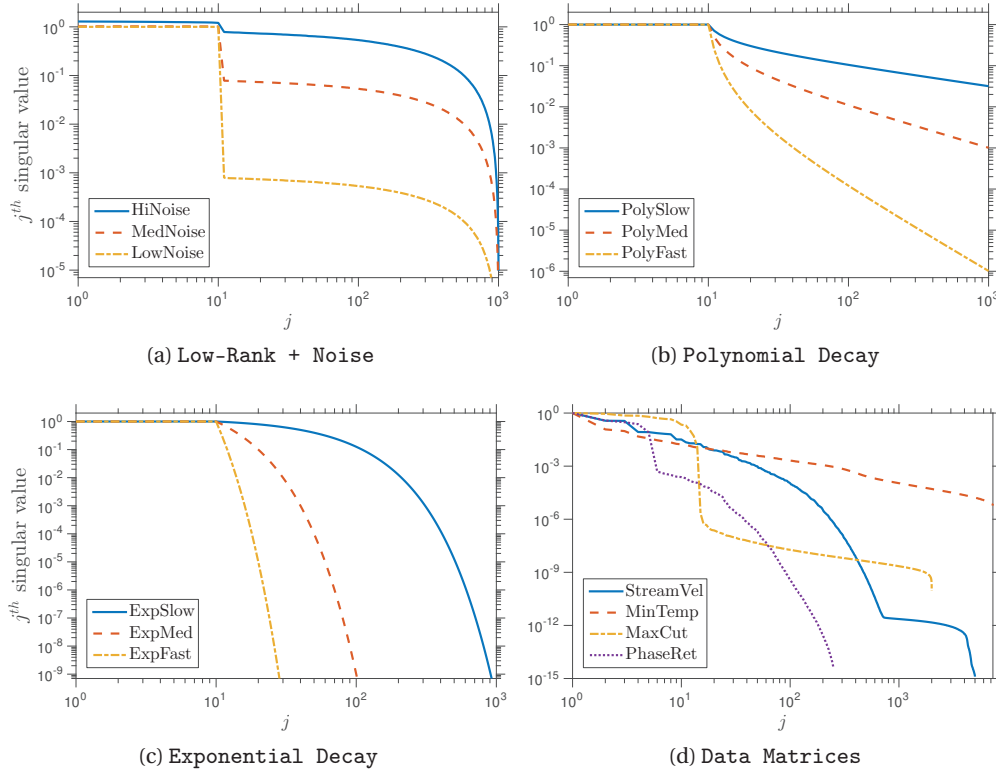


Figure 6.1 – Singular value spectrum for each of the synthetic and real data classes.

**Application Examples.** We present some low-rank data matrices that arise in applications.

- **Navier–Stokes:** This data is courtesy of Beverley McKeon and Sean Symon. The real  $m \times n$  matrix `StreamVel` contains streamwise velocities at  $m = 10,738$  points for each of  $n = 5,001$  time instants.

- **Weather:** This data is courtesy of William North. The real  $m \times n$  matrix `MinTemp` contains the minimum temperature recorded at  $m = 19,264$  stations on each of  $n = 7,305$  days, across the northeastern United States during the years 1981–2016.

We also consider matrices that arise in our optimization problems:

- **MaxCut:** This is a real psd matrix with  $m = n = 2,000$  that gives a high-accuracy solution to the max-cut SDP for a sparse graph [GW95].

- **PhaseRetrieval:** This is a complex psd matrix with  $m = n = 25,921$  that gives a low-accuracy solution to a phase retrieval SDP [HCO<sup>+</sup>15].

#### 6.4.4 Comparison of Formulas for General Low-Rank Matrices

Figure 6.2 presents the results of the following experiment. For synthetic matrices with effective rank  $R = 10$  and truncation rank  $r = 10$ , we compare the relative error eq. (6.36) achieved by each of the four algorithms as a function of storage. We use Gaussian dimension reduction maps in these experiments.

Figure 6.3 contains the results of the following experiment. For each of the four algorithms, we display the relative error eq. (6.36) as a function of storage. We use sparse dimension reduction maps, see Section 3.3 in [TYUC19] for the details.

We plot the oracle error<sup>3</sup> attained by each method. Since the oracle error is not achievable in practice, we also chart the performance of each method at the *a priori* parameter selection.

The numerical work here supports the superiority of the proposed methods.

#### 6.4.5 Comparison of Formulas for Low-Rank PSD Matrices

Figures 6.4 and 6.5 display the performance of the three fixed-rank psd approximation methods. The vertical axis is the Schatten 1-norm relative error (6.36). The variable  $T$  is the storage required for the sketch only. For the Nyström-based approximations (6.32) and (6.33), we have the correspondence  $T/n = k$ . For Algorithm 9 in [TYUC17b], we have  $T/n = k + \ell$ .

Figure 6.6 gives evidence about the numerical challenges involved in implementing Nyström approximations, such as (6.33). Our implementation is based on the Nyström approximation routine `eigenn` released by Tygert [Tyg14] to accompany the paper [LLS<sup>+</sup>17]. We compare with another implementation strategy described in the text of the same paper, Eqn. (13) in [LLS<sup>+</sup>17]. It is surprising to discover very different levels of precision in two implementations designed by professional numerical analysts.

The experiments demonstrate that the proposed method (6.33) has a significant benefit over the alternatives for input matrices that admit a good low-rank approximation.

---

<sup>3</sup>To make fair comparisons among algorithms, we can fix the storage budget and identify the parameter choices that minimize the (average) relative error incurred over the repeated trials. We refer to the minimum as the oracle error for an algorithm. The oracle error is not attainable in practice.

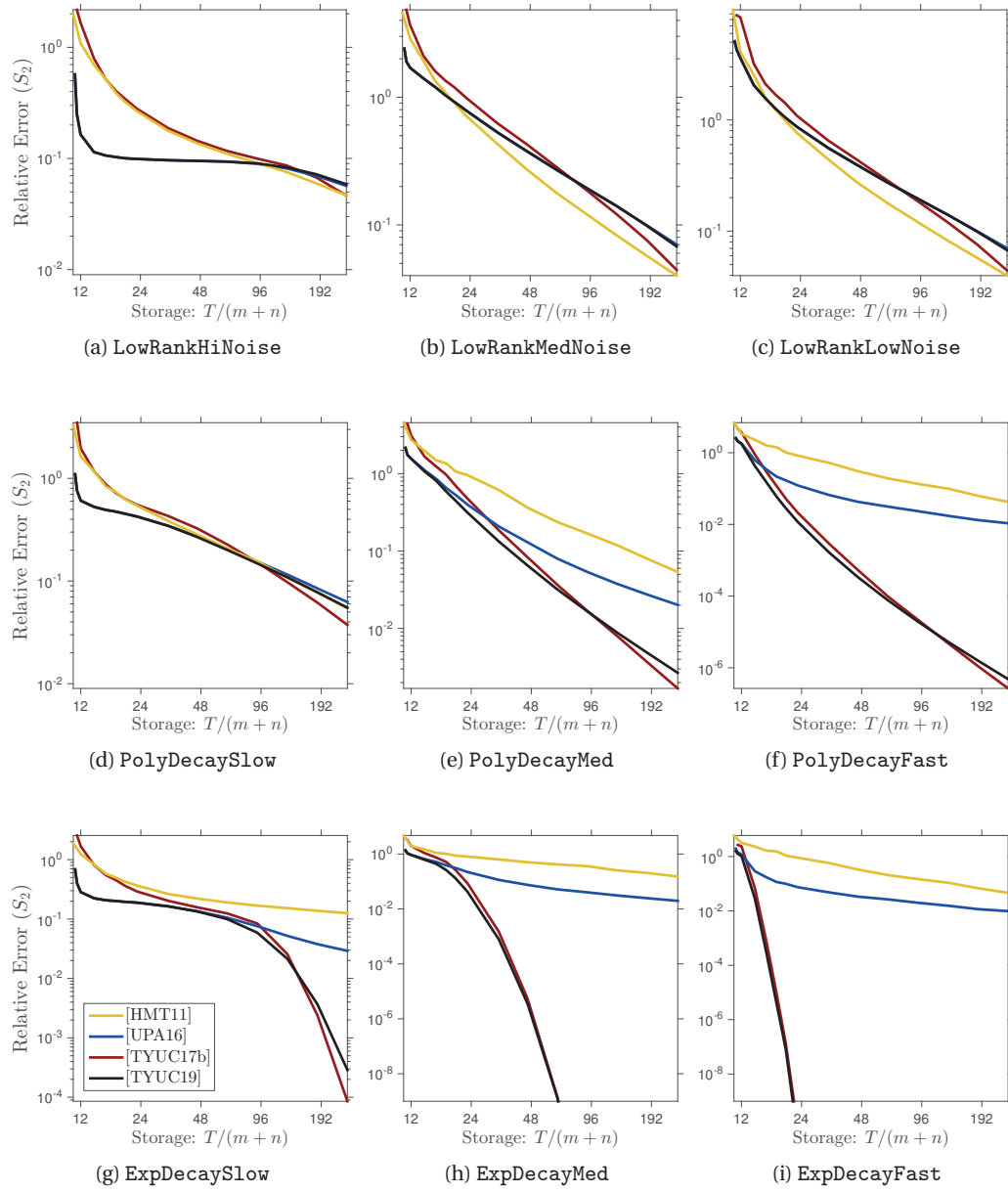


Figure 6.2 – [Comparison of reconstruction formulas: Synthetic examples.] (Gaussian maps, effective rank  $R = 10$ , approximation rank  $r = 10$ , Schatten 2-norm.) We compare the oracle error achieved by different fixed-rank approximation methods.



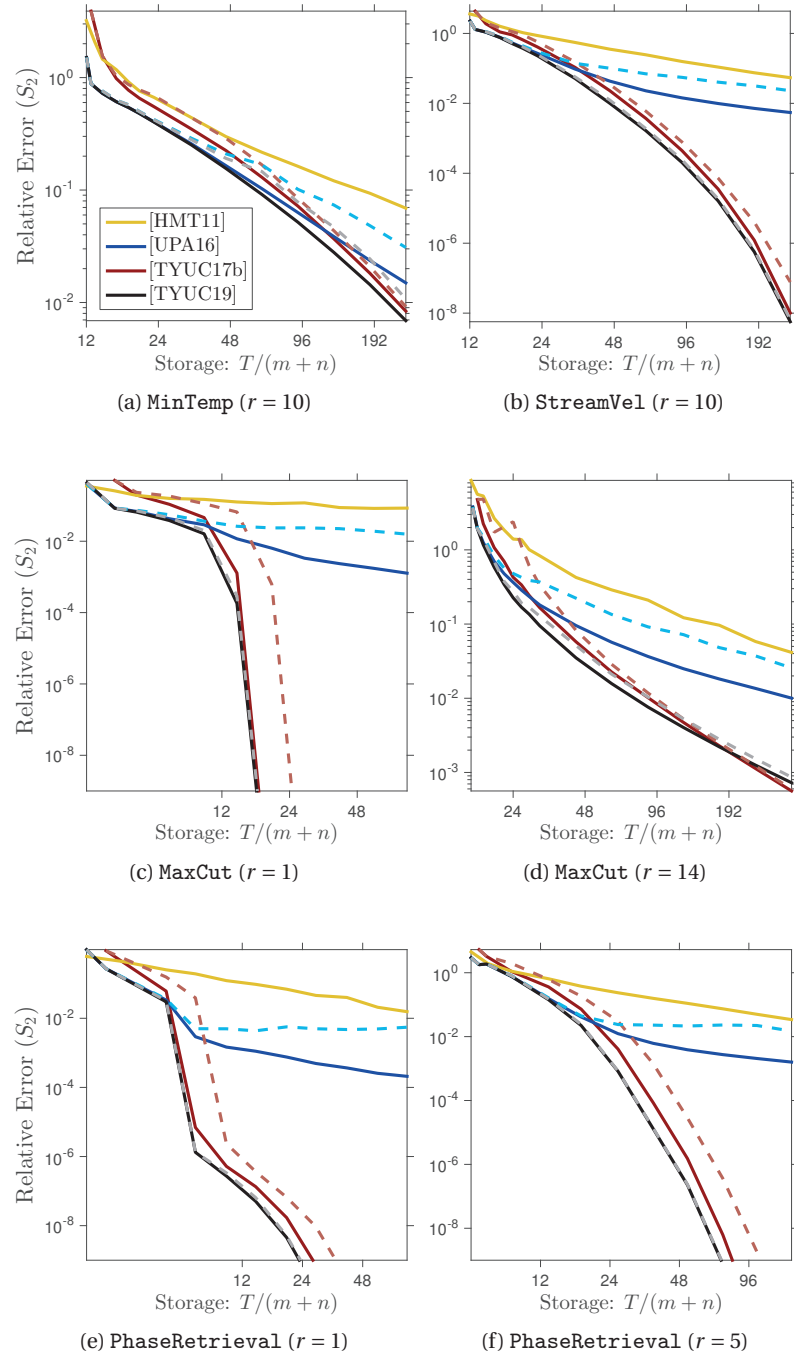


Figure 6.3 – [Comparison of reconstruction formulas: Real data examples.] We compare the relative error achieved by different fixed-rank approximation methods. [Solid lines] are oracle errors; [dashed lines] are errors with “natural” parameter choices. (There is no dashed line for [HMT11].) See Section 5.4 in [TYUC19] for the theoretical guidance on sketch-size parameters.

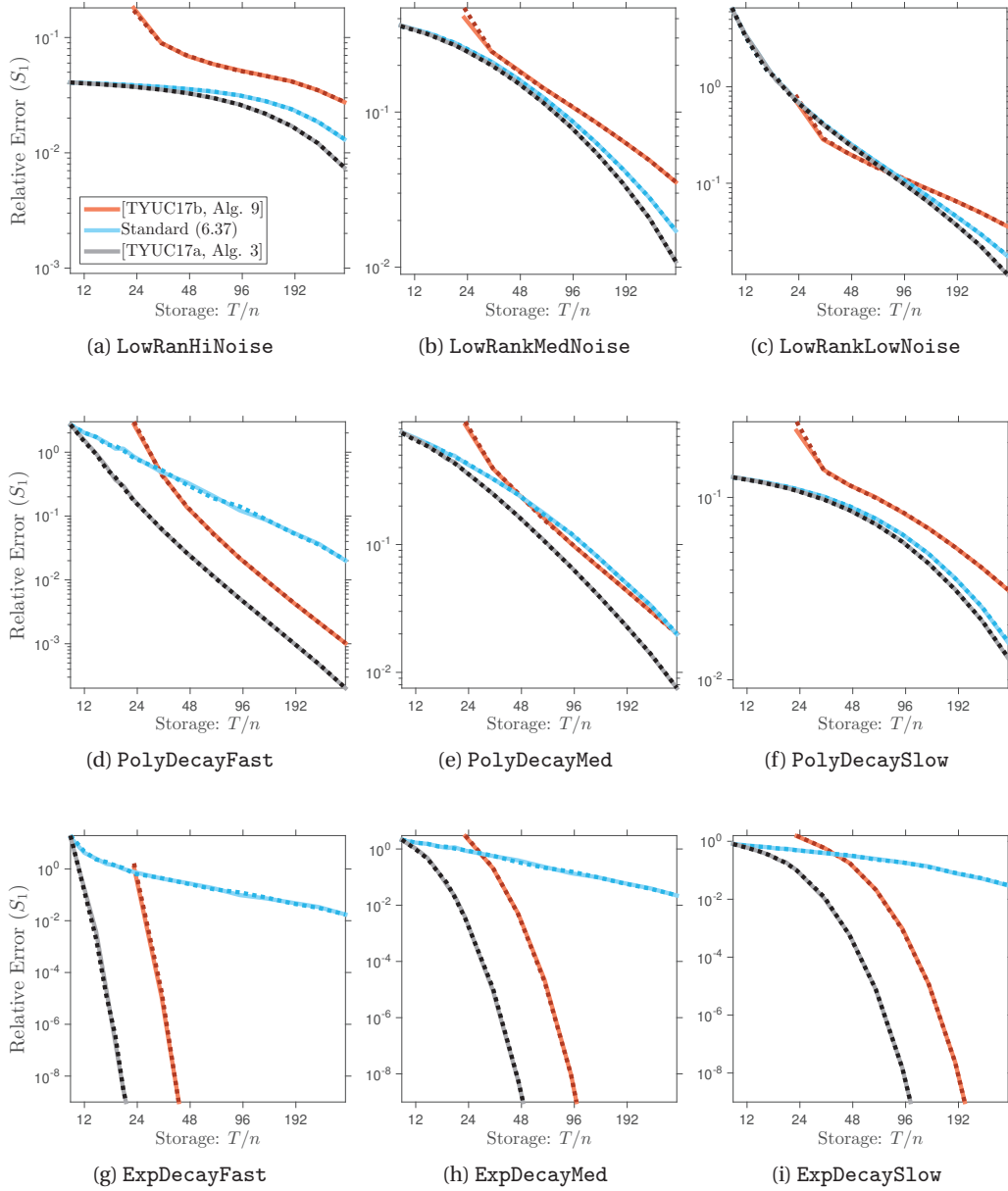


Figure 6.4 – [Comparison of fixed-rank PSD reconstruction formulas: Synthetic examples.] The data series show the performance of three algorithms for rank- $r$  PSD approximation with  $r = 10$ . [Solid lines] are generated from the Gaussian sketch; [dashed lines] are from the scrambled subsampled randomized Fourier transform (SSRFT) sketch. See Section 3 in [TYUC19] for a discussion on the different randomized linear dimension reduction maps.

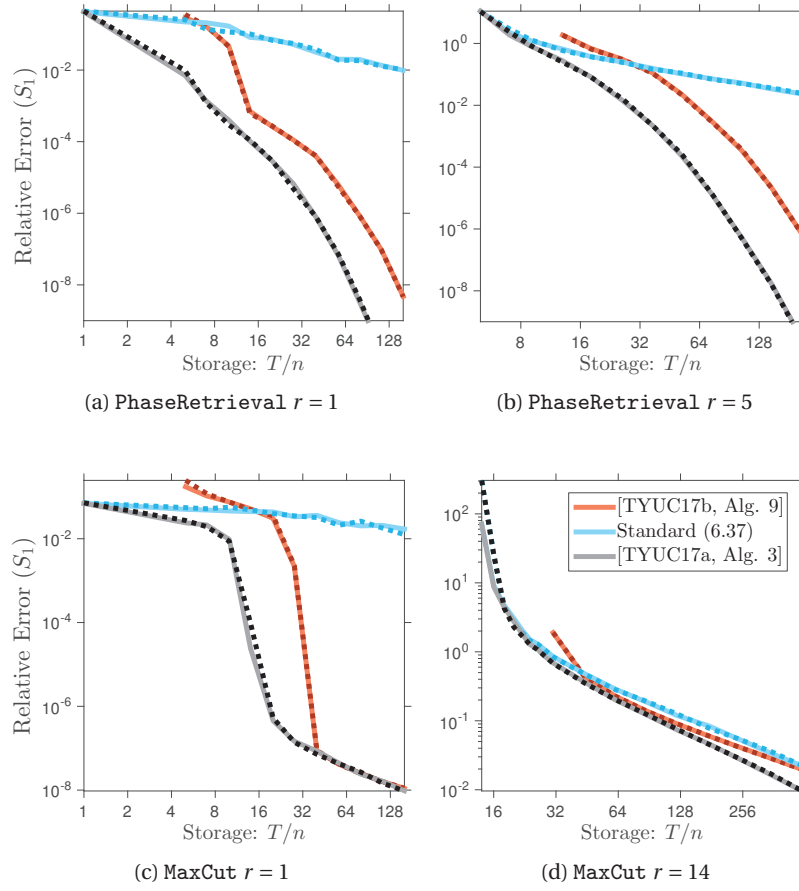


Figure 6.5 – [Comparison of fixed-rank PSD reconstruction formulas: Real data examples.] The data series show the performance of three algorithms for rank- $r$  psd approximation. [Solid lines] are generated from the Gaussian sketch; [dashed lines] are from the scrambled subsampled randomized Fourier transform (SSRFT) sketch. See Section 3 in [TYUC19] for a discussion on different randomized linear dimension reduction maps.

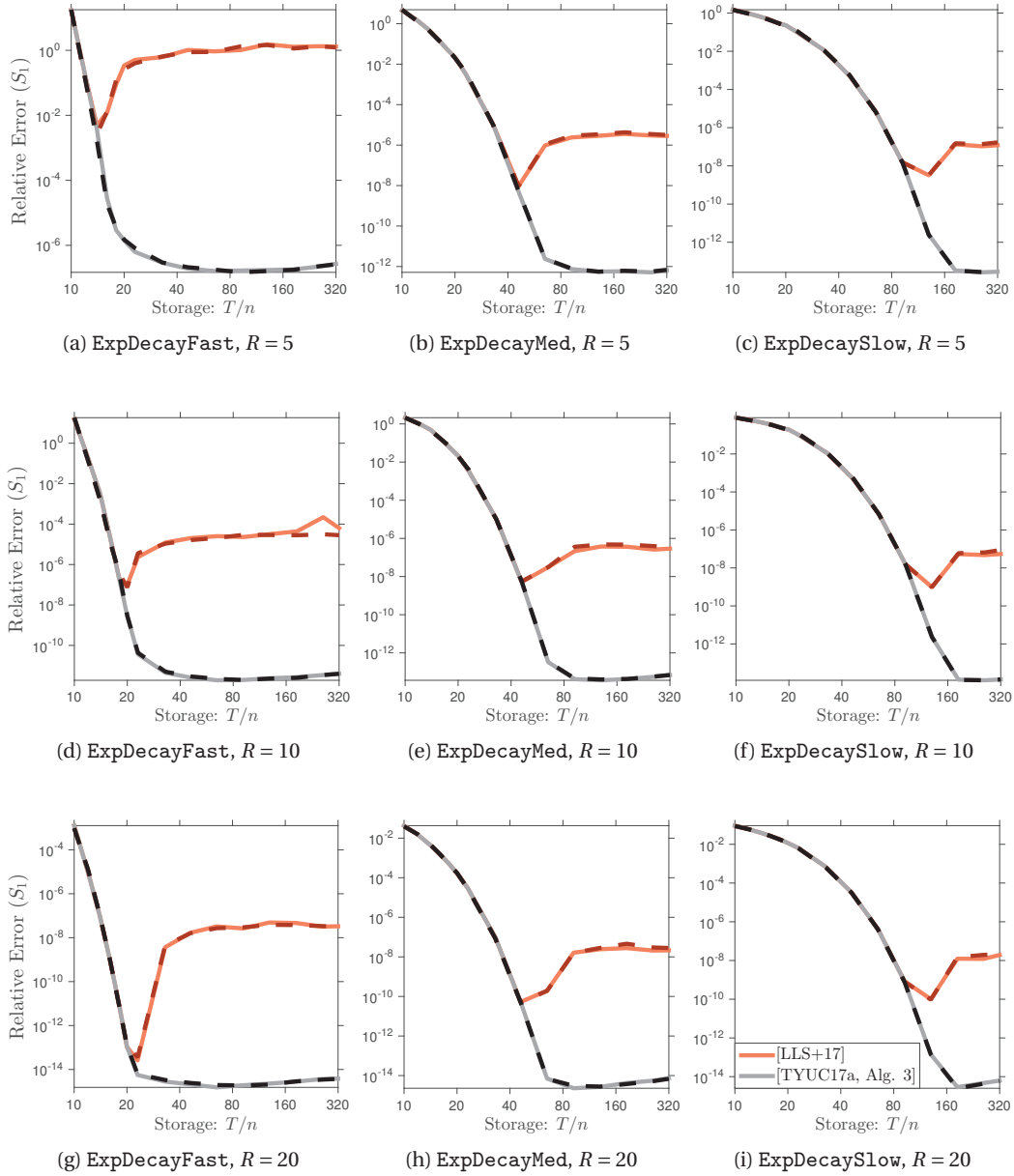


Figure 6.6 – [Numerical stability in implementing Nyström approximation formulas.] The series are generated by two implementations of the fixed-rank PSD approximation (6.33). We compare our implementation with another approach [LLS+17] proposed at Eqn. (13) in [LLS<sup>+</sup>17]. (Approximation Rank  $r = 10$ ) [Solid lines] are generated from the Gaussian sketch; [dashed lines] are from the scrambled subsampled randomized Fourier transform (SSRFT) sketch. See Section 3 in [TYUC19] for a discussion on different randomized linear dimension reduction maps.

