

©Copyright 2025

Ronak Mehta

Statistical Learning from Shifting, Indirect, or Unseen Data: Efficient Algorithms and Theoretical Guarantees

Ronak Mehta

A dissertation
submitted in partial fulfillment of the
requirements for the degree of

Doctor of Philosophy

University of Washington

2025

Reading Committee:

Zaid Harchaoui, Chair

Armeen Taeb

Alex Luedtke

Program Authorized to Offer Degree:

Statistics

University of Washington

Abstract

Statistical Learning from Shifting, Indirect, or Unseen Data: Efficient Algorithms and Theoretical Guarantees

Ronak Mehta

Chair of the Supervisory Committee:

Zaid Harchaoui

Department of Statistics

A fascinating phenomenon underlying statistical machine learning and artificial intelligence is “out-of-distribution” (or OOD) generalization. Data can (and in some settings, must) be used to draw inferences regarding probability distributions other than the one from which they were sampled. Understanding this mystery gives promise to statistical analyses that exhibit a degree of universality, such as clinical trials whose conclusions reflect many subpopulations or pre-defined image/text encodings that can be used to solve many classification tasks simultaneously. This dissertation tackles the theoretical and algorithmic challenges of designing methods that exhibit these modern notions of generalization.

Chapter 2 studies a learning framework called distributionally robust optimization (DRO), which promotes OOD by training models to optimize the worst-case expected loss achievable within a collection of possible training distributions. These maximum-type objectives present challenges for designing stochastic learning algorithms, as unbiased estimates of the gradient are not easily computed. We design an estimator equipped with a progressive bias (and variance) reduction scheme, for which the resulting algorithm is shown to have a linear convergence guarantee. Although our optimization results apply more generally to DRO problems, we focus attention on a subclass of objectives called spectral risk measures, which have appealing statistical and computational properties previously unexplored in machine

learning. We provide theoretical and practical guidance on selecting the various problem parameters, such as the collection of distributions over which to maximize. Finally, we present (among others) extensions to group DRO, a popular extension of the framework amenable to training neural network models.

Chapter 3 takes insights from the DRO application and pursues stochastic algorithms for a more general class of optimization problems, dubbed semilinear min-max problems. These objectives interpolate between the well-understood class of bilinear and relatively less-understood nonbilinear min-max problems, and have applications to problem classes such as convex minimization with functional constraints as special cases. We present the first complexity guarantees for this problem class, using a randomized algorithm with components inspired by the simulation literature (such as adaptive sampling of new data and adaptive averaging of historical data). We prove convergence guarantees in both convex and strongly convex settings with a fine-grained dependence on individual problem constants. The results yield complexity improvements in even specific cases, such as bilinearly coupled problems. We also provide a lower complexity bound on the performance of deterministic algorithms applied to the semilinear problem class.

Chapter 4 shifts focus from the implementation of large-scale learning algorithms to their output. We investigate predictive models that learn via a pre-training procedure with unlabeled data and can then make predictions for downstream classification tasks (without having seen any directly labeled training data from that task). This capability, known as zero-shot prediction, is made possible by three ingredients: 1) massive, carefully curated pre-training datasets, 2) “self-supervised” labels that allow models to learn universal features of structured data (e.g., images/text), and 3) the translation of downstream data into the format seen during pre-training using a technique called prompting. We analyze all three ingredients theoretically by establishing both the sample complexity and the limits of prompting in terms of simple distributional conditions. Inspired by this theory, we explore

variants on the pre-training objective and prompting strategies that show practical benefits such as improved zero-shot classification accuracy.

TABLE OF CONTENTS

	Page
List of Figures	iii
Chapter 1: Introduction	1
1.1 Context and Motivation	3
1.2 Technical Overview	6
1.3 Contributions and Outline	11
Chapter 2: Learning under Distribution Shift with Likelihood Ratios	16
2.1 Introduction	16
2.2 Preliminaries	18
2.3 Properties of Spectral Risk Measures	20
2.4 Smoothing Maximum-Type Objectives	33
2.5 Stochastic Optimization with Bias and Variance Reduction	37
2.6 Convergence Analysis	43
2.7 Uncertainty Sets and Shift Costs	52
2.8 Comparison to Broader Literature	64
2.9 Experiments	75
2.10 Incorporating Group Structure	84
2.11 Possible Extensions	103
2.12 Perspectives & Future Work	113
Chapter 3: Algorithmic Extensions of Distributionally Robust Optimization	115
3.1 Introduction	115
3.2 Preliminaries	120
3.3 Method and General Analysis Template	127
3.4 Stochastic Algorithms for General Objectives	141
3.5 An Algorithm for Dual-Separable Problems	162

3.6	Discussion & Comparisons	174
3.7	Possible Extensions	185
3.8	Perspectives & Future Work	201
Chapter 4:	Generalization Capabilities of Zero-Shot Prediction	206
4.1	Introduction	206
4.2	Preliminaries	208
4.3	Non-Asymptotic Analysis of Variance Reduction	212
4.4	Theoretical Limits of Zero-Shot Prediction	218
4.5	Experiments	229
4.6	Possible Extensions	234
4.7	Perspectives & Future Work	240
Chapter 5:	Conclusion	242
Appendix A:	Appendix to Chapter 2	268
A.1	Technical Background	268
A.2	Convergence Analysis	273
A.3	Implementation Details	290
A.4	Experimental Details	297
Appendix B:	Appendix to Chapter 4	304
B.1	Linear Operators and Variance Reduction	304
B.2	Information Projections	310
B.3	Statistical Analysis of Balancing Estimators	319
B.4	Zero-Shot Prediction	353
B.5	Experimental Details	355

LIST OF FIGURES

Figure Number		Page
1.1	Illustration of Out-of-Distribution Generalization. While traditional statistical learning paradigms would deal only with the data-generating distribution P during evaluation/deployment, this dissertation considers cases in which the evaluation distribution Q may be different. Often, the change/shift from P to Q exhibits additional structure (discussed, e.g., in Section 2.8). Techniques for handling this shift fall into alternative learning algorithms (Chapter 2 and Chapter 3), usage of auxiliary/side information Section 4.3, or methods of querying pre-trained models (Section 4.4).	3
1.2	Illustration of Learning Pipeline. In the illustration above, a predictive model is evaluated on n training examples to produce the histogram of errors (or loss distribution). This loss function is summarized by a single quantity, perhaps by integrating over the empirical distribution P_n , after which this quantity is minimized by invoking an optimization algorithm.	4
1.3	Foundation Model Pre-Training versus Direct Supervision. The purple encoder is optimized to produce an informative representation of input images. In the traditional supervision framework (top), labeled training data is provided so that the encoder produces representations that are informative for the particular task of predicting the given label. Labeled data, often produced using high-quality annotators, is relatively scarce compared to the number of images that are unlabeled. Because unlabeled images may still have accompanying captions, in modern frameworks such as self-supervision, the image encoder may be (pre-)trained to be predictive of its caption (and vice versa).	5
1.4	Linear-Nonlinearly Coupled Objective. Illustration of the objective (1.6), which is possibly nonlinear in θ (the model parameters) but linear in the data weights q	9
1.5	Prompting of Foundation Models. Illustration of the indirect predictor (1.7), which relies on image and text encoders pre-trained via self-supervision (see Figure 1.3)	10

2.1	Example Spectra and Induced Random Variables. Visualization of the superquantile and extremile spectra and the induced distribution function G_s of the likelihood ratio.	23
2.2	Illustration of Spectral Risk and Quantile Functions. The relationship between the order statistics, quantile functions, and the discretized spectrum is shown. Top: Empirical CDF F_n and quantile function F_n^{-1} of X_1, \dots, X_4 . Bottom: Continuous spectra $s(t)$ and their discretization $(\sigma_1, \dots, \sigma_5)$ for various risk measures.	32
2.3	Idealized Visualization of Bias- and Variance-Reduced Algorithms. The Prospect [Mehta et al., 2024b], LSVRG [Mehta et al., 2023], and Loopless LSVRG (Example 2.5.2) algorithms are shown. Left: Expected trajectory over algorithmic randomness. This displays the gradient bias (as compared to full batch gradient descent of each method). LSVRG, unlike the updates in lines 14 and 14 for Algorithm 1, operators in epochs. Right: Observed trajectory of either Prospect/Loopless LSVRG. The variance reduction is interpreted as a control variate correction applied to the initial stochastic gradient estimate.	39
2.4	Regression Benchmarks. The y -axis measures suboptimality as given by (2.49), while the x -axis measures the number of calls to the function value/-gradient oracle divided by n (i.e., passes through the training set). Rows indicate different SRM objectives while columns indicate datasets.	77
2.5	Fairness Benchmarks. Top: Training curves for optimizers on the CVaR and extremile for diabetes (left) and CVaR and extremile for acsincome (right). Bottom: Statistical parity scores for the two classification objectives on diabetes (left) and regression objectives on acsincome . Smaller values indicate better performance for all metrics.	79
2.6	Distribution Shift Benchmarks. Top: Training curves and worst group misclassification error on amazon test. Bottom: Training curves and median group misclassification error on the iwildcam test set. Smaller values indicate better performance for all metrics.	81
2.7	Wall Time and Shift Cost Relationship. Top: Wall time against shift cost for reaching convergence (squared gradient norm less than 10^{-3}) for Prospect. Bottom: Wall time against shift cost for reaching convergence (squared gradient norm less than 10^{-3}) for LSVRG.	82
2.8	Optimal Dual Solutions and Shift Cost Relationship. Top: Visualization of continuous spectra for the superquantile and extremile. Bottom: Sorted optimal dual solution q^* for different values of ν , meant to compare to the superquantile spectrum in the top left panel.	83

2.9	Illustration of Group-Wise Error Evaluation. The top panel depicts a hypothetical probability distribution, where the sample space is partitioned into three groups. When applying a statistical prediction model h to this distribution, the bottom left and right panels show non-uniform and uniform performance conditional on each group. On the right panel specifically, the arrows indicate a slight increase in the average loss across all examples, but a more dramatic decrease in the worst-case group-wise error.	84
2.10	Amazon Reviews Data. The left panel shows the group-wise proportions in the train and test sets, respectively. The right panel shows the average and worst-case group-wise test accuracy. The proposed objectives are the 1.5 or 2-extremile. Even though the group proportions have a minor shift in most product categories, a moderate worst-case accuracy increase is observed for the distributionally robust variants.	101
3.1	Generalized Linear Models as Matrix Games. Visualization to accompany Example 3.1.1, where \mathbf{A} denotes the design matrix, \mathbf{x} is parameter vector, and \mathbf{y} multiplier to adjust the predicted scores \mathbf{Ax}	117
3.2	Moreau Envelope of Bregman Divergences on the Unit Simplex. Visualization of the objective of the proximal operator in the case of the standard ℓ_2 -norm-squared Bregman divergence (left) and the negative entropy-base Bregman divergence (right).	124
3.3	Modified Proximal Step with Historical Regularization. Geometric illustration of the historical regularization penalty applied in the modified primal update (3.11).	130
3.4	Adaptive Sampling for Gradient Estimation Geometric illustration of adaptive sampling for computing the mean of vectors $\bar{\mathbf{g}}_1, \dots, \bar{\mathbf{g}}_N$	135
3.5	Experimental Evaluation of Online Accuracy Certificates. In all plots, the x -axis refers to the iteration count, which may differ between datasets. Each line represents the gradient norm (3.103), certificate (3.101), and the primal-dual gap (3.98).	190
4.1	Data Balancing. Nonlinear and linear operators associated with each iteration of (4.2). Left: Visualization of the exact iterations of (4.2) in the space of probability measures. The purple set contains joint distributions with \mathcal{X} -marginal equal to P_X , whereas the golden set contains joint distributions with \mathcal{Z} -marginal equal to P_Z . Right: Visualization of $\mathbf{L}^2(P)$, the operators defining (4.9), and the singular values given in (4.11).	213

4.2	Graphical Models of Self-Supervised Prediction Paths. Each directed graphical model corresponds to the data types and dependence structures for various SSL pre-training approaches. The variable C represents an unobserved context that determines the observed data-generating distribution. Dotted lines indicate the possibility of presence or absence of the arrow. Methods that are compatible with FSL learn the label Y as a latent variable in the process of solving the pretext task. Methods compatible with ZSP may learn the relationship between X and Z directly, whereas the relationship between Z and Y is estimated via prompting.	219
4.3	Illustration of Prompting Approaches. A hypothetical distribution of embeddings $\beta(Z)$ parametrized by two classes (“cat” and “dog”). Three prompting strategies (template-based, class-conditional, and idealized) are shown with example text and resulting embeddings in \mathbb{R}^d . Prompt bias is the distance of the average of the circular points to the square target points.	221
4.4	Balancing as Data Curation. Depiction of balancing and data curation on ImageNet-Captions dataset, in which \mathcal{X} represents image-caption pairs and \mathcal{Y} represents keywords. Left: Observed marginal $P_{n,Z}$ (orange) and P_Z (blue), which are sorted by order of increasing probability. Right: Zero-shot evaluation of an embedding model trained using the standard CLIP loss on the original versus the balanced training set.	230
4.5	Results: Unbiased Prompting. Pre-trained models are varied along the rows, and sub-tasks (subsets of 50 ImageNet-1k classes) are varied along columns. In all plots, the x -axis denotes the number of prompts sampled for each class embedding, and the y -axis denotes top- k zero-shot classification accuracy. Error bars indicate standard deviations across 10 seeds for prompt sampling. In this setting $P_{Y,Z} = P_{Y,Z}^T$	231
4.6	Results: Class-Conditional Prompting. Pre-trained models are varied along the rows, and evaluation datasets are varied along columns. In all plots, the x -axis denotes the number of prompts sampled for each class embedding, and the y -axis denotes top- k zero-shot classification accuracy. Error bars indicate standard deviations across 10 seeds for prompt sampling.	232
A.1	Harder Risk Parameter Settings. Each row represents a different “hard” variant of the superquantile, extremile, and ESRM spectra. Columns represent different datasets. Suboptimality (2.49) is measured on the y -axis while the x -axis measures the total number of gradient evaluations made divided by n , i.e., the number of passes through the training set.	301

A.2	No Shift Cost Settings. Each row represents a different spectral risk objective with $\nu = 0$ (instead of $\nu = 1$) while each column represents a different datasets. Suboptimality (2.49) is measured on the y -axis while the x -axis measures the total number of gradient evaluations made divided by n , i.e., the number of passes through the training set.	302
A.3	Reduced ℓ_2-regularization settings ($\mu = 1/(10n)$). Each row represents a different spectral risk objective with $\mu = 1/(10n)$ (instead of $\mu = 1/n$) while each column represents a different dataset. Suboptimality (2.49) is measured on the y -axis while the x -axis measures the total number of gradient evaluations made divided by n , i.e., the number of passes through the training set.	303

ACKNOWLEDGMENTS

Throughout my Ph.D., I have been surrounded by the best of the best. I often describe the experience as traveling from dojo to dojo, each time being exposed to a new field and donning the title of “beginner” again and again. Be it my mentors, collaborators, or friends, I could not have asked for a better community in which to learn and grow.

From my research advisor, Zaid Harchaoui, I observed and admired a relentless pursuit of excellence. Thank you for believing that I, too, could join this pursuit. I am now a person who finds joy, perhaps even comfort, in chasing the impossible. Zaid is cultivated in a large number of fields, and is possibly the only one in the world who can use “Rao-Blackwellization” and “large language models” in the same sentence. It is this eclectic blend of classical and cutting-edge disciplines that I originally gravitated toward. Taking pride in one’s work, never cutting corners, and helping others unconditionally are only a few of the lifelong lessons I will carry with me from Zaid’s guidance.

I experienced my first deep dive into the field of optimization with Vincent Roulet. Although we both share an interest in the theoretical aspects of the field, it was from him that I learned to be an experimentalist. From using mathematics to debug code to maximizing information gain with each new experiment, the art and science of numerical methods were extremely challenging for me when I started this journey. Thank you, Vincent, for patiently fostering what I now consider to be my most valuable skill. Beyond this, I also appreciate your perspective on the personal journey, which constitutes the hardest part of a Ph.D.

Having done research in optimization for some time, while I never dared to call myself an expert, I did not realize how far I was until I met Jelena Diakonikolas. When writing our journal paper together, I witnessed her willingness to leave no technical stone unturned,

while still making complicated ideas crystal clear for the reader. In trying to emulate these qualities during late nights in Madison, I finally started to develop my own intuition (and rejoiced at the sight of negative quadratic terms). Thank you, Jelena, for your patience and inspiration.

By working with Lang Liu, my understanding of statistical theory was completely transformed. The two skills of his that I admired the most are the ability to formalize the chaos of applied AI into compelling theoretical problems and the mathematical creativity to invent new tools to solve them. While I attempt to display the first skill in this dissertation, I am not sure I (or anyone, for that matter) can perform the second skill quite like Lang. Thank you, Lang, for spending your personal time even after your graduation to invest in my growth.

At the expense of sounding cliché, as I reflect on my papers over the years, the most successful ones were those that told the best stories. If the technical storytelling of this thesis is to be appreciated, it is really Krishna Pillutla who deserves the credit. For every paper, presentation, rebuttal, or poster, I always struggled with the subtleties and constant decision-making inherent to scientific communication. How should I invite my audience into my problem setting without overwhelming them? Should I surrender to a reviewer's interpretation of my work, or should I firmly push back on their comments? Luckily, life is made easier by simply imagining what Krishna would do in these situations. I must also thank Krishna for staying up until the morning hours during my first paper submissions.

The intellectual community at the University of Washington was a gift that kept on giving. I thank John Thickstun and Alec Greaves-Tunnell for inspiring me during my first year. I am only now realizing the magnitude of the time investments they made during their last year, such as John's generative models course or Alec's contribution to the machine learning course we taught with Zaid. These experiences drew me into the research that became this dissertation. I would also like to recognize and thank my committee members, Alex Luedtke

and Armeen Taeb, as well as my co-authors and collaborators: Azadeh Yazdan-Shahmorad, Noah Stannis, Eric Shea-Brown, Ali Shojaie, Soumik Pal, Ludwig Schmidt, Briana Abrahms, Kasim Rafiq, and Medha Agarwal. To my research peers, Alex Bank, Konstantin Golobokov, Mateus Piovezan Otto, and Facheng Yu, thank you for inspiring me with your enthusiasm and determination. Your futures are beyond bright, and I am looking forward to seeing what they bring.

Essential to the research in this document were two funding initiatives from the National Science Foundation: the Institute for Foundations of Data Science (IFDS) and the Institute for the Foundations of Machine Learning (IFML). Much of my work, which emphasizes theoretical understanding, would not be possible if it were not for these foundations-focused initiatives.

I express my gratitude for the administrative and technical staff of our department: Ellen Reynolds, Vickie Graybeal, Tracy Pham, Kristine Chan, Veronica Bae, and Asa Sourdiffé. I never took for granted your professionalism, efficiency, and care for the students, and have been saved by your emails more times than I should admit.

Looking a bit further into the past, this journey would not have even begun if not for my professors at Johns Hopkins University. I knew that applied mathematics and statistics would become a lifelong endeavor after taking Avanti Athreya's course. Donniell Fishkind took a chance on me and orchestrated my first-ever internship in data science. I inherited a sincere love of research from working with my first advisor, Joshua Vogelstein. Finally, from Carey Priebe, I learned that there is always an exciting lesson to be learned from even the simplest of statistical problems (independent Bernoulli trials would be his case in point).

On a more personal note, I was very fortunate to live in a home that was always filled with loving friends (whether they officially lived there or not). Certain periods of my Ph.D. could not have been traversed without a certain degree of obsession. Yet, I never had to earn this obsession at the price of solitude. Instead, my goals, lifestyle, humor, successes, and

failures were shared with others. Medha, Saksham, Jeff, and Brianna, thank you for your companionship. There is no me without you.

I am grateful for the friends made in my cohort. Shreya and James, thank you for accepting me as I was in the earliest days of the program. Jillian, thank you for teaching me what it means to be an adult and to stand up for myself. Aparna, thank you for your thoughtfulness and reliability over the years. Alana, thank you for sharing with me the interests that no one else does. Finally, I thank Kayla for improving my confidence as a teacher and mentor for others. To you all, I look forward to the next time we have the chance to meet, wherever and whenever it may be.

My final acknowledgement is dedicated to my family. Firstly, to my cousins Aayush, Aditya, Abhishek, Akash, and Priya for traveling across the globe to attend my defense. Secondly, to my brother Kunal for not only supporting me, but also managing our trips, business, and many other responsibilities when I was too overwhelmed to do so myself. I will happily pay back this effort with interest in the coming years. Lastly, to my parents, Deepak and Nisha Mehta. They often like to quip that there must have been a mix-up at the hospital, as they wonder how this Ph.D. kid came from them. I, on the other hand, never saw a reason to wonder. Being kind, hardworking, and loyal like my father, and being disciplined and thoughtful like my mother, is all I ever really tried to do.

DEDICATION

To my parents

Chapter 1

INTRODUCTION

There has never been a more exciting time to work at the intersection of statistical machine learning (ML) and artificial intelligence (AI), as developments in modern AI challenge and redefine the conventional view of “generalization” in statistics. To understand this phenomenon, consider the example of a scientist pursuing knowledge discovery or an ML engineer building a new product. While their goals differ, their processes bear significant similarity; each will frame their goal mathematically using a quantitative *parameter* within a statistical model, and may consider collecting data for the express purpose of learning or inferring this parameter. This reflects Ronald Fisher’s famous Rothamsted experiment in agriculture [Parolini, 2015] or the introduction of CIFAR-10 to study compositional features in image classification [Krizhevsky, 2009]. In both cases, the practitioner seeks the ability to *generalize*, or to use their sample or training set, respectively, to achieve an understanding of possibly unseen data. While this *modus operandi* has been tremendously successful for designing the theory and methods of contemporary statistics, it is increasingly pressing for scientists, engineers, and researchers to use data to learn about populations (or data-generating distributions) *other than* the one from which they were drawn.

Why would such a need for *out-of-distribution (OOD) generalization* arise?

- **Accountability:** Randomized control trials (RCTs) have been the gold standard experimental design for decades. However, when certain interventions are unethical (e.g., exposure to carcinogens) or conducting such a trial is prohibitively expensive, observational data may be used as an alternative option. Managing the discrepancy between interventional and observational data led to the development of the propensity score [Rosenbaum and Rubin, 1983], a mainstay of causal inference.

- **Safety:** AI models are now deployed in critical domain applications such as energy planning [Guigues and Sagastizábal, 2013], materials engineering [Yeh, 2006], and financial regulation [He et al., 2022]. When the deployment environment differs drastically from the training environment, brittle models may lead to catastrophic outcomes such as misdiagnoses, bankruptcies, and mechanical failures. A natural question is whether such shifts in distribution can be simulated or accounted for during the training phase itself.
- **Fairness:** A well-documented phenomenon is the tendency of AI models, such as facial recognition systems, to perform well on majority subgroups in the evaluation data (such as male or lighter-skinned individuals) and exhibit social biases on minority subgroups [Buolamwini and Gebru, 2018]. It is of clear interest whether models can be developed to promote equitable behavior, even if minority subgroups are scarce in the training population, which constitutes one popular notion of algorithmic fairness. As we elaborate on in Chapter 2, fairness by this definition can also be framed as a distribution shift problem, in which minority subgroups appear with higher probability in the evaluation data.
- **Efficiency:** Finally, high-performance models such as large language models (LLMs) may be trained on trillions of tokens and thousands of GPUs in order to be viable for applications such as AI chat assistants [Hoffmann et al., 2022]. Such a magnitude of resources is not accessible to the average individual, motivating research into the reuse of institutionally-trained models with little to no training data from the target distribution provided by the scientist. For this setting in particular, not only may the evaluation data change, but the task itself may change as well, introducing additional complexity.

We argue that in statistics and machine learning, out-of-distribution generalization is the rule rather than the exception. This dissertation contributes to the theory and methods

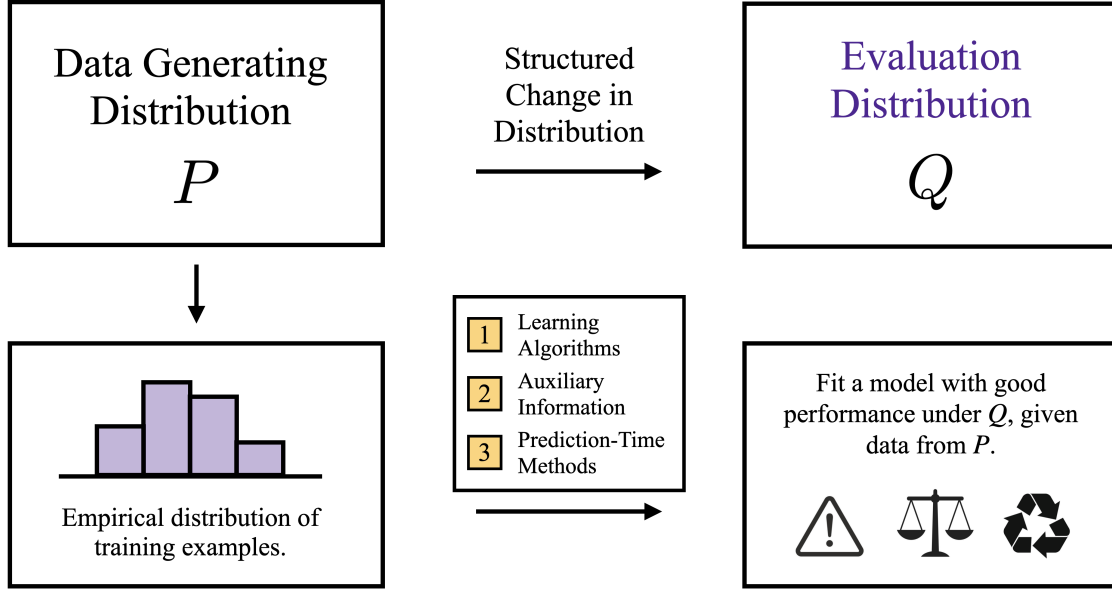


Figure 1.1: **Illustration of Out-of-Distribution Generalization.** While traditional statistical learning paradigms would deal only with the data-generating distribution P during evaluation/deployment, this dissertation considers cases in which the evaluation distribution Q may be different. Often, the change/shift from P to Q exhibits additional structure (discussed, e.g., in Section 2.8). Techniques for handling this shift fall into alternative learning algorithms (Chapter 2 and Chapter 3), usage of auxiliary/side information Section 4.3, or methods of querying pre-trained models (Section 4.4).

of OOD generalization by 1) retrospectively analyzing existing learning algorithms and 2) prospectively developing new ones (see Figure 1.1 for an illustration). In the remainder of this chapter, Section 1.1 introduces existing lines of research at a high-level, Section 1.2 collects technical details necessary to understand the main contributions, and Section 1.3 provides an overview of these contributions and an outline of the document.

1.1 Context and Motivation

1.1.1 Alternative Learning Objectives

Parameters are fit to data by optimizing some performance measure called a *loss function*, and while many losses (such as the squared, hinge, or cross-entropy loss) have been explored

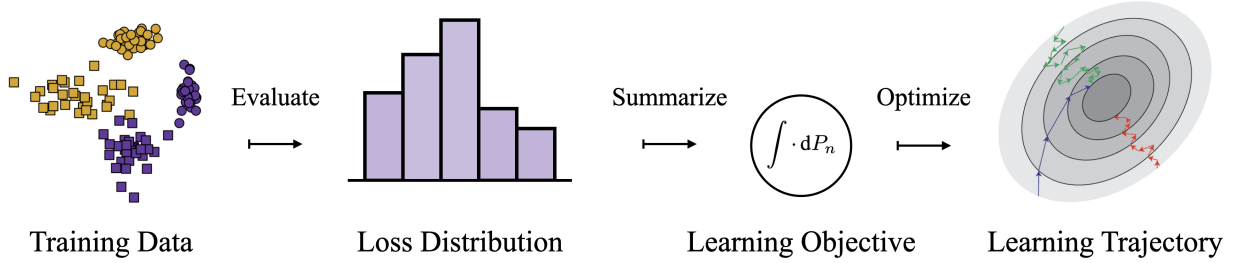


Figure 1.2: **Illustration of Learning Pipeline.** In the illustration above, a predictive model is evaluated on n training examples to produce the histogram of errors (or loss distribution). This loss function is summarized by a single quantity, perhaps by integrating over the empirical distribution P_n , after which this quantity is minimized by invoking an optimization algorithm.

in the literature, one aspect of these objectives has remained consistent: that the loss is averaged over the training data to produce the final summary. Upon reflection, the usage of the simple average is a choice in and of itself, and we may consider other ways of aggregating the loss into a univariate quantity to be optimized (see Figure 1.2). In applications such as finance, alternative *risk measures*, or summaries of the loss distribution, have been explored for decades [He et al., 2022]. Two highly related questions of increasing interest are whether these risk measures can be used to learn models that perform well on (1) the tails of the loss distribution and (2) on other data-generating distributions without including any additional training data.

1.1.2 Data Selection over Model Selection

Statistics pedagogy often describes data in an oracle framework, that is, the analysis starts after a number of independent and identically distributed (i.i.d.) data points are drawn from a probability distribution and supplied to the researcher. In both scientific and industrial applications, the researcher may often have the agency to both influence the data-generating mechanism *and* conduct the analysis simultaneously. While organizations such as Google, Meta, and OpenAI have access to essentially the entire Internet, it is now recognized that

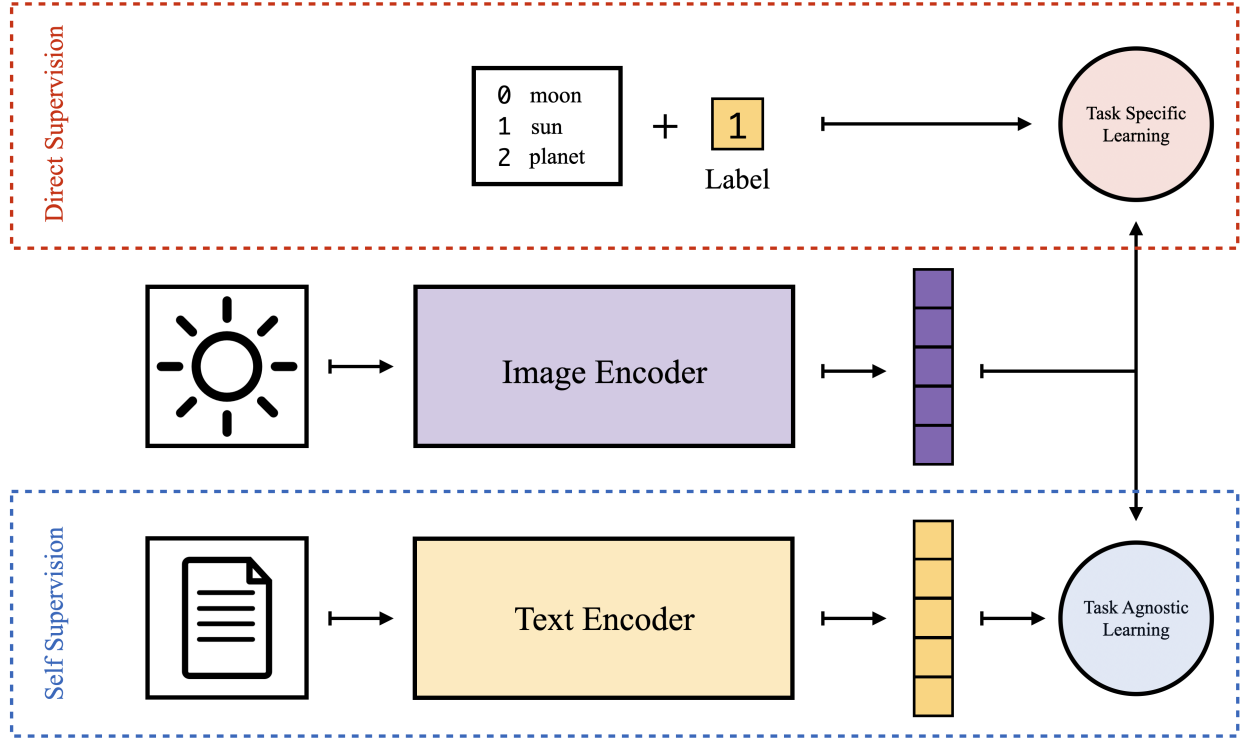


Figure 1.3: **Foundation Model Pre-Training versus Direct Supervision.** The purple encoder is optimized to produce an informative representation of input images. In the traditional supervision framework (top), labeled training data is provided so that the encoder produces representations that are informative for the particular task of predicting the given label. Labeled data, often produced using high-quality annotators, is relatively scarce compared to the number of images that are unlabeled. Because unlabeled images may still have accompanying captions, in modern frameworks such as self-supervision, the image encoder may be (pre-)trained to be predictive of its caption (and vice versa).

a large quantity of data does not produce models that generalize well if much of it is low-quality [Gadre et al., 2023, Li et al., 2024a]. In other words, one may also ask which *data* promotes OOD generalization in addition to which *models*.

1.1.3 Universal Representations of Structured Data

An emerging phenomenon is the advent of *foundation models*, or pre-trained neural networks whose internal layers can be used as fixed or tunable feature mappings for structured data

such as images, text, or audio (see [Devlin et al. \[2019b\]](#), [Chen et al. \[2020\]](#), [Bardes et al. \[2022\]](#), [Oquab et al. \[2024\]](#) and references therein). These reusable feature representations have had a major impact in settings such as data-scarce image or text classification, as they may simplify these problems to simple linear or logistic regression on pre-trained features. Even more surprisingly, they can be used for zero-shot prediction, or generating classifiers without *any* additional training data. This capability has evolved the notion of generalization even further: instead of generalizing to different distributions over the same sample space, these models can now generalize to different input and output spaces entirely. The pre-training process for an image embedding foundation model is shown in [Figure 1.3](#) and compared to the classical notion of (direct) supervised learning. We describe the usage of these models for zero-shot prediction in the upcoming [Section 1.2.3](#). While this technique has forged ahead from the applied perspective [[Radford et al., 2021](#), [Pratt et al., 2023](#), [Xu et al., 2024](#)], we also hope to achieve an improved theoretical understanding of it in [Chapter 4](#) of this dissertation.

1.2 Technical Overview

This dissertation investigates several statistical and computational aspects of these recent perspectives on generalization, from linear models for tabular data to foundation models for images and natural language. Our study is unified by developing and understanding stochastic algorithms for large-scale optimization problems that invariably arise across settings.

1.2.1 Empirical Risk Minimization

We center ourselves in the eminent framework of *empirical risk minimization (ERM)*. Consider a common probability space $(\Omega, \mathcal{F}, \mathbb{P})$, where (Ω, \mathcal{F}) is a measurable space and \mathbb{P} is a probability measure. Consider a statistical model (Ξ, \mathcal{P}) , where \mathcal{P} contains probability measures over Ξ , the space of observable data. Then, let $\xi : \Omega \rightarrow \Xi$ be a random variable that is governed by an unknown probability measure $P \in \mathcal{P}$. Given a parameter space Θ

and *risk functional* $\mathcal{R} : \Theta \times \mathcal{P} \rightarrow \mathbb{R}$, we wish to determine a parameter of the form

$$\boldsymbol{\theta}_0 \in \arg \min_{\boldsymbol{\theta} \in \Theta} \mathcal{R}(\boldsymbol{\theta}, P), \quad (1.1)$$

assuming such a minimizer exists. Although P is unknown, we assume access to i.i.d. samples $\xi_1, \dots, \xi_n \sim P$, often referred to as the *training set*, with associated empirical measure $P_n := \frac{1}{n} \sum_{i=1}^n \delta_{\xi_i}$. The ERM principle defines an estimator $\hat{\boldsymbol{\theta}}_n$ via

$$\boldsymbol{\theta}_n \in \arg \min_{\boldsymbol{\theta} \in \Theta} \mathcal{R}(\boldsymbol{\theta}, P_n). \quad (1.2)$$

A familiar example of (1.2) is maximum likelihood estimation, wherein Θ parametrizes a class of probability measures $\{P_{\boldsymbol{\theta}} : \boldsymbol{\theta} \in \Theta\}$, and $\mathcal{R}(\boldsymbol{\theta}, P_n)$ is the empirical Kullback-Liebler risk between $P_{\boldsymbol{\theta}}$ and P_n . More generally, the user may define a measurable *loss* $\ell : \Theta \times \Xi \rightarrow \mathbb{R}$ and define

$$\mathcal{R}(\boldsymbol{\theta}, P) := \mathbb{E}_{\xi \sim P} [\ell(\boldsymbol{\theta}, \xi)] \implies \mathcal{R}(\boldsymbol{\theta}, P_n) = \frac{1}{n} \sum_{i=1}^n \ell(\boldsymbol{\theta}, \xi_i), \quad (1.3)$$

where $\mathbb{E}[\cdot]$ denotes the expectation functional. While many works may consider the “average loss” format of (1.3) to be synonymous with empirical risk minimization, we consider more general functionals throughout this thesis. The first of two highly relevant examples throughout the thesis is the functional

$$\mathcal{R}(\boldsymbol{\theta}, P) = \max_{Q \in \mathcal{Q}(P)} \mathbb{E}_Q [\ell(\boldsymbol{\theta}, \xi)] \quad (1.4)$$

where $\mathcal{Q}(P)$ is a to-be-specified set of probability measures on Ξ related to P , called the *uncertainty set*. By computing the worst-case expectation over many probability measures, this objective naturally promotes aversion to risk and upweights the “harder” examples in the training set. Letting $P^{\otimes b}$ denote the probability measure governing b i.i.d. draws from P , the second prominent example is

$$\mathcal{R}(\boldsymbol{\theta}, P) = \mathbb{E}_{P^{\otimes b}} [\ell_b(\boldsymbol{\theta}, \xi_1, \dots, \xi_b)], \quad (1.5)$$

where b denotes a batch size and $\ell_b : \Theta \times \Xi^b \rightarrow \mathbb{R}$ is a loss function (dependent on this batch size). This type of objective is often applied in unsupervised problems, where structure within and between data points (similarity, etc.) acts as a replacement for labels. Notably, both (1.4) and (1.5) involve relationships between data points and are not expressible simply as average losses across the observed sample. This is a challenge from both statistical and optimization perspectives, and will guide much of the work in forthcoming sections. We give a brief overview of our topics of study and defer specific details to the individual chapters.

1.2.2 Distribution Shift and Distributional Robustness

A natural application of the risk functional (1.4) is when the practitioner observes training data from P , but expects the estimated parameter to be evaluated on test data from another probability measure Q . For instance, a training and test set of images may differ in distribution due to heterogeneous lighting conditions (a natural shift) or corruption of the test images through blurring (a synthetic shift). While this phenomenon—commonly known as *distribution shift* [Quiñonero-Candela et al., 2022] in the literature—is often treated as a specialized problem setting, we emphasize that virtually all machine learning models are deployed on data that are not distributed identically to the training data. In the empirical setting, by carefully designing the uncertainty set $\mathcal{Q}(P_n)$, one hopes that the parameter is *distributionally robust*, or will have approximately uniform performance across many probability measures that are sufficiently close to P . This framework, aptly named *distributionally robust optimization (DRO)*, differs from classical work in distribution shift [Huang et al., 2006, Sugiyama et al., 2007, 2008], in that the user has no information about the particular distribution Q on which the model will be evaluated (such as unlabeled covariates). Thus, $\mathcal{Q}(\cdot)$ is often defined using qualitative conditions such as existence and boundedness of a likelihood ratio $\frac{dQ}{dP}$ (see Chapter 2).

Under various assumptions on \mathcal{Q} , the maximization problem over probability measures

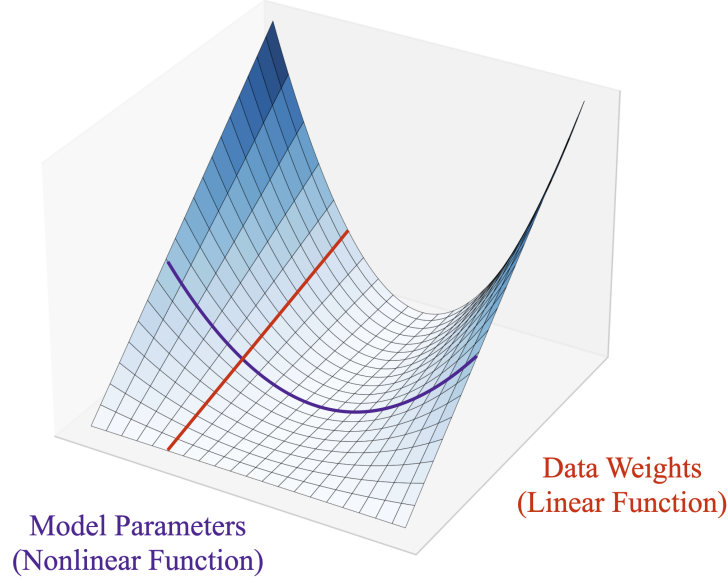
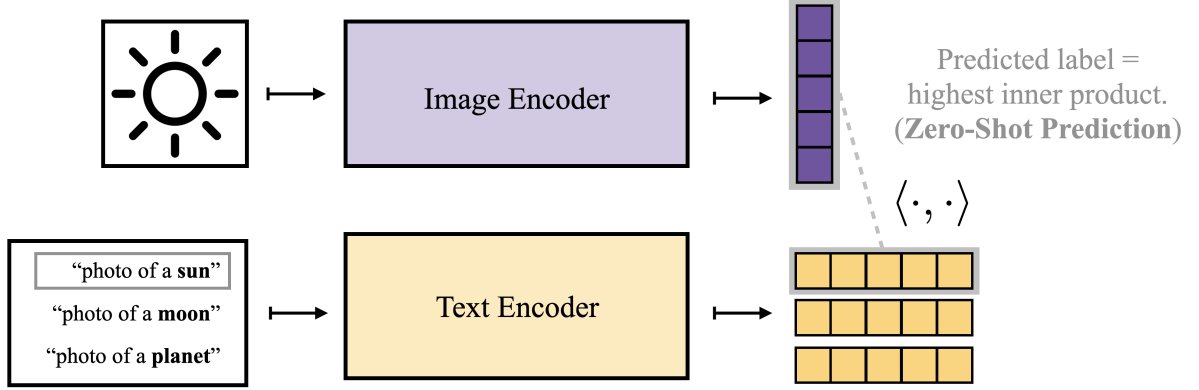


Figure 1.4: **Linear-Nonlinearly Coupled Objective.** Illustration of the objective (1.6), which is possibly nonlinear in $\boldsymbol{\theta}$ (the model parameters) but linear in the data weights \mathbf{q} .

in $\mathcal{Q}(P_n)$ can often be made into a finite-dimensional program of the form

$$\max_{\mathbf{q} \in \mathcal{Q}(P_n)} \mathbb{E}_{\mathbf{Q}} [\ell(\boldsymbol{\theta}, \xi)] = \max_{\mathbf{q} \in \mathcal{Q}} \sum_{i=1}^n q_i \ell(\boldsymbol{\theta}, \xi_i). \quad (1.6)$$

where $\mathcal{Q} \subseteq \mathbb{R}^n$ contains weight vectors $\mathbf{q} = (q_1, \dots, q_n)$ on each training example. However, as we discuss in Chapter 2, this optimization problem will often be intractable for large values of n , and typical stochastic gradient-type algorithms will face challenges that are unseen when minimizing sample averages. Statistical questions include how to construct (sub)gradient estimators for these maximum-type objectives with low bias and variance, and how such objectives behave as n grows (noting that the dimension of the set \mathcal{Q} grows linearly in n). Furthermore, (1.6) generates a min-max problem in which the objective is nonlinear in the minimizing (or primal) variables $\boldsymbol{\theta}$ but is in fact linear in the maximizing (or dual) variables \mathbf{q} (see Figure 1.4). We study this “dual-linear” coupling extensively in Chapter 3.



Prompts: Labels converted into natural language pseudo-captions.

Figure 1.5: **Prompting of Foundation Models.** Illustration of the indirect predictor (1.7), which relies on image and text encoders pre-trained via self-supervision (see Figure 1.3)

1.2.3 Universal Representations and Zero-Shot Prediction

The risk functional (1.5) is based on the objective used to train the CLIP series of neural network models [Radford et al., 2021] and its predecessors [Chen et al., 2020]. Here, we have that $\Xi = \mathcal{X} \times \mathcal{Z}$, where \mathcal{X} and \mathcal{Z} denote observation spaces of two data modalities (most commonly images and text). Then, $\theta = (\alpha, \beta)$, a pair of functions $\alpha : \mathcal{X} \rightarrow \mathbb{R}^d$ and $\beta : \mathcal{Z} \rightarrow \mathbb{R}^d$ which are optimized so that $\alpha(x)$ and $\beta(z)$ are close in \mathbb{R}^d if and only if $x \in \mathcal{X}$ and $z \in \mathcal{Z}$ are semantically similar. Remarkably, these encoders can be used to perform downstream prediction tasks without *any direct labeled data*—generalization of this kind is known as *zero-shot prediction*.

To understand zero-shot prediction, we first contrast it with the related setting of few-shot learning. Let $x \in \mathcal{X}$ be an input that accompanies a label $y \in \mathcal{Y}$, which could be a class label or even a structured object such as a parse tree. Common to both zero-shot prediction and few-shot learning is a pre-training procedure in which a large unlabeled dataset $x_1, \dots, x_n \in \mathcal{X}$ is used to produce α . Pre-training typically occurs through the process of *self-supervised learning (SSL)*, using a pretext task (quantified by an objective

like (1.5)) that can be solved with only instances of \mathbf{x} (e.g., filling in a blank image patch). In few-shot learning, the user may then access a labeled dataset $(\mathbf{x}_1^{\text{lab}}, \mathbf{y}_1^{\text{lab}}), \dots, (\mathbf{x}_{n^{\text{lab}}}^{\text{lab}}, \mathbf{y}_{n^{\text{lab}}}^{\text{lab}})$ from which a predictor can be trained inexpensively. This often takes the form of a linear classifier $\mathbf{x} \mapsto \mathbf{W}\boldsymbol{\alpha}(\mathbf{x}) + \mathbf{b}$ for $\mathbf{W} \in \mathbb{R}^{|\mathcal{Y}| \times d}$ and $\mathbf{b} \in \mathbb{R}^{|\mathcal{Y}|}$. Because such data is not available for zero-shot prediction, the ingenuity of practitioners has yielded the following solution: if (1) each pre-training example \mathbf{x}_i (a web image, say) is paired with another *view* $\mathbf{z}_i \in \mathcal{Z}$ (e.g., a caption in natural language) and (2) if each label $\mathbf{y} \in \mathcal{Y}$ can intelligently be embedded into \mathcal{Z} , then the relationship between each \mathbf{x}_i and \mathbf{z}_i could provide the means to perform prediction. Concretely, one learns the complementary encoder $\boldsymbol{\beta}$ during pre-training and designs a number of *prompts* $\mathbf{z}_j^{\mathbf{y}}$ for $\mathbf{y} \in \mathcal{Y}$ and $j = 1, \dots, M$. Then, the function

$$\mathbf{x} \mapsto \arg \max_{\mathbf{y} \in \mathcal{Y}} \frac{1}{M} \sum_{j=1}^M \langle \boldsymbol{\alpha}(\mathbf{x}), \boldsymbol{\beta}(\mathbf{z}_j^{\mathbf{y}}) \rangle \quad (1.7)$$

is employed for prediction. An example of a prompt is the template text “photo of a ____”, where the blank can be filled by the textual representation of the class (e.g., “cat” or “dog”). The entire procedure is depicted in Figure 1.5.

The zero-shot prediction pipeline, from pre-training to prompt selection, is clearly a wild departure from what is explained by statistical learning theory. Moreover, while some components of these systems have been studied in the context of few-shot learning (such as the reasons why various pre-training objectives result in encoders that provably accelerate learning), unique aspects of zero-shot prediction, such as the role of prompting and the cost of “translating” modalities, have not yet received theoretical treatment. This is the subject of Chapter 4.

1.3 Contributions and Outline

We study these problems of interest across the next three chapters of the dissertation. We conclude in the final chapter with a summary and discussion. Before stating the contributions, we briefly comment on some work that is not included in this dissertation: 1) a

collaborative work on multivariate time series prediction with graph neural networks (featured at the KDD 2023 Workshop on Mining and Learning with Graphs) conducted during an industry internship [Yang et al., 2023], 2) a collaborative work on using black-box prediction sets to classify animal behavior from accelerometry data with uncertainty quantification [Agarwal et al., 2024], 3) a collaborative work studying the relationship of neuropsychiatric symptoms such as depression with Alzheimer’s disease (under review at *Alzheimer’s & Dementia*), and 4) a work for which the dissertation author is the primary contributor that develops stochastic algorithms for solving independent component analysis problems with supervision (under review at ICASSP 2026).

Practical Algorithms for Distributionally Robust Optimization In Chapter 2, we provide scalable, stochastic algorithms for solving distributionally robust optimization (DRO) problems. We prove the first theoretical linear convergence rate for stochastic algorithms on regularized DRO problems and handle several practical questions of interest. These include duality results over the probability simplex that not only allow for the efficient implementation of the algorithm but also give guidance on choosing hyperparameters such as the uncertainty set. We also specialize our results to specific objectives (studied previously in quantitative finance) known as spectral risk measures (SRMs), which we prove can be represented as DRO objectives. We derive several appealing properties of SRMs, including bias bounds for sample estimators under general conditions and computational properties that allow for simpler implementations than previous DRO approaches. The subject of the chapter is not only the theoretical aspects, but also extensive experimentation across tabular, image, and natural language examples. In particular, we test the proposed algorithms in group distribution shift/fairness scenarios and demonstrate performance benefits in terms of worst-case group-wise error.

This chapter is joint work with Vincent Roulet, Krishna Pillutla, Lang Liu, and Zaid Harchaoui. The results span two papers, Mehta et al. [2023] published at AISTATS 2023 and Mehta et al. [2024b] published at ICLR 2024. The latter was accepted as a Spotlight

(top 5% of submissions).

Complexity Guarantees for Semilinear Min-Max Optimization In Chapter 3, we consider general saddle point optimization problems (a.k.a. min-max problems) where the coupled term in the objective is linear in at least one of either the primal or dual variables. This “semilinear” min-max problem includes not only distributionally robust optimization as a special case but also other classes such as convex minimization with functional constraints. We provide a constructive upper bound on the complexity guarantees for controlling the primal-dual gap criterion, in that the algorithm is derived by way of the analysis (contrary to the experimentally-driven methods of Chapter 2). The method improves upon classical approaches to nonbilinearly-coupled min-max problems such as extragradient [Korpelevich, 1976, Nemirovski, 2004], Popov’s method [Popov., 1980], or dual extrapolation [Nesterov, 2007a] by delicately using the non-uniformity Lipschitz and smoothness constants of the nonlinear components of the objective. We also present specialized algorithms for when the linear variable has a rectangular feasible set, i.e., it can be optimized along its coordinates, which yields improved complexities in even the bilinear setting.

This chapter is joint work with Jelena Diakonikolas and Zaid Harchaoui. A portion of the ideas in this work were published in our NeurIPS 2024 paper Mehta et al. [2024a]. They are generalized and improved significantly in the chapter, and a journal manuscript is under review [Mehta et al., 2025].

Generalization Bounds for Data Curation and Zero-Shot Prediction Chapter 4 targets questions such as model reuse and self-supervised learning mentioned in Section 1.2.3. We study the estimation of a joint probability distribution over multimodal data (or a linear functional thereof) when given knowledge of the true marginal distributions of each modality. This problem is motivated by the prominent practice in large-scale machine learning to rebalance a pre-training dataset in order to achieve a particular marginal distribution over subgroups of variables (e.g., text data with equal representation across languages). Con-

sidering the estimation of a linear functional of the data-generating distributions, we prove a non-asymptotic bound on the mean squared error when incorporating the marginals into the estimator. Specifically, a balancing procedure known in statistical contexts as raking ratio estimation, iterative or bi-proportional fitting, or the Sinkhorn algorithm is used to alter the empirical measure into one that satisfies the marginal constraints. The bound is then proven using a recursive formula for iterations of the procedure, which furnishes both a first-order term (which improves upon the variance of the empirical mean estimator) and explicit higher-order terms. Furthermore, when the iteration number is scaled appropriately, we recover the efficient asymptotic variance that was derived in [Bickel et al. \[1991\]](#) using tools of (asymptotic) semiparametric efficiency theory. Our approach, based on the recursive formula, yields a new closed-form expression for this first-order term (which was previously stated variationally via projections) in terms of the conditional mean operator between the modalities. The second part of Chapter 4 casts a zero-shot prediction procedure that is used ubiquitously in modern practice (from pre-training to prompting) as a formal estimator of a function defined on the population distribution. We establish performance limits of zero-shot prediction as compared to the Bayes optimal predictor on downstream tasks. Both theoretical analyses are demonstrated with experiments on language-image pre-training and zero-shot image classification.

This chapter is joint work with Lang Liu and Zaid Harchaoui. It contains some theoretical and experimental results on data curation from our NeurIPS 2024 paper [[Liu et al., 2024](#)] and results on prompting and downstream prediction from our ICML 2025 paper [[Mehta and Harchaoui, 2025](#)]. The latter was accepted as an Oral (top 1% of submissions). Part of this work was conceptualized and developed at the Simons Institute for the Theory of Computing as part of the Modern Paradigms in Generalization (2024) program. The content of previous chapters also benefited greatly from the fruitful scientific interactions at this program.

Software The papers mentioned above accompany software packages to reproduce the code and experiments. These include `1erm` [[Mehta et al., 2023](#)], `prospect` [[Mehta et al.,](#)

2024b], drago [Mehta et al., 2024a], balancing [Liu et al., 2024], and zeroshot [Mehta and Harchaoui, 2025]. Pip-installable packages called `deshift` and `drlearn` collect versions of the DRO algorithms developed in this dissertation for PyTorch and scikit-learn workflows, respectively. All packages are linked at <https://ronakdm.github.io/software>.

Chapter 2

LEARNING UNDER DISTRIBUTION SHIFT WITH LIKELIHOOD RATIOS

2.1 Introduction

In Chapter 1, we introduced alternate notions of risk functionals by which we may define parameters of interest. One such risk, which was (1.4), takes the maximum expected loss achievable over an uncertainty of distributions. When evaluated at an empirical measure P_n , this yields a variant of empirical risk minimization known as *distributionally robust optimization (DRO)*. This chapter is dedicated to studying the properties of such objectives and designing scalable algorithms for optimizing them in practice.

While particular cases of the uncertainty set $\mathcal{Q}(P_n)$ (e.g., a closed ball in f -divergence or Wasserstein distance) have been studied extensively in DRO [Namkoong and Duchi, 2016, Shapiro, 2017, Kuhn et al., 2019, Duchi and Namkoong, 2021], we begin at a broader starting point: placing assumptions on $\mathcal{Q}(\cdot)$ that ensure that the maximization over probability measures can be reformulated into a finite-dimensional program of the form

$$\max_{\mathbf{q} \in \mathcal{Q}} \sum_{i=1}^n q_i \ell(\boldsymbol{\theta}, \xi_i), \quad (2.1)$$

where, as in (1.6), $\mathcal{Q} \subseteq \Delta^{n-1} := \{(p_1, \dots, p_n) : \sum_{i=1}^n p_i = 1, p_i \geq 0 \forall i\}$ contains probability mass vectors. Here, we will call $\boldsymbol{\theta}$ the primal variable and \mathbf{q} the dual variable. Empirical risk minimization (ERM) is recovered by (2.1) when $\mathcal{Q} = \{\mathbf{1}/n\}$ for $\mathbf{1} = (1, \dots, 1)$ (i.e., there is only one feasible dual variable). We call the minimization of (2.1) a *likelihood ratio-based distributionally robust optimization (LR-DRO)* problem, as the finiteness will result from considering $\mathcal{Q}(P_n)$ to contain only distributions that are absolutely continuous with respect to P_n (see Section 2.2 for details). A canonical example of an LR-DRO problem is f -DRO,

in which we define $\mathcal{Q}(P_n)$ as a closed ball about P_n in f -divergence. There is also another class of risk functionals known as *spectral risk measures (SRMs)*, which are lesser-known as DRO objectives in machine learning.

Despite the increasing adoption of objectives of the form (2.1), optimization approaches have relied on using the full-batch or stochastic subgradient method out-of-the-box [Fan et al., 2017, Kawaguchi and Lu, 2020, Laguel et al., 2020, Levy et al., 2020], both enduring considerable limitations. The per-iteration complexity of full-batch methods (those that invoke every first-order oracle $\{\ell(\cdot, \xi_i), \nabla \ell(\cdot, \xi_i)\}_{i=1}^n$) is $\tilde{O}(n)$ function/gradient evaluations. For stochastic¹ variants, unbiased estimates of any subgradient, while needing only $O(1)$ gradient evaluations, still need $O(n)$ function calls to compute the vector of losses $(\ell(\boldsymbol{\theta}, \xi_1), \dots, \ell(\boldsymbol{\theta}, \xi_n))$, as we describe in Section 2.4. This yields the same per-iteration complexity as the full-batch method in automatic differentiation frameworks. A number of methods abandon convergence to the minimal risk altogether and resort to $O(1)$ -time stochastic subgradient updates, but are biased [Kawaguchi and Lu, 2020, Levy et al., 2020], i.e., do *not* converge to the optimal value. In this chapter, we present a class of stochastic algorithms that enjoy a theoretical *linear* convergence guarantee to the solution of a regularized LR-DRO problem while needing only a constant number of calls to the first-order oracles per iteration and have excellent empirical performance compared to baselines. We also provide theoretical and practical guidance on selecting the uncertainty set and other parameters defining the objective. Despite the min-max structure, we adopt a single-hyperparameter primal-only viewpoint to emphasize the importance of practicality, in the spirit of variance-reduced stochastic algorithms used for finite sum objectives [Johnson and Zhang, 2013, Defazio et al., 2014].

The remainder of the chapter is outlined as follows. Section 2.2 introduces likelihood ratio-based DRO and spectral risk measures. Section 2.3 establishes statistical properties such as an upper bound on the bias incurred by stochastic estimates of SRMs. For vanilla

¹We use the term “stochastic” to include both *streaming* algorithms in which fresh samples from the data-generating distribution are provided at each iterate, and *incremental* algorithms, in which multiple passes are made over a fixed dataset.

stochastic gradient-type algorithms that may only converge to a ball about the optimum, these bias bounds allow us to determine the worst-case radius of this ball. In Section 2.4 and Section 2.5, we present the details of our modified form of the objective (2.5) and the proposed algorithm, with its convergence analysis given in Section 2.6. In Section 2.7 we answer questions of practical performance, such as how to select the uncertainty set and how to efficiently solve the associated maximization problem. Important connections to related fields such as empirical likelihood [Owen, 1990] are given in Section 2.8. Numerical benchmarks on regression and classification problems are given in Section 2.9. Finally, we consider several extensions of the algorithms and analyses in Section 2.11 with concluding remarks in Section 2.12.

2.2 Preliminaries

We first characterize LR-DRO as categorization of various DRO objectives, and then introduce the two example objectives used extensively throughout the entire chapter: 1) closed balls in f -divergence (i.e., f -DRO) and 2) spectral risk measures (SRMs). On terminology, we sometimes identify DRO objectives with their uncertainty set, as we have done with f -DRO above. Some results (such as the duality relations in Section 2.7) will be established separately for f -DRO and SRMs. Furthermore, because SRMs are relatively new to machine learning, we dedicate Section 2.3 to some novel results deriving their properties.

In the LR-DRO class, the primary assumption on $\mathcal{Q}(P)$ for any probability measure P over Ξ is that for every $Q \in \mathcal{Q}(P)$, it holds that $Q \ll P$ (i.e., Q is absolutely continuous with respect to P). By narrowing the options for the perturbed distributions in this fashion, any $Q \in \mathcal{Q}(P)$ is exactly specified by the *likelihood ratio*, or Radon-Nikodym derivative $\frac{dQ}{dP} : \Xi \rightarrow [0, \infty)$, that is,

$$\mathcal{R}(\boldsymbol{\theta}, P) = \max_{Q \in \mathcal{Q}(P)} \mathbb{E}_Q [\ell(\boldsymbol{\theta}, \xi)] = \max_{\beta \in \mathcal{B}(P)} \mathbb{E}_P [\beta(\xi) \ell(\boldsymbol{\theta}, \xi)], \quad (2.2)$$

where $\mathcal{B}(P) := \{ \frac{dQ}{dP} : Q \in \mathcal{Q}(P) \}$. The set $\mathcal{B}(P)$ can often be specified by explicitly stating conditions on the likelihood ratio β instead of the probability measure Q , as seen in

the following examples. Below, let $\mathcal{B}_0(P)$ denote the collection of valid likelihood ratios for probability measures under P , or the non-negative measurable functions $\beta(\cdot)$ satisfying $\mathbb{E}_P[\beta(\xi)] = 1$.

Example 2.2.1 (f -DRO). In the case of f -DRO, we have that

$$\mathcal{Q}(P) := \{Q : Q \ll P \text{ and } D_f(Q\|P) \leq \rho\},$$

where $D_f(Q\|P) = \int_{\Xi} f\left(\frac{dQ}{dP}\right) dQ$ is an f -divergence generated by f and ρ is a radius parameter (see Appendix A.1.2 for a review of f -divergences). Then, associated with the f -ball uncertainty set is the feasible set

$$\mathcal{B}(P) = \{\beta \in \mathcal{B}_0(P) : \mathbb{E}_P[f(\beta(\xi))] \leq \rho\}. \quad (2.3)$$

In words, f -DRO places a moment condition on β , which is much weaker than the assumptions placed on spectral risk measures.

Example 2.2.2 (Spectral Risk Measures). For a bounded, measurable, non-decreasing function $s : (0, 1) \rightarrow [0, \infty)$ satisfying $\int_0^1 s(t) dt = 1$ (called the *spectrum*) and a real-valued random variable X with cumulative distribution function (CDF) F , we define the functional \mathbb{L}_s as

$$\mathbb{L}_s[X] = \int_0^1 s(t) F^{-1}(t) dt$$

where $F^{-1}(p) := \inf \{x : F(x) \geq p\}$ is the right generalized inverse CDF or *quantile function* of X . By Hölder's inequality, the quantity above is well-defined when $\mathbb{E}_F|X| < \infty$. We call the function \mathbb{L}_s the *spectral risk measure (SRM)* with spectrum s . In the notation of (2.2), the learning objective associated with the SRM is then

$$\mathcal{R}(\boldsymbol{\theta}, P) := \mathbb{L}_s[\ell(\boldsymbol{\theta}, \xi)] \text{ for } \xi \sim P.$$

It remains to show that there exists an uncertainty set \mathcal{Q} such that

$$\mathbb{L}_s[\ell(\boldsymbol{\theta}, \xi)] = \max_{Q \in \mathcal{Q}(P)} \mathbb{E}_Q[\ell(\boldsymbol{\theta}, \xi)]$$

holds. This is established in the upcoming Proposition 2.3.1 from Section 2.3, where we prove that $\mathcal{Q}(P)$ corresponds to all Q such that $\frac{dQ}{dP}$ follows the same distribution (under P) as $s(U)$ for $U \sim \text{Unif}(0, 1)$. Thus, we immediately have that

$$\mathcal{B}(P) = \left\{ \beta \in \mathcal{B}_0(P) : \beta(\xi) \stackrel{d}{=} s(U) \text{ for } \xi \sim P, U \sim \text{Unif}(0, 1) \right\}. \quad (2.4)$$

Thus, the spectral risk measure uncertainty set constrains the *entire* marginal distribution of the likelihood ratio, as opposed to the moment conditions used in f -DRO.

Given (2.2), in the empirical setting, the learning problem of interest reduces to

$$\min_{\boldsymbol{\theta} \in \Theta} \max_{\beta \in \mathcal{B}(P_n)} \left\{ \mathbb{E}_{P_n} [\beta(\xi) \ell(\boldsymbol{\theta}, \xi)] = \frac{1}{n} \sum_{i=1}^n \beta(\xi_i) \ell(\boldsymbol{\theta}, \xi_i) \right\}. \quad (2.5)$$

Notice that the maximization in (2.5) only depends on the n numbers $(\beta(\xi_1)/n, \dots, \beta(\xi_n)/n)$, which can naturally be interpreted as importance weights for each training example of the objective, as $\frac{1}{n} \sum_{i=1}^n \beta(\xi_i) = \mathbb{E}_{\xi \sim P_n} [\beta(\xi)] = 1$. Even though distributions in the uncertainty set reduce to reweightings of the given training examples, we emphasize that the objective in (2.5) is viewed as a consistent estimator of its population counterpart, which contains distributions that need only be absolutely continuous with respect to P . We proceed to the specific properties of SRMs in Section 2.3 and then discuss algorithms to optimize (2.1) in Section 2.4.

2.3 Properties of Spectral Risk Measures

SRMs have been studied extensively in quantitative finance [Artzner et al., 1999, Föllmer and Schied, 2002, Rockafellar and Uryasev, 2013, Acerbi and Tasche, 2002, Pflug and Ruszczyński, 2005, Kuhn et al., 2019] and convex analysis Rockafellar and Royset [2014], Ben-Tal and Teboulle [2007]. Despite being relatively underexplored in machine learning, one particu-

lar SRM called the *superquantile* or *conditional value-at-risk (CVaR)* has recently received careful attention in the learning setting [Curi et al., 2020, Levy et al., 2020, Laguel et al., 2020, 2021], though it has not been unified with DRO as shown in this chapter. As we will see in Section 2.7, they also have particular computational properties that make them an appealing choice for practitioners.

Because SRMs are really summaries of univariate real-valued random variables, we will derive some relevant properties by only considering a real-valued random variable X representing the loss at a particular parameter θ , i.e., $X = \ell(\theta, \xi)$. The randomness is induced by $\xi \sim P$. We first provide a novel variational form of this functional that is expressed in terms of likelihood ratios, which gives (2.4). We then state a number of examples. In the empirical setting, the maximum can be expressed in terms of an even simpler feasible set, which we discuss alongside optimization details in Section 2.4.

2.3.1 Spectral Risk Measures and the Likelihood Ratio

Before stating the result, note for context that the uncertainty sets in f -DRO (see (2.3)) are based on relatively coarse attributes of β . For example, the χ^2 -divergence is given by $f(x) = x^2 - 1$, indicating that any distribution such that the second moment is bounded as $\mathbb{E}_P[\beta^2(\xi)] \leq \rho + 1$ is feasible. As we will see below, SRMs constrain the entire marginal distribution of β , which is a much stronger condition.

Proposition 2.3.1. *Let s be left-continuous (in addition to being bounded, measurable, non-negative, and non-decreasing). Then, there exists a unique CDF $G_s(v) := \sup \{t : s(t) \leq v\}$ such that*

$$\mathbb{L}_s[\ell(\theta, \xi)] = \max_{\beta \in \mathcal{B}(P)} \mathbb{E}_{\xi \sim P} [\beta(\xi) \ell(\theta, \xi)] \quad (2.6)$$

for

$$\mathcal{B}(P) = \{\beta \in \mathcal{B}_0(P) : \mathbb{P}[\beta(\xi) \leq v] = G_s(v)\}. \quad (2.7)$$

Furthermore, $\beta_\star(\xi)$ maximizing (2.6) can be written as a measurable function of $\ell(\boldsymbol{\theta}, \xi)$.

Proof. Let F be the CDF of $\ell(\boldsymbol{\theta}, \xi)$. Because s is left-continuous and non-decreasing, we have by Bobkov and Ledoux [2019, Proposition A.2] that there exists a unique CDF G_s given by the expression in the statement such that the right generalized inverse of G_s is s . Thus, we have that

$$\begin{aligned} \mathbb{L}_s[\ell(\boldsymbol{\theta}, \xi)] &= \int_0^1 s(t) F^{-1}(t) dt = \int_0^1 G_s^{-1}(t) F^{-1}(t) dt \\ &= \max \{ \mathbb{E} [V \ell(\boldsymbol{\theta}, \xi)] : V : \Omega \rightarrow \mathbb{R} \text{ measurable and } \mathbb{P} [V \leq v] = G_s(v) \} \\ &\geq \max \{ \mathbb{E} [\beta(\xi) \ell(\boldsymbol{\theta}, \xi)] : \beta \in \mathcal{B}(P) \}. \end{aligned} \quad (2.8)$$

To make the inequality above into an equality, we must prove that for V_\star maximizing (2.8), there exists a β_\star in the feasible set above such that $V_\star = \beta_\star(\xi)$. Accordingly, define,

$$\beta_\star(\xi) := G_s^{-1}(F(\ell(\boldsymbol{\theta}, \xi))) = s(F(\ell(\boldsymbol{\theta}, \xi))),$$

which is non-negative because s is non-negative. Measurability follows from the measurability of $\ell(\boldsymbol{\theta}, \cdot)$, F , and s . To show that β_\star is a likelihood ratio, compute the expectation:

$$\mathbb{E}_P [\beta_\star(\xi)] = \mathbb{E}_P [G_s^{-1}(F(\ell(\boldsymbol{\theta}, \xi)))] = \int_0^1 G_s^{-1}(t) dt = 1,$$

because $F(\ell(\boldsymbol{\theta}, \xi)) \sim \text{Unif}(0, 1)$ and $\int_0^1 s(t) dt = 1$. This completes the proof. \square

While Proposition 2.3.1 expresses a measure of tail risk as a distributionally robust objective, the reverse has been done for f -divergences, in that their dual form (see Shapiro [2017, Section 3.2]) is interpreted as a tail error. To understand the CDF G_s , or rather the random variable it corresponds to, we consider two examples of SRMs. Here, we present their continuous spectra $t \mapsto s(t)$ as opposed to their discretizations in (2.13) and (2.14).

- **Superquantile:** For $\tau \in (0, 1]$, the τ -superquantile [Rockafellar and Royset, 2013],

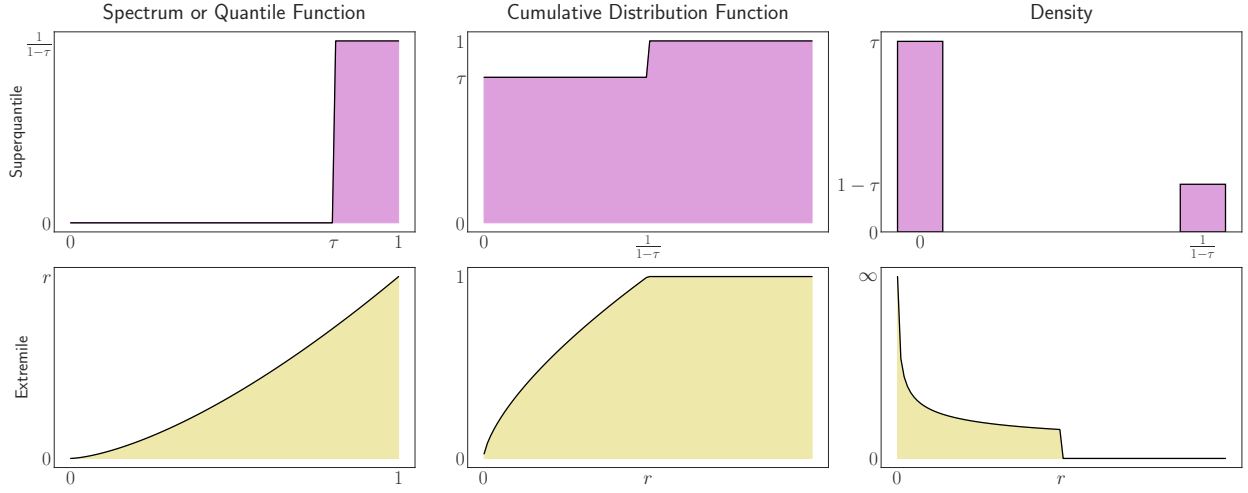


Figure 2.1: **Example Spectra and Induced Random Variables.** Visualization of the superquantile and extremile spectra and the induced distribution function G_s of the likelihood ratio.

a.k.a. conditional value-at-risk, is specified by $s(t) = \frac{1}{1-\tau} \mathbb{1} \{ \tau \leq t \leq 1 \}$. Thus, we have

$$G_s(v) = \tau \mathbb{1} \{ v \geq 0 \} + (1 - \tau) \mathbb{1} \{ v \geq 1/(1 - \tau) \}.$$

In other words, it is the CDF of a Bernoulli-like random variable that is 0 with probability $1 - \tau$ and $1/\tau$ with probability τ .

- **Extremile:** For $r > 1$, the r -extremile [Daouia et al., 2019] is specified by $s(t) = rt^{r-1}$.

Thus, we can compute by direct inversion

$$G_s(v) = \left(\frac{v}{r} \right)^{\frac{1}{r-1}} \mathbb{1} \{ v \in [0, r] \} + \mathbb{1} \{ v > r \}.$$

These examples are visualized alongside their quantile function, CDF, and densities in Figure 2.1.

2.3.2 Bias Bounds for Spectral Risk Measures

Bias bounds that are uniform over $\boldsymbol{\theta} \in \Theta$ help establish the performance guarantees of stochastic gradient descent-type algorithms, as done in [Levy et al. \[2020\]](#), for instance. They are also helpful for establishing uniform convergence of $\mathcal{R}(\cdot, P_n)$ process to $\mathcal{R}(\cdot, P)$, as one may perform the decomposition

$$\begin{aligned} \sup_{\boldsymbol{\theta} \in \Theta} |\mathcal{R}(\boldsymbol{\theta}, P_n) - \mathcal{R}(\boldsymbol{\theta}, P)| \leq \\ \underbrace{\sup_{\boldsymbol{\theta} \in \Theta} |\mathcal{R}(\boldsymbol{\theta}, P_n) - \mathbb{E}_P[\mathcal{R}(\boldsymbol{\theta}, P_n)]|}_{\text{concentration}} + \underbrace{\sup_{\boldsymbol{\theta} \in \Theta} |\mathbb{E}_P[\mathcal{R}(\boldsymbol{\theta}, P_n)] - \mathcal{R}(\boldsymbol{\theta}, P)|}_{\text{bias}}, \end{aligned}$$

and control the first term using standard technical tools, namely, concentration inequalities. As before, we will first consider a real-valued random variable X (which we may think of as $X = \ell(\boldsymbol{\theta}, \xi)$) with CDF F , and n i.i.d. copies denoted $X_1, \dots, X_n \sim F$. Assume in addition that $\mathbb{E}_F[|X_1|] < \infty$. As before, we let F^{-1} denote the quantile function and let s be a spectrum for the spectral risk measure $\mathbb{L}_s[\cdot]$. Similarly, let F_n and F_n^{-1} denote the empirical CDF and empirical quantile function, respectively. We use the notation $\mathbb{L}_s[F]$ below to indicate $\mathbb{L}_s[X]$ for X with CDF F . Under very general settings, namely the existence of greater-than-two moments of X , the bias can be shown to decay at an $n^{-1/2}$ rate. Later in this section, we make stronger assumptions to produce a bound that decays at rate n^{-1} . We will make use of the following theorem.

Theorem 2.3.1 (Theorem 2.10 of [Bobkov and Ledoux \[2019\]](#)). *Consider two probability distributions on \mathbb{R} with associated CDFs F and G , respectively, and quantile functions $F^{-1}(t) := \inf\{x \in \mathbb{R} : F(x) \geq t\}$ and $G^{-1}(t) := \inf\{x \in \mathbb{R} : G(x) \geq t\}$. Given that the random variables associated with F and G are absolutely integrable, we have that*

$$\int_{-\infty}^{\infty} |F(x) - G(x)| dx = \int_0^1 |F^{-1}(t) - G^{-1}(t)| dt,$$

where both the left and right-hand sides are finite.

We proceed to the first bias bound.

Proposition 2.3.2. *Assume that for real-valued random variable X with CDF F , it holds that $\mathbb{E}_F |X|^p < \infty$ for $p > 2$. Then, it holds that*

$$|\mathbb{E}_F[\mathbb{L}_s[F_n]] - \mathbb{L}_s[F]| \leq \sqrt{\frac{2}{n}} \left(\frac{p}{p-2} \right) \sup_{t \in (0,1)} |s(t) - 1| \cdot \|X\|_{\mathbf{L}^p(\mathbb{P})},$$

where $\|X\|_{\mathbf{L}^p(\mathbb{P})}^p := \int_{\Omega} X^p(\omega) d\mathbb{P}(\omega)$. We interpret the supremum above as the essential supremum according to the uniform measure on $(0, 1)$.

Proof. In the notation below, we will sometimes write $F_n(\cdot; \omega)$ to indicate a realisation of $F_n(\cdot)$ at outcome $\omega \in \Omega$, and similarly for $F_n^{-1}(\cdot; \omega)$ and $F^{-1}(\cdot)$. We first observe that

$$\mathbb{E}_F \left[\int_0^1 F_n^{-1}(t) dt \right] = \mathbb{E}_F \left[\frac{1}{n} \sum_{i=1}^n X_i \right] = \mathbb{E}_F [X] = \int_0^1 F^{-1}(t) dt.$$

Thus,

$$\begin{aligned} |\mathbb{E}_F[\mathbb{L}_s[F_n]] - \mathbb{L}_s[F]| &= \left| \mathbb{E}_F \left[\int_0^1 (s(t) - 1)(F_n^{-1}(t) - F^{-1}(t)) dt \right] \right| \\ &\leq \sup_{t \in [0,1]} |s(t) - 1| \cdot \mathbb{E}_F \left[\int_0^1 |F_n^{-1}(t) - F^{-1}(t)| dt \right] \\ &\leq \sup_{t \in [0,1]} |s(t) - 1| \cdot \sqrt{\mathbb{E}_F \left[\left(\int_0^1 |F_n^{-1}(t) - F^{-1}(t)| dt \right)^2 \right]}, \end{aligned}$$

where the second inequality is an application of Jensen's inequality. By Theorem 2.3.1, we also have that $\int_0^1 |F_n^{-1}(t; \omega) - F^{-1}(t)| dt = \int_{-\infty}^{\infty} |F_n(x; \omega) - F(x)| dx$, indicating that the random variable

$$\omega \mapsto \int_{-\infty}^{\infty} |F_n(x; \omega) - F(x)| dx \in \mathbf{L}^2(\mathbb{P}).$$

By the triangle inequality on $\mathbf{L}^2(\mathbb{P})$, we have that

$$\begin{aligned}
\sqrt{\mathbb{E} \left[\left(\int_0^1 |F_n^{-1}(t) - F^{-1}(t)| \, dt \right)^2 \right]} &= \sqrt{\mathbb{E} \left[\left(\int_{-\infty}^{\infty} |F_n(x) - F(x)| \, dx \right)^2 \right]} \\
&= \left\| \int_{-\infty}^{\infty} |F_n(x) - F(x)| \, dx \right\|_{\mathbf{L}^2(\mathbb{P})} \\
&\leq \int_{-\infty}^{\infty} \|F_n(x) - F(x)\|_{\mathbf{L}^2(\mathbb{P})} \, dx \\
&= \int_{-\infty}^{\infty} \sqrt{\mathbb{E} [|F_n(x) - F(x)|^2]} \, dx.
\end{aligned}$$

Next, notice that for fixed $x \in \mathbb{R}$, $nF_n(x) \sim \text{Binom}(n, F(x))$, so that

$$\mathbb{E} [|F_n(x) - F(x)|^2] = \text{Var } F_n(x) = \frac{F(x)(1 - F(x))}{n}.$$

Applying Lemma A.1.1 gives

$$\sqrt{\mathbb{E}_F \left[\left(\int_0^1 |F_n^{-1}(t) - F^{-1}(t)| \, dt \right)^2 \right]} \leq \sqrt{\frac{2}{n}} \left(\frac{p}{p-2} \right) \|X\|_{\mathbf{L}^p(\mathbb{P})},$$

which achieves the desired claim. \square

For context, Proposition 2.3.2 operates in general conditions that are of particular importance in optimization. To put this in context, a number of works provide non-asymptotic uniform learning bounds on spectral (and related) risks [Maurer et al., 2021, Khim et al., 2020, Lee et al., 2020]. However, these approaches require boundedness of the random variable of interest, which eliminates any potential application to heavy-tailed losses. Asymptotic approaches proceed by assuming Lipschitz continuity of the spectrum s [Shao, 1989], the trimming of s (i.e., $s(t) = 0$ for all $t \in [0, \alpha) \cup (1 - \alpha, 1]$ with $0 < \alpha \leq 1$) [Shorack, 2017, Shao, 1989], or bounded derivatives of the population quantile function F^{-1} [Xiang, 1995]. The τ -superquantile does not even have a continuous spectrum, whereas the spectrum of the r -extremile is not Lipschitz for $1 \leq r < 2$. Because s must be non-decreasing to achieve convexity (as we discuss in the upcoming Section 2.4), trimming the upper tail of s

is not reflective of practice. Finally, because losses such as the square loss or logistic loss can grow to infinity for unbounded inputs, the derivative $F^{-1}(t)$ as $t \rightarrow \infty$ cannot be assumed to be bounded. Proposition 2.3.2 only requires that the population loss satisfies a moment condition and holds without trimming or assumptions of boundedness or Lipschitz continuity on the spectrum. Other recent works employ concentration of the empirical measure in Wasserstein distance to give concentration inequalities for spectral risk measures under sub-Gaussian conditions and moment conditions similar to ours [Prashanth and Bhat, 2022, Bhat and Prashanth, 2019, Pandey et al., 2019].

When we are willing to assume boundedness of the random variable in question, we may place smoothness conditions on the inverse CDF to achieve a decay of n^{-1} .

Proposition 2.3.3. *Assume that F^{-1} is twice differentiable on $(0, 1)$, and the derivatives satisfy the conditions $\sup_{t \in (0, 1)} |[F^{-1}]'(t)| \leq M_1$ and $\sup_{t \in (0, 1)} |[F^{-1}]''(t)| \leq M_2$. Then,*

$$|\mathbb{E}_F [\mathbb{L}_s[F_n]] - \mathbb{L}_s[F]| \leq \frac{M_1 s(1)}{2n} + \frac{M_2}{n+1} + \frac{M_2 s(1)}{3(n+2)}.$$

Proof. The proof strategy will be to relate the order statistics of the observed data to the order statistics of uniform random variables via the inverse CDF transform. Then, using the properties of Beta-distributed random variables, we can explicitly compute bias terms where necessary. Denote by $X_{(1)} \leq \dots \leq X_{(n)}$ the order statistics of the sample. Each of these order statistics can be written equivalently as $X_{(i)} = F^{-1}(U_{(i)})$, where $U_{(i)}$ denotes the i -th order statistic of an independently drawn sample $U_1, \dots, U_n \sim \text{Unif}(0, 1)$. Write

$$\begin{aligned} \mathbb{E}_F [\mathbb{L}_s[F_n]] - \mathbb{L}_s[F] &= \mathbb{E}_F \left[\int_0^1 s(t) (F_n^{-1}(t) - F^{-1}(t)) dt \right] \\ &= \sum_{i=1}^n \mathbb{E}_F \left[\int_{(i-1)/n}^{i/n} s(t) (F_n^{-1}(t) - F^{-1}(t)) dt \right] \\ &= \sum_{i=1}^n \mathbb{E}_F \left[\int_{(i-1)/n}^{i/n} s(t) (F^{-1}(U_{(i)}) - F^{-1}(t)) dt \right]. \end{aligned}$$

Due to the boundedness of s and absolute integrability of the random variable X_1 , we may

apply Fubini's theorem so that

$$\mathbb{E}_F [\mathbb{L}_s[F_n]] - \mathbb{L}_s[F] = \int_{(i-1)/n}^{i/n} s(t) \mathbb{E}_F [F^{-1}(U_{(i)}) - F^{-1}(t)] dt.$$

Now, fix any $t \in ((i-1)/n, i/n)$, and apply a second-order Taylor expansion of F^{-1} so that

$$F^{-1}(U_{(i)}) - F^{-1}(t) = \underbrace{[F^{-1}]'(t)(U_{(i)} - t)}_{\text{first-order term}} + \underbrace{\frac{[F^{-1}]''(\tilde{U}_{(i)})}{2}(U_{(i)} - t)^2}_{\text{second-order term}},$$

where $\tilde{U}_{(i)}$ is a real-valued random variable which lies in between $U_{(i)}$ and t . Note that because $U_{(i)}$ follows the $\text{Beta}(i, n-i+1)$ distribution [David and Nagaraja, 2003], we have that $\mathbb{E}_{\text{Unif}(0,1)} [U_{(i)}] = i/(n+1)$, which we use in the upcoming calculations.

Controlling the first-order term. Applying the expectation and the sum, we then upper and lower bound the resulting quantity. Using that $[F^{-1}]'$ is bounded by M_1 and non-negative because F^{-1} is non-decreasing,

$$\begin{aligned} \sum_{i=1}^n \int_{(i-1)/n}^{i/n} s(t) [F^{-1}]'(t) \left(\frac{i}{n+1} - t \right) &\leq \sum_{i=1}^n \int_{(i-1)/n}^{i/n} s(t) [F^{-1}]'(t) \left(\frac{i}{n+1} - \frac{i-1}{n} \right) dt \\ &\leq M_1 \sum_{i=1}^n \int_{(i-1)/n}^{i/n} s(t) \left(\frac{i}{n+1} - \frac{i-1}{n} \right) dt \\ &\leq \frac{M_1 s(1)}{n} \sum_{i=1}^n \left(\frac{i}{n+1} - \frac{i-1}{n} \right). \end{aligned} \quad (2.9)$$

Similarly, we have that

$$\begin{aligned} \sum_{i=1}^n \int_{(i-1)/n}^{i/n} s(t) [F^{-1}]'(t) \left(\frac{i}{n+1} - t \right) &\geq \sum_{i=1}^n \int_{(i-1)/n}^{i/n} s(t) [F^{-1}]'(t) \left(\frac{i}{n+1} - \frac{i}{n} \right) dt \\ &\geq \frac{M_1 s(1)}{n} \sum_{i=1}^n \left(\frac{i}{n+1} - \frac{i}{n} \right), \end{aligned} \quad (2.10)$$

We compute the sums in (2.9) and (2.10) to give the upper bound. Observe that

$$\sum_{i=1}^n \left(\frac{i}{n+1} - \frac{i-1}{n} \right) = \frac{1}{2} \text{ and } \sum_{i=1}^n \left(\frac{i}{n+1} - \frac{i}{n} \right) = -\frac{1}{2},$$

so that

$$\left| \sum_{i=1}^n s(t) [F^{-1}]'(t) \mathbb{E}_{\text{Unif}(0,1)} [U_{(i)} - t] \right| \leq \frac{M_1 s(1)}{2n}.$$

Controlling the second-order term. Here, we use that $[F^{-1}]''$ is bounded by M_2 , so that

$$\begin{aligned} & \mathbb{E}_{\text{Unif}(0,1)} \left[[F^{-1}]''(\tilde{U}_{(i)})(U_{(i)} - t)^2 \right] \\ & \leq M_2 \mathbb{E}_{\text{Unif}(0,1)} [(U_{(i)} - t)^2] \\ & = M_2 \left(\text{Var } U_{(i)} + \left(\frac{i}{n+1} - t \right)^2 \right) \\ & = M_2 \left(\frac{i(n-i+1)}{(n+1)^2(n+2)} + \left(\frac{i}{n+1} - t \right)^2 \right) \\ & = M_2 \left(\frac{i}{(n+1)(n+2)} - \frac{i^2}{(n+1)^2(n+2)} + \left(\frac{i}{n+1} - t \right)^2 \right) \\ & \leq M_2 \left(\frac{1}{n+1} - \frac{i^2}{(n+1)^2(n+2)} + \left(\frac{i}{n+1} - t \right)^2 \right), \end{aligned}$$

where we used the bias-variance decomposition of the mean squared error and the moments of the $\text{Beta}(i, n-i+1)$ distribution. We group the second and third terms and take the integral, so that

$$\begin{aligned} & \sum_{i=1}^n \int_{(i-1)/n}^{i/n} s(t) \mathbb{E}_{\text{Unif}(0,1)} \left[[F^{-1}]''(\tilde{U}_{(i)})(U_{(i)} - t)^2 \right] dt \\ & \leq M_2 \left(\int_0^1 s(t) dt \right) \left(\frac{1}{n+1} \right) + M_2 \sum_{i=1}^n \int_{(i-1)/n}^{i/n} s(t) \left[\left(\frac{i}{n+1} - t \right)^2 - \frac{i^2}{(n+1)^2(n+2)} \right] dt \\ & \leq \frac{M_2}{n+1} + M_2 s(1) \sum_{i=1}^n \left[\int_{(i-1)/n}^{i/n} \left(\frac{i}{n+1} - t \right)^2 dt - \frac{i^2}{n(n+1)^2(n+2)} \right]. \end{aligned} \quad (2.11)$$

It remains to control the second term of (2.11). Expanding the square and integrating, we

generate telescoping terms and sum them to achieve

$$\begin{aligned} \sum_{i=1}^n \int_{(i-1)/n}^{i/n} \left(\frac{i}{n+1} - t \right)^2 dt &= \frac{1}{n} \sum_{i=1}^n \left(\frac{i}{n+1} \right)^2 - 1 + \frac{1}{3} \\ &\leq \frac{1}{n(n+1)^2} \sum_{i=1}^n i^2 \\ &= \frac{(2n+1)}{6(n+1)}. \end{aligned}$$

Then, including the negative term from (2.11), we have that

$$\begin{aligned} &\sum_{i=1}^n \left[\int_{(i-1)/n}^{i/n} \left(\frac{i}{n+1} - t \right)^2 dt - \frac{i^2}{n(n+1)^2(n+2)} \right] \\ &= \frac{2n+1}{6} \left(\frac{1}{n+1} - \frac{1}{n+2} \right) = \frac{2n+1}{6(n+1)(n+2)} \leq \frac{1}{3(n+2)}. \end{aligned}$$

Putting these steps together, the upper bound on the second-order term reads as

$$\mathbb{E}_{\text{Unif}(0,1)} \left[[F^{-1}]''(\tilde{U}_{(i)})(U_{(i)} - t)^2 \right] \leq \frac{M_2}{n+1} + \frac{M_2 s(1)}{3(n+2)},$$

which completes the proof. \square

Note that the Lipschitzness of the inverse CDF was also employed in [Levy et al. \[2020\]](#) to establish guarantees on DRO using the χ^2 -divergence and CVaR uncertainty sets. So far, the results we have shown have relied on the continuous viewpoint of SRMs as functionals of a real-valued random variable. We can derive useful closed-form expressions for the empirical variants of these quantities.

2.3.3 Empirical Spectral Risk Measures as L -Statistics

One key connection that we will exploit is that spectral risk measures can be expressed as a function of the order statistics of the sample in the empirical setting. This is because the empirical quantile function can be written in terms of the order statistics as $F_n^{-1}(t) = X_{(\lceil nt \rceil)}$, as seen in [Figure 2.2](#) (top). Notice in particular that when $t \in (\frac{i-1}{n}, \frac{i}{n})$, we have that

$F_n^{-1}(t) = X_{(i)}$, where end-points are chosen to make F_n^{-1} left continuous. Then, write

$$\begin{aligned}\mathbb{L}_s[F_n] &:= \int_0^1 s(t) \cdot F_n^{-1}(t) dt = \sum_{i=1}^n \left(\int_{(i-1)/n}^{i/n} s(t) \cdot X_{(\lceil nt \rceil)} dt \right) \\ &= \sum_{i=1}^n \left(\int_{(i-1)/n}^{i/n} s(t) dt \right) X_{(i)} = \sum_{i=1}^n \sigma_i X_{(i)}\end{aligned}\tag{2.12}$$

for $\sigma_i := \left(\int_{(i-1)/n}^{i/n} s(t) dt \right)$. The discretized weights $\sigma_1, \dots, \sigma_n$ are shown in Figure 2.2 (bottom). The expression (2.12) is called an *L-estimator* [Shorack, 2017, Maurer et al., 2021] for a generic linear combination of order statistics and an *L-risk* when the ordered elements are losses incurred on a training set [Maurer et al., 2021, Khim et al., 2020]. The σ_i 's allow the practitioner to interpolate between the average-case ($\sigma_i = 1/n \forall i$) and worst-case ($\sigma_n = 1$) performance on the training set. Such objectives have garnered a flurry of recent interest in machine learning [Fan et al., 2017, Williamson and Menon, 2019, Khim et al., 2020, Maurer et al., 2021, Holland and Mehdi Haress, 2022, Leqi et al., 2019, Lee et al., 2020, Kawaguchi and Lu, 2020].

As for particular examples: for $\tau \in (0, 1)$, the τ -CVaR (a.k.a. superquantile) [Rockafellar and Royset, 2013, Kawaguchi and Lu, 2020, Laguel et al., 2021] requires that $k = n(1 - \tau)$ elements of σ be non-zero with equal probability and that the remaining $n - k$ are zero. The *r-extremile* [Daouia et al., 2019] and *γ -exponential spectral risk measure* [Cotter and Dowd, 2006] define their spectra the equations below.

$$\sigma_i = \begin{cases} \frac{1}{n(1-\tau)} & \text{if } i \in \{\lceil n\tau \rceil, \dots, n\} \\ 1 - \frac{\lfloor n(1-\tau) \rfloor}{n(1-\tau)} & \text{if } \lfloor n\tau \rfloor < i < \lceil n\tau \rceil \end{cases} \quad \tau\text{-CVaR}, \tau \in (0, 1) \tag{2.13}$$

$$\sigma_i = \left(\frac{i}{n}\right)^r - \left(\frac{i-1}{n}\right)^r \quad b\text{-extremile}, r \geq 1 \tag{2.14}$$

$$\sigma_i = \frac{e^{\gamma(i/n)} - e^{\gamma(i-1)/n}}{1 - e^{-\gamma}} \quad \gamma\text{-ESRM}, \gamma > 0 \tag{2.15}$$

The multiple cases in the CVaR definition account for the instance in which $n(1 - \tau)$ is not an integer.

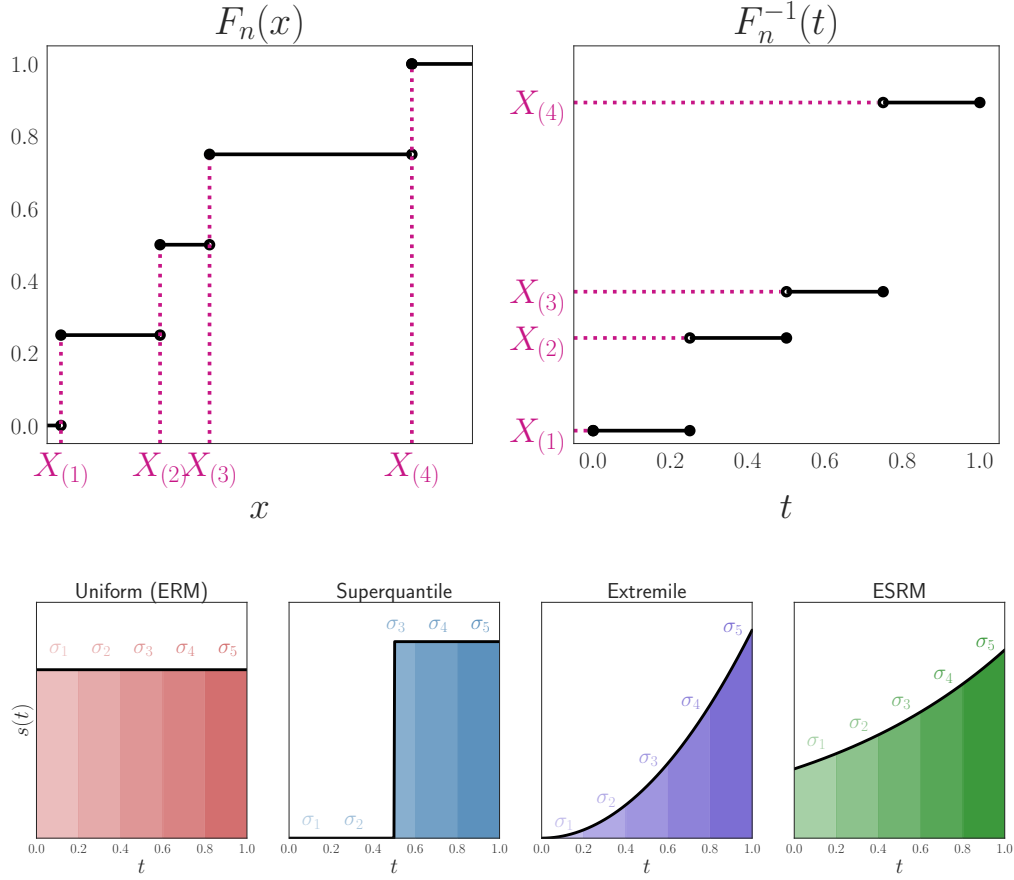


Figure 2.2: **Illustration of Spectral Risk and Quantile Functions.** The relationship between the order statistics, quantile functions, and the discretized spectrum is shown. **Top:** Empirical CDF F_n and quantile function F_n^{-1} of X_1, \dots, X_4 . **Bottom:** Continuous spectra $s(t)$ and their discretization $(\sigma_1, \dots, \sigma_5)$ for various risk measures.

While the bias bounds established above are general-purpose, we return to the optimization problem introduced in Section 2.1 in the next section. We provide a variational form of (2.12) in the empirical setting, which will provide the relationship with distributionally robust optimization by defining a corresponding uncertainty set for SRMs. Observe that because s is non-decreasing, it holds that $\sigma_1 \leq \dots \leq \sigma_n$. Thus, for any vector

$\mathbf{l} = (l_1, \dots, l_n) \in \mathbb{R}^n$ with ordered entries $l_{(1)} \leq \dots \leq l_{(n)}$, we have that

$$\sum_{i=1}^n \sigma_i l_{(i)} = \max_{\text{permutations } \pi} \sum_{i=1}^n \sigma_{\pi(i)} l_i = \max_{\mathbf{q} \in \mathcal{Q}(\sigma)} \sum_{i=1}^n q_i l_i,$$

where

$$\mathcal{Q}(\sigma) := \text{conv} \{ \text{permutations of } \sigma \}, \quad (2.16)$$

also known as the *permutahedron* of σ in \mathbb{R}^n . Using this uncertainty set, we address the optimization problem from Section 2.1.

2.4 Smoothing Maximum-Type Objectives

We now write our objective in the form of (2.1), which is amenable to being studied as an optimization problem. First, for notational ease, define $\ell_i(\boldsymbol{\theta}) := \ell(\boldsymbol{\theta}, \xi_i)$ and the vector $\ell(\boldsymbol{\theta}) = (\ell_1(\boldsymbol{\theta}), \dots, \ell_n(\boldsymbol{\theta}))$. Then,

$$\max_{\beta \in \mathcal{B}(P_n)} \frac{1}{n} \sum_{i=1}^n \beta(\xi_i) \ell(\boldsymbol{\theta}, \xi_i) = \max_{\mathbf{q} \in \mathcal{Q}} \langle \mathbf{q}, \ell(\boldsymbol{\theta}) \rangle, \quad (2.17)$$

where $\mathcal{Q} := \{ \beta(\xi_i)/n : \beta \in \mathcal{B}(P_n) \}$. As shown in Section 2.1, the elements of any $\mathbf{q} = (q_1, \dots, q_n) \in \mathcal{Q}$ are weights of a probability mass function, hence $\sum_{i=1}^n q_i = 1$. As an abuse of terminology, we will also refer to \mathcal{Q} as the *uncertainty set* in this context.

As is common in statistical and machine learning, we will regularize the quantity $\langle \mathbf{q}, \ell(\boldsymbol{\theta}) \rangle$ from (2.17) both on the primal and dual side, with respective hyperparameters $\mu > 0$ and $\nu > 0$, and consider the parameter space $\Theta = \mathbb{R}^d$. To emphasize the statistical context, rather than a primal-dual viewpoint, we consider the objective to be a composition of two functions. The first function \mathcal{A}_ν aggregates the vector (or equivalently, the histogram) of losses $\ell(\boldsymbol{\theta})$ for a particular $\boldsymbol{\theta} \in \Theta$ and produces a univariate summary, given by

$$\mathcal{A}_\nu(\mathbf{l}) := \max_{\mathbf{q} \in \mathcal{Q}} \left\{ \langle \mathbf{q}, \mathbf{l} \rangle - \nu D_f(\mathbf{q} \| \mathbf{1}/n) \right\},$$

where $D_f(\mathbf{q} \| \mathbf{1}/n) = \frac{1}{n} \sum_{i=1}^n f(nq_i)$ is an f -divergence with generator f over empirical mea-

asures. Because ν controls how much we penalize \mathbf{q} from shifting to far from the original uniform weights, we will refer to this parameter as the *shift cost*. The map \mathcal{A}_ν is in fact differentiable (as opposed to $\mathbf{l} \mapsto \max_{\mathbf{q} \in \mathcal{Q}} \langle \mathbf{q}, \mathbf{l} \rangle$), which will be an essential property when deriving the optimization algorithm and establishing the convergence guarantees. Thus, our problem of interest can be written as

$$\min_{\boldsymbol{\theta} \in \mathbb{R}^d} \left[\mathcal{L}(\boldsymbol{\theta}) := \mathcal{A}_\nu(\ell(\boldsymbol{\theta})) + \frac{\mu}{2} \|\boldsymbol{\theta}\|_2^2 \right], \quad (2.18)$$

As introduced in Section 2.2 maintain two prototypical classes of examples for the uncertainty set \mathcal{Q} : f -DRO and SRMs. Our theoretical analyses simply rely on \mathcal{Q} being a closed, convex subset of Δ^{n-1} , whereas implementation details (such as how to compute and differentiate the map \mathcal{A}_ν) will depend heavily on the specific uncertainty set used and are discussed in Section 2.7. Our goal is to design and analyze an iterative algorithm that produces a sequence $(\boldsymbol{\theta}^{(k)})_{k \geq 0}$ converging to the solution of (2.18) with a favorable dependence on the sample size n and constants associated to the loss ℓ . Consider the following standard assumptions.

Assumption 2.4.1. Each ℓ_i is convex, G -Lipschitz continuous w.r.t. $\|\cdot\|_2$, and L -smooth w.r.t. $\|\cdot\|_2$ (i.e., $\nabla \ell_i$ exists and is L -Lipschitz continuous). The generator f of D_f is α_n -strongly convex on the interval $[0, n]$. The constants μ , ν , and α_n are positive.

Common examples of α_n are $2n$ for the χ^2 -divergence and 1 for the KL-divergence. Before proposing an algorithm, we use gradient descent as a baseline. This involves showing that the composition $\boldsymbol{\theta} \mapsto \mathcal{A}_\nu(\ell(\boldsymbol{\theta}))$ is indeed differentiable (in fact, smooth) and characterizing the properties of the objective (2.18).

By Proposition A.1.1, we see that if f is α_n -strongly convex function over $[0, n]$, the f -divergence $D_f(\cdot \| \mathbf{1}/n)$ will be $(\alpha_n n)$ -strongly convex. Given Assumption 2.4.1, we assume that D_f is rescaled to be 1-strongly convex for ease of presentation, define $\text{Reg}(\mathbf{q}) := (\alpha_n n)^{-1} D_f(q \| \mathbf{1}/n)$, and rewrite our risk functional as

$$\mathcal{A}_\nu(\mathbf{l}) := \max_{\mathbf{q} \in \mathcal{Q}} \left\{ \langle \mathbf{q}, \mathbf{l} \rangle - \nu \text{Reg}(\mathbf{q}) \right\}, \quad (2.19)$$

In the following, let $\|\cdot\|$ denote an arbitrary norm on \mathbb{R}^d and let $\|\cdot\|_*$ denote its associated dual norm.

Lemma 2.4.1. *When \mathcal{Q} is closed and convex, and the map Reg is 1-strongly convex over \mathcal{Q} with respect to $\|\cdot\|$, then \mathcal{A}_ν is continuously differentiable with gradient given by*

$$q^{\text{opt}}(\mathbf{l}) := \nabla \mathcal{A}_\nu(\mathbf{l}) = \arg \max_{\mathbf{q} \in \mathcal{Q}} \{\langle \mathbf{q}, \mathbf{l} \rangle - \nu \text{Reg}(\mathbf{q})\} \in \mathbb{R}^n.$$

Moreover, q^{opt} is $(1/\nu)$ -Lipschitz continuous with respect to $\|\cdot\|_*$.

Proof. Because \mathcal{Q} is a closed subset of the compact set Δ^{n-1} , it is also compact. Because \mathcal{Q} is compact and convex, the strongly concave function $\mathbf{q} \mapsto \langle \mathbf{q}, \mathbf{l} \rangle - \nu \text{Reg}(\mathbf{q})$ has a unique maximizer. By Danskin's theorem [Bertsekas, 1997, Proposition B.25], we have that \mathcal{A}_ν is continuously differentiable with the given gradient formula. By Nesterov [2018, Lemma 6.1.2], it is also $(1/\nu)$ -Lipschitz continuous with respect to $\|\cdot\|_*$. \square

Finally, we can pass the smoothness result of Lemma 2.4.1 onto \mathcal{L} through a simple application of the chain rule. Letting $\nabla \ell(\boldsymbol{\theta}) \in \mathbb{R}^{n \times d}$ be the Jacobian of ℓ , we have that the gradient of $\boldsymbol{\theta} \mapsto \mathcal{A}_\nu(\ell(\boldsymbol{\theta}))$ is given by $\nabla \ell(\boldsymbol{\theta})^\top q^{\text{opt}}(\ell(\boldsymbol{\theta}))$. Then, for any $\boldsymbol{\theta}, \boldsymbol{\theta}' \in \mathbb{R}^d$, we have that

$$\begin{aligned} & \left\| \nabla \ell(\boldsymbol{\theta})^\top q^{\text{opt}}(\ell(\boldsymbol{\theta})) - \nabla \ell(\boldsymbol{\theta}')^\top q^{\text{opt}}(\ell(\boldsymbol{\theta}')) \right\|_2 \\ & \leq \left\| \nabla \ell(\boldsymbol{\theta})^\top (q^{\text{opt}}(\ell(\boldsymbol{\theta})) - q^{\text{opt}}(\ell(\boldsymbol{\theta}'))) \right\|_2 + \left\| (\nabla \ell(\boldsymbol{\theta}) - \nabla \ell(\boldsymbol{\theta}'))^\top q^{\text{opt}}(\ell(\boldsymbol{\theta}')) \right\|_2 \\ & \leq \frac{nG}{\nu} \|\ell(\boldsymbol{\theta}) - \ell(\boldsymbol{\theta}')\|_2 + L \|\boldsymbol{\theta} - \boldsymbol{\theta}'\|_2^2 \end{aligned} \tag{2.20}$$

$$\leq \left(\frac{nG^2}{\nu} + L \right) \|\boldsymbol{\theta} - \boldsymbol{\theta}'\|_2^2. \tag{2.21}$$

As a result \mathcal{L} is $(L + nG^2/\nu + \mu)$ -smooth and μ -strongly convex with respect to $\|\cdot\|_2$. By Nesterov [2018, Theorem 2.1.15], using the gradient descent update

$$\boldsymbol{\theta}^{(k+1)} = \boldsymbol{\theta}^{(k)} - \eta \nabla \mathcal{L}(\boldsymbol{\theta}^{(k)})$$

Algorithm 1 Prospect

Inputs: Initial points $\boldsymbol{\theta}^{(0)}$, stepsize $\eta > 0$, number of iterations t .

- 1: Set $\tilde{\boldsymbol{\theta}}_i^{(0)} = \hat{\boldsymbol{\theta}}_i^{(0)} = \boldsymbol{\theta}^{(0)}$ for all $i \in [n]$, with $\tilde{\boldsymbol{\theta}}^{(0)} = (\tilde{\boldsymbol{\theta}}_i^{(0)})_{i=1}^n$ and $\hat{\boldsymbol{\theta}}^{(0)} = (\hat{\boldsymbol{\theta}}_i^{(0)})_{i=1}^n$.
- 2: $\mathbf{q}^{(0)} = q^{\text{opt}}(\ell(\boldsymbol{\theta}^{(0)}))$, $\boldsymbol{\rho}^{(0)} = \mathbf{q}^{(0)}$.
- 3: Set $\mathbf{l}^{(0)} = (\ell_i(\hat{\boldsymbol{\theta}}_i^{(0)}))_{i=1}^n \in \mathbb{R}^n$, $\mathbf{g}^{(0)} = (\nabla r_i(\tilde{\boldsymbol{\theta}}_i^{(0)}))_{i=1}^n \in \mathbb{R}^{d \times n}$.
- 4: **for** $k = 0, \dots, t-1$ **do**
- 5: $i_k \sim \text{Unif}([n])$, $j_k \sim \text{Unif}([n])$.
- 6:
- 7: $\mathbf{v}^{(k)} = nq_{i_k}^{(k)} \nabla r_{i_k}(\boldsymbol{\theta}^{(k)}) - (n\rho_{i_k}^{(k)} \nabla r_{i_k}(\tilde{\boldsymbol{\theta}}_{i_k}^{(k)}) - \sum_{i=1}^n \rho_i^{(k)} \mathbf{g}_i^{(k)})$.
- 8: $\boldsymbol{\theta}^{(k+1)} = \boldsymbol{\theta}^{(k)} - \eta \mathbf{v}^{(k)}$.
- 9:
- 10: $\hat{\boldsymbol{\theta}}^{(k+1)} = \text{UpdateBiasReductionTable}(\boldsymbol{\theta}^{(k)}, \hat{\boldsymbol{\theta}}^{(k)}, i_k)$.
- 11: $\mathbf{l}^{(k+1)} = \ell(\hat{\boldsymbol{\theta}}^{(k+1)})$.
- 12: $\mathbf{q}^{(k+1)} = q^{\text{opt}}(\mathbf{l}^{(k+1)})$.
- 13:
- 14: $\tilde{\boldsymbol{\theta}}^{(k+1)}, \boldsymbol{\rho}^{(k+1)} = \text{UpdateVarianceReductionTable}(\boldsymbol{\theta}^{(k)}, \tilde{\boldsymbol{\theta}}^{(k)}, \boldsymbol{\rho}^{(k)}, i_k)$.
- 15: $\mathbf{g}^{(k+1)} = (\nabla r_i(\tilde{\boldsymbol{\theta}}_i^{(k+1)}))_{i=1}^n$.

Output: Final point $\boldsymbol{\theta}^{(t)}$

would achieve ε -suboptimality in

$$t = O\left(n \left(\frac{L}{\mu} + \frac{nG^2}{\mu\nu} + 1\right) \ln \frac{1}{\varepsilon}\right)$$

calls to the first-order oracles $\{(\ell_i, \nabla \ell_i)\}_{i=1}^n$. Given this baseline, the primary goal in designing the upcoming algorithm is to decouple the term n (representing a full-batch of oracle calls on a single iteration) and the condition number term $\frac{L}{\mu} + \frac{nG^2}{\mu\nu} + 1$ from being multiplied to being summed. This type of improvement has been achieved in the setting of empirical risk minimization via variance reduction [Johnson and Zhang, 2013, Defazio et al., 2014]. We devise a similar class of algorithms in the next section.

2.5 Stochastic Optimization with Bias and Variance Reduction

2.5.1 Algorithm Description

For ease of presentation, we use r_i to denote the regularized losses

$$r_i(\boldsymbol{\theta}) := \ell_i(\boldsymbol{\theta}) + \frac{\mu}{2} \|\boldsymbol{\theta}\|_2^2.$$

We present the Prospect method in Algorithm 1. At a high level, the algorithm stores a table of losses $\mathbf{l}^{(k)} \in \mathbb{R}^n$ and gradients $\mathbf{g}^{(k)} \in \mathbb{R}^{n \times d}$ which are used to approximate $\ell(\boldsymbol{\theta}^{(k)})$ and $\nabla r(\boldsymbol{\theta}^{(k)})$, respectively. In order to accompany the analysis, we use the notation $\hat{\boldsymbol{\theta}}_i^{(k)}$ to denote an element of \mathbb{R}^d such that $l_i^{(k)} = \ell_i(\hat{\boldsymbol{\theta}}_i^{(k)})$ and $\tilde{\boldsymbol{\theta}}_i^{(k)}$ to denote an element of \mathbb{R}^d such that $\mathbf{g}_i^{(k)} = \nabla r_i(\tilde{\boldsymbol{\theta}}_i^{(k)})$. In implementation, however, it may not be necessary to store the tables $\hat{\boldsymbol{\theta}}^{(k)} \in \mathbb{R}^{n \times d}$ or $\tilde{\boldsymbol{\theta}}^{(k)} \in \mathbb{R}^{n \times d}$. In fact, $\mathbf{g}^{(k)}$ itself may not need to be stored either. The algorithm is written in this way because lines 10 and 14 involve updates to these quantities that are left unspecified, which may be random (even conditionally on i_k).

2.5.2 Bias- and Variance-Reduced Updates

These updates may come in different forms, as long as Condition 2.5.1 and Condition 2.5.2 are satisfied. To understand them, we introduce the notation $\mathbb{E}_k[\cdot]$ to denote the conditional expectation given $\boldsymbol{\theta}^{(k)}$. This integrates the randomness of the randomly drawn index i_k as well as any additional randomness resulting from lines 10 and 14.

Condition 2.5.1 (Bias). The update in line 10 of Algorithm 1 satisfies the following. For functions $h_j : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$ for $j \in [n]$, it holds that

$$\sum_{j=1}^n \mathbb{E}_k \left[h_j(\boldsymbol{\theta}^{(k+1)}, \hat{\boldsymbol{\theta}}_j^{(k+1)}) \right] = \frac{1}{n} \sum_{j=1}^n \mathbb{E}_k [h_j(\boldsymbol{\theta}^{(k+1)}, \boldsymbol{\theta}^{(k)})] + \left(1 - \frac{1}{n}\right) \sum_{j=1}^n \mathbb{E}_k \left[h_j(\boldsymbol{\theta}^{(k+1)}, \hat{\boldsymbol{\theta}}_j^{(k)}) \right].$$

Condition 2.5.2 (Variance). The update in line 14 of Algorithm 1 satisfies the following.

For functions $h_i : \mathbb{R}^d \rightarrow \mathbb{R}$ for $i \in [n]$, it holds that

$$\sum_{i=1}^n \mathbb{E}_k [h_i(\rho_i^{(k+1)} \mathbf{g}_i^{(k+1)})] = \frac{1}{n} \sum_{i=1}^n h_i(q_i^{(k)} \nabla \ell_i(\boldsymbol{\theta}^{(k)})) + \left(1 - \frac{1}{n}\right) \sum_{i=1}^n h_i(\rho_i^{(k)} \mathbf{g}_i^{(k)}).$$

After instantiating updates that satisfy these conditions, the implementation of the algorithm may be optimized to ensure an $O(n + d)$ per-iteration time complexity and an $O(nd)$ (or even $O(n + d)$) space complexity. We comment on two particular cases and proceed with the analysis in Section 2.6. Instantiating an algorithm of the form Algorithm 1 involves describing the bias and variance reduction table updates, and then “streamlining” the implementation to ensure that the time complexity of each update is $O(n + d)$. This also depends on which method is used in the computation and recomputation of the weights in line 12.

We present two examples of instantiations of Algorithm 1, and prove that they satisfy Condition 2.5.1 and Condition 2.5.2. In doing so, we gain intuition as to what the role of the tables $\{\hat{\boldsymbol{\theta}}_i^{(k)}\}_{i=1}^n$ and $\{\tilde{\boldsymbol{\theta}}_i^{(k)}\}_{i=1}^n$, that is to contain examples that approximately track $\boldsymbol{\theta}^{(k)}$ in expectation along the course of the optimization.

Example 2.5.1 (Prospect-Style Bias Reduction). In Mehta et al. [2024b], the table $\hat{\boldsymbol{\theta}}^{(k)}$ is updated by drawing an uniformly random index j_k (independent of i_k) and applying

$$\hat{\boldsymbol{\theta}}_j^{(k+1)} = \begin{cases} \boldsymbol{\theta}^{(k)} & \text{if } j = j_k \\ \hat{\boldsymbol{\theta}}_j^{(k)} & \text{if } j \neq j_k \end{cases}.$$

We confirm that Condition 2.5.1 holds by first take the expected value with respect to j_k :

$$\begin{aligned} & \sum_{j=1}^n \mathbb{E}_k [h_j(\boldsymbol{\theta}^{(k+1)}, \hat{\boldsymbol{\theta}}_j^{(k+1)})] \\ &= \frac{1}{n} \sum_{j=1}^n \mathbb{E}_k [h_j(\boldsymbol{\theta}^{(k+1)}, \hat{\boldsymbol{\theta}}_j^{(k+1)}) \mid j_k = j] + \left(1 - \frac{1}{n}\right) \sum_{j=1}^n \mathbb{E}_k [h_j(\boldsymbol{\theta}^{(k+1)}, \hat{\boldsymbol{\theta}}_j^{(k+1)}) \mid j_k \neq j] \\ &= \frac{1}{n} \sum_{j=1}^n \mathbb{E}_k [h_j(\boldsymbol{\theta}^{(k+1)}, \boldsymbol{\theta}^{(k)})] + \left(1 - \frac{1}{n}\right) \sum_{j=1}^n \mathbb{E}_k [h_j(\boldsymbol{\theta}^{(k+1)}, \hat{\boldsymbol{\theta}}_j^{(k)})], \end{aligned}$$

as desired. To perform this update in practice, we need not store $\hat{\boldsymbol{\theta}}^{(k)}$ at all; instead, we store

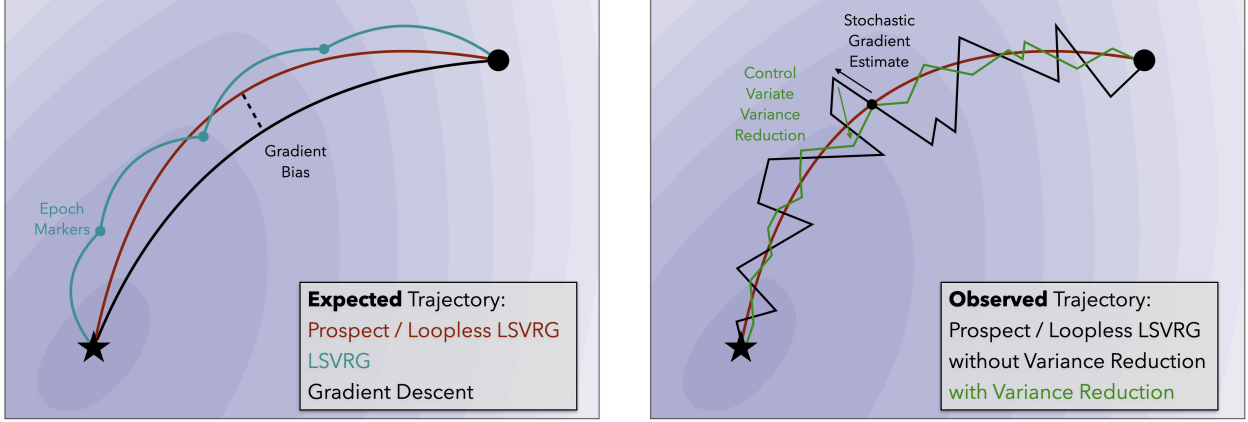


Figure 2.3: **Idealized Visualization of Bias- and Variance-Reduced Algorithms.** The Prospect [Mehta et al., 2024b], LSVRG [Mehta et al., 2023], and Loopless LSVRG (Example 2.5.2) algorithms are shown. **Left:** Expected trajectory over algorithmic randomness. This displays the gradient bias (as compared to full batch gradient descent of each method). LSVRG, unlike the updates in lines 14 and 14 for Algorithm 1, operates in epochs. **Right:** Observed trajectory of either Prospect/Loopless LSVRG. The variance reduction is interpreted as a control variate correction applied to the initial stochastic gradient estimate.

the $O(n)$ sized vector $\mathbf{l}^{(k)}$ and update element j_k on iteration t . The memory requirement for this method is then $O(n)$.

Example 2.5.2 (LSVRG-Style Variance Reduction). A loopless variance of LSVRG [Mehta et al., 2023] is described by the following update:

$$(\tilde{\boldsymbol{\theta}}_i^{(k+1)}, \rho_i^{(k+1)}) = \begin{cases} (\boldsymbol{\theta}^{(k)}, q_i^{(k+1)}) \forall i & \text{w.p. } \frac{1}{n} \\ (\tilde{\boldsymbol{\theta}}_i^{(k)}, \rho_i^{(k)}) \forall i & \text{w.p. } 1 - \frac{1}{n} \end{cases}.$$

It can be seen immediately that Condition 2.5.2 holds. Computationally, we see that $\tilde{\boldsymbol{\theta}}_1^{(k)} = \dots \tilde{\boldsymbol{\theta}}_n^{(k)}$ so we need only store a d -length “checkpoint” vector. Rather than storing $(g_1^{(k)}, \dots, g_n^{(k)})$, we may simply store $\boldsymbol{\rho}^{(k)}$ (the weights) and $\bar{\mathbf{g}}^{(k)} := \sum_{i=1}^n \rho_i^{(k)} \mathbf{g}_i^{(k)}$ (aggregation). When computing line 7, we may simply recompute $\nabla r_{i_k}(\tilde{\boldsymbol{\theta}}_{i_k}^{(k)})$ using the checkpoint at a constant additional cost. Thus, the memory requirement for this step is $O(n + d)$.

2.5.3 Main Results

The convergence analysis given in the next section provides two different convergence rates depending on the choice of the shift cost ν . As shown in (2.21) (Section 2.4), this quantity is inversely proportional to the smoothness constant of the objective, thus directly affecting the conditioning of the objective. When $\nu = \Omega(nG^2/\mu)$ (which we refer to as the *large cost* setting), we see that the constant $(nG^2/\nu + L) = O(\mu + L)$, which is indeed the smoothness constant of the empirical risk minimization (ERM) objective $\frac{1}{n} \sum_{i=1}^n \ell_i(\cdot) + \frac{\mu}{2} \|\cdot\|^2$. Thus, for large values of ν , the rate should be similar to ERM. However, the bias in the gradient estimate remains even when ν is large, which contributes an additional condition number factor

$$\kappa_{\mathcal{P}} = n \max_{i \in [n], q \in \mathcal{P}} q_i$$

in addition to the usual $\kappa = (L + \mu)/\mu$ condition number, and preventing achieving the $(n + \kappa) \ln(1/\varepsilon)$ convergence rate. Note that $\kappa_{\mathcal{P}} = 1$ when $\mathcal{P} = \{\mathbf{1}/n\}$, thus making the large shift cost result a proper generalization of the analysis of ERM. Note that due to strong convexity, there exists a unique minimizer

$$\boldsymbol{\theta}^* = \arg \min_{\boldsymbol{\theta} \in \mathbb{R}^d} \mathcal{L}(\boldsymbol{\theta}) \tag{2.22}$$

which is referenced in the upcoming results. We first present the large cost convergence rate.

Theorem 2.5.1. *Suppose the shift cost satisfies*

$$\nu \geq 8nG^2/\mu.$$

Then, the sequence of iterates produced by Algorithm 1 with $\eta = 1/(12(\mu + M)\kappa_{\mathcal{P}})$ achieves

$$\mathbb{E}_0 \|\boldsymbol{\theta}^{(t)} - \boldsymbol{\theta}^*\|_2^2 \leq (1 + \sigma_n^{-1} + \sigma_n^{-2}) \exp(-t/\tau) \|\boldsymbol{\theta}^{(0)} - \boldsymbol{\theta}^*\|_2^2.$$

with

$$\tau = 2 \max\{n, 24\kappa_{\mathcal{P}}(\kappa + 1)\}.$$

To describe the rate for small costs, define $\delta := nG^2/(\mu\nu)$. The quantity δ captures the effect of the primal regularizer μ and dual regularizer ν as compared to the inherent continuity properties of the underlying losses.

Theorem 2.5.2. *Assume that $n \geq 2$ and that the shift cost $\nu \leq 8nG^2/\mu$. The sequence of iterates produced by Algorithm 1 with*

$$\eta = \frac{1}{16n\mu} \min \left\{ \frac{1}{6[8\delta + (\kappa + 1)\kappa_{\mathcal{P}}]}, \frac{1}{4\delta^2 \max\{2n\kappa^2, \delta\}} \right\}$$

achieves

$$\mathbb{E}_0 \|\boldsymbol{\theta}^{(t)} - \boldsymbol{\theta}^*\|_2^2 \leq \left(5 + 16\delta + \frac{6\kappa^2}{\sigma_n} \right) \exp(-t/\tau) \|\boldsymbol{\theta}^{(0)} - \boldsymbol{\theta}^*\|_2^2$$

for

$$\tau = 32n \max \left\{ 6[8\delta + (\kappa + 1)\kappa_{\mathcal{P}}], 4\delta^2 \max\{2n\kappa^2, \delta\}, 1/16 \right\}.$$

Clearly, the small cost rate is a high-degree polynomial in both n and κ , and does not reduce to the ERM rate. One intuition behind this is that the updates for the sequence $(\mathbf{q}^{(k)})_{t \geq 1}$ cannot be obviously cast as gradient descent steps or even non-Euclidean proximal steps; they are full maximization steps without any regularization to ensure that $\mathbf{q}^{(k+1)}$ is close to $\mathbf{q}^{(k)}$. One benefit of full maximization, however, is that there is no additional dual learning rate to tune.

Note that the guarantees in Theorem 2.5.1 and Theorem 2.5.2 also provide guarantees on the primal-dual gap, which is a common convergence criterion for min-max problems. To describe this quantity, consider the saddle point objective

$$\mathcal{L}(\boldsymbol{\theta}, \mathbf{q}) = \langle \mathbf{q}, \ell(\boldsymbol{\theta}) \rangle - \nu D_f(\mathbf{q} \| \mathbf{1}/n) + \frac{\mu}{2} \|\boldsymbol{\theta}\|_2^2$$

so that $\mathcal{L}(\boldsymbol{\theta}) = \max_{\mathbf{q} \in \mathcal{Q}} \mathcal{L}(\boldsymbol{\theta}, \mathbf{q})$. Due to strong convexity and strong concavity, there exists a unique saddle point $(\boldsymbol{\theta}^*, \mathbf{q}^*)$ satisfying

$$\mathcal{L}(\boldsymbol{\theta}^*, \mathbf{q}) \leq \mathcal{L}(\boldsymbol{\theta}^*, \mathbf{q}^*) \leq \mathcal{L}(\boldsymbol{\theta}, \mathbf{q}^*)$$

for all $\boldsymbol{\theta} \in \mathbb{R}^d$ and $\mathbf{q} \in \mathcal{Q}$, where $\boldsymbol{\theta}^*$ coincides with (2.22). Then, we may the *primal-dual gap* at any point $(\boldsymbol{\theta}, \mathbf{q})$ as

$$\text{Gap}(\boldsymbol{\theta}, \mathbf{q}) = \mathcal{L}(\boldsymbol{\theta}, \mathbf{q}^*) - \mathcal{L}(\boldsymbol{\theta}^*, \mathbf{q}) \geq 0, \quad (2.23)$$

which is uniquely minimized at $(\boldsymbol{\theta}^*, \mathbf{q}^*)$. The following corollary can convert the distance-to-opt bounds into a bound on the primal-dual gap in the case of a smooth dual.

Corollary 2.5.1. *Assume that the rescaled χ^2 -divergence $D_f(\mathbf{q} \parallel \mathbf{1}/n) = \frac{1}{2} \|\mathbf{q} - \mathbf{1}/n\|_2^2$ is used for the shift penalty. Then,*

$$\frac{\mu}{2} \|\boldsymbol{\theta}^{(k)} - \boldsymbol{\theta}^*\|_2^2 \leq \text{Gap}(\boldsymbol{\theta}^{(k)}, q^{\text{opt}}(\ell(\boldsymbol{\theta}^{(k)}))) \leq \frac{1}{2} \left(L + \mu + \frac{nG^2}{\nu} \right) \|\boldsymbol{\theta}^{(k)} - \boldsymbol{\theta}^*\|_2^2.$$

Thus, the expected primal-dual gap $\mathbb{E}_k[\text{Gap}(\boldsymbol{\theta}^{(k)}, q^{\text{opt}}(\ell(\boldsymbol{\theta}^{(k)})))]$ achieves the same convergence rate as $\mathbb{E}_k \|\boldsymbol{\theta}^{(k)} - \boldsymbol{\theta}^*\|_2^2$ under the conditions of Theorem 2.5.1 and Theorem 2.5.2, respectively.

Proof. The gap criterion decomposes as

$$\text{Gap}(\boldsymbol{\theta}, \mathbf{q}) = \underbrace{\mathcal{L}(\boldsymbol{\theta}, \mathbf{q}^*) - \mathcal{L}(\boldsymbol{\theta}^*, \mathbf{q}^*)}_{\geq \frac{\mu}{2} \|\boldsymbol{\theta} - \boldsymbol{\theta}^*\|_2^2} + \underbrace{\mathcal{L}(\boldsymbol{\theta}^*, \mathbf{q}^*) - \mathcal{L}(\boldsymbol{\theta}^*, \mathbf{q})}_{\geq \frac{\nu}{2} \|\mathbf{q} - \mathbf{q}^*\|_2^2},$$

where the inequalities follow by strong convexity in $\boldsymbol{\theta}$ and \mathbf{q} , respectively. This gives the lower bound by using that $\|q^{\text{opt}}(\ell(\boldsymbol{\theta}^{(k)})) - \mathbf{q}^*\|_2^2 \geq 0$. The function $\boldsymbol{\theta} \mapsto \mathcal{L}(\boldsymbol{\theta}, \mathbf{q}^*) - \mathcal{L}(\boldsymbol{\theta}^*, \mathbf{q}^*)$ is $(L + \mu)$ -smooth and is minimized at $\boldsymbol{\theta}^*$, indicating that for any $\boldsymbol{\theta} \in \mathbb{R}^d$,

$$\mathcal{L}(\boldsymbol{\theta}, \mathbf{q}^*) - \mathcal{L}(\boldsymbol{\theta}^*, \mathbf{q}^*) \leq \frac{L + \mu}{2} \|\boldsymbol{\theta} - \boldsymbol{\theta}^*\|_2^2.$$

Similarly, the function $\mathbf{q} \mapsto \mathcal{L}(\boldsymbol{\theta}^*, \mathbf{q}^*) - \mathcal{L}(\boldsymbol{\theta}^*, \mathbf{q})$ is ν -smooth and minimized at \mathbf{q}^* , which

gives that for any $\mathbf{q} \in \mathcal{Q}$

$$\mathcal{L}(\boldsymbol{\theta}^*, \mathbf{q}^*) - \mathcal{L}(\boldsymbol{\theta}^*, \mathbf{q}) \leq \frac{\nu}{2} \|\mathbf{q} - \mathbf{q}^*\|_2^2.$$

Then, we have by $(1/\nu)$ -Lipschitzness (w.r.t. $\|\cdot\|_2$) of q^{opt} and G -Lipschitzness of each ℓ_i that

$$\begin{aligned} \|q^{\text{opt}}(\ell(\boldsymbol{\theta}^{(k)})) - \mathbf{q}^*\|_2^2 &= \|q^{\text{opt}}(\ell(\boldsymbol{\theta}^{(k)})) - q^{\text{opt}}(\ell(\boldsymbol{\theta}^*))\|_2^2 \\ &\leq \frac{1}{\nu^2} \|\ell(\boldsymbol{\theta}^{(k)}) - \ell(\boldsymbol{\theta}^*)\|_2^2 \\ &\leq \frac{nG^2}{\nu^2} \|\boldsymbol{\theta}^{(k)} - \boldsymbol{\theta}^*\|_2^2. \end{aligned}$$

Adding together both bounds gives the desired result. \square

Note that the primal-dual gap, while a convenient theoretical criterion, is still not computable by the algorithm in practice and thus cannot be used as a stopping criterion. However, it can be upper bounded by computable quantities. In Section 3.7.1, we derive an upper bound for a more general convergence criterion, called the *smoothed duality gap*, which contains (2.23) as a special case. In software implementations, we use the smoothed duality gap as an online certificate, or a measure of suboptimality that can be computed by the algorithm as it runs to determine a finite termination point.

The next subsection contains the convergence analysis that yields the two theorems above. The analysis of both methods will follow similarly, but differ in one key step: their usage of a generic bias bound (namely, the upcoming Proposition 2.6.1).

2.6 Convergence Analysis

We first give the broad outline of the analysis. Some key steps are proved in the main text, whereas the remaining proofs can be found in Appendix A.2. The optimum of (2.18) is denoted $\boldsymbol{\theta}^*$ and satisfies

$$\nabla(\mathbf{q}^{*\top} r(\boldsymbol{\theta}^*)) = 0, \text{ for } \mathbf{q}^* = \arg \max_{\mathbf{q} \in \mathcal{Q}} \langle \mathbf{q}, \ell(\boldsymbol{\theta}^*) \rangle - \nu \text{Reg}(\mathbf{q}). \quad (2.24)$$

The optimum $\boldsymbol{\theta}^*$ (or equivalently the saddle point $(\boldsymbol{\theta}^*, \mathbf{q}^*)$) exists due to the respective strong convexity and strong concavity of the objective and because $\mu, \nu > 0$.

When analyzing stochastic gradient methods in the smooth and strongly convex setting, we typically expand

$$\mathbb{E}_k \|\boldsymbol{\theta}^{(k+1)} - \boldsymbol{\theta}^*\|_2^2 = \|\boldsymbol{\theta}^{(k)} - \boldsymbol{\theta}^*\|_2^2 - \underbrace{2\eta \langle \mathbb{E}_k[\mathbf{v}^{(k)}], \boldsymbol{\theta}^{(k)} - \boldsymbol{\theta}^* \rangle}_{\text{descent term}} + \underbrace{\eta^2 \mathbb{E}_k \|\mathbf{v}^{(k)}\|_2^2}_{\text{noise term}}. \quad (2.25)$$

First, note that the expectation of the primal gradient estimate $\mathbf{v}^{(k)}$ is $\nabla r(\boldsymbol{\theta}^{(k)})^\top \mathbf{q}^{(k)}$, where $\mathbf{q}^{(k)} = \mathbf{q}^{\text{opt}}(\mathbf{l}^{(k)})$ and $\mathbf{l}^{(k)} \in \mathbb{R}^n$ denotes the estimate of the full loss vector. Applying standard convex inequalities to the descent term (in particular, Theorem A.1.3) yields

$$\begin{aligned} -\langle \nabla r(\boldsymbol{\theta}^{(k)})^\top \mathbf{q}^{(k)}, \boldsymbol{\theta}^{(k)} - \boldsymbol{\theta}^* \rangle &\leq -\frac{\mu M}{\mu + M} \|\boldsymbol{\theta} - \boldsymbol{\theta}^*\|_2^2 \\ &\quad - \frac{1}{\mu + L} \sum_{i=1}^n q_i^{(k)} \|\nabla r_i(\boldsymbol{\theta}^{(k)}) - \nabla r_i(\boldsymbol{\theta}^*)\|_2^2 \\ &\quad - \langle \nabla r(\boldsymbol{\theta}^*)^\top \mathbf{q}^{(k)}, \boldsymbol{\theta}^{(k)} - \boldsymbol{\theta}^* \rangle. \end{aligned}$$

The first two terms on the right-hand side are negative, which provides a decrease in the expected distance-to-optimum value on every iterate. In the empirical risk minimization setting, the final term on the right-hand side would be zero due to the first-order optimality conditions on $\boldsymbol{\theta}^*$, as $\mathbf{q}^{(k)} = \mathbf{1}/n$, implying the decrease of $\mathbb{E}_k \|\boldsymbol{\theta}^{(k+1)} - \boldsymbol{\theta}^*\|_2^2$ for small enough η . However, because $\mathbf{q}^{(k)}$ is a potentially non-uniform vector estimated using the table of losses $\mathbf{l}^{(k)}$ (as opposed to the loss vector $\ell(\boldsymbol{\theta}^*)$ at optimum), the term $-\langle \nabla r(\boldsymbol{\theta}^*)^\top \mathbf{q}^{(k)}, \boldsymbol{\theta}^{(k)} - \boldsymbol{\theta}^* \rangle$ is non-zero. Additionally, this term is multiplied only by the learning rate η , instead of the noise terms, which are multiplied by η^2 . Thus, this bias term must be bounded carefully in order to achieve the convergence guarantee under this regime. The convergence analysis will proceed in three parts.

1. First, we separately upper bound the descent and noise terms appearing in (2.25). In doing so, we will introduce four non-negative terms, denoted $S^{(k)}$, $T^{(k)}$, $U^{(k)}$, and $R^{(k)}$. These will all be incorporated into a *Lyapunov function* $V^{(k)}$, given by the “main term”

$\|\boldsymbol{\theta}^{(k)} - \boldsymbol{\theta}^*\|_2^2$ and a weighted combination of the other terms, namely,

$$V^{(k)} = \|\boldsymbol{\theta}^{(k)} - \boldsymbol{\theta}^*\|_2^2 + c_1 \textcolor{green}{S}^{(k)} + c_2 \textcolor{red}{T}^{(k)} + c_3 \textcolor{blue}{U}^{(k)} + c_4 \textcolor{brown}{R}^{(k)}. \quad (2.26)$$

The colors are used to track the quantities easily across lemmas.

2. We then study the per-iterate evolution of the additional Lyapunov terms $S^{(k)}$, $T^{(k)}$, $U^{(k)}$, and $R^{(k)}$. We note that in the next step, for certain values of the parameter ν , we need only use $S^{(k)}$ and $T^{(k)}$ in the Lyapunov function.
3. We combine the upper bounds from the previous steps and set the constants (c_1, c_2, c_3, c_4) to achieve a single-iterate progress bound of the form

$$\mathbb{E}_k [V^{(k+1)}] \leq (1 - \tau^{-1})V^{(k)}$$

for $\tau > 1$, which is interpreted as a “half-life” parameter. The final rate will depend on an additional quantity: the condition number $\kappa_{\mathcal{Q}} := n \max_i \max_{q \in \mathcal{Q}} q_i$, which is equal to 1 in the empirical risk minimization setting.

When ν is large, the analysis is simplified significantly by setting $c_3 = c_4 = 0$. However, to achieve unconditional linear convergence of the algorithm under Assumption 2.4.1 using this technique, the terms $U^{(k)}$ and $R^{(k)}$ are needed. We emphasize that this proof technique applies to a class of algorithms, namely those that satisfy Condition 2.5.1 and Condition 2.5.2. The terms are defined as

$$\begin{aligned} \textcolor{green}{S}^{(k)} &= \frac{1}{n} \sum_{i=1}^n \|n\rho_i^{(k)} \nabla r_i(\tilde{\boldsymbol{\theta}}_i^{(k)}) - nq_i^* \nabla r_i(\boldsymbol{\theta}^*)\|_2^2, & \textcolor{red}{T}^{(k)} &= \sum_{i=1}^n \|\hat{\boldsymbol{\theta}}_i^{(k)} - \boldsymbol{\theta}^*\|_2^2, \\ \textcolor{blue}{U}^{(k)} &= \frac{1}{n} \sum_{j=1}^n \|\boldsymbol{\theta}^{(k)} - \hat{\boldsymbol{\theta}}_j^{(k)}\|_2^2, & \textcolor{brown}{R}^{(k)} &= 2\eta n (\textcolor{brown}{q}^{(k)} - \textcolor{brown}{q}^*)^\top (\textcolor{brown}{l}^{(k)} - \textcolor{brown}{l}^*). \end{aligned}$$

Though not included in the Lyapunov function, we will also introduce

$$\textcolor{violet}{Q}^{(k)} = \frac{1}{n} \sum_{i=1}^n \|nq_i^{(k)} \nabla r_i(\boldsymbol{\theta}^{(k)}) - nq_i^* \nabla r_i(\boldsymbol{\theta}^*)\|_2^2,$$

which appears as a term to be canceled. The three steps in the outline comprise the next three subsections. We start with the “main” Lyapunov term $\|\boldsymbol{\theta}^{(k)} - \boldsymbol{\theta}^\star\|_2^2$, whose upper bound will expose the other terms.

2.6.1 Bounding the Distance-to-Optimum

Because the main innovation in this proof technique is the bounding of the descent term with bias, we highlight this step in a result of independent interest.

Proposition 2.6.1 (Bias Bound). *Consider any $\boldsymbol{\theta} \in \mathbb{R}^d$, $\mathbf{l} \in \mathbb{R}^n$, and $\bar{\mathbf{q}} \in \mathcal{Q}$. Define*

$$\mathbf{q} := q^{\text{opt}}(\mathbf{l}) = \arg \max_{\mathbf{p} \in \mathcal{Q}} \langle \mathbf{p}, \mathbf{l} \rangle - \nu \text{Reg}(\mathbf{p}).$$

For any $\alpha_1 \in [0, 1]$,

$$\begin{aligned} & -(\nabla r(\boldsymbol{\theta})^\top \mathbf{q} - \nabla r(\boldsymbol{\theta}^\star)^\top \bar{\mathbf{q}})^\top (\boldsymbol{\theta} - \boldsymbol{\theta}^\star) \\ & \leq -(\mathbf{q} - \bar{\mathbf{q}})^\top (\ell(\boldsymbol{\theta}) - \ell(\boldsymbol{\theta}^\star)) - \frac{\mu}{2} \|\boldsymbol{\theta} - \boldsymbol{\theta}^\star\|_2^2 \\ & - \frac{\alpha_1}{4(M + \mu)\kappa_{\mathcal{Q}}} \frac{1}{n} \sum_{i=1}^n \|nq_i \nabla r_i(\boldsymbol{\theta}) - nq_i^\star \nabla r_i(\boldsymbol{\theta}^\star)\|_2^2 + \frac{2\alpha_1 G^2}{\nu(M + \mu)\kappa_{\mathcal{Q}}} n(\mathbf{q} - \mathbf{q}^\star)^\top (\mathbf{l} - \mathbf{l}^\star). \end{aligned}$$

Notice in Proposition 2.6.1 that when $\bar{\mathbf{q}} = \mathbf{q}^\star$, we have that $\nabla r(\boldsymbol{\theta}^\star)^\top \bar{\mathbf{q}}$ vanishes, as $\mathbf{q}^\star = q^{\text{opt}}(\ell(\boldsymbol{\theta}^\star))$. Armed with Proposition 2.6.1 and a relatively simple bound on the noise, we can upper bound the distance-to-optimum.

Lemma 2.6.1 (Analysis of distance-to-optimum term). *For any constants $\alpha_1 \in [0, 1]$ and*

$\alpha_2 > 0$, and any $\bar{\mathbf{q}} \in \mathcal{Q}$, we have that

$$\begin{aligned} \mathbb{E}_k \|\boldsymbol{\theta}^{(k+1)} - \boldsymbol{\theta}^*\|_2^2 &\leq (1 - \eta\mu) \|\boldsymbol{\theta}^{(k)} - \boldsymbol{\theta}^*\|_2^2 \\ &\quad - 2\eta(\boldsymbol{\theta}^{(k)} - \boldsymbol{\theta}^*)^\top \nabla r(\boldsymbol{\theta}^*) \bar{\mathbf{q}} \\ &\quad - \eta \left(\frac{\alpha_1}{2(M + \mu)\kappa_{\mathcal{Q}}} - \eta(1 + \alpha_2) \right) \mathcal{Q}^{(k)} + \eta^2(1 + \alpha_2^{-1}) \mathcal{S}^{(k)} \\ &\quad + \frac{2\alpha_1 G^2}{\nu(M + \mu)\kappa_{\mathcal{Q}}} \mathcal{R}^{(k)} - 2\eta(\mathbf{q}^{(k)} - \bar{\mathbf{q}})^\top (\ell(\boldsymbol{\theta}) - \ell(\boldsymbol{\theta}^*)). \end{aligned}$$

Proof. Recall the expansion given in (2.25), which is:

$$\mathbb{E}_k \|\boldsymbol{\theta}^{(k+1)} - \boldsymbol{\theta}^*\|_2^2 = \|\boldsymbol{\theta}^{(k)} - \boldsymbol{\theta}^*\|_2^2 - 2\eta \langle \mathbb{E}_k[\mathbf{v}^{(k)}], \boldsymbol{\theta}^{(k)} - \boldsymbol{\theta}^* \rangle + \eta^2 \mathbb{E}_k \|\mathbf{v}^{(k)}\|_2^2. \quad (2.27)$$

Observe that

$$\mathbb{E}_k[\mathbf{v}^{(k)}] = \sum_{i=1}^n \mathbf{q}_i^{(k)} \nabla r(\boldsymbol{\theta}^{(k)}) = \nabla r(\boldsymbol{\theta}^{(k)})^\top \mathbf{q}^{(k)}$$

By Proposition 2.6.1 with $\mathbf{l} = \mathbf{l}^{(k)}$, $\mathbf{q} = \mathbf{q}^{(k)}$, $\boldsymbol{\theta} = \boldsymbol{\theta}^{(k)}$, and multiplying by 2η , we have that

$$\begin{aligned} &- 2\eta(\boldsymbol{\theta}^{(k)} - \boldsymbol{\theta}^*)^\top \nabla r(\boldsymbol{\theta}^{(k)})^\top \mathbf{q}^{(k)} \\ &\leq -2\eta(\boldsymbol{\theta}^{(k)} - \boldsymbol{\theta}^*)^\top \nabla r(\boldsymbol{\theta}^*) \bar{\mathbf{q}} - 2\eta(\mathbf{q}^{(k)} - \bar{\mathbf{q}})^\top (\ell(\boldsymbol{\theta}^{(k)}) - \ell(\boldsymbol{\theta}^*)) \\ &\quad - \mu\eta \|\boldsymbol{\theta}^{(k)} - \boldsymbol{\theta}^*\|_2^2 - \frac{\eta\alpha_1}{2(M + \mu)\kappa_{\mathcal{Q}}} \mathcal{Q}^{(k)} + \frac{2\alpha_1 G^2}{\nu(M + \mu)\kappa_{\mathcal{Q}}} \mathcal{R}^{(k)}. \end{aligned}$$

The noise term is bounded by applying Young's inequality with parameter $\alpha_2 > 0$ and the identity $\mathbb{E}\|X - \mathbb{E}[X]\|_2^2 = \mathbb{E}\|X\|_2^2 - \|\mathbb{E}[X]\|_2^2$ (i.e., the variance is upper bounded by the

second moment). Noting that $\sum_{i=1}^n q_i^* \nabla r_i(\boldsymbol{\theta}^*) = 0$, we get

$$\begin{aligned}
& \mathbb{E}_k [\|\mathbf{v}^{(k)}\|_2^2] \\
&= \mathbb{E}_k [\|\mathbf{v}^{(k)} - \nabla(\mathbf{q}^{*\top} r)(\boldsymbol{\theta}^*)\|_2^2] \\
&= \mathbb{E}_k \left[\|nq_{i_k}^{(k)} \nabla r_{i_k}(\boldsymbol{\theta}^{(k)}) - nq_{i_k}^* \nabla r_{i_k}(\boldsymbol{\theta}^*) \right. \\
&\quad \left. + nq_{i_k}^* \nabla r_{i_k}(\boldsymbol{\theta}^*) - n\rho_{i_k}^{(k)} \nabla r_{i_k}(\tilde{\boldsymbol{\theta}}_{i_k}^{(k)}) - (\nabla(\mathbf{q}^{*\top} r)(\boldsymbol{\theta}^*) - \sum_{i=1}^n \rho_i^{(k)} \nabla r_i(\tilde{\boldsymbol{\theta}}_i^{(k)}))\|_2^2 \right] \\
&\leq (1 + \alpha_2) \mathbb{E}_k [\|nq_{i_k}^{(k)} \nabla r_{i_k}(\boldsymbol{\theta}^{(k)}) - nq_{i_k}^* \nabla r_{i_k}(\boldsymbol{\theta}^*)\|^2] \\
&\quad + (1 + \alpha_2^{-1}) \mathbb{E}_k \left[\|nq_{i_k}^* \nabla r_{i_k}(\boldsymbol{\theta}^*) - n\rho_{i_k}^{(k)} \nabla r_{i_k}(\tilde{\boldsymbol{\theta}}_{i_k}^{(k)}) - (\nabla(\mathbf{q}^{*\top} r)(\boldsymbol{\theta}^*) - \sum_{i=1}^n \rho_i^{(k)} \nabla r_i(\tilde{\boldsymbol{\theta}}_i^{(k)}))\|_2^2 \right] \\
&\leq (1 + \alpha_2) \mathbb{E}_k [\|nq_{i_k}^{(k)} \nabla r_{i_k}(\boldsymbol{\theta}^{(k)}) - nq_{i_k}^* \nabla r_{i_k}(\boldsymbol{\theta}^*)\|^2] \\
&\quad + (1 + \alpha_2^{-1}) \mathbb{E}_k [\|nq_{i_k}^* \nabla r_{i_k}(\boldsymbol{\theta}^*) - n\rho_{i_k}^{(k)} \nabla r_{i_k}(\tilde{\boldsymbol{\theta}}_{i_k}^{(k)})\|_2^2].
\end{aligned}$$

Substituting the definitions of $Q^{(k)}$ and $S^{(k)}$ we have finally that

$$\eta^2 \mathbb{E}_k \|\mathbf{v}^{(k)}\|_2^2 \leq \eta^2 (1 + \alpha_2) Q^{(k)} + \eta^2 (1 + \alpha_2^{-1}) S^{(k)}.$$

Combine the two displays above to get the desired result. \square

2.6.2 Bounding the Evolution of the Lyapunov Function Terms

Notice that the steps above did not depend on the updates to the bias and variance reduction tables. When bounding the evolution of the Lyapunov function terms, we see the conditions being used. We display the proof of Lemma 2.6.2 as an example of the conditions being used, and defer the proofs of the more technical Lemma 2.6.3 and Lemma 2.6.4 to Appendix A.2. Recall that

$$Q^{(k)} = \frac{1}{n} \sum_{i=1}^n \|nq_i^{(k)} \nabla r_i(\boldsymbol{\theta}^{(k)}) - nq_i^* \nabla r_i(\boldsymbol{\theta}^*)\|_2^2.$$

Lemma 2.6.2. *Given Condition 2.5.1 and Condition 2.5.2, we have that*

$$\begin{aligned}\mathbb{E}_k[S^{(k+1)}] &= \frac{1}{n}Q^{(k)} + \left(1 - \frac{1}{n}\right)S^{(k)}, \\ \mathbb{E}_k[T^{(k+1)}] &= \|\boldsymbol{\theta}^{(k)} - \boldsymbol{\theta}^*\|_2^2 + \left(1 - \frac{1}{n}\right)T^{(k)}.\end{aligned}$$

Proof. First, apply Condition 2.5.2 using the collection of functions

$$h_i(\mathbf{x}) := \frac{1}{n}\|n\mathbf{x} - nq_i^*\nabla r_i(\boldsymbol{\theta}^*)\|_2^2$$

in the (*) line below to achieve

$$\begin{aligned}\mathbb{E}_k[S^{(k+1)}] &= \frac{1}{n} \sum_{i=1}^n \mathbb{E}_k \left[\|n\rho_i^{(k+1)}\nabla r_i(\tilde{\boldsymbol{\theta}}_i^{(k+1)}) - nq_i^*\nabla r_i(\boldsymbol{\theta}^*)\|_2^2 \right] \\ &\stackrel{(*)}{=} \frac{1}{n} \sum_{i=1}^n \left[\frac{1}{n} \|nq_i^{(k)}\nabla r_i(\boldsymbol{\theta}^{(k)}) - nq_i^*\nabla r_i(\boldsymbol{\theta}^*)\|_2^2 + \left(1 - \frac{1}{n}\right) \|n\rho_i^{(k)}\nabla r_{i_k}(\tilde{\boldsymbol{\theta}}_i^{(k)}) - nq_i^*\nabla r_i(\boldsymbol{\theta}^*)\|_2^2 \right] \\ &= \frac{1}{n}Q^{(k)} + \left(1 - \frac{1}{n}\right)S^{(k)}.\end{aligned}$$

Similarly, using Condition 2.5.1 with the functions

$$h_i(\mathbf{u}, \mathbf{x}) := \|\mathbf{x} - \boldsymbol{\theta}^*\|_2^2$$

in the (o) line below to achieve

$$\begin{aligned}\mathbb{E}_k[T^{(k+1)}] &= \sum_{i=1}^n \mathbb{E}_k \left[\|\hat{\boldsymbol{\theta}}_i^{(k+1)} - \boldsymbol{\theta}^*\|_2^2 \right] \\ &\stackrel{(o)}{=} \|\boldsymbol{\theta}^{(k)} - \boldsymbol{\theta}^*\|_2^2 + \left(1 - \frac{1}{n}\right) \sum_{i=1}^n \|\hat{\boldsymbol{\theta}}_i^{(k)} - \boldsymbol{\theta}^*\|_2^2 \\ &= \|\boldsymbol{\theta}^{(k)} - \boldsymbol{\theta}^*\|_2^2 + \left(1 - \frac{1}{n}\right)T^{(k)},\end{aligned}$$

completing the proof. □

The remaining two lemmas give the upper bounds for $U^{(k)}$ and $R^{(k)}$.

Lemma 2.6.3. *For any value of $\alpha_2 > 0$, we have that*

$$\begin{aligned} \mathbb{E}_k [U^{(k+1)}] &\leq \eta^2(1 + \alpha_2)Q^{(k)} + \eta^2(1 + \alpha_2^{-1})S^{(k)} \\ &\quad + \frac{\eta M^2}{\mu n} \left(1 - \frac{1}{n}\right) T^{(k)} + \left(1 - \frac{1}{n}\right) \frac{G^2}{2\nu\mu n} R^{(k)} + \left(1 - \frac{1}{n}\right) U^{(k)}. \end{aligned}$$

Lemma 2.6.4. *For any $\alpha_3 > 0$, it holds that*

$$\begin{aligned} \mathbb{E}_k [R^{(k+1)}] &\leq 2\eta(\mathbf{q}^{(k)} - \mathbf{q}^*)^\top (\ell(\boldsymbol{\theta}^{(k)}) - \mathbf{l}^*) + \left(1 - \frac{1}{n}\right) R^{(k)} \\ &\quad + \frac{\eta G^2 n}{2\nu} \alpha_3^{-1} T^{(k)} + \frac{2\eta G^2 n}{\nu} (1 + \alpha_3) U^{(k)}. \end{aligned}$$

2.6.3 Tuning Constants and Final Rate

Finally, we combine all previous upper bounds into an upper bound on $\mathbb{E}_k [V^{(k+1)}]$ in terms of $V^{(k)}$. For Theorem 2.5.1, we may in fact set $c_3 = c_4 = 0$, simplifying the analysis significantly. Of the remaining two terms $S^{(k)}$ and $T^{(k)}$, $S^{(k)}$ is particularly similar to the (only) Lyapunov term in the simplest analysis of the SAGA algorithm (see Bach [2024, Proposition 5.9]). However, in the case of SAGA for ERM, the gradient estimate remains unbiased. In our case, the additional bias requires us to incorporate $T^{(k)}$, which compares the elements in the table used to estimate the dual variables with those at the optimum. For Theorem 2.5.2, the additional terms are necessary to achieve the suboptimal convergence rate, but without any conditions on ν . One weakness of the analysis is the attempt at using a primal-only strategy (analyzing $\|\boldsymbol{\theta}^{(k+1)} - \boldsymbol{\theta}^*\|_2^2$) in a primal-dual setting. In Chapter 3, we employ a strategy based on a primal-dual gap criterion which yields optimal rates in various problem classes. For the remainder of the chapter, we focus on practical questions regarding DRO, such as selecting the uncertainty set and shift cost, solving the maximization problem, and extensions to modern machine learning models such as neural networks.