# 7 CGM with Stochastic Path-Integrated Differential Estimator

So far, we have discussed the problems where data is readily available. Unfortunately, most problems in engineering and machine learning contain randomness to some degree. For instance, the data might be distributed over time and space without being accessible in its entirety at any given time instance.

Extending our results to the stochastic setting is an important piece of future work. This chapter aims to identify and present stochastic CGM variants with the best oracle complexity guarantees as a preliminary study for this direction. The work was done jointly with Suvrit Sra and Volkan Cevher [YSC19].

## Introduction

We study two problem settings in this chapter: The stochastic expectation minimization template, and the so-called finite-sum setting which is an important special case of the former.

$$\underset{\mathbf{x} \in \mathcal{X}}{\text{minimize}} \quad F(\mathbf{x}) := \begin{cases} \mathbb{E}_{\xi} f(\mathbf{x}, \xi) & \text{(expectation)} \\ \frac{1}{n} \sum_{i=1}^n f_i(\mathbf{x}) & \text{(finite-sum)} \end{cases} \quad (7.1)$$

- ▷  $\mathcal{X} \subset \mathbb{R}^p$  is the convex and compact domain;
- ▷  $F$ ,  $f$  and  $f_i$  are differentiable (and *possibly non-convex*);
- ▷  $\xi \sim \mathcal{P}$  is a random variable, supported on  $\Xi \subseteq \mathbb{R}^q$ .

The expectation objective covers a large number of applications in machine learning and statistics. The finite-sum template frequently arises in M-estimation and empirical risk minimization problems. Template (7.1) can be solved by using the well-known projected stochastic gradient descent method (SGD). At each iteration, SGD takes a stochastic gradient step followed by a projection to ensure feasibility of the new point. However, in many applications, projection onto  $\mathcal{X}$  can impose a computational bottleneck (*e.g.*, projection onto the nuclear

norm-ball may require a full singular value decomposition), or it can be even intractable (*e.g.*, dual structural SVMs [LJJSP13]).

There are many applications for stochastic conditional gradient methods, both in convex and non-convex settings. This includes low-rank matrix and tensor factorizations, structured sparse matrix estimation, dictionary learning applications, multi-class classification (considered as a motivating example in [HL16]), constrained deep learning problems (*e.g.*, [RDLS18] presents an application in computer vision), and many more.

To this end, we focus on the theoretical complexity of CGM-type algorithms for stochastic and finite-sum setting, for both convex and non-convex objectives. In the following sections, we identify and present the tightest results known so far.

### Contributions

- We propose a class of novel variance-reduced stochastic optimization algorithms, based on the recent *stochastic path-integrated differential estimator* technique (SPIDER) [FLLZ18]. By combining SPIDER with CGM, we introduce SPIDER-FW (FW stands for Frank-Wolfe).
- We extend our framework with the conditional gradient sliding (CGS) technique [LZ16], and propose SPIDER-CGS. This method further reduces the stochastic first-order oracle complexity of SPIDER-FW in various settings.
- For the non-convex setting, we present a new and compact proof for CGS along with its extension to the stochastic setting. To the best of our knowledge, our analysis is the first to provide convergence results for CGS in terms of the FW-gap (FW-gap is a widely used measure for the convergence of CGM-type methods, see Section 7.2). Although CGS does not seem to provide any improvement upon FW in this setup, we use the proof technique to extend SPIDER-CGS for the non-convex settings
- Finally, we use a learning rate schedule which depends on the global iteration counter, which in turn enables us to provide induction-free proofs. This approach differs from the majority of variance reduced FW methods, as they rely on an induction with respect to the outer loop counter, along with a sufficient improvement condition for each epoch.

**Notation.** We use the notation  $[n] = \{1, 2, \dots, n\}$ .

### 7.1 Related Works

CGM literature in the stochastic setting is much younger compared to the projection-based stochastic gradient methods. We can trace it back to a variant for online learning proposed by Hazan & Kale [HK12]. More recently, Hazan & Luo [HL16] introduced stochastic FW methods with and without variance reduction for finite-sum problems. Very recently, Mokhtari et al. [MHK18] proposed an alternative scheme for the expectation minimization setting.

CGM variants for non-convex stochastic learning are relatively understudied, with the most of the known results being due to Reddi et al. [RSPS16]. We discuss more details on the theoretical aspects of all these CGM variants in Section 7.5.

### 7.1.1 Conditional Gradient Sliding

Lan & Zhou [LZ16] have recently developed the conditional gradient sliding (CGS) method, which uses a modified version of the Nesterov’s acceleration scheme [Nes83] where the projection subproblems are solved inexactly using CGM. In other words, CGS establishes the convergence of an inexact version of the accelerated gradient method. CGS has superior first-order oracle complexity compared to CGM, although they have the same *lmo* complexity. We discuss more details and variants of CGS in Section 7.5.

### 7.1.2 Stochastic Path-Integrated Differential Estimator

There has been an extensive research on the variance reduced stochastic optimization methods, in order to address the needs of big data applications in machine learning. Therefore, various variance reduction techniques are proposed in the last few years such as SAG [RSB12], SVRG [JZ13], SAGA [DBLJ14], and more recently SARA [NLST17], SNVRG [ZXG18], and SPIDER [FLLZ18].

Among these approaches, SARA and SPIDER are closely related, since they use the same sequential update rule for the gradient estimator  $\mathbf{v}^k$ :

$$\mathbf{v}^k = \nabla f_S(\mathbf{x}^k) - \nabla f_S(\mathbf{x}^{k-1}) + \mathbf{v}^{k-1}. \quad (7.2)$$

However, SARA uses this estimator in the classical gradient descent template, while SPIDER adopts a normalized gradient approach. As a consequence, the results and the analyses differ.

As described by Wang et al. [WJZ<sup>+</sup>18], the original SPIDER framework has a restrictive step-size (proportional with the target accuracy  $\varepsilon$ ), which makes the algorithm impractical despite its theoretical appeal. Fortunately, this problem does not translate into the CGM-type analysis.

## 7.2 Preliminaries

As we consider a new template in this chapter, we need to renew some of the definitions we made in Section 1.7, such as the solution set and  $\varepsilon$ -solution.

**Solution.** We denote a solution of problem (7.1) by  $\mathbf{x}_\star$ . Similarly,  $F_\star$  respectively represents the optimal value.

$$\mathbf{x}_\star \in \operatorname{argmin}_{\mathbf{x} \in \mathcal{X}} F(\mathbf{x}) \quad \text{and} \quad F_\star = F(\mathbf{x}_\star). \quad (7.3)$$

**Measure of non-stationarity.** For unconstrained non-convex problems, the typical measure of non-stationarity is the gradient norm, because  $\|\nabla f(\mathbf{x})\| \rightarrow 0$  as  $\mathbf{x}$  converges to a stationary point. However, this measure cannot be used for constrained problems, because  $\|\nabla f(\mathbf{x})\|$  might not converge to 0 when we approach a solution on the boundary.

Instead, we will use a quantity called FW-gap, which naturally appears in the analysis of CGM-type algorithms.

$$\mathcal{G}(\mathbf{x}) := \max_{\mathbf{u} \in \mathcal{X}} \langle \mathbf{u} - \mathbf{x}, -\nabla F(\mathbf{x}) \rangle. \quad (7.4)$$

The FW-gap is always non-negative, and it gets 0 if and only if we are looking at a stationary point or a solution. As a consequence, FW-gap is a meaningful measure of non-stationarity, and it is widely used for CGM-type algorithms, see [LJ16], [RSPS16] and the references therein.

**$\varepsilon$ -solution.** Due to the fundamental difference in the measure of non-stationarity, we use different definitions of approximate solutions for convex and non-convex problems:

▷ If  $F$  is *convex*, we say  $\mathbf{x}_\varepsilon \in \mathcal{X}$  is an  $\varepsilon$ -solution if

$$F(\mathbf{x}_\varepsilon) - F_\star \leq \varepsilon. \quad (7.5)$$

▷ If  $F$  is *non-convex*, we say that a random variable  $\mathbf{x}_\varepsilon$  chosen uniformly from a finite set of points  $\{\mathbf{x}^1, \mathbf{x}^2, \dots, \mathbf{x}^t\}$  is an  $\varepsilon$ -solution if

$$\mathbb{E}[\mathcal{G}(\mathbf{x}_\varepsilon)] \leq \varepsilon. \quad (7.6)$$

It is common to provide convergence guarantees in expectation for a randomly chosen iterate in the non-convex setting. See [RSPS16].

**Oracle models.** We adopt the following black-box oracle model from [RSPS16], to establish a ground for comparing the convergence speed of different algorithms:

- Stochastic first-order oracle (*sfo*)  
For a stochastic function  $\mathbb{E}_\xi f(\cdot, \xi)$  with  $\xi \sim \mathcal{P}$ , *sfo* returns a pair  $(f(\mathbf{x}, \xi'), \nabla f(\mathbf{x}, \xi'))$  where  $\xi'$  is an *iid* sample from  $\mathcal{P}$ . [NY83]
- Incremental first-order oracle (*ifo*)  
For a finite-sum, *ifo* takes an index  $i \in [n]$  and returns  $(f_i(\mathbf{x}), \nabla f_i(\mathbf{x}))$ . [AB15]
- Linear minimization oracle (*lmo*)  
Well-known oracle of FW-type methods.

**Assumptions (finite-sum).** For the finite-sum setting, we assume that  $f_i(\mathbf{x})$  has an averaged  $L$ -Lipschitz gradient:

$$\mathbb{E}\|\nabla f_i(\mathbf{x}) - \nabla f_i(\mathbf{y})\|^2 \leq L^2 \|\mathbf{x} - \mathbf{y}\|^2, \quad \forall (\mathbf{x}, \mathbf{y}) \in \mathcal{X}^2. \quad (7.7)$$



This also implies  $F$  is  $L$ -smooth, since

$$\|\nabla F(\mathbf{x}) - \nabla F(\mathbf{y})\|^2 = \|\mathbb{E}(\nabla f_i(\mathbf{x}) - \nabla f_i(\mathbf{y}))\|^2 \leq \mathbb{E}\|\nabla f_i(\mathbf{x}) - \nabla f_i(\mathbf{y})\|^2 \leq L^2 \|\mathbf{x} - \mathbf{y}\|^2. \quad (7.8)$$

**Assumptions (expectation).** For the expectation minimization, we assume that  $\nabla f(\mathbf{x}, \xi)$  is an unbiased estimate of the gradient:

$$\mathbb{E}\nabla f(\mathbf{x}, \xi) = \nabla F(\mathbf{x}). \quad (7.9)$$

We also assume that the variance is bounded:

$$\mathbb{E}\|\nabla f(\mathbf{x}, \xi) - \nabla F(\mathbf{x})\|^2 \leq \sigma^2 < \infty, \quad \forall \xi \in \Xi, \forall \mathbf{x} \in \mathcal{X}. \quad (7.10)$$

And finally, we assume an averaged  $L$ -Lipschitz gradient condition, *i.e.*, the following condition holds  $\forall \xi \in \Xi$ :

$$\mathbb{E}\|\nabla f(\mathbf{x}, \xi) - \nabla f(\mathbf{y}, \xi)\|^2 \leq L^2 \|\mathbf{x} - \mathbf{y}\|^2, \quad \forall (\mathbf{x}, \mathbf{y}) \in \mathcal{X}^2. \quad (7.11)$$

Similar to the finite-sum, this implies the smoothness of  $F$ .

**Assumptions (non-convex).** Let us denote the initial point by  $\bar{\mathbf{x}}^1$ . Initial suboptimality  $F(\bar{\mathbf{x}}^1) - F_\star$  arises in the convergence bounds for the non-convex setting. For notational convenience, we denote an upper bound on this term by  $\mathcal{E}$ :

$$F(\bar{\mathbf{x}}^1) - F_\star \leq \mathcal{E} \quad (7.12)$$

Suppose that  $F_\star$  is finite, then there exists a finite  $\mathcal{E}$  which satisfies this bound. This is a direct consequence of the smoothness of  $F$  and the boundedness of domain.

All these assumptions are mild and frequently used in the analysis of stochastic methods and CGM-type algorithms.

## 7.3 SPIDER Frank-Wolfe

This section presents SPIDER-FW algorithm and its convergence guarantees under various problem settings.

Our methods have a double loop structure, hence the iterates and the parameters have two different iteration counters  $t$  and  $k$ , such as  $\mathbf{x}^{t,k}$ . For notational brevity, we drop the first counter when there is no ambiguity, such as  $\mathbf{x}^k$ . Throughout,  $s_{t,k}$  denotes the global iteration counter, indexed by the  $k^{\text{th}}$  iteration of the  $t^{\text{th}}$  epoch. In our pseudocodes, *draw samples* means *iid* samples for expectation minimization, and uniform selection with replacement in the finite-sum setting.

---

**Algorithm 7.1** SPIDER Frank-Wolfe

---

```

1: Input:  $\bar{\mathbf{x}}^1 \in \mathcal{X}$ 
2: for  $t = 1, 2, \dots, T$  do
3:   Set  $\mathbf{x}^1 = \bar{\mathbf{x}}^t$ 
4:   Draw  $Q_t$  samples  $\mathcal{Q}_t$ 
5:   Compute  $\mathbf{v}^1 = \nabla f_{\mathcal{Q}_t}(\mathbf{x}^1)$ 
6:   Compute  $\mathbf{w}^1 = \text{lmo}_{\mathcal{X}}(\mathbf{v}^1)$ 
7:   Update  $\mathbf{x}^2 = \mathbf{x}^1 + \eta_{t,1}(\mathbf{w}^1 - \mathbf{x}^1)$ 
8:   for  $k = 2, 3, \dots, K_t$  do
9:     Draw  $S_{t,k}$  samples  $\mathcal{S}_{t,k}$ 
10:    Compute  $\mathbf{v}^k = \nabla f_{\mathcal{S}_{t,k}}(\mathbf{x}^k) - \nabla f_{\mathcal{S}_{t,k}}(\mathbf{x}^{k-1}) + \mathbf{v}^{k-1}$ 
11:    Compute  $\mathbf{w}^k \in \text{lmo}_{\mathcal{X}}(\mathbf{v}^k)$ 
12:    Update  $\mathbf{x}^{k+1} = \mathbf{x}^k + \eta_{t,k}(\mathbf{w}^k - \mathbf{x}^k)$ 
13:   end for
14:   Set  $\bar{\mathbf{x}}^{t+1} = \mathbf{x}^{K_t+1}$ 
15: end for

```

---

### 7.3.1 Convex Finite-Sum

We consider SPIDER-FW with

$$K_t = 2^{t-1} \text{ for } t = 1, 2, \dots, T.$$

We choose the sampling parameters and the learning rate parameter

$$S_{t,k} = K_t \quad \mathcal{Q}_t = [n] \quad \eta_{t,k} = \frac{2}{s_{t,k} + 1} \text{ where } s_{t,k} = K_t + k - 1.$$

**Theorem 7.1.** *Consider the convex finite-sum optimization template, and suppose that the assumptions in Section 7.2 for this template hold. Then, estimate  $\mathbf{x}^{t,k}$  of SPIDER-FW with the parameter setting described above satisfies*

$$\mathbb{E}[F(\mathbf{x}^{t,k})] - F_\star = \mathcal{O}\left(\frac{LD^2}{s_{t,k}}\right) \quad (7.13)$$

**Corollary 7.1.** The *ifo* and *lmo* complexities of SPIDER-FW for achieving  $\epsilon$ -solution in this setting are as follows:

$$\begin{aligned} \#(\text{ifo}) &= \mathcal{O}\left(n \ln\left(\frac{LD^2}{\epsilon}\right) + \frac{L^2 D^4}{\epsilon^2}\right) \\ \#(\text{lmo}) &= \mathcal{O}\left(\frac{LD^2}{\epsilon}\right) \end{aligned} \quad (7.14)$$

### 7.3.2 Convex Expectation Minimization

We consider SPIDER-FW with

$$K_t = 2^{t-1} \text{ for } t = 1, 2, \dots, T.$$

We choose the sampling parameters and the learning rate parameter as

$$S_{t,k} = K_t \quad Q_t = \lceil \frac{\sigma^2 K_t^2}{5L^2 D^2} \rceil \quad \eta_{t,k} = \frac{2}{s_{t,k} + 1} \text{ where } s_{t,k} = K_t + k - 1.$$

**Theorem 7.2.** *Consider the convex expectation minimization template, and suppose that the assumptions in Section 7.2 for this template hold. Then, estimate  $\mathbf{x}^{t,k}$  of SPIDER-FW with the parameter setting described above satisfies*

$$\mathbb{E}[F(\mathbf{x}^{t,k})] - F_\star = \mathcal{O}\left(\frac{LD^2}{s_{t,k}}\right) \quad (7.15)$$

**Corollary 7.2.** The *sfo* and *lmo* complexities of SPIDER-FW for achieving  $\varepsilon$ -solution in this setting are as follows:

$$\begin{aligned} \#(sfo) &= \mathcal{O}\left(\frac{\sigma^2 D^2 + L^2 D^4}{\varepsilon^2}\right) \\ \#(lmo) &= \mathcal{O}\left(\frac{LD^2}{\varepsilon}\right) \end{aligned} \quad (7.16)$$

SPIDER-FW has the same asymptotic oracle complexities as SCGS [LZ16] in this setting.

### 7.3.3 Non-convex Finite-Sum

We consider SPIDER-FW with

$$K_t = K = \lceil \sqrt{n} \rceil.$$

We choose the parameters and the learning rate parameter

$$S_{t,k} = S = \lceil \sqrt{n} \rceil \quad Q_t = [n] \quad \eta_{t,k} = \eta = \frac{1}{\sqrt{s_{T,K}}} \text{ where } s_{T,K} = TK.$$

**Theorem 7.3.** *Consider the non-convex finite-sum template, and suppose that the assumptions in Section 7.2 for this template hold. Denote by  $\mathbf{x}^{out}$  an iterate  $\mathbf{x}^{t,k}$  of SPIDER-FW chosen uniformly random over all  $(t, k)$  pairs up to  $(T, K)$ . Then, the following bound on the FW-gap holds:*

$$\mathbb{E}[\mathcal{G}(\mathbf{x}^{out})] = \mathcal{O}\left(\frac{\mathcal{E} + LD^2}{\sqrt{s_{T,K}}}\right) \quad (7.17)$$

Although it is impractical to store all estimates until the final iteration, all stochastic methods for the non-convex setting shown in Table 7.1 have this type convergence guarantees (see [RSPS16]). More stringently, [LJ16] shows convergence of the non-convex FW in terms of the running best iterate. In the stochastic setting, however, we cannot keep track of the best estimate, simply because we cannot measure the FW-gap.

**Corollary 7.3.** The *ifo* and *lmo* complexities of SPIDER-FW for achieving  $\varepsilon$ -solution in the non-convex finite-sum setting are as follows:

$$\begin{aligned} \#(ifo) &= \mathcal{O}\left(\frac{\sqrt{n}}{\varepsilon^2}(\mathcal{E}^2 + L^2 D^4)\right) \\ \#(lmo) &= \mathcal{O}\left(\frac{1}{\varepsilon^2}(\mathcal{E}^2 + L^2 D^4)\right) \end{aligned} \quad (7.18)$$

SPIDER-FW has better *ifo* complexity than state-of-the-art in the non-convex finite-sum setting, and it improves the dependence on sample size  $n$ . See Table 7.1 for comparison.

#### 7.3.4 Non-convex Expectation Minimization

We consider SPIDER-FW with

$$K_t = K = \lceil \sigma / \varepsilon \rceil.$$

We choose the parameters and the learning rate parameter

$$S_{t,k} = S = \lceil \sigma / \varepsilon \rceil \quad Q_t = Q = \lceil 4(\sigma / \varepsilon)^2 \rceil \quad \eta_{t,k} = \eta = \frac{1}{\sqrt{s_{T,K}}} \text{ where } s_{T,K} = TK.$$

**Theorem 7.4.** Consider the non-convex expectation minimization template, and suppose that the assumptions in Section 7.2 for this template hold. Denote by  $\mathbf{x}^{out}$  an iterate  $\mathbf{x}^{t,k}$  of SPIDER-FW chosen uniformly random over all  $(t, k)$  pairs up to  $(T, K)$ . Then, the following bound holds:

$$\mathbb{E}[\mathcal{G}(\mathbf{x}^{out})] = \mathcal{O}\left(\frac{\mathcal{E} + LD^2}{\sqrt{s_{T,K}}}\right) + \frac{\varepsilon}{2}. \quad (7.19)$$

**Corollary 7.4.** The *sfo* and *lmo* complexities of SPIDER-FW for achieving  $\varepsilon$ -solution in this setting are as follows:

$$\begin{aligned} \#(sfo) &= \mathcal{O}\left(\frac{\sigma}{\varepsilon^3}(\mathcal{E}^2 + L^2 D^4)\right) \\ \#(lmo) &= \mathcal{O}\left(\frac{1}{\varepsilon^2}(\mathcal{E}^2 + L^2 D^4)\right) \end{aligned} \quad (7.20)$$

Once again, SPIDER-FW enjoys superior *sfo* complexity while maintaining the same *lmo* complexity as its competitors. SVRF was the state-of-the-art with  $\mathcal{O}(\varepsilon^{-10/3})$ , see [RSPS16].

## 7.4 SPIDER Conditional Gradient Sliding

This section presents SPIDER-CGS (as shown in Algorithm 7.2) and its convergence guarantees for various settings.

### 7.4.1 Convex Finite-Sum

We consider SPIDER-CGS with

$$K_t = \lceil 2^{t/2} \rceil \text{ for } t = 1, 2, \dots, T.$$

We choose the sampling parameters and the learning rate parameter as

$$S_{t,k} = 9K_t s_{t,K_t}^2 \quad \mathcal{Q}_t = [n] \quad \gamma_{t,k} = \frac{3}{s_{t,k} + 2} \text{ where } s_{t,k} = \sum_{\tau=1}^{t-1} K_\tau + k.$$

Furthermore, we choose the CndG subsolver parameters as

$$\beta_{t,k} = \frac{3}{2} L \gamma_{t,k} \quad \alpha_{t,k} = \frac{2LD^2}{(s_{t,k} + 1)^2}.$$

**Theorem 7.5.** *Consider the convex finite-sum template, and suppose that the assumptions in Section 7.2 for this template hold. Then, estimate  $\mathbf{y}^{t,k}$  of SPIDER-CGS with the parameter setting described above satisfies*

$$\mathbb{E}[F(\mathbf{y}^{t,k}) - F_\star] = \mathcal{O}\left(\frac{LD^2}{s_{t,k}^2}\right) \quad (7.21)$$

**Corollary 7.5.** The *ifo* and *lmo* complexities of SPIDER-CGS for achieving  $\varepsilon$ -solution in this template are as follows:

$$\begin{aligned} \#(ifo) &= \mathcal{O}\left(n \ln\left(\frac{LD^2}{\varepsilon}\right) + \frac{L^2 D^4}{\varepsilon^2}\right) \\ \#(lmo) &= \mathcal{O}\left(\frac{LD^2}{\varepsilon}\right) \end{aligned} \quad (7.22)$$

---

**Algorithm 7.2** SPIDER Conditional Gradient Sliding

---

```

1: Input:  $\bar{\mathbf{x}}^1 = \bar{\mathbf{y}}^1 \in \mathcal{X}$ 
2: for  $t = 1, 2, \dots, T$  do
3:   Set  $\mathbf{x}^1 = \bar{\mathbf{x}}^t$  and  $\mathbf{y}^1 = \bar{\mathbf{y}}^t$ 
4:   Update  $\mathbf{z}^1 = \mathbf{y}^1 + \gamma_{t,1}(\mathbf{x}^1 - \mathbf{y}^1)$ 
5:   Draw  $Q_t$  samples  $\mathcal{Q}_t$ 
6:   Compute  $\mathbf{v}^1 = \nabla f_{\mathcal{Q}_t}(\mathbf{z}^1)$ 
7:    $\mathbf{x}^2 = \text{CNDG}(\mathbf{x}^1, \mathbf{v}^1, \alpha_{t,1}, \beta_{t,1})$ 
8:   Update  $\mathbf{y}^2 = \mathbf{y}^1 + \gamma_{t,1}(\mathbf{x}^2 - \mathbf{y}^1)$ 
9:   for  $k = 2, 3, \dots, K_t$  do
10:    Update  $\mathbf{z}^k = \mathbf{y}^k + \gamma_{t,k}(\mathbf{x}^k - \mathbf{y}^k)$ 
11:    Draw  $S_{t,k}$  samples  $\mathcal{S}_{t,k}$ 
12:    Compute  $\mathbf{v}^k = \nabla f_{\mathcal{S}_{t,k}}(\mathbf{z}^k) - \nabla f_{\mathcal{S}_{t,k}}(\mathbf{z}^{k-1}) + \mathbf{v}^{k-1}$ 
13:     $\mathbf{x}^{k+1} = \text{CNDG}(\mathbf{x}^k, \mathbf{v}^k, \alpha_{t,k}, \beta_{t,k})$ 
14:    Update  $\mathbf{y}^{k+1} = \mathbf{y}^k + \gamma_{t,k}(\mathbf{x}^{k+1} - \mathbf{y}^k)$ 
15:   end for
16:   Set  $\bar{\mathbf{x}}^{t+1} = \mathbf{x}^{K_t+1}$  and  $\bar{\mathbf{y}}^{t+1} = \mathbf{y}^{K_t+1}$ 
17: end for

```

---

```

18: function  $\mathbf{u}^+ = \text{CNDG}(\mathbf{u}, \mathbf{v}, \alpha, \beta)$ 
19:   Set  $\mathbf{u}^1 = \mathbf{u}$ 
20:   for  $k = 1, 2, \dots$  do
21:     Compute  $\mathbf{w}^k = \text{lmo}_{\mathcal{X}}(\mathbf{v} + \beta(\mathbf{u}^k - \mathbf{u}))$ 
22:     Evaluate  $\zeta_k = \langle \mathbf{v} + \beta(\mathbf{u}^k - \mathbf{u}), \mathbf{u}^k - \mathbf{w}^k \rangle$ 
23:     if  $\zeta_k \leq \alpha$  then
24:       break
25:     end if
26:     Set  $\theta_k = \min\{1, \zeta_k / (\beta \|\mathbf{w}^k - \mathbf{u}^k\|^2)\}$ 
27:     Update  $\mathbf{u}^{k+1} = \mathbf{u}^k + \theta_k(\mathbf{w}^k - \mathbf{u}^k)$ 
28:   end for
29:   Set  $\mathbf{u}^+ = \mathbf{u}^k$ 
30: end function

```

---

### 7.4.2 Convex Expectation Minimization

We consider SPIDER-CGS with

$$K_t = \lceil 2^{t/2} \rceil \text{ for } t = 1, 2, \dots, T.$$

We choose the sampling parameters and the learning rate parameter as

$$S_{t,k} = 9K_t s_{t,K_t}^2 \quad Q_t = \lceil \frac{\sigma^2 s_{t,K_t}^4}{L^2 D^2} \rceil \quad \gamma_{t,k} = \frac{2}{s_{t,k} + 1} \text{ where } s_{t,k} = \sum_{\tau=1}^{t-1} K_\tau + k.$$

Furthermore, we choose the CndG subsolver parameters as

$$\beta_{t,k} = \frac{3}{2} L \gamma_{t,k} \quad \alpha_{t,k} = \frac{2LD^2}{(s_{t,k} + 1)^2}.$$

**Theorem 7.6.** *Consider the convex expectation minimization template, and suppose that the assumptions in Section 7.2 for this template hold. Then, estimate  $\mathbf{y}^{t,k}$  of SPIDER-CGS with the parameter setting described above satisfies*

$$\mathbb{E}[F(\mathbf{y}^{t,k}) - F_\star] = \mathcal{O}\left(\frac{LD^2}{s_{t,k}^2}\right) \quad (7.23)$$

**Corollary 7.6.** The *sfo* and *lmo* complexities of SPIDER-CGS for achieving  $\varepsilon$ -solution in convex expectation minimization problems are as follows:

$$\begin{aligned} \#(sfo) &= \mathcal{O}\left(\frac{\sigma^2 D^2 + L^2 D^4}{\varepsilon^2}\right) \\ \#(lmo) &= \mathcal{O}\left(\frac{LD^2}{\varepsilon}\right) \end{aligned} \quad (7.24)$$

### 7.4.3 Non-convex Finite-Sum

We consider SPIDER-CGS with

$$K_t = K = \lceil \sqrt{n} \rceil.$$

We choose the sampling parameters and the learning rate parameter as

$$S_{t,k} = K \quad Q_t = [n] \quad \gamma_{t,k} = \gamma = \frac{1}{\sqrt{s_{T,K}}} \text{ where } s_{T,K} = TK.$$

Furthermore, we choose the CndG subsolver parameters as

$$\beta_{t,k} = \frac{3}{2} L \gamma \quad \alpha_{t,k} = LD^2 \gamma.$$

**Theorem 7.7.** Consider the non-convex finite-sum template, and suppose that the assumptions in Section 7.2 for this template hold. Denote by  $\mathbf{y}^{out}$  an iterate  $\mathbf{y}^{t,k}$  of SPIDER-CGS chosen uniformly random over all  $(t, k)$  pairs up to  $(T, K)$ . Then, the following bound on the FW-gap holds:

$$\mathbb{E}[\mathcal{G}(\mathbf{y}^{out})] = \mathcal{O}\left(\frac{\mathcal{E} + LD^2}{\sqrt{s_{T,K}}}\right) \quad (7.25)$$

**Corollary 7.7.** The *ifo* and *lmo* complexities of SPIDER-CGS for achieving  $\varepsilon$ -solution in non-convex finite-sum are

$$\begin{aligned} \#(ifo) &= \mathcal{O}\left(\frac{\sqrt{n}}{\varepsilon^2}(\mathcal{E}^2 + L^2 D^4)\right) \\ \#(lmo) &= \mathcal{O}\left(\frac{1}{\varepsilon^2}(\mathcal{E}^2 + L^2 D^4)\right) \end{aligned} \quad (7.26)$$

#### 7.4.4 Non-convex Expectation Minimization

We consider SPIDER-CGS with

$$K_t = K = \lceil \sigma / \varepsilon \rceil.$$

We choose the sampling parameters and the learning rate parameter as

$$S_{t,k} = K \quad Q_t = \lceil 4(\sigma / \varepsilon)^2 \rceil \quad \gamma_{t,k} = \gamma = \frac{1}{\sqrt{s_{T,K}}} \text{ where } s_{T,K} = TK.$$

Furthermore, we choose the CndG subsolver parameters as

$$\beta_{t,k} = \frac{3}{2}L\gamma \quad \alpha_{t,k} = LD^2\gamma.$$

**Theorem 7.8.** Consider the non-convex expectation minimization template, and suppose that the assumptions in Section 7.2 for this template hold. Denote by  $\mathbf{y}^{out}$  an iterate  $\mathbf{y}^{t,k}$  of SPIDER-CGS chosen uniformly random over all  $(t, k)$  pairs up to  $(T, K)$ . Then, the following bound holds:

$$\mathbb{E}[\mathcal{G}(\mathbf{y}^{out})] = \mathcal{O}\left(\frac{\mathcal{E} + LD^2}{\sqrt{s_{T,K}}}\right) + \frac{\varepsilon}{2} \quad (7.27)$$

**Corollary 7.8.** The *sfo* and *lmo* complexities of SPIDER-CGS for achieving  $\varepsilon$ -solution in non-convex expectation minimization problems are as follows:

$$\begin{aligned} \#(sfo) &= \mathcal{O}\left(\frac{\sigma}{\varepsilon^3}(\mathcal{E}^2 + L^2 D^4)\right) \\ \#(lmo) &= \mathcal{O}\left(\frac{1}{\varepsilon^2}(\mathcal{E}^2 + L^2 D^4)\right) \end{aligned} \quad (7.28)$$



## 7.5. Comparison & Discussion

	convex				non-convex			
	finite-sum		expectation		finite-sum		expectation	
	(ifo)	(lmo)	(sfo)	(lmo)	(ifo)	(lmo)	(sfo)	(lmo)
FW	$\mathcal{O}(n\epsilon^{-1})$	$\mathcal{O}(\epsilon^{-1})$	-	-	$\mathcal{O}(n\epsilon^{-2})$	$\mathcal{O}(\epsilon^{-2})$	-	-
CGS	$\mathcal{O}(n\epsilon^{-1/2})$	$\mathcal{O}(\epsilon^{-1})$	-	-	$\mathcal{O}(n\epsilon^{-2})$	$\mathcal{O}(\epsilon^{-2})$	-	-
SFW	$\mathcal{O}(\epsilon^{-3})$	$\mathcal{O}(\epsilon^{-1})$	$\mathcal{O}(\epsilon^{-3})$	$\mathcal{O}(\epsilon^{-1})$	$\mathcal{O}(\epsilon^{-4})$	$\mathcal{O}(\epsilon^{-2})$	$\mathcal{O}(\epsilon^{-4})$	$\mathcal{O}(\epsilon^{-2})$
SFW-1	$\mathcal{O}(\epsilon^{-3})$	$\mathcal{O}(\epsilon^{-3})$	$\mathcal{O}(\epsilon^{-3})$	$\mathcal{O}(\epsilon^{-3})$	-	-	-	-
Online-FW	$\mathcal{O}(\epsilon^{-4})$	$\mathcal{O}(\epsilon^{-2})$	$\mathcal{O}(\epsilon^{-4})$	$\mathcal{O}(\epsilon^{-2})$	-	-	-	-
SCGS	$\mathcal{O}(\epsilon^{-2})$	$\mathcal{O}(\epsilon^{-1})$	$\mathcal{O}(\epsilon^{-2})$	$\mathcal{O}(\epsilon^{-1})$	$\mathcal{O}(\epsilon^{-4})$	$\mathcal{O}(\epsilon^{-2})$	$\mathcal{O}(\epsilon^{-4})$	$\mathcal{O}(\epsilon^{-2})$
SVRF / SVFW	$\mathcal{O}(n\ln(\epsilon^{-1}) + \epsilon^{-2})$	$\mathcal{O}(\epsilon^{-1})$	-	-	$\mathcal{O}(n + n^{2/3}\epsilon^{-2})$	$\mathcal{O}(\epsilon^{-2})$	$\mathcal{O}(\epsilon^{-10/3})$	$\mathcal{O}(\epsilon^{-2})$
STORC <sup>†</sup>	$\mathcal{O}(n\ln(\epsilon^{-1}) + \epsilon^{-3/2})$	$\mathcal{O}(\epsilon^{-1})$	-	-	-	-	-	-
<i>SPIDER-FW</i>	$\mathcal{O}(n\ln(\epsilon^{-1}) + \epsilon^{-2})$	$\mathcal{O}(\epsilon^{-1})$	$\mathcal{O}(\epsilon^{-2})$	$\mathcal{O}(\epsilon^{-1})$	$\mathcal{O}(n^{1/2}\epsilon^{-2})$	$\mathcal{O}(\epsilon^{-2})$	$\mathcal{O}(\epsilon^{-3})$	$\mathcal{O}(\epsilon^{-2})$
<i>SPIDER-CGS</i>	$\mathcal{O}(n\ln(\epsilon^{-1}) + \epsilon^{-2})$	$\mathcal{O}(\epsilon^{-1})$	$\mathcal{O}(\epsilon^{-2})$	$\mathcal{O}(\epsilon^{-1})$	$\mathcal{O}(n^{1/2}\epsilon^{-2})$	$\mathcal{O}(\epsilon^{-2})$	$\mathcal{O}(\epsilon^{-3})$	$\mathcal{O}(\epsilon^{-2})$

Table 7.1 – Comparison of conditional gradient methods for stochastic optimization. **Our contribution** is highlighted with red font. See Section 7.5 for more details.

FW [FW56, Jag13], CGS [LZ16], SFW [HL16, RSPS16], SFW-1 [MHK18], Online-FW [HK12], SCGS [LZ16], SVRF / SVFW [HL16, RSPS16], STORC [HL16]

## 7.5 Comparison & Discussion

Next, we present an extensive comparison of the theoretical aspects of the CGM-type algorithms. We compile a summary of this comparison in Table 7.1.

### 7.5.1 Convex Optimization Camp

**Batch setting.** CGM achieves an  $\epsilon$ -solution after  $\mathcal{O}(1/\epsilon)$  iterations. This iteration complexity is optimal for a large class of methods that construct the decision variable through convex combination of *lmo* outputs [Lan14]. CGS, on the other side, enjoys  $\mathcal{O}(1/\sqrt{\epsilon})$  first order oracle complexity while maintaining the optimal  $\mathcal{O}(1/\epsilon)$  *lmo* complexity, by reusing the same gradients over multiple iterations [LZ16].

**Stochastic setting.** Hazan & Kale [HK12] propose Online-FW for an online-learning setting. As mentioned later by Hazan & Luo [HL16], these results can be translated to the stochastic template, by means of a standard online to stochastic conversion approach. This conversion yields  $\mathcal{O}(1/\epsilon^4)$  *sfo* and  $\mathcal{O}(1/\epsilon^2)$  *lmo* complexities.

Algorithm 7.3 presents the most straightforward extension of CGM for the stochastic setting. This method, SFW, converges with  $\mathcal{O}(1/k)$  rate when the sample size  $S_k = \Theta(k^2)$  [HL16]. As a result, it provably achieves an  $\epsilon$ -solution after  $\mathcal{O}(1/\epsilon^3)$  *sfo* and  $\mathcal{O}(1/\epsilon)$  *lmo* calls.

Lan & Zhou introduce an extension of CGS for the stochastic setting (see Section 3 in [LZ16]). This variant, SCGS, achieves  $\mathcal{O}(1/\epsilon^2)$  *sfo* complexity, while maintaining the optimal  $\mathcal{O}(1/\epsilon)$  *lmo* complexity. Under strong convexity, the *lmo* complexity becomes  $\mathcal{O}(1/\epsilon)$ .

---

**Algorithm 7.3** Stochastic Frank-Wolfe

---

```

1: Input:  $\mathbf{x}^1 \in \mathcal{X}$ 
2: for  $k = 1, 2, \dots, K$  do
3:   Draw  $S_k$  samples  $\mathcal{S}_k$ 
4:   Compute  $\mathbf{w}^k = \text{lmo}_{\mathcal{X}}(\nabla f_{\mathcal{S}_k}(\mathbf{x}^k))$ 
5:   Update  $\mathbf{x}^{k+1} = \mathbf{x}^k + \eta_k(\mathbf{w}^k - \mathbf{x}^k)$ 
6: end for

```

---

Hazan & Luo introduce the stochastic variance reduced Frank-Wolfe method, SVRF [HL16]. This method adopts the variance reduction techniques from [JZ13] and [MZJ13]. SVRF is designed specifically for the finite-sum setting, and it requires  $\mathcal{O}(n \ln(1/\epsilon))$  calls of the full gradient oracle, besides  $\mathcal{O}(1/\epsilon^2)$  *ifo* and  $\mathcal{O}(1/\epsilon)$  *lmo* calls, to get an  $\epsilon$ -solution.

In order to improve the *ifo* complexity of the SVRF, Hazan & Luo [HL16] also design a variant based on the CGS. This variant, STORC, also requires  $\mathcal{O}(n \ln(1/\epsilon))$  calls to the full gradient oracle and  $\mathcal{O}(1/\epsilon)$  calls to *lmo*, but it enjoys a reduced number of *ifo* calls, at the order of  $\mathcal{O}(1/\epsilon^{1.5})$ . Compared to the SVRF, however, STORC requires an additional assumption on the Lipschitz continuity of  $F$ . STORC gets better guarantees under strong convexity assumption, however we omit the details.

Lu & Freund [LF18] propose a stochastic CGM variant with  $\mathcal{O}(1/\epsilon)$  *lmo* and  $\mathcal{O}(n + 1/\epsilon)$  *ifo* complexity for the convex finite-sum setting. However, the proposed method relies on a special structure of the objective function, where  $f_i$  are univariate functions of the fitted value  $\langle \mathbf{a}_i, \mathbf{x} \rangle$  for some given data sample  $\mathbf{a}_i$ .

All stochastic FW variants we discussed up to know are based on an increasing mini-batch size. Very recently, Mokhtari et al. [MHK18] have proposed an alternative scheme (SFW-1) for expectation minimization setting, which requires a single *sfo* at each iteration. Nevertheless, SFW-1 has an arguably worse computational complexity compared to SFW, with its  $\mathcal{O}(1/\epsilon^3)$  calls to *sfo* and *lmo*. We emphasize the applications of SFW-1 in submodular maximization, but this is beyond the scope of this dissertation.

For the convex finite-sum setting, SPIDER-FW and SPIDER-CGS share the same complexities as SVRF.

### 7.5.2 Non-convex Optimization Camp

**Batch setting.** CGM converges to a stationary point, see Section 2.2 in [Ber99] for details. This early result guarantees an asymptotical convergence, but does not provide any information about the convergence rate. To our knowledge, the first convergence rate for CGM in the non-convex setting is shown by Yu et al. [YZS14]. More recently, Lacoste-Julien has presented a non-asymptotic  $\mathcal{O}(1/\sqrt{k})$  rate in FW-gap [LJ16].

**Stochastic setting.** The majority of the results known in the literature for this setting are derived by Reddi et al. in [RSPS16]. The basic stochastic CGM extension, SFW, provably gets an  $\varepsilon$ -solution after  $\mathcal{O}(1/\varepsilon^4)$  *sfo* and  $\mathcal{O}(1/\varepsilon^2)$  *lmo* calls. Equipped with a variance reduction technique, SVRF achieves this accuracy with  $\mathcal{O}(1/\varepsilon^{10/3})$  *sfo* and  $\mathcal{O}(1/\varepsilon^2)$  *lmo* complexity. In the finite-sum setting, the *ifo* complexity gets  $\mathcal{O}(n + n^{2/3}/\varepsilon^2)$ . We omit the details of the SAGAFW approach, also present in [RSPS16], due to a technical flaw in its analysis (while telescoping Eq.(14) in page 1249 in [RSPS16]).

Qu et al. [QLX18] show the convergence rate for special instances of CGS and SCGS in the non-convex setting. However, they consider a different convergence criterion based on a proximal gradient mapping rather than the conventional FW-gap. Consequently, their results are incomparable with the rest of the literature. For the fact that we are running a *projection-free* method, the FW-gap is a more natural choice than the projection/proximal gradient norm.

We provide a parameter setting and a compact proof for CGS and SCGS in the supplemental material. Note however this setting simply gets the same guarantees as FW and SFW respectively. Whether CGS can provide improved oracle complexities compared to FW in the non-convex setting is an open problem.

For the non-convex setting, SPIDER-FW and SPIDER-CGS have the same oracle complexities, superior to SVRF (which is the state-of-the-art to our knowledge) for finite-sum and expectation minimization problems.

### 7.5.3 Results from Concurrent Research

By the time we prepared our manuscript [YSC19], the idea of combining SPIDER with CGM analysis was not explored yet. However, there are a few concurrent works that appeared online while our paper was under review. In this section, we discuss these works.

The recent work by Shen et al. [SFZ<sup>+</sup>19] is very closely related to our approach. They propose a class of methods based on CGM and various variance reduction techniques for the non-convex finite-sum setting, including the SPIDER-FW. Besides, they also propose extensions that use second-order approximations. Finally, they provide simulation studies to compare empirical performance of different variants. We refer to this paper for a numerical comparison.

Hassani et al. [HKMS19] introduce a novel variance reduced CGM method, but their work focuses primarily on the submodular maximization. Accordingly, they consider a more general expectation minimization template (the so-called non-oblivious setting) where the probability distribution depends on the decision variable  $\mathbf{x}$  and may change during the optimization procedure. Therefore, the proposed method requires some further assumptions and modifications involving computations with the Hessian approximation. Finally, Zhang et al. [ZCM<sup>+</sup>19] consider a stochastic CGM approach with SPIDER in the distributed and quantized settings.

### Concluding Remarks

We have proposed two novel FW-type methods based on the idea of blending the recent variance reduction technique SPIDER into FW and CGS frameworks. We have shown that the resulting methods enjoy superior oracle complexities in various convex and non-convex optimization templates. Extension of our framework for the strongly convex case is left open. Developing a well-tuned implementation, including one that incorporates parallel optimization, is an important piece of future work.

## Appendix: Proofs

### Preliminaries

This section presents some known results from the existing literature, key to our analysis, for the sake of completeness.