2.6.4 Adaptive Sampling and Non-Uniformity

In Chapter 3, we expand on the optimization aspects of the problem and target a problem class that is more general than DRO. One of the essential ideas will be *adaptive sampling*, which we briefly motivate using the analysis above. For simplicity, consider the non-smooth variant of the objective (2.17), which can be written in the form of a maximum expectation

$$\max_{\mathbf{q} \in \mathcal{Q}} \sum_{i=1}^n q_i \ell(\boldsymbol{\theta}, \xi_i) = \max_{\mathbf{q} \in \mathcal{Q}} \mathbb{E}_{i \sim \text{Unif}[n]} [n q_i \ell(\boldsymbol{\theta}, \xi_i)].$$

We view $n q_i = \beta(\xi_i)$ as the Radon-Nikodym derivative as alluded to before. Technically, the same objective can be considered by considering a weight vector $\mathbf{w} = (w_1, \dots, w_n)$ satisfying $\mathbf{w} > 0$ and $\sum_{i=1}^n w_i = 1$, and writing

$$\max_{\mathbf{q} \in \mathcal{Q}} \mathbb{E}_{i \sim \text{Unif}[n]} [n q_i \ell(\boldsymbol{\theta}, \xi_i)] = \max_{\mathbf{q} \in \mathcal{Q}} \mathbb{E}_{i \sim \mathbf{w}} \left[q_i \frac{\ell(\boldsymbol{\theta}, \xi_i)}{w_i} \right].$$

Minimizing either stochastic formulations in $\boldsymbol{\theta}$ yields the same solution. Furthermore, the stochastic gradient estimate $\mathbf{v}^{(k)}$ from line 7 of Algorithm 1 will retain the same bias properties if $i_k \sim \mathbf{w}$ and if we replace r_i with $r_i/(n w_i)$. While the analysis of the Lyapunov function terms will change due to the change in the sampling scheme, to give intuition as to why a non-uniform sampling scheme can yield complexity improvements, notice that the results in Theorem 2.5.1 and Theorem 2.5.2 depend on the maximum Lipschitz constant G and the maximum smoothness constant L . That is, if each ℓ_i is G_i -Lipschitz continuous and L_i -smooth, then necessarily, we must have $G = \max_i G_i$ and $L = \max_i L_i$. On the other hand, the functions $\{\ell_i/(n w_i)\}_{i=1}^n$ would be governed by Lipschitz and smoothness constants

$$G(\mathbf{w}) = \max_{i=1, \dots, n} \frac{G_i}{n w_i} \text{ and } L(\mathbf{w}) = \max_{i=1, \dots, n} \frac{L_i}{n w_i}.$$

For example, by setting $w_i \propto G_i + L_i$, we get that $\max\{G(\mathbf{w}), L(\mathbf{w})\} \leq \frac{1}{n} \sum_{i=1}^n (G_i + L_i)$, which is the sum of *average* Lipschitz constants, which can be up to a factor n smaller than the maximum Lipschitz constants in the most non-uniform setting. One challenge that exists here, as opposed to a standard primal-only SGD analysis (e.g., as in Needell et al. [2014]),

is that the weights \mathbf{w} seemingly need to adapt to both Lipschitz constants and smoothness constants simultaneously. This will be handled using a primal-dual viewpoint in Chapter 3, in which we view G_1, \dots, G_n not as primal Lipschitz constants, but as *dual smoothness* constants.

2.7 Uncertainty Sets and Shift Costs

Regarding the uncertainty set \mathcal{Q} and shift cost ν , two concrete questions exist: 1) how should we select them, and 2) once selected, how do we compute (2.19) (i.e., solve the dual problem)? We consider the second question and then the first, for the two canonical examples of spectral risk measures, and then for divergence balls.

In both cases, we provide a duality result that converts the maximization into a minimization problem which can be solved by a near-exact iterative procedure. These duality relations also provide guidance on selecting the shift cost ν . As before, we will use the notation $\text{Reg}(\mathbf{q}) = D_f(q \| \mathbf{1}/n)$ to denote an f -divergence as described in Section 2.4. Additionally, while the duality results form the basis of the computational routines, we focus on mathematical details in this section and provide detailed implementation instructions in Appendix A.3.

2.7.1 Duality of Spectral Risk Measures

As in (2.16) we parametrize the uncertainty set as $\mathcal{Q} = \mathcal{Q}(\sigma)$ for spectrum σ .

Proposition 2.7.1. *Let $\mathbf{l} \in \mathbb{R}^n$ be a vector such that $l_1 \leq \dots \leq l_n$. We have the dual relation*

$$\max_{\mathbf{q} \in \mathcal{Q}(\sigma)} \{ \langle \mathbf{q}, \mathbf{l} \rangle - \nu \text{Reg}(\mathbf{q}) \} = \min_{\substack{\mathbf{z} \in \mathbb{R}^n \\ z_1 \leq \dots \leq z_n}} \sum_{i=1}^n \sigma_i z_i + \frac{1}{n} f^* \left(\frac{l_i - z_i}{\nu} \right). \quad (2.28)$$

The optima of both problems, denoted

$$z^{\text{opt}}(\mathbf{l}) = \arg \min_{\substack{\mathbf{z} \in \mathbb{R}^n \\ z_1 \leq \dots \leq z_n}} \sum_{i=1}^n \sigma_i z_i + \frac{1}{n} f^* \left(\frac{l_i - z_i}{\nu} \right), \quad q^{\text{opt}}(\mathbf{l}) = \arg \max_{\mathbf{q} \in \mathcal{Q}(\sigma)} \langle \mathbf{q}, \mathbf{l} \rangle - \nu \text{Reg}(\mathbf{q}),$$

are related as

$$q_i^{\text{opt}}(\mathbf{l}) = \frac{1}{n} [f^*]' \left(\frac{l_i - z_i^{\text{opt}}(\mathbf{l})}{\nu} \right). \quad (2.29)$$

Proof. Let $\iota_{\mathcal{Q}(\sigma)}$ denote the convex indicator function of the permutahedron $\mathcal{Q}(\sigma)$, which is 0 inside $\mathcal{Q}(\sigma)$ and $+\infty$ outside of $\mathcal{Q}(\sigma)$. Its convex conjugate is the support function of the permutahedron, i.e.,

$$\iota_{\mathcal{Q}(\sigma)}^*(\mathbf{l}) = \max_{\mathbf{q} \in \mathcal{Q}(\sigma)} \langle \mathbf{q}, \mathbf{l} \rangle.$$

For two closed convex functions h_1 and h_2 that are bounded from below, the convex conjugate of their sum is the infimal convolution of their conjugate [Hiriart-Urruty and Lemaréchal, 2004, Proposition 6.3.1]:

$$(h_1 + h_2)^*(\mathbf{x}) = \inf_{\mathbf{y} \in \mathbb{R}^d} \{h_1^*(\mathbf{y}) + h_2^*(\mathbf{x} - \mathbf{y})\}.$$

Provided that $h_1 + h_2$ is strictly convex, we have that the maximizer defining the conjugate is unique and equal to the gradient, that is,

$$\nabla(h_1 + h_2)^*(\mathbf{x}) = \arg \max_{\mathbf{z} \in \mathbb{R}^d} \{\langle \mathbf{z}, \mathbf{x} \rangle - (h_1 + h_2)(\mathbf{z})\}.$$

If, in addition, $h_1^* + h_2^*$ is strictly convex and h_2^* is differentiable, we have, by Danskin's theorem [Bertsekas, 1997],

$$\nabla(h_1 + h_2)^*(\mathbf{x}) = \nabla h_2^*(\mathbf{x} - \mathbf{y}^*(\mathbf{x})) \text{ for } \mathbf{y}^*(\mathbf{x}) = \arg \min_{\mathbf{y} \in \mathbb{R}^d} \{h_1^*(\mathbf{y}) + h_2^*(\mathbf{x} - \mathbf{y})\}.$$

Consider then $h_1(\mathbf{q}) = \iota_{\mathcal{Q}(\sigma)}(\mathbf{q})$ and $h_2(\mathbf{q}) = \text{Reg}(\mathbf{q}) := \nu D_f(q \| \mathbf{1}_n/n)$. Provided that f is strictly convex with f^* strictly convex, D_f is also strictly convex with D_f^* strictly convex

since D_f just decomposes as a sum of f on independent variables. We have then

$$\begin{aligned}
\sup_{\mathbf{q} \in \mathcal{Q}(\sigma)} \{ \langle \mathbf{q}, \mathbf{l} \rangle - \text{Reg}(\mathbf{q}) \} &= \sup_{\mathbf{q} \in \mathbb{R}^n} \{ \langle \mathbf{q}, \mathbf{l} \rangle - (\iota_{\mathcal{Q}(\sigma)}(\mathbf{q}) + \text{Reg}(\mathbf{q})) \} \\
&= (\iota_{\mathcal{Q}(\sigma)} + \text{Reg})^*(\mathbf{l}) \\
&= \inf_{\mathbf{y} \in \mathbb{R}^n} \{ \iota_{\mathcal{Q}(\sigma)}^*(\mathbf{y}) + \text{Reg}^*(\mathbf{l} - \mathbf{y}) \} \\
&= \inf_{\mathbf{y} \in \mathbb{R}^n} \left\{ \max_{\mathbf{q} \in \mathcal{Q}(\sigma)} \langle \mathbf{q}, \mathbf{y} \rangle + \text{Reg}^*(\mathbf{l} - \mathbf{y}) \right\} \\
&= \inf_{\mathbf{y} \in \mathbb{R}^n} \left\{ \sum_{i=1}^n \sigma_i y_{(i)} + \text{Reg}^*(\mathbf{l} - \mathbf{y}) \right\}, \tag{2.30}
\end{aligned}$$

where $y_{(1)} \leq \dots \leq y_{(n)}$ are the ordered values of $\mathbf{y} \in \mathbb{R}^n$. Moreover, we have that

$$\arg \max_{\mathbf{q} \in \mathcal{Q}(\sigma)} \{ \langle \mathbf{q}, \mathbf{l} \rangle - \text{Reg}(\mathbf{q}) \} = \nabla \text{Reg}^*(\mathbf{l} - \mathbf{y}^*(\mathbf{l})) \text{ for } \mathbf{y}^*(\mathbf{l}) = \arg \min_{\mathbf{y} \in \mathbb{R}^n} \left\{ \sum_{i=1}^n \sigma_i y_{(i)} + \text{Reg}^*(\mathbf{l} - \mathbf{y}) \right\}.$$

Since for any $\mathbf{x} \in \mathbb{R}^n$, Reg is decomposable into a sum of identical functions evaluated at the coordinates (x_1, \dots, x_n) , that is, $\text{Reg}(\mathbf{x}) = \sum_{i=1}^n \text{Reg}_i(x_i)$, its convex conjugate is $\text{Reg}^*(\mathbf{y}) = \sum_{i=1}^n \text{Reg}_i^*(y_i)$. In our case, $\text{Reg}_i(x_i) = \frac{\nu}{n} f(nx_i)$, so $\text{Reg}_i^*(y_i) = (\nu/n) f^*(y_i/\nu)$.

Next, by convexity of each Reg_i^* , we have that if for scalars l_i, l_j, y_i, y_j such that $l_i \leq l_j$ and $y_i \geq y_j$, then using Lemma A.1.2, we have that

$$\text{Reg}_i^*(l_i - y_i) + \text{Reg}_i^*(l_j - y_j) \geq \text{Reg}_i^*(l_i - y_j) + \text{Reg}_i^*(l_j - y_i).$$

Hence for \mathbf{y} to minimize $\text{Reg}^*(\mathbf{l} - \mathbf{y}) = \sum_{i=1}^n \text{Reg}_i^*(l_i - y_i)$, the coordinates of \mathbf{y} must be ordered as \mathbf{l} . That is, if π is an argsort for \mathbf{l} , s.t. $l_{\pi(1)} \leq \dots \leq l_{\pi(n)}$, then $y_{\pi(1)} \leq \dots \leq y_{\pi(n)}$. Since $\iota_{\mathcal{Q}(\sigma)}^*(\mathbf{y}) = \sum_{i=1}^n \sigma_i y_{(i)}$ does not depend on the ordering of y , the solution of (2.30) must also be ordered as \mathbf{l} such that the dual problem (2.30) can be written as

$$\inf_{\substack{\mathbf{y} \in \mathbb{R}^n \\ y_1 \leq \dots \leq y_n}} \sum_{i=1}^n \sigma_i y_i + \frac{1}{n} f^* \left(\frac{l_i - y_i}{\nu} \right) = \min_{\substack{\mathbf{z} \in \mathbb{R}^n \\ z_1 \leq \dots \leq z_n}} \sum_{i=1}^n \sigma_i z_i + \frac{1}{n} f^* \left(\frac{l_i - z_i}{\nu} \right).$$

By differentiating both sides of (2.28) with respect to l_i and applying Dankin's theorem [Bertsekas, 1997], we have the relation (2.29). \square

Computationally, we see that (2.28) is an instance of an isotonic regression problem, which can be solved by a version of the Pool Adjacent Violators (PAV) algorithm that is adapted to the generator function f . Thus, the dual problem can be solved exactly using an iterative algorithm that terminates in $O(n \log n)$ time.

2.7.2 Duality of f -Divergence Ball

We parametrize the uncertainty as $\mathcal{Q} = \mathcal{Q}(\rho)$ for radius ρ . Recall that Δ^{n-1} denotes the n -dimensional probability simplex. The feasible set can equivalently be written as

$$\mathcal{Q}(\rho) = \{\mathbf{q} \in \Delta^{n-1} : \text{Reg}(\mathbf{q}) \leq \rho\}. \quad (2.31)$$

Observe the following.

Proposition 2.7.2. *Let $\mathbf{l} \in \mathbb{R}^n$ be a vector. We have the dual relation*

$$\max_{\mathbf{q} \in \mathcal{Q}(\rho)} \{\langle \mathbf{q}, \mathbf{l} \rangle - \nu \text{Reg}(\mathbf{q})\} = \min_{t \geq \nu} \{t ((\text{Reg} + \iota_{\Delta^{n-1}})^* (\frac{\mathbf{l}}{t}) + \rho) - \rho \nu\}. \quad (2.32)$$

The optima of both problems, denoted

$$t^{\text{opt}}(\mathbf{l}) = \arg \min_{t \geq \nu} t ((\text{Reg} + \iota_{\Delta^{n-1}})^* (\frac{\mathbf{l}}{t}) + \rho), \quad q^{\text{opt}} = \arg \max_{\mathbf{q} \in \mathcal{Q}(\rho)} \langle \mathbf{q}, \mathbf{l} \rangle - \nu \text{Reg}(\mathbf{q}),$$

are related as

$$q_i^{\text{opt}}(\mathbf{l}) = \nabla (\text{Reg} + \iota_{\Delta^{n-1}})^* (\mathbf{l}/t^{\text{opt}}(\mathbf{l})) = \arg \max_{\mathbf{q} \in \Delta^{n-1}} \{\langle \mathbf{q}, \mathbf{l}/t^{\text{opt}}(\mathbf{l}) \rangle - \text{Reg}(\mathbf{q})\}. \quad (2.33)$$

Proof. First, using (2.31), we represent the left-hand side of (2.32) using a Lagrange multiplier λ applied to the extended real-valued objective

$$\begin{aligned} \max_{\mathbf{q} \in \mathcal{Q}(\rho)} \langle \mathbf{q}, \mathbf{l} \rangle - \nu \text{Reg}(\mathbf{q}) &= \max_{\mathbf{q} \in \Delta^{n-1}} \min_{\lambda \geq 0} \langle \mathbf{q}, \mathbf{l} \rangle - (\nu + \lambda) \text{Reg}(\mathbf{q}) + \lambda \rho \\ &= \min_{\lambda \geq 0} \max_{\mathbf{q} \in \Delta^{n-1}} \langle \mathbf{q}, \mathbf{l} \rangle - (\nu + \lambda) \text{Reg}(\mathbf{q}) + \lambda \rho \\ &= \min_{\lambda \geq 0} (\nu + \lambda) \left(\max_{\mathbf{q} \in \Delta^{n-1}} \langle \mathbf{q}, \mathbf{l}/(\nu + \lambda) \rangle - \text{Reg}(\mathbf{q}) \right) + \lambda \rho. \end{aligned}$$

The second equality follows because $\mathbf{1}/n \in \mathcal{Q}(\rho)$, so Slater's condition is satisfied and strong duality holds. We then reparametrize the problem with $t = \nu + \lambda$, so that

$$\begin{aligned} \max_{\mathbf{q} \in \mathcal{Q}(\rho)} \langle \mathbf{q}, \mathbf{l} \rangle - \nu \text{Reg}(\mathbf{q}) &= \min_{t \geq \nu} t \left(\max_{\mathbf{q} \in \Delta^{n-1}} \langle \mathbf{q}, \mathbf{l}/t \rangle - \text{Reg}(\mathbf{q}) \right) + \rho(t - \nu) \\ &= \min_{t \geq \nu} \left\{ t \left((\text{Reg} + \iota_{\Delta^{n-1}})^* \left(\frac{\mathbf{l}}{t} \right) + \rho \right) - \rho \nu \right\}, \end{aligned}$$

which gives (2.32). Next, we differentiate both sides of (2.32) with respect to l and apply Dankin's theorem [Bertsekas, 1997] to achieve (2.33). \square

To understand the computational implications of Proposition 2.7.2, we define the function

$$h_{\mathbf{l}}(t) = t \left((\text{Reg} + \iota_{\Delta^{n-1}})^* \left(\frac{\mathbf{l}}{t} \right) + \rho(t - \nu) \right),$$

and see that it is the sum of a perspective (in t) of a convex function on \mathbb{R}^n and a linear function in t , hence convex. Its derivative is then given by

$$h'_{\mathbf{l}}(t) = -\frac{1}{t} \underbrace{\left\langle \nabla (\text{Reg} + \iota_{\Delta^{n-1}})^* \left(\frac{\mathbf{l}}{t} \right), \frac{\mathbf{l}}{t} \right\rangle}_{\langle \mathbf{q}^*(t), \mathbf{l}/t \rangle} + \rho \text{ for } \mathbf{q}^*(t) = \arg \max_{\mathbf{q} \in \Delta^{n-1}} \langle \mathbf{q}, \mathbf{l}/t \rangle - \text{Reg}(\mathbf{q}). \quad (2.34)$$

Thus, if $h'_{\mathbf{l}}(t)$ is computable, then $t^{\text{opt}}(\mathbf{l})$ can be computed by binary search. Thus, the strength of the duality result lies on the ability to compute $h'_{\mathbf{l}}(t)$ (or equivalently, $\mathbf{q}^*(t)$).

The quintessential examples of this type of objective are when $\text{Reg}_{\text{KL}}(\mathbf{q}) = \sum_{i=1}^n q_i \log(nq_i)$ (for the KL divergence) and $\text{Reg}_{\chi^2}(\mathbf{q}) = n \|\mathbf{q} - \mathbf{1}/n\|_2^2$ (for the χ^2 -divergence). In this case, we have that

$$\mathbf{q}^*(t) = \arg \max_{\mathbf{q} \in \Delta^{n-1}} \langle \mathbf{q}, \mathbf{l}/t \rangle - \text{Reg}_{\text{KL}}(\mathbf{q}) = \frac{e^{\mathbf{l}/t}}{\sum_{i=1}^n e^{l_i/t}}$$

and

$$\mathbf{q}^*(t) = \arg \max_{\mathbf{q} \in \Delta^{n-1}} \langle \mathbf{q}, \mathbf{l}/t \rangle - \text{Reg}_{\chi^2}(\mathbf{q}) = \arg \min_{\mathbf{q} \in \Delta^{n-1}} \|\mathbf{q} - \mathbf{l}/(2nt)\|_2^2 = \text{proj}_{\Delta^{n-1}}(\mathbf{l}/(2nt)).$$

The first is a closed-form solution, whereas the second is a projection in ℓ_2 -norm onto the probability simplex, which is studied for instance in Condat [2016]. As we saw in Section 2.3,

spectral risk measures have many more constraints on the likelihood ratio than divergence balls. This appears in the case of empirical measures as well, as the dual form of the maximization problem has n Lagrange multipliers for spectral risk measures but only 1 for f -divergence balls.

2.7.3 Setting Problem Parameters

Maximum Uncertainty Set Size The maximum size is governed by the spectrum σ for spectral risk measures and the radius ρ for divergence balls. In either case, the criterion used to determine the maximum size is whether the uncertainty set constraint is inactive (as compared to the standard probability simplex constraint). Recall that both uncertainty sets are permutation invariant, i.e., that $\mathbf{q} \in \mathcal{Q}$ implies that $\mathbf{q}_{\pi(\cdot)} = (q_{\pi(1)}, \dots, q_{\pi(n)}) \in \mathcal{Q}$ for permutation π on $[n]$. We use the following equivalence

$$\mathcal{Q} = \Delta^{n-1} \text{ if and only if } \mathbf{e}_1 := (1, 0, \dots, 0) \in \mathcal{Q}.$$

Observe the following special cases.

- For any spectral risk measure, $(1, 0, \dots, 0) \in \mathcal{Q}$ if and only if $\sigma_n = 1$, which corresponds to the $(1 - 1/n)$ -superquantile.
- For balls in χ^2 -divergence, we measure the length of $n \|\mathbf{e}_1 - \mathbf{1}/n\|_2^2$:

$$n \|\mathbf{e}_1 - \mathbf{1}/n\|_2^2 = n \|\mathbf{e}_1\|_2^2 - 2n \langle \mathbf{e}_1, \mathbf{1}/n \rangle + n \|\mathbf{1}/n\|_2^2 = n - 1.$$

Thus, we have that $\rho = n - 1$ is the upper bound for this uncertainty set.

- For balls in KL-divergence, we argue similarly by computing $\text{KL}(\mathbf{e}_1 \parallel \mathbf{1}/n)$ (using the convention that $0 \log 0 = 0$):

$$\text{KL}(\mathbf{e}_1 \parallel \mathbf{1}/n) = \log n,$$

giving the upper bound $\rho = \log n$.

Minimum Uncertainty Set Size The user might wish to set the uncertainty set size to be at least large enough to achieve some form of robustness guarantee. One such guarantee is determined based on relating the solution of a distributionally robust optimization problem to the expected loss of the resulting procedure. In other words, we may use the minimum value of the distributionally robust objective to create a one-sided confidence interval for the population loss at a particular parameter value $\boldsymbol{\theta} \in \mathbb{R}^d$.

A central assumption that will be needed is that there exists some $B_{\boldsymbol{\theta}}$ such that for $\xi \sim P$,

$$\ell(\boldsymbol{\theta}, \xi) \leq B_{\boldsymbol{\theta}} \text{ almost surely under } P.$$

This is a reasonable assumption when the input space is bounded, as the loss may still be unbounded when w varies as well. Consider i.i.d. samples $\xi_1, \dots, \xi_n \sim P$. We derive a method that will assign to any failure probability $\delta \in (0, 1]$ an uncertainty set $\mathcal{Q}_{\delta} \subseteq \mathbb{R}^{n+1}$ such that with probability at least $1 - \delta$,

$$\max_{\mathbf{q} \in \mathcal{Q}_{\delta}} \left\{ \sum_{i=1}^n q_i \ell(\boldsymbol{\theta}, \xi_i) + q_{n+1} B_{\boldsymbol{\theta}} \right\} \geq \mathbb{E}_{\xi \sim P} [\ell(\boldsymbol{\theta}, \xi)]. \quad (2.35)$$

While this is not a valid confidence bound for a random $\boldsymbol{\theta}$ that is dependent with $\{\xi_i\}_{i=1}^n$ (such as the minimizer of an empirical risk objective), it can still be computed on held-out data and be interpreted conditionally. This is indeed the result of [Coppens and Patrinos \[2023, Proposition III.1\]](#), although we provide an alternate proof below which uses only elementary tools.

To achieve this result, we first prove a technical lemma. Consider the partial order \succeq on \mathbb{R}^{n+1} given by $\mathbf{x} \succeq \mathbf{y}$ if and only if

$$\sum_{i=1}^k x_i \geq \sum_{i=1}^k y_i \quad \forall k \in [n], \text{ and } \sum_{i=1}^{n+1} x_i = \sum_{i=1}^{n+1} y_i.$$

This is related to majorization but is not exactly the same, as we would be checking the inequality for the sorted versions of \mathbf{x} and \mathbf{y} . Then, observe the following lemma.

Lemma 2.7.1. *Consider $\mathbf{x}, \mathbf{y}, \mathbf{z} \in \mathbb{R}^{n+1}$ such that $\mathbf{x} \succeq \mathbf{y}$ and $z_1 \leq \dots \leq z_{n+1}$, i.e., \mathbf{z} is*

sorted. Then,

$$\mathbf{x}^\top \mathbf{z} \leq \mathbf{y}^\top \mathbf{z}.$$

Proof. We prove the result by induction. The case of $n = 0$ is satisfied trivially, as $x_1 = y_1 \implies x_1 z_1 = y_1 z_1$. Assume that the claim holds for vectors of size n ; we show that it holds for vectors of size $n + 1$. First, we define $\bar{\mathbf{x}} \in \mathbb{R}^n$ such that $\bar{x}_i = x_i$ for $i \leq n - 1$, and $\bar{x}_n = \sum_{i=1}^n y_i - \sum_{i=1}^{n-1} x_i$. By construction, it holds that

$$\bar{\mathbf{x}} \succeq (y_1, \dots, y_n).$$

Applying the inductive claim, we are assured that $\sum_{i=1}^n \bar{x}_i z_i \leq \sum_{i=1}^n y_i z_i$. Next, to achieve the desired result, write

$$\begin{aligned} \sum_{i=1}^{n+1} x_i z_i &= \sum_{i=1}^n \bar{x}_i z_i + \left(\sum_{i=1}^{n-1} x_i - \sum_{i=1}^n y_i \right) z_n + x_n z_n + x_{n+1} z_{n+1} \\ &= \sum_{i=1}^n \bar{x}_i z_i + \left(\sum_{i=1}^n x_i - \sum_{i=1}^n y_i \right) z_n + x_{n+1} z_{n+1} \\ &\leq \sum_{i=1}^n y_i z_i + \left(\sum_{i=1}^n x_i - \sum_{i=1}^n y_i \right) z_n + x_{n+1} z_{n+1}. \end{aligned}$$

Thus, the proof is complete if $(\sum_{i=1}^n x_i - \sum_{i=1}^n y_i) z_n + x_{n+1} z_{n+1} \leq y_{n+1} z_{n+1}$. To see this, use $z_n \leq z_{n+1}$ and $\sum_{i=1}^n x_i - \sum_{i=1}^n y_i \geq 0$ to achieve

$$\begin{aligned} \left(\sum_{i=1}^n x_i - \sum_{i=1}^n y_i \right) z_n + x_{n+1} z_{n+1} &\leq \left(\sum_{i=1}^n x_i - \sum_{i=1}^n y_i \right) z_{n+1} + x_{n+1} z_{n+1} \\ &= y_{n+1} z_{n+1} \end{aligned}$$

where we use that $\sum_{i=1}^{n+1} x_i - \sum_{i=1}^{n+1} y_i = 0$ in the last step. \square

We may now present our version of [Coppens and Patrinos \[2023, Proposition III.1\]](#).

Proposition 2.7.3. *Let U be uniformly distributed over Δ^n , the probability simplex in $n + 1$ dimensions. Assume that $\ell(\boldsymbol{\theta}, \xi)$ is a continuous random variable on \mathbb{R} . Then, for uncertainty*

set $\mathcal{Q}_\delta \subseteq \Delta^n$, the event $\{U \succeq \mathbf{q} \text{ for some } \mathbf{q} \in \mathcal{Q}_\delta\}$ is contained within the event (2.35).

Proof. Define the random variables $L_i = \ell(\boldsymbol{\theta}, \xi_i)$ for $i \in [n]$, which are i.i.d. copies of $L := \ell(\boldsymbol{\theta}, \xi)$. Denote by F the (continuous) cumulative distribution function of L . Then, for particular realizations $l_{(1)} \leq \dots \leq l_{(n)}$ of the order statistics of (L_1, \dots, L_n) , define $l_{(0)} := 0$ and $l_{(n+1)} = B_{\boldsymbol{\theta}}$. Write

$$\begin{aligned} \mathbb{E}_P[L] &= \int_0^B x \, dF(\mathbf{x}) = \sum_{i=1}^{n+1} \int_{l_{(i-1)}}^{l_{(i)}} x \, dF(\mathbf{x}) \\ &\leq \sum_{i=1}^{n+1} l_{(i)} \int_{l_{(i-1)}}^{l_{(i)}} dF(\mathbf{x}) \\ &\leq \sum_{i=1}^{n+1} l_{(i)} F(l_{(i)}) - F(l_{(i-1)}). \end{aligned}$$

Thus, for the random variables (L_1, \dots, L_n) , we have that

$$\mathbb{E}_P[L] \leq \sum_{i=1}^{n+1} L_{(i)} U_i,$$

where $U_i = F(L_{(i)}) - F(L_{(i-1)})$. Because L_1, \dots, L_n are continuous, bounded random variables, it follows that U is uniformly distributed over Δ^{n+1} . We now show that

$$U \succeq \mathbf{q} \text{ for some } \mathbf{q} \in \mathcal{Q}(\Delta) \implies \max_{\mathbf{q} \in \mathcal{Q}(\Delta)} \sum_{i=1}^n q_i L_i + q_{n+1} L_{n+1} \geq \mathbb{E}_P[L].$$

Letting $\bar{\mathbf{q}}$ be such that $U \succeq \bar{\mathbf{q}}$ and because $L_{(1)} \leq \dots \leq L_{(n+1)}$, we have by Lemma 2.7.1 that

$$\begin{aligned} \sum_{i=1}^{n+1} L_{(i)} U_i &\leq \sum_{i=1}^{n+1} L_{(i)} \bar{q}_i && \text{Lemma 2.7.1} \\ &\leq \max_{\mathbf{q} \in \mathcal{Q}(\Delta)} \sum_{i=1}^n q_i L_{(i)} + q_{n+1} L_{(n+1)} && \bar{\mathbf{q}} \in \mathcal{Q}(\Delta) \\ &= \max_{\mathbf{q} \in \mathcal{Q}(\Delta)} \sum_{i=1}^n q_i L_i + q_{n+1} L_{n+1}, \end{aligned}$$

where the final equality follows from permutation invariance of $\mathcal{Q}(\Delta)$. □

The significance of Proposition 2.7.3 is that given δ , one can parametrize the uncertainty set using a univariate parameter and compute the probability of $U \succeq \mathbf{q}$ via simulation. The size for which the probability matches $1 - \delta$ can be made the minimum size. This is especially easy for spectral risk measures, given the following result.

Lemma 2.7.2. *Let $\mathcal{Q}(\sigma)$ denote a spectral risk measure uncertainty set with spectrum σ . Then,*

$$\mathbb{P}[U \succeq \mathbf{q} \text{ for some } \mathbf{q} \in \mathcal{Q}(\sigma)] = \mathbb{P}[U \succeq \sigma].$$

Proof. We prove the result by showing the equivalence of the corresponding events. First, because $\sigma \in \mathcal{Q}(\sigma)$, it holds trivially that

$$U \succeq \sigma \implies U \succeq \mathbf{q} \text{ for some } \mathbf{q} \in \mathcal{Q}(\sigma).$$

Conversely, consider $\mathbf{q} \in \mathcal{Q}(\sigma)$ such that $U \succeq \mathbf{q}$. We can write \mathbf{q} as

$$\mathbf{q} = \sum_{\pi \in \Pi_{n+1}} \lambda_{\pi} \sigma_{\pi},$$

where Π_{n+1} is the set of permutations on $[n+1]$, $\sum_{\pi \in \Pi_{n+1}} \lambda_{\pi} = 1$, each $\lambda_{\pi} \geq 0$, and

$$\sigma_{\pi} := (\sigma_{\pi(1)}, \dots, \sigma_{\pi(n+1)}).$$

Then, for any $k \in [n+1]$, it holds that

$$\sum_{i=1}^k U_i \geq \sum_{\pi \in \Pi_{n+1}} \lambda_{\pi} \sum_{i=1}^k \sigma_{\pi(i)} \geq \min_{\pi \in \Pi_{n+1}} \sum_{i=1}^k \sigma_{\pi(i)} = \sum_{i=1}^k \sigma_i,$$

where we use in the last step that $\sigma_1 \leq \dots \leq \sigma_{n+1}$. Thus, $U \succeq \mathbf{q}$, completing the proof. \square

The probabilities can be computed exponentially quickly in the number of samples, rendering the simulation tractable.

Maximum Shift Cost This question is answered by considering uniform approximation bounds between the smoothed ($\nu > 0$) and unsmoothed ($\nu = 0$) objectives. Because the $\nu > 0$ case is typically considered a tool to improve optimization convergence rates in theory and in practice, we would like it to be as small as possible. Noting that

$$\max_{\mathbf{q} \in \mathcal{Q}} \langle \mathbf{q}, \ell(\boldsymbol{\theta}) \rangle \geq \max_{\mathbf{q} \in \mathcal{Q}} \langle \mathbf{q}, \ell(\boldsymbol{\theta}) \rangle - \nu \text{Reg}(\mathbf{q}) \geq \max_{\mathbf{q} \in \mathcal{Q}} \langle \mathbf{q}, \ell(\boldsymbol{\theta}) \rangle - \nu \max_{\mathbf{q}' \in \mathcal{Q}} \text{Reg}(\mathbf{q}'),$$

we compute the $\max_{\mathbf{q}' \in \mathcal{Q}} \text{Reg}(\mathbf{q}')$ for both uncertainty set classes. Because $\text{Reg}(\mathbf{q}) = \frac{1}{n} \sum_{i=1}^n f(nq_i)$ is a strongly convex, rotationally symmetric function centered at $\mathbf{1}/n$, over $\mathcal{Q}(\sigma)$ the maximum is attained at any vertex, yielding $\max_{\mathbf{q}' \in \mathcal{Q}} \text{Reg}(\mathbf{q}') = \text{Reg}(\sigma)$. For divergence balls, we have by definition that $\max_{\mathbf{q}' \in \mathcal{Q}} \text{Reg}(\mathbf{q}') = \rho$. This, if we aim to achieve suboptimality ε for the $\nu = 0$ objective by optimizing the $\nu > 0$ objective, then setting

$$\nu_{\max} = \begin{cases} \frac{\varepsilon}{4 \text{Reg}(\sigma)} & \text{if } \mathcal{Q} = \mathcal{Q}(\sigma) \text{ (spectral risk measure)} \\ \frac{\varepsilon}{4\rho} & \text{if } \mathcal{Q} = \mathcal{Q}(\rho) \text{ (divergence ball)} \end{cases}$$

contributes $\varepsilon/2$ suboptimality by smoothing.

Minimum Shift Cost As before, we address this question separately for spectral risk measures and for f -DRO. To establish notation, let

$$\mathcal{L}^{(\nu)}(\boldsymbol{\theta}) := \max_{\mathbf{q} \in \mathcal{Q}} \langle \mathbf{q}, \ell(\boldsymbol{\theta}) \rangle - \nu \text{Reg}(\mathbf{q}) + \frac{\mu}{2} \|\boldsymbol{\theta}\|_2^2$$

denote the primal objective with shift cost $\nu \geq 0$ with respect to a particular uncertainty set \mathcal{Q} and f -divergence penalty Reg (now indexed with ν). In both cases, we derive a computationally verifiable condition on w under which the following holds: for all $\nu, \bar{\nu} \geq 0$ small enough, we have that $\nabla \mathcal{L}^{(\nu)}(\boldsymbol{\theta}) = \nabla \mathcal{L}^{(\bar{\nu})}(\boldsymbol{\theta})$.

Proposition 2.7.4. *Let $\mathcal{Q} = \mathcal{Q}(\sigma)$ be a spectral risk measure uncertainty set with spectrum σ . Let $z^{(\nu)}$ denote the optimum map z^{opt} of (2.28) under a shift cost $\nu \geq 0$. If for some*

$\bar{\nu} \geq 0$ and $\boldsymbol{\theta} \in \mathbb{R}^d$ we have that

$$z_i^{(\bar{\nu})}(\ell(\boldsymbol{\theta})) = \ell_{(i)}(\boldsymbol{\theta}) - \bar{\nu} f'(n\sigma_i), \quad (2.36)$$

then

$$\nabla \mathcal{L}^{(\bar{\nu})}(\boldsymbol{\theta}) = \nabla \mathcal{L}^{(\nu)}(\boldsymbol{\theta}) \text{ for all } \nu \in (0, \bar{\nu}] \text{ and } \nabla \mathcal{L}^{(\bar{\nu})}(\boldsymbol{\theta}) \in \partial \mathcal{L}^{(0)}(\boldsymbol{\theta}).$$

Furthermore, the condition (2.36) is equivalent to

$$\bar{\nu} (f'(n\sigma_{i+1}) - f'(n\sigma_i)) \leq \ell_{(i+1)}(\boldsymbol{\theta}) - \ell_{(i)}(\boldsymbol{\theta}) \text{ for } i = 1, \dots, n. \quad (2.37)$$

Proof. The left-hand side of (2.36) results from taking the derivative of the right-hand side of (2.28). Thus, (2.36) is equivalent to the unconstrained solution of (2.28) being feasible in \mathbf{z} . Without loss of generality, assume that $\ell(\boldsymbol{\theta})$ is sorted in non-decreasing order. It follows that

$$z_i^{(\nu)}(\ell(\boldsymbol{\theta})) = \ell_i(\boldsymbol{\theta}) - \nu f'(n\sigma_i) \quad (2.38)$$

for all $\nu \leq \bar{\nu}$, as the right-hand side remains feasible for the problem (2.28) (in particular, due to the monotonicity of $z_i^{(\bar{\nu})}(\ell(\boldsymbol{\theta}))$). Defining similarly $q_i^{\text{opt}} = q_i^{(\nu)}$, it follows from Proposition 2.7.1 that

$$q_i^{(\nu)}(\ell(\boldsymbol{\theta})) = \frac{1}{n} [f^*]' \left(\frac{\ell_i(\boldsymbol{\theta}) - z_i^{(\nu)}(\ell(\boldsymbol{\theta}))}{\nu} \right) = \frac{1}{n} [f^*]' (f'(n\sigma_i)) = \sigma_i,$$

which is independent of ν and $\bar{\nu}$. First, consider $\nu > 0$, and notice that

$$\begin{aligned} \nabla \mathcal{L}^{(\bar{\nu})}(\boldsymbol{\theta}) &= \nabla \ell(\boldsymbol{\theta})^\top q^{(\bar{\nu})}(\ell(\boldsymbol{\theta})) + \mu w \\ &= \nabla \ell(\boldsymbol{\theta})^\top q^{(\nu)}(\ell(\boldsymbol{\theta})) + \mu w \\ &= \nabla \mathcal{L}^{(\nu)}(\boldsymbol{\theta}) \end{aligned}$$

For the case of $\nu = 0$, note that $\ell(\boldsymbol{\theta})^\top \sigma + \mu \boldsymbol{\theta}_\nu^*$ is a subgradient of $\mathcal{L}^{(0)}$ (see Mehta et al. [2023, Proposition 2]). The second claim follows from the feasibility (in particular, the

monotonicity) of $z^{(\nu)}(\ell(\boldsymbol{\theta}))$. \square

Proposition 2.7.5. *Let $\mathcal{Q} = \mathcal{Q}(\rho)$ be an f -divergence uncertainty set with radius ρ . Let $t^{(\nu)}$ denote the optimum of t^{opt} of (2.32) under a shift cost $\nu \geq 0$. If for some $\bar{\nu} \geq 0$ and $\boldsymbol{\theta} \in \mathbb{R}^d$, we have that*

$$h'_{\ell(\boldsymbol{\theta})}(t^{(\bar{\nu})}(\ell(\boldsymbol{\theta}))) = 0, \quad (2.39)$$

for h'_i defined in (2.34), then then

$$\nabla \mathcal{L}^{(\bar{\nu})}(\boldsymbol{\theta}) = \nabla \mathcal{L}^{(\nu)}(\boldsymbol{\theta}) \text{ for all } \nu \in (0, \bar{\nu}] \text{ and } \nabla \mathcal{L}^{(\bar{\nu})}(\boldsymbol{\theta}) \in \partial \mathcal{L}^{(0)}(\boldsymbol{\theta}).$$

Furthermore, the condition (2.39) is equivalent to the solution of $h'_{\ell(\boldsymbol{\theta})}(t) = 0$ in t being greater than or equal to $\bar{\nu}$.

Proof. As in the proof of Proposition 2.7.4, the left-hand side of (2.39) results from taking the derivative of the right-hand side of (2.32). Thus, (2.39) is equivalent to the unconstrained solution of (2.32) being feasible in t . The remainder of the proof follows identically to that of Proposition 2.7.4. \square

Before showing numerical benchmarks, we also comment that there are numerous statistical ideas that appear in this chapter that have connections to broader literature, including the method of empirical likelihood. We clarify the similarities and differences between our work and these areas in the next section.

2.8 Comparison to Broader Literature

As reflected by the diversity of topics in this chapter, distributionally robust optimization can be studied from many statistical and computational angles. Accordingly, there are several opportunities to contextualize our results with respect to other subfields within statistics. These include quantitative finance, nonparametric statistics, classical viewpoints on distribution shift, and existing stochastic variance-reduced algorithms for empirical risk minimization.

Risk Measures In Section 2.1, we introduced closed balls f -divergence balls and spectral risk measures as the major classes of distributionally robust objectives. However, spectral risk measures originated from a parallel line of work on *risk measures*, with theoretical roots in convex analysis and applications in econometrics and finance [Rockafellar and Uryasev, 2013, Föllmer and Schied, 2002, He et al., 2022, Ben-Tal and Teboulle, 2007]. Risk measures are functionals of a real-valued random variable (such as the expectation) that quantify some notion of “tail error”. Functionals that satisfied particular axiomatic properties were deemed *coherent* [Artzner et al., 1999], and among them were the class of spectral risk measures [Acerbi and Tasche, 2002]. In other words, spectral risk measures arose conceptually to quantify the tailedness of a random variable but can be viewed in the framework of DRO based on the variational arguments in the upcoming Section 2.3. On the other hand, the etymology of f -divergence uncertainty sets can be seen in reverse: starting as a natural framework for DRO, modern results demonstrated that these also can be represented as measures of tail error (albeit not in closed form). Indeed, by Shapiro [2017, Section 3.2], we have that

$$\sup_{Q \ll P: D_f(Q \| P) \leq \rho} \mathbb{E}_Q [\ell(\boldsymbol{\theta}, \xi)] = \inf_{\lambda \geq 0, \gamma \in \mathbb{R}} \mathbb{E}_P \left[\lambda f^* \left(\frac{\ell(\boldsymbol{\theta}, \xi) - \gamma}{\lambda} \right) + \lambda \rho + \gamma \right]. \quad (2.40)$$

We may actually relate the dual-minimization form to another notion of risk measure, which will help interpret (2.40). Let $\psi : \mathbb{R} \rightarrow \mathbb{R}_+$ be a convex function and define

$$\text{oce}_P^\psi(X) = \inf_{\gamma \in \mathbb{R}} \mathbb{E}_P [\psi(X - \gamma)] + \gamma$$

as the *optimized certainty equivalent (OCE)* with *disutility function* ψ [Leqi et al., 2019, Lee et al., 2020]. This risk measure is often interpreted in quantitative finance as an optimal allocation between current and future losses; the investor is subject to a random loss X and may choose to incur γ of it while subjecting the remaining unknown portion to the function ψ , which represents the value of some loss/return in the future. Then, proceeding from

(2.40), we have that

$$\begin{aligned} \inf_{\lambda \geq 0, \gamma \in \mathbb{R}} \mathbb{E}_P \left[\lambda f^* \left(\frac{\ell(\boldsymbol{\theta}, \xi) - \gamma}{\lambda} \right) + \lambda \rho + \gamma \right] &= \inf_{\lambda \geq 0, \gamma \in \mathbb{R}} \mathbb{E}_P [\psi_\lambda (\ell(\boldsymbol{\theta}, \xi) - \gamma)] + \gamma \\ &= \inf_{\lambda \geq 0} \text{oce}_P^{\psi_\lambda}(\ell(\boldsymbol{\theta}, \xi)), \end{aligned}$$

for $\psi_\lambda(x) = \lambda f^*(x/\lambda) + \lambda \rho$. Thus, the f -DRO objective is the smallest OCE over a class of disutility functions $\{\psi_\lambda : \lambda \geq 0\}$, which themselves may be interpreted as various investments.

A number of recent works study L -risks in the form of (2.12), with a focus on statistical properties. The works [Khim et al. \[2020\]](#) and [Maurer et al. \[2021\]](#) provide classical statistical learning bounds for L -risk objectives, and the latter focuses on unsupervised tasks like clustering. [Holland and Mehdi Haress \[2022\]](#) present a derivative-free learning procedure for general L -risk problems in the fully stochastic/streaming setting. As for optimizing these risk measures, [Fan et al. \[2017\]](#) and [Kawaguchi and Lu \[2020\]](#) study batch and stochastic algorithms, respectively, for the “average top- k ” loss, which is exactly equivalent to the superquantile. We instead focus on developing incremental algorithms, akin to those for ERM [\[Mairal, 2014, Le Roux et al., 2012, Shalev-Shwartz and Zhang, 2013, Johnson and Zhang, 2013, Defazio et al., 2014\]](#), which apply to all L -risks.

Empirical Likelihood We review the method of empirical likelihood (EL) from nonparametric statistics [\[Owen, 1990\]](#) while comparing it to distributionally robust optimization (DRO) from both the statistical and computational perspectives. We show the following.

- The distributionally robust objective with the reverse-KL ball uncertainty set is equal to the upper limit of an empirical likelihood confidence region for mean estimation. The confidence level can be expressed in terms of the radius of the ball.
- The maximum over the weights q occurs under different constraints in DRO versus EL. In EL, the constraint is that the distribution induced by q has a particular mean (an affine constraint), as opposed to the small divergence constraint in DRO. In other

words, DRO maximizes a linear objective under a nonlinear constraint, whereas EL maximizes a nonlinear objective with a linear constraint.

EL is used to derive nonparametric confidence intervals for an arbitrary functional of a probability measure. For now, we consider the special case mean estimation from i.i.d. real-valued observations X_1, \dots, X_n with cumulative distribution function F . Rather than specifying a distributional class for F , the idea is to construct a likelihood ratio of a distribution with mean μ dominated by the empirical CDF F_n itself and reject the proposed mean μ when this likelihood ratio is small. Letting $\mathbf{q} = (q_1, \dots, q_n)$ be the weights on the samples, we consider *all* valid likelihood ratios for such a distribution, and take their highest value, which we denote by $\mathcal{R}(\mu)$.

$$\mathcal{R}(\mu) := \max \left\{ \prod_{i=1}^n (nq_i) : \sum_{i=1}^n q_i X_i = \mu, \mathbf{q} \in \Delta^{n-1} \right\}. \quad (2.41)$$

The objective is log-strictly concave and the feasible set is compact in \mathbb{R}^n , so a unique maximizer assuredly exists. The confidence or acceptance region for μ is then given by

$$\{\mu : \mathcal{R}(\mu) \geq r\},$$

where r is the to-be-calibrated threshold on the likelihood ratio. To address the first bullet, observe the following.

Proposition 2.8.1. *Let $f(t) = -\log(t)$ for $t > 0$. Then,*

$$\max \left\{ \sum_{i=1}^n q_i X_i : D_f(\mathbf{q} \| \mathbf{1}/n) \leq \rho, \mathbf{q} \in \Delta^{n-1} \right\} = \max \{ \mu \in \mathbb{R} : \mathcal{R}(\mu) \geq e^{-n\rho} \}.$$

Proof. We start with the quantity on the left-hand side and derive the right-hand side as a result. For notational ease, set $T_n(\mathbf{q}) := \sum_{i=1}^n q_i X_i$. Introduce the auxiliary variable μ to

the DRO problem:

$$\begin{aligned}
& \max \{T_n(\mathbf{q}) : D_f(\mathbf{q} \parallel \mathbf{1}/n) \leq \rho, \mathbf{q} \in \Delta^{n-1}\} \\
&= \max_{\mu \in \mathbb{R}} \max \{ \mu : T_n(\mathbf{q}) = \mu, D_f(\mathbf{q} \parallel \mathbf{1}/n) \leq \rho, \mathbf{q} \in \Delta^{n-1} \} \\
&= \max_{\mu \in \mathbb{R}} \max_{\substack{\mathbf{q} \in \Delta^{n-1} \\ T_n(\mathbf{q}) = \mu}} \min_{\xi \geq 0} \mu - \xi(D_f(\mathbf{q} \parallel \mathbf{1}/n) - \rho),
\end{aligned}$$

where we introduced the Lagrange multiplier ξ . Because of the existence of a primal-dual strictly feasible point $(\mathbf{q}, \xi) = (\mathbf{1}/n, 1)$, we apply Slater's condition and strong duality to the inner max-min problem to achieve

$$\begin{aligned}
& \max_{\mathbf{q} \in \Delta^{n-1}} \{T_n(\mathbf{q}) : D_f(\mathbf{q} \parallel \mathbf{1}/n) \leq \rho\} \\
&= \max_{\mu \in \mathbb{R}} \min_{\xi \geq 0} \max_{\substack{\mathbf{q} \in \Delta^{n-1} \\ T_n(\mathbf{q}) = \mu}} \mu - \xi(D_f(\mathbf{q} \parallel \mathbf{1}/n) - \rho) \\
&= \max_{\mu \in \mathbb{R}} \min_{\xi \geq 0} \max_{\substack{\mathbf{q} \in \Delta^{n-1} \\ T_n(\mathbf{q}) = \mu}} \mu - \frac{\xi}{n} \left(- \sum_{i=1}^n \log(nq_i) - n\rho \right) \\
&= \max_{\mu \in \mathbb{R}} \min_{\xi \geq 0} \mu + \frac{\xi}{n} \left(\max_{\substack{\mathbf{q} \in \Delta^{n-1} \\ T_n(\mathbf{q}) = \mu}} \sum_{i=1}^n \log(nq_i) + n\rho \right) \\
&= \max_{\mu \in \mathbb{R}} \min_{\xi/n \geq 0} \mu - \frac{\xi}{n} (-\log(\mathcal{R}(\mu)) - n\rho).
\end{aligned}$$

Next, we interpret ξ/n as a Lagrange multiplier for the constraint $-\log(\mathcal{R}(\mu)) \leq n\rho$, which is equivalent to $\mathcal{R}(\mu) \geq e^{-n\rho}$. Rewriting the above in constraint form completes the argument. \square

Next, we consider the computational aspect. As stated before, DRO maximizes a linear objective under a nonlinear constraint, whereas EL maximizes a nonlinear objective with a

linear constraint:

$$\max_{\mathbf{q} \in \Delta^{n-1}} \{T_n(\mathbf{q}) : D_f(\mathbf{q} \| \mathbf{1}/n) \leq \rho\}, \quad (\text{DRO})$$

$$\max_{\mathbf{q} \in \Delta^{n-1}} \left\{ \sum_{i=1}^n \log(nq_i) : T_n(\mathbf{q}) = \mu \right\}. \quad (\text{EL})$$

Furthermore, EL does not compute a primal-optimal point in closed form; instead, it computes a saddle point of the Lagrangian when introducing a Lagrange multiplier for the constraint $T_n(\mathbf{q}) = \mu$. If λ denotes this multiplier, then we have that

$$q_i = \frac{1}{n} \frac{1}{1 + \lambda(X_i - \mu)} \quad \forall i \in \{1, \dots, n\}.$$

The proof of the univariate mean empirical likelihood theorem, or ELT [Owen, 2001, Theorem 2.2], i.e., the asymptotic χ_1^2 distribution of $-2 \log \mathcal{R}(\mu)$ does actually compute λ (hence, does not have a closed form for q) but shows that $\lambda = O_p(n^{-1/2})$, from which asymptotic arguments follow. A potentially different computational scheme is required for every choice of T_n , and a subsequent analysis of the Lagrange multiplier. That is why there are different ELT variants for each function (quantiles, variance, regression parameters, etc.).

DRO for In-Distribution Generalization While typically motivated through considerations such as distribution shift from P to Q , there are applications of DRO methods even for traditional generalization guarantees from P_n to P . In particular, consider the two distributional parameters (which we assume to uniquely exist for the sake of discussion) as

$$\boldsymbol{\theta}_{\text{in}}^* := \arg \min_{\boldsymbol{\theta} \in \Theta} \underbrace{\mathbb{E}_{\xi \sim P} [\ell(\boldsymbol{\theta}, \xi)]}_{\mathcal{R}_{\text{in}}(\boldsymbol{\theta}, P)}, \text{ and } \boldsymbol{\theta}_{\text{out}}^* := \arg \min_{\boldsymbol{\theta} \in \Theta} \underbrace{\sup_{Q \in \mathcal{Q}(P)} \mathbb{E}_{\xi \sim Q} [\ell(\boldsymbol{\theta}, \xi)]}_{\mathcal{R}_{\text{out}}(\boldsymbol{\theta}, P)}.$$

The statistical and optimization techniques developed for DRO are often in service of estimating $\boldsymbol{\theta}_{\text{out}}^*$. We now describe scenarios in which these same techniques can in fact be used to achieve desirable guarantees for the estimation of $\boldsymbol{\theta}_{\text{in}}^*$ as well. In the statements below, all

probabilities are taken with respect to an independently drawn sample $\xi_1, \dots, \xi_n \sim P$, with

$$\boldsymbol{\theta}_n \equiv \boldsymbol{\theta}_n(\xi_1, \dots, \xi_n) := \arg \min_{\boldsymbol{\theta} \in \Theta} \underbrace{\sup_{Q \in \mathcal{Q}(P_n)} \mathbb{E}_{\xi \sim Q} [\ell(\boldsymbol{\theta}, \xi)]}_{\mathcal{R}_{\text{in}}(\boldsymbol{\theta}, P_n)}.$$

The first type of guarantee concerns *certification*, or claiming that with high probability,

$$\max_{Q \in \mathcal{Q}(P_n)} \mathbb{E}_{\xi \sim Q} [\ell(\boldsymbol{\theta}_n, \xi) | P_n] \geq \mathbb{E}_{\xi \sim P} [\ell(\boldsymbol{\theta}_n, \xi) | P_n]. \quad (2.42)$$

In other words, the optimal value of the distributionally robust empirical risk is, with high probability, an upper bound for the expected risk of the minimizer conditioned on the training data. Thus, while $\boldsymbol{\theta}_n$ is computed with a distributionally robust empirical objective, we still consider its performance on the data-generating distribution. Clearly, we may satisfy (2.42) if the sufficient condition

$$P \in \mathcal{Q}(P_n). \quad (2.43)$$

holds. However, this is only possible when $\mathcal{Q}(P_n)$ contains distributions that may not be absolutely continuous with respect to P_n . This precludes the likelihood ratio-based framework from Section 2.1. For these, we turn to methods that define $\mathcal{Q}(P_n)$ using a ball in the Wasserstein metric \mathbf{W} . In this case, if ε is the radius of the Wasserstein ball, then (2.43) follows if and only if $\mathbf{W}(P, P_n) \leq \varepsilon$. Thus, the guarantees will follow by appealing to concentration results of $\mathbf{W}(P, P_n)$ (see Kuhn et al. [2019]). In fact, (2.43) implies the stronger condition that (2.42) holds uniformly, i.e.,

$$\max_{Q \in \mathcal{Q}(P_n)} \mathbb{E}_{\xi \sim Q} [\ell(\boldsymbol{\theta}, \xi)] \geq \mathbb{E}_{\xi \sim P} [\ell(\boldsymbol{\theta}, \xi)] \quad \forall \boldsymbol{\theta} \in \Theta. \quad (2.44)$$

Thus, if ε is fixed, we may say that if n is sufficiently large, then (2.42) holds.

For likelihood ratio-based DRO, as in the case of f -divergences, results typically operate in reverse: for a fixed sample size n , if the radius ρ is large enough, then we may achieve

(see [Coppens and Patrinos \[2023\]](#)) the following result:

$$\forall \boldsymbol{\theta} \in \Theta, \exists \rho : \max_{Q \in \mathcal{Q}_\rho(P_n)} \mathbb{E}_{\xi \sim Q} [\ell(\boldsymbol{\theta}, \xi)] \geq \mathbb{E}_{\xi \sim P} [\ell(\boldsymbol{\theta}, \xi)] \text{ w.h.p.} \quad (2.45)$$

Notice that the statement above is a random event with respect to the empirical measure P_n . Moreover, for continuous loss distributions, the selected ρ can be chosen independently of $\boldsymbol{\theta}$. We apply a similar result to select a minimum uncertainty set size in Section 2.7.

A second type of guarantee concerns *confidence*, or to say that with any prespecified probability δ , we may define a random upper bound U_n^δ based on observations ξ_1, \dots, ξ_n such that

$$\liminf_{n \rightarrow \infty} \mathbb{P} \left[U_n^\delta \geq \min_{\boldsymbol{\theta} \in \Theta} \mathbb{E}_{\xi \sim P} [\ell(\boldsymbol{\theta}, \xi)] \right] \geq 1 - \delta. \quad (2.46)$$

This is simply a one-sided asymptotic $(1 - \delta)$ -confidence interval for the minimal value of $\boldsymbol{\theta} \mapsto \mathbb{E}_{\xi \sim P} [\ell(\boldsymbol{\theta}, \xi)]$. This setting can be interpreted as estimating a distributional parameter which is defined as the minimum of a particular functional of P defined for each $\boldsymbol{\theta}$, with associated uncertainty quantification. Notably, one possible choice of $U_n^\delta = \max_{Q \in \mathcal{Q}(P_n)} \mathbb{E}_{\xi \sim Q} [\ell(\boldsymbol{\theta}_n, \xi) | P_n]$, because by (2.42), we have that

$$\max_{Q \in \mathcal{Q}(P_n)} \mathbb{E}_{\xi \sim Q} [\ell(\boldsymbol{\theta}_n, \xi) | P_n] \geq \mathbb{E}_{\xi \sim P} [\ell(\boldsymbol{\theta}_n, \xi) | P_n] \geq \min_{\boldsymbol{\theta} \in \Theta} \mathbb{E}_{\xi \sim P} [\ell(\boldsymbol{\theta}, \xi)] = \mathbb{E}_{\xi \sim P} [\ell(\boldsymbol{\theta}_{\text{in}}^*, \xi)].$$

However, this may result in a loose interval, which can be seen explicitly in the upcoming expression for U_n^δ . As shown in [Duchi et al. \[2021\]](#), we have that

$$U_n^\delta = \min_{\boldsymbol{\theta} \in \Theta} \max_{Q \in \mathcal{Q}(P_n)} \mathbb{E}_{\xi \sim Q} [\ell(\boldsymbol{\theta}, \xi)],$$

$$\mathcal{Q}(P_n) = \left\{ Q \ll P_n : D_f(Q \| P_n) \leq \frac{\chi_{1,1-2\delta}^2}{n} \right\} \quad (2.47)$$

satisfies (2.46), where $\chi_{1,1-2\delta}^2$ denotes the $(1 - 2\delta)$ -quantile of the χ^2 -distribution with 1 degree of freedom. To comment on the looseness of the left-hand side of (2.42) for this purpose, we notice that the radius in (2.47) is shrinking at a rate of $O(n^{-1})$. Thus, as expected from a confidence interval, its size shrinks to zero as $n \rightarrow \infty$. In contrast, when

specifically pursuing distributional robustness, the radius is kept constant. Furthermore, the functional $\mathbb{E}_{\xi \sim P} [\ell(\boldsymbol{\theta}, \xi)]$ can be replaced by a Hadamard differentiable functional of P . This method is known as *generalized empirical likelihood*, with respect to Owen’s classical empirical likelihood approach for uncertainty quantification for Fréchet differentiable functionals [Owen \[1990\]](#). [Lam and Zhou \[2017\]](#) also achieved (2.46) for particular f -divergences.

Direct Likelihood Ratio Estimation Our choice of objective is motivated by the fact that the user does not know what type of evaluation distribution Q they may observe during deployment. However, similar techniques can be applied in the setting of *domain adaptation*. This setting posits that $P \neq Q$ in some structured way, and this structure may be surmised from a potentially small number of examples from the shifted distribution. For instance, a training and test set of images may differ in distribution due to heterogeneous lighting conditions (a natural shift) or corruption of the test images through blurring (a synthetic shift). Here, we do in fact assume that $Q \ll P$, so there is a to-be-estimated population likelihood ratio $\beta : (\mathbf{x}, \mathbf{y}) \mapsto \frac{dQ}{dP}(\mathbf{x}, \mathbf{y})$ for $(\mathbf{x}, \mathbf{y}) \in \mathcal{X} \times \mathcal{Y}$. Given that such a population quantity exists, we may naturally attempt to estimate it, so that we may apply (2.2) to the sample.

One formal notion of structure is *covariate shift*, in which we assume that $P_{Y|X} = Q_{Y|X}$ so that $P \neq Q$ if and only if $P_X \neq Q_X$. As an abuse of notation, we consider the true β to only be a function of $\mathbf{x} \in \mathcal{X}$. We produce an estimate $\beta_n : \mathcal{X} \rightarrow \mathbb{R}$ which may be a function-valued estimator and not necessarily a list of n weights. Continuing within the covariate shift framework, because the shift only occurs in P_X , the statistical analysis will be *transductive*, in that we consider fixed (non-random) covariates $\{\mathbf{x}_i^{\text{tr}}\}_{i=1}^n$, random responses $Y_1, \dots, Y_n \sim P_{Y|X}$, and fixed, unlabeled examples $\{\mathbf{x}_i^{\text{te}}\}_{i=1}^m$ from the target domain. We discuss two seminal examples of methods in this setting. Both methods are kernel-based, in that the user first specifies a statistical model \mathcal{P} such that $P_X, Q_X \in \mathcal{P}$ and a reproducing kernel Hilbert space (RKHS) \mathcal{H} over \mathbb{R} containing functions on \mathcal{X} . We denote by $\phi : \mathcal{X} \rightarrow \mathcal{H}$ the feature/lifting map and by $\mu_X : \mathcal{P} \rightarrow \mathcal{H}$ the mean map of \mathcal{H} , i.e.,

$\mu_X(P_X)(\mathbf{x}) = \mathbb{E}_{X \sim P_X} [\phi(X)(\mathbf{x})]$ for any $\mathbf{x} \in \mathcal{X}$. Consider the empirical measure of the examples from the target domain.

$$Q_{m,X} := \frac{1}{m} \sum_{i=1}^m \delta_{\mathbf{x}_i^{\text{te}}}. \quad (2.48)$$

Note that while Q is absolutely continuous with respect to P , there is no guarantee that $Q_{m,X} \ll (P_n)_X$. In kernel mean matching [Gretton et al., 2008], β_n is estimated by solving

$$\min_{\beta \in \mathcal{B}(P_n)} \left\| \mu_X(Q_{m,X}) - \mathbb{E}_{(X,Y) \sim P_n} [\beta(X,Y)\phi(X)] \right\|_{\mathcal{H}}^2$$

which can be written as a finite-dimensional quadratic program. In practice, β is optimized with a slight relaxation of the constraint $\mathbb{E}_{(X,Y) \sim P_n} [\beta(X,Y)] = 1$, which is standard for a likelihood ratio. Note that this approach, while employing RKHS techniques, does not actually assume that $\beta \in \mathcal{H}$. On the other hand, kernelized unconstrained least-squares importance fitting [Kanamori et al., 2013] selects β_n as an element of the RKHS using a clever objective. Observe that for any β , the squared loss criterion against $\frac{dQ}{dP}$ can be written as

$$\left\| \beta - \frac{dQ}{dP} \right\|_{\mathbf{L}^2(P)}^2 = \int_{\mathcal{X}} \beta^2(\mathbf{x}) dP(\mathbf{x}) - 2 \int_{\mathcal{X}} \beta(\mathbf{x}) dQ(\mathbf{x}) + \int_{\mathcal{X}} \frac{dQ}{dP}(\mathbf{x}) dQ(\mathbf{x})$$

with the last term independent of β . Thus, the optimization problem is written as

$$\min_{\beta \in \mathcal{H}} \mathbb{E}_{\xi \sim P_n} [\beta^2(Z)] - 2\mathbb{E}_{\xi \sim Q_{m,X}} [\beta(Z)],$$

Standard techniques to convert problems over RKHS to finite-dimensional programs via Tikhonov regularization can be applied to the problem above.

Importance-Weighted Algorithms for ERM Related to the likelihood ratio estimation literature, one may view a DRO routine as simultaneously estimating importance weights for training examples while actually solving the resulting importance-weighted empirical risk minimization (ERM) problem. The idea of combining ERM with a weighting function for generic losses is often credited to Shimodaira [2000], in the context of (weighted) maximum

likelihood estimation. In fact, this principle can be traced even further to pseudolikelihood estimation used in the analysis of survey data [Binder, 1983]. From then, importance-weighted ERM became a prominent approach in statistical learning for correcting sample bias and handling covariate shift [Huang et al., 2006, Sugiyama et al., 2007, 2008, Wen et al., 2014], with the likelihood ratio estimation problem studied in its own right [Kanamori et al., 2009]. Note that these applications of likelihood ratios for *bias reduction* methods differ from their use in stochastic simulation for *variance reduction* (e.g., in Monte Carlo estimation) [Kahn and Marshall, 1953].

To be precise, the use of importance weighting in a Monte Carlo application implies that there is a *known* target distribution Q , and we may select the distribution P from which we sample in order to generate the best estimate of a parameter of Q . In our setting, that of bias reduction, we may not change the data-generating distribution P , yet still estimate a parameter of Q . If either P or Q were unknown, this leads to an intermediate problem of estimating the likelihood function. Even if both were known, then we may devise unbiased estimators but still may wish to characterize the rates of decay in the variance [Ma et al., 2023]. Finally, in either bias or variance-reducing reweighting, if the estimand is some population risk function, there is the subsequent question of how empirical risk minimization performs on a rebalanced objective.

After the initial surge of importance weighted ERM for covariate shift in supervised learning, these methods appeared in the context of active learning [Mahmood et al., 2014, Swaminathan and Joachims, 2015, Wang et al., 2021], as a major source of bias comes from sampling from a particular policy (as opposed to “expert” demonstrations). More recently, the role of importance weighting in overparametrized statistical models (e.g., neural networks) has been studied aggressively in theory and practice [Byrd and Lipton, 2019, Xu et al., 2021]. As before, the task is typically thought of as a bilevel optimization problem consisting of 1) estimating the likelihood ratio, and 2) applying the learned likelihood ratio to importance-weighted ERM. The likelihood ratio may itself be parametrized by a neural network [Kato and Teshima, 2021] or related to an optimal transport map [Gong et al., 2016],

and even be learned in an online manner [Zhang et al., 2023]. A number of “end-to-end” approaches (i.e., learning the weighted estimator with a single optimization problem) have also been proposed [Fang et al., 2020, Zhang et al., 2020]. Finally, from the perspective of domain generalization more broadly, likelihood ratio estimation is seen as a form of *domain alignment* Zhang et al. [2013], such as moment/distribution matching between the source and target distributions.

Fairness and Subpopulation Shift DRO objectives, which are maxima over reweightings of the observed training data, are a special case of *subpopulation shift*, wherein the data-generating distribution is modeled as a mixture of subpopulations, and the distribution shift stems from changes in the mixture. In our case, the subpopulations are point masses at the observed data points. In the context of *algorithmic fairness*, the subpopulations may represent data conditioned on some protected attribute (e.g., race, gender, age range), and common notations of fairness such as *demographic/statistical parity* [Agarwal et al., 2018, 2019] impose (informally) that model performance with respect to each subpopulation should be roughly equal. As such, robustness to reweighting and algorithmic fairness are often aligned notions [Williamson and Menon, 2019], with recent research arguing that distributionally robust models are more fair [Hashimoto et al., 2018, Vu et al., 2022] and that fair models are more distributionally robust [Mukherjee et al., 2022].

2.9 Experiments

We compare Prospect against baselines in a variety of learning tasks. While we focus attention on its performance as an optimizer of its training objective, we also highlight metrics of interest on the test set in fairness and distribution shift benchmarks. To clarify terminology, we use “Prospect” and “LSVRG” to refer to the methods as described in Figure 2.3. In comparisons, we include stochastic algorithms that either are single-hyperparameter “out-of-the-box” methods, such as stochastic gradient descent and stochastic regularized dual averaging [Xiao, 2009], or multi-hyperparameter methods that converge linearly on strongly

convex SRM-based objectives such as stochastic saddle-point SAGA [Palaniappan and Bach, 2016]. The algorithm implementation and data preparation code are made publicly available online: <https://github.com/ronakdm/prospect>.

Setting, Baselines, Evaluation. We consider supervised learning tasks with input-label example (\mathbf{x}_i, y_i) . Losses are of the form $\ell_i(\boldsymbol{\theta}) := \text{err}(y_i, \langle \boldsymbol{\theta}, \varphi(\mathbf{x}_i) \rangle)$, with a fixed feature embedding ϕ , and err measuring prediction loss. Uncertainty sets considered are the CVaR/superquantile, extremile, and ESRM. We compare against four baselines: minibatch stochastic gradient descent (SGD), stochastic regularized dual averaging (SRDA) [Xiao, 2009], Saddle-SAGA [Palaniappan and Bach, 2016], and LSVRG [Mehta et al., 2023]. For SGD and SRDA, we use a batch size of 64, and for LSVRG we use an epoch length of n . For Saddle-SAGA, we find that allowing different learning rates for the primal and dual variables improves experimental performance, so we compare against an improved heuristic (setting the dual stepsize as $10n$ times smaller than the primal stepsize). We plot

$$\text{Suboptimality}(\boldsymbol{\theta}) = (\mathcal{L}(\boldsymbol{\theta}) - \mathcal{L}(\boldsymbol{\theta}^*)) / (\mathcal{L}(\boldsymbol{\theta}^{(0)}) - \mathcal{L}(\boldsymbol{\theta}^*)), \quad (2.49)$$

where $\boldsymbol{\theta}^*$ is approximated by running LBFGS [Nocedal and Wright, 1999] on the objective until convergence. The x -axis displays the number of calls to any first-order oracle $\boldsymbol{\theta} \mapsto (\ell_i(\boldsymbol{\theta}), \nabla \ell_i(\boldsymbol{\theta}))$ divided by n , i.e., the number of passes through the training set. We fix the shift cost $\nu = 1$ and regularization parameter $\mu = 1/n$. Experimental details and additional experiments with various hyperparameters are contained in Appendix A.4.

2.9.1 Tabular Least-Squares Regression

We consider five tabular regression benchmarks under square loss. The datasets used are `yacht` ($n = 244$) [Tsanas and Xifara, 2012], `energy` ($n = 614$) [Baressi Segota et al., 2020], `concrete` ($n = 824$) [Yeh, 2006], `kin8nm` ($n = 6553$) [Akujuobi and Zhang, 2017], and `power` ($n = 7654$) [Tüfekci, 2014]. The training curves are shown in Figure 2.4.

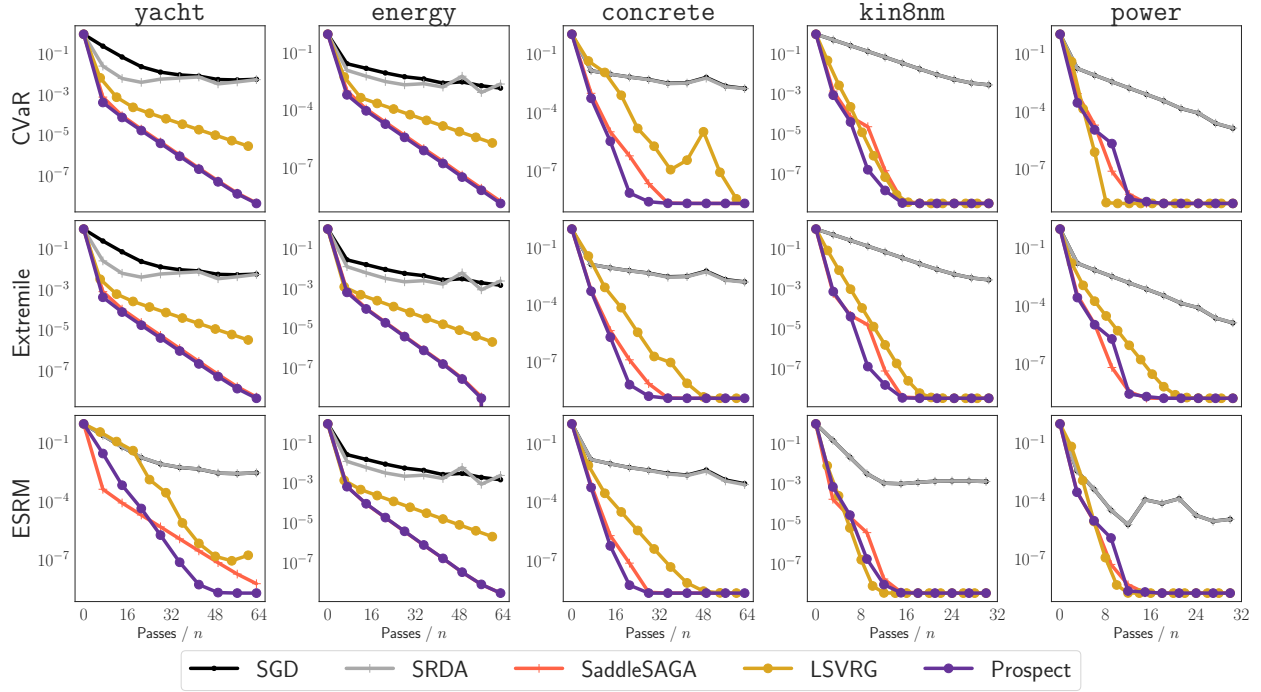


Figure 2.4: **Regression Benchmarks.** The y -axis measures suboptimality as given by (2.49), while the x -axis measures the number of calls to the function value/gradient oracle divided by n (i.e., passes through the training set). Rows indicate different SRM objectives while columns indicate datasets.

Results. Across datasets and objectives, we find that Prospect exhibits linear convergence at a rate no worse than SaddleSAGA and LSVRG, but that is often much better. For example, Prospect converges to precision 10^{-8} for the CVaR on *concrete* and the extremile on *power* within half the number of passes that LSVRG takes for the same suboptimality. Similarly, for the ESRM on *yacht*, SaddleSAGA requires 64 epochs to reach the same precision as Prospect at 40 epochs. The direct stochastic methods, SGD and SRDA, are biased and fail to converge for any learning rate.

2.9.2 Fair Classification and Regression

Inspired by [Williamson and Menon \[2019\]](#), we explore the relationship between distributional robustness and group fairness on 2 common tabular benchmarks. **Diabetes 130-Hospitals** (`diabetes`) is a classification task of predicting readmission for diabetes patients based on clinical data from US hospitals [\[Rizvi et al., 2014\]](#). **Adult Census** (`acsincome`) is a regression task of predicting income of US adults from data compiled by the American Community Survey [\[Ding et al., 2021\]](#).

Evaluation. We evaluate fairness with the *statistical parity score*, which compares predictive distributions of a model given different values of a particular protected attribute [Agarwal et al. \[2018, 2019\]](#). Letting $Z = (X, Y, A)$ denote a random (input, label, metadata attribute) triplet, a model g is said to satisfy statistical parity (SP) if the conditional distribution of $g(X)$ over predictions given $A = a$ is equal for any value a . Intuitively, SP scores measure the maximum deviation between these distributions for any over a , so values close to zero indicate SP-fairness. In `diabetes`, we use gender as the protected attribute A , whereas in `acsincome` we use race as the protected attribute. Note that the protected attributes are not supplied to the models. Results are given in Figure 2.5.

Results. Firstly, we note that Prospect converges rapidly on both datasets while LSVRG fails to converge on `diabetes` and SaddleSAGA fails to converge on `acsincome`. Secondly, LSVRG does not stabilize with respect to classification SP, showing a mean/std SP score of $1.38 \pm 0.25\%$ within the final ten passes on the `diabetes` CVaR, whereas Prospect gives $0.82 \pm 0.00\%$, i.e., a 40% relative improvement with greater stability. While SaddleSAGA does stabilize in SP on `diabetes`, it fails to qualitatively decrease at all on the `acsincome`. Interestingly, while suboptimality and SP-fairness are correlated for Prospect, SGD (reaching only 10^{-1} suboptimality with respect to the CVaR objectives on `acsincome`) achieves a lower fairness score. Again, across both suboptimality and fairness, Prospect is either the best or close to the best.