The following Lemma from [FLLZ18] provides an error bound for the estimator  $\mathbf{v}^k$  obtained by the SPIDER approach.

**Lemma 7.9** (Lemma 1 from [FLLZ18], more specifically Eqn.(A.3) in its supplements). Suppose that  $\mathcal{S}_{t,k}$  is a subset that samples  $S_{t,k}$  *iid* realizations from the distribution  $\mathcal{P}$ . Let the stochastic estimator  $\nabla f_{\mathcal{S}_{t,k}}$  satisfy the averaged  $L$ -Lipschitz gradients condition from Section 7.2. Set the estimator  $\mathbf{v}^k$  as

$$\mathbf{v}^k = \nabla f_{\mathcal{S}_{t,k}}(\mathbf{x}^k) - \nabla f_{\mathcal{S}_{t,k}}(\mathbf{x}^{k-1}) + \mathbf{v}^{k-1}. \quad (7.29)$$

Then, the following bound holds:

$$\mathbb{E}\|\nabla F(\mathbf{x}^k) - \mathbf{v}^k\|^2 \leq \frac{L^2}{S_{t,k}} \|\mathbf{x}^k - \mathbf{x}^{k-1}\|^2 + \|\nabla F(\mathbf{x}^{k-1}) - \mathbf{v}^{k-1}\|^2. \quad (7.30)$$

The next Lemma draws a well-known bound on the variance in terms of the mini-batch size.

**Lemma 7.10** (Eqn.(3.5) from [LZ16], or Lemma 2 from [RSPS16]). Suppose that  $\mathcal{S}_{t,k}$  is a subset that samples  $S_{t,k}$  *iid* realizations from the distribution  $\mathcal{P}$ . Let the stochastic estimator  $\nabla f_{\mathcal{S}_{t,k}}$  satisfy the bounded variance condition from Section 7.2. Then, the following bound holds:

$$\mathbb{E}\|\nabla f_{\mathcal{S}_{t,k}}(\mathbf{x}) - \nabla F(\mathbf{x})\|^2 \leq \frac{\sigma^2}{S_{t,k}} \quad \forall \mathbf{x} \in \mathcal{X}. \quad (7.31)$$

Finally, we recall the convergence guarantees for the CndG procedure of CGS-type methods.

**Lemma 7.11** (Similar to Theorem 2.2 part (c) from [LZ16], or more generally Theorem 2 from [Jag13]). Remark that CndG procedure simply applies CGM (with exact line-search) for the following projection subproblem:

$$\underset{\mathbf{x} \in \mathcal{X}}{\text{minimize}} \quad \frac{\beta}{2} \|\mathbf{x} - \mathbf{u} + \frac{1}{\beta} \mathbf{v}\|^2. \quad (7.32)$$

The objective function in this subproblem is  $\beta$ -smooth, hence CGM requires at most  $\mathcal{O}(\frac{4\beta D^2}{\alpha})$  iterations to satisfy the convergence criterion.

We refer to the associated references for the proofs.

### Non-convex Conditional Gradient Sliding

In this section, we prove convergence of a CGS instance, and derive its oracle complexities in the non-convex settings. We also extend our results for the stochastic variant SCGS.

#### Proof of convergence for non-convex CGS

---

##### Algorithm 7.4 Conditional Gradient Sliding

---

```

1: Input:  $\mathbf{x}^1 \in \mathcal{X}$ 
2: Set:  $\alpha = \gamma LD^2$ ,  $\beta = \gamma L/2$ ,  $\gamma = 1/\sqrt{K}$ 
3: for  $k = 1, 2, \dots, K$  do
4:   Update  $\mathbf{z}^k = \mathbf{y}^k + \gamma(\mathbf{x}^k - \mathbf{y}^k)$ 
5:    $\mathbf{x}^{k+1} = \text{CndG}(\mathbf{x}^k, \nabla F(\mathbf{z}^k), \alpha, \beta)$ 
6:   Update  $\mathbf{y}^{k+1} = \mathbf{y}^k + \gamma(\mathbf{x}^{k+1} - \mathbf{y}^k)$ 
7: end for

```

---

**Theorem 7.12.** Consider the CGS algorithm with the parameters as described in Algorithm 7.4 (in the batch setting). Denote by  $\mathbf{y}^{out}$  a random iterate  $\mathbf{y}^k$  drawn uniformly random over all iterates of the CGS. Then, the following bound holds:

$$\mathbb{E}[\mathcal{G}(\mathbf{y}^{out})] = \frac{\mathcal{E} + 3LD^2}{\sqrt{K}} \quad (7.33)$$

**Corollary 7.9.** The *ifo* and *lmo* complexities of CGS for achieving an  $\varepsilon$ -solution in the non-convex minimization setting are

$$\begin{aligned} \#(ifo) &= \mathcal{O}\left((\mathcal{E}^2 + L^2 D^4) \frac{n}{\varepsilon^2}\right) \\ \#(lmo) &= \mathcal{O}\left((\mathcal{E}^2 + L^2 D^4) \frac{1}{\varepsilon^2}\right) \end{aligned} \quad (7.34)$$

*Proof.* We start by the Taylor expansion and smoothness:

$$\begin{aligned} F(\mathbf{y}^{k+1}) &\leq F(\mathbf{y}^k) + \langle \nabla F(\mathbf{y}^k), \mathbf{y}^{k+1} - \mathbf{y}^k \rangle + \frac{L}{2} \|\mathbf{y}^{k+1} - \mathbf{y}^k\|^2 \\ &= F(\mathbf{y}^k) + \gamma \langle \nabla F(\mathbf{y}^k), \mathbf{x}^{k+1} - \mathbf{y}^k \rangle + \gamma^2 \frac{L}{2} \|\mathbf{x}^{k+1} - \mathbf{y}^k\|^2 \\ &\leq F(\mathbf{y}^k) + \gamma \langle \nabla F(\mathbf{y}^k), \mathbf{x}^{k+1} - \mathbf{y}^k \rangle + \gamma^2 \frac{L}{2} D^2 \\ &= F(\mathbf{y}^k) + \gamma \langle \nabla F(\mathbf{y}^k), \mathbf{w}_\star^k - \mathbf{y}^k \rangle + \gamma \langle \nabla F(\mathbf{y}^k), \mathbf{x}^{k+1} - \mathbf{w}_\star^k \rangle + \gamma^2 \frac{L}{2} D^2 \\ &= F(\mathbf{y}^k) - \gamma \mathcal{G}(\mathbf{y}^k) + \gamma \langle \nabla F(\mathbf{y}^k), \mathbf{x}^{k+1} - \mathbf{w}_\star^k \rangle + \gamma^2 \frac{L}{2} D^2 \end{aligned} \quad (7.35)$$

where we define  $\mathbf{w}_\star^k = \arg \max_{\mathbf{x} \in \mathcal{X}} \langle \mathbf{x}, -\nabla F(\mathbf{y}^k) \rangle$ .

We can equivalently write this inequality as

$$\begin{aligned} F(\mathbf{y}^{k+1}) &\leq F(\mathbf{y}^k) - \gamma \mathcal{G}(\mathbf{y}^k) + \gamma \langle \nabla F(\mathbf{y}^k) - \nabla F(\mathbf{z}^k), \mathbf{x}^{k+1} - \mathbf{w}_\star^k \rangle \\ &\quad + \gamma \langle \nabla F(\mathbf{z}^k), \mathbf{x}^{k+1} - \mathbf{w}_\star^k \rangle + \gamma^2 \frac{L}{2} D^2 \end{aligned} \quad (7.36)$$

Focus on the last inner-product term

$$\begin{aligned} \gamma \langle \nabla F(\mathbf{z}^k), \mathbf{x}^{k+1} - \mathbf{w}_\star^k \rangle &= \gamma \langle \nabla F(\mathbf{z}^k) + \beta(\mathbf{x}^{k+1} - \mathbf{x}^k), \mathbf{x}^{k+1} - \mathbf{w}_\star^k \rangle - \gamma \beta \langle \mathbf{x}^{k+1} - \mathbf{x}^k, \mathbf{x}^{k+1} - \mathbf{w}_\star^k \rangle \\ &\leq \gamma \alpha - \gamma \beta \langle \mathbf{x}^{k+1} - \mathbf{x}^k, \mathbf{x}^{k+1} - \mathbf{w}_\star^k \rangle \\ &\leq \gamma \alpha + \gamma \beta D^2 \end{aligned} \quad (7.37)$$

where the first inequality follows from the role of  $\alpha$  in CndG, and the second one from the Cauchy-Schwarz inequality.

Combining these two inequalities, we obtain

$$\begin{aligned} F(\mathbf{y}^{k+1}) &\leq F(\mathbf{y}^k) - \gamma \mathcal{G}(\mathbf{y}^k) + \gamma \langle \nabla F(\mathbf{y}^k) - \nabla F(\mathbf{z}^k), \mathbf{x}^{k+1} - \mathbf{w}_\star^k \rangle + \gamma \alpha + \gamma \beta D^2 + \gamma^2 \frac{L}{2} D^2 \\ &\leq F(\mathbf{y}^k) - \gamma \mathcal{G}(\mathbf{y}^k) + \gamma D \|\nabla F(\mathbf{y}^k) - \nabla F(\mathbf{z}^k)\| + \gamma \alpha + \gamma \beta D^2 + \gamma^2 \frac{L}{2} D^2 \\ &\leq F(\mathbf{y}^k) - \gamma \mathcal{G}(\mathbf{y}^k) + \gamma^2 L D \|\mathbf{x}^k - \mathbf{y}^k\| + \gamma \alpha + \gamma \beta D^2 + \gamma^2 \frac{L}{2} D^2 \\ &\leq F(\mathbf{y}^k) - \frac{1}{\sqrt{K}} \mathcal{G}(\mathbf{y}^k) + \frac{3LD^2}{K} \end{aligned} \quad (7.38)$$

Taking expectation of both sides, rearranging, and summing over all iterations, we obtain

$$\frac{1}{\sqrt{K}} \sum_{k=1}^K \mathbb{E}[\mathcal{G}(\mathbf{y}^k)] \leq F(\bar{\mathbf{x}}^1) - \mathbb{E}[F(\mathbf{y}^K)] + \sum_{k=1}^K \frac{3LD^2}{K} \leq F(\bar{\mathbf{x}}^1) - F(\mathbf{x}_\star) + 3LD^2 \quad (7.39)$$

Hence, by definition of  $\mathbf{y}^{\text{out}}$ , we get

$$\mathbb{E}[\mathcal{G}(\mathbf{y}^{\text{out}})] \leq \frac{F(\bar{\mathbf{x}}^1) - F(\mathbf{x}_\star)}{\sqrt{K}} + \frac{3LD^2}{\sqrt{K}} = \frac{\mathcal{E} + 3LD^2}{\sqrt{K}} \quad (7.40)$$

This completes the convergence rate proof.

To get  $\varepsilon$ -solution, we set the number of iterations  $K_\varepsilon$  such that

$$\mathbb{E}[\mathcal{G}(\mathbf{y}^{\text{out}})] \leq \frac{\mathcal{E} + 3LD^2}{\sqrt{K_\varepsilon}} \leq \varepsilon. \quad (7.41)$$

Hence, we can calculate the  $lmo$  complexity by using Lemma 7.11 as

$$\#(lmo) = \mathcal{O}\left(K_\varepsilon \frac{4\beta D^2}{\alpha}\right) = \mathcal{O}\left((\mathcal{E}^2 + L^2 D^4) \frac{1}{\varepsilon^2}\right) \quad (7.42)$$

which completes the proof.  $\square$

### Proof of convergence for the non-convex SCGS

---

#### Algorithm 7.5 Stochastic Conditional Gradient Sliding

---

```

1: Input:  $\mathbf{x}^1 \in \mathcal{X}$ 
2: Set:  $\alpha = \gamma L D^2$ ,  $\beta = \gamma L/2$ ,  $\gamma = 1/\sqrt{K}$ 
3: for  $k = 1, 2, \dots, K$  do
4:   Update  $\mathbf{z}^k = \mathbf{y}^k + \gamma(\mathbf{x}^k - \mathbf{y}^k)$ 
5:   Draw  $K$  samples  $\mathcal{S}_k$ 
6:    $\mathbf{x}^{k+1} = \text{CndG}(\mathbf{x}^k, \nabla f_{\mathcal{S}_k}(\mathbf{z}^k), \alpha, \beta)$ 
7:   Update  $\mathbf{y}^{k+1} = \mathbf{y}^k + \gamma(\mathbf{x}^{k+1} - \mathbf{y}^k)$ 
8: end for

```

---

**Theorem 7.13.** Consider the SCGS algorithm with the parameters as described in Algorithm 7.5. Assume that the conditions for the expectation minimization from Section 7.2 hold. Denote by  $\mathbf{y}^{out}$  a random iterate  $\mathbf{y}^k$  drawn uniformly random over all iterates of the SCGS. Then, the following bound holds:

$$\mathbb{E}[\mathcal{G}(\mathbf{y}^{out})] = \frac{\mathcal{E} + \sigma D + 3LD^2}{\sqrt{K}} \quad (7.43)$$

**Corollary 7.10.** The  $sfo$  and  $lmo$  complexities of SCGS for achieving  $\varepsilon$ -solution in the non-convex minimization setting are

$$\#(sfo) = \mathcal{O}\left((\mathcal{E} + \sigma D + LD^2)^4 \frac{1}{\varepsilon^4}\right) \quad \text{and} \quad \#(lmo) = \mathcal{O}\left((\mathcal{E} + \sigma D + LD^2)^2 \frac{1}{\varepsilon^2}\right) \quad (7.44)$$

*Proof.* From (7.35), and similar to (7.36) and (7.37), we can show

$$F(\mathbf{y}^{k+1}) \leq F(\mathbf{y}^k) - \gamma \mathcal{G}(\mathbf{y}^k) + \gamma \langle \nabla F(\mathbf{y}^k) - \nabla f_{\mathcal{S}_k}(\mathbf{z}^k), \mathbf{x}^{k+1} - \mathbf{w}_\star^k \rangle + \gamma \alpha + \gamma \beta D^2 + \gamma^2 \frac{L}{2} D^2 \quad (7.45)$$

where  $\mathbf{w}_\star^k = \arg \max_{\mathbf{x} \in \mathcal{X}} \langle \mathbf{x}, -\nabla F(\mathbf{y}^k) \rangle$ .



Focusing on the inner-product term, we get the following bound:

$$\begin{aligned}
 \gamma \langle \nabla F(\mathbf{y}^k) - \nabla f_{S_k}(\mathbf{z}^k), \mathbf{x}^{k+1} - \mathbf{w}_\star^k \rangle &\leq \gamma \|\nabla F(\mathbf{y}^k) - \nabla f_{S_k}(\mathbf{z}^k)\| \|\mathbf{x}^{k+1} - \mathbf{w}_\star^k\| \\
 &\leq \gamma D \|\nabla F(\mathbf{y}^k) - \nabla f_{S_k}(\mathbf{z}^k)\| \\
 &\leq \gamma D \|\nabla F(\mathbf{y}^k) - \nabla F(\mathbf{z}^k)\| + \gamma D \|\nabla F(\mathbf{z}^k) - \nabla f_{S_k}(\mathbf{z}^k)\| \\
 &\leq \gamma LD \|\mathbf{y}^k - \mathbf{z}^k\| + \gamma D \|\nabla F(\mathbf{z}^k) - \nabla f_{S_k}(\mathbf{z}^k)\| \\
 &= \gamma^2 LD \|\mathbf{x}^k - \mathbf{y}^k\| + \gamma D \|\nabla F(\mathbf{z}^k) - \nabla f_{S_k}(\mathbf{z}^k)\| \\
 &\leq \gamma^2 LD^2 + \gamma D \|\nabla F(\mathbf{z}^k) - \nabla f_{S_k}(\mathbf{z}^k)\|
 \end{aligned} \tag{7.46}$$

Substituting back and taking expectations, we obtain

$$\begin{aligned}
 \mathbb{E}[F(\mathbf{y}^{k+1})] &\leq \mathbb{E}[F(\mathbf{y}^k)] - \frac{1}{\sqrt{K}} \mathbb{E}[\mathcal{G}(\mathbf{y}^k)] + \frac{D}{\sqrt{K}} \mathbb{E}[\|\nabla F(\mathbf{z}^k) - \nabla f_{S_k}(\mathbf{z}^k)\|] + 3\gamma^2 LD^2 \\
 &= \mathbb{E}[F(\mathbf{y}^k)] - \frac{1}{\sqrt{K}} \mathbb{E}[\mathcal{G}(\mathbf{y}^k)] + \frac{D}{\sqrt{K}} \mathbb{E}[\|\nabla F(\mathbf{z}^k) - \nabla f_{S_k}(\mathbf{z}^k)\|] + \frac{3LD^2}{K}
 \end{aligned} \tag{7.47}$$

Now, we use Lemma 7.10 with the Jensen's inequality to obtain

$$\mathbb{E}[F(\mathbf{y}^{k+1})] \leq \mathbb{E}[F(\mathbf{y}^k)] - \frac{1}{\sqrt{K}} \mathbb{E}[\mathcal{G}(\mathbf{y}^k)] + \frac{\sigma D + 3LD^2}{K} \tag{7.48}$$

From here, we follow the same steps as in the proof of CGS and get (7.43).

Then, to achieve an  $\varepsilon$ -solution, we can calculate *sfo* complexity as

$$\#(sfo) = \sum_{k=1}^{K_\varepsilon} K_\varepsilon = K_\varepsilon^2 = \mathcal{O}\left((\varepsilon + \sigma D + LD^2)^4 \frac{1}{\varepsilon^4}\right) \tag{7.49}$$

Finally, *lmo* complexity can be found using Lemma 7.11

$$\#(lmo) = \mathcal{O}\left(K_\varepsilon \frac{4\beta D^2}{\alpha}\right) = \mathcal{O}\left((\varepsilon + \sigma D + LD^2)^2 \frac{1}{\varepsilon^2}\right) \tag{7.50}$$

This completes the proof.  $\square$

### Proofs for SPIDER-FW

**Lemma 7.14.** Suppose that the assumptions listed in Section 7.2 hold. Then, for  $k = 1, \dots, K_t$ , we have the following bounds:

$$\text{Convex finite-sum} \quad \mathbb{E}\|\nabla F(\mathbf{x}^k) - \mathbf{v}^k\| \leq 2LD/K_t \tag{7.51}$$

$$\text{Convex expectation} \quad \mathbb{E}\|\nabla F(\mathbf{x}^k) - \mathbf{v}^k\| \leq 3LD/K_t \tag{7.52}$$

$$\text{Non-convex finite-sum} \quad \mathbb{E}\|\nabla F(\mathbf{x}^k) - \mathbf{v}^k\| \leq LD/\sqrt{TK} \tag{7.53}$$

$$\text{Non-convex expectation} \quad \mathbb{E}\|\nabla F(\mathbf{x}^k) - \mathbf{v}^k\| \leq LD/\sqrt{TK} + \varepsilon/2 \tag{7.54}$$

## Chapter 7. CGM with Stochastic Path-Integrated Differential Estimator

---

*Proof.* From Lemma 7.9, we have the following inequality for  $k = 2, 3, \dots, K_t$ :

$$\mathbb{E}\|\nabla F(\mathbf{x}^k) - \mathbf{v}^k\|^2 \leq \frac{L^2}{S_{t,k}} \|\mathbf{x}^k - \mathbf{x}^{k-1}\|^2 + \|\nabla F(\mathbf{x}^{k-1}) - \mathbf{v}^{k-1}\|^2. \quad (7.55)$$

By definition,  $\|\mathbf{x}^{k-1} - \mathbf{x}^k\|^2 = \|\eta_{t,k-1}(\mathbf{w}^{k-1} - \mathbf{x}^{k-1})\|^2 \leq \eta_{t,k-1}^2 D^2$ . Hence, we get

$$\mathbb{E}\|\nabla F(\mathbf{x}^k) - \mathbf{v}^k\|^2 \leq \frac{\eta_{t,k-1}^2 L^2 D^2}{S_{t,k}} + \|\nabla F(\mathbf{x}^{k-1}) - \mathbf{v}^{k-1}\|^2. \quad (7.56)$$

Convex finite-sum:

We take the telescopic sum of (7.56) from  $i = 2$  to  $k$

$$\mathbb{E}\|\nabla F(\mathbf{x}^k) - \mathbf{v}^k\|^2 \leq \sum_{i=2}^k \frac{\eta_{t,i-1}^2 L^2 D^2}{S_{t,i}} + \underbrace{\|\nabla F(\mathbf{x}^1) - \mathbf{v}^1\|^2}_{0, \text{ since we take full batch}} \leq \frac{4L^2 D^2}{K_t} \sum_{i=2}^k \frac{1}{(s_{t,i-1} + 1)^2} \quad (7.57)$$

By definition of  $s_{t,k}$ , for any  $i \geq 2$  we have

$$s_{t,i-1} + 1 = K_t + i - 1 \geq K_t. \quad (7.58)$$

Hence, we get

$$\mathbb{E}\|\nabla F(\mathbf{x}^k) - \mathbf{v}^k\|^2 \leq \frac{4L^2 D^2}{K_t^3} \sum_{i=2}^k 1 \leq \frac{4L^2 D^2}{K_t^3} k \leq \frac{4L^2 D^2}{K_t^2}. \quad (7.59)$$

By using Jensen's inequality, we get (7.51).

Convex expectation minimization:

We take the telescopic sum of (7.56) from  $i = 2$  to  $k$

$$\mathbb{E}\|\nabla F(\mathbf{x}^k) - \mathbf{v}^k\|^2 \leq \frac{4L^2 D^2}{K_t^2} + \|\nabla F(\mathbf{x}^1) - \mathbf{v}^1\|^2 \leq \frac{9L^2 D^2}{K_t^2}, \quad (7.60)$$

where the bound on  $\|\nabla F(\mathbf{x}^1) - \mathbf{v}^1\|^2$  follows from Lemma 7.10 as

$$\|\nabla F(\mathbf{x}^1) - \mathbf{v}^1\|^2 \leq \frac{\sigma^2}{Q_t} = \frac{5L^2 D^2}{K_t^2}. \quad (7.61)$$

We get (7.52) by using Jensen's inequality.

*Non-convex finite-sum:*

We take the telescopic sum of (7.56) from  $i = 2$  to  $k$

$$\mathbb{E}\|\nabla F(\mathbf{x}^k) - \mathbf{v}^k\|^2 \leq \sum_{i=2}^k \frac{\eta_{t,i}^2 L^2 D^2}{S_{t,i}} + \underbrace{\|\nabla F(\mathbf{x}^1) - \mathbf{v}^1\|^2}_{0, \text{ since we take full batch}} \leq \frac{L^2 D^2}{TK^2} \sum_{i=2}^k 1 \leq \frac{L^2 D^2}{TK} \quad (7.62)$$

We get (7.53) by using Jensen's inequality.