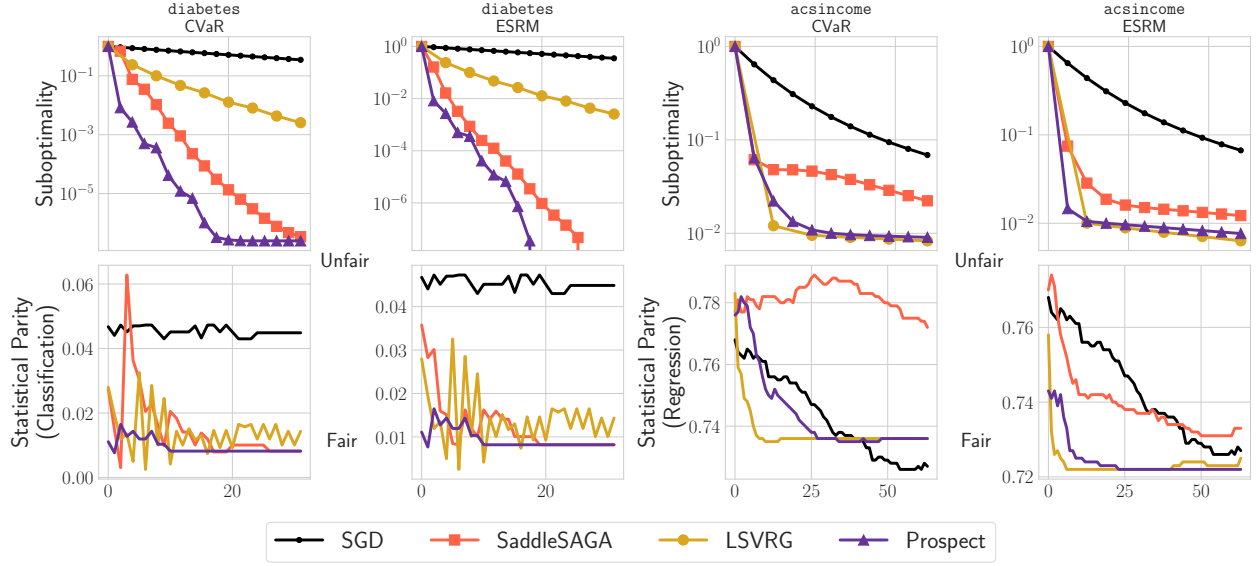


Figure 2.5: **Fairness Benchmarks.** **Top:** Training curves for optimizers on the CVaR and extremile for **diabetes** (left) and CVaR and extremile for **acsincome** (right). **Bottom:** Statistical parity scores for the two classification objectives on **diabetes** (left) and regression objectives on **acsincome**. Smaller values indicate better performance for all metrics.

2.9.3 Image and Text Classification under Distribution Shift

We consider two tasks from the WILDS distribution shift benchmark [Koh et al., 2021]. The **Amazon Reviews** (**amazon**) task [Ni et al., 2019] consists of classifying text reviews of products to a rating of 1-5, with disjoint train and test reviewers. The **iWildCam** (**iwildcam**) image classification challenge [Beery et al., 2020] contains labeled images of animals, flora, and backgrounds from cameras placed in wilderness sites. Shifts are due to changes in camera angles, locations, lighting... We use $n = 10000$ and $n = 20000$ examples respectively. For both datasets, we train a *linear probe classifier*, i.e., a linear model over a frozen deep representation. For **amazon**, we use a pretrained BERT model [Devlin et al., 2019a] fine-tuned on a held-out subset of the Amazon Reviews training set for 2 epochs. For **iwildcam**, we use a ResNet50 pretrained on ImageNet (without fine-tuning).

Evaluation. Apart from the training suboptimality, we evaluate the spectral risk objectives on their robustness to subpopulation shifts. We define each subpopulation group based on the true label. For `amazon`, we use the *worst group misclassification error* on the test set [Sagawa et al., 2020]. For `iwildcam`, we use the *median group error* owing to its larger number of classes.

Results. For both `amazon` and `iwildcam`, Prospect and SaddleSAGA (with our heuristic) outperform LSVRG in training suboptimality. We hypothesize that this phenomenon is due to checkpoints of LSVRG getting stale over the n -length epochs for these datasets with large n (leading to a slow reduction of bias). In contrast, Prospect and SaddleSAGA avoid this issue by dynamically updating the running estimates of the importance weights. For the worst group error for `amazon`, Prospect and SaddleSAGA outperform LSVRG. Prospect has a mean/std worst group error of $77.38 \pm 0.00\%$ over the last ten passes on the extremile, whereas SaddleSAGA has a slightly worse $77.53 \pm 1.57\%$. Interestingly, on `iwildcam`, LSVRG and Prospect give stronger generalization performance, nearly 1pp better, than SaddleSAGA in terms of median group misclassification rate. In summary, across tasks and objectives, Prospect demonstrates best or close to best performance.

2.9.4 Scaling Laws and Shift Cost

We aim to disentangle the many effects of the shift cost parameter ν in the following experiments. While ν has the statistical interpretation of penalizing values of \mathbf{q} that stray far from the original uniform weights $\mathbf{1}/n$, we choose to control this distributional robustness property by instead using the uncertainty set \mathcal{Q} . This allows users of the non-smooth ($\nu = 0$) and smooth ($\nu > 0$) variants of the objective to use similar intuition for designing the uncertainty set. Accordingly, we view ν purely from an optimization lens, that is, we use it to allow for the algorithm to converge quickly. This involves trading off the per-iteration cost and the number of iterations.

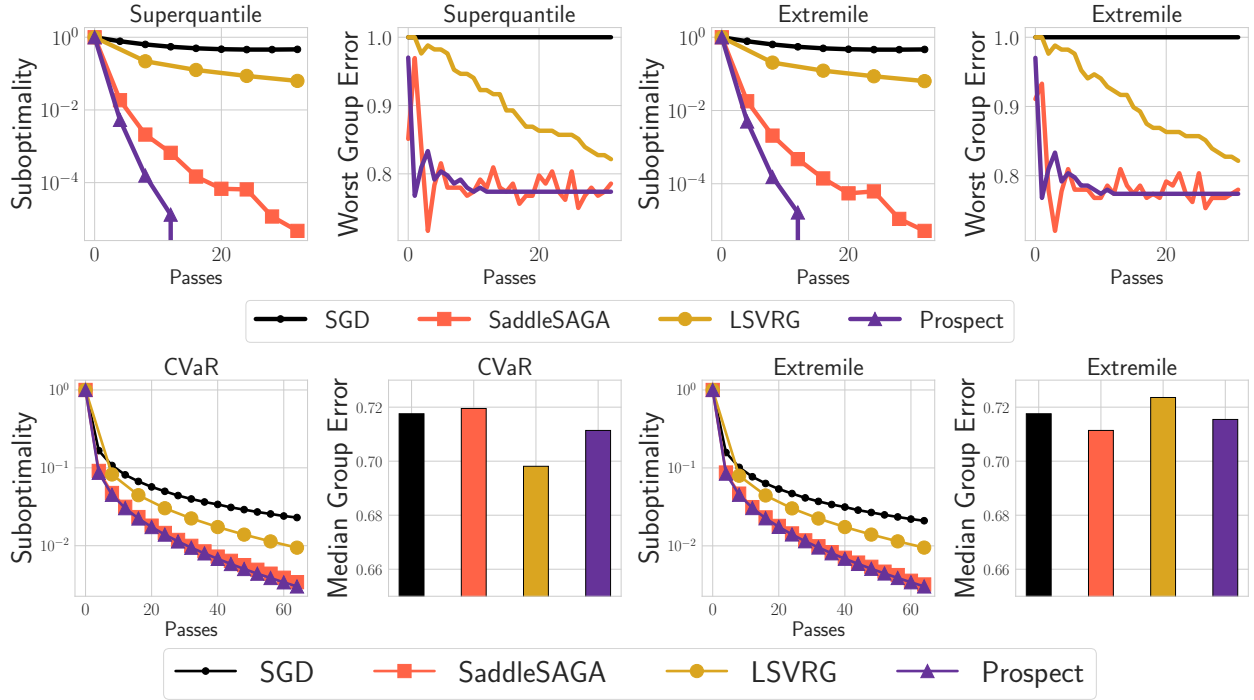


Figure 2.6: **Distribution Shift Benchmarks.** **Top:** Training curves and worst group misclassification error on `amazon` test. **Bottom:** Training curves and median group misclassification error on the `iwildcam` test set. Smaller values indicate better performance for all metrics.

Evaluation. Consider the following competing effects.

- **Number of Iterations:** Let $\varepsilon > 0$ be the desired suboptimality. Ignoring other parameters, the number of iterations has an $O(1/\nu)$ dependence on $\nu > 0$ (for any $\mu \geq 0$). More intuitively, ν regularizes the dual problem, leading to better conditioning of the objective. Thus, increasing ν will generally *decrease* the number of steps needed for convergence (but for a different objective).
- **Per-Iteration Cost:** The per-iteration cost has a subtle dependence on ν . Each iteration performs full maximization over q , which relies on subroutines such as sorting (e.g., the Pool Adjacent Violators algorithm in Appendix A.3). In particular, the map $l \mapsto \max_{q \in \mathcal{Q}} \langle q, l \rangle - \nu D_f(q \| \mathbf{1}/n)$ can often be computed by fitting a monotonic function

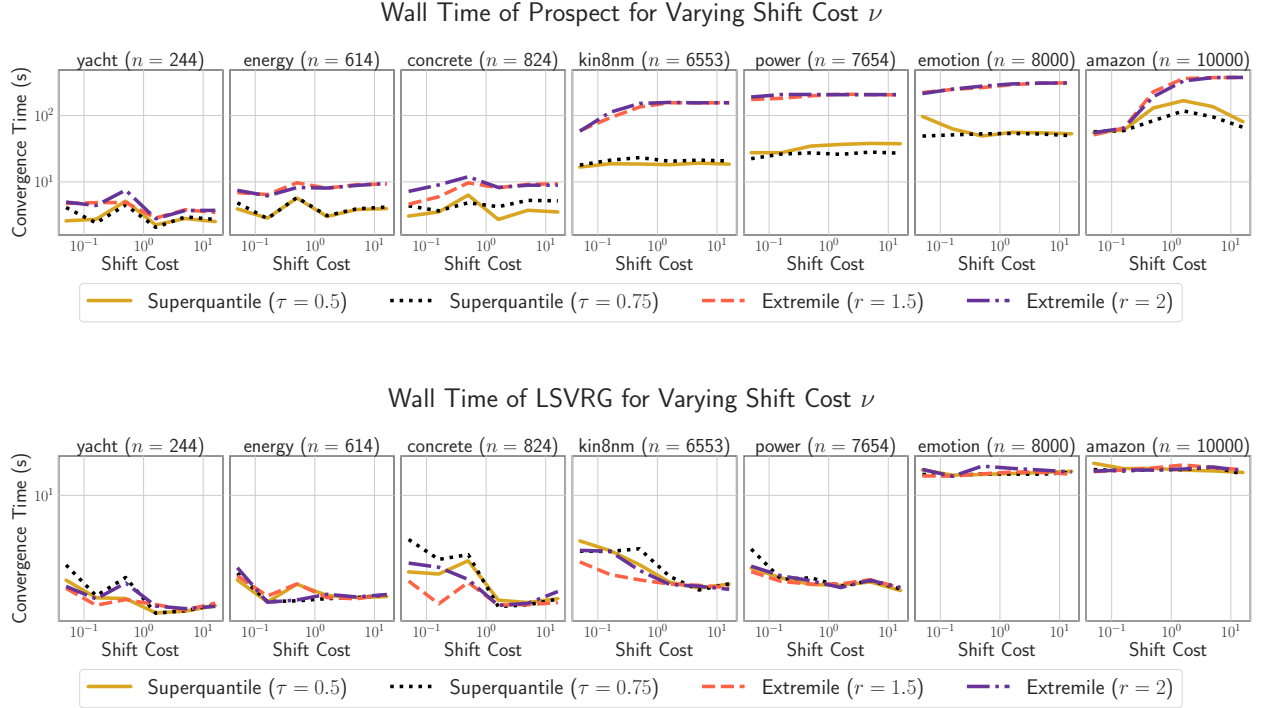


Figure 2.7: **Wall Time and Shift Cost Relationship.** **Top:** Wall time against shift cost for reaching convergence (squared gradient norm less than 10^{-3}) for Prospect. **Bottom:** Wall time against shift cost for reaching convergence (squared gradient norm less than 10^{-3}) for LSVRG.

to the sorted vector $(l_{(1)}, \dots, l_{(n)})$ where $l_{(1)} \leq \dots \leq l_{(n)}$. The parameter ν perturbs this sorted vector, perhaps breaking monotonicity and hence making it more “difficult” to fit a monotonic function. This difficulty renders as more inner loop iterations, so that the per-iteration cost may *increase* with increasing ν .

- **Altering the Objective/Solution:** Finally, changing ν also changes the objective, and in turn changes the optimal primal-dual pair (θ^*, \mathbf{q}^*) . If the practitioner designs *Qcal* using a superquantile, for example, they imagine \mathbf{q}^* to be a “top- k weights” like vector. However, as we see below, increasing ν may alter the dual solution to the point of resembling the uniform weights with a few large “spikes”. Thus, we would not like to bias our objective significantly.

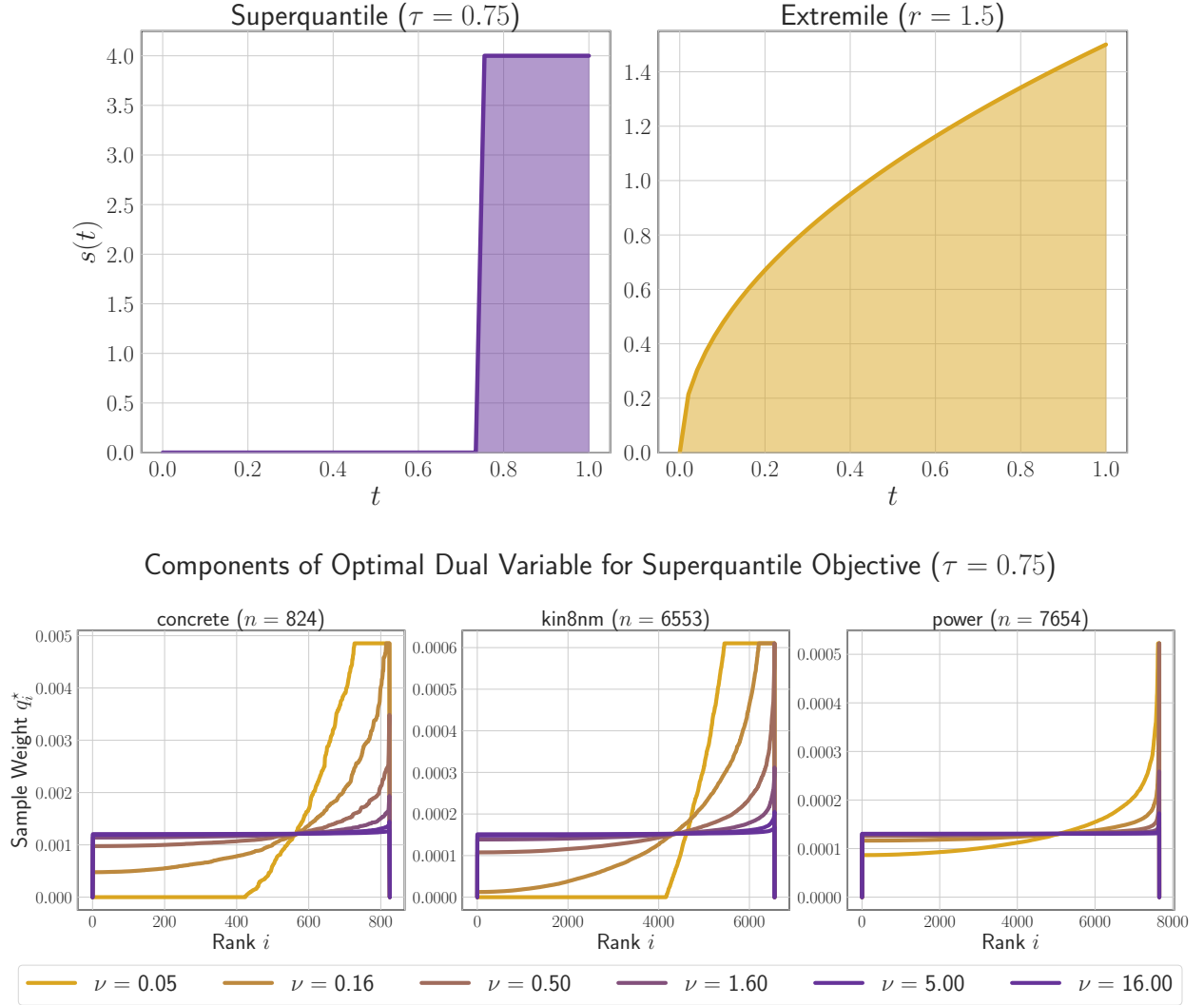


Figure 2.8: **Optimal Dual Solutions and Shift Cost Relationship.** **Top:** Visualization of continuous spectra for the superquantile and extremile. **Bottom:** Sorted optimal dual solution \mathbf{q}^* for different values of ν , meant to compare to the superquantile spectrum in the top left panel.

Results. Ultimately, we find that setting ν to be as small as possible without harming convergence (e.g., a dual value of $\nu = 0.05$) generally works well, as 1) convergence times are generally *faster* for small shift costs (Figure 2.7), and 2) the dual solution remains close to that of $\nu = 0$ and depends only on the uncertainty set (Figure 2.8).

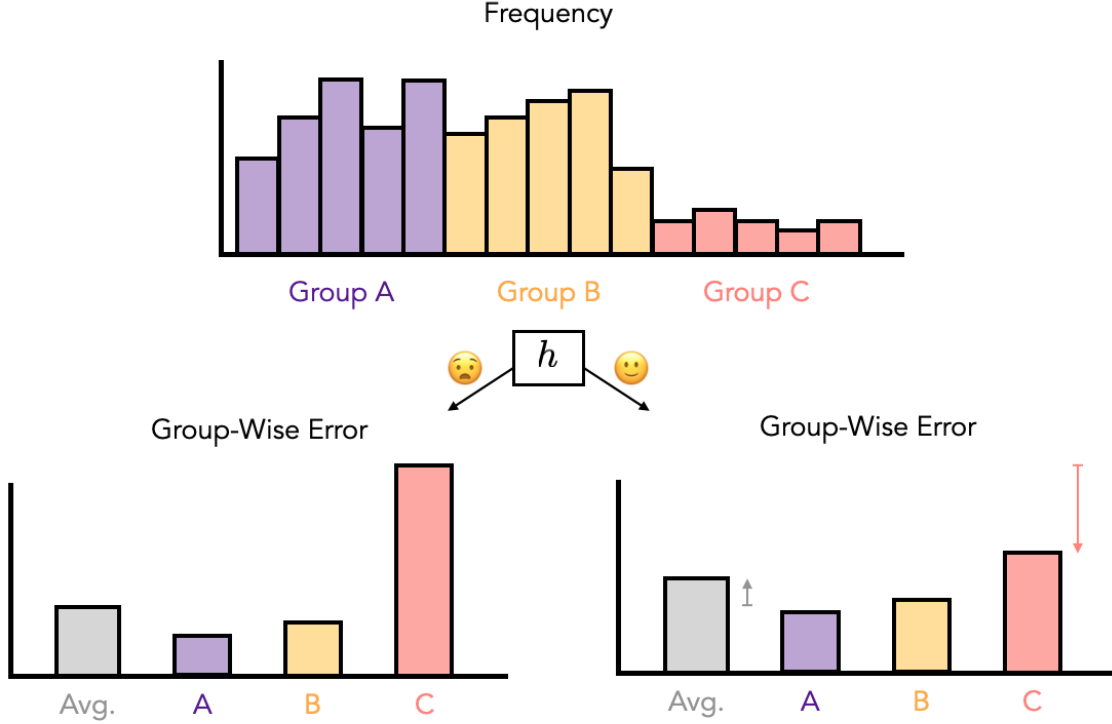


Figure 2.9: **Illustration of Group-Wise Error Evaluation.** The top panel depicts a hypothetical probability distribution, where the sample space is partitioned into three groups. When applying a statistical prediction model h to this distribution, the bottom left and right panels show non-uniform and uniform performance conditional on each group. On the right panel specifically, the arrows indicate a slight increase in the average loss across all examples, but a more dramatic decrease in the worst-case group-wise error.

We consider extensions to both the theoretical analysis and algorithmic aspects in the next section.

2.10 Incorporating Group Structure

While the stochastic algorithms introduced in this chapter make the DRO problem tractable for moderately large sample sizes, the $O(n)$ cost of each update of the dual variables becomes a major limitation when attempting to adapt methods to neural network applications such as computer vision and natural language processing. The seminal work of [Sagawa et al. \[2020\]](#) applied DRO at the level of groups or subpopulations within the data. That is, while different

groups may be upweighted or downweighted in the objective, the overall performance of a model on a particular group was measured by the simple average of within-group errors. In this section, we refer to the setup of Section 2.1 as *instance-level DRO*, to be contrasted with *group DRO* as described below. At a high level, the practitioner wishes to train a model that controls *worst case group-wise error* across groups by changing the training objective in a manner similar to that of Section 2.4. Unfavorable and favorable group-wise error distributions are depicted in Figure 2.9.

To describe the setting formally, we rely on the likelihood ratio-based framework introduced in Chapter 1 and Section 2.8. We assume that the data-generating distribution P and the shifted distribution Q are over the augmented domain $\Xi \times \mathcal{A}$, where Ξ denotes observable data (possibly feature-label pairs) and $\mathcal{A} = \{\mathbf{a}_1, \dots, \mathbf{a}_M\}$ denotes an additional component of the datum, often called the *group label* or *protected attribute* in various contexts. Denote by P_ξ and P_A the marginal distributions of P on Ξ and \mathcal{A} , respectively. Because A is discrete, we will use the notation $P_A(\mathbf{a}) \equiv P_A(\{\mathbf{a}\})$ to indicate the probability mass function in this section. To proceed, we collect the assumptions required for a well-defined empirical risk objective below.

Assumption 2.10.1. The following statements hold.

- We have that $P_A(\mathbf{a}) > 0$ for all $\mathbf{a} \in \mathcal{A}$.
- An i.i.d. sample $(\xi_1, A_1), \dots, (\xi_n, A_n) \sim P$ is observed at train time. In particular, training examples are accompanied by group labels (but test examples may not be).
- Consider the empirical measure P_n with marginals $P_{n,\xi}$ and $P_{n,A}$ on Ξ and \mathcal{A} , respectively. Then, $P_{n,A}(\mathbf{a}) > 0$ for all $\mathbf{a} \in \mathcal{A}$ (i.e., all groups are observed in the data).

The first part of the Assumption 2.10.1 is purely technical, so that group conditional means are well-defined. The second is methodological, implying that the user may observe the group information when training models. The third part is in fact a random event for

which it makes sense to consider an objective that depends on group-wise means. In a statistical analysis, we would return another estimator in the case that $P_{n,A}(\mathbf{a}) = 0$ for some $a \in \mathcal{A}$, and upper bound the probability of this occurring in terms of n . See Proposition B.3.3, used to prove the main results of Section 4.3, for a direct example of such an argument.

We first introduce the population objective. While the format will resemble the likelihood ratio-based reweighted objective (2.2) from Section 2.2, it involves structural assumptions similar to the covariate shift example from Section 2.8 and the distributional assumption from (2.3). We first construct

$$\mathcal{B}^{\text{struct}}(P) := \left\{ \beta \in \mathcal{B}_0(P) : \beta(\mathbf{z}, \mathbf{a}) = \beta(\mathbf{z}', \mathbf{a}) \ \forall (\mathbf{z}, \mathbf{z}', \mathbf{a}) \in \Xi \times \Xi \times \mathcal{A} \right\},$$

indicating that the likelihood ratio can only be a function of the attribute input $\mathbf{a} \in \mathcal{A}$. This is analogous to covariate shift, in which case the likelihood ratio was only a function of the feature component of ξ . Then, using f -DRO as an example, we also introduce

$$\mathcal{B}^{\text{DRO}}(P) := \{ \beta \in \mathcal{B}_0(P) : \mathbb{E}_P[f(\beta(\xi, A))] \leq \rho \},$$

and finally set $\mathcal{B}(P) := \mathcal{B}^{\text{struct}}(P) \cap \mathcal{B}^{\text{DRO}}(P)$. In words, we consider distribution shifts that only occur on the weights of P_A (via $\mathcal{B}^{\text{struct}}(P)$), subject to a constraint on uncertainty (via $\mathcal{B}^{\text{DRO}}(P)$). The objective then reads as

$$\min_{\theta \in \mathbb{R}^d} \max_{\beta \in \mathcal{B}(P)} \left\{ \mathbb{E}_{(\xi, A) \sim P} [\beta(\xi, A) \ell(\theta, \xi)] := \mathbb{E}_{A \sim P_A} [\beta(\cdot, A) \cdot \mathbb{E}_P[\ell(\theta, \xi)] | A] \right\}. \quad (2.50)$$

Notice that another implicit choice is made here, which is that while we observe realizations of A , they are not included as input to the prediction function or the instance-level loss function $\ell : \mathbb{R}^d \times \Xi \rightarrow \mathbb{R}$. Next, we describe the empirical counterpart of this objective along with the conceptual and practical differences between group DRO and the instance-level setting.

2.10.1 Learning Objective and Algorithm Description

Learning Objective By replacing P with P_n in (2.50), develop an objective similar to (2.18) at the group level. Given observed data $\{(\xi_i, A_i)\}_{i=1}^n$, we maintain the notation $\ell_i(\boldsymbol{\theta}) := \ell(\boldsymbol{\theta}, \xi_i)$ for parameter vector $\boldsymbol{\theta} \in \mathbb{R}^d$ and define groups of indices $\mathcal{I}_1, \dots, \mathcal{I}_M$ such that $A_i = \mathbf{a}_m \iff i \in \mathcal{I}_m$. Letting $n_m := |\mathcal{I}_m| \geq 2$ be the group-level sample size, we denote the empirical group means as

$$\bar{\ell}_m(\boldsymbol{\theta}) := \frac{1}{n_m} \sum_{i \in \mathcal{I}_m} \ell_i(\boldsymbol{\theta}).$$

We then naturally define the objective

$$\min_{\boldsymbol{\theta} \in \mathbb{R}^d} \left[\mathcal{L}^{\text{gr}}(\boldsymbol{\theta}) := \max_{\mathbf{q} \in \mathcal{Q}} \sum_{m=1}^M q_m \cdot \bar{\ell}_m(\boldsymbol{\theta}) - \nu \text{Reg}(\mathbf{q}) + \frac{\mu}{2} \|\boldsymbol{\theta}\|_2^2 \right], \quad (2.51)$$

where $\mathcal{Q} \subseteq \Delta^{M-1}$ is an uncertainty set and $\text{Reg} : \Delta^{M-1} \rightarrow \mathbb{R}$ is a regularization function that is strongly convex with respect to a particular norm. While (2.51) might seem like a direct analog of (2.18), there are a number of differences to be considered in the group DRO setting.

- There is no canonical choice of Reg based on a statistical divergence, as there are multiple choices justified for different reasons. At the instance level, the empirical measure of the data places uniform weight on all examples. While we may use $\text{Reg}(\mathbf{q}) := D_f(\mathbf{q} \parallel \mathbf{1}_M/M)$ as a regularizer, the empirically observed group weights $\mathbf{p}^{\text{gr}} := (n_1/n, \dots, n_M/n)$ could be completely different. We may consider using $\text{Reg}(\mathbf{q}) := D_f(\mathbf{q} \parallel \mathbf{p}^{\text{gr}})$ to promote similarity to the observed distribution.
- Unlike in (2.18), where we may compute each $(\ell_i, \nabla \ell_i)$ at cost $\tilde{O}(d)$, we require $\tilde{O}(nd/M)$ operations to query $(\bar{\ell}_i, \nabla \bar{\ell}_i)$ on average. In a practical setting, we may only access unbiased estimates of the $(\bar{\ell}_i, \nabla \bar{\ell}_i)$ oracles, which introduces another source of statistical error to be considered in the analysis.

Algorithm 2 Group-DRO with Prospect Bias Reduction

Inputs: Initial points $\boldsymbol{\theta}^{(0)}$, stepsize $\eta > 0$, per-group batch sizes $b_1 \in [n_1], \dots, b_M \in [n_M]$, primal regularization (weight decay) parameter μ , dual regularization (shift cost) parameter ν , and number of iterations t .

- 1: Set initialization $\boldsymbol{\theta}^{(0)}$.
- 2: Compute $\mathbf{l}^{(0)} = (l_1^{(0)}, \dots, l_n^{(0)}) = \ell(\boldsymbol{\theta}^{(0)})$.
- 3: Compute $\bar{\mathbf{l}}^{(0)} = (\bar{l}_1^{(0)}, \dots, \bar{l}_M^{(0)})$ where $\bar{l}_m^{(0)} = \frac{1}{n_m} \sum_{i \in \mathcal{I}_m} l_i^{(0)}$ for $m = 1, \dots, M$.
- 4: **for** $k = 0, \dots, t - 1$ **do**
- 5: Sample b_m points from each group, yielding indices $\{(i_{m,1}, \dots, i_{m,b_m})\}_{m=1}^M$.
- 6:
- 7: Compute the group-wise average losses via

$$\ell^{(k+1)}(\boldsymbol{\theta}^{(k)}) = \left(\frac{1}{b_1} \sum_{j=1}^{b_1} \ell_{i_{1,j}}(\boldsymbol{\theta}^{(k)}), \dots, \frac{1}{b_M} \sum_{j=1}^{b_M} \ell_{i_{M,j}}(\boldsymbol{\theta}^{(k)}) \right) \in \mathbb{R}^M.$$

- 8: Compute the group weights $\mathbf{q}^{(k)} = q^{\text{opt}}(\bar{\mathbf{l}}^{(k)})$, which incorporates the shift cost ν .
- 9: Compute the gradient estimate $\mathbf{g}^{(k+1)} = \sum_{m=1}^M q_m^{(k)} \nabla \ell_m^{(k+1)}(\boldsymbol{\theta}^{(k)})$. In an auto-differentiation framework, this can be accomplished by backpropagating through the computation graph of $\boldsymbol{\theta} \mapsto \langle \text{StopGrad}(\mathbf{q}^{(k)}), \ell^{(k+1)}(\boldsymbol{\theta}) \rangle$.
- 10: $\boldsymbol{\theta}^{(k+1)} = \text{Step}(\boldsymbol{\theta}^{(k)}, \mathbf{g}^{(k+1)}, \eta, \mu)$.
- 11:
- 12: Update $l_{i_{m,j}}^{(k+1)} = \ell_{i_{m,j}}(\boldsymbol{\theta}^{(k+1)})$ for $m \in [M]$ and $j \in [b_m]$.
- 13: Set $l_i^{(k+1)} = l_i^{(k)}$ for any $i \in [n]$ such that $i \neq i_{m,j}$ for all (m, j) .
- 14: Update $\bar{l}_m^{(k+1)} = \frac{1}{n_m} \sum_{i \in \mathcal{I}_m} l_i^{(k+1)}$ for $m = 1, \dots, M$.

Output: Final point $\boldsymbol{\theta}^{(t)}$

Algorithm Description On the algorithmic side, we apply an optimization scheme in a similar spirit to Algorithm 1, described in Algorithm 2 below. In this case, we eschew the variance reduction component because such algorithms are generally used for large-scale applications with neural network models, so storing historical parameter vectors may be impractical (even at the group level). Instead, we will capture the variance of the group mean estimates described above in an SGD-style analysis. Algorithm 2 operates by maintaining an $O(n)$ -sized table $\mathbf{l}^{(k)}$ and changing $\sum_m b_m$ coordinates at each iteration. This table aggregates past information and is used to estimate the dual variables (the sample weights), which would otherwise cost $O(n)$ oracles to compute. However, even in the case of bias reduction, we find

Algorithm 3 Group-DRO with Stratified Sampling

- Inputs:** Initial points $\boldsymbol{\theta}^{(0)}$, stepsize $\eta > 0$, per-group batch sizes $b_1 \in [n_1], \dots, b_M \in [n_M]$, primal regularization (weight decay) parameter μ , dual regularization (shift cost) parameter ν , and number of iterations t .
- 1: Set initialization $\boldsymbol{\theta}^{(0)}$.
 - 2: **for** $k = 0, \dots, t - 1$ **do**
 - 3: Sample b_m points from each group, yielding indices $\{(i_{m,1}, \dots, i_{m,b_m})\}_{m=1}^M$.
 - 4: Compute $\ell^{(k+1)}(\boldsymbol{\theta}^{(k)}) = \left(\frac{1}{b_1} \sum_{j=1}^{b_1} \ell_{i_{1,j}}(\boldsymbol{\theta}^{(k)}), \dots, \frac{1}{b_M} \sum_{j=1}^{b_M} \ell_{i_{M,j}}(\boldsymbol{\theta}^{(k)}) \right) \in \mathbb{R}^M$.
 - 5: Compute $\mathbf{q}^{(k)} = q^{\text{opt}}(\ell^{(k+1)}(\boldsymbol{\theta}^{(k)}))$ and the gradient estimate $\mathbf{g}^{(k+1)} = \sum_{m=1}^M q_m^{(k)} \nabla \ell_m^{(k+1)}(\boldsymbol{\theta}^{(k)})$ (see Algorithm 2).
 - 6: $\boldsymbol{\theta}^{(k+1)} = \text{Step}(\boldsymbol{\theta}^{(k)}, \mathbf{g}^{(k+1)}, \eta, \mu)$.
- Output:** Final point $\boldsymbol{\theta}^{(t)}$
-

in experiments that this approach generally does not change optimization or classification performance, as the average within-group loss can be estimated with low bias using moderate batch sizes (between the orders of 10^1 and 10^2). Thus, we also describe a more standard stochastic gradient algorithm Algorithm 3 that can be easily plugged into any deep learning workflow with minimal additional code. Algorithm 3 still differs from previous approaches in the literature, as we may apply the techniques for handling spectral risk measure uncertainty sets and smoothing developed throughout this chapter. Finally, we allow for the choice of a general update function in the description of the algorithms, namely in line 10 of Algorithm 2 and line 6 of Algorithm 3. While we typically consider the standard stochastic gradient update

$$\boldsymbol{\theta}^{(k+1)} = \text{Step}(\boldsymbol{\theta}^{(k)}, \mathbf{g}^{(k+1)}, \eta, \mu) = \boldsymbol{\theta}^{(k)} - \eta \mathbf{v}^{(k+1)}, \text{ for } \mathbf{v}^{(k+1)} = \mathbf{g}^{(k+1)} + \mu \boldsymbol{\theta}^{(k)}, \quad (2.52)$$

the algorithm works well empirically by using update steps that are attuned to deep learning models, such as the Adam optimizer [Kingma and Ba, 2015].

Sampling Details In both Algorithm 2 and Algorithm 3, the data is sampled in a *stratified* manner; that is, we select batch sizes for each group and sample data points from each group on every iteration [Singh and Mangat, 1996, pages 102–144]. This ensures that the

user may compute estimates of the average loss on each group and control their statistical precision. While this can be implemented in practice by using multiple data loaders, one trick to approximate this with minimal changes to a typical PyTorch workflow is to sample the data using *weighted sampling with replacement*. In this scheme, one specifies sample weights (p_1, \dots, p_n) (or an unnormalized version thereof) associated with each data point, and samples from this distribution with replacement when producing a mini-batch. If the weights are selected so that for all $m \in [M]$, the conditions

$$p_i = p_j \text{ for all } i, j \in \mathcal{I}_m \text{ and } \sum_{i \in \mathcal{I}_m} p_i = \frac{1}{M},$$

then we may expect approximately b observations from each group in any mini-batch of size bM . In PyTorch, we simply supply the sample weight vector when initializing the data loader. In the code snippet below, we use `group_labels` to refer to the vector (A_1, \dots, A_n) , where the $\mathcal{A} = \{0, \dots, M-1\}$, and `batch_size` refers to bM . The reason for setting `drop_last=True` is so that if bM does not divide n , we do not get a smaller mini-batch that does not have enough samples to compute each group mean.

```

import numpy as np

from torch.utils.data import DataLoader, WeightedRandomSample

...

group_label_count = np.bincount(group_labels)

sample_weight = (1 / group_label_count)[group_labels] / len(group_label_count)

train_dataloader = DataLoader(
    train.dataset,
    sampler=WeightedRandomSampler(
        sample_weight,
        len(sample_weight),
        replacement=True
    ),
    batch_size=batch_size,
    drop_last=True
)

```

2.10.2 Convergence Analysis

The convergence analysis will employ similar tools to Section 2.6, but incorporate the noise resulting from imperfect estimates of the average loss of the parameter on each group. Using a smoothed group DRO objective will be instrumental in providing bias and variance control in terms of the per-group batch sizes. We introduce some notation that is used in the proof, which is analogous to the one used in the convergence analysis of Section 2.6. Let $\bar{r}_m(\cdot) := \bar{\ell}_m(\cdot) + \frac{\mu}{2} \|\cdot\|_2^2$ denote the regularized group-wise loss, defining the function $\bar{r} : \mathbb{R}^d \rightarrow \mathbb{R}^n$ via $\bar{r}(\boldsymbol{\theta}) := (\bar{r}_1(\boldsymbol{\theta}), \dots, \bar{r}_M(\boldsymbol{\theta}))$. We also define the iterates $\hat{\boldsymbol{\theta}}_i^{(k)} \in \mathbb{R}^d$ as those satisfying $l_i^{(k)} = \ell_i(\hat{\boldsymbol{\theta}}_i^{(k)})$ in Algorithm 2. Due to strong convexity and strong concavity, we

have that (2.51) has a unique saddle point $(\boldsymbol{\theta}^*, \mathbf{q}^*) \in \mathbb{R}^d \times \mathcal{Q}$ (recalling that $\mathcal{Q} \subseteq \mathbb{R}^M$). We maintain the following assumption throughout.

Assumption 2.10.2. The unregularized losses ℓ_i are convex, G_m -Lipschitz, and $(\bar{M} - \mu)$ -smooth for $\bar{M} \geq \mu$, $i \in \mathcal{I}_m$, and $m \in [M]$. The shift penalty $\text{Reg} : \mathcal{Q} \rightarrow \mathbb{R}$ is 1-strongly convex with respect to $\|\cdot\|_2$.

We emphasize the main difference between this setting and the instance-level incremental setting is that here, we observe a noisy version *every* element of the vector $\bar{\ell}(\boldsymbol{\theta}^{(k)})$ on each iteration, whereas before, we observed an exact version of *one* element of the vector $\ell(\boldsymbol{\theta}^{(k)})$. To handle this, we will require an additional variance assumption. Recall the notation

$$\nabla \ell_m^{(k+1)}(\boldsymbol{\theta}^{(k)}) = \frac{1}{b_m} \sum_{j=1}^{b_m} \nabla \ell_{i_{m,j}}(\boldsymbol{\theta}^{(k)})$$

from Algorithm 2. Conditioned on $\boldsymbol{\theta}^{(k)}$, the randomness in $\nabla \ell_m^{(k+1)}$ is governed by sampling the indices $i_{m,1}, \dots, i_{m,b_m}$ without replacement from $[n_m]$, motivating the assumption below. Let \mathbb{E}_k denote the conditional expectation given $\boldsymbol{\theta}^{(k)}$.

Assumption 2.10.3. There exists constants $\bar{\sigma}_1^2, \dots, \bar{\sigma}_M^2 \geq 0$ such that for all $k \geq 0$ and $m = 1, \dots, M$, it holds that

$$\mathbb{E}_k \left\| \nabla \ell_m^{(k+1)}(\boldsymbol{\theta}^{(k)}) - \nabla \bar{\ell}_m(\boldsymbol{\theta}^{(k)}) \right\|_2^2 \leq \frac{(n_m - b_m) \bar{\sigma}_m^2}{(n_m - 1) b_m}$$

Crucially, this assumption is made on the gradient of the *unregularized* losses. Under Assumption 2.10.2, these are bounded random variables which allow us to satisfy Assumption 2.10.3 uniformly over $\boldsymbol{\theta} \in \mathbb{R}^d$. To see this, fix $\boldsymbol{\theta}^{(k)} \in \mathbb{R}^d$ and compute directly (using the

finite population correction) the expression

$$\begin{aligned}\mathbb{E}_k \left\| \nabla \ell_m^{(k+1)}(\boldsymbol{\theta}^{(k)}) - \nabla \bar{\ell}_m(\boldsymbol{\theta}^{(k)}) \right\|_2^2 &= \frac{(n_m - b_m)}{(n_m - 1)b_m} \left(\frac{1}{n_m} \sum_{i \in \mathcal{I}_m} \left\| \nabla \ell_i(\boldsymbol{\theta}^{(k)}) - \nabla \bar{\ell}_m(\boldsymbol{\theta}^{(k)}) \right\|_2^2 \right) \\ &\leq \frac{(n_m - b_m)}{(n_m - 1)b_m} \left(\frac{1}{n_m} \sum_{i \in \mathcal{I}_m} \left\| \nabla \ell_i(\boldsymbol{\theta}^{(k)}) \right\|_2^2 \right) \\ &\leq \frac{(n_m - b_m)}{(n_m - 1)b_m} G_m^2,\end{aligned}$$

where the second line follows because the second moment upper bounds the variance, and the third line follows by Assumption 2.10.2. Thus, $\bar{\sigma}_m^2$ is always upper bounded by G_m^2 , but could be much smaller if the gradients within a group are large but concentrated about their mean. The last assumption is on the update function.

Assumption 2.10.4. Consider $\mathbf{v}^{(k+1)}$ as defined in (2.52). The Step function from line 10 of Algorithm 2 and line 6 of Algorithm 3 satisfies the decomposition

$$\mathbb{E}_k \left\| \boldsymbol{\theta}^{(k+1)} - \boldsymbol{\theta}^* \right\|_2^2 \leq \left\| \boldsymbol{\theta}^{(k)} - \boldsymbol{\theta}^* \right\|_2^2 - 2\eta C_1 \left\langle \bar{r}(\boldsymbol{\theta}^{(k)})^\top \mathbf{q}^{(k)}, \boldsymbol{\theta}^{(k)} - \boldsymbol{\theta}^* \right\rangle + \eta^2 C_2 \mathbb{E}_k \left\| \mathbf{v}^{(k+1)} \right\|_2^2.$$

In the case of the standard stochastic gradient update (2.52), Assumption 2.10.4 is satisfied as an equality with $C_1 = C_2 = 1$.

As before, we organize the analysis into a bound on the descent term (or bias analysis) and on the noise term (or variance analysis), and derive conditions on the learning rate $\eta > 0$. The following result is similar to Proposition 2.6.1.

Lemma 2.10.1 (Bias Analysis). *For any $\mathbf{q}, \mathbf{q}^* \in \Delta^{M-1}$ and $\boldsymbol{\theta}, \boldsymbol{\theta}^* \in \mathbb{R}^d$, it holds that*

$$\begin{aligned}-(\nabla \bar{r}(\boldsymbol{\theta})^\top \mathbf{q} - \nabla \bar{r}(\boldsymbol{\theta}^*)^\top \mathbf{q}^*)^\top (\boldsymbol{\theta} - \boldsymbol{\theta}^*) &\leq -(\mathbf{q} - \mathbf{q}^*)^\top (\bar{\ell}(\boldsymbol{\theta}) - \bar{\ell}(\boldsymbol{\theta}^*)) - \frac{\mu}{2} \left\| \boldsymbol{\theta} - \boldsymbol{\theta}^* \right\|_2^2 \\ &\quad - \frac{1}{2(\bar{M} + \mu)} \left\| (\nabla \bar{r}(\boldsymbol{\theta}) - \nabla \bar{r}(\boldsymbol{\theta}^*))^\top \mathbf{q} \right\|_2^2.\end{aligned}$$

Proof. The functions $\boldsymbol{\theta} \mapsto \bar{r}(\boldsymbol{\theta})^\top \mathbf{q}$ and $\boldsymbol{\theta} \mapsto \bar{r}(\boldsymbol{\theta})^\top \mathbf{q}^*$ are each \bar{M} -smooth and μ -strongly convex, as \mathbf{q} and \mathbf{q}^* are both contained in the probability simplex. By applying Lemma A.1.1,

we achieve the inequalities

$$\begin{aligned} \mathbf{q}^\top \bar{r}(\boldsymbol{\theta}^*) &\geq \mathbf{q}^\top \bar{r}(\boldsymbol{\theta}) + \mathbf{q}^\top \nabla \bar{r}(\boldsymbol{\theta})(\boldsymbol{\theta}^* - \boldsymbol{\theta}) \\ &\quad + \frac{1}{2(\bar{M} + \mu)} \|(\nabla \bar{r}(\boldsymbol{\theta}) - \nabla \bar{r}(\boldsymbol{\theta}^*))^\top \mathbf{q}\|_2^2 + \frac{\mu}{4} \|\boldsymbol{\theta} - \boldsymbol{\theta}^*\|_2^2 \end{aligned}$$

and

$$\begin{aligned} (\mathbf{q}^*)^\top \bar{r}(\boldsymbol{\theta}) &\geq (\mathbf{q}^*)^\top \bar{r}(\boldsymbol{\theta}^*) + (\mathbf{q}^*)^\top \nabla \bar{r}(\boldsymbol{\theta}^*)(\boldsymbol{\theta} - \boldsymbol{\theta}^*) \\ &\quad + \frac{1}{2(\bar{M} + \mu)} \|(\nabla \bar{r}(\boldsymbol{\theta}^*) - \nabla \bar{r}(\boldsymbol{\theta}))^\top \mathbf{q}^*\|_2^2 + \frac{\mu}{4} \|\boldsymbol{\theta}^* - \boldsymbol{\theta}\|_2^2. \end{aligned}$$

Summing the two inequalities gives

$$\begin{aligned} -(\nabla \bar{r}(\boldsymbol{\theta}))^\top \mathbf{q} - \nabla \bar{r}(\boldsymbol{\theta}^*)^\top \mathbf{q}^*)^\top (\boldsymbol{\theta} - \boldsymbol{\theta}^*) &\leq -(\mathbf{q} - \mathbf{q}^*)^\top (\bar{r}(\boldsymbol{\theta}) - \bar{r}(\boldsymbol{\theta}^*)) - \frac{\mu}{2} \|\boldsymbol{\theta} - \boldsymbol{\theta}^*\|_2^2 \\ &\quad - \frac{1}{2(\bar{M} + \mu)} \|(\nabla \bar{r}(\boldsymbol{\theta}) - \nabla \bar{r}(\boldsymbol{\theta}^*))^\top \mathbf{q}\|_2^2 \\ &\quad - \frac{1}{2(\bar{M} + \mu)} \|(\nabla \bar{r}(\boldsymbol{\theta}) - \nabla \bar{r}(\boldsymbol{\theta}^*))^\top \mathbf{q}^*\|_2^2. \end{aligned}$$

Drop the final non-positive term and mirror the argument of Proposition 2.6.1 to achieve

$$(\mathbf{q} - \mathbf{q}^*)^\top (\bar{r}(\boldsymbol{\theta}) - \bar{r}(\boldsymbol{\theta}^*)) = (\mathbf{q} - \mathbf{q}^*)^\top (\bar{\ell}(\boldsymbol{\theta}) - \bar{\ell}(\boldsymbol{\theta}^*)).$$

and complete the proof. \square

In the variance bound, we will also use the spectral norm, or largest singular value $s_{\max}(\nabla \bar{r}(\boldsymbol{\theta}^*))$ of the Jacobian of the regularized loss at the optimum. We only use this constant at the optimum, as the quantity can be unbounded over all of $\boldsymbol{\theta} \in \mathbb{R}^d$ due to regularization.

Lemma 2.10.2 (Variance Analysis). *It holds that*

$$\begin{aligned} \mathbb{E}_k \|\mathbf{v}^{(k+1)}\|_2^2 &\leq \max_{\mathbf{q} \in \mathcal{Q}} \|\mathbf{q}\|_2^2 \sum_{m=1}^M \frac{(n_m - b_m) \bar{\sigma}_m^2}{(n_m - 1) b_m} + 2 \|(\nabla \bar{r}(\boldsymbol{\theta}^{(k)}) - \nabla \bar{r}(\boldsymbol{\theta}^*))^\top \mathbf{q}^{(k)}\|_2^2 \\ &\quad + 2 s_{\max}(\nabla \bar{r}(\boldsymbol{\theta}^*))^2 \|\mathbf{q}^{(k)} - \mathbf{q}^*\|_2^2. \end{aligned}$$

Proof. Using the variance identity $\mathbb{E} \|X\|_2^2 = \mathbb{E} \|X - \mathbb{E}[X]\|_2^2 + \|\mathbb{E}[X]\|_2^2$, we first write

$$\begin{aligned}
& \mathbb{E}_k \|\mathbf{v}^{(k+1)}\|_2^2 \\
&= \mathbb{E}_k \|\mathbf{g}^{(k+1)} + \mu \boldsymbol{\theta}^{(k)}\|_2^2 \\
&= \mathbb{E}_k \|\mathbf{g}^{(k+1)} + \mu \boldsymbol{\theta}^{(k)} - \nabla \bar{r}(\boldsymbol{\theta}^{(k)})^\top \mathbf{q}^{(k)}\|_2^2 + \|\nabla \bar{r}(\boldsymbol{\theta}^{(k)})^\top \mathbf{q}^{(k)}\|_2^2 \\
&= \mathbb{E}_k \|\mathbf{g}^{(k+1)} - \nabla \bar{\ell}(\boldsymbol{\theta}^{(k)})^\top \mathbf{q}^{(k)}\|_2^2 + \|\nabla \bar{r}(\boldsymbol{\theta}^{(k)})^\top \mathbf{q}^{(k)}\|_2^2 \\
&\leq \mathbb{E}_k \|\mathbf{g}^{(k+1)} - \nabla \bar{\ell}(\boldsymbol{\theta}^{(k)})^\top \mathbf{q}^{(k)}\|_2^2 + 2 \|(\nabla \bar{r}(\boldsymbol{\theta}^{(k)}) - \nabla \bar{r}(\boldsymbol{\theta}^*))^\top \mathbf{q}^{(k)}\|_2^2 + 2 \|\nabla \bar{r}(\boldsymbol{\theta}^*)^\top \mathbf{q}^{(k)}\|_2^2.
\end{aligned}$$

We control the first and the third term above. For the first term, using the notation of Algorithm 2, use the vector-valued Cauchy-Schwarz inequality to achieve

$$\begin{aligned}
\mathbb{E}_k \|\mathbf{g}^{(k+1)} - \nabla \bar{\ell}(\boldsymbol{\theta}^{(k)})^\top \mathbf{q}^{(k)}\|_2^2 &= \mathbb{E}_k \left\| \sum_{m=1}^M q_m^{(k)} (\nabla \ell_m^{(k+1)}(\boldsymbol{\theta}^{(k)}) - \nabla \bar{\ell}_m(\boldsymbol{\theta}^{(k)})) \right\|_2^2 \\
&\leq \|\mathbf{q}^{(k)}\|_2^2 \sum_{m=1}^M \mathbb{E}_k \|\nabla \ell_m^{(k+1)}(\boldsymbol{\theta}^{(k)}) - \nabla \bar{\ell}_m(\boldsymbol{\theta}^{(k)})\|_2^2 \\
&\leq \|\mathbf{q}^{(k)}\|_2^2 \sum_{m=1}^M \frac{(n_m - b_m) \bar{\sigma}_m^2}{(n_m - 1) b_m},
\end{aligned}$$

where the last line follows from Assumption 2.10.3. For the third term, using that $\nabla \bar{r}(\boldsymbol{\theta}^*)^\top \mathbf{q}^* = \mathbf{0}$ at the optimum, we have that

$$\begin{aligned}
\|\nabla \bar{r}(\boldsymbol{\theta}^*)^\top \mathbf{q}^{(k)}\|_2^2 &= \|\nabla \bar{r}(\boldsymbol{\theta}^*)^\top (\mathbf{q}^{(k)} - \mathbf{q}^*)\|_2^2 \\
&\leq \|\nabla \bar{r}(\boldsymbol{\theta}^*)^\top\|_{2,2}^2 \|\mathbf{q}^{(k)} - \mathbf{q}^*\|_2^2,
\end{aligned}$$

and by using $\|\nabla \bar{r}(\boldsymbol{\theta}^*)^\top\|_{2,2} = \|\nabla \bar{r}(\boldsymbol{\theta}^*)\|_{2,2} = s_{\max}(\bar{r}(\boldsymbol{\theta}^*))$, we complete the proof. \square

Having established all the required bounds, we return to the overall analysis. We will adopt a Lyapunov stability argument for large shift costs, as in Section 2.6. The Lyapunov function will be

$$V^{(k)} = \|\boldsymbol{\theta}^{(k)} - \boldsymbol{\theta}^*\|_2^2 + \sum_{m=1}^M c_m \sum_{i \in \mathcal{I}_m} \|\hat{\boldsymbol{\theta}}_i^{(k)} - \boldsymbol{\theta}^*\|_2^2, \quad (2.53)$$

where c_1, \dots, c_M are to-be-specified constants. We can easily bound the evolution of the M

additional terms for a given batch size sequence, as for any $i \in \mathcal{I}_m$,

$$\mathbb{E}_k \left[\|\hat{\boldsymbol{\theta}}_i^{(k+1)} - \boldsymbol{\theta}^*\|_2^2 \right] = \frac{b_m}{n_m} \|\boldsymbol{\theta}^{(k)} - \boldsymbol{\theta}^*\|_2^2 + \left(\frac{n_m - b_m}{n_m} \right) \|\hat{\boldsymbol{\theta}}_i^{(k)} - \boldsymbol{\theta}^*\|_2^2$$

which then implies that

$$\mathbb{E}_k \left[\sum_{i \in \mathcal{I}_m} \|\hat{\boldsymbol{\theta}}_i^{(k+1)} - \boldsymbol{\theta}^*\|_2^2 \right] = b_m \|\boldsymbol{\theta}^{(k)} - \boldsymbol{\theta}^*\|_2^2 + \left(\frac{n_m - b_m}{n_m} \right) \sum_{i \in \mathcal{I}_m} \|\hat{\boldsymbol{\theta}}_i^{(k)} - \boldsymbol{\theta}^*\|_2^2. \quad (2.54)$$

The result below is a “large shift cost” regime analysis, where the condition on the smoothing parameter can in fact be optimized with respect to the per-group batch size parameter. The convergence will not be exact, but will have a convergence radius that scales with the learning rate.

Proposition 2.10.1. *Recalling the constants C_1 and C_2 from Assumption 2.10.4, assume that*

$$\eta \leq \min \left\{ \frac{C_1}{2C_2(\bar{M} + \mu)}, \frac{2C_1^2 \sum_{m=1}^M G_m^2}{C_2 s_*^2 \mu} \right\} \text{ and } \nu \geq \frac{8 \sum_{m'=1}^M G_{m'}^2}{C_1 \mu} \sqrt{\max_m n_m}.$$

Define the half-life

$$\tau := 2 \max \left\{ \frac{1}{\eta \mu}, \frac{n_1}{b_1}, \dots, \frac{n_M}{b_M} \right\}. \quad (2.55)$$

Then, by setting $c_m = \frac{\eta \mu n_m}{4M}$ for $m = 1, \dots, M$ in (2.53), it holds that

$$\mathbb{E}_k [V^{(k+1)}] \leq (1 - \tau^{-1}) V^{(k)} + \eta^2 \left(\max_{\mathbf{q} \in \mathcal{Q}} \|\mathbf{q}\|_2^2 \right) \sum_{m=1}^M \frac{(n_m - b_m) \bar{\sigma}_m^2}{(n_m - 1) b_m}.$$

Proof. By combining Assumption 2.10.4, Lemma 2.10.1, Lemma 2.10.2, and (2.54), we have

that

$$\begin{aligned} \mathbb{E}_k [V^{(k+1)}] &\leq \left(1 - \eta\mu + \sum_{m=1}^M c_m b_m\right) \|\boldsymbol{\theta}^{(k)} - \boldsymbol{\theta}^*\|_2^2 \\ &\quad - 2\eta C_1 (\mathbf{q}^{(k)} - \mathbf{q}^*)^\top (\bar{\ell}(\boldsymbol{\theta}^{(k)}) - \bar{\ell}(\boldsymbol{\theta}^*)) \end{aligned} \quad (2.56)$$

$$- \eta \left(\frac{C_1}{\bar{M} + \mu} - 2\eta C_2 \right) \|(\nabla \bar{r}(\boldsymbol{\theta}^{(k)}) - \nabla \bar{r}(\boldsymbol{\theta}^*))^\top \mathbf{q}^{(k)}\|_2^2 \quad (2.57)$$

$$\begin{aligned} &+ \eta^2 C_2 \left(\max_{\mathbf{q} \in \mathcal{Q}} \|\mathbf{q}\|_2^2 \right) \sum_{m=1}^M \frac{(n_m - b_m) \bar{\sigma}_m^2}{(n_m - 1) b_m} \\ &+ 2\eta^2 C_2 s_{\max} (\nabla \bar{r}(\boldsymbol{\theta}^*))^2 \|\mathbf{q}^{(k)} - \mathbf{q}^*\|_2^2 \\ &+ \sum_{m=1}^M c_m \left(\frac{n_m - b_m}{n_m} \right) \sum_{i \in \mathcal{I}_m} \|\hat{\boldsymbol{\theta}}_i^{(k)} - \boldsymbol{\theta}^*\|_2^2. \end{aligned} \quad (2.58)$$

The condition $\eta \leq \frac{C_1}{2C_2(\bar{M} + \mu)}$ implies that the term (2.57) is non-positive. To handle (2.56) and (2.58), we use the smoothness of $q^{\text{opt}}(\cdot)$ (via Lemma 2.4.1) and Assumption 2.10.2. First, apply Young's inequality with parameter $\alpha > 0$ to achieve

$$\begin{aligned} 2\eta(\mathbf{q}^{(k)} - \mathbf{q}^*)^\top (\bar{\ell}(\boldsymbol{\theta}^{(k)}) - \bar{\ell}(\boldsymbol{\theta}^*)) &\leq \eta\alpha \|\mathbf{q}^{(k)} - \mathbf{q}^*\|_2^2 + \frac{\eta}{\alpha} \|\bar{\ell}(\boldsymbol{\theta}^{(k)}) - \bar{\ell}(\boldsymbol{\theta}^*)\|_2^2 \\ &\leq \eta\alpha \|\mathbf{q}^{(k)} - \mathbf{q}^*\|_2^2 + \frac{\eta}{\alpha} (\sum_m G_m^2) \|\boldsymbol{\theta}^{(k)} - \boldsymbol{\theta}^*\|_2^2. \end{aligned}$$

We also have that

$$\|\mathbf{q}^{(k)} - \mathbf{q}^*\|_2^2 \leq \frac{1}{\nu^2} \sum_{m=1}^M G_m^2 \sum_{i \in \mathcal{I}_m} \|\hat{\boldsymbol{\theta}}_i^{(k)} - \boldsymbol{\theta}^*\|_2^2.$$

Combining these bounds together and letting $s_{\max}(\nabla \bar{r}(\boldsymbol{\theta}^*)) = s_*$, we have ultimately that

$$\begin{aligned} \mathbb{E}_k [V^{(k+1)}] &\leq \left[1 - \eta\mu + \sum_{m=1}^M \left(c_m b_m + \frac{\eta C_1 G_m^2}{\alpha} \right) \right] \|\boldsymbol{\theta}^{(k)} - \boldsymbol{\theta}^*\|_2^2 \\ &\quad + \sum_{m=1}^M c_m \left[1 + \frac{\eta G_m^2}{c_m \nu^2} (C_1 \alpha + 2\eta C_2 s_*^2) - \frac{b_m}{n_m} \right] \sum_{i \in \mathcal{I}_m} \|\hat{\boldsymbol{\theta}}_i^{(k)} - \boldsymbol{\theta}^*\|_2^2 \\ &\quad + \eta^2 C_2 \left(\max_{\mathbf{q} \in \mathcal{Q}} \|\mathbf{q}\|_2^2 \right) \sum_{m=1}^M \frac{(n_m - b_m) \bar{\sigma}_m^2}{(n_m - 1) b_m}. \end{aligned}$$

The convergence rate will be determined by setting the constants in the square brackets.

When setting $c_m = \frac{\eta G_m^2 \mu}{4b_m \sum_{m'} G_{m'}^2}$ and $\alpha = \frac{4C_1 \sum_m G_m^2}{\mu}$, we have that

$$1 - \eta\mu + \sum_{m=1}^M \left(c_m b_m + \frac{\eta C_1 G_m^2}{\alpha} \right) = 1 - \frac{\eta\mu}{2}$$

and

$$\frac{\eta G_m^2}{c_m \nu^2} (C_1 \alpha + 2\eta C_2 s_*^2) = \frac{8b_m \sum_{m'} G_{m'}^2}{\mu \nu^2} \left(\frac{2C_1^2 \sum_{m'} G_{m'}^2}{\mu} + \eta C_2 s_*^2 \right) \stackrel{\text{want}}{\leq} \frac{b_m}{2n_m}.$$

The inequality is accomplished when the conditions

$$\eta \leq \frac{2C_1^2 \sum_m G_m^2}{C_2 s_*^2 \mu} \text{ and } \nu \geq \frac{8 \sum_{m'=1}^M G_{m'}^2}{C_1 \mu} \sqrt{\max_m n_m}$$

hold, which achieves the desired rate and completes the proof. \square

We highlight the smoothness condition, which, by leveraging non-uniformity of G_1, \dots, G_M , scales as $\sqrt{\max_m n_m}$ times the sum of squared Lipschitz constants (as opposed to the total sample size times maximum Lipschitz constant squared seen in the instance-level analysis earlier in this chapter). We will observe a similar theme in Chapter 3, where we group coordinates of the loss vector and are able to sample entire blocks of coordinates on each iteration. To simplify the upcoming discussion, we assume that $C_1 = C_2 = 1$.

Next, iterating the result of Proposition 2.10.1 for a total iteration bound of t , we have

that

$$\begin{aligned}\mathbb{E}_0[V^{(t)}] &\leq (1 - \tau^{-1})^t V^{(0)} + \left(\sum_{k=0}^{t-1} (1 - \tau^{-1})^k \right) \eta^2 \left(\max_{\mathbf{q} \in \mathcal{Q}} \|\mathbf{q}\|_2^2 \right) \sum_{m=1}^M \frac{(n_m - b_m) \bar{\sigma}_m^2}{(n_m - 1) b_m} \\ &\leq (1 - \tau^{-1})^t V^{(0)} + \tau \eta^2 \left(\max_{\mathbf{q} \in \mathcal{Q}} \|\mathbf{q}\|_2^2 \right) \sum_{m=1}^M \frac{(n_m - b_m) \bar{\sigma}_m^2}{(n_m - 1) b_m}.\end{aligned}$$

While τ depends on the per-group batch sizes, if the dominant condition in (2.55) is the first, then the radius scales as

$$\frac{\eta}{\mu} \left(\max_{\mathbf{q} \in \mathcal{Q}} \|\mathbf{q}\|_2^2 \right) \sum_{m=1}^M \frac{(n_m - b_m) \bar{\sigma}_m^2}{(n_m - 1) b_m},$$

which, in terms of the $O(\eta/\mu)$ dependence, resembles state-of-the-art stochastic gradient descent analyses in similar settings [Cutler et al., 2023]. At a high level, the notion of training examples partitioned into logical groups and adapting to their non-uniformity is a theme that will be discussed further in Chapter 3. Here, we leveraged the per-group batch sizes and operated in a primal-only SGD-like framework, pursuing a practical and simple algorithm. As alluded to above, in Chapter 3, we pursue other ideas such as adaptive sampling and block coordinate-wise updates to achieve improved complexities in the primal-dual setting.

2.10.3 Experiments

We provide a numerical benchmark to apply Algorithm 3 in practice. As mentioned before, while Algorithm 2 and Algorithm 3 perform almost equivalently in practice, Algorithm 3 is much easier to implement for practitioners operating within an existing PyTorch workflow, as the indices of data points may not be provided by the sampling mechanism. Consider the following setup, which resembles the instance-level distribution shift benchmark from Figure 2.6 (Section 2.9).

- **Data:** A subset of 100,000 points from the Amazon Reviews dataset from the WILD Distribution Shift Benchmark. They are split into 50,000 train and 50,000 test examples, respectively. Each data point represents a review of an Amazon product written

in natural language, along with a categorical label indicating the number of stars given from 1 to 5. The review is augmented with metadata, including the product index, product category, and year. The text is discretized using the `bert-base-uncased` tokenizer and truncated to 100 tokens. The group labels are given by the product category, and may shift in proportion between the train and test sets.

- **Model:** The model architecture is a 2-layer transformer network, where each transformer block uses 8 heads, an embedding dimension of 512, and pre-normalization.
- **Optimizer:** Algorithm 3 with the base optimizer Adam with default moment and variance parameters (β_1, β_2) are used without weight decay and a fixed learning rate of 3×10^{-4} . The spectrum is chosen to be either the 1.5 or 2-extremile compared to a group-wise empirical risk minimization baseline (or $\mathcal{Q} = \{\mathbf{1}/M\}$).

The results are shown in Figure 2.10. We observe that the average performance over the test set is indistinguishable between the various objectives, which include group-wise empirical risk minimization, as well as the 1.5-extremile and 2-extremile. However, the worst group-wise accuracy for each of the spectral risk measure-based group DRO models achieves a 5% increase over the ERM baseline.

To motivate the upcoming discussion on the relationship between worst group-wise error and fairness, compare the experiment above to Figure 2.5. We used two datasets, and in `diabetes`, we use gender as the protected attribute A , whereas in `acsincome` we use race as the protected attribute. It is natural to consider the protected attribute to be a group label, and for statistical parity (SP) to be a measure of worst-case group-wise performance. However, the instance-level DRO algorithms, while evaluated on SP, were never supplied the protected attribute (i.e., the group label) at training time, which is a common constraint in fairness benchmarks. This differs from the models explored in this section, for which observing the group label is an essential component of the algorithm. If the protected attribute is available to the model, then group DRO presents an alternative formulation that

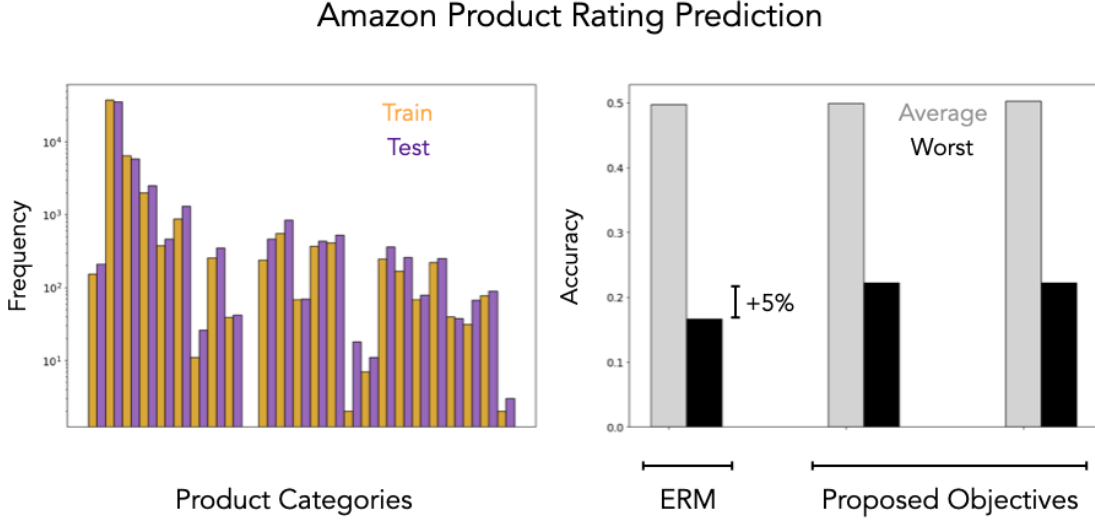


Figure 2.10: **Amazon Reviews Data.** The left panel shows the group-wise proportions in the train and test sets, respectively. The right panel shows the average and worst-case group-wise test accuracy. The proposed objectives are the 1.5 or 2-extremile. Even though the group proportions have a minor shift in most product categories, a moderate worst-case accuracy increase is observed for the distributionally robust variants.

explicitly promotes fairness via the group reweighting structure. We provide details on these group-wise metrics and others in the next subsection.

2.10.4 Other Group-Wise Performance Metrics

We conclude this section by commenting on our group DRO formulation in the context of other group-level performance measures, such as class-specific performance and algorithmic fairness. Consider the setting of supervised learning, where $\Xi = \mathcal{X} \times \mathcal{Y}$ for feature space \mathcal{X} and label space \mathcal{Y} . Furthermore, the sake of exposition, assume that $\mathcal{Y} = \{0, 1\}$, so that for a parametrized function $f_{\theta} : \mathcal{X} \rightarrow (0, 1)$ and $\xi = (X, Y)$, we define

$$\ell(\theta, \xi) := Y \log(f_{\theta}(X)) + (1 - Y) \log(1 - f_{\theta}(X)).$$

We interpret $f_{\theta}(\mathbf{x})$ as the predicted probability that $Y = 1$ given $X = \mathbf{x}$. Accordingly, define the predicted value $\hat{Y}_{\theta}(X) = \mathbb{1}\{f_{\theta}(X) \geq 0.5\}$. One could, in principle, set $A_1 =$

$Y_1, \dots, A_n = Y_n$, i.e., the group labels are themselves the class labels. Thus, if a class y is harder under θ , or

$$\mathbb{P}_{(X,Y) \sim P} [\hat{Y}_\theta(X) = y | Y = y] \leq \mathbb{P}_{(X,Y) \sim P} [\hat{Y}_\theta(X) = 1 - y | Y = 1 - y],$$

then reweighting the group labels toward this label would balance performance across classes. The quantities on the left and right-hand side above are, in fact, the *recall* for classes y and $1 - y$, respectively. Typically, the reweighting parameter might be selected using a held-out validation set. However, to adapt to unknown class imbalances, the group DRO formulation of (2.51) offers a principled solution that simultaneously incorporates many possible class distributions (or *label shifts*) that may be observed during deployment.

Relatedly, in the fairness literature, a common claim of success is that some notion of error is approximately equal across groups. As mentioned above, the group label is more commonly referred to as the “protected attribute”. Under $(X, Y) \sim P$, two examples of parameter choice $\theta \in \mathbb{R}^d$ that are fair under *equalized odds* are those that satisfy either false positive or false negative invariance to the attribute A :

$$\begin{aligned} \mathbb{P} [\hat{Y}_\theta(X) = 1 | Y = 0, A = \mathbf{a}_1] &\approx \dots \approx \mathbb{P} [\hat{Y}_\theta(X) = 1 | Y = 0, A = \mathbf{a}_M] \\ \mathbb{P} [\hat{Y}_\theta(X) = 0 | Y = 1, A = \mathbf{a}_1] &\approx \dots \approx \mathbb{P} [\hat{Y}_\theta(X) = 0 | Y = Y, A = \mathbf{a}_M]. \end{aligned}$$

This motivates the use of (2.51), as a parameter with approximately equal performance on the group conditional risk $\mathbb{E}_P [\ell(\theta, \xi) | A = \mathbf{a}]$ across all $\mathbf{a} \in \mathcal{A}$ cannot be largely increased by reweighting. Note that we explored this connection experimentally at the instance-level in Section 2.9.

Finally, there is an alternative utilization of group information which is similar to (2.51) algorithmically but distinct in terms of the underlying statistical problem. Observe that at an instance level, (2.51) applies equal weight to all examples in the group, where the distribution over groups is learned. The assumption that the within-group distribution on Ξ is uniform regardless of the distributional shift. Consider instead having *known* target

weights q_1^*, \dots, q_M^* over groups, but no reasonable justification that the per-group distribution over Ξ is uniform. This setting can be formulated as an instance-level DRO problem with uncertainty set

$$\mathcal{B} := \left\{ \beta \in \mathcal{B}^{\text{DRO}} : \int_{\Xi \times \{\mathbf{a}_m\}} \beta(\mathbf{x}, \mathbf{a}) \, dP_n(\mathbf{x}, \mathbf{a}) = q_m^* \, \forall m \in [M] \right\}.$$

This viewpoint is a combination of the “known” distribution framework direct likelihood estimation (see Section 2.8) and the “unknown” target distribution. We may know our target distribution up to a partition of the space Ξ of resolution M with probabilities $\{q_i^*\}_{i=1}^M$, but wish to account for many possible distributions within each cell of the partition.

2.11 Possible Extensions

2.11.1 Changes in the Support

We consider here a version of the objective in which the learned parameter can account for changes in the support of the input distribution, instead of just changes in the relative weights on each training example. Consider the squared loss

$$\ell(\boldsymbol{\theta}, (\mathbf{x}, y)) := \frac{1}{2} (y - \langle \boldsymbol{\theta}, \mathbf{x} \rangle)^2,$$

which is not Lipschitz continuous unless $\boldsymbol{\theta}$ is restricted to a compact domain. We perturb \mathbf{x} by a random vector $\boldsymbol{\varepsilon}$ realized in \mathbb{R}^d , whose mean is denoted $\boldsymbol{\mu}$ and whose covariance matrix is denoted $\boldsymbol{\Sigma} \in \mathbb{R}^{d \times d}$. We consider the maximum expected loss achievable by perturbation distributions whose mean and covariance are bounded in norm-balls of radius r and Δ , respectively. The $\boldsymbol{\mu}$ is bounded in the Euclidean norm $\|\cdot\|_2$ and the covariance is bounded in spectral norm $\|\cdot\|_{2,2}$. While the specific distribution of $\boldsymbol{\varepsilon}$ does not affect the objective beyond the first-two moments, we let $\boldsymbol{\varepsilon} \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ for simplicity. Given data points $\{\xi_i\}_{i=1}^n$

for $\xi_i = (\mathbf{x}_i, y_i) \in \mathbb{R}^d \times \mathbb{R}$, this gives the objective

$$\mathcal{L}_{r,\Delta}(\boldsymbol{\theta}) := \max_{\mathbf{q} \in \mathcal{Q}} \max_{\substack{\boldsymbol{\mu}: \|\boldsymbol{\mu}\|_2 \leq r, \\ \boldsymbol{\Sigma}: \|\boldsymbol{\Sigma}\|_{2,2} \leq \Delta}} \sum_{i=1}^n \frac{q_i}{2} \mathbb{E}_{\boldsymbol{\varepsilon}_i \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})} [(y_i - \langle \boldsymbol{\theta}, \mathbf{x}_i + \boldsymbol{\varepsilon}_i \rangle)^2], \quad (2.59)$$

which is still convex in $\boldsymbol{\theta}$. By expanding this term, we can in fact remove the optimization over $\boldsymbol{\Sigma}$, whereas the dependence on $\boldsymbol{\mu}$ can be simplified using a univariate optimization problem.

Proposition 2.11.1. *For any $r, \Delta > 0$, it holds that*

$$\mathcal{L}_{r,\Delta}(\boldsymbol{\theta}) := \max_{\mathbf{q} \in \mathcal{Q}} \left[\langle \mathbf{q}, \ell(\boldsymbol{\theta}) \rangle + \max \left\{ 0, r (\langle \boldsymbol{\theta}, \bar{\mathbf{x}}_{\mathbf{q}} \rangle - \bar{y}_{\mathbf{q}}) \|\boldsymbol{\theta}\|_2 + \frac{r^2}{2} \|\boldsymbol{\theta}\|_2^2 \right\} \right] + \frac{\Delta}{2} \|\boldsymbol{\theta}\|_2^2 \quad (2.60)$$

for $\ell_i(\boldsymbol{\theta}) = \frac{1}{2} (y_i - \langle \boldsymbol{\theta}, \mathbf{x}_i \rangle)$, $\bar{y}_{\mathbf{q}} = \sum_{i=1}^n q_i y_i$ and $\bar{\mathbf{x}}_{\mathbf{q}} = \sum_{i=1}^n q_i \mathbf{x}_i$. Furthermore, the values of $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma}$ achieving (2.60) for $\boldsymbol{\theta} \neq \mathbf{0}$ are given by

$$\begin{aligned} \boldsymbol{\mu}^*(\boldsymbol{\theta}) &= \left(\arg \max_{a \in \{0, r\}} a (\langle \boldsymbol{\theta}, \bar{\mathbf{x}}_{\mathbf{q}} \rangle - \bar{y}_{\mathbf{q}}) \|\boldsymbol{\theta}\|_2 + \frac{a^2}{2} \|\boldsymbol{\theta}\|_2^2 \right) \cdot \frac{\boldsymbol{\theta}}{\|\boldsymbol{\theta}\|_2} \\ \boldsymbol{\Sigma}^*(\boldsymbol{\theta}) &= \frac{\Delta}{\|\boldsymbol{\theta}\|_2^2} \boldsymbol{\theta} \boldsymbol{\theta}^\top + \boldsymbol{\Sigma}_0, \end{aligned}$$

where $\boldsymbol{\Sigma}_0$ is any positive semidefinite matrix satisfying $\text{range}(\boldsymbol{\Sigma}_0) = \text{span}\{\boldsymbol{\theta}\}^\perp$.