*Non-convex expectation minimization:*

We take the telescopic sum of (7.56) from  $i = 2$  to  $k$

$$\mathbb{E}\|\nabla F(\mathbf{x}^k) - \mathbf{v}^k\|^2 \leq \frac{L^2 D^2}{TK} + \|\nabla F(\mathbf{x}^1) - \mathbf{v}^1\|^2 \leq \frac{L^2 D^2}{TK} + \frac{\varepsilon^2}{4} \quad (7.63)$$

where the bound on  $\|\nabla F(\mathbf{x}^1) - \mathbf{v}^1\|^2$  follows from Lemma 7.10 as

$$\|\nabla F(\mathbf{x}^1) - \mathbf{v}^1\|^2 \leq \frac{\sigma^2}{Q_t} \leq \frac{\varepsilon^2}{4} \quad (7.64)$$

We get (7.54) by using Jensen's inequality.  $\square$

### Proof for Theorem 7.1 and Corollary 7.1

We start by the Taylor expansion and smoothness:

$$\begin{aligned} F(\mathbf{x}^{k+1}) &\leq F(\mathbf{x}^k) + \langle \nabla F(\mathbf{x}^k), \mathbf{x}^{k+1} - \mathbf{x}^k \rangle + \frac{L}{2} \|\mathbf{x}^{k+1} - \mathbf{x}^k\|^2 \\ &\leq F(\mathbf{x}^k) + \eta_{t,k} \langle \nabla F(\mathbf{x}^k), \mathbf{w}^k - \mathbf{x}^k \rangle + \eta_{t,k}^2 \frac{L}{2} D^2 \\ &= F(\mathbf{x}^k) + \eta_{t,k} \langle \mathbf{v}^k, \mathbf{w}^k - \mathbf{x}^k \rangle + \eta_{t,k} \langle \nabla F(\mathbf{x}^k) - \mathbf{v}^k, \mathbf{w}^k - \mathbf{x}^k \rangle + \eta_{t,k}^2 \frac{L}{2} D^2 \end{aligned} \quad (7.65)$$

By definition of  $\mathbf{w}^k$ , we have

$$\langle \mathbf{v}^k, \mathbf{w}^k - \mathbf{x}^k \rangle = \min_{\mathbf{x} \in \mathcal{X}} \langle \mathbf{v}^k, \mathbf{x} - \mathbf{x}^k \rangle \leq \langle \mathbf{v}^k, \mathbf{x}_\star - \mathbf{x}^k \rangle \quad (7.66)$$

Substituting this inequality back, and rearranging, we get

$$\begin{aligned} F(\mathbf{x}^{k+1}) &\leq F(\mathbf{x}^k) + \eta_{t,k} \langle \mathbf{v}^k, \mathbf{x}_\star - \mathbf{x}^k \rangle + \eta_{t,k} \langle \nabla F(\mathbf{x}^k) - \mathbf{v}^k, \mathbf{w}^k - \mathbf{x}^k \rangle + \eta_{t,k}^2 \frac{L}{2} D^2 \\ &= F(\mathbf{x}^k) + \eta_{t,k} \langle \nabla F(\mathbf{x}^k), \mathbf{x}_\star - \mathbf{x}^k \rangle + \eta_{t,k} \langle \nabla F(\mathbf{x}^k) - \mathbf{v}^k, \mathbf{w}^k - \mathbf{x}_\star \rangle + \eta_{t,k}^2 \frac{L}{2} D^2 \end{aligned} \quad (7.67)$$

From convexity of  $F$ , we know

$$\langle \nabla F(\mathbf{x}^k), \mathbf{x}_\star - \mathbf{x}^k \rangle \leq F_\star - F(\mathbf{x}^k) \quad (7.68)$$

and by using Cauchy-Schwarz, we have

$$\langle \nabla F(\mathbf{x}^k) - \mathbf{v}^k, \mathbf{w}^k - \mathbf{x}_\star \rangle \leq \|\nabla F(\mathbf{x}^k) - \mathbf{v}^k\| \|\mathbf{w}^k - \mathbf{x}_\star\| \leq \|\nabla F(\mathbf{x}^k) - \mathbf{v}^k\| D \quad (7.69)$$

Putting (7.68) and (7.69) back into (7.67), and subtracting  $F_\star$  from both sides, we obtain:

$$F(\mathbf{x}^{k+1}) - F_\star \leq (1 - \eta_{t,k})(F(\mathbf{x}^k) - F_\star) + \eta_{t,k} D \|\nabla F(\mathbf{x}^k) - \mathbf{v}^k\| + \eta_{t,k}^2 \frac{L}{2} D^2 \quad (7.70)$$

Then, we take the expectation of both sides and use (7.51) to get

$$\begin{aligned} \mathbb{E}[F(\mathbf{x}^{k+1})] - F_\star &\leq (1 - \eta_{t,k})(\mathbb{E}[F(\mathbf{x}^k)] - F_\star) + \eta_{t,k} D \mathbb{E}\|\nabla F(\mathbf{x}^k) - \mathbf{v}^k\| + \eta_{t,k}^2 \frac{L}{2} D^2 \\ &\leq (1 - \eta_{t,k})(\mathbb{E}[F(\mathbf{x}^k)] - F_\star) + \eta_{t,k} \frac{2LD^2}{K_t} + \eta_{t,k}^2 \frac{L}{2} D^2 \end{aligned} \quad (7.71)$$

Telescopic sum of this inequality over  $(t, k)$  pairs gives

$$\begin{aligned} \mathbb{E}[F(\mathbf{x}^{k+1})] - F_\star &\leq \sum_{(\tau, i)} \left( \eta_{\tau, i} \frac{2LD^2}{K_\tau} + \eta_{\tau, i}^2 \frac{L}{2} D^2 \right) \prod_{(\tau', j)=(\tau, i)}^{(t, k)} (1 - \eta_{\tau', j}) \\ &\quad + (\mathbb{E}[F(\mathbf{x}^{1,1})] - F_\star) \prod_{(\tau, i)} (1 - \eta_{\tau, i}) \end{aligned} \quad (7.72)$$

The last term vanishes due to 0 factor ( $\eta_{1,1} = 1$ ). Remark that

$$\prod_{(\tau', j)=(\tau, i)}^{(t, k)} (1 - \eta_{\tau', j}) = \prod_{r=i}^{K_\tau} \frac{s_{\tau, r} - 1}{s_{\tau, r} + 1} \prod_{\tau'=\tau+1}^{t-1} \prod_{j=1}^{K_{\tau'}} \frac{s_{\tau', j} - 1}{s_{\tau', j} + 1} \prod_{\ell=1}^k \frac{s_{t, \ell} - 1}{s_{t, \ell} + 1} = \frac{(s_{\tau, i} - 1)s_{\tau, i}}{s_{t, k}(s_{t, k} + 1)} \quad (7.73)$$

Combining these, we get

$$\mathbb{E}[F(\mathbf{x}^{k+1})] - F_\star \leq \sum_{(\tau, i)} \left( \eta_{\tau, i} \frac{2LD^2}{K_\tau} + \eta_{\tau, i}^2 \frac{L}{2} D^2 \right) \frac{(s_{\tau, i} - 1)s_{\tau, i}}{s_{t, k}(s_{t, k} + 1)}. \quad (7.74)$$

We focus on the individual terms:

$$\sum_{(\tau, i)} \eta_{\tau, i} \frac{2LD^2}{K_\tau} \frac{(s_{\tau, i} - 1)s_{\tau, i}}{s_{t, k}(s_{t, k} + 1)} \leq \frac{8LD^2}{s_{t, k}(s_{t, k} + 1)} \sum_{(\tau, i)} 1 \leq \frac{8LD^2}{s_{t, k} + 1} \quad (7.75)$$

$$\sum_{(\tau, i)} \eta_{\tau, i}^2 \frac{L}{2} D^2 \frac{(s_{\tau, i} - 1)s_{\tau, i}}{s_{t, k}(s_{t, k} + 1)} \leq \frac{2LD^2}{s_{t, k}(s_{t, k} + 1)} \sum_{(\tau, i)} 1 \leq \frac{2LD^2}{K_t + k} \quad (7.76)$$

We proved the convergence rate:

$$\mathbb{E}[F(\mathbf{x}^{k+1})] - F_\star \leq \frac{10LD^2}{s_{t, k} + 1} \quad (7.77)$$

To get  $\varepsilon$ -solution, we set the number of outer iterations  $T_\varepsilon$  such that

$$\mathbb{E}[F(\bar{\mathbf{x}}^{T_\varepsilon})] - F_\star \leq \frac{10LD^2}{K_{T_\varepsilon}} \leq \varepsilon. \quad (7.78)$$

Then, it is sufficient to choose

$$T_\varepsilon = \log_2 \left( \frac{10LD^2}{\varepsilon} \right) + 1. \quad (7.79)$$

Then, to achieve an  $\varepsilon$ -solution, we can calculate the *ifo* complexity as

$$\begin{aligned} \#(ifo) &= \sum_{t=1}^{T_\varepsilon} \left( Q_t + \sum_{k=2}^{K_t} S_{t,k} \right) = \sum_{t=1}^{T_\varepsilon} \left( n + \sum_{k=2}^{K_t} 2^{t-1} \right) \\ &\leq \sum_{t=1}^{T_\varepsilon} (n + 2^{2(t-1)}) \\ &= \mathcal{O}(nT_\varepsilon + 2^{2T_\varepsilon}) = \mathcal{O} \left( n \ln \left( \frac{LD^2}{\varepsilon} \right) + \frac{L^2 D^4}{\varepsilon^2} \right) \end{aligned} \quad (7.80)$$

and the *lmo* complexity as

$$\#(lmo) = \sum_{t=1}^{T_\varepsilon} K_t \leq 2K_{T_\varepsilon} = 2^{T_\varepsilon} = \mathcal{O} \left( \frac{LD^2}{\varepsilon} \right). \quad (7.81)$$

### Proof for Theorem 7.2 and Corollary 7.2

Proof is similar to that for finite-sum setting, but we use (7.52) instead of (7.51) at (7.71), hence the constants change:

$$\mathbb{E}[F(\mathbf{x}^{k+1})] - F_\star \leq \frac{14LD^2}{K_t + k}. \quad (7.82)$$

To get  $\varepsilon$ -solution, we set the number of outer iterations  $T_\varepsilon$  as

$$T_\varepsilon = \log_2 \left( \frac{14LD^2}{\varepsilon} \right) + 1. \quad (7.83)$$

Then, to achieve  $(1 - \varepsilon)$  accuracy, we can calculate *sfo* complexity as

$$\#(sfo) = \sum_{t=1}^{T_\varepsilon} \left( Q_t + \sum_{k=2}^{K_t} S_{t,k} \right) \leq \sum_{t=1}^{T_\varepsilon} \left( \left\lceil \frac{\sigma^2 K_t^2}{5L^2 D^2} \right\rceil + K_t^2 \right) = \mathcal{O} \left( \frac{\sigma^2 D^2}{\varepsilon^2} + \frac{L^2 D^4}{\varepsilon^2} \right). \quad (7.84)$$

The *lmo* complexity is same as the finite-sum case.

### Proof for Theorem 7.3 and Corollary 7.3

We start by the Taylor expansion and smoothness:

$$\begin{aligned}
 F(\mathbf{x}^{k+1}) &\leq F(\mathbf{x}^k) + \langle \nabla F(\mathbf{x}^k), \mathbf{x}^{k+1} - \mathbf{x}^k \rangle + \frac{L}{2} \|\mathbf{x}^{k+1} - \mathbf{x}^k\|^2 \\
 &\leq F(\mathbf{x}^k) + \eta \langle \nabla F(\mathbf{x}^k), \mathbf{w}^k - \mathbf{x}^k \rangle + \eta^2 \frac{L}{2} D^2 \\
 &= F(\mathbf{x}^k) + \eta \langle \mathbf{v}^k, \mathbf{w}^k - \mathbf{x}^k \rangle + \eta \langle \nabla F(\mathbf{x}^k) - \mathbf{v}^k, \mathbf{w}^k - \mathbf{x}^k \rangle + \eta^2 \frac{L}{2} D^2 \\
 &\leq F(\mathbf{x}^k) + \eta \langle \mathbf{v}^k, \mathbf{w}_\star^k - \mathbf{x}^k \rangle + \eta \langle \nabla F(\mathbf{x}^k) - \mathbf{v}^k, \mathbf{w}^k - \mathbf{x}^k \rangle + \eta^2 \frac{L}{2} D^2 \\
 &= F(\mathbf{x}^k) + \eta \langle \nabla F(\mathbf{x}^k), \mathbf{w}_\star^k - \mathbf{x}^k \rangle + \eta \langle \nabla F(\mathbf{x}^k) - \mathbf{v}^k, \mathbf{w}^k - \mathbf{w}_\star^k \rangle + \eta^2 \frac{L}{2} D^2 \\
 &\leq F(\mathbf{x}^k) - \eta \mathcal{G}(\mathbf{x}^k) + \eta D \|\nabla F(\mathbf{x}^k) - \mathbf{v}^k\| + \eta^2 \frac{L}{2} D^2
 \end{aligned} \tag{7.85}$$

where  $\mathbf{w}_\star^k = \arg\max_{\mathbf{x} \in \mathcal{X}} \langle \mathbf{x}, -\nabla F(\mathbf{x}^k) \rangle$ . Taking expectation of both sides and using (7.53),

$$\mathbb{E}[F(\mathbf{x}^{k+1})] \leq \mathbb{E}[F(\mathbf{x}^k)] - \frac{1}{\sqrt{TK}} \mathbb{E}[\mathcal{G}(\mathbf{x}^k)] + \frac{3LD^2}{2TK} \tag{7.86}$$

Rearranging, and summing over all  $(t, k)$  pairs up to  $(T, K)$ , we obtain

$$\frac{1}{\sqrt{TK}} \sum_{(\tau, i)}^{(T, K)} \mathbb{E}[\mathcal{G}(\mathbf{x}^{\tau, i})] \leq F(\bar{\mathbf{x}}^1) - \mathbb{E}[F(\mathbf{x}^{T, K})] + \sum_{(\tau, i)}^{(T, K)} \frac{3LD^2}{2TK} \leq F(\bar{\mathbf{x}}^1) - F(\mathbf{x}_\star) + \frac{3LD^2}{2}. \tag{7.87}$$

Hence, by definition of  $\mathbf{x}^{\text{out}}$ , we have

$$\mathbb{E}[\mathcal{G}(\mathbf{x}^{\text{out}})] \leq \frac{F(\bar{\mathbf{x}}^1) - F(\mathbf{x}_\star)}{\sqrt{TK}} + \frac{3LD^2}{2\sqrt{TK}} = \frac{2\mathcal{E} + 3LD^2}{2\sqrt{TK}} \tag{7.88}$$

This completes the convergence rate proof.

To get  $\varepsilon$ -solution, we set the number of outer iterations  $T_\varepsilon$  such that

$$\mathbb{E}[\mathcal{G}(\mathbf{x}^{\text{out}})] \leq \frac{2\mathcal{E} + 3LD^2}{2\sqrt{T_\varepsilon K}} \leq \varepsilon. \tag{7.89}$$

Hence, it suffices to choose

$$T_\varepsilon = \left( \frac{2\mathcal{E} + 3LD^2}{2\varepsilon\sqrt{K}} \right)^2 = \frac{(2\mathcal{E} + 3LD^2)^2}{4\varepsilon^2 K}. \tag{7.90}$$

Then, to achieve an  $\varepsilon$ -solution, we can calculate *ifo* complexity as

$$\begin{aligned} \#(ifo) &= \sum_{t=1}^{T_\varepsilon} \left( Q_t + \sum_{k=2}^K S_{t,k} \right) = \sum_{t=1}^{T_\varepsilon} \left( n + \sum_{k=2}^K K \right) \leq (n + K^2) T_\varepsilon \leq \sqrt{n} \frac{3(2\mathcal{E} + 3LD^2)^2}{4\varepsilon^2} \\ &= \mathcal{O} \left( (\mathcal{E}^2 + L^2 D^4) \frac{\sqrt{n}}{\varepsilon^2} \right) \end{aligned} \quad (7.91)$$

and the *lmo* complexity as

$$\#(lmo) = \sum_{t=1}^{T_\varepsilon} K = K T_\varepsilon = \frac{(2\mathcal{E} + 3LD^2)^2}{4\varepsilon^2} = \mathcal{O} \left( (\mathcal{E}^2 + L^2 D^4) \frac{1}{\varepsilon^2} \right) \quad (7.92)$$

#### Proof for Theorem 7.4 and Corollary 7.4

Proof is similar to that for non-convex finite-sum setting, but we use (7.54) instead of (7.53) at (7.85), and we get

$$\mathbb{E}[F(\mathbf{x}^{k+1})] \leq \mathbb{E}[F(\mathbf{x}^k)] - \frac{1}{\sqrt{TK}} \mathbb{E}[\mathcal{G}(\mathbf{x}^k)] + \frac{3LD^2}{2TK} + \frac{\varepsilon}{2\sqrt{TK}} \quad (7.93)$$

Rearranging, and summing over all  $(t, k)$  pairs up to  $(T, K)$ , we get

$$\begin{aligned} \frac{1}{\sqrt{TK}} \sum_{(t,i)}^{(T,K)} \mathbb{E}[\mathcal{G}(\mathbf{x}^{t,i})] &\leq F(\bar{\mathbf{x}}^1) - \mathbb{E}[F(\mathbf{x}^{T,K})] + \sum_{(t,i)}^{(T,K)} \left( \frac{3LD^2}{2TK} + \frac{\varepsilon}{2\sqrt{TK}} \right) \\ &\leq F(\bar{\mathbf{x}}^1) - F(\mathbf{x}_\star) + \frac{3LD^2}{2} + \frac{\varepsilon\sqrt{TK}}{2}. \end{aligned} \quad (7.94)$$

Hence, by definition of  $\mathbf{x}^{\text{out}}$ , we have

$$\mathbb{E}[\mathcal{G}(\mathbf{x}^{\text{out}})] \leq \frac{F(\bar{\mathbf{x}}^1) - F(\mathbf{x}_\star)}{\sqrt{TK}} + \frac{3LD^2}{2\sqrt{TK}} + \frac{\varepsilon}{2} = \frac{2\mathcal{E} + 3LD^2}{2\sqrt{TK}} + \frac{\varepsilon}{2} \quad (7.95)$$

This completes the convergence proof.

To get  $\varepsilon$ -solution, we set the number of outer iterations  $T_\varepsilon$  such that

$$\mathbb{E}[\mathcal{G}(\mathbf{x}^{\text{out}})] \leq \frac{2\mathcal{E} + 3LD^2}{2\sqrt{T_\varepsilon K}} + \frac{\varepsilon}{2} \leq \varepsilon. \quad (7.96)$$

Then, to achieve an  $\varepsilon$ -solution, we can calculate *sfo* complexity as

$$\begin{aligned} \#(sfo) &= \sum_{t=1}^{T_\varepsilon} \left( Q_t + \sum_{k=2}^K S_{t,k} \right) = \sum_{t=1}^{T_\varepsilon} \left( n + \sum_{k=2}^K K \right) \leq ([4(\sigma/\varepsilon)^2] + ([\sigma/\varepsilon])^2) T_\varepsilon \\ &= \mathcal{O} \left( (\mathcal{E}^2 + L^2 D^4) \frac{\sigma}{\varepsilon^3} \right) \end{aligned} \quad (7.97)$$

Finally, *lmo* complexity is same as the non-convex finite-sum case.

### Proofs for SPIDER-CGS

**Lemma 7.15.** Suppose that the assumptions listed in Section 7.2 hold. Then, for  $k = 1, \dots, K_t$ , we have the following bounds:

$$\text{Convex finite-sum} \quad \mathbb{E} \|\nabla F(\mathbf{z}^k) - \mathbf{v}^k\| \leq 2\sqrt{2}LD/(s_{t,k} + 1)^2 \quad (7.98)$$

$$\text{Convex expectation} \quad \mathbb{E} \|\nabla F(\mathbf{z}^k) - \mathbf{v}^k\| \leq 3LD/(s_{t,k} + 1)^2 \quad (7.99)$$

$$\text{Non-convex finite-sum} \quad \mathbb{E} \|\nabla F(\mathbf{z}^k) - \mathbf{v}^k\| \leq 2LD/\sqrt{TK} \quad (7.100)$$

$$\text{Non-convex expectation} \quad \mathbb{E} \|\nabla F(\mathbf{z}^k) - \mathbf{v}^k\| \leq 2LD/\sqrt{TK} + \varepsilon/2 \quad (7.101)$$

*Proof.* From Lemma 7.9, we have the following inequality for all  $k = 2, \dots, K_t$ :

$$\mathbb{E} \|\nabla F(\mathbf{z}^k) - \mathbf{v}^k\|^2 \leq \frac{L^2}{S_{t,k}} \|\mathbf{z}^k - \mathbf{z}^{k-1}\|^2 + \|\nabla F(\mathbf{z}^{k-1}) - \mathbf{v}^{k-1}\|^2. \quad (7.102)$$

By definition,

$$\begin{aligned} \mathbf{z}^k &= \mathbf{y}^k + \gamma_{t,k}(\mathbf{x}^k - \mathbf{y}^k) \\ &= \mathbf{y}^{k-1} + \gamma_{t,k-1}(\mathbf{x}^k - \mathbf{y}^{k-1}) + \gamma_{t,k}(\mathbf{x}^k - \mathbf{y}^k) \\ \& \quad \mathbf{z}^{k-1} = \mathbf{y}^{k-1} + \gamma_{t,k-1}(\mathbf{x}^{k-1} - \mathbf{y}^{k-1}) \\ \implies \|\mathbf{z}^k - \mathbf{z}^{k-1}\|^2 &= \|\gamma_{t,k-1}(\mathbf{x}^k - \mathbf{x}^{k-1}) + \gamma_{t,k}(\mathbf{x}^k - \mathbf{y}^k)\|^2 \\ &= \gamma_{t,k-1}^2 \|\mathbf{x}^k - \mathbf{x}^{k-1}\|^2 + \gamma_{t,k}^2 \|\mathbf{x}^k - \mathbf{y}^k\|^2 + 2\gamma_{t,k-1}\gamma_{t,k} \langle \mathbf{x}^k - \mathbf{y}^k, \mathbf{x}^k - \mathbf{x}^{k-1} \rangle \\ &\leq \gamma_{t,k-1}^2 \|\mathbf{x}^k - \mathbf{x}^{k-1}\|^2 + \gamma_{t,k-1}^2 \|\mathbf{x}^k - \mathbf{y}^k\|^2 + 2\gamma_{t,k-1}^2 \|\mathbf{x}^k - \mathbf{y}^k\| \|\mathbf{x}^k - \mathbf{x}^{k-1}\| \\ &\leq 4\gamma_{t,k-1}^2 D^2 \end{aligned} \quad (7.103)$$

Substituting into (7.102), we get

$$\mathbb{E} \|\nabla F(\mathbf{z}^k) - \mathbf{v}^k\|^2 \leq \frac{4\gamma_{t,k-1}^2 L^2 D^2}{S_{t,k}} + \|\nabla F(\mathbf{z}^{k-1}) - \mathbf{v}^{k-1}\|^2. \quad (7.104)$$

Convex finite-sum:

We take the telescopic sum of (7.104) from  $i = 2$  to  $k$

$$\begin{aligned} \mathbb{E} \|\nabla F(\mathbf{z}^k) - \mathbf{v}^k\|^2 &\leq \sum_{i=2}^k \frac{4\gamma_{t,i-1}^2 L^2 D^2}{S_{t,i}} + \underbrace{\|\nabla F(\mathbf{z}^1) - \mathbf{v}^1\|^2}_{0, \text{ since we take full batch}} \\ &= \frac{4L^2 D^2}{9K_t(s_{t,K_t} + 1)^2} \sum_{i=2}^k \frac{9}{(s_{t,i-1} + 2)^2} = \frac{4L^2 D^2}{K_t(s_{t,K_t} + 1)^2} \sum_{i=2}^k \frac{1}{(s_{t,i} + 1)^2} \end{aligned} \quad (7.105)$$



Clearly,  $s_{t,i} + 1 \geq \frac{s_{t,K_t} + 1}{\sqrt{2}}$  for all  $2 \leq i \leq K_t$ . Hence, we get

$$\mathbb{E} \|\nabla F(\mathbf{z}^k) - \mathbf{v}^k\|^2 \leq \frac{8L^2 D^2}{K_t(s_{t,K_t} + 1)^2} \sum_{i=2}^k \frac{1}{(s_{t,K_t} + 1)^2} \leq \frac{8L^2 D^2}{(s_{t,k} + 1)^4} \quad (7.106)$$

We get (7.98) by using Jensen's inequality.