Proof. First, expand the summands of the objective (2.59) to achieve

$$\begin{aligned} & \frac{1}{2} \mathbb{E}_{\boldsymbol{\varepsilon}_i \sim \mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})} [(y_i - \langle \boldsymbol{\theta}, \mathbf{x}_i + \boldsymbol{\varepsilon}_i \rangle)^2] \\ &= \frac{1}{2} (y_i - \langle \boldsymbol{\theta}, \mathbf{x}_i \rangle)^2 - (y_i - \langle \boldsymbol{\theta}, \mathbf{x}_i \rangle) \mathbb{E}_{\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})} [\langle \boldsymbol{\theta}, \boldsymbol{\varepsilon}_i \rangle] + \frac{1}{2} \mathbb{E}_{\mathcal{N}(\boldsymbol{\mu}, \boldsymbol{\Sigma})} [\langle \boldsymbol{\theta}, \boldsymbol{\varepsilon}_i \rangle^2] \\ &= \frac{1}{2} (y_i - \langle \boldsymbol{\theta}, \mathbf{x}_i \rangle)^2 - (y_i - \langle \boldsymbol{\theta}, \mathbf{x}_i \rangle) \langle \boldsymbol{\theta}, \boldsymbol{\mu} \rangle + \frac{1}{2} \boldsymbol{\theta}^\top \boldsymbol{\Sigma} \boldsymbol{\theta} + \frac{1}{2} \langle \boldsymbol{\theta}, \boldsymbol{\mu} \rangle^2. \end{aligned}$$

Thus, it holds that $\mathcal{L}_{r,\Delta}(\boldsymbol{\theta})$ is equal to

$$\max_{\mathbf{q} \in \mathcal{Q}} \left\{ \langle \mathbf{q}, \ell(\boldsymbol{\theta}) \rangle + \max_{\boldsymbol{\mu}: \|\boldsymbol{\mu}\|_2 \leq r} (\langle \boldsymbol{\theta}, \bar{\mathbf{x}}_{\mathbf{q}} \rangle - \bar{y}_{\mathbf{q}}) \langle \boldsymbol{\theta}, \boldsymbol{\mu} \rangle + \frac{1}{2} \langle \boldsymbol{\theta}, \boldsymbol{\mu} \rangle^2 \right\} + \frac{1}{2} \max_{\boldsymbol{\Sigma}: \|\boldsymbol{\Sigma}\|_{2,2} \leq \Delta} \boldsymbol{\theta}^\top \boldsymbol{\Sigma} \boldsymbol{\theta}$$

The final term above can be computed in closed form, as the objective depends on $\boldsymbol{\Sigma}$ only based on evaluation on the span of $\boldsymbol{\theta}$. Thus, one solution (in the case that $\boldsymbol{\theta} \neq \mathbf{0}$) of the

maximization problem is $\Sigma := \frac{\Delta}{\|\theta\|_2^2} \theta \theta^\top$, yielding

$$\frac{1}{2} \max_{\Sigma: \|\Sigma\|_{2,2} \leq \Delta} \theta^\top \Sigma \theta = \frac{\Delta}{2} \|\theta\|_2^2,$$

which is simply a regularization term. The equality above also holds when $\theta = \mathbf{0}$. Note that even though the expectation was computed under a Gaussian distribution (which requires Σ to be full rank), we may make it full rank by simply completing the singular value decomposition with arbitrary positive singular values that are strictly less than Δ . On the other hand, we see that the objective also depends on μ only through $\langle \theta, \mu \rangle$ (also assuming that $\theta \neq 0$), we may reparametrize the optimization problem by considering

$$\mu = a \cdot \frac{\theta}{\|\theta\|_2}$$

for $a \in [0, r]$. To compute μ , we must then solve the univariate problem

$$\max_{0 \leq a \leq r} a (\langle \theta, \bar{x}_q \rangle - \bar{y}_q) \|\theta\|_2 + \frac{a^2}{2} \|\theta\|_2^2 = \max \left\{ 0, r (\langle \theta, \bar{x}_q \rangle - \bar{y}_q) \|\theta\|_2 + \frac{r^2}{2} \|\theta\|_2^2 \right\},$$

where the solution follows by simply considering the two endpoints, as the objective it is the maximum of a strongly convex function over an interval. \square

In order to preserve the dual linearity of the saddle-point objective in q , one could alternatively maximize over q and μ in the inner-most objective by considering $\mathcal{L}_{r,\Delta}(\theta) = \max_{q,\mu} \tilde{\mathcal{L}}(\theta, q, \mu)$ for

$$\begin{aligned} \tilde{\mathcal{L}}(\theta, q, \mu) &= \langle q, \ell(\theta) \rangle + (\langle \theta, \bar{x}_q \rangle - \bar{y}_q) \langle \theta, \mu \rangle + \frac{1}{2} \langle \theta, \mu \rangle^2 + \frac{\Delta}{2} \|\theta\|_2^2 \\ &= \langle q, \ell(\theta) \rangle - \bar{y}_q \mu^\top \theta + \frac{1}{2} \theta^\top ((2\bar{x}_q + \mu) \mu^\top + \Delta I) \theta \end{aligned}$$

as the solution for μ is computable in closed form. Thus, by modifying the optimization problem only slightly, we can account for some degree of robustness to input noise, effectively changing the support of the input distribution.

2.11.2 A Shrinking Uncertainty Set

As mentioned in Section 2.8, there are several applications (such as empirical likelihood) in which the size of the uncertainty set is dependent on the sample size n . If one wishes to use an estimator learned from a distributionally robust objective that also decreases the amount of uncertainty as more data arrives, it is worthwhile to know the necessary conditions on this decay to match various rates.

To do so, we change notation slightly to better agree with standard generalization bounds (e.g., Wainwright [2019]) and consider a generic class $f \in \mathcal{F}$ of measurable functions $f : \Xi \rightarrow \mathbb{R}$ (indicating losses) and consider an uncertainty set $\mathcal{Q}_n \subseteq \Delta^{n-1}$ dependent on n . Finally, consider an empirical DR risk minimizer

$$\hat{f}_n \in \arg \min_{f \in \mathcal{F}} \left[R_n(f) := \max_{\mathbf{q} \in \mathcal{Q}_n} \sum_{i=1}^n q_i f(\xi_i) - \nu \text{Reg}(\mathbf{q}) \right].$$

We also employ the standard empirical process notation $P(f) := \mathbb{E}_{\xi \sim P} [f(\xi)]$ and let P_n be the empirical measure of the data. We have the following guarantee.

Lemma 2.11.1. *Assume that $f(\xi) \leq B$ with P -probability one for every $f \in \mathcal{F}$ and that $\text{Reg}(\mathbf{1}/n) = 0$. Then, for any $f \in \mathcal{F}$, it holds that*

$$P(\hat{f}_n) - P(f) \leq (P - P_n)(\hat{f}_n - f) + B \max_{\mathbf{q} \in \mathcal{Q}_n} \|\mathbf{q} - \mathbf{1}/n\|_1. \quad (2.61)$$

Proof. First, write the decomposition

$$\begin{aligned} P(\hat{f}_n) - P(f) &= (P - P_n)(\hat{f}_n) + \underbrace{(P_n - R_n)(\hat{f}_n)}_{\leq 0} + \underbrace{R_n(\hat{f}_n) - R_n(f)}_{\leq 0} \\ &\quad + (R_n - P_n)(f) + (P_n - P)(f) \\ &\leq (P - P_n)(\hat{f}_n - f) + (R_n - P_n)(f), \end{aligned}$$

where we used in the first line that $(P_n - R_n)(f) \leq 0$ for all $f \in \mathcal{F}$ because $\text{Reg}(\mathbf{1}/n) = 0$,

whereas $R_n(\hat{f}_n) - R_n(f) \leq 0$ due to the optimality of \hat{f}_n . Next, write

$$\begin{aligned} (R_n - P_n)(f) &= \max_{\mathbf{q} \in \mathcal{Q}_n} \sum_{i=1}^n (q_i - 1/n) f(\xi_i) - \nu \text{Reg}(\mathbf{q}) \\ &\leq \max_{\mathbf{q} \in \mathcal{Q}_n} \sum_{i=1}^n (q_i - 1/n) f(\xi_i) \\ &\leq B \max_{\mathbf{q} \in \mathcal{Q}_n} \|\mathbf{q} - \mathbf{1}/n\|_1, \end{aligned}$$

completing the proof. \square

Notice that (2.61) contains the standard empirical process term, which can be bounded in a variety of ways, which do not necessarily depend on the fact that \hat{f}_n is the minimizer of the average loss. For example, the standard “slow rate” approach would immediately apply

$$(P - P_n)(\hat{f}_n - f) \leq 2 \sup_{f \in \mathcal{F}} |(P_n - P)f|,$$

where we assume that the right-hand side remains measurable. For a localization-based approach, first assume that there exists a unique minimizer

$$f_\star := \arg \min_{f \in \mathcal{F}} P(f).$$

Define the following *excess risk* term as

$$\mathcal{E}_n(f) := (P - P_n)(f - f_\star).$$

This term differs from the usual excess risk process $f \mapsto P(f - f_\star)$, but is defined in this way to apply the arguments of, for instance, [Ohn and Kim \[2025, Section 2\]](#) to control it. Note that these do not depend on the fact that \hat{f}_n was generated by optimizing the average loss, so they can be applied for this setting. In various settings, this can be shown to be bounded by an $O(1/n)$ term with high probability. Thus, it only remains to control the second term of (2.61). In the case of spectral risk measures, we may achieve the following bounds in terms of Kullback-Liebler (KL) or χ^2 -divergence.

Lemma 2.11.2. *Let \mathcal{Q}_n be an ambiguity set determined by a spectral risk measure with spectrum s_n on $(0, 1)$ (see Section 2.3). Then, it holds that*

$$\max_{\mathbf{q} \in \mathcal{Q}_n} \|\mathbf{q} - \mathbf{1}/n\|_1 \leq \sqrt{\frac{1}{2} \int_0^1 s_n(t) \log s_n(t) dt}.$$

Proof. First, the maximizer of the strongly convex function $\mathbf{q} \mapsto \|\mathbf{q} - \mathbf{1}/n\|_1^2$ occurs on the boundary of the closed, convex permutahedron $\mathcal{Q}_n \equiv \mathcal{Q}(\sigma)$. Furthermore, because this objective is permutation invariant, it holds that

$$\max_{\mathbf{q} \in \mathcal{Q}_n} \|\mathbf{q} - \mathbf{1}/n\|_1 = \|\sigma - \mathbf{1}/n\|_1.$$

Then, it holds that

$$\begin{aligned} \|\sigma - \mathbf{1}/n\|_1 &= \sum_{i=1}^n \left| \int_{(i-1)/n}^{i/n} s_n(t) dt - 1/n \right| \\ &= \sum_{i=1}^n \left| \int_{(i-1)/n}^{i/n} (s_n(t) - 1) dt \right| \\ &\leq \sum_{i=1}^n \int_{(i-1)/n}^{i/n} |s_n(t) - 1| dt \\ &= \|s_n - 1\|_1. \end{aligned}$$

By an immediate application of Pinsker's inequality,

$$\|s_n - 1\|_1 \leq \sqrt{\frac{1}{2} \text{KL}(s_n \| 1)},$$

where the number 1 indicates the uniform distribution on $(0, 1)$. Then, compute the KL-divergence term and use that $\log(1) = 0$ to complete the proof. \square

We may use the upper bound in the lemma above to compute the decay of the term to zero as $n \rightarrow \infty$, and ensure that it matches the first term of (2.61). To see this, let us recall the examples from Section 2.3 and determine what values of the risk parameter are necessary to enforce $\max_{\mathbf{q} \in \mathcal{Q}_n} \|\mathbf{q} - \mathbf{1}/n\|_1 = O(n^{-\alpha})$.

- **Superquantile:** For $\tau_n \in (0, 1]$, the τ_n -superquantile [Rockafellar and Royset, 2013], is specified by $s_n(t) = \frac{1}{1-\tau_n} \mathbb{1} \{ \tau_n \leq t \leq 1 \}$. Then, it holds that

$$\sqrt{\int_0^1 s_n(t) \log s_n(t) dt} = \sqrt{\log \left(\frac{1}{1-\tau_n} \right)} \implies \tau_n = O \left(1 - e^{-(1/n^{2\alpha})} \right).$$

- **Extremile:** For $r_n > 1$, the r_n -extremile [Daouia et al., 2019] is specified by $s_n(t) = r_n t^{r_n-1}$. Then, it holds that

$$\sqrt{\int_0^1 s_n(t) \log s_n(t) dt} = \sqrt{\log r_n - \left(\frac{r_n-1}{r_n} \right)} \implies r_n = O \left(e^{(1/n^{2\alpha})} \right).$$

For $\alpha = 1$, these indicate extremely fast convergence of $\tau_n \rightarrow 0$ or $r_n \rightarrow 1$ in order to match the “fast” rate of $O(1/n)$ for the excess risk term.

One avenue for future work in this area is to establish fast rates for a fixed uncertainty set, describing convergence to the minimizer of the population (distributionally robust) risk. There is work in this area showing matching upper and lower bounds in the case of a particular family of f -divergences [Duchi and Namkoong, 2021], but similar results for spectral risk measure uncertainty sets remain an open question.

2.11.3 Concentration under Heavy Tails

In risk-averse applications, a common consideration is the heavy-tailedness of the distribution being analyzed. As commented in Mikosch [1999], while the term “heavy-tailed” does not have a universal definition, two classes of distributions are ubiquitous cases of interest: random variables with regularly varying tails, and sub-exponential random variables. They are defined below.

Definition 2.11.1 (Regularly Varying Tails). For a real-valued random variable X with CDF F , we say that X has *regularly varying tails* with index $\alpha \in \mathbb{R}$ if

$$\lim_{x \rightarrow \infty} \frac{1 - F(tx)}{1 - F(x)} = t^\alpha \text{ for all } t > 0.$$

Definition 2.11.2 (Sub-Exponential R.V.). We say that a real-valued, non-negative random variable X is *sub-exponential* if there exist positive constants a and b such that

$$\log \mathbb{P}[X \geq x] \leq \log a - bx.$$

The relationship between the two classes of random variables is the following. Let X_1, \dots, X_n be independent and identically distributed copies of a real-valued, non-negative random variable X . For function f, g , we write $f(x) \sim g(x)$ to mean that $\lim_{x \rightarrow \infty} f(x)/g(x) = 1$. It holds that both random variables with regularly varying tails and sub-exponential random variables satisfy

$$\mathbb{P}[\sum_{i=1}^n X_i > x] \sim n\mathbb{P}[X > x] \sim \mathbb{P}[\max_i X_i > x] \text{ as } x \rightarrow \infty. \quad (2.62)$$

In fact, (2.62) is an equivalent (albeit qualitative) condition to being sub-exponential [Mikosch, 1999]. In words, (2.62) captures the intuition that the maximum of a collection of i.i.d. heavy-tailed random variables contributes nearly all of the value in the sum. As such, the sub-exponential distribution is seen as a *boundary* between light- and heavy-tailed distributions.

Recently, Vladimirova et al. [2020] and Kuchibhotla and Chakraborty [2022] independently proposed a class of random variables with Weibull-like tails, which has connections to both notions of heavy-tailed distributions given above.

Definition 2.11.3 (Sub-Weibull R.V.). We say that a real-valued, non-negative random variable X is *sub-exponential* if there exist positive constants a , b , and ζ such that

$$\log \mathbb{P}[X \geq x] \leq \log a - bx^\zeta.$$

Equivalently, we may have parameters (K, ζ) such that

$$\|X\|_{\mathbf{L}^p(P)} \leq Kp^{1/\zeta}.$$

It can clearly be seen that by setting $\zeta = 1$, we recover sub-exponentiality. On the other

hand, when X is bounded by B , we have that if

$$\lim_{x \rightarrow \infty} \frac{1 - F(B - (tx)^{-1})}{1 - F(B - x^{-1})} = t^{-\alpha} \text{ for all } t > 0.$$

for $\alpha \geq 0$, then we have that the distribution of X is in the maximum domain of attraction of the Weibull distribution with parameter α (often denoted Ψ_α in the extreme value theory literature). Thus, the random variables satisfying Definition 2.11.3 for $\zeta \geq 1$ are relevant representatives of heavy-tailed distributions.

We use two steps for concentration results for spectral risk measures applied to collections of sub-Weibull random variables. To do so, we first relate spectral risk measures to scalings of standard uniform averages. Then, we appeal to concentration bounds for sums of independent sub-Weibull random variables.

Lemma 2.11.3. *Let $x_1, \dots, x_n \in [0, +\infty)$ define empirical CDF $F_n(x) := \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{(-\infty, x_i]}(x)$ and empirical quantile function $F_n^{-1}(p) := \inf \{x : F_n(x) \geq p\}$. Then, for any integrable spectrum $s : [0, 1] \rightarrow [0, +\infty)$, we have that*

$$\mathbb{L}_s[F_n] = \sum_{i=1}^n \sigma_i x_{(i)} \leq \frac{s(1)}{n} \sum_{i=1}^n x_i,$$

where $\sigma_i := \int_{(i-1)/n}^{i/n} s(p) dp$ and $x_{(1)} \leq \dots \leq x_{(n)}$ are the order statistics of (x_1, \dots, x_n) .

Proof. Note that we may express the quantile function of an empirical measure using the order statistics of the sample. Then,

$$\mathbb{L}_s[F_n] = \int_0^1 s(p) F_n^{-1}(p) dp = \sum_{i=1}^n \left(\int_{(i-1)/n}^{i/n} s(p) x_{(\lceil np \rceil)} dp \right) = \sum_{i=1}^n \left(\int_{(i-1)/n}^{i/n} s(t) dt \right) x_{(i)}.$$

. To prove the first claim, we compute the integral

$$\sum_{i=1}^n \left(\int_{(i-1)/n}^{i/n} s(t) dt \right) x_{(i)} = \sum_{i=1}^n \sigma_i x_{(i)}.$$

For the second claim, we apply monotonicity of s and non-negativity of x_1, \dots, x_n to upper

bound the integral via

$$\sum_{i=1}^n \left(\int_{(i-1)/n}^{i/n} s(t) \, dp \right) x_{(i)} \leq \frac{s(1)}{n} \sum_{i=1}^n x_{(i)} = \frac{s(1)}{n} \sum_{i=1}^n x_i.$$

□

Thus, spectral risk measures are at most a constant factor larger than the simple average. For the superquantile, we have that $s(1) = \frac{1}{1-\tau}$, whereas for the extremile, we have that $s(1) = r$. From Lemma 2.11.3, we can derive a tail bound for sub-Weibull random variables.

Next, we present a concentration inequality for sub-Weibull random variables proved in Vladimirova et al. [2020]. Before stating the result, we note that tail bounds proved for sums of sub-Weibull random variables in Vladimirova et al. [2020] does not yield concentration of averages, while that of Kuchibhotla and Chakraborty [2022] does with confidence band of size $O(n^{-1/2})$. In fact, Zhang and Wei [2022] eventually proved sharper concentration rates in terms of problem constants (and the same dependence on n), but with less interpretable constant. Consider the following claim.

Proposition 2.11.2. *Let $X_1, \dots, X_n \sim P$ be i.i.d. non-negative sub-Weibull random variables with parameters K and $\zeta < 1$. Then, for any $t \geq s(1) \max\{eK, 1\}$, it holds that*

$$\mathbb{P} [\mathbb{L}_s[F_n] \geq t] \leq Ke \exp \left(-\frac{n}{eK} \left(\frac{t}{s(1)} \right)^{\min\{2, 1/\zeta\}} \right).$$

Proof. First, apply Lemma 2.11.3 so that for any $t \geq 0$, we have that

$$\mathbb{P} [\mathbb{L}_s[F_n] \geq t] \leq \mathbb{P} \left[\sum_{i=1}^n X_i \geq nt/s(1) \right].$$

We may then apply Vladimirova et al. [2020, Eq. (7)] (note that K_ζ in their notation is equal to eK in ours) and the condition that $t \geq eKs(1)$ to achieve

$$\mathbb{P} \left[\sum_{i=1}^n X_i \geq nt/s(1) \right] \leq Ke \exp \left(-\frac{n}{eK} \min \left\{ \frac{t^2}{s^2(1)}, \frac{t^{1/\zeta}}{s^{1/\zeta}(1)} \right\} \right).$$

We use that $t \geq s(1)$ to pass the minimum to the exponent to complete the proof. □

From Proposition 2.11.2, we see that the spectrum can scale the confidence region by a factor of $s(1)$. The restriction $\zeta < 1$ is made to rule out sub-Gaussian and sub-exponential random variables.

2.12 Perspectives & Future Work

In this chapter, we proposed a class of practical algorithms for the distributionally robust optimization (DRO) and proved their linear convergence guarantee for smoothed DRO problems. We paid particular attention to spectral risk measure objectives by first formulating them as DRO objectives and using their quantile-based representation to bound their bias under general conditions (e.g. the loss may be unbounded). Unlike previous DRO formulations, which were typically based on closed balls in KL-divergence [Kumar et al., 2024] or χ^2 -divergence [Duchi and Namkoong, 2019], the dual problem for spectral risk measures is easy to solve by subroutines such as sorting or isotonic regression. Furthermore, the hyperparameter $\sigma = (\sigma_1, \dots, \sigma_n)$ gives the user both visibility and control concerning the optimal dual solution (i.e. the solution will resemble a permutation of σ).

Several directions for future work were outlined in Section 2.11. One direct path would be to generalize the results on squared error loss and changes in the data support from Proposition 2.11.1 to generalized linear model (GLM) loss functions, i.e., those that depend on the parameter θ through convex functions of the quantities $\langle x_i, \theta \rangle$ for $i = 1, \dots, n$. Importantly, this addresses a major limitation of the likelihood ratio-based approaches in this chapter, that is, the absolute continuity constraint of the shifted distributions. The main competing DRO approach is based on Wasserstein metric-based uncertainty sets, which allow for changes in the support but are much less tractable computationally and may only be formulated into finite-dimensional programs under strict conditions [Kuhn et al., 2019]. Incorporating input noise into the likelihood ratio-based approaches of this chapter might yield a “best of both worlds” method in this regard.

Secondly, while the perturbations discussed in this chapter are based on reweighting the observed training data, another interpretation of “changing data” is direct transformation/cor-

ruption of data instances. While this idea is explored using Gaussian noise for vector-valued inputs in Section 2.11.1, formulating the problem using a finite set of naturally-occurring distortions (such as image blur) would be an interesting line of work for applications in vision and language.

In the next chapter, we expand on the ideas mentioned in Section 2.6.4. We recognize that the structure of the DRO objective applies to multiple problems in the data sciences and aim to derive theoretically optimal algorithms using a primal-dual approach. Note that generalizing in this manner does not supplant the work of this chapter; instead, questions such as the large-sample properties of the objective (Section 2.3), uncertainty set selection and computing the maximization map (Section 2.7), and viability on realistic distribution shift scenarios (Section 2.9 and Section 2.11), remain exclusive subjects of Chapter 2.

Chapter 3

ALGORITHMIC EXTENSIONS OF DISTRIBUTIONALLY ROBUST OPTIMIZATION

3.1 Introduction

In Chapter 2, we introduced the distributionally robust optimization problem (2.18) and provided an algorithm and convergence analysis from a “primal-only” viewpoint. Concretely, this refers to the insistence on having a single hyperparameter $\eta > 0$ as a learning rate and using a proof technique called the Lyapunov stability argument (see Section 2.6), which is largely inspired by variance-reduced stochastic algorithms for finite-sum minimization [Johnson and Zhang, 2013, Defazio et al., 2014]. Moreover, the class of algorithms studied in Chapter 2 was motivated primarily by their excellent performance in numerical benchmarks, as observed in Section 2.9. Nevertheless, we would be remiss not to recognize the primal-dual structure of this saddle-point problem, leading to a possibly sharper theoretical convergence guarantee and applicability to more general problems.

Indeed, the DRO problem is an instance of a more general one, which has been a subject of fundamental research in optimization for decades. Consider the saddle point (or min-max) problem

$$\min_{\mathbf{x} \in \mathcal{X}} \max_{\mathbf{y} \in \mathcal{Y}} c(\mathbf{x}, \mathbf{y}) - \psi(\mathbf{y}) + \phi(\mathbf{x}), \quad (3.1)$$

with respect to primal variable $\mathbf{x} \in \mathcal{X} \subseteq \mathbb{R}^d$ and dual variable $\mathbf{y} \in \mathcal{Y} \subseteq \mathbb{R}^n$. We refer to $c(\mathbf{x}, \mathbf{y})$ as the *coupled* component, whereas $\phi(\mathbf{x})$ and $\psi(\mathbf{y})$ are the *individual* components of the objective. Such problems are further categorized as *bilinearly coupled* if $c(\mathbf{x}, \mathbf{y}) = \mathbf{y}^\top \mathbf{A} \mathbf{x}$ for a matrix $\mathbf{A} \in \mathbb{R}^{n \times d}$, and are called *nonbilinearly coupled* otherwise. While primal-dual algorithms for saddle point problems are still an active area of modern machine learning and

optimization research, nonbilinearly coupled objectives have received much less attention than their bilinearly coupled counterparts. This is due in part to bilinearly coupled objectives being amenable to various innovations such as Chambolle-Pock-style [Chambolle and Pock, 2011] and/or stochastic coordinate-wise updates [Song et al., 2021] that yield an improved runtime. Interestingly, the objective (2.18) has more structure than in a general nonbilinearly coupled problem.

The explorations in this chapter are motivated by the observation that numerous non-bilinearly coupled objectives, including DRO, are in fact linear with respect to one of the two decision variables. Accordingly, insights from the bilinear setting can be used to design efficient algorithms for these “dual-linear” min-max problems¹. Formally, we introduce the optimization problem

$$\min_{\mathbf{x} \in \mathcal{X}} \max_{\mathbf{y} \in \mathcal{Y}} [\mathcal{L}(\mathbf{x}, \mathbf{y}) := \langle \mathbf{y}, f(\mathbf{x}) \rangle - \psi(\mathbf{y}) + \phi(\mathbf{x})], \quad (3.2)$$

in which $f = (f_1, \dots, f_n)$ has convex components, ϕ is μ -strongly convex ($\mu \geq 0$), and ψ is ν -strongly convex ($\nu \geq 0$). We further require that $\mathcal{Y} \subseteq \{\mathbf{y} \in \mathbb{R}^n : y_j \geq 0 \text{ if } f_j \text{ is non-linear}\}$, so that (3.2) constitutes a legitimate convex-concave saddle point problem. Beyond DRO, we consider the following examples, which appear frequently in statistical learning applications.

Example 3.1.1 (Generalized Linear Models (GLMs)). Consider a design matrix $\mathbf{A} \in \mathbb{R}^{n \times d}$, where each row \mathbf{A}_i contains a d -dimensional feature vector. Fitting a generalized linear model (GLM) via the maximum likelihood principle results in a problem of the form

$$\min_{\mathbf{x} \in \mathcal{X}} \sum_{i=1}^n \psi_i^*(\langle \mathbf{A}_i, \mathbf{x} \rangle) + \phi(\mathbf{x}),$$

where \mathcal{X} is a parameter space, $\psi_i^* : \mathbb{R} \rightarrow \mathbb{R} \cup \{+\infty\}$ is a convex loss function (e.g., the negative log-likelihood of the i -th output conditioned on \mathbf{A}_i), and ϕ is a regularizer such as the ℓ_2 -norm squared or elastic net penalty. When taking the Fenchel conjugate of ψ_i^*

¹While we fix the convention that the coupled term is linear in the dual variables, our methods extend by analogy to primal-linear objectives.

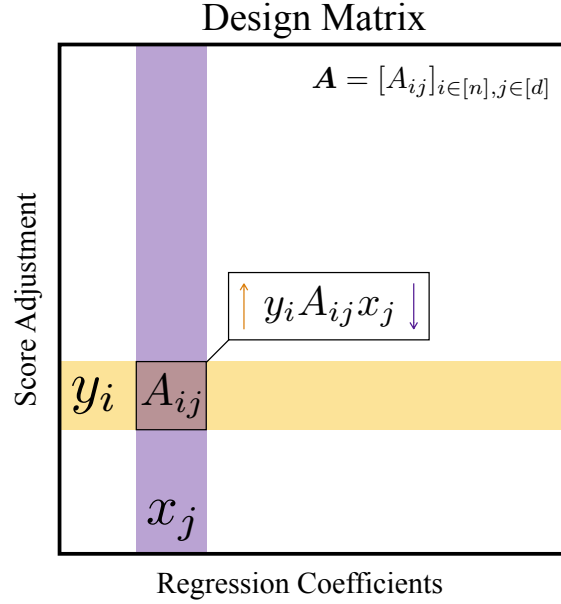


Figure 3.1: **Generalized Linear Models as Matrix Games.** Visualization to accompany Example 3.1.1, where \mathbf{A} denotes the design matrix, \mathbf{x} is parameter vector, and \mathbf{y} multiplier to adjust the predicted scores \mathbf{Ax} .

(denoted ψ_i), we uncover the min-max problem

$$\min_{\mathbf{x} \in \mathcal{X}} \max_{\mathbf{y} \in \mathbb{R}^n} \langle \mathbf{y}, \mathbf{Ax} \rangle - \underbrace{\sum_{i=1}^n \psi_i(y_i)}_{\psi(\mathbf{y})} + \phi(\mathbf{x}), \quad (3.3)$$

which is an example of (3.2) with $f_i(\mathbf{x}) = \langle \mathbf{A}_i, \mathbf{x} \rangle$. This bilinearly coupled min-max problem, depicted in Figure 3.1, has interesting features beyond the classical setup. Notice that in terms of the dual variables, the objective of (3.3) is the sum of functions that depend only on each $y_i \in \mathbb{R}$ individually. This “separable” structure inspired recent methods for bilinearly coupled min-max problems with applications to statistical learning [Song et al., 2021], and is one that we pay specific attention to in Section 3.5.

Example 3.1.2 (Fully Composite Optimization). The fully composite optimization problem (see Cui et al. [2018], Drusvyatskiy and Paquette [2019], Doikov and Nesterov [2022],

Vladarean et al. [2023] and references therein) writes as

$$\min_{\mathbf{x} \in \mathcal{X}} F(h(\mathbf{x}), \mathbf{x}),$$

where $h : \mathcal{X} \rightarrow \mathbb{R}^{n-d}$ (for $n > d$) is component-wise convex and $F : \mathbb{R}^n \times \mathcal{X} \rightarrow \mathbb{R}$ is closed and convex. It is assumed that h is smooth and “hard” to compute, whereas F may be non-differentiable but “easy” to compute. This formulation can be viewed as a generalization of the typical additive composite problem in which $n - d = 1$ and $F(h(\mathbf{x}), \mathbf{x}) = h(\mathbf{x}) + g(\mathbf{x})$ for a smooth component h and non-smooth component g . By taking the Fenchel conjugate $F^*(\mathbf{y}) := \sup_{\mathbf{z} \in \mathbb{R}^n \times \mathcal{X}} \langle \mathbf{z}, \mathbf{y} \rangle - F(\mathbf{z})$ of F , we achieve the formulation (3.2) with $f(\mathbf{x}) := (h(\mathbf{x}), \mathbf{x}) \in \mathbb{R}^n$, $\psi(\mathbf{y}) = F^*(\mathbf{y})$, and $\phi(\mathbf{x}) = 0$. To ensure that the overall problem is convex-concave we also assume that $F(\cdot, \mathbf{x})$ is monotone in that $\mathbf{u} \leq \mathbf{v}$ element-wise $\implies F(\mathbf{u}, \mathbf{x}) \leq F(\mathbf{v}, \mathbf{x})$.

Example 3.1.3 (Problems with Functional Constraints). Consider the classical convex minimization problem with constraints defined by sublevel sets of convex functions

$$\min_{\mathbf{x} \in \mathcal{X}} \phi(\mathbf{x}) \text{ s.t. } f_j(\mathbf{x}) \leq 0 \text{ for all } j = 1, \dots, n.$$

Then, the Lagrangian formulation yields the expression (3.2) by letting $\mathbf{y} \in \mathcal{Y} = \mathbb{R}_+^n$ denote the Lagrange multipliers and setting $\psi \equiv 0$. The objective (3.2) also encompasses the related setting of “soft” functional constraints, where we set ψ as any ν -strongly convex function ψ with $\nu > 0$ to produce a faster convergence rate at an approximation cost governed by the parameter ν . In this case, the primal solution resulting from this smoothed problem may only approximately satisfy the functional constraints. A classical example in statistical learning is the support vector machine problem [Cortes and Vapnik, 1995].

Example 3.1.4 (Maximal eigenvalue minimization). Given a collection of d symmetric matrices $\mathbf{A}_1, \dots, \mathbf{A}_d \in \mathbb{R}^{m \times m}$, the classical problem of minimizing the maximal eigenvalue $\lambda_{\max}(\sum_{i=1}^d x_i \mathbf{A}_i)$ in $\mathbf{x} = (x_1, \dots, x_d) \in \mathbb{R}^d$ can be formulated [Nesterov, 2007b, Baes et al.,

2013] as the saddle-point problem

$$\min_{\mathbf{x} \in \mathcal{X}} \max_{\mathbf{y} \in \mathcal{Y}} \sum_{i=1}^d \text{tr}((\mathbf{A}_i x_i) \mathbf{y}),$$

where \mathcal{Y} is the set of positive semi-definite matrices satisfying $\text{tr}(\mathbf{y}) = 1$ and \mathcal{X} is any convex, compact subset of \mathbb{R}^d . Here, $n = m^2$ depends quadratically on the height/width of the matrices $(\mathbf{A}_i)_{i=1}^d$, and blocks may naturally correspond to matrix structure, such as rows or columns. This constitutes a non-separable matrix game, a specific problem class discussed in Section 3.6.

Even though some of the examples above move from statistical learning problems to more general optimization problems in applied mathematics, the field of statistics is not simply a source of example problems that fit our setting; instead, ideas such as using adaptive sampling and carefully handling complex dependence structures that arise will be essential for our improved complexity guarantees. In order to make statements of complexity, we also note that we will use a different convergence criterion in this chapter than in the previous one (another concrete facet of the “primal-dual” mindset). This *primal-dual gap* criterion will be amenable to an analysis that does not require a Lyapunov function and will be able to handle cases of non-strong convexity in a unified manner. Furthermore, the analysis of the primal-dual gap done in this chapter is of mathematical value; our proof techniques (outlined Section 3.3) are in line with prior work providing *constructive* arguments for the analysis of optimization methods [Diakonikolas and Orecchia, 2019, Mehta et al., 2024a, Li et al., 2024b, Diakonikolas, 2025]. In other words, the theoretical analysis provides guiding principles for deriving optimization algorithms that may otherwise be unsuspected. Applying this method, we present an algorithm that employs a unique combination of randomized updates and a “historical regularization” technique that is conceptually novel and interesting in its own right. This contrasts with the retrospective, empirically motivated approach of Chapter 2. By furnishing both perspectives, we argue that a deeper understanding of the problem of interest is developed. Let us outline the rest of the chapter.

In Section 3.2, we state the assumptions that fully specify the problem class and review the complexities of classical methods for nonbilinearly coupled saddle point problems. We also define *block separable* problems, whose structure we will exploit in later sections. Section 3.3 we introduce the high-level steps of the analysis and technical lemmas used throughout the proofs in various settings. In Section 3.4 and Section 3.5, we realize the template from the previous section for general and block separable objectives, respectively. Extensive comparisons are made to recent literature in Section 3.6. Extensions such as online convergence certificates and lower bounds are discussed in Section 3.7, and future work is commented on in Section 3.8.

3.2 Preliminaries

To initiate the discussion of computational complexity, we first make precise our (strong) convexity assumption, where we interpret values of the strong convexity modulus being zero ($\mu = 0$ or $\nu = 0$) as the corresponding function being simply convex. Let $\text{ri}(\cdot)$ denote the relative interior of a set, and let $\|\cdot\|_{\mathcal{X}}$ and $\|\cdot\|_{\mathcal{Y}}$ be norms on \mathbb{R}^d and \mathbb{R}^n , respectively.

Assumption 3.2.1. Assume that ϕ is proper with $\mathcal{X} \subseteq \text{dom}(\phi) := \{\mathbf{x} \in \mathbb{R}^d : \phi(\mathbf{x}) < +\infty\}$, closed (i.e., has a closed epigraph in \mathbb{R}^{d+1}), and μ -strongly convex ($\mu \geq 0$) with respect to $\|\cdot\|_{\mathcal{X}}$, that is, for any $\mathbf{s} \in \partial\phi(\mathbf{u})$, we have that $\phi(\mathbf{z}) \geq \phi(\mathbf{u}) + \langle \mathbf{s}, \mathbf{z} - \mathbf{u} \rangle + \frac{\mu}{2} \|\mathbf{z} - \mathbf{u}\|_{\mathcal{X}}^2$. Assume in addition that $\mathcal{X} \cap \text{ri}(\text{dom}(\phi))$ is non-empty. Similarly, assume that ψ is proper (with $\mathcal{Y} \subseteq \text{dom}(\psi)$ and $\mathcal{Y} \cap \text{ri}(\text{dom}(\psi))$ non-empty), closed, and ν -strongly convex ($\nu \geq 0$) with respect to $\|\cdot\|_{\mathcal{Y}}$.

Recall that when $\mu > 0$, we aimed to control $\|\boldsymbol{\theta} - \boldsymbol{\theta}^*\|_2^2$ in the results of Section 2.5.3, where the $\boldsymbol{\theta}^* \in \mathbb{R}^d$ was the unique optimum of the objective. That analog of this criterion for (3.2) would clearly be $\|\mathbf{x} - \mathbf{x}^*\|_{\mathcal{X}}^2$ with $\mathbf{x}^* \in \mathcal{X}$ defined analogously. Other than the possibly non-Euclidean geometry induced by $\|\cdot\|_{\mathcal{X}}$, note that \mathbf{x}^* may not exist or be unique for $\mu = 0$, which is also true of the similarly defined $\mathbf{y}^* \in \mathcal{Y}$. Instead, as is standard for primal-dual algorithm of the same type [Chambolle and Pock, 2011, Alacaoglu et al., 2020,

[Song et al., 2021], the optimality criterion will be the *primal-dual gap* of $(\mathbf{x}, \mathbf{y}) \in \mathcal{X} \times \mathcal{Y}$ at $\mathbf{u} \in \mathcal{X}$ and $\mathbf{v} \in \mathcal{Y}$:

$$\text{Gap}^{\mathbf{u}, \mathbf{v}}(\mathbf{x}, \mathbf{y}) := \mathcal{L}(\mathbf{x}, \mathbf{v}) - \mathcal{L}(\mathbf{u}, \mathbf{y}). \quad (3.4)$$

To relate this criterion to our previous one in the strongly convex case (when $\mu > 0$ and $\nu > 0$), we set $\mathbf{u} = \mathbf{x}_\star$ and $\mathbf{v} = \mathbf{y}_\star$, where $(\mathbf{x}_\star, \mathbf{y}_\star)$ also forms the unique saddle point of the objective (3.2), satisfying $\mathcal{L}(\mathbf{x}_\star, \mathbf{y}) \leq \mathcal{L}(\mathbf{x}_\star, \mathbf{y}_\star) \leq \mathcal{L}(\mathbf{x}, \mathbf{y}_\star)$. In fact, by strong convexity, we have that

$$\text{Gap}^{\mathbf{x}^\star, \mathbf{y}^\star}(\mathbf{x}, \mathbf{y}) - \frac{\mu}{2} \|\mathbf{x} - \mathbf{x}^\star\|_{\mathcal{X}}^2 - \frac{\nu}{2} \|\mathbf{y} - \mathbf{y}^\star\|_{\mathcal{Y}}^2 \geq 0,$$

which can be seen by adding the two non-negative terms $\mathcal{L}(\mathbf{x}, \mathbf{y}^\star) - \mathcal{L}(\mathbf{x}^\star, \mathbf{y}^\star) - \frac{\mu}{2} \|\mathbf{x} - \mathbf{x}^\star\|_{\mathcal{X}}^2 \geq 0$ and $\mathcal{L}(\mathbf{x}^\star, \mathbf{y}^\star) - \mathcal{L}(\mathbf{x}^\star, \mathbf{y}) - \frac{\nu}{2} \|\mathbf{y} - \mathbf{y}^\star\|_{\mathcal{Y}}^2 \geq 0$. Thus, a quantitative convergence guarantee on (3.4) immediately yields such a guarantee on the distance-to-optimum terms $\|\mathbf{x} - \mathbf{x}^\star\|_{\mathcal{X}}^2$ and $\|\mathbf{y} - \mathbf{y}^\star\|_{\mathcal{Y}}^2$. Importantly, the primal-dual gap can be used when $\mu = 0$ or $\nu = 0$ by taking a supremum of (3.4) over $\mathbf{u} \in \mathcal{U}$ or $\mathbf{v} \in \mathcal{V}$, where $\mathcal{U} \subseteq \mathcal{X}$ and $\mathcal{V} \subseteq \mathcal{Y}$ are compact sets in which the iterates are contained.

Having described the convergence criterion formally, we may now motivate the usefulness of specialized algorithms for (3.2). We will compare various approaches in terms of the total arithmetic complexity or runtime needed to guarantee that (3.4) is smaller than a suboptimality parameter $\varepsilon > 0$. Consider treating (3.2) generically as a nonbilinearly coupled saddle-point problem (or more generally still as a variational inequality (VI) problem). Let λ be the Lipschitz parameter of $(\mathbf{x}, \mathbf{y}) \mapsto (\nabla_{\mathbf{x}} \mathcal{L}, -\nabla_{\mathbf{y}} \mathcal{L})$, when it exists. Classical approaches such as Korpelevich's extragradient [Korpelevich, 1976], Popov's method [Popov., 1980], Nemirovski's mirror-prox [Nemirovski, 2004], Nesterov's dual extrapolation [Nesterov, 2007a] will achieve the complexity $O(nd\lambda/\varepsilon)$ when $\mu = \nu = 0$, whereas their linearly convergent variants [Nesterov and Scramali, 2006, Marcotte, 1991, Tseng, 1995, Mokhtari et al., 2020, Le Thi Thanh Hai and Vuong, 2025] achieve $O(nd(\lambda/\min\{\mu, \nu\}) \ln(1/\varepsilon))$ in the strongly

convex-strongly concave setting.

Furthermore, using generic acceleration Catalyst schemes [Lin et al., 2018] adapted for min-max problems, improved complexity bounds can be obtained. Yang et al. [2020] achieve an $\tilde{O}(nd\lambda/\sqrt{\mu\varepsilon})$ runtime in the strongly convex-concave regime. Recently Lan and Li [2023] obtained a complexity of $\tilde{O}(nd\lambda/\sqrt{\mu\nu} \ln(\frac{1}{\varepsilon}))$ for strongly convex-strongly concave problems. A disadvantage of classical approaches that is overcome by the Catalyst approach is the dependence on the minimum of the two strong convexity constants. In the examples above, ψ often plays the role of a strongly convex smoothing penalty whose strong convexity constant may be near-zero (e.g., $O(\varepsilon)$ for some applications). Thus, guarantees depending on $(\mu \wedge \nu)$ are undesirable. A second disadvantage incurred by all of the approaches above is the dependence on a coarse Lipschitz constant λ for the entire vector field. When it is only known that each component function f_j is G_{\max} -Lipschitz continuous and M_{\max} -smooth, and that $|y_j| \leq D_{\max}$ for any $\mathbf{y} = (y_1, \dots, y_n) \in \mathcal{Y}$, then λ can be estimated as $n(G_{\max} + D_{\max} M_{\max})$. Combined with the order- nd per-iteration cost, the total arithmetic complexity can have an $O(n^2d)$ dependence on the dimensions (n, d) . These two issues highlight the price to pay when using generic acceleration schemes: the loss of adaptation to the structure of the problem and the related constants. In contrast, by leveraging non-uniformity in the various Lipschitz/smoothness constants for each f_j and randomized updates, we may achieve complexities that are linear in $(n + d)$.

In contrast, the (possibly non-uniform) Lipschitz continuity and smoothness parameters of the component functions represent quantities of practical interest in applications. In Example 3.1.1, the Lipschitz continuity properties of each f_j will depend on the norm of the corresponding row $\mathbf{A}_{j\cdot}$. In Example 3.1.2, we have that $f(\mathbf{x}) = (h(\mathbf{x}), \mathbf{x})$ so that f_{n-d+1}, \dots, f_n are 1-Lipschitz continuous and 0-smooth; a complexity result that only depends on the n times the maxima of these constants over the components of f could be overly pessimistic.

In our analysis, we handle every combination of $\mu = 0$ versus $\mu > 0$ and $\nu = 0$ versus $\nu > 0$ to achieve the same dependence on ε in a unified way with a constant that enjoys a

transparent dependence on these component-wise Lipschitz and smoothness constants. The complexity benefits of using component-wise problem constants have been shown in recent works on stochastic variance-reduced methods for variational inequality problems [Alacaoglu and Malitsky, 2022, Cai et al., 2024, Pichugin et al., 2024, Alizadeh et al., 2024, Diakonikolas, 2025].

We now introduce the mathematical objects and assumptions used in this chapter. As before, we let $\|\cdot\|_{\mathcal{X}}$ denote a norm on \mathbb{R}^d . Let $\|\cdot\|_{\mathcal{Y}}$ denote an ℓ_p -norm on \mathbb{R}^n for $p \in [1, 2]$. The associated dual norms are denoted by $\|\cdot\|_{\mathcal{X}^*}$ and $\|\cdot\|_{\mathcal{Y}^*}$ and defined in the usual way as $\|\mathbf{w}\|_{\mathcal{X}^*} = \sup_{\mathbf{x}: \|\mathbf{x}\|_{\mathcal{X}} \leq 1} \langle \mathbf{w}, \mathbf{x} \rangle$, $\|\mathbf{z}\|_{\mathcal{Y}^*} = \sup_{\mathbf{y}: \|\mathbf{y}\|_{\mathcal{Y}} \leq 1} \langle \mathbf{z}, \mathbf{y} \rangle$. The following assumptions about the objective in (3.2) are made throughout the chapter. We employ block coordinate-wise updates in the upcoming stochastic algorithms, of which coordinate-wise updates are a special case with block size one. To do so, we introduce a partitioning of the n components of f into N blocks and define the relevant Lipschitz constants for each one. In discussions of arithmetic complexity, we may assume a uniform block size n/N , but the analysis natively handles blocks of possibly non-uniform size.

Assumption 3.2.2. Assume that each $f_j : \mathbb{R}^d \rightarrow \mathbb{R}$ is convex and the indices $[n] = \{1, \dots, n\}$ are partitioned into *blocks* (B_1, \dots, B_N) . There exist constants $\mathbf{G}_1, \dots, \mathbf{G}_N \geq 0$ and $\mathbf{L}_1, \dots, \mathbf{L}_N \geq 0$ such that for each $J = 1, \dots, N$:

$$\sup_{\mathbf{x} \in \mathcal{X}} \left\| \sum_{j \in B_J} z_j \nabla f_j(\mathbf{x}) \right\|_{\mathcal{X}^*} \leq \mathbf{G}_J \|\mathbf{z}\|_2, \quad \forall \mathbf{z} \in \mathbb{R}^{|B_J|}, \quad (3.5)$$

$$\sup_{\mathbf{y} \in \mathcal{Y}} \left\| \sum_{j \in B_J} y_j (\nabla f_j(\mathbf{x}) - \nabla f_j(\mathbf{x}')) \right\|_{\mathcal{X}^*} \leq \mathbf{L}_J \|\mathbf{x} - \mathbf{x}'\|_{\mathcal{X}}, \quad \forall \mathbf{x}, \mathbf{x}' \in \mathcal{X}. \quad (3.6)$$

In addition, for $\nabla f(\mathbf{x}) := (\nabla f_1(\mathbf{x}), \dots, \nabla f_n(\mathbf{x})) \in \mathbb{R}^{n \times d}$, there exist $G \geq 0$ and $L \geq 0$ such that

$$\sup_{\mathbf{x} \in \mathcal{X}} \left\| \nabla f(\mathbf{x})^\top \mathbf{z} \right\|_{\mathcal{X}^*} \leq G \|\mathbf{z}\|_{\mathcal{Y}}, \quad \forall \mathbf{z} \in \mathbb{R}^n, \quad (3.7)$$

$$\sup_{\mathbf{y} \in \mathcal{Y}} \left\| (\nabla f(\mathbf{x}) - \nabla f(\mathbf{x}'))^\top \mathbf{y} \right\|_{\mathcal{X}^*} \leq L \|\mathbf{x} - \mathbf{x}'\|_{\mathcal{X}}, \quad \forall \mathbf{x}, \mathbf{x}' \in \mathcal{X}. \quad (3.8)$$

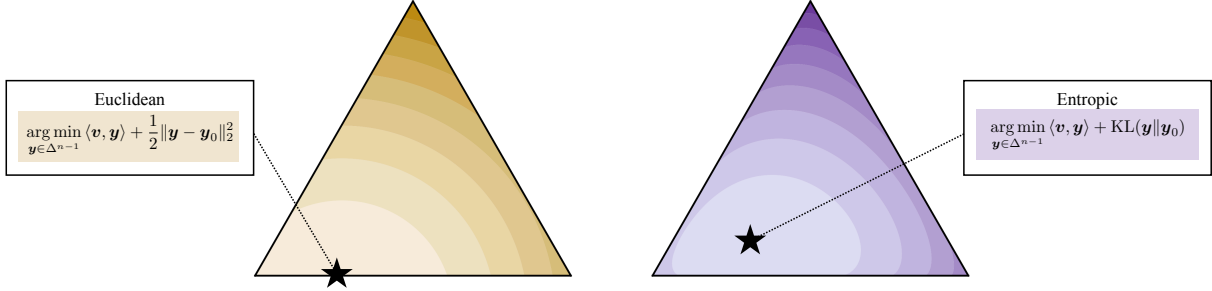


Figure 3.2: **Moreau Envelope of Bregman Divergences on the Unit Simplex.** Visualization of the objective of the proximal operator in the case of the standard ℓ_2 -norm-squared Bregman divergence (left) and the negative entropy-base Bregman divergence (right).

Next, we address the possible non-Euclidean structure of \mathcal{X} and \mathcal{Y} . A classic example of such a feasible set is the probability simplex $\Delta^{n-1} = \{\mathbf{y} \in \mathbb{R}_{\geq 0}^n : \langle \mathbf{y}, \mathbf{1} \rangle = 1\}$, visualized in Figure 3.2. Notice that using the ℓ_2 -geometry results in a constrained optimization problem (Euclidean projection on the simplex), whereas the entropic proximal step objective has a unique minimizer contained within the set. Arguing analytically, consider the case in which $\|\cdot\|_{\mathcal{X}} = \|\cdot\|_2$ and we wish to upper bound G from (3.7) in terms of the individual Lipschitz constants of the component functions f_1, \dots, f_n . In the case of $\|\cdot\|_{\mathcal{Y}} = \|\cdot\|_2$ we have that

$$\sup_{\mathbf{x} \in \mathcal{X}} \|\nabla f(\mathbf{x})^\top \mathbf{y}\|_{\mathcal{X}^*} \leq \sup_{\mathbf{x} \in \mathcal{X}} \underbrace{\sqrt{\sum_{j=1}^n \|\nabla f_j(\mathbf{x})\|_2^2}}_{\|\nabla f_j(\mathbf{x})\|_{\text{Fro}}} \|\mathbf{y}\|_2$$

whereas for $\|\cdot\|_{\mathcal{Y}} = \|\cdot\|_1$, we have

$$\sup_{\mathbf{x} \in \mathcal{X}} \|\nabla f(\mathbf{x})^\top \mathbf{y}\|_{\mathcal{X}^*} \leq \sup_{\mathbf{x} \in \mathcal{X}} \underbrace{\max_{j \in [n]} \|\nabla f_j(\mathbf{x})\|_2}_{\|\nabla f(\mathbf{x})^\top\|_{1,2}} \|\mathbf{y}\|_1.$$

For the second inequality, the resulting bound is up to \sqrt{n} smaller than for the first. To take advantage of non-Euclidean norms, we rely upon the following definition of Bregman divergences for possibly non-differentiable distance-generating functions.

Consider a convex subset $\mathcal{Z} \subseteq \mathbb{R}^m$, and let φ be a proper, closed, and 1-strongly convex

function satisfying $\mathcal{Z} \subseteq \text{dom}(\varphi)$. Define the *Bregman divergence* $\Delta_\varphi : \text{dom}(\varphi) \times \text{ri}(\text{dom}(\varphi)) \rightarrow \mathbb{R}$ as

$$\Delta_\varphi(\mathbf{z}, \mathbf{z}') := \varphi(\mathbf{z}) - \varphi(\mathbf{z}') - \langle \nabla \varphi(\mathbf{z}'), \mathbf{z} - \mathbf{z}' \rangle.$$

The notation $\nabla \varphi(\mathbf{z}') \in \partial \varphi(\mathbf{z}')$ denotes an arbitrary, but consistently chosen subgradient at \mathbf{z}' when applied to a convex but possibly non-differentiable function. This slight modification is made for purely technical reasons, as we perform a mirror descent-style analysis with Bregman divergences generated by the (possibly non-smooth) component functions ϕ and ψ . Lemma 3.2.1 provides a modified three-point inequality used in both the upper and lower bounds on the objective used in bounding the initial gap estimates.

Lemma 3.2.1. *Let h , g , and φ be proper, closed, and convex functions whose domains contain \mathcal{Z} and that map to $\mathbb{R} \cup \{+\infty\}$. Assume that g is relatively γ -strongly convex with respect to φ on \mathcal{Z} , i.e., $\Delta_g(\mathbf{u}, \mathbf{z}) \geq \gamma \Delta_\varphi(\mathbf{u}, \mathbf{z})$ for $\gamma \geq 0$ and $\mathbf{u}, \mathbf{z} \in \mathcal{Z}$. Let $A \geq 0$, $a > 0$, $\gamma_0 > 0$ be constants, and let $\mathbf{z}_1, \dots, \mathbf{z}_r \in \text{ri}(\text{dom}(\varphi))$. Let*

$$\mathbf{z}^+ = \arg \min_{\mathbf{u} \in \mathcal{Z}} \left\{ m(\mathbf{u}) := h(\mathbf{u}) + ag(\mathbf{u}) + \frac{A\gamma + \gamma_0}{2} \sum_{i=1}^r w_i \Delta_\varphi(\mathbf{u}, \mathbf{z}_i) \right\}, \quad (3.9)$$

where each $w_i \geq 0$ and $\sum_{i=1}^r w_i = 1$. Then, for any $\mathbf{u} \in \mathcal{Z}$,

$$m(\mathbf{u}) \geq m(\mathbf{z}^+) + \left(\frac{(A+a)\gamma + \gamma_0}{2} \right) \Delta_\varphi(\mathbf{u}, \mathbf{z}^+) + \frac{a\gamma}{2} \Delta_\varphi(\mathbf{u}, \mathbf{z}^+).$$

Proof. By the definition of the Bregman divergence generated by m , we have that

$$\begin{aligned} m(\mathbf{u}) &= m(\mathbf{z}^+) + \langle \nabla m(\mathbf{z}^+), \mathbf{u} - \mathbf{z}^+ \rangle + \Delta_m(\mathbf{u}, \mathbf{z}^+), \\ &\geq m(\mathbf{z}^+) + \Delta_m(\mathbf{u}, \mathbf{z}^+), \end{aligned}$$

where we use that $\langle \nabla m(\mathbf{z}^+), \mathbf{u} - \mathbf{z}^+ \rangle \geq 0$ for any subgradient $\nabla m(\mathbf{z}^+)$ as $\mathbf{z}^+ \in \arg \min_{\mathbf{u} \in \mathcal{Z}} m(\mathbf{u})$.

Then, by using the definition of m , and that $\Delta_{\Delta_\varphi(\cdot, \mathbf{z})} = \Delta_\varphi$ for any fixed $\mathbf{z} \in \text{ri}(\text{dom}(\varphi))$, we

have that

$$\begin{aligned}
m(\mathbf{u}) &\geq m(\mathbf{z}^+) + \Delta_m(\mathbf{u}, \mathbf{z}^+) \\
&= m(\mathbf{z}^+) + \Delta_h(\mathbf{u}, \mathbf{z}^+) + a\Delta_g(\mathbf{u}, \mathbf{z}^+) + \frac{A\gamma + \gamma_0}{2} \sum_{i=1}^r w_i \Delta_{\Delta_\varphi(\cdot, \mathbf{z}_i)}(\mathbf{u}, \mathbf{z}^+) \\
&= m(\mathbf{z}^+) + \Delta_h(\mathbf{u}, \mathbf{z}^+) + a\Delta_g(\mathbf{u}, \mathbf{z}^+) + \frac{A\gamma + \gamma_0}{2} \Delta_\varphi(\mathbf{u}, \mathbf{z}^+).
\end{aligned}$$

Use $\Delta_h(\mathbf{u}, \mathbf{z}^+) \geq 0$ and then relative strong convexity $a\Delta_g(\mathbf{u}, \mathbf{z}^+) \geq a\gamma\Delta_\varphi(\mathbf{u}, \mathbf{z}^+)$ to prove the desired result. \square

Henceforth, we use the notation $\Delta_{\mathcal{X}}(\cdot, \cdot)$ to denote the Bregman divergence on \mathcal{X} that is both 1-strongly convex with respect to $\|\cdot\|_{\mathcal{X}}$ and that satisfies $\Delta_\phi \geq \mu\Delta_{\mathcal{X}}$ (i.e., ϕ is relatively μ -strongly convex with respect to $\Delta_{\mathcal{X}}$). As an example, this is satisfied by $\Delta_{\mathcal{X}} := \Delta_{\phi/\mu}$ for $\mu > 0$. We define $\Delta_{\mathcal{Y}}(\cdot, \cdot)$ analogously. As a technical consideration, we also assume that the unique solution of (3.9) lies in $\text{ri}(\text{dom}(\varphi))$, which is satisfied for common choices of the generator φ . Finally, we introduce additional structure on the problem (3.2) which can be exploited when available.

Definition 3.2.1. We call \mathcal{L} a *dual-separable objective* if the dual component decomposes as $\psi(\mathbf{y}) = \sum_{J=1}^N \psi_J(\mathbf{y}_J)$ where \mathbf{y}_J denotes the components of $\mathbf{y} \in \mathcal{Y}$ corresponding to the indices in block B_J . We call $\mathcal{X} \times \mathcal{Y}$ a *dual-separable feasible set* if $\mathcal{Y} = \mathcal{Y}_1 \times \dots \times \mathcal{Y}_N$ and $\mathbf{y}_J \in \mathcal{Y}_J$ for $J = 1, \dots, N$. We call the problem (3.2) a *dual-separable problem* if its objective and feasible set are both dual-separable.

Dual-separability of the objective is commonly satisfied, such as when ψ represents an ℓ_2 or negative entropy penalty. In this case, we have that $\Delta_{\mathcal{Y}}(\mathbf{y}, \mathbf{y}') := \sum_{J=1}^N \Delta_J(\mathbf{y}_J, \mathbf{y}'_J)$, where each $\Delta_J(\cdot, \cdot)$ is a Bregman divergence on $\text{dom}(\psi_J) \times \text{ri}(\text{dom}(\psi_J))$. As before, ψ is relatively ν -strongly convex with respect to $\Delta_{\mathcal{Y}}$. In later sections, we may also use the subscript on $\mathbf{y}_k \in \mathcal{Y}$ to denote a particular time index of an algorithm, as opposed to the block index J on $\mathbf{y}_J \in \mathcal{Y}_J$; the difference will be clear from context. Dual-separability of the feasible set is a less common assumption. It is not satisfied, for instance, on simplicial domains such as

the one in Example 1. Based on this observation, Section 3.4 and Section 3.5 are dedicated to the proposed algorithms for non-separable and separable problems, respectively.

3.3 Method and General Analysis Template

We consider an algorithm to be a sequence of primal-dual iterates $(\mathbf{x}_k, \mathbf{y}_k)_{k \geq 0}$, with fixed initial point $(\mathbf{x}_0, \mathbf{y}_0)$. This sequence may be random, in which case the relevant probabilistic information is introduced when we analyze stochastic algorithms.

Recalling the gap function (3.4), we first fix \mathbf{u}, \mathbf{v} and aim to show that $\limsup_{t \rightarrow \infty} \text{Gap}^{\mathbf{u}, \mathbf{v}}(\mathbf{x}_t, \mathbf{y}_t) \leq 0$ (possibly in expectation), with a convergence rate in terms of the problem constants from Section 3.2 and iteration count $t \geq 0$. To this end, we introduce an averaging sequence $(a_k)_{k \geq 1}$ of positive constants with $a_0 = 0$, and their aggregation $A_t := \sum_{k=0}^t a_k$, and then pursue an upper bound of the form

$$\sum_{k=1}^t a_k \text{Gap}^{\mathbf{u}, \mathbf{v}}(\mathbf{x}_k, \mathbf{y}_k) \leq G_0(\mathbf{u}, \mathbf{v}), \quad (3.10)$$

where $G_0(\mathbf{u}, \mathbf{v})$ is a constant independent of t , so that when dividing by A_t , the average expected gap decays at rate A_t^{-1} . Accordingly, we wish for A_t to grow as fast as possible with t . By way of convexity, we have that $\text{Gap}^{\mathbf{u}, \mathbf{v}}(\tilde{\mathbf{x}}_t, \tilde{\mathbf{y}}_t) \leq G_0(\mathbf{u}, \mathbf{v})/A_t$ for $(\tilde{\mathbf{x}}_t, \tilde{\mathbf{y}}_t) := A_t^{-1} \sum_{k=1}^t a_k(\mathbf{x}_k, \mathbf{y}_k)$, which can be returned by the algorithm to realize the gap bound (3.10). We may conclude by taking the supremum of $G_0(\mathbf{u}, \mathbf{v})$. In the randomized scenarios, we also discuss in Section 3.6 how the analysis can be adapted to an even stronger convergence criterion for which the supremum is taken prior to the expectation. Because many of the technical ideas remain similar when proving convergence with respect to this stronger criterion, we describe only the parts that change in Section 3.6.

For any algorithm we consider, the analysis will proceed by constructing an upper bound on $a_k \text{Gap}^{\mathbf{u}, \mathbf{v}}(\mathbf{x}_k, \mathbf{y}_k)$ containing telescoping and non-positive terms, by first lower bounding $a_k \mathcal{L}(\mathbf{u}, \mathbf{y}_k)$ and upper bounding $a_k \mathcal{L}(\mathbf{x}_k, \mathbf{v})$. As we will see, the update for \mathbf{x}_k will be used to produce the lower bound while the update for \mathbf{y}_k will be used to produce the upper bound.

Algorithm 4 Template Method

- 1: **Input:** Initial point $(\mathbf{x}_0, \mathbf{y}_0)$, averaging sequence $(a_k)_{k=0}^t$, non-negative weights $(\gamma_I)_{I=1}^N$ that sum to one, balancing sequences $(w_k^P)_{k=1}^t$ and $(w_k^D)_{k=1}^t$, functions SUBROUTINE₁, SUBROUTINE₂, and SUBROUTINE₃.
- 2: Initialize the comparison points $\hat{\mathbf{x}}_{0,I} = \mathbf{x}_0$ for all $I \in [N]$ and $\hat{\mathbf{y}}_0 = \mathbf{y}_0$.
- 3: **for** $k = 1$ **to** t **do**
- 4: SUBROUTINE₁: Compute $\bar{\mathbf{g}}_{k-1}$ using stored information and oracle calls to $\nabla f_i(\mathbf{x}_{k-1})$, $i \in [n]$
- 5: Perform the primal update

$$\begin{aligned} \mathbf{x}_k = \arg \min_{\mathbf{x} \in \mathcal{X}} \bigg\{ & a_k \langle \bar{\mathbf{g}}_{k-1}, \mathbf{x} \rangle + a_k \phi(\mathbf{x}) \\ & + (A_{k-1}\mu + \mu_0) \left(\underbrace{\frac{1-w_k^P}{2} \Delta \mathcal{X}(\mathbf{x}, \mathbf{x}_{k-1})}_{\text{standard proximity term}} + \underbrace{\frac{w_k^P}{2} \sum_{I=1}^N \gamma_I \Delta \mathcal{X}(\mathbf{x}, \hat{\mathbf{x}}_{k-1,I})}_{\text{primal historical regularization}} \right) \bigg\}. \end{aligned} \quad (3.11)$$

- 6: SUBROUTINE₂: Compute $\bar{\mathbf{f}}_{k-1/2}$ using stored information and some calls to $f_1(\mathbf{x}_k), \dots, f_n(\mathbf{x}_k)$.
- 7: Perform the dual update

$$\begin{aligned} \mathbf{y}_k = \arg \max_{\mathbf{y} \in \mathcal{Y}} \bigg\{ & a_k \langle \mathbf{y}, \bar{\mathbf{f}}_{k-1/2} \rangle - a_k \psi(\mathbf{y}) \\ & - (A_{k-1}\nu + \nu_0) \left(\underbrace{\frac{1-w_k^D}{2} \Delta \mathcal{Y}(\mathbf{y}, \mathbf{y}_{k-1})}_{\text{standard proximity term}} + \underbrace{\frac{w_k^D}{2} \Delta \mathcal{Y}(\mathbf{y}, \hat{\mathbf{y}}_{k-1})}_{\text{dual historical regularization}} \right) \bigg\}. \end{aligned} \quad (3.12)$$

- 8: SUBROUTINE₃: Update comparison points $(\hat{\mathbf{x}}_{k,I})_{I=1}^N$ and $\hat{\mathbf{y}}_k$.
 - return** $A_t^{-1} \sum_{k=1}^t a_k(\mathbf{x}_k, \mathbf{y}_k)$.
-

The update rules are motivated directly by the analysis. Similar steps will take place in the stochastic setting, except using the expectation of (3.4) under algorithmic randomness. We start with an arbitrary point $(\mathbf{x}_0, \mathbf{y}_0) \in \text{ri}(\text{dom}(\phi)) \times \text{ri}(\text{dom}(\psi))$. The parameters $\mu_0 > 0$ and $\nu_0 > 0$ appearing below are employed in order to handle the strongly convex and non-strongly convex settings in a unified manner.

This general, high-level idea and analysis template are in line with prior work providing constructive arguments for the analysis of optimization methods [Diakonikolas and Orecchia,

2019, Mehta et al., 2024a, Li et al., 2024b, Diakonikolas, 2025]. As such, it provides a clear guiding principle for the analysis and motivation for the algorithmic choices. It is of note, however, that while the general principle is common to all these works, the specifics of the analysis and associated algorithms differ significantly, as the technical obstacles they need to address are problem-specific. For instance, the stochastic algorithms in the present work have a unique combination of randomized updates and historical regularization that are conceptually novel and interesting in their own right. It is the combination of *both* non-uniform sampling and non-uniform regularization that leads to complexity improvements even in more specific problem classes such as bilinearly coupled problems. For the separable case, our proof technique relies on an auxiliary sequence of dual variables that offers an elegant extension of previous meta-analyses when using coordinate-wise updates.

3.3.1 Algorithm Template and a First Gap Bound

We fix the convention that the update for \mathbf{x}_k occurs before the update for \mathbf{y}_k . Thus, we require that only information available up to and including time $k - 1$ is used in the update. In the stochastic setting, this requirement will be formalized in the language of measurability. Both the primal and dual updates will resemble those of a proximal gradient-type algorithm, wherein \mathbf{x}_k and \mathbf{y}_k are defined by minimizing or maximizing an approximation of (3.2) with a proximity term (i.e., a Bregman divergence). In the primal update, the proximity term promotes \mathbf{x}_k being close to not only \mathbf{x}_{k-1} , but several additional to-be-specified comparison points $\hat{\mathbf{x}}_{k-1,1}, \dots, \hat{\mathbf{x}}_{k-1,N}$. Similarly, \mathbf{y}_k will be made close to \mathbf{y}_{k-1} along with a single comparison point $\hat{\mathbf{y}}_{k-1}$. The use of multiple comparison points is in fact the motivation for the modified three-point inequality (Lemma 3.2.1). This “historical regularization” is visualized in Figure 3.3. The remaining components to specify are $\bar{\mathbf{g}}_{k-1} \in \mathbb{R}^d$, a vector which will be used to linearize the objective (3.2) in the primal update, and $\bar{\mathbf{f}}_{k-1/2} \in \mathbb{R}^n$, an analogous vector used in the dual update. The subscript $k - 1/2$ indicates that “half” of the information in iteration k (namely, the value of \mathbf{x}_k) can be used in the update. The components introduced so far generate a template algorithm which can combine them in various ways;

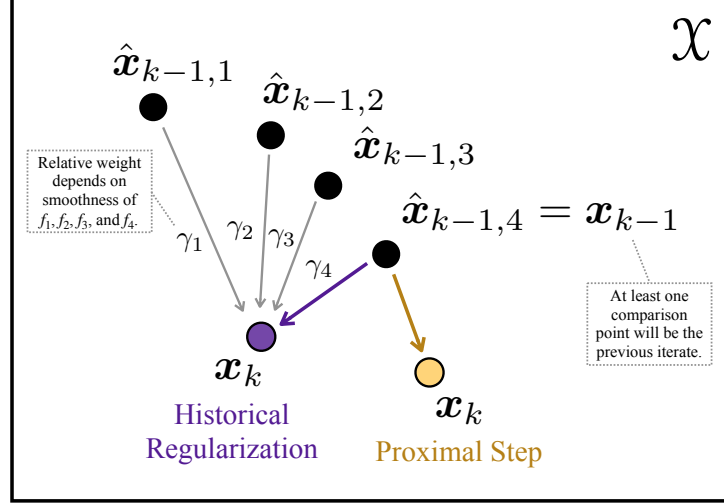


Figure 3.3: **Modified Proximal Step with Historical Regularization.** Geometric illustration of the historical regularization penalty applied in the modified primal update (3.11).

the pseudocode for this method is shown in Algorithm 4, in which the algorithm-specific content is abstracted into three subroutines.

The hyperparameters are shown explicitly to better accompany the theoretical analysis. They are set to specific values over the course of the proofs. We motivate the updates (3.11) and (3.12) using a lower bound on $a_k \mathcal{L}(\mathbf{u}, \mathbf{y}_k)$ and an upper bound on $a_k \mathcal{L}(\mathbf{x}_k, \mathbf{v})$. Several terms will appear that either telescope when summed or are used to cancel errors incurred at each iteration. For the reader's convenience, we summarize this notation below.

Notation: Throughout the analyses for each algorithm, the primal and dual “distance-to-opt” terms will be written as

$$\mathcal{T}_k^P := (A_k\mu + \mu_0)\Delta_{\mathcal{X}}(\mathbf{u}, \mathbf{x}_k), \quad \mathcal{T}_k^D := (A_k\nu + \nu_0)\Delta_{\mathcal{Y}}(\mathbf{u}, \mathbf{y}_k), \quad (3.13)$$

indicating that they telescope when summed. The bounds also produce the negation of the terms

$$\mathcal{C}_k^P := (A_{k-1}\mu + \mu_0)\Delta_{\mathcal{X}}(\mathbf{x}_k, \mathbf{x}_{k-1}), \quad \mathcal{C}_k^D := (A_{k-1}\nu + \nu_0)\Delta_{\mathcal{Y}}(\mathbf{y}_k, \mathbf{y}_{k-1}), \quad (3.14)$$

which appear in the gap function bound and are used for canceling errors that appear when controlling the primal-dual gap. Analogous terms appear based on the comparison points instead of the iterates. That is, consider for each $I \in [N]$ the additional telescoping terms

$$\hat{\mathcal{T}}_{k,I}^P := (A_k\mu + \mu_0)\Delta_{\mathcal{X}}(\mathbf{u}, \hat{\mathbf{x}}_{k,I}), \quad \hat{\mathcal{T}}_k^D := (A_k\nu + \nu_0)\Delta_{\mathcal{Y}}(\mathbf{v}, \hat{\mathbf{y}}_k), \quad (3.15)$$

along with similar cancellation terms

$$\hat{\mathcal{C}}_{k,I}^P := (A_{k-1}\mu + \mu_0)\Delta_{\mathcal{X}}(\mathbf{x}_k, \hat{\mathbf{x}}_{k-1,I}), \quad \hat{\mathcal{C}}_k^D := (A_{k-1}\nu + \nu_0)\Delta_{\mathcal{Y}}(\mathbf{y}_k, \hat{\mathbf{y}}_{k-1}). \quad (3.16)$$

Finally, we will also define the inner product terms

$$\mathcal{I}_k^P = a_k \langle \nabla f(\mathbf{x}_k)^\top \mathbf{y}_k - \bar{\mathbf{g}}_{k-1}, \mathbf{u} - \mathbf{x}_k \rangle, \quad \mathcal{I}_k^D = a_k \langle f(\mathbf{x}_k) - \bar{f}_{k-1/2}, \mathbf{v} - \mathbf{y}_k \rangle, \quad (3.17)$$

which will comprise the errors that are cancelled by the terms above.

To achieve the lower bound, we compare the objective to the one minimized by \mathbf{x}_k and use properties of Bregman divergences to produce telescoping terms akin to a mirror descent-style analysis.

Lemma 3.3.1. *For any $k \geq 1$, let $\bar{\mathbf{g}}_{k-1} \in \mathbb{R}^d$ and $w_k^P, w_{k-1}^P \in [0, 1)$ such that $w_{k-1}^P \leq w_k^P$.*

For $k \geq 1$, if \mathbf{x}_k is defined via the update (3.11), then it holds that

$$a_k \mathcal{L}(\mathbf{u}, \mathbf{y}_k) \geq a_k \mathcal{L}(\mathbf{x}_k, \mathbf{y}_k) + \mathcal{I}_k^P \quad (3.18)$$

$$\begin{aligned} &+ \left(\frac{1-w_k^P}{2} \mathcal{T}_k^P - \frac{1-w_{k-1}^P}{2} \mathcal{T}_{k-1}^P \right) + \frac{w_k^P}{2} \left(\mathcal{T}_k^P - \sum_{I=1}^N \gamma_I \hat{\mathcal{T}}_{k-1,I}^P \right) \\ &+ \frac{1-w_k^P}{2} \mathcal{C}_k^P + \frac{w_k^P}{2} \sum_{I=1}^N \gamma_I \hat{\mathcal{C}}_{k,I}^P + \frac{a_k \mu}{2} \Delta_{\mathcal{X}}(\mathbf{u}, \mathbf{x}_k). \end{aligned} \quad (3.19)$$

Proof. Because \mathbf{y}_k is observed in \mathcal{Y} , we have that $\mathbf{x} \mapsto \langle \mathbf{y}_k, f(\mathbf{x}) \rangle$ is convex and differentiable. As a result,

$$\begin{aligned} a_k \mathcal{L}(\mathbf{u}, \mathbf{y}_k) &= a_k (\mathbf{y}_k^\top f(\mathbf{u}) - \psi(\mathbf{y}_k) + \phi(\mathbf{u})) \\ &\geq a_k \langle \mathbf{y}_k, f(\mathbf{x}_k) + \nabla f(\mathbf{x}_k)(\mathbf{u} - \mathbf{x}_k) \rangle - a_k \psi(\mathbf{y}_k) + a_k \phi(\mathbf{u}). \end{aligned} \quad (3.20)$$

Then, add and subtract terms from the objective defining (3.11) and apply Lemma 3.2.1 with $A = A_{k-1}$, $a = a_k$, $\gamma = \mu$, and $\gamma_0 = \mu_0$ to achieve

$$\begin{aligned} a_k \mathcal{L}(\mathbf{u}, \mathbf{y}_k) &\geq a_k \langle \bar{\mathbf{g}}_{k-1}, \mathbf{u} - \mathbf{x}_k \rangle + a_k \phi(\mathbf{u}) + a_k \langle \nabla f(\mathbf{x}_k)^\top \mathbf{y}_k - \bar{\mathbf{g}}_{k-1}, \mathbf{u} - \mathbf{x}_k \rangle \\ &+ \frac{1-w_k^P}{2} (A_{k-1} \mu + \mu_0) \Delta_{\mathcal{X}}(\mathbf{u}, \mathbf{x}_{k-1}) - \frac{1-w_k^P}{2} \underbrace{(A_{k-1} \mu + \mu_0) \Delta_{\mathcal{X}}(\mathbf{u}, \mathbf{x}_{k-1})}_{\mathcal{T}_{k-1}^P \text{ from (3.13)}} \\ &+ \frac{w_k^P}{2} (A_{k-1} \mu + \mu_0) \sum_{I=1}^N \gamma_I \Delta_{\mathcal{X}}(\mathbf{u}, \hat{\mathbf{x}}_{k-1,I}) - \frac{w_k^P}{2} (A_{k-1} \mu + \mu_0) \underbrace{\sum_{I=1}^N \gamma_I \Delta_{\mathcal{X}}(\mathbf{u}, \hat{\mathbf{x}}_{k-1,I})}_{\sum_I \gamma_I \hat{\mathcal{T}}_{k-1,I}^P \text{ from (3.15)}} \\ &+ a_k \langle \mathbf{y}_k, f(\mathbf{x}_k) \rangle - a_k \psi(\mathbf{y}_k). \end{aligned}$$

When applying Lemma 3.2.1 and using the definitions of \mathcal{C}_k^P from (3.14) and $\hat{\mathcal{C}}_{k,I}^P$ from (3.16),

we achieve the inequality

$$\begin{aligned}
a_k \mathcal{L}(\mathbf{u}, \mathbf{y}_k) &\geq a_k \phi(\mathbf{x}_k) + a_k \langle \nabla f(\mathbf{x}_k)^\top \mathbf{y}_k - \bar{\mathbf{g}}_{k-1}, \mathbf{u} - \mathbf{x}_k \rangle \\
&\quad + \frac{1-w_k^P}{2} \mathcal{C}_k^P - \frac{1-w_k^P}{2} \mathcal{T}_{k-1}^P \\
&\quad + \frac{w_k^P}{2} \sum_{I=1}^N \gamma_I \hat{\mathcal{C}}_{k,I}^P - \frac{w_k^P}{2} \sum_I \gamma_I \hat{\mathcal{T}}_{k-1,I}^P \\
&\quad + a_k \langle \mathbf{y}_k, f(\mathbf{x}_k) \rangle - a_k \psi(\mathbf{y}_k).
\end{aligned}$$

Finally, use $-\frac{1-w_k^P}{2} \mathcal{T}_{k-1}^P \geq -\frac{1-w_{k-1}^P}{2} \mathcal{T}_{k-1}^P$, the substitution $\mathcal{L}(\mathbf{x}_k, \mathbf{y}_k) = \langle \mathbf{y}_k, f(\mathbf{x}_k) \rangle - \psi(\mathbf{y}_k) + \phi(\mathbf{x}_k)$, and the definition of \mathcal{I}_k^P from (3.17) to complete the proof. \square

The upper bound is proved in a nearly identical fashion, except the step that employs convexity in (3.20) is not used; the proof is omitted for brevity.

Lemma 3.3.2. *For any $k \geq 1$, let $\bar{\mathbf{f}}_{k-1/2} \in \mathbb{R}^n$ and $w_k^D, w_{k-1}^D \in (0, 1)$ such that $w_{k-1}^D \leq w_k^D$. For $k \geq 1$, if \mathbf{y}_k is defined via the update (3.12), then it holds that*

$$a_k \mathbb{E}[\mathcal{L}(\mathbf{x}_k, \mathbf{v})] \leq a_k \mathcal{L}(\mathbf{x}_k, \mathbf{y}_k) + \mathcal{I}_k^D \quad (3.21)$$

$$+ \left(\frac{1-w_{k-1}^D}{2} \mathcal{T}_{k-1}^D - \frac{1-w_k^D}{2} \mathcal{T}_k^D \right) + \frac{w_k^D}{2} \left(\hat{\mathcal{T}}_{k-1}^D - \mathcal{T}_k^D \right) \quad (3.22)$$

$$- \frac{1-w_k^D}{2} \mathcal{C}_k^D - \frac{w_k^D}{2} \hat{\mathcal{C}}_k^D - \frac{a_k \nu}{2} \Delta_{\mathbf{y}}(\mathbf{v}, \mathbf{y}_k). \quad (3.23)$$

Combining the derived upper and lower bounds and canceling matching terms, we claim the following result. This bound will be the starting point for every subsequent result in the chapter, so we reference it heavily in the sequel.

Claim 3.3.1. *Using the notation from (3.13), (3.14), (3.15), (3.16), and (3.17), we have*

$$\begin{aligned} \sum_{k=1}^t a_k \text{Gap}^{\mathbf{u}, \mathbf{v}}(\mathbf{x}_k, \mathbf{y}_k) &\leq \left(\frac{1-w_0^{\text{P}}}{2} \mathcal{T}_0^{\text{P}} - \frac{1-w_t^{\text{P}}}{2} \mathcal{T}_t^{\text{P}} \right) + \left(\frac{1-w_0^{\text{D}}}{2} \mathcal{T}_0^{\text{D}} - \frac{1-w_t^{\text{D}}}{2} \mathcal{T}_t^{\text{D}} \right) \\ &\quad + \sum_{k=1}^t [\mathcal{I}_k^{\text{D}} - \mathcal{I}_k^{\text{P}}] - \left[\frac{1-w_k^{\text{P}}}{2} \mathcal{C}_k^{\text{P}} + \frac{1-w_k^{\text{D}}}{2} \mathcal{C}_k^{\text{D}} + \frac{w_k^{\text{P}}}{2} \sum_J \gamma_I \hat{\mathcal{C}}_{k,J}^{\text{P}} + \frac{w_k^{\text{D}}}{2} \hat{\mathcal{C}}_k^{\text{D}} \right] \end{aligned} \quad (3.24)$$

$$+ \frac{1}{2} \sum_{k=1}^t \left[w_k^{\text{P}} \left(\sum_{I=1}^N \gamma_I \hat{\mathcal{T}}_{k-1,I}^{\text{P}} - \mathcal{T}_k^{\text{P}} \right) - a_k \mu \Delta_{\mathcal{X}}(\mathbf{u}, \mathbf{x}_k) \right] \quad (3.25)$$

$$+ \frac{1}{2} \sum_{k=1}^t \left[w_k^{\text{D}} \left(\hat{\mathcal{T}}_{k-1}^{\text{D}} - \mathcal{T}_k^{\text{D}} \right) - a_k \nu \Delta_{\mathcal{Y}}(\mathbf{v}, \mathbf{y}_k) \right] \quad (3.26)$$

The remaining work will be to bound the inner product terms \mathcal{I}_k^{P} and \mathcal{I}_k^{D} by quantities that either telescope or can be cancelled by the remaining terms within (3.24). Then, if $w_k^{\text{P}} > 0$ and $w_k^{\text{D}} > 0$ for any k , we will also control (3.25) and (3.26). For these two lines, we will bound the entire sum over k by a term that does not grow with t . For all three lines, this will rely both on the specific form of $\bar{\mathbf{g}}_{k-1}$ and $\bar{\mathbf{f}}_{k-1/2}$ (which we call the primal and dual gradient estimates) and the growth conditions placed on the sequence $(a_k)_{k \geq 1}$. While these steps may also be modified slightly in the separable setting (see Definition 3.2.1), the format of the analysis remains the same. The growth of the $(a_k)_{k \geq 1}$ sequence (e.g., constant, polynomial, exponential) is determined if the user knows whether the objective is convex or strongly convex and concave or strongly concave in the primal and dual variables, respectively. Regarding hyperparameters, we consider some variants that do not use the historical regularization by setting w_k^{P} and w_k^{D} to zero in (3.11) and (3.12), meaning that $(\gamma_I)_{I=1}^N$ is no longer a hyperparameter that needs to be set. Thus, the number of hyperparameters decreases considerably for each of the cases in Section 3.4 and Section 3.5.

The comparison points (when used) are reflective of SAGA-style variance reduction methods [Defazio et al., 2014, Palaniappan and Bach, 2016]. In general, they are snapshots of previous iterates and may be used not only to define proximity terms but in the definitions of $\bar{\mathbf{g}}_{k-1}$ and $\bar{\mathbf{f}}_{k-1/2}$ as well. The gradient estimates may be computed with adaptive sam-

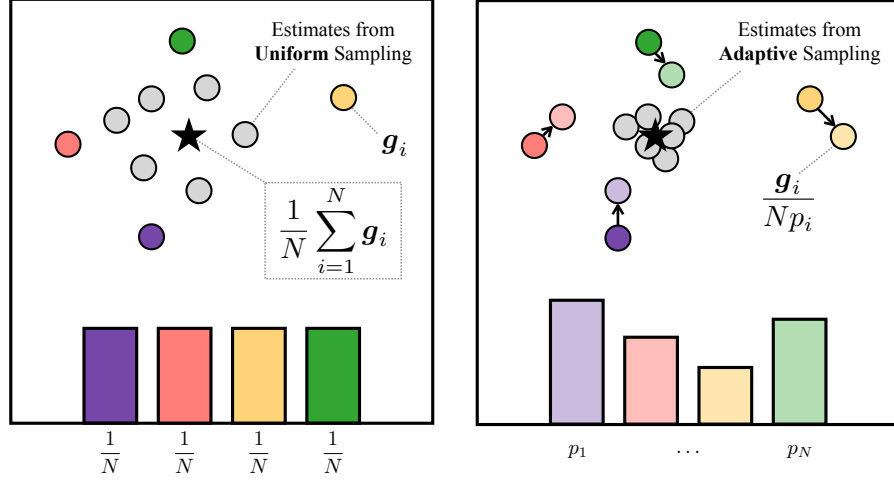


Figure 3.4: **Adaptive Sampling for Gradient Estimation** Geometric illustration of adaptive sampling for computing the mean of vectors $\bar{g}_1, \dots, \bar{g}_N$.

pling, as mentioned in Section 3.1 and visualized in Figure 3.4. By convention, we take any time-dependent element at a negative index to be equal to its initial value (indexed by zero). The number of points stored for the primal updates is equal to the number of blocks N , which may be much smaller than n (avoiding the $O(nd)$ complexity of SAGA). There is only a single additional comparison point in the dual update, incurring a storage cost of $O(n)$. Finally, although the updates include the strong convexity parameters μ and ν , this choice is for readability when proceeding through the analysis. Practically, the growth of the sequence $(A_k)_{k \geq 1}$ will be derived in terms of k with an unknown constant scaling. This constant is a hyperparameter to be searched by the algorithm. For example, when $\mu > 0$ and $\nu > 0$, we have that $a_{k+1} = \alpha A_k = \alpha \sum_{i=1}^k a_i$ for $k \geq 1$, where $\alpha > 0$ will be a tunable hyperparameter whose optimal value depends on μ and ν . We conclude this section with a deterministic full vector update method to provide intuition and work with Claim 3.3.1, whereas more advanced algorithms are presented in Section 3.4 and Section 3.5.

3.3.2 Warm-up: deterministic algorithm

Here, we allow the algorithm to access all first-order oracles $\{(f_i, \nabla f_i)\}_{i=1}^n$ in each iteration, for a total per-iteration cost of $O(nd)$. To observe how the gradient estimates can be set to control (3.17), first consider

$$\bar{\mathbf{g}}_{k-1} = \nabla f(\mathbf{x}_{k-1})^\top \mathbf{y}_{k-1} + \frac{a_{k-1}}{a_k} (\nabla f(\mathbf{x}_{k-1})^\top \mathbf{y}_{k-1} - \nabla f(\mathbf{x}_{k-2})^\top \mathbf{y}_{k-2}). \quad (3.27)$$

Then, by substituting (3.27) into \mathcal{I}_k^P , we have that

$$\begin{aligned} \mathcal{I}_k^P &= a_k \langle \nabla f(\mathbf{x}_k)^\top \mathbf{y}_k - \nabla f(\mathbf{x}_{k-1})^\top \mathbf{y}_{k-1}, \mathbf{u} - \mathbf{x}_k \rangle \\ &\quad - a_{k-1} \langle \nabla f(\mathbf{x}_{k-1})^\top \mathbf{y}_{k-1} - \nabla f(\mathbf{x}_{k-2})^\top \mathbf{y}_{k-2}, \mathbf{u} - \mathbf{x}_k \rangle \\ &= a_k \langle \nabla f(\mathbf{x}_k)^\top \mathbf{y}_k - \nabla f(\mathbf{x}_{k-1})^\top \mathbf{y}_{k-1}, \mathbf{u} - \mathbf{x}_k \rangle \\ &\quad - a_{k-1} \langle \nabla f(\mathbf{x}_{k-1})^\top \mathbf{y}_{k-1} - \nabla f(\mathbf{x}_{k-2})^\top \mathbf{y}_{k-2}, \mathbf{u} - \mathbf{x}_{k-1} \rangle - \mathcal{E}_k^P \end{aligned} \quad (3.28)$$

where we define the error term

$$\mathcal{E}_k^P = a_{k-1} \langle \nabla f(\mathbf{x}_{k-1})^\top \mathbf{y}_{k-1} - \nabla f(\mathbf{x}_{k-2})^\top \mathbf{y}_{k-2}, \mathbf{x}_{k-1} - \mathbf{x}_k \rangle. \quad (3.29)$$

We further set $w_k^P = w_k^D = 0$ to simplify the result of Claim 3.3.1 significantly to become

$$\begin{aligned} \sum_{k=1}^t a_k \text{Gap}^{\mathbf{u}, \mathbf{v}}(\mathbf{x}_k, \mathbf{y}_k) &\leq \mathcal{T}_0^P - \mathcal{T}_t^P + \mathcal{T}_0^D - \mathcal{T}_t^D \\ &\quad + \sum_{k=1}^{t-1} \mathcal{E}_k^P - a_k \langle \nabla f(\mathbf{x}_k)^\top \mathbf{y}_k - \nabla f(\mathbf{x}_{k-1})^\top \mathbf{y}_{k-1}, \mathbf{u} - \mathbf{x}_k \rangle - \sum_{k=1}^t (\mathcal{C}_k^P + \mathcal{C}_k^D). \end{aligned} \quad (3.30)$$

Evidently, the goal is for the sum of terms in (3.30) to be a quantity that does not grow with t . A key step in this proof (and the proofs of upcoming results) will be to bound above \mathcal{E}_k^P such that the resulting terms can be canceled by \mathcal{C}_k^P and \mathcal{C}_k^D from (3.14). While we may encounter a similar term when bounding \mathcal{I}_k^D , because \mathbf{x}_k can be used in the definition of $\bar{\mathbf{f}}_{k-1/2}$ and $O(nd)$ operations are permitted, we may simply set $\bar{\mathbf{f}}_{k-1/2} = f(\mathbf{x}_k)$ to make \mathcal{I}_k^D vanish. Because the terms \mathcal{E}_k^P and \mathcal{E}_k^D appear in some form for every algorithm, we can use

their definitions (and the definitions of subroutines from Algorithm 4) as an identity card for each algorithm, in the form below.

Identity Card 1: Full vector update method	
SUBROUTINE ₁ :	$\bar{\mathbf{g}}_{k-1} = \nabla f(\mathbf{x}_{k-1})^\top \mathbf{y}_{k-1} + \frac{a_{k-1}}{a_k} (\nabla f(\mathbf{x}_{k-1})^\top \mathbf{y}_{k-1} - \nabla f(\mathbf{x}_{k-2})^\top \mathbf{y}_{k-2})$
SUBROUTINE ₂ :	$\bar{\mathbf{f}}_{k-1/2} = f(\mathbf{x}_k)$
SUBROUTINE ₃ :	None
<hr/>	
<i>Primal error:</i>	$\mathcal{E}_k^P = a_{k-1} \langle \nabla f(\mathbf{x}_{k-1})^\top \mathbf{y}_{k-1} - \nabla f(\mathbf{x}_{k-2})^\top \mathbf{y}_{k-2}, \mathbf{x}_{k-1} - \mathbf{x}_k \rangle$
<i>Dual error:</i>	$\mathcal{E}_k^D = 0$

The necessary steps are performed in the following proposition, which also establishes the growth rate of $(a_k)_{k \geq 1}$.

Proposition 3.3.1. *Let $(\mathbf{x}_0, \mathbf{y}_0) \in \text{ri}(\text{dom}(\phi)) \times \text{ri}(\text{dom}(\psi))$ and $\{(\mathbf{x}_k, \mathbf{y}_k)\}_{k \geq 1}$ be generated by the updates (3.11) and (3.12), with $\bar{\mathbf{g}}_{k-1}$ given by (3.27) and $\bar{\mathbf{f}}_{k-1/2} = f(\mathbf{x}_k)$. Select $(a_k)_{k \geq 1}$ to satisfy*

$$a_k \leq \min \left\{ \frac{\sqrt{(A_k \mu + \mu_0)(A_{k-1} \nu + \nu_0)}}{\sqrt{2}G}, \frac{\sqrt{(A_k \mu + \mu_0)(A_{k-1} \mu + \mu_0)}}{2L} \right\}, \quad (3.31)$$

Recalling the notation of (3.13), we have that for any $\mathbf{u} \in \mathcal{X}$ and $\mathbf{v} \in \mathcal{Y}$,

$$\sum_{k=1}^t a_k \text{Gap}^{\mathbf{u}, \mathbf{v}}(\mathbf{x}_k, \mathbf{y}_k) + \frac{1}{2} \mathcal{T}_t^P + \mathcal{T}_t^D \leq \mathcal{T}_0^P + \mathcal{T}_0^D. \quad (3.32)$$

Proof. We proceed from steps leading up to the gap bound (3.30). Consider $k \geq 2$ (as $\mathcal{E}_k^P = 0$ for $k \leq 1$). Apply Young's inequality with parameter $(A_{k-1} \mu + \mu_0)/2$ and the strong convexity of Bregman divergences to write

$$\mathcal{E}_k^P \leq \frac{1}{2} \mathcal{C}_k^P + \frac{a_{k-1}^2}{A_{k-1} \mu + \mu_0} \left\| \nabla f(\mathbf{x}_{k-1})^\top \mathbf{y}_{k-1} - f(\mathbf{x}_{k-2})^\top \mathbf{y}_{k-2} \right\|_{\mathcal{X}^*}^2.$$

To bound the second term above, we first decompose it via Young's inequality. Write

$$\begin{aligned}
& \left\| \nabla f(\mathbf{x}_{k-1})^\top \mathbf{y}_{k-1} - f(\mathbf{x}_{k-2})^\top \mathbf{y}_{k-2} \right\|_{\mathcal{X}^*}^2 \\
& \leq 2 \left\| (\nabla f(\mathbf{x}_{k-1}) - \nabla f(\mathbf{x}_{k-2}))^\top \mathbf{y}_{k-1} \right\|_{\mathcal{X}^*}^2 \\
& \quad + 2 \left\| \nabla f(\mathbf{x}_{k-1})^\top (\mathbf{y}_{k-1} - \mathbf{y}_{k-2}) \right\|_{\mathcal{X}^*}^2 \\
& \leq 4L^2 \Delta_{\mathcal{X}}(\mathbf{x}_{k-1}, \mathbf{x}_{k-2}) + 4G^2 \Delta_{\mathcal{Y}}(\mathbf{y}_{k-1}, \mathbf{y}_{k-2}),
\end{aligned} \tag{3.33}$$

where the last inequality follows by (3.7) and (3.8) from Assumption 3.2.2. Recall that $\mathcal{C}_k^{\text{P}} = (A_{k-1}\mu + \mu_0)\Delta_{\mathcal{X}}(\mathbf{x}_k, \mathbf{x}_{k-1})$ and $\mathcal{C}_k^{\text{D}} = (A_{k-1}\nu + \nu_0)\Delta_{\mathcal{Y}}(\mathbf{y}_k, \mathbf{y}_{k-1})$ (see (3.14)). Combining the steps above for $k \geq 2$ and the condition (3.31) yields

$$\begin{aligned}
\mathcal{E}_k^{\text{P}} & \leq \frac{1}{2}\mathcal{C}_k^{\text{P}} + \frac{4a_{k-1}^2}{(A_{k-1}\mu + \mu_0)} \left[\frac{L^2\mathcal{C}_{k-1}^{\text{P}}}{A_{k-2}\mu + \mu_0} + \frac{G^2\mathcal{C}_{k-1}^{\text{D}}}{A_{k-2}\nu + \nu_0} \right] \\
& \leq \frac{1}{2}\mathcal{C}_k^{\text{P}} + \frac{1}{2}\mathcal{C}_{k-1}^{\text{P}} + \mathcal{C}_{k-1}^{\text{D}}.
\end{aligned} \tag{3.34}$$

Note that for the case of $k = 2$, our choice of a_1 satisfies (3.31) as well. Summing up the current gap bound over $k = 1, \dots, t$ and dropping non-positive terms yields

$$\begin{aligned}
\sum_{k=1}^t \text{Gap}^{u,v}(\mathbf{x}_k, \mathbf{y}_k) & \leq \mathcal{T}_0^{\text{P}} - \mathcal{T}_t^{\text{P}} + \mathcal{T}_0^{\text{D}} - \mathcal{T}_t^{\text{D}} \\
& \quad - \left(\frac{1}{2}\mathcal{C}_t^{\text{P}} + \mathcal{C}_t^{\text{D}} \right) - a_t \langle \nabla f(\mathbf{x}_t)^\top \mathbf{y}_t - \nabla f(\mathbf{x}_{t-1})^\top \mathbf{y}_{t-1}, \mathbf{u} - \mathbf{x}_t \rangle.
\end{aligned} \tag{3.35}$$

For the remaining inner product term, we apply Young's inequality with parameter $(A_t\mu + \mu_0)/2$ and apply a similar argument as for (3.34):

$$\begin{aligned}
& a_t \langle \nabla f(\mathbf{x}_t)^\top \mathbf{y}_t - \nabla f(\mathbf{x}_{t-1})^\top \mathbf{y}_{t-1}, \mathbf{u} - \mathbf{x}_t \rangle \\
& = -a_t \langle \nabla f(\mathbf{x}_t)^\top \mathbf{y}_t - \nabla f(\mathbf{x}_{t-1})^\top \mathbf{y}_{t-1}, \mathbf{u} - \mathbf{x}_t \rangle \\
& \leq \frac{a_t^2}{A_t\mu + \mu_0} \left\| \nabla f(\mathbf{x}_t)^\top \mathbf{y}_t - \nabla f(\mathbf{x}_{t-1})^\top \mathbf{y}_{t-1} \right\|_{\mathcal{X}^*}^2 + \frac{1}{2}\mathcal{T}_t^{\text{P}} \\
& \leq \frac{1}{2}\mathcal{C}_t^{\text{P}} + \mathcal{C}_t^{\text{D}} + \frac{1}{2}\mathcal{T}_t^{\text{P}},
\end{aligned}$$

which will each be cancelled by terms in (3.35), leading to the claimed bound. \square

To convert the convergence guarantee in Proposition 3.3.1 into a complexity result, we consider the four possible cases for whether $\mu = 0$ and/or $\nu = 0$, which proves Theorem 3.3.1.

Theorem 3.3.1. *Under Assumption 3.2.2 and Assumption 3.2.1, consider any $\mathbf{u} \in \mathcal{X}$, $\mathbf{v} \in \mathcal{Y}$, and precision $\varepsilon > 0$. Define the initial distance term*

$$D_0 = \sqrt{\frac{\mu_0}{\nu_0}} \Delta_{\mathcal{X}}(\mathbf{u}, \mathbf{x}_0) + \sqrt{\frac{\nu_0}{\mu_0}} \Delta_{\mathcal{Y}}(\mathbf{v}, \mathbf{y}_0). \quad (3.36)$$

There exists a choice of the sequence $(a_k)_{k=1}^t$ such that Algorithm 4 with Identity Card 1 produces an output point $(\tilde{\mathbf{x}}_t, \tilde{\mathbf{y}}_t) \in \mathcal{X} \times \mathcal{Y}$ satisfying $\text{Gap}^{\mathbf{u}, \mathbf{v}}(\tilde{\mathbf{x}}_t, \tilde{\mathbf{y}}_t) \leq \varepsilon$ for t that depends on ε according to the following iteration complexities.

Case	Iteration Complexity
$\mu > 0$ and $\nu > 0$	$O\left(\left(\frac{L}{\mu} + \frac{G}{\sqrt{\mu\nu}}\right) \ln\left(\frac{(L\sqrt{\nu_0/\mu_0} + G)D_0}{\varepsilon}\right)\right)$
$\mu > 0$ and $\nu = 0$	$O\left(\frac{L}{\mu} \ln\left(\frac{(L\sqrt{\nu_0/\mu_0} + G)D_0}{\varepsilon}\right) + G\sqrt{\frac{\sqrt{\mu_0/\nu_0}D_0}{\mu\varepsilon}}\right)$
$\mu = 0$ and $\nu > 0$	$O\left(\frac{L\sqrt{\nu_0/\mu_0}D_0}{\varepsilon} + G\sqrt{\frac{\sqrt{\nu_0/\mu_0}D_0}{\nu\varepsilon}}\right)$
$\mu = 0$ and $\nu = 0$	$O\left(\frac{(L\sqrt{\nu_0/\mu_0} + G)D_0}{\varepsilon}\right)$

Proof. We first determine the growth of the sequence A_t , so that A_t^{-1} gives the convergence rate in terms of the number of iterations. The growth rate can be derived by providing a sequence $(a_k)_{k \geq 0}$ such that (3.31) is satisfied. For the dependence of the required number of iterations on the suboptimality parameter ε , we write $G_0 A_t^{-1} \leq \varepsilon$ and solve for t , noting that $G_0 = \mu_0 \Delta_{\mathcal{X}}(\mathbf{u}, \mathbf{x}_0) + \nu_0 \Delta_{\mathcal{Y}}(\mathbf{v}, \mathbf{y}_0)$. In all cases, set $a_1 = \min\{\sqrt{\mu_0\nu_0}/(\sqrt{2}G), \mu_0/(2L)\}$ so that the condition (3.31) is satisfied, and

$$\frac{G_0}{a_1} = \max\left\{2L\sqrt{\frac{\nu_0}{\mu_0}}, \sqrt{2}G\right\} D_0 = O\left(\left(L\sqrt{\frac{\nu_0}{\mu_0}} + G\right) D_0\right).$$

Although the $(a_k)_{k \geq 1}$ values are unitless, the relative quantity can help determine the optimal values for μ_0 and ν_0 .

Case 1: $\mu > 0, \nu > 0$. Let $\alpha = \min \left\{ \frac{\sqrt{\mu\nu}}{\sqrt{2G}}, \frac{\mu}{2L} \right\}$. For $k \geq 2$, write $A_k - A_{k-1} = a_k = \alpha \sqrt{A_k A_{k-1}} \geq \alpha A_{k-1}$, which implies that $A_t \geq (1 + \alpha)^t a_1$. Then,

$$\frac{G_0}{A_t} \leq \frac{G_0}{a_1(1 + \alpha)^t} \lesssim (1 + \alpha)^{-t} \left(L \sqrt{\frac{\nu_0}{\mu_0}} + G \right) D_0 \stackrel{\text{want}}{\leq} \varepsilon,$$

which is satisfied for t at the given big- O order.

Case 2: $\mu > 0, \nu = 0$. We have that $a_k = \min \left\{ \frac{\mu\nu_0}{4G^2} k, \frac{\mu}{2L} A_{k-1} \right\}$ for $k \geq 2$ satisfies the rate condition. Then, there exists a $k^* \geq 0$ such that $A_t \geq (1 + \frac{\mu}{2L})^t a_1$ for all $t < k^*$ and $A_t \geq \frac{\mu\nu_0}{4G^2} \sum_{k=k^*+1}^t k + (1 + \frac{\mu}{2L})^{k^*} a_1$ for $t \geq k^*$. To compute the complexity, we consider when either term is dominant. For $(1 + \frac{\mu}{2L})^{k^*} a_1$, we apply the same argument as Case 1. Otherwise, we have that

$$\frac{G_0}{A_t} \lesssim \frac{G^2}{\mu\nu_0 t^2} (\mu_0 \Delta_x(\mathbf{u}, \mathbf{x}_0) + \nu_0 \Delta_y(\mathbf{v}, \mathbf{y}_0)) = \frac{G^2}{\mu t^2} \sqrt{\frac{\mu_0}{\nu_0}} D_0 \stackrel{\text{want}}{\leq} \varepsilon.$$

Case 3: $\mu = 0, \nu > 0$. For $k \geq 2$, we have that $a_k = \min \left\{ \frac{\mu_0\nu}{4G^2} k, \frac{\mu_0}{2L} \right\}$ satisfies the rate condition by direct computation. Thus, for $k \leq k^* = \frac{G^2}{2\nu L}$, we have that $A_t \geq \frac{\mu_0\nu t(t+1)}{8G^2}$ (for which we argue similarly to Case 2 above), and otherwise, $A_t \geq \frac{\mu_0(t-k^*)}{2L} + \frac{\mu_0\nu k^*(k^*+1)}{8G^2}$ (for which we argue similarly to Case 4 below).

Case 4: $\mu = 0, \nu = 0$. Here, a_k is equal to a constant, so $A_t = a_1 t$. Then, arguing similarly to Case 1,

$$\frac{G_0}{A_t} \leq \frac{G_0}{a_1 t} \lesssim t^{-1} \left(L \sqrt{\frac{\nu_0}{\mu_0}} + G \right) D_0 \stackrel{\text{want}}{\leq} \varepsilon,$$

which is satisfied for t at the given big- O order, completing the proof. \square

We assume that the cost of querying the first-order oracle $\mathbf{x} \mapsto (f_j(\mathbf{x}), \nabla f_j(\mathbf{x}))$ is $O(d)$ for any $j = 1, \dots, n$, and that the optimization problems (3.11) defining the primal update

and (3.12) defining the dual update can be solved at $\tilde{O}(d)$ and $\tilde{O}(n)$ cost, respectively. Each iteration requires querying all first-order oracles and furthermore computes matrix-vector products using $n \times d$ matrices, thus the overall arithmetic complexity of each iteration of Algorithm 4 with Identity Card 1 is $\tilde{O}(nd)$.

3.4 Stochastic Algorithms for General Objectives

In the case of a randomized algorithm, we allow the vectors $\bar{\mathbf{g}}_{k-1}$ and $\bar{\mathbf{f}}_{k-1/2}$ to be defined based on a randomly chosen subset of the indices in $[n]$. This amounts to accessing first-order information $(f_j, \nabla f_j)$ for only some $j \in [n]$. The expressions for $\bar{\mathbf{g}}_{k-1}$ and $\bar{\mathbf{f}}_{k-1/2}$ will depend on historical values of the primal and dual iterations, in the spirit of variance reduction or random extrapolation for convex minimization (see, e.g., Gower et al. [2020]). We first describe precisely which comparison points are stored by the algorithm and then specify how they are used to compute $\bar{\mathbf{g}}_{k-1}$ and $\bar{\mathbf{f}}_{k-1/2}$. We store N previous primal *iterates* $\hat{\mathbf{x}}_{k,1}, \dots, \hat{\mathbf{x}}_{k,N} \in \mathcal{X}$ and a collection of past dual *coordinate blocks* $\hat{\mathbf{y}}_k = (\hat{\mathbf{y}}_{k,1}, \dots, \hat{\mathbf{y}}_{k,N}) \in \mathcal{Y}$ associated to each block index. We define a primal gradient table $\hat{\mathbf{g}}_k = (\hat{\mathbf{g}}_{k,1}, \dots, \hat{\mathbf{g}}_{k,n}) \in \mathbb{R}^{n \times d}$ and dual gradient table $\hat{\mathbf{f}}_k = (\hat{f}_{k,1}, \dots, \hat{f}_{k,n}) \in \mathbb{R}^n$ constructed via

$$(\hat{f}_{k,i}, \hat{\mathbf{g}}_{k,i}) = (f_i(\hat{\mathbf{x}}_{k,I}), \nabla f_i(\hat{\mathbf{x}}_{k,I})) \text{ for all } i \in B_I. \quad (3.37)$$

In other words, the tables contain the first-order information of each f_i in block B_I at $\hat{\mathbf{x}}_{k,I}$. Note that we do not necessarily need to store the $O(nd)$ -sized gradient table, and need only store the $O(Nd)$ -sized table of comparison points, which can be much smaller. We update these tables randomly at the end of each iteration by independently sampling block indices R_k and S_k (possibly non-uniformly) and setting each block to

$$\hat{\mathbf{x}}_{k,I} = \begin{cases} \mathbf{x}_k & \text{if } I = R_k \\ \hat{\mathbf{x}}_{k-1,I} & \text{otherwise} \end{cases} \text{ and } \hat{\mathbf{y}}_{k,I} = \begin{cases} \mathbf{y}_{k,I} & \text{if } I = S_k \\ \hat{\mathbf{y}}_{k-1,I} & \text{otherwise} \end{cases}. \quad (3.38)$$

As mentioned in Section 3.3, we define $(\hat{\mathbf{x}}_{k,I}, \hat{\mathbf{y}}_{k,I}) = (\mathbf{x}_0, \mathbf{y}_{0,I})$ for any block B_I and iteration $k < 0$. The probability weights that govern the randomness in R_k and S_k are denoted as $\mathbf{r} = (r_1, \dots, r_N)$, and $\mathbf{s} = (s_1, \dots, s_N)$, respectively.

Next, for computing the primal and dual gradient estimates, we sample two more block indices P_k and Q_k with associated probability mass vectors $\mathbf{p} = (p_1, \dots, p_N)$ and $\mathbf{q} = (q_1, \dots, q_N)$, respectively. Letting \mathbf{e}_j denote the j -th standard basis vector in \mathbb{R}^n , we construct

$$\bar{\mathbf{g}}_{k-1} = \hat{\mathbf{g}}_{k-1}^\top \hat{\mathbf{y}}_{k-1} + \frac{a_{k-1}}{p_{P_k} a_k} \sum_{i \in B_{P_k}} (y_{k-1,i} \nabla f_i(\mathbf{x}_{k-1}) - \hat{y}_{k-2,i} \hat{\mathbf{g}}_{k-2,i}) \quad (3.39)$$

$$\bar{\mathbf{f}}_{k-1/2} = \hat{\mathbf{f}}_k + \frac{a_{k-1}}{q_{Q_k} a_k} \sum_{j \in B_{Q_k}} (f_j(\mathbf{x}_{k-1}) - \hat{f}_{k-1,j}) \mathbf{e}_j. \quad (3.40)$$

Even though \mathbf{x}_k is known during the update of \mathbf{y}_k , notice that we do not set $\bar{\mathbf{f}}_{k-1/2} = f(\mathbf{x}_k)$ in this setting to avoid an $\tilde{O}(nd)$ per-iteration complexity. Also notice that $\hat{\mathbf{g}}_{k-1}^\top \hat{\mathbf{y}}_{k-1}$ can be maintained at an $O(nd/N)$ cost per iteration on average (as opposed to $O(nd)$), because both $\hat{\mathbf{y}}_k$ and $\hat{\mathbf{g}}_k$ only change within a single coordinate block each (R_k and S_k , respectively). This is discussed in detail alongside the per-iteration complexities of specific algorithms.

We collect here the probabilistic notation used in this chapter. We have introduced four random variables for any iteration k : $P_k \sim \mathbf{p}$, $Q_k \sim \mathbf{q}$, $R_k \sim \mathbf{r}$, and $S_k \sim \mathbf{s}$. To formally analyze the resulting algorithm, consider a filtered probability space $\mathcal{P} = (\Omega, (\mathcal{F}_k)_{k \geq 0}, \mathbb{P})$, where we use the natural filtration $(\mathcal{F}_k)_{k \geq 0}$ (with $\mathcal{F}_0 = \{\emptyset, \Omega\}$ and \mathcal{F}_k being the σ -algebra generated by the collection $\{(P_\kappa, Q_\kappa, R_\kappa, S_\kappa)\}_{\kappa=1}^k$). We also denote by $\mathcal{F}_{k-1/2}$ the σ -algebra generated by $\{(P_\kappa, Q_\kappa, R_\kappa, S_\kappa)\}_{\kappa=1}^{k-1} \cup \{P_k\}$, which captures information up until and including the computation of \mathbf{x}_k (but not \mathbf{y}_k). Thus, in the language of probability theory, we will say that \mathbf{x}_k is $\mathcal{F}_{k-1/2}$ -measurable and \mathbf{y}_k is \mathcal{F}_k -measurable. The full or marginal expectation on \mathcal{P} is given by \mathbb{E} , whereas the conditional expectation given \mathcal{F}_k is denoted by \mathbb{E}_k . We let $\mathbf{z}_{k,J} = (z_{k,j})_{j \in B_J}$ be the block B_J coordinates of a time-indexed vector $\mathbf{z}_k \in \mathbb{R}^n$. In the arguments below, we will always consider (\mathbf{u}, \mathbf{v}) that is independent of $\{(P_\kappa, Q_\kappa, R_\kappa, S_\kappa)\}_{\kappa \geq 1}$.

In Section 3.6.1, we then precisely describe how this assumption can be relaxed to achieve the same complexity guarantees for a stronger convergence criterion than the one in Theorem 3.4.1 and Theorem 3.4.2.

We proceed to the convergence analysis. The primal-dual sequence $(\mathbf{x}_k, \mathbf{y}_k)_{k \geq 0}$ is now a stochastic process; we do not distinguish random variables and realizations when clear from context. We will use the following fact throughout this section: because $\|\cdot\|_{\mathbf{y}}$ is an ℓ_p -norm with $p \in [1, 2]$ (see Section 3.2), it holds that $\|\cdot\|_{\mathbf{y}} \geq \|\cdot\|_2$ and $\|\cdot\|_{\mathbf{y}^*} \leq \|\cdot\|_2$. We will use similar techniques as before to upper bound $\sum_{k=1}^t a_k \mathbb{E}[\text{Gap}^{\mathbf{u}, \mathbf{v}}(\mathbf{x}_k, \mathbf{y}_k)]$. First, either by using Lemma 3.3.1 and Lemma 3.3.2, we produce a lower bound for $a_k \mathcal{L}(\mathbf{u}, \mathbf{y}_k)$ and an upper bound for $a_k \mathcal{L}(\mathbf{x}_k, \mathbf{v})$. Before stating these results, we describe the aspects of the analysis that are similar to the algorithm from Section 3.3.2. Recall the notation box from Section 3.3.

As in Section 3.3.2, we need to use the structure of $\bar{\mathbf{g}}_{k-1}$ (defined in (3.39)) and $\bar{\mathbf{f}}_{k-1/2}$ (defined in (3.40)) and conditions on $(a_k)_{k \geq 1}$ to control (the expected value of) the inner product terms that appear in (3.42) and (3.43) below. We describe the analogous argument to the one leading to (3.29), except for the stochastic setting. Using $\bar{\mathbf{g}}_{k-1}$ as an example, we take the conditional expectation of \mathcal{I}_k^{P} from (3.17) given \mathcal{F}_{k-1} for $k \geq 1$ (recalling that \mathbf{u} is independent of the algorithm randomness) to write

$$\begin{aligned} \mathbb{E}_{k-1}[\mathcal{I}_k^{\text{P}}] &= a_k \mathbb{E}_{k-1} \langle \nabla f(\mathbf{x}_k)^\top \mathbf{y}_k - \hat{\mathbf{g}}_{k-1}^\top \hat{\mathbf{y}}_{k-1}, \mathbf{u} - \mathbf{x}_k \rangle \\ &\quad - a_{k-1} \langle \nabla f(\mathbf{x}_{k-1})^\top \mathbf{y}_{k-1} - \hat{\mathbf{g}}_{k-2}^\top \hat{\mathbf{y}}_{k-2}, \mathbf{u} - \mathbf{x}_{k-1} \rangle - \mathbb{E}_{k-1}[\mathcal{E}_k^{\text{P}}] \end{aligned} \quad (3.41)$$

for the error term

$$\mathcal{E}_k^{\text{P}} = a_{k-1} \left\langle \frac{1}{p_{P_k}} \sum_{i \in B_{P_k}} (y_{k-1,i} \nabla f_i(\mathbf{x}_{k-1}) - \hat{y}_{k-2,i} \hat{\mathbf{g}}_{k-2,i}), \mathbf{x}_{k-1} - \mathbf{x}_k \right\rangle. \quad (3.42)$$

Note that the telescoping occurs when taking the marginal expectation \mathbb{E} over the entire sequence. The term \mathcal{E}_k^{D} is defined analogously by substituting $\bar{\mathbf{f}}_{k-1/2}$ from (3.40) into the

expression for \mathcal{I}_k^{D} , yielding

$$\mathcal{E}_k^{\text{D}} = a_{k-1} \left\langle \frac{1}{q_{Q_k}} \sum_{j \in B_{Q_k}} (f_j(\mathbf{x}_{k-1}) - f_j(\hat{\mathbf{x}}_{k-1, Q_k})) \mathbf{e}_j, \mathbf{y}_k - \mathbf{y}_{k-1} \right\rangle. \quad (3.43)$$

We summarize the above using an identity card.

Identity Card 2: Stochastic update method for general objectives

$$\text{SUBROUTINE}_1: \bar{\mathbf{g}}_{k-1} = \hat{\mathbf{g}}_{k-1}^\top \hat{\mathbf{y}}_{k-1} + \frac{a_{k-1}}{p_{P_k} a_k} \sum_{i \in B_{P_k}} (y_{k-1,i} \nabla f_i(\mathbf{x}_{k-1}) - \hat{y}_{k-2,i} \hat{\mathbf{g}}_{k-2,i})$$

$$\text{SUBROUTINE}_2: \bar{\mathbf{f}}_{k-1/2} = \hat{\mathbf{f}}_k + \frac{a_{k-1}}{q_{Q_k} a_k} \sum_{j \in B_{Q_k}} (f_j(\mathbf{x}_{k-1}) - \hat{f}_{k-1,j}) \mathbf{e}_j$$

$$\text{SUBROUTINE}_3: \text{Update } (\hat{\mathbf{x}}_{k,I}, \hat{\mathbf{y}}_{k,I}) \text{ for all } I \text{ using (3.38) and } (\hat{\mathbf{g}}_k, \hat{\mathbf{f}}_k) \text{ using (3.37).}$$

$$\text{Primal error: } \mathcal{E}_k^P = a_{k-1} \left\langle \frac{1}{p_{P_k}} \sum_{i \in B_{P_k}} (y_{k-1,i} \nabla f_i(\mathbf{x}_{k-1}) - \hat{y}_{k-2,i} \hat{\mathbf{g}}_{k-2,i}), \mathbf{x}_{k-1} - \mathbf{x}_k \right\rangle$$

$$\text{Dual error: } \mathcal{E}_k^D = a_{k-1} \left\langle \frac{1}{q_{Q_k}} \sum_{j \in B_{Q_k}} (f_j(\mathbf{x}_{k-1}) - f_j(\hat{\mathbf{x}}_{k-1,Q_k})) \mathbf{e}_j, \mathbf{y}_k - \mathbf{y}_{k-1} \right\rangle$$

Notice that the terms \mathcal{E}_k^P and \mathcal{E}_k^D that appear in (3.42) and (3.43) are measurements of “table bias”, or how stale the elements in the tables are compared to the current iterates \mathbf{x}_k (for \mathcal{E}_k^P) and \mathbf{y}_k (for \mathcal{E}_k^D). The algorithms below provide two different strategies for achieving convergence while controlling these errors. Because terms of the form \mathcal{E}_k^P will appear in multiple analyses, we collect a repeatedly used bound below.

Lemma 3.4.1. *Consider \mathcal{F}_{k-1} -measurable random vectors \mathbf{x} and $\hat{\mathbf{x}}_1, \dots, \hat{\mathbf{x}}_N$ realized in \mathcal{X} , and similarly, let \mathcal{F}_{k-1} -measurable \mathbf{y} and $\hat{\mathbf{y}}$ be realized in \mathcal{Y} . For any collection of positive constants $(b_I)_{I=1}^N$ and $(c_I)_{I=1}^N$, we have that*

$$\begin{aligned} & \mathbb{E} \left\| \frac{1}{p_{P_k}} \sum_{i \in B_{P_k}} (y_i \nabla f_i(\mathbf{x}) - \hat{y}_i \nabla f_i(\hat{\mathbf{x}}_{P_k})) \right\|_{\mathcal{X}^*}^2 \\ & \leq 2 \left(\max_J \frac{\mathbf{L}_J^2}{p_J b_J} \right) \sum_{I=1}^N b_I \mathbb{E} \|\mathbf{x} - \hat{\mathbf{x}}_I\|_{\mathcal{X}}^2 + 2 \left(\max_J \frac{\mathbf{G}_J^2}{p_J c_J} \right) \sum_{I=1}^N c_I \mathbb{E} \|\mathbf{y}_I - \hat{\mathbf{y}}_I\|_2^2. \end{aligned}$$

Proof. First, take the conditional expectation given \mathcal{F}_{k-1} so that

$$\begin{aligned} & \mathbb{E}_{k-1} \left\| \frac{1}{p_{P_k}} \sum_{i \in B_{P_k}} (y_i \nabla f_i(\mathbf{x}) - \hat{y}_i \nabla f_i(\hat{\mathbf{x}}_{P_k})) \right\|_{\mathcal{X}^*}^2 \\ & = \sum_{I=1}^N \frac{1}{p_I} \left\| \sum_{i \in B_I} (y_i \nabla f_i(\mathbf{x}) - \hat{y}_i \nabla f_i(\hat{\mathbf{x}}_I)) \right\|_{\mathcal{X}^*}^2 \end{aligned}$$

Applying Young's inequality, the quantity above is upper bounded via

$$\begin{aligned}
& \sum_{I=1}^N \frac{2}{p_I} \left[\left\| \sum_{i \in B_I} y_i (\nabla f_i(\mathbf{x}) - \nabla f_i(\hat{\mathbf{x}}_I)) \right\|_{\mathcal{X}^*}^2 + \left(\sum_{i \in B_I} |y_i - \hat{y}_i| \|\nabla f_i(\hat{\mathbf{x}}_I)\|_{\mathcal{X}^*} \right)^2 \right] \\
& \stackrel{(\circ)}{\leq} \sum_{I=1}^N \frac{2}{p_I} \left[\mathbf{L}_I^2 \|\mathbf{x} - \hat{\mathbf{x}}_I\|_{\mathcal{X}}^2 + \mathbf{G}_I^2 \|\mathbf{y}_I - \hat{\mathbf{y}}_I\|_2^2 \right] \\
& \leq 2 \left(\max_J \frac{\mathbf{L}_J^2}{p_J b_J} \right) \sum_{I=1}^N b_I \|\mathbf{x} - \hat{\mathbf{x}}_I\|_{\mathcal{X}}^2 + 2 \left(\max_J \frac{\mathbf{G}_J^2}{p_J c_J} \right) \sum_{I=1}^N c_I \|\mathbf{y}_I - \hat{\mathbf{y}}_I\|_2^2.
\end{aligned}$$

where we used Assumption 3.2.2 in (◦). Take the marginal expectation to complete the proof. \square

3.4.1 Strategy 1: Non-Uniform Historical Regularization

Here, the goal will be to select the balancing sequences $(w_k^P)_{k \geq 1}$ and $(w_k^D)_{k \geq 1}$ along with the weights $\boldsymbol{\gamma} = (\gamma_I)_{I=1}^N$ from (3.11) and (3.12) to achieve the desired complexity guarantee. This is the subject of Proposition 3.4.1. The rate conditions will be stated in terms of three quantities that largely depend on the sampling schemes \mathbf{p} and \mathbf{q} (which are used to define $\bar{\mathbf{g}}_{k-1}$ and $\bar{\mathbf{f}}_{k-1/2}$) and the primal regularization weights $\boldsymbol{\gamma}$. Those are

$$\mathbf{G}_{\mathbf{p}} := \sqrt{\max_I \frac{\mathbf{G}_I^2}{p_I}}, \mathbf{L}_{\mathbf{p}, \boldsymbol{\gamma}} := \sqrt{\max_I \frac{\mathbf{L}_I^2}{p_I \gamma_I}}, \text{ and } \mathbf{G}_{\mathbf{q}, \boldsymbol{\gamma}} := \sqrt{\max_I \frac{\mathbf{G}_I^2}{q_I \gamma_I}}. \quad (3.44)$$

Recall that the vectors $\mathbf{r} = (r_1, \dots, r_N)$ and $\mathbf{s} = (s_1, \dots, s_N)$ contain the probabilities by which the primal and dual table blocks are updated at each iteration. We will set these probabilities to the uniform vectors $\mathbf{r} = \mathbf{1}/N$ and $\mathbf{s} = \mathbf{1}/N$ as they will not affect the convergence rates in this analysis. Moreover, we assume that the objective (3.2) is dual separable (but not necessarily the feasible set). We discuss how this assumption can be avoided by a minor modification of the algorithm after the proof.

Proposition 3.4.1. *Let $(\mathbf{x}_0, \mathbf{y}_0) \in \text{ri}(\text{dom}(\phi)) \times \text{ri}(\text{dom}(\psi))$ and $\{(\mathbf{x}_k, \mathbf{y}_k)\}_{k \geq 1}$ be generated by the update from Lemma 3.3.1 with non-increasing sequences of weights $w_k^P \in [0, 1/2]$ and Lemma 3.3.2 with non-increasing sequence $w_k^D \in [0, 1/2]$, with $\bar{\mathbf{g}}_{k-1}$ and $\bar{\mathbf{f}}_{k-1/2}$ given*

by (3.39) and (3.40), respectively. Define $a_1 = \min \left\{ \frac{\sqrt{w_0^D \mu_0 \nu_0}}{4\mathbf{G}_P}, \frac{\sqrt{w_0^P \mu_0}}{4\sqrt{2}\mathbf{L}_{P,\gamma}}, \frac{\sqrt{w_0^P \mu_0 \nu_0}}{4\mathbf{G}_{q,\gamma}} \right\}$ and select $(a_k)_{k \geq 2}$ such that the conditions

$$a_k \leq \left\{ \frac{\sqrt{w_{k-1}^D (A_k \mu + \mu_0) (A_{k-1} \nu + \nu_0)}}{4\mathbf{G}_P}, \frac{\sqrt{w_{k-1}^P (A_k \mu + \mu_0) (A_{k-1} \mu + \mu_0)}}{4\sqrt{2}\mathbf{L}_{P,\gamma}}, \frac{\sqrt{w_{k-1}^P (A_k \nu + \nu_0) (A_{k-1} \mu + \mu_0)}}{4\mathbf{G}_{q,\gamma}} \right\} \quad (3.45)$$

are satisfied. In addition, impose that for any $\ell = 1, \dots, t-1$, it holds that

$$\frac{\mu}{N} \sum_{k'=0}^{\infty} w_{\ell+k'+1}^P A_{\ell+k'} (1 - 1/N)^{k'} \leq (w_{\ell}^P A_{\ell} + a_{\ell}) \mu \quad (3.46)$$

$$\frac{\nu}{N} \sum_{k'=0}^{\infty} w_{\ell+k'+1}^D A_{\ell+k'} (1 - 1/N)^{k'} \leq (w_{\ell}^D A_{\ell} + a_{\ell}) \nu \quad (3.47)$$

and that $A_k \leq \left(1 + \frac{1}{2N}\right)^k a_1$ when $\mu > 0$ or $\nu > 0$. We have that for any $\mathbf{u} \in \mathcal{X}$ and $\mathbf{v} \in \mathcal{Y}$,

$$\sum_{k=1}^t a_k \mathbb{E}[\text{Gap}^{\mathbf{u}, \mathbf{v}}(\mathbf{x}_k, \mathbf{y}_k)] + \frac{1}{4} \mathbb{E}[\mathcal{T}_t^P] + \frac{1}{4} \mathbb{E}[\mathcal{T}_t^D] \leq (a_1 \mu + \mu_0) \Delta_{\mathcal{X}}(\mathbf{u}, \mathbf{x}_0) + (a_1 \nu + \nu_0) \Delta_{\mathcal{Y}}(\mathbf{v}, \mathbf{y}_0).$$

Proof. Our starting point is Claim 3.3.1, after which we must show that the terms in (3.24), (3.25), and (3.26) are bounded by a quantity that does not grow with t . Recall that for (3.24), we used the three-term decomposition (3.41), which generated two telescoping terms and one error term. The first part of the proof uses this argument and bounds the error term.

1. Controlling (3.24): We follow the arguments at the beginning of this section to produce \mathcal{E}_k^P in (3.42) and \mathcal{E}_k^D in (3.43). We first upper bound the error terms $\mathbb{E}[\mathcal{E}_k^P]$ for $k = 2, \dots, t-1$ (noting that $\mathcal{E}_1^P = \mathcal{E}_1^D = 0$) and the last term of the telescoping inner products in (3.41). By Young's inequality with parameter $(1 - w_k^P)(A_{k-1}\mu + \mu_0)/4$, we have

$$\begin{aligned} \mathbb{E}[\mathcal{E}_k^P] &= a_{k-1} \mathbb{E} \left\langle \frac{1}{p_{P_k}} \sum_{i \in B_{P_k}} (y_{k-1,i} \nabla f_i(\mathbf{x}_{k-1}) - \hat{y}_{k-2,i} \hat{\mathbf{g}}_{k-2,i}), \mathbf{x}_k - \mathbf{x}_{k-1} \right\rangle \\ &\leq \frac{1 - w_k^P}{4} \mathbb{E}[\mathcal{C}_k^P] + \frac{2a_{k-1}^2}{(1 - w_k^P)(A_{k-1}\mu + \mu_0)} \mathbb{E} \left\| \frac{1}{p_{P_k}} \sum_{i \in B_{P_k}} (y_{k-1,i} \nabla f_i(\mathbf{x}_{k-1}) - \hat{y}_{k-2,i} \hat{\mathbf{g}}_{k-2,i}) \right\|_{\mathcal{X}^*}^2 \end{aligned} \quad (3.48)$$

and for the last term of the decomposition (3.41),

$$\begin{aligned}
& -\mathbb{E}[a_t \langle \nabla f(\mathbf{x}_t)^\top \mathbf{y}_t - \hat{\mathbf{g}}_{t-1}^\top \hat{\mathbf{y}}_{t-1}, \mathbf{u} - \mathbf{x}_t \rangle] \\
& \leq \frac{1 - w_t^P}{4} \mathbb{E}[\mathcal{T}_t^P] + \frac{2a_t^2}{(1 - w_t^P)(A_t\mu + \mu_0)} \mathbb{E} \left\| \frac{1}{p_{P_t}} \sum_{i \in B_{P_t}} (y_{t,i} \nabla f_i(\mathbf{x}_t) - \hat{y}_{t-1,i} \hat{\mathbf{g}}_{t-1,i}) \right\|_{\mathcal{X}^*}^2 \quad (3.49)
\end{aligned}$$

To handle the second term in (3.48) and (3.49) for $k \in \{2, \dots, t\}$, apply Lemma 3.4.1 with $b_I = \gamma_I$ and $c_I = 1$ to achieve

$$\begin{aligned}
& \mathbb{E} \left\| \frac{1}{p_{P_k}} \sum_{i \in B_{P_k}} (y_{k-1,i} \nabla f_i(\mathbf{x}_{k-1}) - \hat{y}_{k-2,i} \hat{\mathbf{g}}_{k-2,i}) \right\|_{\mathcal{X}^*}^2 \\
& \leq 2 \left(\max_I \frac{\mathbf{L}_J^2}{p_J \gamma_J} \right) \cdot \sum_{I=1}^N \gamma_I \mathbb{E} \|\mathbf{x}_{k-1} - \hat{\mathbf{x}}_{k-2,I}\|_{\mathcal{X}}^2 + 2 \left(\max_J \frac{\mathbf{G}_J^2}{p_J} \right) \mathbb{E} \|\mathbf{y}_{k-1} - \hat{\mathbf{y}}_{k-2}\|_2^2 \\
& \leq 4 \underbrace{\left(\max_J \frac{\mathbf{L}_J^2}{p_J \gamma_J} \right)}_{\mathbf{L}_{p,\gamma}^2} \cdot \sum_{I=1}^N \gamma_I \mathbb{E}[\Delta_{\mathcal{X}}(\mathbf{x}_{k-1}, \hat{\mathbf{x}}_{k-2,I})] + 4 \underbrace{\left(\max_J \frac{\mathbf{G}_J^2}{p_J} \right)}_{\mathbf{G}_p^2} \mathbb{E}[\Delta_{\mathcal{Y}}(\mathbf{y}_{k-1}, \hat{\mathbf{y}}_{k-2,I})],
\end{aligned}$$

where in the last line we applied $\|\cdot\|_2 \leq \|\cdot\|_{\mathcal{Y}}$ and the strong convexity of Bregman divergences. Recall that $\hat{\mathcal{C}}_{k,I}^P := (A_{k-1}\mu + \mu_0)\Delta_{\mathcal{X}}(\mathbf{x}_k, \hat{\mathbf{x}}_{k-1,I})$ and $\hat{\mathcal{C}}_k^D := (A_{k-1}\nu + \nu_0)\Delta_{\mathcal{Y}}(\mathbf{y}_k, \hat{\mathbf{y}}_{k-1})$. Combining the steps above, using that $1/(1 - w_k^P) \leq 2$, and applying the condition (3.45), we have the upper bound

$$\begin{aligned}
\mathbb{E}[\mathcal{E}_k^P] & \leq \frac{1 - w_k^P}{4} \mathbb{E}[\mathcal{C}_k^P] + \frac{8a_{k-1}^2 \mathbf{G}_p^2 \mathbb{E}[\hat{\mathcal{C}}_{k-1}^D]}{(1 - w_k^P)(A_{k-1}\mu + \mu_0)(A_{k-2}\nu + \nu_0)} \\
& \quad + \frac{8a_{k-1}^2 \mathbf{L}_{p,\gamma}^2 \mathbb{E} \left[\sum_{I=1}^N \gamma_I \hat{\mathcal{C}}_{k-1,I}^P \right]}{(1 - w_k^P)(A_{k-1}\mu + \mu_0)(A_{k-2}\mu + \mu_0)} \\
& \leq \frac{1 - w_{k-1}^P}{2} \mathbb{E}[\mathcal{C}_k^P] + \frac{w_{k-1}^D}{2} \mathbb{E}[\hat{\mathcal{C}}_{k-1}^D] + \frac{w_{k-1}^P}{4} \mathbb{E} \left[\sum_{I=1}^N \gamma_I \hat{\mathcal{C}}_{k-1,I}^P \right]
\end{aligned}$$

with a similar bound holding for (3.49). These terms will cancel with the corresponding non-positive terms in (3.24).

The upper bounds for $\mathbb{E}[\mathcal{E}_k^D]$ and $-\mathbb{E}[a_k \langle \nabla f(\mathbf{x}_k) - \hat{\mathbf{f}}_k, \mathbf{v} - \mathbf{y}_k \rangle]$ follow by very similar arguments as above. Applying Young's inequality with parameter $(1 - w_k^D)(A_{k-1}\nu + \nu_0)/4$

we upper bound $\mathbb{E}_{k-1/2}[\mathcal{E}_k^D]$ via

$$\begin{aligned} & a_{k-1} \mathbb{E}_{k-1/2} \left\langle \frac{1}{q_{Q_k}} \sum_{j \in B_{Q_k}} (f_j(\mathbf{x}_{k-1}) - \hat{f}_{k-2,j}) \mathbf{e}_j, \mathbf{y}_k - \mathbf{y}_{k-1} \right\rangle \\ & \leq \frac{1 - w_k^D}{4} \mathbb{E}_{k-1/2}[\mathcal{C}_k^D] + \frac{2a_{k-1}^2(1 - w_k^D)^{-1}}{(A_{k-1}\nu + \nu_0)} \sum_{J=1}^N \frac{1}{q_J} \left\| \sum_{j \in B_J} (f_j(\mathbf{x}_{k-1}) - \hat{f}_{k-1,j}) \mathbf{e}_j \right\|_{\mathbf{y}_*}^2 \end{aligned} \quad (3.50)$$

Recall that the index R_{k-1} determines which block of the primal table is updated. Thus, it holds in conditional expectation that

$$\begin{aligned} \mathbb{E}_{k-2} [\|\mathbf{x}_{k-1} - \hat{\mathbf{x}}_{k-1,J}\|_{\mathcal{X}}^2] &= \mathbb{E}_{k-2} [\|\mathbf{x}_{k-1} - \hat{\mathbf{x}}_{k-1,J}\|_{\mathcal{X}}^2 \mathbb{1}_{J=R_{k-1}}] + \mathbb{E}_{k-2} [\|\mathbf{x}_{k-1} - \hat{\mathbf{x}}_{k-1,J}\|_{\mathcal{X}}^2 \mathbb{1}_{J \neq R_{k-1}}] \\ &= 0 + \mathbb{E}_{k-2} [\|\mathbf{x}_{k-1} - \hat{\mathbf{x}}_{k-2,J}\|_{\mathcal{X}}^2 \mathbb{1}_{J \neq R_{k-1}}] \\ &\leq \mathbb{E}_{k-2} [\|\mathbf{x}_{k-1} - \hat{\mathbf{x}}_{k-2,J}\|_{\mathcal{X}}^2]. \end{aligned}$$

Taking the marginal expectation, the second term of (3.50) can be upper bounded as

$$\begin{aligned} \sum_{J=1}^N \frac{1}{q_J} \mathbb{E} \left\| \sum_{j \in B_J} (f_j(\mathbf{x}_{k-1}) - \hat{f}_{k-1,j}) \mathbf{e}_j \right\|_{\mathbf{y}_*}^2 &\leq \sum_{J=1}^N \frac{G_J^2}{q_J} \mathbb{E} \|\mathbf{x}_{k-1} - \hat{\mathbf{x}}_{k-1,J}\|_{\mathcal{X}}^2 \\ &\leq \sum_{J=1}^N \frac{G_J^2}{q_J} \mathbb{E} \|\mathbf{x}_{k-1} - \hat{\mathbf{x}}_{k-2,J}\|_{\mathcal{X}}^2 \\ &\leq \underbrace{2 \max_I \frac{G_I^2}{q_I \gamma_I}}_{G_{q,\gamma}^2} \sum_{J=1}^N \gamma_J \Delta_{\mathcal{X}}(\mathbf{x}_{k-1}, \hat{\mathbf{x}}_{k-2,J})_{\mathcal{X}}^2. \end{aligned} \quad (3.51)$$

Invoking condition (3.45) once again and taking the marginal expectation, we have

$$\mathbb{E}[\mathcal{E}_k^D] \leq \frac{1 - w_k^D}{4} \mathbb{E}[\mathcal{C}_k^D] + \frac{4a_{k-1}^2 G_{q,\gamma}^2}{(1 - w_k^D)(A_{k-1}\nu + \nu_0)(A_{k-2}\mu + \mu_0)} \mathbb{E} \left[\sum_{I=1}^N \gamma_I \hat{\mathcal{C}}_{k-1,I}^P \right] \quad (3.52)$$

$$\leq \frac{1 - w_k^D}{4} \mathbb{E}[\mathcal{C}_k^D] + \frac{w_{k-1}^P}{4} \mathbb{E} \left[\sum_{I=1}^N \gamma_I \hat{\mathcal{C}}_{k-1,I}^P \right]. \quad (3.53)$$

Similarly, $-\mathbb{E}[a_k \langle \nabla f(\mathbf{x}_k) - \hat{\mathbf{f}}_k, \mathbf{v} - \mathbf{y}_k \rangle] \leq \frac{1 - w_t^D}{4} \mathbb{E}[\mathcal{T}_t^D] + \frac{w_{t-1}^P}{4} \mathbb{E} \left[\sum_{I=1}^N \gamma_I \hat{\mathcal{C}}_{t,I}^P \right]$. All these terms will cancel with corresponding terms in (3.24).

2. Controlling (3.25): For this step, we will express the $\hat{\mathcal{T}}_k^{\text{P}}$ terms as a function of the \mathcal{T}_k^{P} terms by analyzing the random sampling that governs the block updates. For any $k \geq 1$, write

$$\begin{aligned} \sum_{I=1}^N \gamma_I \mathbb{E}[\Delta_{\mathcal{X}}(\mathbf{u}, \hat{\mathbf{x}}_{k,I})] &= \sum_{I=1}^N \gamma_I ((1/N) \mathbb{E}[\Delta_{\mathcal{X}}(\mathbf{u}, \mathbf{x}_k)] + (1 - 1/N) \mathbb{E}[\Delta_{\mathcal{X}}(\mathbf{u}, \hat{\mathbf{x}}_{k-1,I})]) \\ &= (1/N) \sum_{k'=0}^k (1 - 1/N)^{k'} \mathbb{E}[\Delta_{\mathcal{X}}(\mathbf{u}, \mathbf{x}_{k-k'})] \end{aligned} \quad (3.54)$$

Using (3.54), the expression (3.25) can be expanded to

$$\begin{aligned} &\sum_{k=0}^{t-1} w_{k+1}^{\text{P}} (A_k \mu + \mu_0) (1/N) \sum_{k'=0}^k (1 - 1/N)^{k'} \mathbb{E}[\Delta_{\mathcal{X}}(\mathbf{u}, \mathbf{x}_{k-k'})] \\ &= \sum_{k'=0}^{t-1} \sum_{k=k'}^{t-1} w_{k+1}^{\text{P}} (A_k \mu + \mu_0) (1/N) (1 - 1/N)^{k'} \mathbb{E}[\Delta_{\mathcal{X}}(\mathbf{u}, \mathbf{x}_{k-k'})] \quad \text{exchange sums} \\ &= \sum_{k'=0}^{t-1} \sum_{\ell=0}^{t-1-k'} w_{\ell+k'+1}^{\text{P}} (A_{\ell+k'} \mu + \mu_0) \mathbb{E}[\Delta_{\mathcal{X}}(\mathbf{u}, \mathbf{x}_{\ell})] (1/N) (1 - 1/N)^{k'}. \quad \text{reparameterize } \ell = k - k' \end{aligned}$$

Using the identity above and $t - 1 - k' < t - 1 < \infty$, we can decompose (3.25) as

$$\begin{aligned} &\frac{1}{2} \sum_{k=1}^t \mathbb{E} \left[w_k^{\text{P}} \left(\sum_{I=1}^N \gamma_I \hat{\mathcal{T}}_{k-1,I}^{\text{P}} - \mathcal{T}_k^{\text{P}} \right) - a_k \mu \Delta_{\mathcal{X}}(\mathbf{u}, \mathbf{x}_k) \right] \\ &\leq \frac{1}{2} \sum_{\ell=0}^{t-1} \underbrace{\left(\sum_{k'=0}^{\infty} w_{\ell+k'+1}^{\text{P}} (A_{\ell+k'} \mu + \mu_0) (1/N) (1 - 1/N)^{k'} \right)}_{\leq w_{\ell}^{\text{P}} (A_{\ell} \mu + \mu_0) + a_{\ell} \mu \text{ for } \ell \geq 1} \mathbb{E}[\Delta_{\mathcal{X}}(\mathbf{u}, \mathbf{x}_{\ell})] \\ &\quad - \frac{1}{2} \sum_{\ell=0}^{t-1} (w_{\ell+1}^{\text{P}} (A_{\ell+1} \mu + \mu_0) + a_{\ell+1} \mu) \mathbb{E}[\Delta_{\mathcal{X}}(\mathbf{u}, \mathbf{x}_{\ell+1})] \end{aligned}$$

where the inequality under the braces follows from the theorem assumptions. By telescoping the resulting terms, (3.25) is upper bounded by

$$\frac{1}{2} \left(\sum_{k'=0}^{\infty} w_{k'+1}^{\text{P}} (A_{k'} \mu + \mu_0) (1/N) (1 - 1/N)^{k'} \right) \Delta_{\mathcal{X}}(\mathbf{u}, \mathbf{x}_0).$$

Finally, we upper bound the leading constant on $\Delta_{\mathcal{X}}(\mathbf{u}, \mathbf{x}_0)$ above. We separate the $(A_{k'}\mu + \mu_0)$ constants into two terms each and bound the resulting sums. Because $w_{k'+1}^{\mathbf{P}} \leq 1/2$, we have that $\frac{\mu_0}{2N} (\sum_{k'=0}^{\infty} w_{k'+1}^{\mathbf{P}} (1 - 1/N)^{k'}) \leq \frac{\mu_0}{4}$ and because $A_{k'} \leq (1 + 1/(2N))^{k'} a_1$, we have

$$\frac{\mu}{2N} \left(\sum_{k'=0}^{\infty} w_{k'+1}^{\mathbf{P}} A_{k'} (1 - 1/N)^{k'} \right) \leq \frac{a_1 \mu}{4N} \sum_{k'=0}^{\infty} (1 - 1/(2N))^{k'} = \frac{a_1 \mu}{2}.$$

Summing the two gives the leading constant of $\Delta_{\mathcal{X}}(\mathbf{u}, \mathbf{x}_0)$ appearing in the statement, as $(a_1 \mu / 2 + \mu_0 / 4 + \mu_0 / 2) \leq a_1 \mu + \mu_0$.

3. Controlling (3.26): This will follow from similar steps as those shown above, but will rely upon using the probabilistic arguments on each coordinate block of $\hat{\mathbf{y}}_k$. Recall the notation $\Delta_I(\cdot, \cdot)$ from Section 3.2. Using dual-separability of the objective, write

$$\begin{aligned} \mathbb{E}[\Delta_{\mathcal{Y}}(\mathbf{v}, \hat{\mathbf{y}}_k)] &= \sum_{I=1}^N \mathbb{E}[\Delta_I(\mathbf{v}_I, \hat{\mathbf{y}}_{k,I})] \\ &= \sum_{I=1}^N ((1/N) \mathbb{E}[\Delta_I(\mathbf{v}_I, \mathbf{y}_{k,I})] + (1 - 1/N) \mathbb{E}[\Delta_I(\mathbf{v}_I, \hat{\mathbf{y}}_{k-1,I})]) \\ &= (1/N) \sum_{k'=0}^k (1 - 1/N)^{k'} \underbrace{\sum_{I=1}^N \mathbb{E}[\Delta_I(\mathbf{v}_I, \mathbf{y}_{k-k',I})]}_{\mathbb{E}[\Delta_{\mathcal{Y}}(\mathbf{v}, \mathbf{y}_{k-k'})]}. \end{aligned}$$

The remainder of the argument follows identically to Step 2 and produces the leading constant of $\Delta_{\mathcal{Y}}(\mathbf{v}, \mathbf{y}_0)$ appearing in the statement. \square

Note that the separability of $\Delta_{\mathcal{Y}}(\cdot, \cdot)$ is only used in Step 3, which could be eschewed by updating the entirety of $\hat{\mathbf{y}}_{k-1}$ with probability $1/N$ as opposed to the block-wise updates currently being used. Recall the initial distance D_0 from (3.36).

Theorem 3.4.1. *Under Assumption 3.2.2 and Assumption 3.2.1, consider any $\mathbf{u} \in \mathcal{X}$, $\mathbf{v} \in \mathcal{Y}$ and precision $\varepsilon > 0$. There exists a choice of the sequence $(a_k)_{k=1}^t$, and the parameters $w_k^{\mathbf{P}}$ and $w_k^{\mathbf{D}}$ such that Algorithm 4 with Identity Card 2 produces an output point $(\tilde{\mathbf{x}}_t, \tilde{\mathbf{y}}_t) \in \mathcal{X} \times \mathcal{Y}$ satisfying $\mathbb{E}[\text{Gap}^{\mathbf{u}, \mathbf{v}}(\tilde{\mathbf{x}}_t, \tilde{\mathbf{y}}_t)] \leq \varepsilon$ for t that depends on ε according to*

the following iteration complexities. They depend logarithmically on the constants $C_1 := \sqrt{N}(\mathbf{L}_{\mathbf{p},\gamma}\sqrt{\nu_0/\mu_0} + (\mathbf{G}_{\mathbf{p}} \vee \mathbf{G}_{\mathbf{q},\gamma}))D_0 + \mu\Delta_{\mathcal{X}}(\mathbf{u}, \mathbf{x}_0) + \nu\Delta_{\mathcal{Y}}(\mathbf{v}, \mathbf{y}_0)$, $C_2 := (\sqrt{N}\mathbf{L}_{\mathbf{p},\gamma}\sqrt{\nu_0/\mu_0} + (\mathbf{G}_{\mathbf{p}} \vee N\mathbf{G}_{\mathbf{q},\gamma}))D_0 + \mu\Delta_{\mathcal{X}}(\mathbf{u}, \mathbf{x}_0)$, and $C_3 := (\mathbf{L}_{\mathbf{p},\gamma}\sqrt{\nu_0/\mu_0} + (N\mathbf{G}_{\mathbf{p}} \vee \mathbf{G}_{\mathbf{q},\gamma}))D_0 + \nu\Delta_{\mathcal{Y}}(\mathbf{v}, \mathbf{y}_0)$.

Case	Iteration Complexity
$\mu > 0$ and $\nu > 0$	$O\left(\left(N + \frac{\sqrt{N}\mathbf{L}_{\mathbf{p},\gamma}}{\mu} + \frac{\sqrt{N}(\mathbf{G}_{\mathbf{p}} \vee \mathbf{G}_{\mathbf{q},\gamma})}{\sqrt{\mu\nu}}\right) \ln\left(\frac{C_1}{\varepsilon}\right)\right)$
$\mu > 0$ and $\nu = 0$	$O\left(\left(N + \frac{\sqrt{N}\mathbf{L}_{\mathbf{p},\gamma}}{\mu}\right) \ln\left(\frac{C_2}{\varepsilon}\right) + (\mathbf{G}_{\mathbf{p}} \vee N\mathbf{G}_{\mathbf{q},\gamma}) \sqrt{\frac{\sqrt{\mu_0/\nu_0}D_0 + (a_1\mu/\nu_0)\Delta_{\mathcal{X}}(\mathbf{u}, \mathbf{x}_0)}{\mu\varepsilon}}\right)$
$\mu = 0$ and $\nu > 0$	$O\left(N \ln\left(\frac{C_3}{\varepsilon}\right) + \frac{\mathbf{L}_{\mathbf{p},\gamma}\sqrt{\nu_0/\mu_0}D_0}{\varepsilon} + (N\mathbf{G}_{\mathbf{p}} \vee \mathbf{G}_{\mathbf{q},\gamma}) \sqrt{\frac{\sqrt{\nu_0/\mu_0}D_0 + (a_1\nu/\mu_0)\Delta_{\mathcal{Y}}(\mathbf{v}, \mathbf{y}_0)}{\nu\varepsilon}}\right)$
$\mu = 0$ and $\nu = 0$	$O\left(\frac{(\mathbf{L}_{\mathbf{p},\gamma}\sqrt{\nu_0/\mu_0} + (\mathbf{G}_{\mathbf{p}} \vee \mathbf{G}_{\mathbf{q},\gamma}))D_0}{\varepsilon}\right)$

Proof. We split the proof into the same case-by-case strategy as employed in Theorem 3.3.1. While we may match those arguments exactly for most of the conditions of Proposition 3.4.1, we need to additionally set the correct values of the sequences $(w_k^{\mathbf{P}})_{k \geq 1}$ and $(w_k^{\mathbf{D}})_{k \geq 1}$ to complete the analysis. In all cases, the requirement that $A_k \leq (1 + \frac{1}{2N})^k a_1$ introduces a term of the form $N \ln\left(\frac{G_0}{a_1\varepsilon}\right)$ to the iteration complexity, where $G_0 = (a_1\mu + \mu_0)\Delta_{\mathcal{X}}(\mathbf{u}, \mathbf{x}_0) + (a_1\nu + \nu_0)\Delta_{\mathcal{Y}}(\mathbf{v}, \mathbf{y}_0)$. When substituting the particular values of $(w_0^{\mathbf{P}}, w_0^{\mathbf{D}})$ in each case, we derive the constants (C_1, C_2, C_3) .

Case 1: $\mu > 0, \nu > 0$. We consider here a constant choice of the sequences and $w_k^{\mathbf{P}} = w_0^{\mathbf{P}}$ and $w_k^{\mathbf{D}} = w_0^{\mathbf{D}}$. Then, all conditions on the growth of $(a_k)_{k \geq 1}$ can be satisfied using $a_k = \alpha A_{k-1}$, where

$$\alpha \lesssim \min \left\{ \frac{\sqrt{w_0^{\mathbf{D}}\mu\nu}}{\mathbf{G}_{\mathbf{p}}}, \frac{\sqrt{w_0^{\mathbf{P}}\mu\nu}}{\mathbf{G}_{\mathbf{q},\gamma}}, \frac{\sqrt{w_0^{\mathbf{P}}\mu}}{\mathbf{L}_{\mathbf{p},\gamma}} \right\}.$$

This implies that $A_k = (1 + \alpha)A_{k-1}$. Then, considering (3.46), we have that

$$\frac{1}{N} \sum_{k'=0}^{\infty} w_0^P A_{\ell+k'} (1 - 1/N)^{k'} = \frac{1}{N} w_0^P A_{\ell} \sum_{k'=0}^{\infty} (1 + \alpha)^{k'} (1 - 1/N)^{k'}.$$

Consider a setting of α and a constant $c > 0$ (which may depend on N) such that

$$\frac{1}{N} \sum_{k'=0}^{\infty} (1 + \alpha)^{k'} (1 - 1/N)^{k'} \leq 1 + \frac{c\alpha}{1 + \alpha}, \quad (3.55)$$

where the right-hand side can be tightened to $1 + c\alpha/2$ when $\alpha \leq 1$. Then, by setting $w_0^P \leq 2c^{-1}$, the condition (3.46) is satisfied (with an identical argument holding for (3.47)). When $N = 1$, the term above vanishes, so we may consider $N \geq 2$. To satisfy (3.55), we need that $\alpha \leq \frac{N-c/2-1}{1-N} = \frac{1}{N-1}$ for $c = 2N$, which imposes the condition that $\alpha \leq \frac{1}{N-1}$.

Case 2: $\mu > 0, \nu = 0$. We set $w_k^P = 1/2$ for all $k \geq 0$ and need only set w_k^P , which will be piece-wise constant. For all $\ell \leq k^*$ such that the second condition of (3.45) is the dominant condition, that is, $a_{\ell} = \alpha_{\ell} A_{\ell-1}$ for $\alpha_{\ell} \lesssim \frac{\sqrt{w_{\ell-1}^P} \mu}{L_{p,\gamma}}$, we may set $w_{\ell}^P = 1/N$ and $\alpha_{\ell} \equiv \alpha_{k^*}$ for all $\ell \in \{0, \dots, k^*\}$ and follow the argument of Case 1 to achieve the first term in the complexity. For $\ell \geq k^* + 1$, we will derive the value of $w_{k^*+1}^P$. We have that $a_{\ell} = c\mu\nu_0 \min \left\{ \frac{1}{G_p^2}, \frac{w_{k^*+1}^P}{G_{q,\gamma}^2} \right\} \ell$ for an absolute constant $c > 0$, and moreover, that

$$A_{\ell} \leq c\mu\nu_0 \min \left\{ \frac{1}{G_p^2}, \frac{w_{k^*+1}^P}{G_{q,\gamma}^2} \right\} \frac{\ell(\ell+1)}{2}.$$

Furthermore, the condition (3.46) is satisfied if

$$\frac{1}{N} w_{k^*+1}^P \sum_{k'=0}^{\infty} \frac{(\ell+k')(\ell+1+k')}{2} (1 - 1/N)^{k'} \leq w_{k^*+1}^P \frac{\ell(\ell+1)}{2} + \ell$$

which in turn is satisfied when

$$w_{k^*+1}^P \underbrace{\frac{1}{N} \sum_{k'=0}^{\infty} \frac{k'(2\ell+k'+1)}{2} (1 - 1/N)^{k'}}_{\leq N\ell + N^2 + 1/2} \leq \ell,$$

where the upper bound follows by summing over k' and applying $(1 - 1/N) \leq 1$. Set

$$w_{k^*+1}^P = \frac{1}{N + N^2/k^* + 1/(2k^*)},$$

which introduces the term $\frac{N\mathbf{G}_{q,\gamma}}{\sqrt{\mu\varepsilon}} \cdot \sqrt{\frac{G_0}{\nu_0}}$ in the given complexity.