Convex expectation minimization:

We take the telescopic sum of (7.104) from  $i = 2$  to  $k$

$$\mathbb{E} \|\nabla F(\mathbf{z}^k) - \mathbf{v}^k\|^2 \leq \frac{8L^2 D^2}{(s_{t,k} + 1)^4} + \|\nabla F(\mathbf{z}^1) - \mathbf{v}^1\|^2 \leq \frac{8L^2 D^2}{(s_{t,k} + 1)^4} + \frac{\sigma^2}{Q_t} \leq \frac{9L^2 D^2}{(s_{t,k} + 1)^4} \quad (7.107)$$

where the bound on  $\|\nabla F(\mathbf{z}^1) - \mathbf{v}^1\|^2$  follows from Lemma 7.10.

We get (7.99) by using Jensen's inequality.

Non-convex finite-sum:

We take the telescopic sum of (7.104) from  $i = 2$  to  $k$

$$\mathbb{E} \|\nabla F(\mathbf{z}^k) - \mathbf{v}^k\|^2 \leq \sum_{i=2}^k \frac{4\gamma^2 L^2 D^2}{S} + \underbrace{\|\nabla F(\mathbf{z}^1) - \mathbf{v}^1\|^2}_{0, \text{ since we take full batch}} \leq \sum_{i=2}^k \frac{4L^2 D^2}{TK^2} \leq \frac{4L^2 D^2}{TK} \quad (7.108)$$

We get (7.100) by using Jensen's inequality.

Non-convex expectation minimization:

We take the telescopic sum of (7.104) from  $i = 2$  to  $k$

$$\mathbb{E} \|\nabla F(\mathbf{z}^k) - \mathbf{v}^k\|^2 \leq \sum_{i=2}^k \frac{4\gamma^2 L^2 D^2}{S} + \|\nabla F(\mathbf{z}^1) - \mathbf{v}^1\|^2 \leq \frac{4L^2 D^2}{TK} + \frac{\varepsilon^2}{4} \quad (7.109)$$

where the bound on  $\|\nabla F(\mathbf{z}^1) - \mathbf{v}^1\|^2$  follows from Lemma 7.10. We get (7.101) by using Jensen's inequality.  $\square$

**Proof for Theorem 7.5 and Corollary 7.5**

We start by Taylor expansion and smoothness:

$$\begin{aligned}
 F(\mathbf{y}^{k+1}) &\leq F(\mathbf{z}^k) + \langle \nabla F(\mathbf{z}^k), \mathbf{y}^{k+1} - \mathbf{z}^k \rangle + \frac{L}{2} \|\mathbf{y}^{k+1} - \mathbf{z}^k\|^2 \\
 &= F(\mathbf{z}^k) + \langle \nabla F(\mathbf{z}^k), \mathbf{y}^k - \mathbf{z}^k \rangle + \gamma_{t,k} \langle \nabla F(\mathbf{z}^k), \mathbf{x}^{k+1} - \mathbf{x}_\star \rangle \\
 &\quad + \gamma_{t,k} \langle \nabla F(\mathbf{z}^k), \mathbf{x}_\star - \mathbf{y}^k \rangle + \frac{L\gamma_{t,k}^2}{2} \|\mathbf{x}^{k+1} - \mathbf{x}^k\|^2 \quad (7.110) \\
 &= (1 - \gamma_{t,k})(F(\mathbf{z}^k) + \langle \nabla F(\mathbf{z}^k), \mathbf{y}^k - \mathbf{z}^k \rangle) + \gamma_{t,k}(F(\mathbf{z}^k) + \langle \nabla F(\mathbf{z}^k), \mathbf{x}_\star - \mathbf{z}^k \rangle) \\
 &\quad + \gamma_{t,k} \langle \nabla F(\mathbf{z}^k), \mathbf{x}^{k+1} - \mathbf{x}_\star \rangle + \frac{L\gamma_{t,k}^2}{2} \|\mathbf{x}^{k+1} - \mathbf{x}^k\|^2
 \end{aligned}$$

From the convexity of  $F$ , we have

$$F(\mathbf{y}^k) \geq F(\mathbf{z}^k) + \langle \nabla F(\mathbf{z}^k), \mathbf{y}^k - \mathbf{z}^k \rangle \quad \text{and} \quad F_\star \geq F(\mathbf{z}^k) + \langle \nabla F(\mathbf{z}^k), \mathbf{x}_\star - \mathbf{z}^k \rangle \quad (7.111)$$

Hence, we get

$$\begin{aligned}
 F(\mathbf{y}^{k+1}) &\leq (1 - \gamma_{t,k})F(\mathbf{y}^k) + \gamma_{t,k}F_\star + \gamma_{t,k} \langle \nabla F(\mathbf{z}^k), \mathbf{x}^{k+1} - \mathbf{x}^k \rangle + \frac{L\gamma_{t,k}^2}{2} \|\mathbf{x}^{k+1} - \mathbf{x}^k\|^2 \\
 &= (1 - \gamma_{t,k})F(\mathbf{y}^k) + \gamma_{t,k}F_\star + \gamma_{t,k} \langle \mathbf{v}^k + \beta_{t,k}(\mathbf{x}^{k+1} - \mathbf{x}^k), \mathbf{x}^{k+1} - \mathbf{x}_\star \rangle \\
 &\quad - \gamma_{t,k}\beta_{t,k} \langle \mathbf{x}^{k+1} - \mathbf{x}^k, \mathbf{x}^{k+1} - \mathbf{x}_\star \rangle + \gamma_{t,k} \langle \nabla F(\mathbf{z}^k) - \mathbf{v}^k, \mathbf{x}^{k+1} - \mathbf{x}_\star \rangle \\
 &\quad + \frac{L\gamma_{t,k}^2}{2} \|\mathbf{x}^{k+1} - \mathbf{x}^k\|^2 \\
 &\leq (1 - \gamma_{t,k})F(\mathbf{y}^k) + \gamma_{t,k}F_\star + \gamma_{t,k}\alpha_{t,k} - \gamma_{t,k}\beta_{t,k} \langle \mathbf{x}^{k+1} - \mathbf{x}^k, \mathbf{x}^{k+1} - \mathbf{x}_\star \rangle \\
 &\quad + \gamma_{t,k} \langle \nabla F(\mathbf{z}^k) - \mathbf{v}^k, \mathbf{x}^{k+1} - \mathbf{x}_\star \rangle + \frac{L\gamma_{t,k}^2}{2} \|\mathbf{x}^{k+1} - \mathbf{x}^k\|^2 \quad (7.112) \\
 &= (1 - \gamma_{t,k})F(\mathbf{y}^k) + \gamma_{t,k}F_\star + \gamma_{t,k}\alpha_{t,k} + \frac{\beta_{t,k}\gamma_{t,k}}{2} \left( \|\mathbf{x}^k - \mathbf{x}_\star\|^2 - \|\mathbf{x}^{k+1} - \mathbf{x}_\star\|^2 \right) \\
 &\quad + \gamma_{t,k} \left( \frac{L\gamma_{t,k} - \beta_{t,k}}{2} \|\mathbf{x}^{k+1} - \mathbf{x}^k\|^2 + \langle \nabla F(\mathbf{z}^k) - \mathbf{v}^k, \mathbf{x}^{k+1} - \mathbf{x}_\star \rangle \right) \\
 &\leq (1 - \gamma_{t,k})F(\mathbf{y}^k) - \gamma_{t,k}F_\star + \gamma_{t,k}\alpha_{t,k} + \frac{\beta_{t,k}\gamma_{t,k}}{2} \left( \|\mathbf{x}^k - \mathbf{x}_\star\|^2 - \|\mathbf{x}^{k+1} - \mathbf{x}_\star\|^2 \right) \\
 &\quad + \gamma_{t,k}D\|\nabla F(\mathbf{z}^k) - \mathbf{v}^k\|
 \end{aligned}$$

where the second inequality follows from the definition of  $\alpha_{t,k}$ , and the last inequality from the fact that  $\beta_{t,k} \geq L\gamma_{t,k}$  together with Cauchy-Schwarz inequality. Then, we subtract  $F_\star$  from both sides, take the expectation of both sides, and compute the telescopic sum over  $(t, k)$ . The first term vanishes due to the  $(1 - \gamma_{1,1}) = 0$  factor.

Denoting by  $\mathcal{E}_{t,k} := \mathbb{E}\|\mathbf{x}^k - \mathbf{x}_\star\|^2$ , and  $\mathcal{E}_{t,k+} := \mathbb{E}\|\mathbf{x}^{k+1} - \mathbf{x}_\star\|^2$ , we get

$$\begin{aligned} & \mathbb{E}[F(\mathbf{y}^{k+1}) - F_\star] \\ & \leq \sum_{(\tau,i)} \left[ \gamma_{\tau,i} \alpha_{\tau,i} + \frac{\beta_{\tau,i} \gamma_{\tau,i}}{2} (\mathcal{E}_{\tau,i} - \mathcal{E}_{\tau,i+}) + \gamma_{t,k} D \mathbb{E}\|\nabla F(\mathbf{z}^k) - \mathbf{v}^k\| \right] \prod_{(\tau',j)=(\tau,i)}^{(t,k)} (1 - \gamma_{\tau',j}). \end{aligned} \quad (7.113)$$

Remark that

$$\prod_{(\tau',j)=(\tau,i)}^{(t,k)} (1 - \gamma_{\tau',j}) = \prod_{r=i}^{K_\tau} \frac{s_{\tau,r} - 1}{s_{\tau,r} + 2} \prod_{\tau'=\tau+1}^{t-1} \prod_{j=1}^{K_{\tau'}} \frac{s_{\tau',j} - 1}{s_{\tau',j} + 2} \prod_{\ell=1}^k \frac{s_{t,\ell} - 1}{s_{t,\ell} + 2} = \frac{(s_{\tau,i} - 1) s_{\tau,i} (s_{\tau,i} + 1)}{s_{t,k} (s_{t,k} + 1) (s_{t,k} + 2)} \quad (7.114)$$

Now we focus on the individual terms

$$\sum_{(\tau,i)} \gamma_{\tau,i} \alpha_{\tau,i} \frac{(s_{\tau,i} - 1) s_{\tau,i} (s_{\tau,i} + 1)}{s_{t,k} (s_{t,k} + 1) (s_{t,k} + 2)} \leq \sum_{(\tau,i)} \frac{3}{s_{\tau,i} + 2} \frac{2LD^2}{s_{\tau,i}^2} \frac{(s_{\tau,i} - 1) s_{\tau,i} (s_{\tau,i} + 1)}{s_{t,k} (s_{t,k} + 1) (s_{t,k} + 2)} \quad (7.115)$$

Now, show that

$$\begin{aligned} & \sum_{(\tau,i)} \frac{\beta_{\tau,i} \gamma_{\tau,i}}{2} (\mathcal{E}_{\tau,i} - \mathcal{E}_{\tau,i+}) \prod_{(\tau',j)=(\tau,i)}^{(t,k)} (1 - \gamma_{\tau',j}) \\ & = \sum_{(\tau,i)} \frac{27L}{4} \frac{1}{(s_{\tau,i} + 2)^2} \frac{(s_{\tau,i} - 1) s_{\tau,i} (s_{\tau,i} + 1)}{s_{t,k} (s_{t,k} + 1) (s_{t,k} + 2)} (\mathcal{E}_{\tau,i} - \mathcal{E}_{\tau,i+}) \\ & = \frac{27L}{4s_{t,k} (s_{t,k} + 1) (s_{t,k} + 2)} \sum_{(\tau,i)} \frac{(s_{\tau,i} - 1) s_{\tau,i} (s_{\tau,i} + 1)}{(s_{\tau,i} + 2)^2} (\mathcal{E}_{\tau,i} - \mathcal{E}_{\tau,i+}) \end{aligned} \quad (7.116)$$

Remark that

$$\begin{aligned} & \sum_{s_{\tau,i}=1}^{s_{t,k}} \frac{(s_{\tau,i} - 1) s_{\tau,i} (s_{\tau,i} + 1)}{(s_{\tau,i} + 2)^2} \mathcal{E}_{\tau,i} - \sum_{s_{\tau,i}=1}^{s_{t,k}} \frac{(s_{\tau,i} - 1) s_{\tau,i} (s_{\tau,i} + 1)}{(s_{\tau,i} + 2)^2} \mathcal{E}_{\tau,i+} \\ & = \sum_{s_{\tau,i}=1}^{s_{t,k}-1} \frac{s_{\tau,i} (s_{\tau,i} + 1) (s_{\tau,i} + 2)}{(s_{\tau,i} + 3)^2} \mathcal{E}_{\tau,i+} - \sum_{s_{\tau,i}=1}^{s_{t,k}} \frac{(s_{\tau,i} - 1) s_{\tau,i} (s_{\tau,i} + 1)}{(s_{\tau,i} + 2)^2} \mathcal{E}_{\tau,i+} \\ & \leq \sum_{s_{\tau,i}=1}^{s_{t,k}} \frac{s_{\tau,i} (s_{\tau,i} + 1) (s_{\tau,i} + 2)}{(s_{\tau,i} + 3)^2} \mathcal{E}_{\tau,i+} - \sum_{s_{\tau,i}=1}^{s_{t,k}} \frac{(s_{\tau,i} - 1) s_{\tau,i} (s_{\tau,i} + 1)}{(s_{\tau,i} + 2)^2} \mathcal{E}_{\tau,i+} \\ & \leq D^2 \sum_{s_{\tau,i}=1}^{s_{t,k}} \left( \frac{s_{\tau,i} (s_{\tau,i} + 1) (s_{\tau,i} + 2)}{(s_{\tau,i} + 3)^2} - \frac{(s_{\tau,i} - 1) s_{\tau,i} (s_{\tau,i} + 1)}{(s_{\tau,i} + 2)^2} \right) \\ & \leq D^2 s_{t,k} \end{aligned} \quad (7.117)$$

Hence, we get

$$\sum_{(\tau,i)} \frac{\beta_{\tau,i} \gamma_{\tau,i}}{2} (\mathcal{E}_{\tau,i} - \mathcal{E}_{\tau,i+}) \prod_{(\tau',j)=(\tau,i)}^{(t,k)} (1 - \gamma_{\tau',j}) \leq \frac{27LD^2}{4(s_{t,k} + 1)(s_{t,k} + 2)} \quad (7.118)$$

Finally, we focus on the last term:

$$\begin{aligned} & \sum_{(\tau,i)} \gamma_{\tau,i} DE \|\nabla F(\mathbf{z}^{\tau,i}) - \mathbf{v}^{\tau,i}\| \frac{(s_{\tau,i} - 1)s_{\tau,i}(s_{\tau,i} + 1)}{s_{t,k}(s_{t,k} + 1)(s_{t,k} + 2)} \\ & \leq LD^2 \sum_{s_{\tau,i}=1}^{s_{t,k}} \frac{6\sqrt{2}}{(s_{\tau,i} + 2)(s_{\tau,i} + 1)^2} \frac{(s_{\tau,i} - 1)s_{\tau,i}(s_{\tau,i} + 1)}{s_{t,k}(s_{t,k} + 1)(s_{t,k} + 2)} \\ & \leq \frac{6\sqrt{2}LD^2}{(s_{t,k} + 1)(s_{t,k} + 2)} \end{aligned} \quad (7.119)$$

Combining these bounds, we obtain

$$\mathbb{E}[F(\mathbf{y}^{k+1}) - F_{\star}] = \mathcal{O}\left(\frac{LD^2}{(s_{t,k} + 1)(s_{t,k} + 2)}\right) \quad (7.120)$$

Easy to verify by induction that  $K_t \leq s_{t,k} \leq 4K_t$ . Hence,  $s_{t,k} = \Theta(K_t) = \Theta(2^{t/2})$ .

Hence,  $\mathbb{E}[F(\mathbf{y}^{k+1}) - F_{\star}] = \mathcal{O}(LD^2 2^{-t})$ . Therefore, to get  $\epsilon$ -solution, we set  $T_{\epsilon}$  as

$$T_{\epsilon} = \Theta\left(\log_2\left(\frac{LD^2}{\epsilon}\right)\right) \quad (7.121)$$

Then, to achieve this accuracy, we can calculate *ifo* complexity as

$$\begin{aligned} \#(ifo) &= \sum_{t=1}^{T_{\epsilon}} \left( Q_t + \sum_{k=2}^{K_t} S_{t,k} \right) = \mathcal{O}\left(\sum_{t=1}^{T_{\epsilon}} \left( n + \sum_{k=2}^{K_t} 2^{3t/2} \right)\right) = \mathcal{O}\left(nT_{\epsilon} + \sum_{t=1}^{T_{\epsilon}} 2^{2t}\right) \\ &= \mathcal{O}(nT_{\epsilon} + 2^{2T_{\epsilon}}) \\ &= \mathcal{O}\left(n \ln\left(\frac{LD^2}{\epsilon}\right) + \frac{L^2 D^4}{\epsilon^2}\right) \end{aligned} \quad (7.122)$$

Finally, we can find *lmo* complexity by using Lemma 7.11:

$$\#(lmo) = \mathcal{O}\left(\sum_{t=1}^{T_{\epsilon}} \sum_{k=1}^{K_t} \frac{4\beta_{t,k} D^2}{\alpha_{t,k}}\right) = \mathcal{O}\left(\sum_{t=1}^{T_{\epsilon}} K_t^2\right) = \mathcal{O}(2^{T_{\epsilon}}) = \mathcal{O}\left(\frac{LD^2}{\epsilon}\right) \quad (7.123)$$

**Proof for Theorem 7.6 and Corollary 7.6**

Similar to the proof of finifte-sum setting, but we use (7.99) instead of (7.98), hence the constants at (7.119) change.

To get  $\epsilon$ -solution, we can calculate  $sfo$  complexity as

$$\begin{aligned} \#(sfo) &= \sum_{t=1}^{T_\epsilon} \left( Q_t + \sum_{k=2}^{K_t} S_{t,k} \right) = \mathcal{O} \left( \sum_{t=1}^{T_\epsilon} \left( \frac{\sigma^2 K_t^4}{L^2 D^2} + \sum_{k=2}^{K_t} 2^{3t/2} \right) \right) = \mathcal{O} \left( \left( \frac{\sigma^2}{L^2 D^2} + 1 \right) \sum_{t=1}^{T_\epsilon} 2^{2t} \right) \\ &= \mathcal{O} \left( \left( \frac{\sigma^2}{L^2 D^2} + 1 \right) 2^{T_\epsilon} \right) = \mathcal{O} \left( \frac{\sigma^2 D^2 + L^2 D^4}{\epsilon^2} \right) \end{aligned} \quad (7.124)$$

The  $lmo$  complexity is same as the convex finite-sum case.

**Proof for Theorem 7.7 and Corollary 7.7**

We start by (7.35), and rearrange it to obtain

$$F(\mathbf{y}^{k+1}) \leq F(\mathbf{y}^k) - \gamma \mathcal{G}(\mathbf{y}^k) + \gamma \langle \nabla F(\mathbf{y}^k) - \mathbf{v}^k, \mathbf{x}^{k+1} - \mathbf{w}_\star^k \rangle + \gamma \langle \mathbf{v}^k, \mathbf{x}^{k+1} - \mathbf{w}_\star^k \rangle + \gamma^2 \frac{L}{2} D^2 \quad (7.125)$$

where  $\mathbf{w}_\star^k = \arg \max_{\mathbf{x} \in \mathcal{X}} \langle \mathbf{x}, -\nabla F(\mathbf{y}^k) \rangle$ .

Focus on the last inner-product term

$$\begin{aligned} \gamma \langle \mathbf{v}^k, \mathbf{x}^{k+1} - \mathbf{w}_\star^k \rangle &= \gamma \langle \mathbf{v}^k + \beta(\mathbf{x}^{k+1} - \mathbf{x}^k), \mathbf{x}^{k+1} - \mathbf{w}_\star^k \rangle - \gamma \beta \langle \mathbf{x}^{k+1} - \mathbf{x}^k, \mathbf{x}^{k+1} - \mathbf{w}_\star^k \rangle \\ &\leq \gamma \alpha - \gamma \beta \langle \mathbf{x}^{k+1} - \mathbf{x}^k, \mathbf{x}^{k+1} - \mathbf{w}_\star^k \rangle \\ &\leq \gamma \alpha + \gamma \beta D^2 \end{aligned} \quad (7.126)$$

The first inequality follows from the role of  $\alpha$  in CndG, and the second from Cauchy-Schwarz.

Now we use

$$\begin{aligned} \gamma \langle \nabla F(\mathbf{y}^k) - \mathbf{v}^k, \mathbf{x}^{k+1} - \mathbf{w}_\star^k \rangle &\leq \gamma D \|\nabla F(\mathbf{y}^k) - \mathbf{v}^k\| \\ &\leq \gamma D \|\nabla F(\mathbf{y}^k) - \nabla F(\mathbf{z}^k)\| + \gamma D \|\nabla F(\mathbf{z}^k) - \mathbf{v}^k\| \\ &\leq \gamma L D \|\mathbf{y}^k - \mathbf{z}^k\| + \gamma D \|\nabla F(\mathbf{z}^k) - \mathbf{v}^k\| \\ &= \gamma^2 L D \|\mathbf{x}^k - \mathbf{y}^k\| + \gamma D \|\nabla F(\mathbf{z}^k) - \mathbf{v}^k\| \\ &= \gamma^2 L D^2 + \gamma D \|\nabla F(\mathbf{z}^k) - \mathbf{v}^k\| \end{aligned} \quad (7.127)$$

Combining these bounds, we obtain

$$\begin{aligned} F(\mathbf{y}^{k+1}) &\leq F(\mathbf{y}^k) - \gamma \mathcal{G}(\mathbf{y}^k) + \gamma D \|\nabla F(\mathbf{z}^k) - \mathbf{v}^k\| + \gamma \alpha + \gamma \beta D^2 + \frac{3}{2} \gamma^2 L D^2 \\ &= F(\mathbf{y}^k) - \frac{1}{\sqrt{TK}} \mathcal{G}(\mathbf{y}^k) + \frac{D}{\sqrt{TK}} \|\nabla F(\mathbf{z}^k) - \mathbf{v}^k\| + \frac{4LD^2}{TK} \end{aligned} \quad (7.128)$$

## Chapter 7. CGM with Stochastic Path-Integrated Differential Estimator

---

Now we take expectation of both sides and use (7.100), and we obtain

$$\mathbb{E}[F(\mathbf{y}^{k+1})] \leq \mathbb{E}[F(\mathbf{y}^k)] - \frac{1}{\sqrt{TK}} \mathbb{E}[\mathcal{G}(\mathbf{y}^k)] + \frac{6LD^2}{TK} \quad (7.129)$$

Rearranging, and summing over all  $(t, k)$  pairs up to  $(T, K)$ , we get

$$\frac{1}{\sqrt{TK}} \sum_{(\tau, l)} \mathbb{E}[\mathcal{G}(\mathbf{y}^{\tau, l})] \leq F(\bar{\mathbf{x}}^1) - F(\mathbf{x}_\star) + 6LD^2 \quad (7.130)$$

Hence, by definition of  $\mathbf{y}^{\text{out}}$ , we get

$$\mathbb{E}[\mathcal{G}(\mathbf{y}^{\text{out}})] \leq \frac{\mathcal{E} + 6LD^2}{\sqrt{TK}} \quad (7.131)$$

This completes the convergence rate proof.

Proof for the *ifo* complexity follows similarly to the one for SPIDER-FW.

To show *lmo* complexity, we use

$$\#(lmo) = \mathcal{O} \left( \sum_{t=1}^{T_\epsilon} \sum_{k=1}^K \frac{4\beta D^2}{\alpha} \right) = \mathcal{O} \left( \sum_{t=1}^{T_\epsilon} \sum_{k=1}^K 6 \right) = \mathcal{O}(6KT_\epsilon) = \mathcal{O} \left( (\mathcal{E}^2 + L^2 D^4) \frac{1}{\epsilon^2} \right) \quad (7.132)$$

### Proof for Theorem 7.8 and Corollary 7.8

Follows similarly as in the non-convex finite-sum case, but we use (7.101) instead of (7.100), hence we have additional  $\epsilon/2$  term on the right-hand-side of (7.131). *sfo* complexity follows similarly as in SPIDER-FW. *lmo* complexity is the same as the finite-sum case.

## 8 Conclusion & Future Directions

This chapter presents an overview of our main contributions along with some future directions.

In this dissertation, we considered a fundamental class of constrained convex optimization problem (1.1) with broad applications. In particular, we focused on an SDP formulation (1.4), which can be written as a special instance of this problem template. As we described in Chapter 1, storing the matrix decision variable is the critical bottleneck that prevents us from solving many important applications of this SDP formulation at scale.

- In Chapter 2, we introduced a novel convex optimization paradigm for designing algorithms with optimal storage cost, which constructs an approximate solution with rigorous guarantees to large-scale problems with structured solutions. The two key ingredients are to maintain only a sketch of the large decision variable and to use algorithmic templates that are compatible with these sketches. A natural target for this approach is the class of matrix optimization problems with low-rank solutions. We presented the first storage-optimal algorithm for this problem class, SketchyCGM, which modifies the standard conditional gradient method (CGM) to work with sketches. In spite of recent enthusiasm for nonconvex heuristics, our research justifies that the convex optimization did not reach yet its limits of scalability and that it still remains as a viable option at large-scale.

Based on this paradigm, our recent research is focused on designing a practical convex optimization solver for a broad set of semidefinite programming (SDP) instances. Some main ingredients for this design are presented in the subsequent sections. The detailed analysis and implementation of the resulting method for this specific template, with potent numerical demonstrations in cutting-edge applications, is left for future work.

Our paradigm is beyond low-rank optimization. As an immediate future direction, we may consider convex optimization problems with other types of structured solutions. There are different dimensionality reduction techniques to store and retrieve structured objects. One example in this line is designing a storage-optimal convex optimization algorithm for problems with sparse or group-sparse solutions, with potential applications in optimal transport.

- In Chapter 3, we proposed and studied new primal-dual optimization algorithms for solving a generic constrained convex optimization template, which covers our model SDP formulation (1.4) as a special instance. The proposed methods, UPD and its accelerated variant AUPD, have several key advantages compared to the existing primal-dual methods in the literature. As opposed to the smoothing (of the dual problem), we preserve the original structure of the dual problem. Then, we apply a (sub)gradient method to this dual problem, and recover the primal variable as a weighted sum of the solutions to the primal subproblem. This primal averaging sequence is crucial for our storage-optimal optimization paradigm. Moreover, our algorithms are adaptive to the unknown smoothness level (Hölder-smoothness order) of the dual problem, as our analysis leads to optimal iteration complexity guarantees in the sense of first-order black-box models [NY83].

Unfortunately, these algorithms have some substantial limitations at the current stage. Our algorithms borrow a line-search strategy introduced by Nesterov [Nes15] for his universal gradient methods in the unconstrained setting, to achieve adaptivity in the dual problem. Unfortunately, this line-search procedure, hence also our algorithms explicitly depend on the target accuracy input  $\varepsilon$ , and the methods only guarantee convergence towards an  $\varepsilon$ -suboptimal set. As a consequence, it is difficult to tune these algorithms without a rough knowledge of the optimal value. Also note that oracle errors can cause these algorithms to get stuck in the line-search procedure.

The literature on adaptive optimization has recently evolved, following the increased interest in stochastic and online learning settings. AdaGrad and its variants attempt to adapt the learning rate using an accumulation of the gradient norms. These developments in the stochastic setting can be translated to the deterministic setting via some online to offline conversion mechanisms to design first-order algorithms that adapt to the smoothness, see [CBCG04, Lev17]. In this line, we introduced an accelerated adaptive method [LYC18]. Our method, AcceleGrad, can be viewed as an accelerated variant of AdaGrad in the deterministic setting. In contrast to AdaGrad, AcceleGrad achieves optimal rates for both smooth and non-smooth problems. It also has convergence guarantees in the stochastic setting.

These results, so far, remained in the unconstrained optimization setting. An important future direction with significant potential impact is extending these results to the constrained setting. We might expect the resulting algorithms to maintain the key benefits of our universal primal-dual convex optimization framework, while addressing most of its seeming limitations. In various problem instances, we observed significant empirical improvement with a heuristic primal-dual variant of AcceleGrad, against UPD and AUPD.

- In Chapter 4, we presented an extension of CGM for solving a generic composite convex minimization template over a convex and compact domain. Our method (HCGM) combines the ideas of smoothing (of the primal problem) and homotopy for the smoothness parameter under the standard CGM template. By means of set-indicator functions, HCGM can also be applied to the problems with affine constraints, and in particular to our model SDP formu-



---

lation (1.4). In this case, we can view HCGM as a quadratic penalty method. We proved that the method achieves an optimal  $\mathcal{O}(1/\sqrt{t})$  convergence rate, which holds even when the linear minimization oracles are noisy with additive or multiplicative errors.

The numerical evaluation of HCGM for the clustering SDP (in Section 4.4.1) also reveals an important phenomenon that supports our claims about trading accuracy for scalability. It is clear from Figure 4.1 that the overall error after rounding (the misclassification rate in this case) saturates after a low-to-medium accuracy level (that we achieved after  $\sim 1000$  iterations in this example), and solving the optimization problem beyond this accuracy does not lead to any further useful information. This observation can be generalized for the majority of the SDP formulations obtained as the relaxation of a combinatorial decision problem, excluding the examples where the empirical study suggests that the relaxation is exact.

- The slow empirical convergence of HCGM strains its practical impact in solving large-scale SDP instances. In Chapter 5, we address this issue by introducing a natural extension of HCGM, going from quadratic penalty to an augmented Lagrangian formulation. The new method, CGAL, retains the strong theoretical guarantees of HCGM, but it exhibits significantly superior empirical performance. This improvement is due to the dual updates, which can be viewed as shifting the center point of the quadratic penalty in a favorable direction adaptively. In a selection of numerical test setups with SDP formulations, we observed empirical  $\mathcal{O}(1/t)$  convergence rate for CGAL, compared against  $\mathcal{O}(1/\sqrt{t})$  rate of HCGM.

An immediate future direction is finding a theoretical explanation for this performance improvement, beyond an educated intuition. Remark that we empirically observed  $\mathcal{O}(1/t)$  convergence rate for CGAL in various numerical experiments. It might be possible to prove this rate, maybe under more stringent assumptions such as a quadratic growth condition.

HCGM and CGAL preserve the key features of the standard CGM, such as the cheap linear minimization oracles and the structured updates on the decision variable. As a result, they are also amenable to sketching for storage-optimal extensions. Considering its superior empirical performance; CGAL, in particular, might be the key ingredient for the next phase of large-scale SDP solvers.

- Our storage-optimal optimization paradigm is based on a sketching model to truncate the storage costs. Hence, designing practical sketching algorithms for low-rank approximation is a fundamental aspect of our research. In Chapter 6, we summarized our contributions.

Based on the specific role of sketching in our approach, we focus on minimizing the sketch size. Moreover, our sketching algorithms are specifically intended for environments where the data matrix is presented as a stream of linear updates. Our aim is to address questions that arise when using sketching algorithms in practice, that had been often neglected in the literature. Unlike the majority of the existing works, our algorithms are accompanied by informative error bounds, which provide guidance when we set the sketch-size parameters. In addition, we also address the numerical stability problems.

We designed a benchmark for numerical evaluation of streaming low-rank matrix approximation, with synthetic and real datasets. We presented extensive numerical experiments to compare our methods and the other methods in the literature, in order to identify methods that produce the best approximations in practice, under various scenarios. We implemented our methods as well as the other methods that we consider in these experiments, treating the sketch as an abstract data type using ideas from object-oriented programming. Beyond optimization, we also presented concrete applications on “on-the-fly compression” of large-scale scientific data. We omit most of these details in this dissertation for brevity, we refer to our original manuscripts for these details. We published our code online, which is available at “<https://github.com/alpyurtsever/SKETCH>” as a software toolbox for MATLAB. This toolbox also includes various dimensionality reduction maps, optional *a posteriori* error estimator that allows the user to validate the quality of the reconstructed matrix, as well as the codes to regenerate the results of our numerical experiments.

See the follow-up paper [SGL<sup>+</sup>19] of our collaborators on an algorithm for low-rank Tucker approximation of a tensor from streaming data as an extension in this line of research.

- Lastly, in Chapter 7, we studied the theoretical complexity of the CGM variants for stochastic expectation minimization and finite-sum formulations, in order to identify and present the tightest results known so far. In this chapter, we also introduced a class of novel variance-reduced stochastic CGM variants, based on the recent stochastic path-integrated differential estimator (SPIDER) technique [FLLZ18].

An interesting future direction is to extend these algorithms for solving problems with affine constraints, in a similar way to HCGM or CGAL. In this line, we have already taken an initial step in our recent work [LYFC19], by combining our homotopy smoothing approach with a stochastic CGM framework by Mokhtari et al. [MHK18], and demonstrated the application of the proposed method for solving some basic stochastic SDP problems. While the guaranteed rate of our method is not very fast, this is the first stochastic CGM variant for solving problems with affine constraints. We can expect to achieve better rates using our SPIDER-CGM approach as the base method. Developing a well-tuned implementation is another important piece of future work.

*SDP formulations are key in signal processing, machine learning, and other engineering applications. Most aspects of my research are immediately applicable in addition to their long term promise. By greatly enhancing the scalability of optimization methods for solving SDP problems, we can potentially unlock new results in important scientific applications.*

# Bibliography

- [AB15] A. Agarwal and L. Bottou. A lower bound for the optimization of finite sums. In *Proc. 32nd Int. Conf. Machine Learning*, 2015.
- [Ali91] F. Alizadeh. *Combinatorial optimization with interior point methods and semi-definite matrices*. PhD thesis, University of Minnesota, 1991.
- [Ali93] F. Alizadeh. Interior point methods in semidefinite programming with applications to combinatorial optimization. *SIAM J. Optim.*, 5(1):13–51, 1993.
- [ARR14] A. Ahmed, B. Recht, and J. Romberg. Blind deconvolution using convex programming. *IEEE Trans. on Inf. Theory*, 60(3):1711–1732, 2014.
- [AT06] A. Auslender and M. Teboulle. Interior gradient and proximal methods for convex and conic optimization. *SIAM J. Optim.*, 16(3):697–725, 2006.
- [Bac15] F. Bach. Duality between subgradient and conditional gradient methods. *SIAM J. Optim.*, 25(1):115–129, 2015.
- [Bar95] A. Barvinok. Problems of distance geometry and convex properties of quadratic maps. *Discrete & Computational Geometry*, 13(2):189–202, 1995.
- [BBV16] A.S. Bandeira, N. Boumal, and V. Voroninski. On the low-rank approach for semidefinite programs arising in synchronization and community detection. In *29th Annual Conference on Learning Theory*, 2016.
- [BC11] Heinz H. Bauschke and Patrick L. Combettes. *Convex analysis and monotone operator theory in Hilbert spaces*. Springer, New York, NY, 2011.
- [Ber76] D. P. Bertsekas. On penalty and multiplier methods for constrained minimization. *SIAM J. Control Optim.*, 14(2):216–235, 1976.
- [Ber99] D. Bertsekas. *Nonlinear Programming*. Athena Scientific, 2nd edition, 1999.
- [BKS16] S. Bhojanapalli, A. Kyrillidis, and S. Sanghavi. Dropping convexity for faster semi-definite optimization. *J. Mach. Learn. Res.*, 49:1–53, 2016.

## Bibliography

---

- [BM03] S. Burer and R. D. C. Monteiro. A nonlinear programming algorithm for solving semidefinite programs via low-rank factorization. *Math. Program.*, 95(2, Ser. B):329–357, 2003.
- [BNPS17] J. Bolte, T. P. Nguyen, J. Peypouquet, and B. W. Suter. From error bounds to the complexity of first-order descent methods for convex functions. *Math. Program.*, 165(2):471–507, 2017.
- [Bou15] N. Boumal. A riemannian low-rank method for optimization over semidefinite matrices with block-diagonal constraints. arXiv:1506.00575v2, 2015.
- [BPC<sup>+</sup>11] S. Boyd, N. Parikh, E. Chu, B. Peleato, and J. Eckstein. Distributed optimization and statistical learning via the alternating direction method of multipliers. *Found. Trends Machine Learning*, 3(1):1–122, 2011.
- [Bra06] M. Brand. Fast low-rank modifications of the thin singular value decomposition. *Linear Algebra Appl.*, 415(1):20–30, 2006.
- [BT09] A. Beck and M. Teboulle. A fast iterative shrinkage-thresholding algorithm for linear inverse problems. *SIAM J. Imaging Sciences*, 2(1):183–202, 2009.
- [BVB18] N. Boumal, V. Voroninski, and A.S. Bandeira. Deterministic guarantees for Burer–Monteiro factorizations of smooth semidefinite programs. arXiv:1804.02008v1, 2018.
- [BWZ16] C. Boutsidis, D. Woodruff, and P. Zhong. Optimal principal component analysis in distributed and streaming models. In *Proc. 48th ACM Symp. Theory of Computing*, 2016.
- [CBCG04] N. Cesa-Bianchi, A. Conconi, and C. Gentile. On the generalization ability of on-line learning algorithms. *IEEE Trans. on Inf. Theory*, 50(9):2050–2057, 2004.
- [CD13] J. Chiu and L. Demanet. Sublinear randomized algorithms for skeleton decompositions. *SIAM J. Matrix Anal. Appl.*, 34(3):1361–1383, 2013.
- [CESV13] E. J. Candès, Y. C. Eldar, T. Strohmer, and V. Voroninski. Phase retrieval via matrix completion. *SIAM J. Imaging Sciences*, 6(1):199–225, 2013.
- [CJN17] B. Cox, A. Juditsky, and A. Nemirovski. Decomposition techniques for bilinear saddle point problems and variational inequalities with affine monotone operators. *J. Optim. Theory Appl.*, 2(402–435), 2017.
- [CKS15] K. N. Chaudhury, Y. Khoo, and A. Singer. Global registration of multiple point clouds using semidefinite programming. *SIAM J. Optim.*, 25(1):468–501, 2015.
- [Cla10] K. L. Clarkson. Coresets, sparse greedy approximation, and the Frank-Wolfe algorithm. *ACM Trans. Algorithms*, 6(4), 2010.

- 
- [CLDM18] H. Chang, Y. Lou, Y. Duan, and S. Marchesini. Total variation-based phase retrieval for poisson noise removal. *SIAM J. Imaging Sciences*, 11(1):24–55, 2018.
  - [CLS15a] E. J. Candès, X. Li, and M. Soltanolkotabi. Phase retrieval from coded diffraction patterns. *Applied and Computational Harmonic Analysis*, 39(2):277–299, 2015.
  - [CLS15b] E. J. Candès, X. Li, and M. Soltanolkotabi. Phase retrieval via Wirtinger Flow: Theory and algorithms. *IEEE Trans. Inform. Theory*, 61(4):1985–2007, 2015.
  - [CP08] P. L. Combettes and J.-C. Pesquet. A proximal decomposition method for solving convex variational inverse problems. *Inverse Problems*, 24(6):065014, 2008.
  - [CR12] E. Candès and B. Recht. Exact matrix completion via convex optimization. *Communications of the ACM*, 55(6):111–119, 2012.
  - [CSV12] E. J. Candès, T. Strohmer, and V. Voroninski. Phaselift: Exact and stable signal recovery from magnitude measurements via convex programming. *Communications on Pure and Applied Mathematics*, 2012.
  - [CVY18] V. Cevher, B. C. Vu, and A. Yurtsever. Stochastic forward Douglas-Rachford splitting method for monotone inclusions. In P. Giselsson and A. Rantzer, editors, *Large-Scale and Distributed Optimization*, chapter 7, pages 149–179. Springer International Publishing, 2018.
  - [CVY19] V. Cevher, B. C. Vu, and A. Yurtsever. Inertial three-operator splitting method and applications. arXiv:1904.12980v1, 2019.
  - [CZH<sup>+</sup>17] R. Chandra, Z. Zhong, J. Hontz, V. McCulloch, C. Studer, and T. Goldstein. Phasepack: A phase retrieval library. *Asilomar Conference on Signals, Systems, and Computers*, 2017.
  - [DBLJ14] A. Defazio, F. Bach, and S. Lacoste-Julien. SAGA: A fast incremental gradient method with support for non-strongly convex composite objectives. In *Advances in Neural Information Processing Systems 28*, 2014.
  - [DFTJ16] C. Dünner, S. Forte, M. Takác, and M. Jaggi. Primal–dual rates and certificates. In *Proc. 33rd Int. Conf. Machine Learning*, 2016.
  - [dGJL07] A. d’Aspremont, L.E. Ghaoui, M.I. Jordan, and G.R.G. Lanckriet. A direct formulation for sparse PCA using semidefinite programming. *SIAM Review*, 49(3):434–448, 2007.
  - [DH78] J. Dunn and S. Harshbarger. Conditional gradient algorithms with open loop step size rules. *J. Math. Anal. Appl.*, 62(2):432–444, 1978.
  - [DHS11] J. Duchi, E. Hazan, and Y. Singer. Adaptive subgradient methods for online learning and stochastic optimization. *J. Mach. Learn. Res.*, 12:2121–2159, 2011.

## Bibliography

---

- [DM05] P. Drineas and M. W. Mahoney. On the Nystrom method for approximating a Gram matrix for improved kernel-based learning. *J. Mach. Learn. Res.*, 6:2153–2175, 2005.
- [Dun79] J. Dunn. Rates of convergence for conditional gradient algorithms near singular and nonsingular extremals. *SIAM J. Control Optim.*, 17(2):187–211, 1979.
- [Dun80] J. Dunn. Convergence rates for conditional gradient sequences generated by implicit step length rules. *SIAM J. Control Optim.*, 18(5):473–487, 1980.
- [DYC<sup>+</sup>19] L. Ding, A. Yurtsever, V. Cevher, J. A. Tropp, and M. Udell. An optimal-storage approach to semidefinite programming using approximate complementarity. arXiv:1902.03373v1, 2019.
- [Faz02] M. Fazel. *Matrix rank minimization with applications*. PhD thesis, Stanford Univ., 2002.
- [FBCM04] C. Fowlkes, S. Belongie, F. Chung, and J. Malik. Spectral grouping using the Nystrom method. *IEEE Trans. Pattern Anal. Mach. Intell.*, 26(2):214–225, Jan. 2004.
- [FGM17] R. M. Freund, P. Grigas, and R. Mazumder. An extended Frank–Wolfe method with "in-face" directions and its application to low-rank matrix completion. *SIAM J. Optim.*, 27(1):319–346, 2017.
- [FLLZ18] C. Fang, C. J. Li, Z. Lin, and T. Zhang. Spider: Near-optimal non-convex optimization via stochastic path-integrated differential estimator. In *Advances in Neural Information Processing Systems 32*, 2018.
- [FW56] M. Frank and P. Wolfe. An algorithm for quadratic programming. *Naval Research Logistics Quarterly*, 3:95–110, 1956.
- [Gar12] D. Garber. *Approximating Semidefinite Programs in Sublinear Time*. PhD thesis, Israel Institute of Technology, 2012.
- [GEB13] T. Goldstein, E. Esser, and R. Baraniuk. Adaptive primal-dual hybrid gradient methods for saddle point problems. arXiv:1305.0546v1, 2013.
- [GH11] D. Garber and E. Hazan. Approximating semidefinite programs in sublinear time. In *Advances in Neural Information Processing Systems 25*, 2011.
- [GH16] D. Garber and E. Hazan. Sublinear time algorithms for approximate semidefinite programming. *Math. Program., Ser. A*, (158):329–361, 2016.
- [Git11] A. Gittens. The spectral norm error of the naïve Nystrom extension. arXiv:1110.5305, 2011.