Case 3: $\mu = 0, \nu > 0$. Here, we may set $w_k^P = 1/2$ for all $k \geq 0$ and derive the required setting for $(w_k^D)_{k \geq 1}$. We may repeat the argument above and for Case 3 of Theorem 3.3.1 to set $w_k^D \sim 1/N^2$, which achieves the given complexity.

Case 4: $\mu = 0, \nu = 0$. The conditions (3.46) and (3.47) vanish, so we reuse the sequence $A_t = O(\min\{\sqrt{\mu_0\nu_0}/G, \mu_0/L\}t)$ as before to complete the proof. \square

We instantiate the problem constants by selecting a sampling scheme. Recall from Section 3.2 that $\boldsymbol{\lambda} = (\boldsymbol{\lambda}_1, \dots, \boldsymbol{\lambda}_N)$ for $\boldsymbol{\lambda}_I := \sqrt{\mathbf{G}_I^2 + \mathbf{L}_I^2}$ along with the constants from (3.44). The non-uniform sampling complexity given below follows by letting $p_I \propto \boldsymbol{\lambda}_I$, $\gamma_I \propto \boldsymbol{\lambda}_I$, and $q_I \propto \mathbf{G}_I$. The constants appearing in Theorem 3.4.1 are tabulated below.

Constant	Uniform Sampling	Non-Uniform Sampling
$\mathbf{G}_p \vee \mathbf{G}_{q,r}$	$N \ \mathbf{G}\ _\infty$	$\ \boldsymbol{\lambda}\ _1^{1/2} \ \mathbf{G}\ _1^{1/2}$
$\mathbf{L}_{p,\gamma}$	$N \ \mathbf{L}\ _\infty$	$\ \boldsymbol{\lambda}\ _1$

We discuss memory and per-iteration complexity of the method. If we show that the optimization problem (3.11) can be solved at $\tilde{O}(d)$ cost, then the total arithmetic complexity is given by $\tilde{O}(n(d+N)/N)$, where we recall that we assume uniform block sizes for arithmetic complexity discussions. The relevant terms are the matrix-vector product $\hat{\mathbf{g}}_k^\top \hat{\mathbf{y}}_k$ and the sum of weighted Bregman divergences.

For the former, let $\hat{\mathbf{g}}_{k,I} \in \mathbb{R}^{n/N \times d}$ denote the matrix containing the primal gradients for the elements in block I . Noting that R_k is the block of the primal table that is updated at

each iteration and S_k plays the same role for the dual table, it holds that

$$\begin{aligned}
\hat{\mathbf{g}}_k^\top \hat{\mathbf{y}}_k &= \hat{\mathbf{g}}_{k-1}^\top \hat{\mathbf{y}}_k + (\hat{\mathbf{g}}_{k,R_k} - \hat{\mathbf{g}}_{k-1,R_k})^\top \hat{\mathbf{y}}_{k,R_k} \\
&= \hat{\mathbf{g}}_{k-1}^\top \left(\hat{\mathbf{y}}_{k-1} + \sum_{j \in B_{S_k}} (\hat{\mathbf{y}}_{k,j} - \hat{\mathbf{y}}_{k-1,j}) \mathbf{e}_j \right) + (\hat{\mathbf{g}}_{k,R_k} - \hat{\mathbf{g}}_{k-1,R_k})^\top \hat{\mathbf{y}}_{k,R_k} \\
&= \hat{\mathbf{g}}_{k-1}^\top \hat{\mathbf{y}}_{k-1} + \hat{\mathbf{g}}_{k-1,S_k}^\top (\hat{\mathbf{y}}_{k,S_k} - \hat{\mathbf{y}}_{k-1,S_k}) + (\hat{\mathbf{g}}_{k,R_k} - \hat{\mathbf{g}}_{k-1,R_k})^\top \hat{\mathbf{y}}_{k,R_k}. \tag{3.56}
\end{aligned}$$

Assuming that $\hat{\mathbf{g}}_{k-1}^\top \hat{\mathbf{y}}_{k-1} \in \mathbb{R}^d$ is already stored, everything above can be computed with $\hat{\mathbf{g}}_{k,R_k}$, $\hat{\mathbf{y}}_{k,S_k}$, and past information at cost $O(nd/N)$. Furthermore, we need not retain the entire table $\hat{\mathbf{g}}_k \in \mathbb{R}^{n \times d}$, as $\hat{\mathbf{g}}_{k-1,R_k}$ and $\hat{\mathbf{g}}_{k-1,S_k}$ above can be recomputed from $\hat{\mathbf{x}}_{k-1,R_k}$ and $\hat{\mathbf{x}}_{k-1,S_k}$. The entire memory footprint is $O(Nd + n)$, which could be much smaller than $O(nd)$ (for instance, when $N = d$).

For the latter, letting φ be the generator of $\Delta_{\mathcal{X}}(\cdot, \cdot)$, we write

$$\sum_{I=1}^N \gamma_I \Delta_{\mathcal{X}}(\mathbf{x}, \hat{\mathbf{x}}_{k,I}) = \varphi(\mathbf{x}) + \left\langle \sum_{I=1}^N \gamma_I \nabla \varphi(\hat{\mathbf{x}}_{k,I}), \mathbf{x} \right\rangle + \text{const}(\mathbf{x}),$$

where the term $\text{const}(\mathbf{x})$ does not vary with respect to \mathbf{x} . It then holds that

$$\sum_{I=1}^N \gamma_I \nabla \varphi(\hat{\mathbf{x}}_{k,I}) = \sum_{I=1}^N \gamma_I \nabla \varphi(\hat{\mathbf{x}}_{k-1,I}) + \gamma_{R_k} (\nabla \varphi(\hat{\mathbf{x}}_{k,R_k}) - \nabla \varphi(\hat{\mathbf{x}}_{k-1,R_k})),$$

so we need only compute $\nabla \varphi(\hat{\mathbf{x}}_{k,R_k})$ at each iteration. Retaining the $(\nabla \varphi(\hat{\mathbf{x}}_{k,I}))_{I=1}^N$ comes at an $O(Nd)$ storage cost, which is the same as the table itself. The total per-iteration complexity is then $\tilde{O}(n(d + N)/N)$.

3.4.2 Strategy 2: Non-Uniform Block Replacement Probabilities

In the previous approach, we relied on the non-uniform weights $(\gamma_I)_{I=1}^N$ in order to achieve complexities that were independent of the number of blocks N . Here, rather than relying on the historical regularization, we will instead tune the sampling probabilities $\mathbf{r} = (r_1, \dots, r_N)$ and $\mathbf{s} = (s_1, \dots, s_N)$, which govern the element R_k of $\hat{\mathbf{x}}_{k-1,1}, \dots, \hat{\mathbf{x}}_{k-1,N}$ and which coordinate block S_k of $\hat{\mathbf{y}}_{k-1}$ gets updates at each iteration k (see (3.38)).

We bound the same terms as in Proposition 3.4.1, although (3.25) and (3.26) are immediately non-positive as we will set $w_k^P = 0$ and $w_k^D = 0$ for all k . This argument is employed in Proposition 3.4.2. Similarly to (3.44), the resulting complexity will depend on the sampling probabilities \mathbf{p} , \mathbf{q} , \mathbf{r} , and \mathbf{s} through the following constants:

$$\mathbf{G}_{\mathbf{p},\mathbf{s}} := \sqrt{\max_{I \in [N]} \frac{\mathbf{G}_I^2}{p_I s_I^2}}, \mathbf{L}_{\mathbf{p},\mathbf{r}} := \sqrt{\sum_{I=1}^N \frac{\mathbf{L}_I^2}{p_I r_I^2}}, \text{ and } \mathbf{G}_{\mathbf{q},\mathbf{r}} := \sqrt{\sum_{I=1}^N \frac{\mathbf{G}_I^2}{q_I r_I^2}}.$$

Observe the following.

Proposition 3.4.2. *Let $(\mathbf{x}_0, \mathbf{y}_0) \in \text{ri}(\text{dom}(\phi)) \times \text{ri}(\text{dom}(\psi))$ and $\{(\mathbf{x}_k, \mathbf{y}_k)\}_{k \geq 1}$ be generated using $\bar{\mathbf{g}}_{k-1}$ and $\bar{\mathbf{f}}_{k-1/2}$ given by (3.39) and (3.40), respectively. Define $a_1 = \frac{1}{15} \min \left\{ \frac{\sqrt{\mu_0 \nu_0}}{\mathbf{G}_{\mathbf{p},\mathbf{s}} \vee \mathbf{G}_{\mathbf{q},\mathbf{r}}}, \frac{\mu_0}{\mathbf{L}_{\mathbf{p},\mathbf{r}}} \right\}$ and select $(a_k)_{k \geq 2}$ such that both the conditions*

$$\frac{a_k^2}{A_k \mu + \mu_0} \leq \min_I (1 + (s_I \wedge r_I)/5) \frac{a_{k-1}^2}{A_{k-1} \mu + \mu_0} \quad (3.57)$$

$$\frac{a_k^2}{A_k \nu + \nu_0} \leq \min_I (1 + r_I/5) \frac{a_{k-1}^2}{A_{k-1} \nu + \nu_0} \quad (3.58)$$

and

$$a_k \leq \min \left\{ \frac{\sqrt{(A_k \mu + \mu_0)(A_{k-1} \nu + \nu_0)}}{10 \mathbf{G}_{\mathbf{p},\mathbf{s}}}, \frac{\sqrt{(A_k \mu + \mu_0)(A_{k-1} \mu + \mu_0)}}{10 \mathbf{L}_{\mathbf{p},\mathbf{r}}}, \frac{\sqrt{(A_k \nu + \nu_0)(A_{k-1} \mu + \mu_0)}}{15 \mathbf{G}_{\mathbf{q},\mathbf{r}}} \right\} \quad (3.59)$$

are satisfied. We have that for any $\mathbf{u} \in \mathcal{X}$ and $\mathbf{v} \in \mathcal{Y}$,

$$\sum_{k=1}^t a_k \mathbb{E}[\text{Gap}^{\mathbf{u},\mathbf{v}}(\mathbf{x}_k, \mathbf{y}_k)] + \frac{1}{2} \mathbb{E}[\mathcal{T}_t^P] + \frac{1}{2} \mathbb{E}[\mathcal{T}_t^D] \leq \mathcal{T}_0^P + \mathcal{T}_0^D.$$

Proof. As stated before, we aim to show that the sum of terms in (3.24) is upper bounded by a constant independent of t . This is composed of the terms $\mathbb{E}[\mathcal{I}_k^P]$ and $\mathbb{E}[\mathcal{I}_k^D]$. We divide this task into bounding the primal and dual components separately. As before, by applying the argument leading to the expressions (3.42) and (3.43), bounding the inner product terms $\mathbb{E}[\mathcal{I}_k^P]$ and $\mathbb{E}[\mathcal{I}_k^D]$ reduces to bounding the error terms $\mathbb{E}[\mathcal{E}_k^P]$ and $\mathbb{E}[\mathcal{E}_k^D]$ and the final element of the telescoping inner product terms $a_t \langle \nabla f(\mathbf{x}_t)^\top \mathbf{y}_t - \hat{\mathbf{g}}_{t-1}^\top \hat{\mathbf{y}}_{t-1}, \mathbf{u} - \mathbf{x}_t \rangle$

and $a_t \langle \nabla f(\mathbf{x}_t) - \hat{\mathbf{f}}_t, \mathbf{v} - \mathbf{y}_t \rangle$.

1. Controlling $\mathbb{E}[\mathcal{I}_k^P]$: The first step follows similarly to Step 1 from the proof of Proposition 3.4.1. By Young's inequality with parameter $(A_{k-1}\mu + \mu_0)/2$, we have

$$\begin{aligned} \mathbb{E}[\mathcal{E}_k^P] &= a_{k-1} \mathbb{E} \left\langle \frac{1}{p_{P_k}} \sum_{i \in B_{P_k}} (y_{k-1,i} \nabla f_i(\mathbf{x}_{k-1}) - \hat{y}_{k-2,i} \hat{\mathbf{g}}_{k-2,i}), \mathbf{x}_k - \mathbf{x}_{k-1} \right\rangle \\ &\leq \frac{1}{2} \mathbb{E}[\mathcal{C}_k^P] + \frac{a_{k-1}^2}{A_{k-1}\mu + \mu_0} \mathbb{E} \left\| \frac{1}{p_{P_k}} \sum_{i \in B_{P_k}} (y_{k-1,i} \nabla f_i(\mathbf{x}_{k-1}) - \hat{y}_{k-2,i} \hat{\mathbf{g}}_{k-2,i}) \right\|_{\mathcal{X}^*}^2 \end{aligned} \quad (3.60)$$

and

$$\begin{aligned} &- \mathbb{E}[a_k \langle \nabla f(\mathbf{x}_k)^\top \mathbf{y}_k - \hat{\mathbf{g}}_{k-1}^\top \hat{\mathbf{y}}_{k-1}, \mathbf{u} - \mathbf{x}_k \rangle] \\ &\leq \frac{1}{2} \mathbb{E}[\mathcal{T}_t^P] + \frac{a_{t-1}^2}{A_{t-1}\mu + \mu_0} \mathbb{E} \left\| \frac{1}{p_{P_t}} \sum_{i \in B_{P_t}} (y_{t,i} \nabla f_i(\mathbf{x}_t) - \hat{y}_{t-1,i} \hat{\mathbf{g}}_{t-1,i}) \right\|_{\mathcal{X}^*}^2. \end{aligned} \quad (3.61)$$

Then, we apply Lemma 3.4.1 with $b_I = 1$ and $c_I = 1$ to handle the second term in (3.60) (and (3.61)), and write

$$\begin{aligned} &\mathbb{E} \left\| \frac{1}{p_{P_k}} \sum_{i \in B_{P_k}} (y_{k-1,i} \nabla f_i(\mathbf{x}_{k-1}) - \hat{y}_{k-2,i} \hat{\mathbf{g}}_{k-2,i}) \right\|_{\mathcal{X}^*}^2 \\ &\leq 2 \left(\max_J \frac{\mathbf{L}_J^2}{p_J} \right) \sum_{I=1}^N \mathbb{E} \|\mathbf{x}_{k-1} - \hat{\mathbf{x}}_{k-2,I}\|_{\mathcal{X}}^2 + 2 \left(\max_J \frac{\mathbf{G}_J^2}{p_J} \right) \sum_{I=1}^N \mathbb{E} \|\mathbf{y}_{k-1,I} - \hat{\mathbf{y}}_{k-2,I}\|_2^2. \end{aligned}$$

Individually, the colored table bias terms can be further upper bounded by applying Lemma 3.4.2 (stated after this proof), so that the sums of (3.60) and (3.61) can be further developed to

$$\begin{aligned} \sum_{k=2}^t \mathbb{E}[\mathcal{I}_k^P] &\leq \frac{1}{2} \sum_{k=2}^t \mathbb{E}[\mathcal{C}_k^P] + \frac{1}{2} \mathbb{E}[\mathcal{T}_t^P] \\ &\quad + 10 \sum_{k=1}^t \frac{a_k^2}{A_k\mu + \mu_0} \sum_{I=1}^N \frac{\mathbf{G}_I^2}{p_I s_I} \sum_{k'=1}^k (1 - s_I/2)^{k-k'} \mathbb{E} \|\mathbf{y}_{k',I} - \mathbf{y}_{k'-1,I}\|_2^2 \end{aligned} \quad (3.62)$$

$$+ 10 \sum_{k=1}^t \frac{a_k^2}{A_k\mu + \mu_0} \sum_{I=1}^N \frac{\mathbf{L}_I^2}{p_I r_I} \sum_{k'=1}^k (1 - r_I/2)^{k-k'} \mathbb{E} \|\mathbf{x}_{k'} - \mathbf{x}_{k'-1}\|_{\mathcal{X}}^2. \quad (3.63)$$

To control the resulting sums (3.62) and (3.63), we exchange the order of summation to

compute them. First, (3.62) can be written as

$$\begin{aligned} & 10 \sum_{k'=1}^t \sum_{I=1}^N \frac{\mathbf{G}_I^2}{p_I s_I} \mathbb{E} \|\mathbf{y}_{k',I} - \mathbf{y}_{k'-1,I}\|_2^2 \cdot \sum_{k=k'}^t \frac{a_k^2}{A_k \mu + \mu_0} (1 - s_I/2)^{k-k'} \\ & \leq 50 \sum_{k'=1}^t \frac{a_{k'}^2}{A_{k'} \mu + \mu_0} \sum_{I=1}^N \frac{\mathbf{G}_I^2}{p_I s_I^2} \mathbb{E} \|\mathbf{y}_{k',I} - \mathbf{y}_{k'-1,I}\|_2^2 \end{aligned} \quad (3.64)$$

$$\leq 50 \underbrace{\left(\max_{I \in [N]} \frac{\mathbf{G}_I^2}{p_I s_I^2} \right)}_{\mathbf{G}_{p,s}^2} \sum_{k'=1}^t \frac{a_{k'}^2}{A_{k'} \mu + \mu_0} \mathbb{E} \|\mathbf{y}_{k'} - \mathbf{y}_{k'-1}\|_{\mathbf{y}}^2, \quad (3.65)$$

where the inequality (3.64) follows by the given assumption (3.57) that $\frac{a_k^2}{A_k \mu + \mu_0} \leq \min_I (1 + s_I/5) \frac{a_{k-1}^2}{A_{k-1} \mu + \mu_0}$ and the sequence of steps

$$\begin{aligned} \frac{a_k^2}{A_k \mu + \mu_0} \sum_{k=k'}^t (1 - s_I/2)^{k-k'} & \leq \sum_{k=k'}^t [(1 - s_I/2)(1 + s_I/5)]^{k-k'} \frac{a_{k'}^2}{A_{k'} \mu + \mu_0} \\ & \leq \sum_{k=k'}^t (1 - s_I/5)^{k-k'} \frac{a_{k'}^2}{A_{k'} \mu + \mu_0} \leq \frac{5}{s_I} \frac{a_{k'}^2}{A_{k'} \mu + \mu_0}. \end{aligned}$$

In (3.65), we also used that $\|\cdot\|_2 \leq \|\cdot\|_{\mathbf{y}}$. This argument is the most technical part of the analysis and is repeated two more times in the remainder of the proof. The first of the two is to upper bound (3.63) by the quantity

$$50 \underbrace{\left(\sum_{I=1}^N \frac{\mathbf{L}_I^2}{p_I r_I^2} \right)}_{\mathbf{L}_{p,r}^2} \sum_{k'=1}^t \frac{a_{k'}^2}{A_{k'} \mu + \mu_0} \mathbb{E} \|\mathbf{x}_{k'} - \mathbf{x}_{k'-1}\|_{\mathbf{x}}^2, \quad (3.66)$$

whereas the second appears in the steps used to bound the dual error terms below.

2. Controlling $\mathbb{E}[\mathcal{I}_k^D]$: The following steps are the dual analog of ones shown above for the primal terms. Applying Young's inequality with parameter $(A_{k-1}\nu + \nu_0)/2$ we have that

$$\begin{aligned}
\mathbb{E}_{k-1/2}[\mathcal{E}_k^D] &= a_{k-1} \mathbb{E}_{k-1/2} \left\langle \frac{1}{q_{Q_k}} \sum_{j \in B_{Q_k}} (f_j(\mathbf{x}_{k-1}) - \hat{f}_{k-2,j}) \mathbf{e}_j, \mathbf{y}_k - \mathbf{y}_{k-1} \right\rangle \\
&\leq \frac{1}{2} \mathbb{E}_{k-1/2}[\mathcal{C}_k^D] + \frac{a_{k-1}^2}{A_{k-1}\nu + \nu_0} \sum_{J=1}^N \frac{1}{q_J} \left\| \sum_{j \in B_J} (f_j(\mathbf{x}_{k-1}) - \hat{f}_{k-1,j}) \mathbf{e}_j \right\|_{y_*}^2 \\
&\leq \frac{1}{2} \mathbb{E}_{k-1/2}[\mathcal{C}_k^D] + \frac{a_{k-1}^2}{A_{k-1}\nu + \nu_0} \sum_{J=1}^N \frac{1}{q_J} \left\| \sum_{j \in B_J} (f_j(\mathbf{x}_{k-1}) - \hat{f}_{k-1,j}) \mathbf{e}_j \right\|_2^2 \quad (3.67) \\
&\leq \frac{1}{2} \mathbb{E}_{k-1/2}[\mathcal{C}_k^D] + \frac{a_{k-1}^2}{A_{k-1}\nu + \nu_0} \sum_{J=1}^N \frac{\mathbf{G}_J^2}{q_J} \|\mathbf{x}_{k-1} - \hat{\mathbf{x}}_{k-1,J}\|_{\mathcal{X}}^2
\end{aligned}$$

where in (3.67) we used that $\|\cdot\|_{y_*} \leq \|\cdot\|_2$. Summing over k and taking the marginal expectation,

$$\sum_{k=2}^t \mathbb{E}[\mathcal{I}_k^D] \leq \frac{1}{2} \sum_{k=2}^t \mathbb{E}[\mathcal{C}_k^D] + \frac{1}{2} \mathbb{E}[\mathcal{T}_t^D] + \frac{a_{k-1}^2}{A_{k-1}\nu + \nu_0} \sum_{J=1}^N \frac{\mathbf{G}_J^2}{q_J} \|\mathbf{x}_{k-1} - \hat{\mathbf{x}}_{k-1,J}\|_{\mathcal{X}}^2, \quad (3.68)$$

where (3.68) can be upper bounded using the same arguments leading to (3.66) under the given assumption (3.58), yielding

$$25 \underbrace{\left(\sum_{J=1}^N \frac{\mathbf{G}_J^2}{q_J r_J^2} \right)}_{\mathbf{G}_{q,r}^2} \sum_{k'=1}^t \frac{a_{k'}^2}{A_{k'}\nu + \nu_0} \mathbb{E} \|\mathbf{x}_{k'} - \mathbf{x}_{k'-1}\|_{\mathcal{X}}^2. \quad (3.69)$$

To summarize progress thus far, we must cancel the terms (3.65), (3.66), and (3.69) to complete the proof, which requires setting the appropriate conditions on the sequence $(a_k)_{k \geq 1}$.

3. Deriving the rate conditions: Under the condition (3.59), we may bound (3.65) by $\frac{1}{2} \sum_{k'=1}^t \mathbb{E}[\mathcal{C}_{k'}^D]$, (3.66) by $\frac{1}{4} \sum_{k'=1}^t \mathbb{E}[\mathcal{C}_{k'}^P]$, and (3.69) by $\frac{1}{4} \sum_{k'=1}^t \mathbb{E}[\mathcal{C}_{k'}^P]$. All terms of (3.24) now cancel, completing the proof. \square

The following technical lemma was used to express the terms that quantified the table bias terms (distance between the iterates and their counterparts in the respective tables)

with terms that can be canceled by quantities in (3.24).

Lemma 3.4.2. *[Diakonikolas, 2025, Lemma 2] For any $k \geq 1$ and $I \in [N]$, the following hold:*

$$\begin{aligned}\mathbb{E} \|\mathbf{x}_k - \hat{\mathbf{x}}_{k-1,I}\|_{\mathcal{X}}^2 &\leq \frac{5}{r_I} \sum_{k'=1}^k (1 - r_I/2)^{k-k'} \mathbb{E} \|\mathbf{x}_{k'} - \mathbf{x}_{k'-1}\|_{\mathcal{X}}^2, \\ \mathbb{E} \|\mathbf{y}_{k,I} - \hat{\mathbf{y}}_{k-1,I}\|_2^2 &\leq \frac{5}{s_I} \sum_{k'=1}^k (1 - s_I/2)^{k-k'} \mathbb{E} \|\mathbf{y}_{k',I} - \mathbf{y}_{k'-1,I}\|_2^2.\end{aligned}$$

We now convert the Proposition 3.4.2 into a complexity guarantee, using again the constant D_0 from Theorem 3.3.1.

Theorem 3.4.2. *Under Assumption 3.2.2 and Assumption 3.2.1, consider any $\mathbf{u} \in \mathcal{X}$, $\mathbf{v} \in \mathcal{Y}$ and precision $\varepsilon > 0$. Assume that $\min_I r_I \geq 1/(2N)$ and $\min_I s_I \geq 1/(2N)$. Define $(a_k)_{k=1}^t$ such that Algorithm 4 with Identity Card 2 produces an output point $(\tilde{\mathbf{x}}_t, \tilde{\mathbf{y}}_t) \in \mathcal{X} \times \mathcal{Y}$ satisfying $\mathbb{E} [\text{Gap}^{\mathbf{u},\mathbf{v}}(\tilde{\mathbf{x}}_t, \tilde{\mathbf{y}}_t)] \leq \varepsilon$ for t that depends on ε according to the following iteration complexities. Below, we use the constant $C_0 := (\mathbf{L}_{\mathbf{p},\mathbf{r}} \sqrt{\nu_0/\mu_0} + (\mathbf{G}_{\mathbf{p},\mathbf{s}} \vee \mathbf{G}_{\mathbf{q},\mathbf{r}})) D_0$.*

Case	Iteration Complexity
$\mu > 0$ and $\nu > 0$	$O\left(\left(N + \frac{\mathbf{L}_{\mathbf{p},\mathbf{r}}}{\mu} + \frac{\mathbf{G}_{\mathbf{p},\mathbf{s}} \vee \mathbf{G}_{\mathbf{q},\mathbf{r}}}{\sqrt{\mu\nu}}\right) \ln\left(\frac{C_0}{\varepsilon}\right)\right)$
$\mu > 0$ and $\nu = 0$	$O\left(\left(N + \frac{\mathbf{L}_{\mathbf{p},\mathbf{r}}}{\mu}\right) \ln\left(\frac{C_0}{\varepsilon}\right) + (\mathbf{G}_{\mathbf{p},\mathbf{s}} \vee \mathbf{G}_{\mathbf{q},\mathbf{r}}) \sqrt{\frac{\sqrt{\mu_0/\nu_0} D_0}{\mu\varepsilon}}\right)$
$\mu = 0$ and $\nu > 0$	$O\left(N \ln\left(\frac{C_0}{\varepsilon}\right) + \frac{\mathbf{L}_{\mathbf{p},\mathbf{r}} \sqrt{\nu_0/\mu_0} D_0}{\varepsilon} + (\mathbf{G}_{\mathbf{p},\mathbf{s}} \vee \mathbf{G}_{\mathbf{q},\mathbf{r}}) \sqrt{\frac{\sqrt{\nu_0/\mu_0} D_0}{\nu\varepsilon}}\right)$
$\mu = 0$ and $\nu = 0$	$O\left(N \ln\left(\frac{C_0}{\varepsilon}\right) + \frac{(\mathbf{L}_{\mathbf{p},\mathbf{r}} \sqrt{\nu_0/\mu_0} + (\mathbf{G}_{\mathbf{p},\mathbf{s}} \vee \mathbf{G}_{\mathbf{q},\mathbf{r}})) D_0}{\varepsilon}\right)$

Proof. The exact case-by-case strategy of Theorem 3.3.1 can be applied (ignoring absolute constant factors) by setting $G \leftarrow \mathbf{G}_{\mathbf{p},\mathbf{s}} \vee \mathbf{G}_{\mathbf{q},\mathbf{r}}$ and $L \leftarrow \mathbf{L}_{\mathbf{p},\mathbf{r}}$. The only additional conditions that need to be incorporated are (3.57) and (3.58). Under the given assumption that

$\min_I r_I \geq 1/(2N)$ and $\min_I s_I \geq 1/(2N)$, these conditions will be satisfied when, for all $k \geq 2$, it holds that $\frac{a_k^2}{A_k\mu+\mu_0} \leq (1 + 1/(10N)) \frac{a_{k-1}^2}{A_{k-1}\mu+\mu_0}$ and $\frac{a_k^2}{A_k\nu+\nu_0} \leq (1 + 1/(10N)) \frac{a_{k-1}^2}{A_{k-1}\nu+\nu_0}$.

Taking the first condition as an example, it can be rewritten as

$$\frac{a_k^2}{a_{k-1}^2} \leq \left(1 + \frac{1}{10N}\right) \frac{A_k\mu + \mu_0}{A_{k-1}\mu + \mu_0}.$$

When $\mu = 0$ and $\nu = 0$, this condition is satisfied automatically as the ratio on the left-hand side is equal to 1 (because a_k is a constant sequence). Otherwise, because A_k is an increasing sequence, we can reduce the condition to $a_k^2 \leq (1 + 1/(10N))a_{k-1}^2$. The fastest growth condition on $(a_k)_{k \geq 1}$ that is possible under the constraint is $a_k \leq (1 + \alpha)a_{k-1}$, where $(1 + \alpha) \leq \sqrt{1 + 1/(10N)}$. Then,

$$\sqrt{1 + 1/(10N)} \geq \sqrt{1 + 69/(900N)} \geq 1 + 1/(30N),$$

so the imposition $\alpha \leq \frac{1}{30N}$ suffices. This adds an $O(N \ln(C_0/\varepsilon))$ term to the resulting complexities and completes the proof. \square

We provide similar computations as those used in Theorem 3.4.1 to uncover the dependence on the sampling scheme. We again use $\boldsymbol{\lambda} = (\boldsymbol{\lambda}_1, \dots, \boldsymbol{\lambda}_N)$ where $\boldsymbol{\lambda}_I := \sqrt{\mathbf{G}_I^2 + \mathbf{L}_I^2}$ from Section 3.2. The non-uniform sampling complexity given below follows by letting $p_I \propto \boldsymbol{\lambda}_I^{1/2}$, $r_I \propto \boldsymbol{\lambda}_I^{1/2}$, $s_I \propto \mathbf{G}_I^{1/2}$ and $q_I \propto \mathbf{G}_I^{1/2}$.

Constant	Uniform Sampling	Non-Uniform Sampling
$\mathbf{G}_{p,s} \vee \mathbf{G}_{q,r}$	$N^{3/2} \ \mathbf{G}\ _2$	$\ \boldsymbol{\lambda}\ _{1/2}^{1/2} \ \mathbf{G}\ _{1/2}^{1/2}$
$\mathbf{L}_{p,r}$	$N^{3/2} \ \mathbf{L}\ _2$	$\ \boldsymbol{\lambda}\ _{1/2}^{3/4} \ \mathbf{L}\ _{1/2}^{1/4}$

Notice that the complexities in Theorem 3.4.1 (as opposed to the ones shown in Theorem 3.4.2) depend on additional factors in N . Thus, taking $\mu > 0$ and $\nu > 0$ as an example, in the case of uniform sampling, the method of Theorem 3.4.1 has an $N^{3/2}(\|\mathbf{G}\|_\infty + \|\mathbf{L}\|_\infty)$ dependence on problem constants, whereas Theorem 3.4.2 has an $N^{3/2}(\|\mathbf{G}\|_2 + \|\mathbf{L}\|_2)$ depen-

dence. These are the same in the case of highly non-uniform Lipschitz constants, but have a \sqrt{N} factor difference for approximately uniform Lipschitz constants.

This method inherits the exact per-iteration complexity as the one analyzed in Theorem 3.4.1. Maintaining the matrix-vector multiplication $\hat{\mathbf{g}}_k^\top \hat{\mathbf{y}}_k \in \mathbb{R}^d$ operates just as in (3.56). The total arithmetic complexity is given by $\tilde{O}(n(d + N)/N)$. The main difference in the algorithms from the upcoming Section 3.5, from a per-iteration complexity viewpoint, is that the full $\tilde{O}(n)$ cost update of \mathbf{y}_k is replaced by a single block update of cost $\tilde{O}(n/N)$. Thus, we aim to improve the per-iteration complexity to $\tilde{O}(nd/N)$ in the separable case.

3.5 An Algorithm for Dual-Separable Problems

Given Definition 3.2.1, we design an algorithm variant that performs stochastic block-wise updates in the dual variable, akin to similar strategies applied to bilinearly coupled objectives [Song et al., 2021]. Precisely, we will only update a single randomly chosen block Q_k on each iteration k , in an effort to achieve an improved complexity guarantee. Updates in this form introduce additional technical challenges because different blocks of the dual iterate $\mathbf{y}_k = (\mathbf{y}_{k,1}, \dots, \mathbf{y}_{k,N})$ have different dependences on the block index Q_k . As such, we carefully handle the expectations computations in the upcoming Lemma 3.5.1. Furthermore, a key difference in the proof structure of this section is that we will track an auxiliary sequence $(\bar{\mathbf{y}}_k)_{k \geq 1}$ of *return values*, such that the algorithm returns $(\mathbf{x}_t, \bar{\mathbf{y}}_t)$ in the final iteration instead of $(\mathbf{x}_t, \mathbf{y}_t)$. Conceptually, each block $\bar{\mathbf{y}}_{k,J}$ represents the J -th block of \mathbf{y}_k if $J = Q_k$, or if block J was the one updated at time k . In other words, it stores all possible block updates that could have occurred from step $k - 1$ to step k in one vector. Similar to before, we will define our update sequences in the process of deriving upper and lower bounds on $a_k \mathbb{E}[\mathcal{L}(\mathbf{x}_k, \mathbf{v})]$ and $a_k \mathbb{E}[\mathcal{L}(\mathbf{u}, \bar{\mathbf{y}}_k)]$.

Crucially, we do not need to compute the elements of the sequence $(\bar{\mathbf{y}}_k)_{k \geq 1}$, as doing so would defeat the purpose of considering coordinate-wise updates. Instead, we may realize (3.10) in expectation by computing only one instance of $\bar{\mathbf{y}}_k$ with the following trick [Alacaoglu et al., 2022]: we randomly draw an index \hat{t} (independent of all other randomness

in the algorithm and of (\mathbf{u}, \mathbf{v}) from $\{1, \dots, t\}$, where $\hat{t} = k$ with probability a_k/A_t . Thus, by computing the conditional expectation over \hat{t} given the sequence of iterates,

$$\mathbb{E}[\text{Gap}^{\mathbf{u}, \mathbf{v}}(\mathbf{x}_{\hat{t}}, \bar{\mathbf{y}}_{\hat{t}}) | \mathcal{F}_t] = \sum_{k=1}^t a_k \text{Gap}^{\mathbf{u}, \mathbf{v}}(\mathbf{x}_k, \bar{\mathbf{y}}_k).$$

In practice, we may simply run the algorithm to iteration \hat{t} . We use the following technical lemma to provide expectation formulas regarding \mathbf{y}_k and $\bar{\mathbf{y}}_k$.

Lemma 3.5.1. *Let $h : \mathcal{Y} \times \mathcal{X} \rightarrow \mathbb{R}$ be block separable in its first argument, i.e., $h(\mathbf{y}, \mathbf{x}) = \sum_{J=1}^N h_J(\mathbf{y}_J, \mathbf{x})$ for $J \in [N]$ and $\mathbf{y} \in \mathcal{Y}$. Assume that $\mathbf{y}_{k,J} = \mathbf{y}_{k-1,J}$ for all $J \neq Q_k$, that $\bar{\mathbf{y}}_{k,Q_k} = \mathbf{y}_{k,Q_k}$, and that $\bar{\mathbf{y}}_k$ is $\mathcal{F}_{k-1/2}$ -measurable. Then, if Q_k is sampled uniformly on $[N]$, it holds that*

$$N\mathbb{E}[h_{Q_k}(\mathbf{y}_{k,Q_k}, \mathbf{x}_k)] = N\mathbb{E}[h(\mathbf{y}_k, \mathbf{x}_k)] - (N-1)\mathbb{E}[h(\mathbf{y}_{k-1}, \mathbf{x}_k)] = \mathbb{E}[h(\bar{\mathbf{y}}_k, \mathbf{x}_k)]. \quad (3.70)$$

Proof. Write

$$\begin{aligned} \mathbb{E}_{k-1/2}[h_{Q_k}(\mathbf{y}_{k,Q_k}, \mathbf{x}_k)] &= \mathbb{E}_{k-1/2}[h(\mathbf{y}_k, \mathbf{x}_k)] - \mathbb{E}_{k-1/2}\left[\sum_{J \neq Q_k} h_J(\mathbf{y}_{k,J}, \mathbf{x}_k)\right] \\ &= \mathbb{E}_{k-1/2}[h(\mathbf{y}_k, \mathbf{x}_k)] - \mathbb{E}_{k-1/2}\left[\sum_{J \neq Q_k} h_J(\mathbf{y}_{k-1,J}, \mathbf{x}_k)\right] \\ &= \mathbb{E}_{k-1/2}[h(\mathbf{y}_k, \mathbf{x}_k)] - \frac{N-1}{N}h(\mathbf{y}_{k-1}, \mathbf{x}_k). \end{aligned}$$

Take the marginal expectation of both terms to prove the first result. For the second,

$$\begin{aligned} \mathbb{E}_{k-1/2}[h(\mathbf{y}_k, \mathbf{x}_k)] &= \frac{1}{N} \sum_{J=1}^N \mathbb{E}_{k-1/2}[h_J(\bar{\mathbf{y}}_{k,J}, \mathbf{x}_k) | J = Q_k] \\ &\quad + \frac{N-1}{N} \sum_{J=1}^N \mathbb{E}_{k-1/2}[h_J(\mathbf{y}_{k-1,J}, \mathbf{x}_k) | J \neq Q_k] \\ &= \frac{1}{N} \sum_{J=1}^N h_J(\bar{\mathbf{y}}_{k,J}, \mathbf{x}_k) + \frac{N-1}{N} \sum_{J=1}^N h_J(\mathbf{y}_{k-1,J}, \mathbf{x}_k) \\ &= \frac{1}{N} h(\bar{\mathbf{y}}_k, \mathbf{x}_k) + \frac{N-1}{N} h(\mathbf{y}_{k-1}, \mathbf{x}_k). \end{aligned}$$

Rearrange terms and apply the marginal expectation to achieve the second result. \square

The condition that $\bar{\mathbf{y}}_k$ is $\mathcal{F}_{k-1/2}$ -measurable reflects the viewpoint that $\bar{\mathbf{y}}_k$ can be computed via a deterministic update given \mathbf{x}_k , after which \mathbf{y}_k can be computed exactly from the random triple $(\mathbf{y}_{k-1}, \bar{\mathbf{y}}_k, Q_k)$.

We proceed to the details of the convergence analysis. The formula for $\bar{\mathbf{g}}_{k-1} \in \mathbb{R}^d$ is given after rigorously introducing the sequence $(\bar{\mathbf{y}}_k)_{k \geq 0}$, on which $\bar{\mathbf{g}}_{k-1}$ depends. We will first specify the upper bound and dual update. We define first the update for \mathbf{y}_k (and by extension, the update for $\bar{\mathbf{y}}_k$). Because of separability, we may perform this update at $O(b)$ for $b := n/N$ cost on average across blocks. We do so by including only one additive component of ψ in the objective that \mathbf{y}_k maximizes. We notice that the strong concavity constant of the objective defining \mathbf{y}_k may change if we only use some components of ψ . To account for this, we design a slightly different proximity term from the one in Algorithm 4. Recall the definition of the Bregman divergences $\Delta_J(\cdot, \cdot)$ from Section 3.2. The upper bound is stated in expectation, as the results may not hold for all realizations (as was the case for previous versions of the initial gap bounds).

Lemma 3.5.2. *For all $k \geq 1$, consider the update*

$$\bar{\mathbf{y}}_k = \arg \max_{\mathbf{y} \in \mathcal{Y}} \left\{ a_k \langle \mathbf{y}, \bar{\mathbf{f}}_{k-1/2} \rangle - a_k \psi(\mathbf{y}) - \frac{A_{k-1}\nu + \nu_0}{2} \Delta_{\mathcal{Y}}(\mathbf{y}, \mathbf{y}_{k-1}) \right\}, \quad (3.71)$$

followed by setting $\mathbf{y}_{k,Q_k} = \bar{\mathbf{y}}_{k,Q_k}$ and $\mathbf{y}_{k,J} = \mathbf{y}_{k-1,J}$ for $J \neq Q_k$. Then, it holds that

$$\begin{aligned} a_k \mathbb{E}[\mathcal{L}(\mathbf{x}_k, \mathbf{v})] &\leq a_k \mathbb{E}[\mathcal{L}(\mathbf{x}_k, \bar{\mathbf{y}}_k)] + \frac{N}{2} \mathbb{E}[\mathcal{T}_{k-1}^D] - \frac{N}{2} \mathbb{E}[\mathcal{T}_k^D] - \frac{N}{2} \mathbb{E}[\mathcal{C}_k^D] \\ &\quad + a_k \mathbb{E} \langle \mathbf{v} - \bar{\mathbf{y}}_k, f(\mathbf{x}_k) - \bar{\mathbf{f}}_{k-1/2} \rangle. \end{aligned} \quad (3.72)$$

Proof. As in Lemma 3.3.2, write

$$\begin{aligned}
a_k \mathcal{L}(\mathbf{x}_k, \mathbf{v}) &= a_k \langle \mathbf{v}, \bar{\mathbf{f}}_{k-1/2} \rangle - a_k \psi(\mathbf{v}) + a_k \phi(\mathbf{x}_k) - \frac{A_{k-1}\nu+\nu_0}{2} \Delta_{\mathbf{y}}(\mathbf{v}, \mathbf{y}_{k-1}) \\
&\quad + \frac{A_{k-1}\nu+\nu_0}{2} \Delta_{\mathbf{y}}(\mathbf{v}, \mathbf{y}_{k-1}) + a_k \langle \mathbf{v}, f(\mathbf{x}_k) - \bar{\mathbf{f}}_{k-1/2} \rangle \\
&\leq a_k \langle \bar{\mathbf{y}}_k, \bar{\mathbf{f}}_{k-1/2} \rangle - a_k \psi(\bar{\mathbf{y}}_k) + a_k \phi(\mathbf{x}_k) - \frac{A_{k-1}\nu+\nu_0}{2} \Delta_{\mathbf{y}}(\bar{\mathbf{y}}_k, \mathbf{y}_{k-1}) \\
&\quad + \frac{A_{k-1}\nu+\nu_0}{2} \Delta_{\mathbf{y}}(\mathbf{v}, \mathbf{y}_{k-1}) + a_k \langle \mathbf{v}, f(\mathbf{x}_k) - \bar{\mathbf{f}}_{k-1/2} \rangle - \frac{A_{k-1}\nu+\nu_0}{2} \Delta_{\mathbf{y}}(\mathbf{v}, \bar{\mathbf{y}}_k) \\
&= a_k \mathcal{L}(\mathbf{x}_k, \bar{\mathbf{y}}_k) - \frac{A_{k-1}\nu+\nu_0}{2} \Delta_{\mathbf{y}}(\bar{\mathbf{y}}_k, \mathbf{y}_{k-1}) \\
&\quad + \frac{A_{k-1}\nu+\nu_0}{2} \Delta_{\mathbf{y}}(\mathbf{v}, \mathbf{y}_{k-1}) + a_k \langle \mathbf{v} - \bar{\mathbf{y}}_k, f(\mathbf{x}_k) - \bar{\mathbf{f}}_{k-1/2} \rangle - \frac{A_{k-1}\nu+\nu_0}{2} \Delta_{\mathbf{y}}(\mathbf{v}, \bar{\mathbf{y}}_k),
\end{aligned}$$

where the inequality follows by Lemma 3.2.1 (dropping the non-positive term $\frac{a_k \nu}{2} \Delta_{\mathbf{y}}(\mathbf{v}, \bar{\mathbf{y}}_k)$).

We take the marginal expectation and use Lemma 3.5.1 to achieve

$$\begin{aligned}
\frac{A_{k-1}\nu+\nu_0}{2} \mathbb{E}[\Delta_{\mathbf{y}}(\mathbf{v}, \mathbf{y}_{k-1})] - \frac{A_{k-1}\nu+\nu_0}{2} \mathbb{E}[\Delta_{\mathbf{y}}(\mathbf{v}, \bar{\mathbf{y}}_k)] &= \frac{N}{2} \mathbb{E}[\mathcal{T}_{k-1}^{\mathbf{D}}] - \frac{N}{2} \mathbb{E}[\mathcal{T}_k^{\mathbf{D}}] \\
\frac{A_{k-1}\nu+\nu_0}{2} \mathbb{E}[\Delta_{\mathbf{y}}(\bar{\mathbf{y}}_k, \mathbf{y}_{k-1})] &= \frac{N}{2} \mathbb{E}[\mathcal{C}_k^{\mathbf{D}}],
\end{aligned}$$

completing the proof. \square

Notice that we need only compute $\bar{\mathbf{y}}_{k,Q_k}$ in order to define \mathbf{y}_k . While the update (3.71) is written for the purpose of the proof, notice that for any $J \in [N]$, under Definition 3.2.1,

$$\bar{\mathbf{y}}_{k,J} = \arg \max_{\mathbf{y}_J \in \mathcal{Y}_J} \left\{ a_k \langle \mathbf{y}_J, \bar{\mathbf{f}}_{k-1/2,J} \rangle - a_k \psi_J(\mathbf{y}_J) - \frac{A_{k-1}\nu+\nu_0}{2} \Delta_J(\mathbf{y}_J, \mathbf{y}_{k-1,J}) \right\}.$$

Defining the primal update will reflect two different strategies—namely, those used in Section 3.4. For both cases, recall the table $\hat{\mathbf{x}}_{k,1}, \dots, \hat{\mathbf{x}}_{k,N}$ introduced in (3.37) and (3.38). We will use the primal gradient estimate

$$\bar{\mathbf{g}}_{k-1} = \hat{\mathbf{g}}_{k-1}^{\top} \mathbf{y}_{k-1} + \frac{N a_{k-1}}{a_k} \sum_{i \in B_{P_k}} (\bar{y}_{k-1,i} \nabla f_i(\mathbf{x}_{k-1}) - y_{k-2,i} \hat{\mathbf{g}}_{k-2,i}). \quad (3.73)$$

Notice that we do not need to use the table $\hat{\mathbf{y}}_k$ introduced in (3.39) from Section 3.4, as the matrix-vector product $\hat{\mathbf{g}}_{k-1}^{\top} \mathbf{y}_{k-1}$ above can be maintained in $O(bd)$ on average because only $b = n/N$ components of \mathbf{y}_k change each iteration. When using these tables of iterates,

we will encounter the familiar terms $\hat{\mathcal{T}}_{k,I}^P = (A_k\mu + \mu_0)\Delta_{\mathcal{X}}(\mathbf{u}, \hat{\mathbf{x}}_{k,I})$ from (3.15) and $\hat{\mathcal{C}}_{k,I}^P = (A_{k-1}\mu + \mu_0)\Delta_{\mathcal{X}}(\mathbf{x}_k, \hat{\mathbf{x}}_{k-1,I})$ from (3.16) for $I = 1, \dots, N$.

With these elements in hand, we produce the following lower bound, which follows from identical steps to Lemma 3.3.1 and so has its proof omitted. Recall the use of probability weights $\gamma_1, \dots, \gamma_N$ in the update (3.11).

Lemma 3.5.3. *For any $k \geq 1$, let $\bar{\mathbf{g}}_{k-1} \in \mathbb{R}^d$ and $w_k^P \in [0, 1)$, consider the update*

$$\begin{aligned} \mathbf{x}_k = \arg \min_{\mathbf{x} \in \mathcal{X}} \Big\{ & a_k \langle \bar{\mathbf{g}}_{k-1}, \mathbf{x} \rangle + a_k \phi(\mathbf{x}) + \frac{1-w_k^P}{2} (A_{k-1}\mu + \mu_0) \Delta_{\mathcal{X}}(\mathbf{x}, \mathbf{x}_{k-1}) \\ & + \frac{w_k^P}{2} (A_{k-1}\mu + \mu_0) \sum_{I=1}^N \gamma_I \Delta_{\mathcal{X}}(\mathbf{x}, \hat{\mathbf{x}}_{k-1,I}) \Big\}. \end{aligned} \quad (3.74)$$

For $k \geq 1$, it holds that

$$a_k \mathcal{L}(\mathbf{u}, \bar{\mathbf{y}}_k) \geq a_k \mathcal{L}(\mathbf{x}_k, \bar{\mathbf{y}}_k) + a_k \langle \nabla f(\mathbf{x}_k)^\top \bar{\mathbf{y}}_k - \bar{\mathbf{g}}_{k-1}, \mathbf{u} - \mathbf{x}_k \rangle \quad (3.75)$$

$$+ \left(\frac{1-w_k^P}{2} \mathcal{T}_k^P - \frac{1-w_{k-1}^P}{2} \mathcal{T}_{k-1}^P \right) + \frac{w_k^P}{2} \left(\mathcal{T}_k^P - \sum_{I=1}^N \gamma_I \hat{\mathcal{T}}_{k-1,I}^P \right) \quad (3.76)$$

$$+ \frac{1-w_k^P}{2} \mathcal{C}_k^P + \frac{w_k^P}{2} \sum_{I=1}^N \gamma_I \hat{\mathcal{C}}_{k-1,I}^P + \frac{a_k \mu}{2} \Delta_{\mathcal{X}}(\mathbf{u}, \mathbf{x}_k). \quad (3.77)$$

The only difference between Lemma 3.5.3 and Lemma 3.3.1 is the replacement of \mathbf{y}_k by $\bar{\mathbf{y}}_k$, which is also an element of \mathcal{Y} . On the dual side, we set $\bar{\mathbf{f}}_{k-1/2} = f(\mathbf{x}_k)$, so the inner product terms in (3.72) vanish in expectation (in other words, \mathcal{I}_k^D can be thought of as zero). Following the same argument used to produce \mathcal{E}_k^P in (3.42), we may now build the identity card.

Identity Card 3: Stochastic update method for separable objectives

SUBROUTINE₁: $\hat{\mathbf{g}}_{k-1}^\top \mathbf{y}_{k-1} + \frac{Na_{k-1}}{a_k} \sum_{i \in B_{P_k}} (\bar{y}_{k-1,i} \nabla f_i(\mathbf{x}_{k-1}) - y_{k-2,i} \hat{\mathbf{g}}_{k-2,i})$

SUBROUTINE₂: $\bar{\mathbf{f}}_{k-1/2} = f(\mathbf{x}_k)$

SUBROUTINE₃: Update $\hat{\mathbf{x}}_{k,I}$ for all I using (3.38) and $(\hat{\mathbf{g}}_k, \hat{\mathbf{f}}_k)$ using (3.37).

Primal error: $\mathcal{E}_k^P = a_{k-1} \left\langle \frac{1}{p_{P_k}} \sum_{i \in B_{P_k}} (\bar{y}_{k-1,i} \nabla f_i(\mathbf{x}_{k-1}) - y_{k-2,i} \hat{\mathbf{g}}_{k-2,i}), \mathbf{x}_k - \mathbf{x}_{k-1} \right\rangle$

Dual error: $\mathcal{E}_k^D = 0$

Using Lemma 3.5.2 and Lemma 3.5.3, we may produce a version of Claim 3.3.1:

$$\begin{aligned} \sum_{k=1}^t a_k \mathbb{E}[\text{Gap}^{\mathbf{u}, \mathbf{v}}(\mathbf{x}_k, \mathbf{y}_k)] &\leq \frac{1-w_k^P}{2} (\mathcal{T}_0^P - \mathbb{E}[\mathcal{T}_t^P]) + \frac{N}{2} (\mathcal{T}_0^D - \mathbb{E}[\mathcal{T}_t^D]) \\ &+ \sum_{k=1}^t \mathbb{E}[\mathcal{I}_k^P] - \underbrace{\mathbb{E} \left[\frac{1-w_k^P}{2} \mathcal{C}_k^P + \frac{N}{2} \mathcal{C}_k^D + \frac{w_k^P}{2} \sum_J \gamma_I \hat{\mathcal{C}}_{k,J}^P \right]}_{\text{cancellation terms from Lemma 3.5.2 and Lemma 3.5.3}} \end{aligned} \quad (3.78)$$

$$\begin{aligned} &+ \underbrace{\frac{1}{2} \sum_{k=1}^t \mathbb{E} \left[w_k^P \left(\sum_{I=1}^N \gamma_I \hat{\mathcal{T}}_{k-1,I}^P - \mathcal{T}_k^P \right) - a_k \mu \Delta_{\mathcal{X}}(\mathbf{u}, \mathbf{x}_k) \right]}_{\text{primal table terms from Lemma 3.5.3}}. \end{aligned} \quad (3.79)$$

We proceed to the individual analyses to cancel the lines (3.78) and (3.79).

3.5.1 Strategy 1: Non-Uniform Historical Regularization

Here, we do not assume that $w_k^P = 0$, and so will cancel both the lines (3.78) and (3.79). The resulting complexity will depend on the sampling probabilities \mathbf{p} and regularization weights γ through the following constants:

$$\mathbf{G}_{\mathbf{p}} := \sqrt{\max_{I \in [N]} \frac{\mathbf{G}_I^2}{p_I}} \text{ and } \mathbf{L}_{\mathbf{p}, \gamma} := \sqrt{\max_{I \in [N]} \frac{\mathbf{L}_I^2}{p_I \gamma_I}}.$$

The proof follows very similar steps to the proof of Proposition 3.4.1.

Proposition 3.5.1. *Let $(\mathbf{x}_0, \mathbf{y}_0) \in \text{ri}(\text{dom}(\phi)) \times \text{ri}(\text{dom}(\psi))$ and $\{(\mathbf{x}_k, \mathbf{y}_k)\}_{k \geq 1}$ be generated by the update from Lemma 3.5.3 with sequence $w_k^P \in [0, 1)$ and Lemma 3.5.2, with $\bar{\mathbf{g}}_{k-1}$ and $\bar{\mathbf{f}}_{k-1/2}$ given by (3.73), and $f(\mathbf{x}_k)$, respectively. Define $a_1 = \min \left\{ \frac{\sqrt{(1-w_0^P)\mu_0\nu_0}}{2\sqrt{2}\mathbf{G}_P}, \frac{\sqrt{w_0^P(1-w_0^P)\mu_0}}{4\mathbf{L}_{P,\gamma}} \right\}$ and select $(a_k)_{k \geq 2}$ such that the conditions*

$$a_k \leq \min \left\{ \frac{\sqrt{(1-w_k^P)(A_k\mu+\mu_0)(A_{k-1}\nu+\nu_0)}}{2\sqrt{2}\mathbf{G}_P}, \frac{\sqrt{w_k^P(1-w_k^P)(A_k\mu+\mu_0)(A_{k-1}\mu+\mu_0)}}{4\mathbf{L}_{P,\gamma}} \right\}, \quad (3.80)$$

are satisfied. In addition, impose that for any $\ell = 1, \dots, t-1$, it holds that

$$\frac{\mu}{N} \sum_{k'=0}^{\infty} w_{\ell+k'+1}^P A_{\ell+k'} (1-1/N)^{k'} \leq (w_\ell^P A_\ell + a_\ell) \mu \quad (3.81)$$

and that $A_k \leq \left(1 + \frac{1}{2N}\right)^k a_1$. We have that for any $\mathbf{u} \in \mathcal{X}$ and $\mathbf{v} \in \mathcal{Y}$,

$$\sum_{k=1}^t a_k \mathbb{E}[\text{Gap}^{\mathbf{u}, \mathbf{v}}(\mathbf{x}_k, \bar{\mathbf{y}}_k)] + \frac{1-w_k^P}{4} \mathbb{E}[\mathcal{T}_t^P] + \frac{1}{4} \mathbb{E}[\mathcal{T}_t^D] \leq (a_1\mu + \mu_0) \Delta_{\mathcal{X}}(\mathbf{u}, \mathbf{x}_0) + \frac{1}{2} \mathcal{T}_0^D.$$

Proof. Just as in the proof of Proposition 3.4.1, we aim to show that the sum of terms in (3.78) and (3.79) is a constant independent of t . These lines are exactly analogous to (3.24) and (3.25), and so we control them in a similar fashion.

1. Controlling (3.78): By Young's inequality with parameter $(1-w_k^P)(A_{k-1}\mu + \mu_0)/4$,

$$\begin{aligned} \mathbb{E}[\mathcal{E}_k^P] &= a_{k-1} \mathbb{E} \left\langle \frac{1}{p_{P_k}} \sum_{i \in B_{P_k}} (\bar{y}_{k-1,i} \nabla f_i(\mathbf{x}_{k-1}) - y_{k-2,i} \hat{\mathbf{g}}_{k-2,i}), \mathbf{x}_k - \mathbf{x}_{k-1} \right\rangle \\ &\leq \frac{1-w_k^P}{4} \mathbb{E}[\mathcal{C}_k^P] + \frac{2a_{k-1}^2}{(1-w_k^P)(A_{k-1}\mu + \mu_0)} \mathbb{E} \left\| \frac{1}{p_{P_k}} \sum_{i \in B_{P_k}} (\bar{y}_{k-1,i} \nabla f_i(\mathbf{x}_{k-1}) - y_{k-2,i} \hat{\mathbf{g}}_{k-2,i}) \right\|_{\mathcal{X}^*}^2 \end{aligned} \quad (3.82)$$

and for the last term

$$\begin{aligned} & - \mathbb{E}[a_t \langle \nabla f(\mathbf{x}_t)^\top \bar{\mathbf{y}}_t - \hat{\mathbf{g}}_{t-1}^\top \mathbf{y}_{t-1}, \mathbf{u} - \mathbf{x}_t \rangle] \\ & \leq \frac{1-w_t^P}{4} \mathbb{E}[\mathcal{T}_t^P] + \frac{2a_t^2}{(1-w_k^P)(A_t\mu + \mu_0)} \mathbb{E} \left\| \frac{1}{p_{P_t}} \sum_{i \in B_{P_t}} (\bar{y}_{t,i} \nabla f_i(\mathbf{x}_t) - y_{t-1,i} \hat{\mathbf{g}}_{t-1,i}) \right\|_{\mathcal{X}^*}^2 \end{aligned} \quad (3.83)$$

For the second term in (3.82) and (3.83), apply Lemma 3.4.1 with $b_I = \gamma_I$ and $c_I = 1$ to achieve

$$\begin{aligned}
& \mathbb{E} \left\| \frac{1}{p_{P_k}} \sum_{i \in B_{P_k}} (\bar{y}_{k-1,i} \nabla f_i(\mathbf{x}_{k-1}) - y_{k-2,i} \hat{\mathbf{g}}_{k-2,i}) \right\|_{\mathcal{X}^*}^2 \\
& \leq 2 \left(\max_I \frac{\mathbf{L}_J^2}{p_J \gamma_J} \right) \cdot \sum_{I=1}^N \gamma_I \mathbb{E} \|\mathbf{x}_{k-1} - \hat{\mathbf{x}}_{k-2,I}\|_{\mathcal{X}}^2 + 2 \left(\max_J \frac{\mathbf{G}_J^2}{p_J} \right) \mathbb{E} \|\bar{\mathbf{y}}_{k-1} - \mathbf{y}_{k-2}\|_2^2 \\
& \stackrel{(*)}{\leq} 2 \left(\max_I \frac{\mathbf{L}_J^2}{p_J \gamma_J} \right) \cdot \sum_{I=1}^N \gamma_I \mathbb{E} \|\mathbf{x}_{k-1} - \hat{\mathbf{x}}_{k-2,I}\|_{\mathcal{X}}^2 + 2N \left(\max_J \frac{\mathbf{G}_J^2}{p_J} \right) \mathbb{E} \|\mathbf{y}_{k-1} - \mathbf{y}_{k-2}\|_{\mathcal{Y}}^2 \\
& \stackrel{(\circ)}{\leq} \underbrace{4 \left(\max_J \frac{\mathbf{L}_J^2}{p_J \gamma_J} \right) \cdot \sum_{I=1}^N \gamma_I \mathbb{E} [\Delta_{\mathcal{X}}(\mathbf{x}_{k-1}, \hat{\mathbf{x}}_{k-2,I})]}_{\mathbf{L}_{\mathbf{p},\gamma}^2} + \underbrace{4N \left(\max_J \frac{\mathbf{G}_J^2}{p_J} \right) \mathbb{E} [\Delta_{\mathcal{Y}}(\mathbf{y}_{k-1}, \mathbf{y}_{k-2})]}_{\mathbf{G}_{\mathbf{p}}^2}
\end{aligned}$$

where in $(*)$ we applied Lemma 3.5.1 and $\|\cdot\|_2 \leq \|\cdot\|_{\mathcal{Y}}$ to the second term and in (\circ) we applied $\|\cdot\|_2 \leq \|\cdot\|_{\mathcal{Y}}$ and the strong convexity of Bregman divergences.

Combining the steps above, recalling that $\hat{\mathcal{C}}_{k,I}^{\mathbf{P}} := (A_{k-1}\mu + \mu_0)\Delta_{\mathcal{X}}(\mathbf{x}_k, \hat{\mathbf{x}}_{k-1,I})$, and applying the condition (3.80), we have the upper bound

$$\begin{aligned}
\mathbb{E}[\mathcal{E}_k^{\mathbf{P}}] & \leq \frac{1 - w_k^{\mathbf{P}}}{4} \mathbb{E}[\mathcal{C}_k^{\mathbf{P}}] + \frac{4Na_{k-1}^2 \mathbf{G}_{\mathbf{p}}^2 \mathbb{E}[\mathcal{C}_{k-1}^{\mathbf{D}}]}{(1 - w_k^{\mathbf{P}})(A_{k-1}\mu + \mu_0)(A_{k-2}\nu + \nu_0)} \\
& \quad + \frac{4a_{k-1}^2 \mathbf{L}_{\mathbf{p},\gamma}^2 \mathbb{E} \left[\sum_{I=1}^N \gamma_I \hat{\mathcal{C}}_{k-1,I}^{\mathbf{P}} \right]}{(1 - w_k^{\mathbf{P}})(A_{k-1}\mu + \mu_0)(A_{k-2}\mu + \mu_0)} \\
& \leq \frac{1 - w_k^{\mathbf{P}}}{2} \mathbb{E}[\mathcal{C}_k^{\mathbf{P}}] + \frac{N}{2} \mathbb{E}[\mathcal{C}_{k-1}^{\mathbf{D}}] + \frac{w_k^{\mathbf{P}}}{4} \mathbb{E} \left[\sum_{I=1}^N \gamma_I \hat{\mathcal{C}}_{k-1,I}^{\mathbf{P}} \right]
\end{aligned}$$

with a similar bound holding for $-\mathbb{E}[a_t \langle \nabla f(\mathbf{x}_t)^\top \bar{\mathbf{y}}_t - \hat{\mathbf{g}}_{t-1}^\top \mathbf{y}_{t-1}, \mathbf{u} - \mathbf{x}_t \rangle]$ except with $\mathbb{E}[\mathcal{C}_k^{\mathbf{P}}]$ replaced by $\mathbb{E}[\mathcal{T}_t^{\mathbf{P}}]$. These terms will cancel with the corresponding non-positive terms in (3.78).

2. Controlling (3.79): This follows from an identical argument to the one used to bound (3.25) in the proof of Proposition 3.4.1, appealing to the condition (3.81). \square

The following result mirrors the logic of Theorem 3.4.1; the proof is omitted.

Theorem 3.5.1. *Under Assumption 3.2.2 and Assumption 3.2.1, consider any $\mathbf{u} \in \mathcal{X}$, $\mathbf{v} \in \mathcal{Y}$*

and precision $\varepsilon > 0$. Define the initial distance constant

$$D_{0,N} = \sqrt{\frac{\mu_0}{\nu_0}} \Delta_{\mathcal{X}}(\mathbf{u}, \mathbf{x}_0) + \sqrt{\frac{\nu_0}{\mu_0}} N \Delta_{\mathcal{Y}}(\mathbf{v}, \mathbf{y}_0) \quad (3.84)$$

There exists a choice of the sequence $(a_k)_{k=1}^t$ and parameters $(w_k^{\mathbf{P}})_{k=0}^t$ such that Algorithm 4 with Identity Card 3 produces an output point $(\tilde{\mathbf{x}}_t, \tilde{\mathbf{y}}_t) \in \mathcal{X} \times \mathcal{Y}$ satisfying $\mathbb{E} [\text{Gap}^{\mathbf{u}, \mathbf{v}}(\tilde{\mathbf{x}}_t, \tilde{\mathbf{y}}_t)] \leq \varepsilon$ for t that depends on ε according to the following iteration complexities. They depend logarithmically on the constants $C_1 := \sqrt{N}(\mathbf{L}_{\mathbf{p}, \gamma} \sqrt{\nu_0/\mu_0} + \mathbf{G}_{\mathbf{p}}) D_{0,N} + \mu \Delta_{\mathcal{X}}(\mathbf{u}, \mathbf{x}_0)$, and $C_2 := (\sqrt{N} \mathbf{L}_{\mathbf{p}, \gamma} \sqrt{\nu_0/\mu_0} + \mathbf{G}_{\mathbf{p}}) D_{0,N} + \mu \Delta_{\mathcal{X}}(\mathbf{u}, \mathbf{x}_0)$.

Case	Iteration Complexity
$\mu > 0$ and $\nu > 0$	$O\left(\left(N + \frac{\sqrt{N} \mathbf{L}_{\mathbf{p}, \gamma}}{\mu} + \frac{\sqrt{N} \mathbf{G}_{\mathbf{p}}}{\sqrt{\mu\nu}}\right) \ln\left(\frac{C_1}{\varepsilon}\right)\right)$
$\mu > 0$ and $\nu = 0$	$O\left(\left(N + \frac{\sqrt{N} \mathbf{L}_{\mathbf{p}, \gamma}}{\mu}\right) \ln\left(\frac{C_2}{\varepsilon}\right) + \mathbf{G}_{\mathbf{p}} \sqrt{\frac{\sqrt{\mu_0/\nu_0} D_{0,N} + (a_1 \mu/\nu_0) \Delta_{\mathcal{X}}(\mathbf{u}, \mathbf{x}_0)}{\mu\varepsilon}}\right)$
$\mu = 0$ and $\nu > 0$	$O\left(\frac{\mathbf{L}_{\mathbf{p}, \gamma} \sqrt{\nu_0/\mu_0} D_{0,N}}{\varepsilon} + N \mathbf{G}_{\mathbf{p}} \sqrt{\frac{\sqrt{\nu_0/\mu_0} D_{0,N}}{\nu\varepsilon}}\right)$
$\mu = 0$ and $\nu = 0$	$O\left(\frac{(\mathbf{L}_{\mathbf{p}, \gamma} \sqrt{\nu_0/\mu_0} + \mathbf{G}_{\mathbf{p}}) D_{0,N}}{\varepsilon}\right)$

In the discussions in Section 3.6, we set $\nu_0 \sim \mu_0/N$, so the $D_{0,N}$ term appearing in Theorem 3.5.1 is interpreted as \sqrt{N} times a constant.

We recall from Section 3.2 that $\boldsymbol{\lambda} = (\boldsymbol{\lambda}_1, \dots, \boldsymbol{\lambda}_N)$ for $\boldsymbol{\lambda}_I := \sqrt{\mathbf{G}_I^2 + \mathbf{L}_I^2}$. The non-uniform sampling complexity given below follows by letting $p_I \propto \boldsymbol{\lambda}_I$ and $\gamma_I \propto \mathbf{L}_I$.

Constant	Uniform Sampling	Non-Uniform Sampling
$\mathbf{G}_{\mathbf{p}}$	$\sqrt{N} \ \mathbf{G}\ _{\infty}$	$\ \boldsymbol{\lambda}\ _1^{1/2} \ \mathbf{G}\ _{\infty}^{1/2}$
$\mathbf{L}_{\mathbf{p}, \gamma}$	$N \ \mathbf{L}\ _{\infty}$	$\ \boldsymbol{\lambda}\ _1^{1/2} \ \mathbf{L}\ _1^{1/2}$

There are strict advantages both in the dependence on \mathbf{G} and \mathbf{L} in the separable case over, say, the complexities given in Theorem 3.4.1. This primarily results from the freedom to

select the constants $\gamma_1, \dots, \gamma_N$ based only on the smoothness constants \mathbf{L} . Thus, in the case of non-uniform sampling, the dependencies from Theorem 3.4.1 (using the same historical regularization strategy) reduce in Theorem 3.5.1 from $\|\boldsymbol{\lambda}\|_1^{1/2} \|\mathbf{G}\|_1^{1/2}$ to $\|\boldsymbol{\lambda}\|_1^{1/2} \|\mathbf{G}\|_\infty^{1/2}$ and from $\|\boldsymbol{\lambda}\|_1$ to $\|\boldsymbol{\lambda}\|_1^{1/2} \|\mathbf{L}\|_1^{1/2}$.

Regarding per-iteration complexity, the exact arguments of Section 3.4.1 apply, except that the update defining \mathbf{y}_k occurs at cost $\tilde{O}(n/N)$ instead of $\tilde{O}(n)$. Thus, the total per-iteration complexity is $\tilde{O}(nd/N)$.

3.5.2 Strategy 2: Non-Uniform Block Replacement Probabilities

Like its counterpart Section 3.4.2, this section will rely on choosing the parameter $\mathbf{r} = (r_1, \dots, r_N)$, which governs the update probabilities of the primal table $\hat{\mathbf{x}}_{k,1}, \dots, \hat{\mathbf{x}}_{k,N}$. Accordingly, we may set $w_k^P = 0$ to simplify our gap bound to

$$\begin{aligned} \sum_{k=1}^t a_k \mathbb{E}[\text{Gap}^{u,v}(\mathbf{x}_k, \mathbf{y}_k)] &\leq \frac{1}{2} (\mathcal{T}_0^P - \mathbb{E}[\mathcal{T}_t^P]) + \frac{N}{2} (\mathcal{T}_0^D - \mathbb{E}[\mathcal{T}_t^D]) \\ &+ \sum_{k=1}^{t-1} \mathbb{E}[\mathcal{I}_k^P] - \frac{1}{2} \sum_{k=1}^t \mathbb{E}[\mathcal{C}_k^P + N\mathcal{C}_k^D] \end{aligned} \quad (3.85)$$

By this point, all arguments used in the analysis have been seen before, in that ideas related to separability were employed in Section 3.5.1 and ideas related to block replacement probabilities were employed in Section 3.4.2. Thus, the proofs are relatively short in this section. We cancel (3.85) in Proposition 3.5.2. The resulting complexity will depend on the sampling probabilities \mathbf{p} and \mathbf{r} through the constants

$$\mathbf{G}_{\mathbf{p}} := \sqrt{\max_{I \in [N]} \frac{\mathbf{G}_I^2}{p_I}} \text{ and } \mathbf{L}_{\mathbf{p}, \mathbf{r}} := \sqrt{\sum_{I=1}^N \frac{\mathbf{L}_I^2}{p_I r_I^2}},$$

as described in Proposition 3.5.2.

Proposition 3.5.2. *Let $(\mathbf{x}_0, \mathbf{y}_0) \in \text{ri}(\text{dom}(\phi)) \times \text{ri}(\text{dom}(\psi))$ and $\{(\mathbf{x}_k, \mathbf{y}_k)\}_{k \geq 1}$ be generated using $\bar{\mathbf{g}}_{k-1}$ and $\bar{\mathbf{f}}_{k-1/2}$ given by (3.73) and $f(\mathbf{x}_k)$, respectively. Define $a_1 = \min \left\{ \frac{\sqrt{\mu_0 \nu_0}}{5\sqrt{2}\mathbf{G}_{\mathbf{p}}}, \frac{\mu_0}{10\mathbf{L}_{\mathbf{p}, \mathbf{r}}} \right\}$*

and select $(a_k)_{k \geq 2}$ such that the conditions

$$a_k \leq \min \left\{ \frac{\sqrt{(A_k \mu + \mu_0)(A_{k-1} \nu + \nu_0)}}{5\sqrt{2}\mathbf{G}_p}, \frac{\sqrt{(A_k \mu + \mu_0)(A_{k-1} \mu + \mu_0)}}{10\mathbf{L}_{p,r}} \right\} \quad (3.86)$$

and $\frac{a_k^2}{A_k \mu + \mu_0} \leq \min_I(1 + r_I/5) \frac{a_{k-1}^2}{A_{k-1} \mu + \mu_0}$ hold. We have that for any $\mathbf{u} \in \mathcal{X}$ and $\mathbf{v} \in \mathcal{Y}$,

$$\sum_{k=1}^t a_k \mathbb{E}[\text{Gap}^{\mathbf{u}, \mathbf{v}}(\mathbf{x}_k, \mathbf{y}_k)] + \frac{1}{2} \mathbb{E}[\mathcal{T}_t^{\text{P}}] + \frac{N}{2} \mathbb{E}[\mathcal{T}_t^{\text{D}}] \leq \frac{1}{2} \mathcal{T}_0^{\text{P}} + \frac{N}{2} \mathcal{T}_0^{\text{D}}.$$

Proof. Mirroring the proof of Proposition 3.4.2, by Young's inequality with parameter $(A_{k-1} \mu + \mu_0)/2$, we have for $k = 2, \dots, t-1$ that

$$\begin{aligned} \mathbb{E}[\mathcal{E}_k^{\text{P}}] &= a_{k-1} \mathbb{E} \left\langle \frac{1}{p_{P_k}} \sum_{i \in B_{P_k}} (\bar{y}_{k-1,i} \nabla f_i(\mathbf{x}_{k-1}) - y_{k-2,i} \hat{\mathbf{g}}_{k-2,i}), \mathbf{x}_k - \mathbf{x}_{k-1} \right\rangle \\ &\leq \frac{1}{2} \mathbb{E}[\mathcal{C}_k^{\text{P}}] + \frac{a_{k-1}^2}{A_{k-1} \mu + \mu_0} \mathbb{E} \left\| \frac{1}{p_{P_k}} \sum_{i \in B_{P_k}} (\bar{y}_{k-1,i} \nabla f_i(\mathbf{x}_{k-1}) - y_{k-2,i} \hat{\mathbf{g}}_{k-2,i}) \right\|_{\mathcal{X}^*}^2 \end{aligned} \quad (3.87)$$

$$(3.88)$$

and

$$\begin{aligned} &- a_k \mathbb{E}[\langle \nabla f(\mathbf{x}_k)^\top \bar{\mathbf{y}}_k - \hat{\mathbf{g}}_{k-1}^\top \mathbf{y}_{k-1}, \mathbf{u} - \mathbf{x}_k \rangle] \\ &\leq \frac{1}{2} \mathbb{E}[\mathcal{T}_t^{\text{P}}] + \frac{a_{t-1}^2}{A_{t-1} \mu + \mu_0} \mathbb{E} \left\| \frac{1}{p_{P_t}} \sum_{i \in B_{P_t}} (\bar{y}_{t-1,i} \nabla f_i(\mathbf{x}_{t-1}) - y_{t-2,i} \hat{\mathbf{g}}_{t-2,i}) \right\|_{\mathcal{X}^*}^2. \end{aligned} \quad (3.89)$$

Apply Lemma 3.4.1 to achieve

$$\begin{aligned} &\mathbb{E} \left\| \frac{1}{p_{P_k}} \sum_{i \in B_{P_k}} (\bar{y}_{k-1,i} \nabla f_i(\mathbf{x}_{k-1}) - y_{k-2,i} \hat{\mathbf{g}}_{k-2,i}) \right\|_{\mathcal{X}^*}^2 \\ &\leq 2 \sum_{I=1}^N \frac{\mathbf{G}_I^2}{p_I} \mathbb{E} \|\bar{\mathbf{y}}_{k-1,I} - \mathbf{y}_{k-2,I}\|_2^2 + 2 \sum_{I=1}^N \frac{\mathbf{L}_I^2}{p_I} \mathbb{E} \|\mathbf{x}_{k-1} - \hat{\mathbf{x}}_{k-2,I}\|_{\mathcal{X}}^2, \end{aligned}$$

where the last line follows from Young's inequality. For the first term, we compute its

expectation using Lemma 3.5.1,

$$\begin{aligned}
2 \sum_{I=1}^N \frac{G_I^2}{p_I} \mathbb{E} \|\bar{\mathbf{y}}_{k-1,I} - \mathbf{y}_{k-2,I}\|_2^2 &\leq 2 \max_I \frac{G_I^2}{p_I} \mathbb{E} \|\bar{\mathbf{y}}_{k-1} - \mathbf{y}_{k-2}\|_2^2 \\
&\leq 2 \max_I \frac{G_I^2}{p_I} \mathbb{E} \|\bar{\mathbf{y}}_{k-1} - \mathbf{y}_{k-2}\|_{\mathcal{Y}}^2 \\
&= 2N \max_I \frac{G_I^2}{p_I} \mathbb{E} \|\mathbf{y}_{k-1} - \mathbf{y}_{k-2}\|_{\mathcal{Y}}^2 \\
&\leq \underbrace{4N \max_I \frac{G_I^2}{p_I}}_{G_p^2} \mathbb{E}[\Delta_{\mathcal{Y}}(\mathbf{y}_{k-1}, \mathbf{y}_{k-2})],
\end{aligned}$$

where we used $\|\cdot\|_2 \leq \|\cdot\|_{\mathcal{Y}}$ in the second inequality. For the second term, use the same argument leading to (3.66) in the proof of Proposition 3.4.2 to achieve

$$\begin{aligned}
2 \sum_{I=1}^N \frac{L_I^2}{p_I} \mathbb{E} \|\mathbf{x}_{k-1} - \hat{\mathbf{x}}_{k-2,I}\|_{\mathcal{X}}^2 &\leq 50 \left(\sum_{I=1}^N \frac{L_I^2}{p_I r_I^2} \right) \sum_{k'=1}^t \frac{a_{k'}^2}{A_{k'} \mu + \mu_0} \mathbb{E} \|\mathbf{x}_{k'} - \mathbf{x}_{k'-1}\|_{\mathcal{X}}^2 \\
&\leq \underbrace{100 \left(\sum_{I=1}^N \frac{L_I^2}{p_I r_I^2} \right)}_{L_{p,r}^2} \sum_{k'=1}^t \frac{a_{k'}^2}{A_{k'} \mu + \mu_0} \mathbb{E}[\Delta_{\mathcal{X}}(\mathbf{x}_{k'}, \mathbf{x}_{k'-1})].
\end{aligned}$$

Thus, under the conditions (3.86), it holds that

$$\sum_{k=1}^{t-1} \mathbb{E} [\mathcal{I}_k^{\mathcal{P}}] \leq \frac{1}{2} \sum_{k=1}^t \mathbb{E} [\mathcal{C}_k^{\mathcal{P}} + N \mathcal{C}_k^{\mathcal{D}}],$$

which completes the proof. \square

As in Section 3.5.1, because the following result follows the same argument as Theorem 3.4.2, the proof is omitted. We use the initial distance quantity $D_{0,N}$ from Theorem 3.5.1.

Theorem 3.5.2. *Under Assumption 3.2.2 and Assumption 3.2.1, consider any $\mathbf{u} \in \mathcal{X}$, $\mathbf{v} \in \mathcal{Y}$ and precision $\varepsilon > 0$. There exists a choice of the sequences $(a_k)_{k=1}^t$ and $(w_k^{\mathcal{P}})_{k=1}^t$ such that Algorithm 4 with Identity Card 3 produces an output point $(\tilde{\mathbf{x}}_t, \tilde{\mathbf{y}}_t) \in \mathcal{X} \times \mathcal{Y}$ satisfying $\mathbb{E} [\text{Gap}^{\mathbf{u}, \mathbf{v}}(\tilde{\mathbf{x}}_t, \tilde{\mathbf{y}}_t)] \leq \varepsilon$ for t that depends on ε according to the following iteration complexities. Let $C_{0,N} := (L_{p,r} \sqrt{\nu_0/\mu_0} + G_p) D_{0,N}$, for $D_{0,N}$ defined in (3.84).*

Case	Iteration Complexity
$\mu > 0$ and $\nu > 0$	$O\left(\left(N + \frac{\mathbf{L}_{p,r}}{\mu} + \frac{\mathbf{G}_p}{\sqrt{\mu\nu}}\right) \ln\left(\frac{C_{0,N}}{\varepsilon}\right)\right)$
$\mu > 0$ and $\nu = 0$	$O\left(\left(N + \frac{\mathbf{L}_{p,r}}{\mu}\right) \ln\left(\frac{C_{0,N}}{\varepsilon}\right) + \mathbf{G}_p \sqrt{\frac{\sqrt{\mu_0/\nu_0} D_{0,N}}{\mu\varepsilon}}\right)$
$\mu = 0$ and $\nu > 0$	$O\left(N \ln\left(\frac{C_{0,N}}{\varepsilon}\right) + \frac{\mathbf{L}_{p,r} \sqrt{\nu_0/\mu_0} D_{0,N}}{\varepsilon} + \mathbf{G}_p \sqrt{\frac{\sqrt{\nu_0/\mu_0} D_{0,N}}{\nu\varepsilon}}\right)$
$\mu = 0$ and $\nu = 0$	$O\left(N \ln\left(\frac{C_{0,N}}{\varepsilon}\right) + \frac{(\mathbf{L}_{p,r} \sqrt{\nu_0/\mu_0} + \mathbf{G}_p) D_{0,N}}{\varepsilon}\right)$

As with Theorem 3.5.1, by setting $\nu_0 \sim \mu_0/N$, we interpret $D_{0,N}$ as \sqrt{N} times a constant. The non-uniform sampling complexity given below follows by letting $p_I \propto \boldsymbol{\lambda}_I^{1/2}$ and $r_I \propto \mathbf{L}_I^{1/2}$.

Constant	Uniform Sampling	Non-Uniform Sampling
\mathbf{G}_p	$\sqrt{N} \ \mathbf{G}\ _\infty$	$\ \boldsymbol{\lambda}\ _{1/2}^{1/4} \ \mathbf{G}\ _\infty^{3/4}$
$\mathbf{L}_{p,r}$	$N^{3/2} \ \mathbf{L}\ _2$	$\ \boldsymbol{\lambda}\ _{1/2}^{1/4} \ \mathbf{L}\ _{1/2}^{3/4}$

We compare the resulting complexity to the non-separable analog in Theorem 3.4.2. There are improvements both in the dependence on \mathbf{G} and \mathbf{L} when the sampling scheme can be tuned. Theorem 3.5.2 improves the dependence on \mathbf{G} from $\|\boldsymbol{\lambda}\|_{1/2}^{1/2} \|\mathbf{G}\|_\infty^{1/4}$ to $\|\boldsymbol{\lambda}\|_{1/2}^{1/4} \|\mathbf{G}\|_\infty^{3/4}$. As for the dependence on \mathbf{L} , this improves from $\|\boldsymbol{\lambda}\|_{1/2}^{3/4} \|\mathbf{L}\|_{1/2}^{1/4}$ to $\|\boldsymbol{\lambda}\|_{1/2}^{1/4} \|\mathbf{L}\|_{1/2}^{3/4}$. The mechanism is analogous to the improvement from Theorem 3.4.1 to Theorem 3.5.1; because $\mathcal{I}_k^D = 0$, the constants do not have to adapt in order to control an error term of the form \mathcal{E}_k^D .

Finally, on per-iteration complexity, the exact arguments of Section 3.4.2 apply, except that the update defining \mathbf{y}_k may now occur at cost $\tilde{O}(n/N)$ instead of $\tilde{O}(n)$. Thus, the total per-iteration complexity is $\tilde{O}(nd/N)$.

3.6 Discussion & Comparisons

Our discussion covers internal comparisons between full vector methods and stochastic methods, as well as external comparisons to contemporary methods for solving saddle point and

variational inequality problems. We briefly comment on our chosen convergence criterion, as it may differ slightly from those used in comparisons.

3.6.1 Stronger Convergence Criteria

The results discussed in this section are expressed in terms of *global complexity* or *arithmetic complexity*, which is computed by multiplying the number of iterations shown in Theorem 3.3.1 to 3.5.2 by $\tilde{O}(nd)$ for full vector update methods; for block-wise methods with N blocks of size n/N , we use $\tilde{O}(n(d/N + 1))$ for full updates of \mathbf{y}_k and $\tilde{O}(n(d/N))$ for partial updates of \mathbf{y}_k . In order to effectively compare to methods using our block coordinate-wise Lipschitz and smoothness constant from Section 3.2, we will apply a particular finite sum decomposition (see (3.93)) that will allow us to directly apply methods from the finite sum variational inequality literature.

Importantly, in the case of randomized algorithms, the complexity in terms of the number of iterations is itself determined by the number t such that the algorithm may output a point $(\tilde{\mathbf{x}}_t, \tilde{\mathbf{y}}_t)$ satisfying $\mathbb{E}[\text{Gap}^{\mathbf{u}, \mathbf{v}}(\tilde{\mathbf{x}}_t, \tilde{\mathbf{y}}_t)] \leq \varepsilon$, where the expectation is taken over all algorithm randomness with $\mathbf{u} \in \mathcal{X}$ and $\mathbf{v} \in \mathcal{Y}$ fixed (i.e., (\mathbf{u}, \mathbf{v}) is independent of the algorithm randomness). When $\mu > 0$, the criterion is made meaningful by setting $\mathbf{u} = \mathbf{x}_\star$ as the unique minimizer of the strongly convex objective $\mathbf{x} \mapsto \max_{\mathbf{y} \in \mathcal{Y}} \mathcal{L}(\mathbf{x}, \mathbf{y})$. Otherwise, we choose a compact set $\mathcal{U} \subseteq \mathcal{X}$ and consider $\sup_{\mathbf{u} \in \mathcal{U}} \mathbb{E}[\text{Gap}^{\mathbf{u}, \mathbf{v}}(\tilde{\mathbf{x}}_t, \tilde{\mathbf{y}}_t)] \leq \varepsilon$, where \mathbf{v} is replaced by the unique maximizer \mathbf{y}_\star when $\nu > 0$ or another supremum is taken over $\mathbf{v} \in \mathcal{V} \subseteq \mathcal{Y}$ for \mathcal{V} compact otherwise. As described in Alacaoglu et al. [2022, Example 1], the “supremum of expected gap” criterion is weaker than the “expected supremum of gap” criterion, as algorithms with divergent behavior can still converge according to the first criterion. We render guarantees for the stronger criterion as a technical detail, as in light of previous work, largely similar steps can be applied to achieve the same complexity guarantee for the expected supremum of gap. To not overcomplicate the proofs, we only highlight the parts of the analysis that change. Consider the argument used to derive (3.41), in which the expectation is applied to the terms $a_k \langle \nabla f(\mathbf{x}_k)^\top \mathbf{y}_k - \hat{\mathbf{g}}_{k-1}^\top \hat{\mathbf{y}}_{k-1}, \mathbf{u} - \mathbf{x}_k \rangle$, after which they telescope. If the supremum

is to be taken, we can no longer apply the expectation to these terms directly. Instead,

$$\begin{aligned}
a_k \langle \nabla f(\mathbf{x}_k)^\top \mathbf{y}_k - \bar{\mathbf{g}}_{k-1}, \mathbf{u} - \mathbf{x}_k \rangle &= a_k \langle \nabla f(\mathbf{x}_k)^\top \mathbf{y}_k - \hat{\mathbf{g}}_{k-1}^\top \hat{\mathbf{y}}_{k-1}, \mathbf{u} - \mathbf{x}_k \rangle \\
&\quad - \underbrace{\left\langle \frac{a_{k-1}}{p_{P_k}} \sum_{i \in B_{P_k}} (y_{k-1,i} \nabla f_i(\mathbf{x}_{k-1}) - \hat{y}_{k-2,i} \hat{\mathbf{g}}_{k-2,i}), \mathbf{u} - \mathbf{x}_{k-1} \right\rangle}_{\mathbf{z}_k} \\
&\quad - \underbrace{a_{k-1} \left\langle \frac{1}{p_{P_k}} \sum_{i \in B_{P_k}} (y_{k-1,i} \nabla f_i(\mathbf{x}_{k-1}) - \hat{y}_{k-2,i} \hat{\mathbf{g}}_{k-2,i}), \mathbf{x}_{k-1} - \mathbf{x}_k \right\rangle}_{\mathcal{E}_k^P},
\end{aligned}$$

and by taking the expectation over P_k yields the identity

$$\begin{aligned}
a_k \langle \nabla f(\mathbf{x}_k)^\top \mathbf{y}_k - \bar{\mathbf{g}}_{k-1}, \mathbf{u} - \mathbf{x}_k \rangle &= a_k \langle \nabla f(\mathbf{x}_k)^\top \mathbf{y}_k - \hat{\mathbf{g}}_{k-1}^\top \hat{\mathbf{y}}_{k-1}, \mathbf{u} - \mathbf{x}_k \rangle \\
&\quad - a_{k-1} \langle \nabla f(\mathbf{x}_{k-1})^\top \mathbf{y}_{k-1} - \hat{\mathbf{g}}_{k-2}^\top \hat{\mathbf{y}}_{k-2}, \mathbf{u} - \mathbf{x}_{k-1} \rangle \\
&\quad - \langle \mathbf{z}_k - \mathbb{E}_{k-1}[\mathbf{z}_k], \mathbf{u} \rangle - \mathcal{E}_k^P.
\end{aligned} \tag{3.90}$$

The familiar term \mathcal{E}_k^P does not depend on \mathbf{u} and can be bounded using the same techniques as in the proofs of Theorem 3.4.1 and Theorem 3.4.2, whereas $\mathbf{z}_k - \mathbb{E}_{k-1}[\mathbf{z}_k]$ is zero-mean conditional on \mathcal{F}_{k-1} (i.e., a martingale difference sequence), but is not necessarily independent of \mathbf{u} . Next, we may apply Diakonikolas [2025, Lemma 4] (adapted from Alacaoglu and Malitsky [2022, Lemma 3.5]) to achieve

$$\begin{aligned}
-\mathbb{E} \left[\sum_{k=1}^t \langle \mathbf{z}_k - \mathbb{E}_{k-1}[\mathbf{z}_k], \mathbf{u}_\star \rangle \right] &\leq \mathbb{E}[\Delta_{\mathcal{X}}(\mathbf{u}_\star, \mathbf{x}_0)] + \frac{1}{2} \sum_{k=1}^t \mathbb{E} \|\mathbf{z}_k - \mathbb{E}_{k-1}[\mathbf{z}_k]\|_2^2 \\
&\leq \mathbb{E}[\Delta_{\mathcal{X}}(\mathbf{u}_\star, \mathbf{x}_0)] + \frac{1}{2} \sum_{k=1}^t \mathbb{E} \|\mathbf{z}_k\|_2^2,
\end{aligned}$$

where \mathbf{u}_\star is the element of $\mathbf{u} \in \mathcal{U}$ that achieves the supremum in the gap criterion (recalling that \mathcal{U} is chosen to be compact). The term $\mathbb{E}[\Delta_{\mathcal{X}}(\mathbf{u}_\star, \mathbf{x}_0)]$ is upper bounded by a constant, whereas the $\mathbb{E} \|\mathbf{z}_k\|_2^2$ terms are bounded using the exact same techniques used to bound the \mathcal{E}_k^P terms. We choose to describe the argument in the manner above (as opposed to including it formally in the proofs) as it changes neither the other technical ideas nor the resulting complexities; we comment on this subtlety for the sake of completeness.

Algorithm Type	Global Complexity (big-O)
Full vector (Theorem 3.3.1)	$nd \left(\frac{L}{\mu} + \frac{G}{\sqrt{\mu\nu}} \right) \ln \left(\frac{1}{\varepsilon} \right)$
Stochastic (Theorem 3.4.1)	$\left(\frac{nd}{N} + n \right) \left(N + \frac{\sqrt{N}\ \boldsymbol{\lambda}\ _1}{\mu} + \frac{\sqrt{N}\ \boldsymbol{\lambda}\ _1^{1/2}\ \mathbf{G}\ _1^{1/2}}{\sqrt{\mu\nu}} \right) \ln \left(\frac{1}{\varepsilon} \right)$
Stochastic (Theorem 3.4.2)	$\left(\frac{nd}{N} + n \right) \left(N + \frac{\ \boldsymbol{\lambda}\ _{1/2}^{3/4}\ \mathbf{L}\ _{1/2}^{1/4}}{\mu} + \frac{\ \boldsymbol{\lambda}\ _{1/2}^{1/2}\ \mathbf{G}\ _{1/2}^{1/2}}{\sqrt{\mu\nu}} \right) \ln \left(\frac{1}{\varepsilon} \right)$
Block Coordinate-wise (Theorem 3.5.1)	$\frac{nd}{N} \left(N + \frac{\sqrt{N}\ \boldsymbol{\lambda}\ _1^{1/2}\ \mathbf{L}\ _1^{1/2}}{\mu} + \frac{\sqrt{N}\ \boldsymbol{\lambda}\ _1^{1/2}\ \mathbf{G}\ _\infty^{1/2}}{\sqrt{\mu\nu}} \right) \ln \left(\frac{1}{\varepsilon} \right)$
Block Coordinate-wise (Theorem 3.5.2)	$\frac{nd}{N} \left(N + \frac{\ \boldsymbol{\lambda}\ _{1/2}^{1/4}\ \mathbf{L}\ _{1/2}^{3/4}}{\mu} + \frac{\ \boldsymbol{\lambda}\ _{1/2}^{1/4}\ \mathbf{G}\ _\infty^{3/4}}{\sqrt{\mu\nu}} \right) \ln \left(\frac{1}{\varepsilon} \right)$

Table 3.1: **Complexity Bounds for Full Vector and Stochastic Methods for the case $\mu, \nu > 0$.** Arithmetic or global complexity (i.e., the total number of elementary operations required to compute (\mathbf{x}, \mathbf{y}) satisfying $\mathbb{E} [\text{Gap}^{\mathbf{u}, \mathbf{v}}(\mathbf{x}, \mathbf{y}) \leq \varepsilon]$ for fixed $(\mathbf{u}, \mathbf{v}) \in \mathcal{X} \times \mathcal{Y}$, with the expectation taken over all algorithmic randomness.

3.6.2 Full Vector Update versus Stochastic Methods

To compare methods both within this chapter and alternatives for solving nonbilinearly coupled min-max problems, we first state some relationships between the constants introduced in Assumption 3.2.2. We then proceed with fine-grained comparisons to alternatives in the existing literature on min-max optimization and monotone variational inequalities. To simplify some comparisons, we assume that $L/\mu \geq 1$ and $G/\sqrt{\mu\nu} \geq 1$, so that they may be interpreted as “primal” and “mixed” condition numbers, respectively. First, observe that by the triangle inequality, in terms of the constants $\mathbf{G} = (\mathbf{G}_1, \dots, \mathbf{G}_N)$ and $\mathbf{L} = (\mathbf{L}_1, \dots, \mathbf{L}_N)$, the constants G and L can be upper bounded as

$$G \leq \|\mathbf{G}\|_1 \text{ and } L \leq \|\mathbf{L}\|_1.$$

Comparing the complexities for $\mu, \nu > 0$ in Table 3.1, the findings are summarized as follows.

Remark 3.6.1. *We observe the following in Table 3.1.*

- *The strongly convex case highlights a limitation in the historical regularization method*

(Theorem 3.4.1 and Theorem 3.5.1) as a factor \sqrt{N} is gained in terms of the dependence on \mathbf{G} and \mathbf{L} , which is not observed for the full vector method from Theorem 3.3.1. When comparing Theorem 3.4.1 to Theorem 3.3.1, the dependence on the smoothness constants renders as L versus $\|\boldsymbol{\lambda}\|_1 / \sqrt{N}$ and on the Lipschitz constants as G versus $\|\boldsymbol{\lambda}\|_1^{1/2} \|\mathbf{G}\|_1^{1/2} / \sqrt{N}$. In both cases, when the constants are highly non-uniform, we are still afforded an up to \sqrt{n} improvement in complexity from the stochastic method. Note that Theorem 3.5.1 is a strict improvement over Theorem 3.4.1 due to separability.

- On the other hand, in the highly non-uniform setting, we may gain an up to d factor (resp. n factor) of improvement in terms of complexity using the methods of Theorem 3.4.2 (resp. Theorem 3.5.2) over the full vector method. The results of Theorem 3.4.1 and Theorem 3.4.2 (and by analogy, the results of Theorem 3.5.1 and Theorem 3.5.2) do not have a uniformly dominating method, as the extra factor of \sqrt{N} may be on par with the improvement of the $\|\cdot\|_1$ norm over the $\|\cdot\|_{1/2}$ in terms of the dependence on $(\mathbf{G}, \mathbf{L}, \boldsymbol{\lambda})$.

3.6.3 Alternative Methods for Min-Max and Variational Inequality Problems

In some comparisons, we must access the smoothness constants of the function $(\mathbf{x}, \mathbf{y}) \mapsto \mathcal{L}(\mathbf{x}, \mathbf{y})$ both with respect to \mathbf{x} and \mathbf{y} separately, and additionally, with respect to the pair (\mathbf{x}, \mathbf{y}) . For the sake of comparison, assume that ϕ and ψ are differentiable and

$$\|\nabla\phi(\mathbf{x}) - \nabla\phi(\mathbf{x}')\|_{\mathbf{x}^*} \leq \mu \|\mathbf{x} - \mathbf{x}'\|_{\mathbf{x}} \text{ and } \|\nabla\psi(\mathbf{y}) - \nabla\psi(\mathbf{y}')\|_{\mathbf{y}^*} \leq \nu \|\mathbf{y} - \mathbf{y}'\|_{\mathbf{y}}.$$

Then, the following relations hold:

$$\begin{aligned}
\sup_{\mathbf{y} \in \mathcal{Y}} \|\nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \mathbf{y}) - \nabla_{\mathbf{x}'} \mathcal{L}(\mathbf{x}', \mathbf{y})\|_{\mathcal{X}^*} &\leq (L + \mu) \|\mathbf{x} - \mathbf{x}'\|_{\mathcal{X}}, \\
\sup_{\mathbf{x} \in \mathcal{X}} \|\nabla_{\mathbf{y}} \mathcal{L}(\mathbf{x}, \mathbf{y}) - \nabla_{\mathbf{y}'} \mathcal{L}(\mathbf{x}, \mathbf{y}')\|_{\mathcal{X}^*} &\leq \nu \|\mathbf{y} - \mathbf{y}'\|_{\mathcal{Y}}, \\
\sup_{\mathbf{y} \in \mathcal{Y}} \|\nabla_{\mathbf{y}} \mathcal{L}(\mathbf{x}, \mathbf{y}) - \nabla_{\mathbf{y}} \mathcal{L}(\mathbf{x}', \mathbf{y})\|_{\mathcal{X}^*} &\leq G \|\mathbf{x} - \mathbf{x}'\|_{\mathcal{X}}, \\
\sup_{\mathbf{x} \in \mathcal{X}} \|\nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \mathbf{y}) - \nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \mathbf{y}')\|_{\mathcal{X}^*} &\leq G \|\mathbf{y} - \mathbf{y}'\|_{\mathcal{Y}}.
\end{aligned} \tag{3.91}$$

As is done in the case of the ℓ_2 -norm, define the norms

$$\|(\mathbf{x}, \mathbf{y})\|^2 := \|\mathbf{x}\|_{\mathcal{X}}^2 + \|\mathbf{y}\|_{\mathcal{Y}}^2 \quad \text{and} \quad \|\nabla \mathcal{L}(\mathbf{x}, \mathbf{y})\|_*^2 := \|\nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}, \mathbf{y})\|_{\mathcal{X}^*}^2 + \|\nabla_{\mathbf{y}} \mathcal{L}(\mathbf{x}, \mathbf{y})\|_{\mathcal{Y}^*}^2.$$

Then, we finally have that

$$\|\nabla \mathcal{L}(\mathbf{x}, \mathbf{y}) - \nabla \mathcal{L}(\mathbf{x}', \mathbf{y}')\|_* \leq \sqrt{3 \max\{L^2 + G^2 + \mu^2, G^2 + \nu^2\}} \|(\mathbf{x}, \mathbf{y}) - (\mathbf{x}', \mathbf{y}')\|. \tag{3.92}$$

Remark 3.6.2. *We do not assume that ϕ and ψ are smooth in our convergence guarantees, as the individual components ϕ and ψ are allowed to be non-differentiable in our framework. The additional assumptions above are made only for the sake of comparison. To this end, we will assume that $\nu \leq G$ (and maintain $\mu \leq L$). Then the constants in (3.91) and (3.92) can be simplified to absolute constants times L and $\sqrt{L^2 + G^2}$, respectively. Thus, we define the vector field $(\mathbf{x}, \mathbf{y}) \mapsto \nabla \mathcal{L}(\mathbf{x}, \mathbf{y})$, which is $C\sqrt{L^2 + G^2}$ -Lipschitz (for an absolute constant $C > 0$) and $(\mu \wedge \nu)$ -strongly monotone.*

Given the above, we use these constants to compare to algorithms designed for solving variational inequality problems. In the upcoming remarks, we discuss the convex-concave (monotone) and strongly convex-strongly concave (strongly monotone) settings. First, consider the case in which $\mu \wedge \nu > 0$, so that we may observe the dependence on all problem constants in the full vector update settings.

Remark 3.6.3. *We observe the following in Table 3.2.*

- *When using classical methods for monotone variational inequalities, such as dual extrapolation [Nesterov and Scrimali, 2006], we highlight the $(\mu \wedge \nu)$ term, which may suffer when there is a large asymmetry in the strong convexity and strong concavity constants. As emphasized in Section 3.1, ν is often chosen as a small approximation or smoothing parameter, meaning that the theoretical complexity will not necessarily improve for large values of the primal strong convexity constant μ . Observe that all other methods will improve with increasing μ .*
- *In recent works [Jin et al., 2022, Li et al., 2023] and ours, the Lipschitz constants from different components of the objective function are separated in the complexity. In particular, the constants L and ν are decoupled, which is especially helpful in scenarios in which L is much larger than G . Examples include losses that are themselves smooth approximations of non-smooth losses, such as Huber approximations of the mean absolute error function. We improve over Jin et al. [2022], Li et al. [2023], which are designed for general nonbilinearly-coupled objectives, by a logarithmic factor in iteration complexity and achieve the same result in global complexity.*

When viewed as a saddle point or variational inequality problem, notice that (3.2) has a finite sum structure. In order to make direct comparisons, we decompose the objective block-wise, that is, $\mathcal{L}(\mathbf{x}, \mathbf{y}) = \sum_{j=1}^N \mathcal{L}_J(\mathbf{x}, \mathbf{y})$, where

$$\mathcal{L}_J(\mathbf{x}, \mathbf{y}) = \begin{cases} \sum_{j \in B_J} y_j f_j(\mathbf{x}) - \psi_J(\mathbf{y}_J) + \frac{1}{N} \phi(\mathbf{x}) & \text{if Definition 3.2.1 (separability) is satisfied} \\ \sum_{j \in B_J} y_j f_j(\mathbf{x}) - \frac{1}{N} \psi(\mathbf{y}) + \frac{1}{N} \phi(\mathbf{x}) & \text{otherwise} \end{cases} \quad (3.93)$$

Thus, when comparing to methods designed for finite sum objectives, we may consider the overall complexity of querying the oracle $(\mathcal{L}_J, \nabla \mathcal{L}_J)$ to be $O(nd/N)$ if the objective is separable or $O(n(d/N + 1))$ if it is non-separable. Even if $(\mathcal{L}_J, \nabla \mathcal{L}_J)$ is of cost $O(nd/N)$ to

compute, each oracle call may be associated with a single step of the algorithm, which is still $O(n)$ if it updates the primal and dual variables in their entirety. This is the first consideration when computing the global complexities in Table 3.3 and Table 3.4. For the second consideration, we compute the Lipschitzness and smoothness assumption of the individual component functions *on average* with uniform or non-uniform sampling, as is commonly used in analyses of methods for sum-decomposable objectives. Many contemporary results are stated in terms of “on average” smoothness (for min-max problems) and Lipschitzness (for variational inequality problems). Using the same norms defined in (3.92), we say that $\mathcal{L}_1, \dots, \mathcal{L}_N$ are L_{avg} -smooth on average according to sampling weights $\mathbf{p} = (p_1, \dots, p_N)$ if

$$\mathbb{E}_{J \sim \mathbf{p}} \|(1/p_J) (\nabla \mathcal{L}_J(\mathbf{x}, \mathbf{y}) - \nabla \mathcal{L}_J(\mathbf{x}', \mathbf{y}'))\|_*^2 \leq L_{\text{avg}}^2 \|(\mathbf{x}, \mathbf{y}) - (\mathbf{x}', \mathbf{y}')\|^2.$$

Recall the constants $\boldsymbol{\lambda} = (\boldsymbol{\lambda}_1, \dots, \boldsymbol{\lambda}_N)$, where $\boldsymbol{\lambda}_J = \sqrt{\mathbf{G}_J^2 + \mathbf{L}_J^2}$ are the Lipschitz constants of each \mathcal{L}_J with respect to the norm $\|\cdot\|$ defined above. The prototypical sampling schemes are the uniform and importance-weighted schemes

$$\begin{aligned} \mathbb{E}_{J \sim \text{unif}[N]} \|n (\nabla \mathcal{L}_J(\mathbf{x}, \mathbf{y}) - \nabla \mathcal{L}_J(\mathbf{x}', \mathbf{y}'))\|_*^2 &\leq \underbrace{N \|\boldsymbol{\lambda}\|_2^2}_{\lambda_{\text{unif}}^2} \|(\mathbf{x}, \mathbf{y}) - (\mathbf{x}', \mathbf{y}')\|^2 \\ \mathbb{E}_{J \sim \boldsymbol{\lambda}} \|(1/\boldsymbol{\lambda}_J) (\nabla \mathcal{L}_J(\mathbf{x}, \mathbf{y}) - \nabla \mathcal{L}_J(\mathbf{x}', \mathbf{y}'))\|_*^2 &\leq \underbrace{\|\boldsymbol{\lambda}\|_1^2}_{\lambda_{\text{imp}}^2} \|(\mathbf{x}, \mathbf{y}) - (\mathbf{x}', \mathbf{y}')\|^2. \end{aligned}$$

Note that under the same sampling scheme, L_{avg} -smoothness on average is also implied by $(1/L_{\text{avg}})$ -cocoercivity on average [Cai et al., 2024, Assumption 3]. Finally, for $\varepsilon \searrow 0$, we have that $O\left(N + \sqrt{N} \lambda_{\text{unif}}/\varepsilon\right) = O\left(\sqrt{N} \lambda_{\text{unif}}/\varepsilon\right)$ and $O\left(N + \sqrt{N} \lambda_{\text{imp}}/\varepsilon\right) = O\left(\sqrt{N} \lambda_{\text{imp}}/\varepsilon\right)$ which, combined with the oracle cost for (3.93), leads to the results in Table 3.3 and Table 3.4.

Remark 3.6.4. We observe the following in Table 3.3 and Table 3.4.

- For the improved dependence on the problem constants (\mathbf{G}, \mathbf{L}) and strong convexity constants (μ, ν) there are two main themes, which both involve “decoupling” of the two

constants. Indeed, all results except for ours depend only on the aggregate Lipschitz constants $\boldsymbol{\lambda}$ and the minimum of the strong convexity constants $(\mu \wedge \nu)$, whereas we may (partially) separate these into dependences on “ \mathbf{L} over μ ” terms and “ \mathbf{G} over $\sqrt{\mu\nu}$ ” terms.

- In Table 3.3, fixing (n, d, ε) and without considering the differences between the aggregate constants $\boldsymbol{\lambda}$ and decoupled constants (\mathbf{G}, \mathbf{L}) , we notice a dependence on the ℓ_2 -norm in the results of [Alacaoglu and Malitsky \[2022\]](#), [Cai et al. \[2024\]](#), and [Pichugin et al. \[2024\]](#), as opposed to the $(1/N)$ times the ℓ_1 -norm dependence in Theorem 3.4.1. Because $N^{-1} \|\cdot\|_1 \leq N^{-1/2} \|\cdot\|_2$, we observe a \sqrt{N} improvement in complexity. Comparing the analogous results in Table 3.4, the complexity result of Theorem 3.4.1 scales as $N^{-1/2} \|\cdot\|_1$ improves over, but may still be on par with, the $\|\cdot\|_2$ scaling of [Alacaoglu and Malitsky \[2022\]](#) and [Cai et al. \[2024\]](#).
- In comparison to the best result of [Diakonikolas \[2025\]](#) in Table 3.3 for the non-separable case, the improvement of Theorem 3.4.1 comes from the use of the ℓ_1 -norm over the $\ell_{1/2}$ -norm, which may be up to \sqrt{N} smaller.

We also mention the work of [Boob and Khalafi \[2024\]](#), which solves a functionally constrained variational inequality formulation akin to Example 3 from Section 3.1. They operate under a completely different set of assumptions, largely to handle the possible unboundedness of the domain \mathcal{Y} of the Lagrange multipliers. Furthermore, the gradient operator and the functional constraints satisfy non-standard deviation control (as opposed to Lipschitzness) inequalities (see [Boob and Khalafi \[2024\]](#), Eq. (1.3) and (1.4)), and thus this work is not directly comparable to ours.

3.6.4 Bilinearly Coupled Problems

While originally motivated by nonbilinearly-coupled min-max problems, the bilinearly-coupled setting constitutes an important special case of (3.2), defined via $f(\mathbf{x}) = \mathbf{A}\mathbf{x}$ for $\mathbf{A} \in \mathbb{R}^{n \times d}$.

For the sake of discussion, we consider the convex-concave case (which includes matrix games for which $\phi \equiv 0$ and $\psi \equiv 0$). Randomized algorithms (such as randomized mirror-prox) for such problems were explored in Juditsky et al. [2011, 2013] to achieve complexity guarantees of the form $O(\sigma^2 \varepsilon^{-2} + \|\mathbf{A}\|_{\mathcal{X}, \mathcal{Y}^*} \varepsilon^{-1})$, where σ^2 is a measurement of noise arising from using stochastic estimates of matrix-vector multiplications and $\|\mathbf{A}\|_{\mathcal{X}, \mathcal{Y}^*} = \sup \{\|\mathbf{A}\mathbf{x}\|_{\mathcal{Y}^*} : \|\mathbf{x}\|_{\mathcal{X}} = 1\}$ is the induced matrix norm. For general problems, the approaches of Chambolle and Pock [2011], Alacaoglu and Malitsky [2022] achieve

$$\min \left\{ O \left(nd \|\mathbf{A}\|_{2,2} \varepsilon^{-1} \right), O \left(nd + \sqrt{nd(n+d) \|\mathbf{A}\|_{\text{Fro}} \varepsilon^{-1}} \right) \right\}$$

whereas these are reduced by Song et al. [2021] and Alacaoglu et al. [2022]² to

$$O \left(nd + d \frac{n \max_j \|\mathbf{A}_{j\cdot}\|_2}{\varepsilon} \right), \quad (3.94)$$

when the objective function is *separable*, where \mathbf{A}_i denotes the i -th row of \mathbf{A} . In our notation, $G_i = \|\mathbf{A}_i\|_\infty$ when we equip \mathcal{X} with the ℓ_1 -norm. Thus, due to separability, we may apply Theorem 3.5.1 with $N = n$ to achieve a global complexity of

$$O \left(nd + d \frac{\sqrt{(\sum_{i=1}^n \|\mathbf{A}_i\|_\infty) \cdot (n \max_j \|\mathbf{A}_{j\cdot}\|_\infty)}}{\varepsilon} \right). \quad (3.95)$$

Because

$$n \max_j \|\mathbf{A}_{j\cdot}\|_2 \stackrel{(\circ)}{\geq} n \max_j \|\mathbf{A}_{j\cdot}\|_\infty \stackrel{(*)}{\geq} \sum_{i=1}^n \|\mathbf{A}_i\|_\infty,$$

and (\circ) offers an up to a \sqrt{d} -factor improvement and $(*)$ offers an up to n -factor improvement, our result can offer up to an order- \sqrt{nd} improvement overall. This improvement is realized when within-row entries are highly uniform and within-column entries are highly non-uniform, leading to highly non-uniform infinity norms of the rows. For linearly constrained problems, Alacaoglu et al. [2022] achieve $O(nd + d \sum_{i=1}^n \|\mathbf{A}_i\|_2 / \varepsilon)$. It is relevant

²Guarantees might be for the expected supremum of gap or supremum of expected gap. We compare them side-by-side due to the argument in the earlier part of the section.

to note that our improvement relies heavily on the non-uniform sampling and non-uniform historical regularization applied in Section 3.4.1 and Section 3.5.1. Non-uniform sampling has been applied in the works above, as well as earlier works such as Alacaoglu et al. [2017] and Chambolle et al. [2018], but our findings show that sampling strategies may not be sufficient to remove the extraneous dimension factors in the iteration complexity.

We also consider a particular application to (non-separably constrained) matrix games, that is, the problem

$$\min_{\mathbf{x} \in \Delta^{d-1}} \max_{\mathbf{y} \in \Delta^{n-1}} [\mathcal{L}(\mathbf{x}, \mathbf{y}) := \mathbf{y}^\top \mathbf{A} \mathbf{x}],$$

for a matrix $\mathbf{A} \in \mathbb{R}^{n \times d}$ and probability simplices Δ^{d-1} and Δ^{n-1} , meaning that $\phi \equiv 0$ and $\psi \equiv 0$. Let A_{ij} denote the element (i, j) -th entry of \mathbf{A} and let $\text{nnz}(\mathbf{A})$ be the number of such entries that are non-zero. For this example problem, the results of Carmon et al. [2019, Thm. 1] and Alacaoglu and Malitsky [2022, Coro. 9 & Pg. 37] achieve a global complexity of

$$\tilde{O} \left(\text{nnz}(\mathbf{A}) + \sqrt{\text{nnz}(\mathbf{A}) \cdot (n + d)^{\frac{\max_{i,j} |A_{ij}|}{\varepsilon}}} \right) \quad (3.96)$$

whereas the method of Diakonikolas [2025, Thm. 1 & Pg. 19] achieves

$$\tilde{O} \left(\text{nnz}(\mathbf{A}) + (n + d) \frac{(\sum_i \|\mathbf{A}_i\|_\infty^{2/3} + \sum_j \|\mathbf{A}_{\cdot j}\|_\infty^{2/3})^{3/2}}{\varepsilon} \right), \quad (3.97)$$

where $\mathbf{A}_{\cdot j}$ is the j -th column of \mathbf{A} and \mathbf{A}_i is defined as in (3.94). Due to the non-separability of Δ^{n-1} , we may apply the result of Theorem 3.4.1 with $N = n$ and $G_i = \|\mathbf{A}_i\|_\infty$ to achieve a global complexity of

$$\tilde{O} \left(\text{nnz}(\mathbf{A}) + \frac{(n + d) \sum_i \|\mathbf{A}_i\|_\infty}{\varepsilon} \right).$$

Notice that this is a direct improvement over (3.97) by replacing the $\ell_{2/3}$ -norms with the ℓ_1 -norm. On the other hand, for cases when \mathbf{A} is dense but if the infinity norms of the rows are highly non-uniform, the complexity above will improve over (3.96) by a factor of

$\sqrt{\text{nnz}(\mathbf{A})/(n+d)}$ (which is order \sqrt{n} for dense square matrices). If these row norms are not highly non-uniform, but the matrix is still dense, then our complexity will be worse than (3.96) by an $n \cdot \sqrt{(n+d)/\text{nnz}(\mathbf{A})}$ -factor (which is also order- \sqrt{n} for dense square matrices).

3.7 Possible Extensions

3.7.1 Certificates of Suboptimality

The goal of this section is to provide an *online accuracy certificate* (or simply *certificate*) for the dual-linear min-max problem (3.2). In other words, we wish to find a continuous function of the primal-dual pair that is zero if and only if it is evaluated at the optimum, and can be computed by the user without necessarily having knowledge of the minimum. This is especially useful in software packages, as it can be used to determine a stopping criterion for the algorithm.

For any $(\mathbf{x}, \mathbf{y}) \in \mathcal{X} \times \mathcal{Y}$, define the functions

$$\Phi(\mathbf{x}) = \max_{\mathbf{v} \in \mathcal{Y}} \mathcal{L}(\mathbf{x}, \mathbf{v}) \text{ and } \Psi(\mathbf{y}) = \min_{\mathbf{u} \in \mathcal{X}} \mathcal{L}(\mathbf{u}, \mathbf{y}).$$

One such certificate is a direct upper bound on the primal-dual gap

$$\text{Gap}(\mathbf{x}, \mathbf{y}) = \max_{\mathbf{v} \in \mathcal{Y}} \{\mathcal{L}(\mathbf{x}, \mathbf{v}) - \mathcal{L}(\mathbf{x}, \mathbf{y})\} - \min_{\mathbf{u} \in \mathcal{X}} \{\mathcal{L}(\mathbf{u}, \mathbf{y}) - \mathcal{L}(\mathbf{x}, \mathbf{y})\}. \quad (3.98)$$

While the quantity above has the appealing property of agreeing with our notion of suboptimality tracked in the proofs, it is not necessarily the best quantity to track in practice when the strong convexity constants are low (or even zero). Instead, we pursue an upper bound on the *smoothed duality gap* [Walwil and Fercoq, 2025] given by a hyperparameter $\beta = (\beta_{\mathbf{x}}, \beta_{\mathbf{y}})$ and defined as

$$\begin{aligned} \text{Gap}_{\beta}(\mathbf{x}, \mathbf{y}) = & \max_{\mathbf{v} \in \mathcal{Y}} \left\{ \mathcal{L}(\mathbf{x}, \mathbf{v}) - \mathcal{L}(\mathbf{x}, \mathbf{y}) - \frac{\beta_{\mathbf{y}}}{2} \|\mathbf{v} - \mathbf{y}\|_{\mathcal{Y}}^2 \right\} \\ & - \min_{\mathbf{u} \in \mathcal{X}} \left\{ \mathcal{L}(\mathbf{u}, \mathbf{y}) - \mathcal{L}(\mathbf{x}, \mathbf{y}) + \frac{\beta_{\mathbf{x}}}{2} \|\mathbf{u} - \mathbf{x}\|_{\mathcal{X}}^2 \right\}, \end{aligned} \quad (3.99)$$

for a potential return value (\mathbf{x}, \mathbf{y}) (possibly computed using past iterates of the algorithm), where $\beta_{\mathbf{x}} \geq 0$ and $\beta_{\mathbf{y}} \geq 0$. This may also be called the *self-centered* smoothed duality gap because the primal regularizer is centered at \mathbf{x} and the dual regularizer is centered at \mathbf{y} . See Fercoq [2023, Proposition 15] for a proof of its status as a convergence certificate, in that $\text{Gap}_{\beta}(\mathbf{x}, \mathbf{y}) \geq 0$ and that $(\mathbf{x}^*, \mathbf{y}^*)$ is a saddle point of the objective (assuming that one exists) if and only if $\text{Gap}_{\beta}(\mathbf{x}^*, \mathbf{y}^*) = 0$. Notice in addition that for any choice of β and any $(\mathbf{x}, \mathbf{y}) \in \mathcal{X} \times \mathcal{Y}$, it holds that

$$\frac{\mu + \beta_{\mathbf{x}}}{2} \|\mathbf{x} - \mathbf{u}^*(\mathbf{x})\|_{\mathcal{X}}^2 + \frac{\nu + \beta_{\mathbf{y}}}{2} \|\mathbf{y} - \mathbf{v}^*(\mathbf{y})\|_{\mathcal{Y}}^2 \leq \text{Gap}_{\beta}(\mathbf{x}, \mathbf{y}) \leq \text{Gap}(\mathbf{x}, \mathbf{y}),$$

where $\mathbf{u}^*(\mathbf{x})$ and $\mathbf{v}^*(\mathbf{y})$ denote the minimizer and maximizer of (3.99), respectively. Thus, $\text{Gap}_{\beta}(\mathbf{x}, \mathbf{y})$ enjoys the same convergence rate as the primal-dual gap.

We proceed to upper bound $\text{Gap}_{\beta}(\mathbf{x}, \mathbf{y})$ in a form that may not depend on unknown quantities such as the saddle-point $(\mathbf{x}^*, \mathbf{y}^*)$ in the strong convex-strongly concave setting. Instead, we operate with the following toolkit. For any $\mathbf{u} \in \mathbb{R}^d$ and $\mathbf{v} \in \mathbb{R}^n$, we assume access to the “proximal gradient” oracles

$$\mathbf{x} \mapsto \min_{\mathbf{u} \in \mathcal{X}} \langle \mathbf{x}, \mathbf{u} \rangle + \phi(\mathbf{u}) + \frac{\beta_{\mathbf{x}}}{2} \|\mathbf{u} - \mathbf{x}\|_{\mathcal{X}}^2 \quad \text{and} \quad \mathbf{y} \mapsto \max_{\mathbf{v} \in \mathcal{Y}} \langle \mathbf{y}, \mathbf{v} \rangle - \psi(\mathbf{v}) - \frac{\beta_{\mathbf{y}}}{2} \|\mathbf{v} - \mathbf{y}\|_{\mathcal{Y}}^2, \quad (3.100)$$

for which the output of either map is finite when $\mu + \beta_{\mathbf{x}} > 0$ and $\nu + \beta_{\mathbf{y}} > 0$. With a slight abuse of notation, consider a sequence of primal-dual iterates $(\mathbf{x}_k, \mathbf{y}_k)_{k \geq 1}$, where each $(\mathbf{x}_k, \mathbf{y}_k)$ represents the estimated solution *after* $(n \vee d)k$ iterations of an underlying algorithm. Thus, we evaluate the algorithm at every $(n \vee d)$ iterations, so that if we make a full batch gradient computation (at $O(nd)$ cost) the complexity of computing the certificate costs $O(n + d)$ complexity when amortized over the sequence. Finally, our bound can be applied to *any* algorithm, or sequence of primal-dual iterates, including the methods of Chapter 2.

We prove an upper bound on (3.99) that can be computed in practice and supports averaging iterates. We first state a generic upper bound that can aggregate many points and specify practical options after the statement of the result. Below, the notation $i : k$ is used

to index an element of the k -length row of a triangular array.

Lemma 3.7.1. *Consider probability mass weights $(\lambda_{i:k})_{i=1}^k$ and evaluation points $\bar{\mathbf{x}}_k, \hat{\mathbf{x}}_{1:k}, \dots, \hat{\mathbf{x}}_{k:k} \in \mathcal{X}$, and $\hat{\mathbf{y}}_{1:k}, \dots, \hat{\mathbf{y}}_{k:k} \in \mathcal{Y}$. Define*

$$\bar{\mathbf{y}}_k = \sum_{i=1}^k \lambda_{i:k} \hat{\mathbf{y}}_{i:k}.$$

Assume that f_1, \dots, f_n are bounded below by zero. Then, we have that

$$\text{Gap}_\beta(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k) \leq \max_{\mathbf{v} \in \mathcal{Y}} \left\{ \langle f(\bar{\mathbf{x}}_k), \mathbf{v} \rangle - \psi(\mathbf{v}) - \frac{\beta_{\mathbf{y}}}{2} \|\mathbf{v} - \bar{\mathbf{y}}_k\|_{\mathcal{Y}}^2 \right\} + \phi(\bar{\mathbf{x}}_k) + \psi(\bar{\mathbf{y}}_k) - \max \{M_k, 0\}, \quad (3.101)$$

where

$$\begin{aligned} M_k = \min_{\mathbf{u} \in \mathcal{X}} & \left\{ \sum_{i=1}^k \lambda_{i:k} \langle \nabla f(\hat{\mathbf{x}}_{i:k})^\top \hat{\mathbf{y}}_{i:k}, \mathbf{u} \rangle + \phi(\mathbf{u}) + \frac{\beta_{\mathbf{x}}}{2} \|\mathbf{u} - \bar{\mathbf{x}}_k\|_{\mathcal{X}}^2 \right\} \\ & + \sum_{i=1}^k \lambda_{i:k} \langle \hat{\mathbf{y}}_{i:k}, f(\hat{\mathbf{x}}_{i:k}) - \nabla f(\hat{\mathbf{x}}_{i:k}) \hat{\mathbf{x}}_{i:k} \rangle. \end{aligned}$$

Proof. First, we expand the two terms in (3.99) and observe

$$\begin{aligned} & \min_{\mathbf{u} \in \mathcal{X}} \left\{ \mathcal{L}(\mathbf{u}, \bar{\mathbf{y}}_k) - \mathcal{L}(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k) + \frac{\beta_{\mathbf{x}}}{2} \|\mathbf{u} - \bar{\mathbf{x}}_k\|_{\mathcal{X}}^2 \right\} \\ &= \langle f(\bar{\mathbf{x}}_k), \bar{\mathbf{y}}_k \rangle + \phi(\bar{\mathbf{x}}_k) - \min_{\mathbf{u} \in \mathcal{X}} \left\{ \langle f(\mathbf{u}), \bar{\mathbf{y}}_k \rangle + \phi(\mathbf{u}) + \frac{\beta_{\mathbf{x}}}{2} \|\mathbf{u} - \bar{\mathbf{x}}_k\|_{\mathcal{X}}^2 \right\} \\ & \max_{\mathbf{v} \in \mathcal{Y}} \left\{ \mathcal{L}(\bar{\mathbf{x}}_k, \mathbf{v}) - \mathcal{L}(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k) - \frac{\beta_{\mathbf{y}}}{2} \|\mathbf{v} - \bar{\mathbf{y}}_k\|_{\mathcal{Y}}^2 \right\} \\ &= \max_{\mathbf{v} \in \mathcal{Y}} \left\{ \langle f(\bar{\mathbf{x}}_k), \mathbf{v} \rangle - \psi(\mathbf{v}) - \frac{\beta_{\mathbf{y}}}{2} \|\mathbf{v} - \bar{\mathbf{y}}_k\|_{\mathcal{Y}}^2 \right\} - \langle f(\bar{\mathbf{x}}_k), \bar{\mathbf{y}}_k \rangle + \psi(\bar{\mathbf{y}}_k), \end{aligned}$$

which sum to

$$\begin{aligned} \text{Gap}(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k) &= \max_{\mathbf{v} \in \mathcal{Y}} \left\{ \langle f(\bar{\mathbf{x}}_k), \mathbf{v} \rangle - \psi(\mathbf{v}) - \frac{\beta_{\mathbf{y}}}{2} \|\mathbf{v} - \bar{\mathbf{y}}_k\|_{\mathcal{Y}}^2 \right\} + \phi(\bar{\mathbf{x}}_k) + \psi(\bar{\mathbf{y}}_k) \\ & \quad - \min_{\mathbf{u} \in \mathcal{X}} \left\{ \langle \bar{\mathbf{y}}_k, f(\mathbf{u}) \rangle + \phi(\mathbf{u}) + \frac{\beta_{\mathbf{x}}}{2} \|\mathbf{u} - \bar{\mathbf{x}}_k\|_{\mathcal{X}}^2 \right\}. \end{aligned} \quad (3.102)$$

Note that all terms are computable besides the minimization over $u \in \mathcal{X}$. Thus, the rest of the

proof relies on lower bounding the minimization term $\min_{\mathbf{u} \in \mathcal{X}} \{ \langle \bar{\mathbf{y}}, f(\mathbf{u}) \rangle + \phi(\mathbf{u}) + \frac{\beta_{\mathbf{x}}}{2} \|\mathbf{u} - \bar{\mathbf{x}}_k\|_{\mathcal{X}}^2 \}$, which is non-negative by assumption. Write

$$\begin{aligned} & \langle \bar{\mathbf{y}}_k, f(\mathbf{u}) \rangle + \phi(\mathbf{u}) + \frac{\beta_{\mathbf{x}}}{2} \|\mathbf{u} - \bar{\mathbf{x}}_k\|_{\mathcal{X}}^2 \\ &= \sum_{i=1}^k \lambda_{i:k} \langle \hat{\mathbf{y}}_{i:k}, f(\mathbf{u}) \rangle + \phi(\mathbf{u}) + \frac{\beta_{\mathbf{x}}}{2} \|\mathbf{u} - \bar{\mathbf{x}}_k\|_{\mathcal{X}}^2 \\ &\geq \sum_{i=1}^k \lambda_{i:k} \left(\langle \hat{\mathbf{y}}_{i:k}, f(\hat{\mathbf{x}}_{i:k}) \rangle + \langle \nabla f(\hat{\mathbf{x}}_{i:k})^\top \hat{\mathbf{y}}_{i:k}, \mathbf{u} - \hat{\mathbf{x}}_{i:k} \rangle \right) + \phi(\mathbf{u}) + \frac{\beta_{\mathbf{x}}}{2} \|\mathbf{u} - \bar{\mathbf{x}}_k\|_{\mathcal{X}}^2, \end{aligned}$$

by convexity so that

$$\begin{aligned} & \min_{\mathbf{u} \in \mathcal{X}} \left\{ \langle \bar{\mathbf{y}}, f(\mathbf{u}) \rangle + \phi(\mathbf{u}) + \frac{\beta_{\mathbf{x}}}{2} \|\mathbf{u} - \bar{\mathbf{x}}_k\|_{\mathcal{X}}^2 \right\} \\ &\geq \min_{\mathbf{u} \in \mathcal{X}} \left\{ \sum_{i=1}^k \lambda_{i:k} \langle \nabla f(\hat{\mathbf{x}}_{i:k})^\top \hat{\mathbf{y}}_{i:k}, \mathbf{u} \rangle + \phi(\mathbf{u}) + \frac{\beta_{\mathbf{x}}}{2} \|\mathbf{u} - \bar{\mathbf{x}}_k\|_{\mathcal{X}}^2 \right\} \\ &\quad + \sum_{i=1}^k \lambda_{i:k} \langle \hat{\mathbf{y}}_{i:k}, f(\hat{\mathbf{x}}_{i:k}) \rangle - \sum_{i=1}^k \lambda_{i:k} \langle \nabla f(\hat{\mathbf{x}}_{i:k})^\top \hat{\mathbf{y}}_{i:k}, \hat{\mathbf{x}}_{i:k} \rangle. \end{aligned}$$

Substituting the definition of M_k completes the proof. \square

One practical note is that the term that usually blows up when $\mu \approx 0$ is M_k —the minimization over $\mathbf{u} \in \mathcal{X}$ —whenever the primal gradient is not nearly zero.

To convert this upper bound into an algorithm, we make the choices

$$\hat{\mathbf{x}}_{i:k} = \mathbf{x}_i, \quad \hat{\mathbf{x}}_{i:k} = \mathbf{x}_i, \quad \text{and} \quad \bar{\mathbf{x}}_k = \sum_{i=1}^k \lambda_{i:k} \mathbf{x}_i.$$

For the weights $(\lambda_{i:k})_{i=1}^k$, we use an exponential moving average. Notice that all terms in (3.101) can be computed easily if every term that is indexed by k is maintained at every iteration of the optimization problem. See Algorithm 5 for an implementation-friendly description.

Figure 3.5 compares the certificate from Algorithm 5 to the standard (primal) gradient norm criterion on the distributionally robust classification and regression benchmarks from

Algorithm 5 Dual-Linear Online Accuracy Certificate

Inputs: Initial points $(\mathbf{x}_0, \mathbf{y}_0)$, total epochs K , averaging constant $\lambda \in [0, 1]$, hyperparameters $\beta = (\beta_{\mathbf{x}}, \beta_{\mathbf{y}})$.

- 1: $\mathbf{z}_0 = \nabla f(\mathbf{x}_0)^\top \mathbf{y}_0 \in \mathbb{R}^d$, $r_0 = \langle \mathbf{y}_0, f(\mathbf{x}_0) - \nabla f(\mathbf{x}_0) \mathbf{x}_0 \rangle \in \mathbb{R}$, $(\bar{\mathbf{x}}_0, \bar{\mathbf{y}}_0) = (\mathbf{x}_0, \mathbf{y}_0)$.
- 2: **for** $k = 1, \dots, K$ **do**
- 3: $(\mathbf{x}_k, \mathbf{y}_k) = \text{RunEpoch}(\mathbf{x}_{k-1}, \mathbf{y}_{k-1})$.
- 4: $(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k) = (1 - \lambda) \cdot (\bar{\mathbf{x}}_{k-1}, \bar{\mathbf{y}}_{k-1}) + \lambda \cdot (\mathbf{x}_k, \mathbf{y}_k)$.
- 5: $\mathbf{z}_k = (1 - \lambda) \cdot \mathbf{z}_{k-1} + \lambda \cdot \nabla f(\mathbf{x}_k)^\top \mathbf{y}_k$.
- 6: $r_k = (1 - \lambda) \cdot r_{k-1} + \lambda \cdot \langle \mathbf{y}_k, f(\mathbf{x}_k) - \nabla f(\mathbf{x}_k) \mathbf{x}_k \rangle$.
- 7: $M_k = \min_{\mathbf{u} \in \mathcal{X}} \left\{ \langle \mathbf{z}_k, \mathbf{u} \rangle + \phi(\mathbf{u}) + \frac{\beta_{\mathbf{x}}}{2} \|\mathbf{u} - \bar{\mathbf{x}}_k\|_{\mathcal{X}}^2 \right\} + r_k$.
- 8: $\varepsilon_k = \max_{\mathbf{v} \in \mathcal{Y}} \left\{ \langle f(\bar{\mathbf{x}}_k), \mathbf{v} \rangle - \psi(\mathbf{v}) - \frac{\beta_{\mathbf{y}}}{2} \|\mathbf{v} - \bar{\mathbf{y}}_k\|_{\mathcal{Y}}^2 \right\} - \max\{M_k, 0\} + \phi(\bar{\mathbf{x}}_k) + \psi(\bar{\mathbf{y}}_k)$.

Output: Sequence $(\varepsilon_k)_{k=1}^K$.

Section 2.9 under the same experimental setup. Precisely, the gradient norm criterion is computed as

$$\|\nabla \Phi(\mathbf{x}_k)\|_{\mathcal{X}^*}^2 \text{ for } \Phi(\mathbf{x}) := \max_{\mathbf{v} \in \mathcal{Y}} \mathcal{L}(\mathbf{x}, \mathbf{v}) \quad (3.103)$$

at epoch k . For the certificate parameters, we set $\beta_{\mathbf{x}} = 1$, $\beta_{\mathbf{y}} = 0$, and $\nu = 0.05$ in all experiments. In contrast, μ is tuned via training an empirical risk minimization model and choosing the value that yields the best generalization performance. Note that for some datasets, this may be as small as $\mu \approx 10^{-5}$, making a positive value of $\beta_{\mathbf{x}}$ essential. We find that the certificate, while not guaranteed to be an upper bound to the primal-dual gap, remains close in experimental settings using a flat hyperparameter choice. Thus, the tolerance can be set at approximately the same magnitude as would be required of the primal-dual gap/suboptimality.

3.7.2 Lower Bounds

In this section, we show that under particular parameter regimes, the full vector update method analyzed in Theorem 3.3.1 achieves a matching lower bound on the number of iterations k required to achieve an ε -suboptimal primal-dual pair $(\mathbf{x}_k, \mathbf{y}_k)$. This helps calibrate

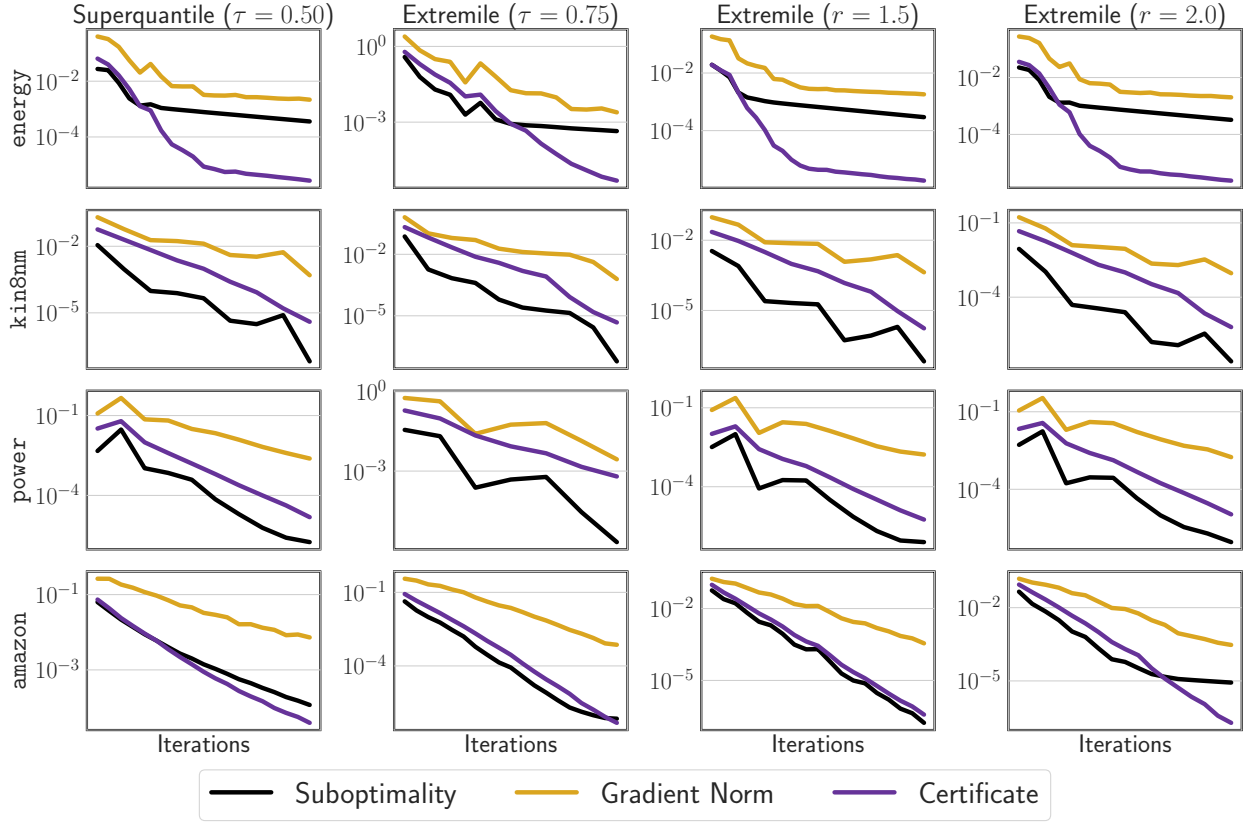


Figure 3.5: **Experimental Evaluation of Online Accuracy Certificates.** In all plots, the x -axis refers to the iteration count, which may differ between datasets. Each line represents the gradient norm (3.103), certificate (3.101), and the primal-dual gap (3.98).

the upper bounds stated in previous sections to the true problem hardness of semilinear min-max programs. We assume that $\mu > 0$ and $\nu > 0$, guaranteeing the existence of a unique solution $(\mathbf{x}^*, \mathbf{y}^*)$ satisfying

$$\max_{\mathbf{y} \in \mathcal{Y}} \mathcal{L}(\mathbf{x}^*, \mathbf{y}) \leq \mathcal{L}(\mathbf{x}^*, \mathbf{y}^*) \leq \min_{\mathbf{x} \in \mathcal{X}} \mathcal{L}(\mathbf{x}, \mathbf{y}^*),$$

i.e. a saddle point of the objective \mathcal{L} . We use the gap criterion (3.4) by setting $(\mathbf{u}, \mathbf{v}) = (\mathbf{x}^*, \mathbf{y}^*)$, and recognize by strong convexity that

$$\text{Gap}^{\mathbf{x}^*, \mathbf{y}^*}(\mathbf{x}, \mathbf{y}) := \mathcal{L}(\mathbf{x}, \mathbf{y}^*) - \mathcal{L}(\mathbf{x}^*, \mathbf{y}) \geq \frac{\mu}{2} \|\mathbf{x} - \mathbf{x}^*\|_{\mathcal{X}}^2 + \frac{\nu}{2} \|\mathbf{y} - \mathbf{y}^*\|_{\mathcal{Y}}^2 \geq 0.$$

The lower bound will be applied to the term $\|\mathbf{x} - \mathbf{x}^*\|_{\mathbf{x}}^2$, and as a result apply to the gap criterion. The proof follows similar techniques as those used in Zhang et al. [2022] for bilinearly coupled objectives:

1. We describe precisely the class of algorithms considered, and verify that the method of Theorem 3.3.1 is such an algorithm.
2. We then introduce a “hard” instance which satisfies Assumption 3.2.2 and Assumption 3.2.1, and then characterize particular subspaces for which the primal-dual iterates are guaranteed to be members.
3. Finally, we compute by hand an arbitrarily close approximation $\hat{\mathbf{x}}^*$ of the optimal primal solution \mathbf{x}^* , and show that for d sufficiently large and $\sqrt{\mu\nu}/G \leq 1/2$,

$$\|\mathbf{x}_{2k} - \hat{\mathbf{x}}^*\|_{\mathbf{x}}^2 \geq C(1 - 2\sqrt{\mu\nu}/G)^{2k} \|\mathbf{x}_0 - \hat{\mathbf{x}}^*\|_{\mathbf{x}}^2.$$

This describes the parameter regime

$$\mu G^2 \geq \nu L^2, \tag{3.104}$$

indicating that the mixed condition number in Theorem 3.3.1 is dominant. This reflects practice in the case of the examples from Section 3.2. In distributionally robust optimization, ν is often chosen as a smoothing parameter for a problem that is originally non-smooth in the primal. In fully composite optimization, we have that $(1/\nu)$ is the smoothness constant of F , which could be large for F “close to non-smooth”. Finally, in the case of optimization with functional constraints, ν refers to a penalty applied to the Lagrange multipliers, which would originally be unpenalized in order to enforce the constraint exactly. Remarkably, in the parameter regime (3.104), we may use a *bilinear* example as our hard instance. We fix the dual-linear convention in this section as well, and handle the primal-linear case at the end of the section.

Algorithm Class and Hard Instance This section contains the portions of the proof that generalize and improve on the result of Zhang et al. [2022]. We define our class of methods to satisfy a linear span assumption involving proximal oracles. This definition will reflect [Zhang et al., 2022, Definition 2.1], but account for a number of changes: 1) possible nonbilinearity, 2) differing problem dimensions $n \neq d$, 3) non-Euclidean proximal oracles, 4) extrapolated gradient estimates, and 5) the use of \mathbf{x}_k in the update for \mathbf{y}_k . Define the non-Euclidean proximal operators

$$\begin{aligned} p_{\eta_k, \phi}(\tilde{\mathbf{x}}, \tilde{\mathbf{g}}) &:= \arg \min_{\mathbf{x} \in \mathcal{X}} \langle \tilde{\mathbf{g}}, \mathbf{x} \rangle + \phi(\mathbf{x}) + \eta_k \Delta_{\mathcal{X}}(\mathbf{x}, \tilde{\mathbf{x}}) & (\eta_k \geq 0, \tilde{\mathbf{x}} \in \mathcal{X}, \tilde{\mathbf{g}} \in \mathcal{X}^*) \\ p_{\delta_k, \psi}(\tilde{\mathbf{y}}, \tilde{\mathbf{f}}) &:= \arg \max_{\mathbf{y} \in \mathcal{Y}} \langle \tilde{\mathbf{f}}, \mathbf{y} \rangle - \psi(\mathbf{y}) - \delta_k \Delta_{\mathcal{Y}}(\mathbf{y}, \tilde{\mathbf{y}}) & (\delta_k \geq 0, \tilde{\mathbf{y}} \in \mathcal{Y}, \tilde{\mathbf{f}} \in \mathcal{Y}^*) \end{aligned}$$

and consider the following definition.

Definition 3.7.1 (Dual-Linear Proximal Algorithm Class). We define a (dual-linear) *deterministic proximal gradient algorithm* to be a sequence of primal-dual iterates $(\mathbf{x}_k, \mathbf{y}_k)_{k \geq 0}$ satisfying $(\mathbf{x}_k, \mathbf{y}_k) \in \mathcal{H}_k^{\mathcal{X}} \times \mathcal{H}_k^{\mathcal{Y}}$, where $\mathcal{H}_k^{\mathcal{X}} \subseteq \mathcal{X}$ and $\mathcal{H}_k^{\mathcal{Y}} \subseteq \mathcal{Y}$ are constructed from the following steps. First, for $k \geq 0$, define the subspaces

$$\begin{aligned} \mathcal{G}_k^{\mathcal{X}} &= \text{span} \{ \nabla f(\mathbf{x}_i)^\top \mathbf{y}_i : 0 \leq i \leq k \} \subseteq \mathcal{X}^* \\ \mathcal{G}_k^{\mathcal{Y}} &= \text{span} \{ f(\mathbf{x}_i) : 0 \leq i \leq k+1 \} \subseteq \mathcal{Y}^*. \end{aligned}$$

Define $\mathcal{H}_0^{\mathcal{X}} = \text{span}\{\mathbf{x}_0\}$ and $\mathcal{H}_0^{\mathcal{Y}} = \text{span}\{\mathbf{y}_0\}$, and for $k \geq 1$, we recursively define

$$\begin{aligned} \mathcal{H}_k^{\mathcal{X}} &= \text{span} \left(\{ p_{\eta_k, \phi}(\tilde{\mathbf{x}}_{k-1}, \tilde{\mathbf{g}}_{k-1}) : \tilde{\mathbf{x}}_{k-1} \in \mathcal{H}_{k-1}^{\mathcal{X}}, \tilde{\mathbf{g}}_{k-1} \in \mathcal{G}_{k-1}^{\mathcal{X}} \} \cup \mathcal{H}_{k-1}^{\mathcal{X}} \right), \\ \mathcal{H}_k^{\mathcal{Y}} &= \text{span} \left(\{ p_{\delta_k, \psi}(\tilde{\mathbf{y}}_{k-1}, \tilde{\mathbf{f}}_k) : \tilde{\mathbf{y}}_{k-1} \in \mathcal{H}_{k-1}^{\mathcal{Y}}, \tilde{\mathbf{f}}_k \in \mathcal{G}_{k-1}^{\mathcal{Y}} \} \cup \mathcal{H}_{k-1}^{\mathcal{Y}} \right). \end{aligned}$$

We must confirm that the algorithm analyzed in Theorem 3.3.1 adheres to this definition.

Lemma 3.7.2. *The primal-dual updates (3.11) and (3.12) define a deterministic proximal gradient algorithm in the sense of Definition 3.7.1.*

Proof. First for any $k \geq 1$, notice the primal update (3.11) and (3.12) can be written as

$$\mathbf{x}_k = p_{\eta_k \phi}(\mathbf{x}_{k-1}, \bar{\mathbf{g}}_{k-1}) \text{ for } \eta_k = \frac{A_{k-1}\mu + \mu_0}{a_k} \geq 0, \quad (3.105)$$

$$\mathbf{y}_k = p_{\delta_k \psi}(\mathbf{y}_{k-1}, f(\mathbf{x}_k)) \text{ for } \delta_k = \frac{A_{k-1}\nu + \nu_0}{a_k} \geq 0. \quad (3.106)$$

By the definition of the primal gradient estimate $\bar{\mathbf{g}}_{k-1}$ from (3.27), we also have that

$$(\bar{\mathbf{g}}_{k-1}, f(\mathbf{x}_k)) \in \mathcal{G}_{k-1}^{\mathbf{x}} \times \mathcal{G}_{k-1}^{\mathbf{y}}.$$

By construction, we have that initial iterates $(\mathbf{x}_0, \mathbf{y}_0) \in \mathcal{H}_0^{\mathbf{x}} \times \mathcal{H}_0^{\mathbf{y}}$. Combine this base case with the recursions (3.105), and (3.106) to prove by induction that $\mathbf{x}_k \in \mathcal{H}_k^{\mathbf{x}}$ and $\mathbf{y}_k \in \mathcal{H}_k^{\mathbf{y}}$. \square

Next, we proceed to define a special case of the objective \mathcal{L} and establish a *zero-chain* property for the sequence of iterates. Consider a setting in which $n \geq d$, $\|\cdot\|_{\mathbf{x}} = \|\cdot\|_2$ and $\|\cdot\|_{\mathbf{y}} = \|\cdot\|_2$ (the ℓ_2 -norms on \mathbb{R}^d and \mathbb{R}^n), respectively. The objective is written

$$\mathcal{L}(\mathbf{x}, \mathbf{y}) = \underbrace{\langle \mathbf{y}, \frac{G}{2} \mathbf{A} \mathbf{x} \rangle}_{\langle \mathbf{y}, f(\mathbf{x}) \rangle} - \underbrace{\frac{\nu}{2} \|\mathbf{y}\|_2^2}_{\psi(\mathbf{y})} + \underbrace{\mathbf{c}^\top \mathbf{x} + \frac{\mu}{2} \|\mathbf{x}\|_2^2}_{\phi(\mathbf{x})}, \quad (3.107)$$

where $\mathbf{A} \in \mathbb{R}^{n \times d}$ and $\mathbf{c} \in \mathbb{R}^d$ are to-be-specified. As a result, we have that

$$\Delta_{\mathbf{x}}(\mathbf{x}, \mathbf{x}') = \frac{1}{2} \|\mathbf{x} - \mathbf{x}'\|_2^2 \text{ and } \Delta_{\mathbf{y}}(\mathbf{y}, \mathbf{y}') = \frac{1}{2} \|\mathbf{y} - \mathbf{y}'\|_2^2.$$

We characterize the subspaces $\mathcal{H}_k^{\mathbf{x}}$ and $\mathcal{H}_k^{\mathbf{y}}$ from Definition 3.7.1 on this particular example, which will reveal the structure of \mathbf{x}_k that is used in the next step. First, notice that the proximal oracles can be computed directly for (3.107). Recalling that we defined $(\mathbf{x}_k, \mathbf{y}_k) = (\mathbf{x}_0, \mathbf{y}_0)$ for $k < 0$, we have for all $k \geq 1$,

$$\mathbf{x}_k = \frac{\eta_k, \mathbf{x}_{k-1} - \bar{\mathbf{g}}_{k-1}}{\mu + \eta_k} \subseteq \text{span} \{ \mathbf{x}_{k-1}, \mathbf{c}, \mathbf{A}^\top \mathbf{y}_{k-1}, \mathbf{A}^\top \mathbf{y}_{k-2}, \dots, \mathbf{A}^\top \mathbf{y}_0 \} \quad (3.108)$$

$$\mathbf{y}_k = \frac{\delta_k, \mathbf{y}_{k-1} - \bar{\mathbf{f}}_{k-1}}{\nu + \delta_k} \subseteq \text{span} \{ \mathbf{y}_{k-1}, \mathbf{A} \mathbf{x}_k, \mathbf{A} \mathbf{x}_{k-1}, \dots, \mathbf{A} \mathbf{x}_0 \}. \quad (3.109)$$

We can then present the key result of this step.

Lemma 3.7.3. *For any deterministic proximal gradient algorithm in the sense of Definition 3.7.1 with $\mathbf{x}_0 = \mathbf{0}_d$ and $\mathbf{y}_0 = \mathbf{0}_n$ applied to the problem (3.107), we have that $\mathcal{H}_1^{\mathbf{x}} = \{\mathbf{c}\}$, $\mathcal{H}_1^{\mathbf{y}} = \text{span}\{\mathbf{Ac}\}$, and for all $k \geq 2$,*

$$\mathcal{H}_k^{\mathbf{x}} \subseteq \text{span}\{(\mathbf{A}^\top \mathbf{A})^i \mathbf{c} : 0 \leq i \leq k-1\}, \quad (3.110)$$

$$\mathcal{H}_k^{\mathbf{y}} \subseteq \text{span}\{(\mathbf{AA}^\top)^i (\mathbf{Ac}) : 0 \leq i \leq k-1\}. \quad (3.111)$$

Proof. Using that $\mathbf{x}_0 = \mathbf{0}_d$ and $\mathbf{y}_0 = \mathbf{0}_n$, and directly applying the formulas (3.108) and (3.109) we achieve the base cases

$$\mathcal{H}_1^{\mathbf{x}} = \{\mathbf{c}\} \text{ and } \mathcal{H}_1^{\mathbf{y}} \subseteq \text{span}\{\mathbf{Ac}\}.$$

We argue the general case by induction. Fix $k \geq 2$ and assume that for $\{\kappa : 0 \leq \kappa \leq k-1\}$, we have that (3.110) and (3.111) hold. Then, we may rewrite the conclusion of (3.108) and (3.109) as

$$\mathbf{x}_k \subseteq \text{span}\{\mathbf{x}_{k-1}, \mathbf{c}, \mathbf{A}^\top \mathbf{y}_{k-1}\} \text{ and } \mathbf{y}_k \subseteq \text{span}\{\mathbf{y}_{k-1}, \mathbf{Ax}_k\},$$

and we may apply these inclusions in an alternating manner to claim

$$\begin{aligned} \mathbf{x}_k &\subseteq \text{span}\{\mathbf{x}_{k-1}, \mathbf{c}, \mathbf{A}^\top \mathbf{y}_{k-1}\} \subseteq \text{span}\{\mathbf{x}_{k-1}, \mathbf{c}, \mathbf{A}^\top \mathbf{Ax}_{k-1}, \mathbf{A}^\top \mathbf{y}_{k-2}\} \\ &\subseteq \text{span}\{\mathbf{x}_{k-1}, \mathbf{c}, (\mathbf{A}^\top \mathbf{A})\mathbf{x}_{k-1}, \dots, (\mathbf{A}^\top \mathbf{A})\mathbf{x}_1\} \\ &= (\mathbf{A}^\top \mathbf{A}) \left(\text{span}\{(\mathbf{A}^\top \mathbf{A})^i \mathbf{c} : 0 \leq i \leq k-2\} \right) \\ &= \text{span}\{(\mathbf{A}^\top \mathbf{A})^i \mathbf{c} : 0 \leq i \leq k-1\}. \end{aligned}$$

Arguing similarly for the sequence $(\mathbf{y}_k)_{k \geq 1}$, we have that

$$\begin{aligned} \mathbf{y}_k &\subseteq \text{span}\{\mathbf{y}_{k-1}, \mathbf{Ac}, (\mathbf{AA}^\top)\mathbf{y}_{k-1}, \dots, (\mathbf{AA}^\top)\mathbf{y}_1\} \\ &= \text{span}\{(\mathbf{AA}^\top)^i (\mathbf{Ac}) : 0 \leq i \leq k-1\}. \end{aligned}$$

This completes the proof. □

Note that $\mathbf{x}_0 = \mathbf{0}_d$ and $\mathbf{y}_0 = \mathbf{0}_n$ can hold without loss of generality, as we may shift the input space to satisfy this condition. Now, we are ready to partially specify \mathbf{A} and \mathbf{c} . We will set \mathbf{A} to recover a variant of Nesterov's tridiagonal matrix (see [Nesterov, 2018, Section 2.3]), so that for

$$\mathbf{A} = \begin{bmatrix} 1 & -1 & & \\ & 1 & -1 & \\ & & \ddots & \ddots \\ & & & 1 \end{bmatrix}, \quad \mathbf{B} = \begin{bmatrix} 2 & -1 & & \\ -1 & \ddots & -1 & \\ & -1 & 2 & -1 \\ & & -1 & 1 \end{bmatrix}, \quad (3.112)$$

we have that

$$\mathbf{A}\mathbf{A}^\top = \begin{bmatrix} \mathbf{B} & \mathbf{0}_{d \times n-d} \\