- 
- [Git13] A. Gittens. *Topics in Randomized Numerical Linear Algebra*. PhD thesis, California Institute of Technology, 2013.
  - [GJLJ17] G. Gidel, T. Jebara, and S. Lacoste-Julien. Frank-Wolfe algorithms for saddle point problems. In *Proc. 20th Int. Conf. Artificial Intelligence and Statistics*, 2017.
  - [GLF<sup>+</sup>10] D. Gross, Y.-K. Liu, S. T. Flammia, S. Becker, and J. Eisert. Quantum state tomography via compressed sensing. *Phys. Rev. Lett.*, 105, 2010.
  - [GM13] A. Gittens and M. W. Mahoney. Revisiting the Nyström method for improved large-scale machine learning. arXiv:1303.1849, 2013.
  - [GM16] A. Gittens and M. W. Mahoney. Revisiting the Nyström method for improved large-scale machine learning. *J. Mach. Learn. Res.*, 17:Paper No. 117, 65, 2016.
  - [GPLJ18] G. Gidel, F. Pedregosa, and S. Lacoste-Julien. Frank-Wolfe splitting via augmented Lagrangian method. In *Proc. 21st Int. Conf. Artificial Intelligence and Statistics*, 2018.
  - [GW95] M. X. Goemans and D. P. Williamson. Improved approximation algorithms for maximum cut and satisfiability problems using semidefinite programming. *J. Assoc. Comput. Mach.*, 42(6):1115–1145, 1995.
  - [Ham84] H. J. Hammond. *Solving asymmetric variational inequality problems and systems of equations with generalized nonlinear programming algorithms*. PhD thesis, MIT, 1984.
  - [Haz08] E. Hazan. Sparse approximate solutions to semidefinite programs. In *Proc. 8th Latin American Conf. Theoretical Informatics*, pages 306–316, 2008.
  - [HCO<sup>+</sup>15] R. Horstmeyer, R. Y. Chen, X. Ou, B. Ames, J. A. Tropp, and C. Yang. Solving ptychography with a convex relaxation. *New J. Physics*, 17(5):053044, 2015.
  - [HH15] N. He and Z. Harchaoui. Semi-proximal mirror-prox for nonsmooth composite minimization. In *Advances in Neural Information Processing Systems 28*, 2015.
  - [HJL96] C. L. Hamilton-Jester and C.-K. Li. Extreme vectors of doubly nonnegative matrices. *Rocky Mountain J. Math.*, 26(4):1371–1383, 1996.
  - [HJN15] Z. Harchaoui, A. Juditsky, and A. Nemirovski. Conditional gradient algorithms for norm-regularized smooth convex optimization. *Math. Program., Ser. A*, (152):75–112, 2015.
  - [HK12] E. Hazan and S. Kale. Projection-free online learning. In *Proc. 29th Int. Conf. Machine Learning*, 2012.
  - [HK16] F. M. Harper and J. A. Konstan. The movielens datasets: History and context. *Acm Trans. Interactive Intelligent Systems*, 5(4):19, 2016.



## Bibliography

---

- [HKMS19] H. Hassani, A. Karbasi, A. Mokhtari, and Z. Shen. Stochastic conditional gradient++. *arXiv:1902.06992*, 2019.
- [HL16] E. Hazan and H. Luo. Variance-reduced and projection-free stochastic optimization. In *Proc. 33rd Int. Conf. Machine Learning*, 2016.
- [HMT11] N. Halko, P. G. Martinsson, and J. A. Tropp. Finding structure with randomness: probabilistic algorithms for constructing approximate matrix decompositions. *SIAM Rev.*, 53(2):217–288, 2011.
- [HUS10] J.-B. Hiriart-Urruty and A. Seeger. A variational approach to copositive matrices. *SIAM Review*, 52(4):593–629, 2010.
- [Jag13] M. Jaggi. Revisiting Frank–Wolfe: Projection–free sparse convex optimization. In *Proc. 30th Int. Conf. Machine Learning*, 2013.
- [JN16] A. Juditsky and A. Nemirovski. Solving variational inequalities with monotone operators on domains given by Linear Minimization Oracles. *Math. Program.*, 156(1-2):221–256, 2016.
- [JNS13] P. Jain, P. Netrapalli, and S. Sanghavi. Low-rank matrix completion using alternating minimization. In *Proc. 45th Ann. ACM Symp. Theory of Computing*, 2013.
- [JS10] M. Jaggi and M. Sulovský. A Simple Algorithm for Nuclear Norm Regularized Problems. In *Proc. 27th Int. Conf. Machine Learning*, 2010.
- [JZ13] R. Johnson and T. Zhang. Accelerating stochastic gradient descent using predictive variance reduction. In *Advances in Neural Information Processing Systems 26*, 2013.
- [KN12] S. Khot and A. Naor. Grothendieck-type inequalities in combinatorial optimization. *Communications on Pure and Applied Mathematics*, 65(7):992–1035, 2012.
- [Lan14] G. Lan. The complexity of large-scale convex programming under a linear optimization oracle. *arXiv:1309.5550v2*, 2014.
- [LCG<sup>+</sup>04] G.R.G. Lanckriet, N. Cristianini, L.E. Ghaoui, P. Bartlett, and M.I. Jordan. Learning the kernel matrix with semidefinite programming. *J. Mach. Learn. Res.*, 5:27–72, 2004.
- [Lev17] K. Y. Levy. Online to offline conversions, universality and adaptive minibatch sizes. In *Advances in Neural Information Processing Systems*, pages 1612–1621, 2017.
- [LF18] H. Lu and R. M. Freund. Generalized stochastic Frank-wolfe algorithm with stochastic" substitute"gradient for structured convex optimization. *arXiv:1807.07680*, 2018.



- 
- [LJ16] S. Lacoste-Julien. Convergence rate of frank-wolfe for non-convex objectives. arXiv:1607.00345, 2016.
  - [LJJP13] S. Lacoste-Julien, M. Jaggi, M. Schmidt, and P. Pletscher. Block-coordinate Frank-Wolfe optimization for structural SVMs. In *Proc. 30th Int. Conf. Machine Learning*, 2013.
  - [LKTJ17] F. Locatello, R. Khanna, M. Tschannen, and M. Jaggi. A unified optimization view on generalized matching pursuit and Frank-Wolfe. In *Proc. 20th Int. Conf. Artificial Intelligence and Statistics*, 2017.
  - [LL12] J. Lavaei and H.L. Low. Zero duality gap in optimal power flow problem. *IEEE Trans. on Power Syst.*, 27(1):92–107, February 2012.
  - [LLM19] Y.-F. Liu, X. Liu, and S. Ma. On the non-ergodic convergence rate of an inexact augmented Lagrangian framework for composite convex programming. *Mathematics of Operations Research*, 44(2):632–650, 2019.
  - [LLS<sup>+</sup>17] H. Li, G. C. Linderman, A. Szlam, K. P. Stanton, Y. Kluger, and M. Tygert. Algorithm 971: An implementation of a randomized algorithm for principal component analysis. *ACM Trans. Math. Softw.*, 43(3):28:1–28:14, Jan. 2017.
  - [LM15] G. Lan and R. D. C. Monteiro. Iteration-Complexity of First-Order Augmented Lagrangian Methods for Convex Programming. *Math. Program.*, pages 1–37, 2015.
  - [LMWY13] J. Liu, P. Musialski, P. Wonka, and J. Ye. Tensor completion for estimating missing values in visual data. *IEEE Trans. Pattern Anal. Mach. Intell.*, 35(1):208–230, 2013.
  - [Lov03] L. Lovász. Semidefinite programs and combinatorial optimization. In *Recent advances in algorithms and combinatorics*, pages 137–194. Springer, 2003.
  - [LP66] E. Levitin and B. Polyak. Constrained minimization methods. *USSR Comput. Math. & Math. Phys.*, 6(5):1–50, 1966.
  - [LPZZ17] G. Lan, S. Pokutta, Y. Zhou, and D. Zink. Conditional accelerated lazy stochastic gradient descent. In *Proc. 34th Int. Conf. Machine Learning*, 2017.
  - [LTRJ17] F. Locatello, M. Tschannen, G. Rätsch, and M. Jaggi. Greedy algorithms for cone constrained optimization with convergence guarantees. In *Advances in Neural Information Processing Systems 30*, 2017.
  - [LYC18] K. Y. Levy, A. Yurtsever, and V. Cevher. Online adaptive methods, universality and acceleration. In *Advances in Neural Information Processing Systems*, 2018.
  - [LYFC19] F. Locatello, A. Yurtsever, O. Fercoq, and V. Cevher. Stochastic conditional gradient method for composite convex minimization. arXiv:1901.10348v2, 2019.

## Bibliography

---

- [LZ16] G. Lan and Y. Zhou. Conditional gradient sliding for convex optimization. *SIAM J. Optim.*, 26(2):1379–1409, 2016.
- [MHK18] A. Mokhtari, H. Hassani, and A. Karbasi. Stochastic conditional gradient methods: From convex minimization to submodular maximization. arXiv:1804.09554, 2018.
- [MNS15] E. Mossel, J. Neeman, and A. Sly. Consistency thresholds for the planted bisection model. In *Proceedings of the Forty-Seventh Annual ACM on Symposium on Theory of Computing*, pages 69–75. ACM, 2015.
- [Mot52] T. S. Motzkin. Copositive quadratic forms. Technical report, National Bureau of Standards Report No. 1818, 1952.
- [MVW17] D. G. Mixon, S. Villar, and R. Ward. Clustering subgaussian mixtures by semidefinite programming. *Information and Inference: A Journal of the IMA*, 6(4):389–415, 2017.
- [MZJ13] M. Mahdavi, L. Zhang, and R. Jin. Mixed optimization for smooth functions. In *Advances in Neural Information Processing Systems 26*, 2013.
- [Nes83] Y. Nesterov. A method of solving a convex programming problem with convergence rate  $O(1/k^2)$ . *Soviet Mathematics Doklady*, 27(2):372–376, 1983.
- [Nes05] Y. Nesterov. Smooth minimization of non-smooth functions. *Math. Program.*, 103:127–152, 2005.
- [Nes09] Y. Nesterov. Primal-dual subgradient methods for convex problems. *Math. Program.*, 120(1, Ser. B):221–259, 2009.
- [Nes15] Y. Nesterov. Universal gradient methods for convex optimization problems. *Math. Program.*, 152(1-2):381–404, 2015.
- [Nes17] Y. Nesterov. Complexity bounds for primal-dual methods minimizing the model of objective function. *Math. Program.*, 2017.
- [NLST17] L. M. Nguyen, J. Liu, K. Scheinberg, and M. Takác. SARAH: A novel method for machine learning problems using stochastic recursive gradient. In *Proc. 34th Int. Conf. Machine Learning*, 2017.
- [NN89] Y. Nesterov and A. Nemirovski. *Self-concordant functions and polynomial time methods in convex programming*. USSR Academy of Sciences, Central Economic & Mathematic Institute, 1989.
- [NN94] Y. Nesterov and A. Nemirovski. *Interior-point polynomial algorithms in convex programming*. Society for Industrial and Applied Mathematics, 1994.
- [NY83] A. Nemirovski and D. Yudin. *Problem Complexity and Method Efficiency in Optimization*. John Wiley & Sons, 1983.

- [OLY<sup>+</sup>16] G. Odor, Y.-H. Li, A. Yurtsever, Y.-P. Hsieh, Q. Tran-Dinh, M. El Halabi, and V. Cevher. Frank-Wolfe works for non-lipschitz continuous gradient objectives: Scalable poisson phase retrieval. In *41st IEEE Int. Conf. Acoustics, Speech & Signal Processing*, 2016.
- [Pat98] G. Pataki. On the rank of extreme matrices in semidefinite programs and the multiplicity of optimal eigenvalues. *Mathematics of Operations Research*, 23(2):338–358, 1998.
- [Pla05] J. C. Platt. FastMap, MetricMap, and Landmark MDS are all Nyström algorithms. In *Proc. 10th Int. Workshop Artificial Intelligence and Statistics*, 2005.
- [PW07] J. Peng and Y. Wei. Approximating K-means-type clustering via semidefinite programming. *SIAM J. Optim.*, 18(1):186–205, 2007.
- [QLX18] C. Qu, Y. Li, and H. Xu. Non-convex conditional gradient sliding. In *Proc. 35th Int. Conf. Machine Learning*, 2018.
- [Rag08] P. Raghavendra. Optimal algorithms and inapproximability results for every CSP? In *Proceedings of the 40th Annual ACM Symposium on Theory of Computing*, pages 245–254, 2008.
- [RDLS18] S. N. Ravi, T. Dinh, V. S. R. Lokhande, and V. Singh. Constrained deep learning using conditional gradient and applications in computer vision. arXiv:1803.06453v1, 2018.
- [Ren14] J. Renegar. Efficient first-order methods for linear programming and semidefinite programming. arXiv:1409.5832, 2014.
- [Roc70] R. T. Rockafellar. *Convex analysis*. Princeton Mathematical Series. Princeton University Press, 1970.
- [Roc76] R. T. Rockafellar. Monotone operators and the proximal point algorithm. *SIAM J. Control Optim.*, 14(5):877–898, 1976.
- [RSB12] N. L. Roux, M. Schmidt, and F. Bach. A stochastic gradient method with an exponential convergence rate for finite training sets. In *Advances in Neural Information Processing Systems 26*, 2012.
- [RSPS16] S. J. Reddi, S. Sra, B. Póczos, and A. Smola. Stochastic Frank-Wolfe methods for nonconvex optimization. In *54th Annual Allerton Conf. Communication, Control, and Computing*, 2016.
- [SEC<sup>+</sup>15] Y. Shechtman, Y. C. Eldar, O. Cohen, H. N. Chapman, J. Miao, and M. Segev. Phase retrieval with application to optical imaging. *IEEE Signal Process. Mag.*, 32(3):87–109, 2015.

## Bibliography

---

- [SFZ<sup>+</sup>19] Z. Shen, C. Fang, P. Zhao, J. Huang, and H. Qian. Complexities in projection-free stochastic non-convex minimization. In *Proc. 22nd Int. Conf. Artificial Intelligence and Statistics*, 2019.
- [SGL<sup>+</sup>19] Y. Sun, Y. Guo, C. Luo, J. A. Tropp, and M. Udell. Low-rank Tucker approximation of a tensor from streaming data. arXiv:1904.10951v1, 2019.
- [Sin11] A. Singer. Angular synchronization by eigenvectors and semidefinite programming. *Applied and computational harmonic analysis*, 30(1):20–36, 2011.
- [Sio58] M. Sion. On general minimax theorems. *Pacific J. Math.*, 8(1):171–176, 1958.
- [SRJ04] N. Srebro, J. Rennie, and T. Jaakkola. Maximum-margin matrix factorizations. In *Adv. Neural Information Processing Systems 17*, Vancouver, Dec. 2004.
- [SS11] A. Singer and Y. Shkolnisky. Three-dimensional structure determination from common lines in cryo-EM by eigenvectors and semidefinite programming. *SIAM J. Imaging Sciences*, 4(2):543–572, 2011.
- [SSGB07] L. Song, A. Smola, A. Gretton, and K. M. Borgwardt. A dependence maximization view of clustering. In *Proc. 20th Int. Conf. Machine Learning*, 2007.
- [ST14] R. Shefi and M. Teboulle. Rate of Convergence Analysis of Decomposition Methods Based on the Proximal Method of Multipliers for Convex Minimization. *SIAM J. Optim.*, 24(1):269–297, 2014.
- [TDC14] Q. Tran-Dinh and V. Cevher. Constrained convex minimization via model-based excessive gap. In *Advances in Neural Information Processing Systems 27*, 2014.
- [TDFC18] Q. Tran-Dinh, O. Fercoq, and V. Cevher. A smooth primal-dual optimization framework for nonsmooth composite convex minimization. *SIAM J. Optim.*, 28(1):96–134, 2018.
- [TLJJ06] B. Taskar, S. Lacoste-Julien, and M.I. Jordan. Structured prediction, dual extragradient and Bregman projections. *J. Mach. Learn. Res.*, 2006.
- [Tyg14] M. Tygert. Beta versions of Matlab routines for principal component analysis. Available at <http://tygert.com/software.html>, 2014.
- [TYUC17a] J. A. Tropp, A. Yurtsever, M. Udell, and V. Cevher. Fixed-rank approximation of a positive-semidefinite matrix from streaming data. In *Advances in Neural Information Processing Systems 30*, 2017.
- [TYUC17b] J. A. Tropp, A. Yurtsever, M. Udell, and V. Cevher. Practical sketching algorithms for low-rank matrix approximation. *SIAM J. Matrix Anal. Appl.*, 38(4):1454–1485, 2017.

- 
- [TYUC17c] J. A. Tropp, A. Yurtsever, M. Udell, and V. Cevher. Randomized single-view algorithms for low-rank matrix approximation. ACM Report 2017-01, Caltech, Jan. 2017.
  - [TYUC18] J. A. Tropp, A. Yurtsever, M. Udell, and V. Cevher. More practical sketching algorithms for low-rank matrix approximation. ACM Report 2018-01, Caltech, Oct. 2018.
  - [TYUC19] J. A. Tropp, A. Yurtsever, M. Udell, and V. Cevher. Streaming low-rank matrix approximation with an application to scientific simulation. arXiv:1902.08651, 2019.
  - [Upa16] J. Upadhyay. Fast and space-optimal low-rank factorization in the streaming model with application in differential privacy. arXiv:1604.01429v3, 2016.
  - [VNM44] J. Von Neumann and O. Morgenstern. *Theory of games and economic behavior*. Princeton press, 1944.
  - [WJZ<sup>+</sup>18] K. Wang, K. Ji, Y. Zhou, Y. Liang, and V. Tarokh. SpiderBoost: A class of faster variance-reduced algorithms for nonconvex optimization. arXiv:1810.10690, 2018.
  - [WLRT08] F. Woolfe, E. Liberty, V. Rokhlin, and M. Tygert. A fast randomized algorithm for the approximation of matrices. *Appl. Comput. Harmon. Anal.*, 25(3):335–366, 2008.
  - [WS00] C. K. I. Williams and M. Seeger. Using the Nyström method to speed up kernel machines. In *Advances in Neural Information Processing Systems 13*, 2000.
  - [WW18] I. Waldspurger and A. Waters. Rank optimality for the Burer-Monteiro factorization. arXiv:1812.03046v1, 2018.
  - [Xu17] H.-K. Xu. Convergence analysis of the Frank–Wolfe algorithm and its generalization in Banach spaces. arXiv:1710.07367v1, 2017.
  - [YFC19] A. Yurtsever, O. Fercoq, and V. Cevher. A conditional gradient-based augmented lagrangian framework. In *Proc. 36th Int. Conf. Machine Learning*, 2019.
  - [YFLC18] A. Yurtsever, O. Fercoq, F. Locatello, and V. Cevher. A conditional gradient framework for composite convex minimization with applications to semidefinite programming. In *Proc. 35th Int. Conf. Machine Learning*, 2018.
  - [YHC15] A. Yurtsever, Y.-P. Hsieh, and V. Cevher. Scalable convex methods for phase retrieval. In *6th IEEE Intl. Workshop on Computational Advances in Multi-Sensor Adaptive Processing*, 2015.

## Bibliography

---

- [YM10] A. Yoshise and Y. Matsukawa. On optimization over the doubly nonnegative cone. In *IEEE International Symposium on Computer-Aided Control System Design*, 2010.
- [YSC19] A. Yurtsever, S. Sra, and V. Cevher. Conditional gradient methods via stochastic path-integrated differential estimator. In *Proc. 36th Int. Conf. Machine Learning*, 2019.
- [YST15] L. Yang, D. Sun, and K.-C. Toh. SDPNAL+: A majorized semismooth Newton–CG augmented Lagrangian method for semidefinite programming with nonnegative constraints. *Math. Program. Comput.*, 7(3):331–366, 2015.
- [YTDC15a] A. Yurtsever, Q. Tran-Dinh, and V. Cevher. A universal primal-dual convex optimization framework. In *Advances in Neural Information Processing Systems 28*, 2015.
- [YTDC15b] A. Yurtsever, Q. Tran-Dinh, and V. Cevher. Universal primal-dual proximal-gradient methods. arXiv:1502.03123v1, 2015.
- [YUTC17] A. Yurtsever, M. Udell, J.A. Tropp, and V. Cevher. Sketchy decisions: Convex low-rank matrix optimization with optimal storage. In *Proc. 20th Int. Conf. Artificial Intelligence and Statistics*, 2017.
- [YVC16] A. Yurtsever, B. C. Vu, and V. Cevher. Stochastic three-composite convex minimization. In *Advances in Neural Information Processing Systems 29*, 2016.
- [YZS14] Y. Yu, X. Zhang, and D. Schuurmans. Generalized conditional gradient for sparse estimation. arXiv:1410.4828v1, 2014.
- [ZCM<sup>+</sup>19] M. Zhang, L. Chen, A. Mokhtari, H. Hassani, and A. Karbasi. Quantized Frank-Wolfe: Communication-efficient distributed optimization. arXiv:1902.06332, 2019.
- [ZS18] W.-J. Zeng and H. C. So. Outlier-robust matrix completion via  $\ell_p$ -minimization. *IEEE Trans. on Sig. Process*, 66(5):1125–1140, 2018.
- [ZXG18] D. Zhou, P. Xu, and Q. Gu. Stochastic nested variance reduction for nonconvex optimization. In *Advances in Neural Information Processing Systems 32*, 2018.

Alp Yurtsever  
EPFL STI IEL LIONS  
ELD 244 (Bâtiment ELD)  
Station 11, 1015 Lausanne  
Switzerland

Website: [alpyurtsever.github.io](https://alpyurtsever.github.io)  
E-mail: [alpyurtsever@gmail.com](mailto:alpyurtsever@gmail.com)

## Education

---

Sep 2013 – June 2019	<b>PhD</b> École Polytechnique Fédérale de Lausanne, Lausanne Computer and Communication Sciences (EDIC) Laboratory of Information and Inference Systems (LIONS) Thesis topic : Scalable convex optimization methods for semidefinite programming Thesis supervisor : Prof. Volkan Cevher
Sep 2018 – Dec 2018	<b>Visitor</b> Massachusetts Institute of Technology, Cambridge Institute for Data, Systems and Society (IDSS) Laboratory for Information and Decision Systems (LIDS) Hosting faculty : Prof. Suvrit Sra
Sep 2009 – June 2013	<b>BSc</b> Middle East Technical University, Ankara Electrical and Electronics Engineering Physics (Double Major)

## Awards and Fellowships

---

- SNF Early PostDoc.Mobility  
*The fellowship includes a grant towards living costs, awarded by Swiss National Science Foundation for 18 months to successful early-career postdocs who wish to enhance their scientific profile by working at a research institution abroad.*
- IEEE CAMSAP 2015 student paper award  
*Three best student papers of the workshop among the papers by graduate or undergraduate students are identified by the committee, in Cancun, Mexico, in December 2015.*
- EDIC fellowship *One year fellowship program that is granted for the promising applicants with strong academic record, by the doctoral school of IC department at EPFL.*

## Service to the community

---

### Reviewing activities

OXFORD Academic: IMA Journal of Numerical Analysis (IMAJNA)  
IEEE Transactions on Signal and Information Processing over Networks (SIPN)  
IEEE Transactions on Signal Processing (T-SP)  
Journal of Selected Topics in Signal Processing (J-STSP)  
International Journal of Control, Automation and Systems (IJCAS)  
Operations Research Letters  
International Conference on Machine Learning (ICML)  
Neural Information Processing Systems (NeurIPS)

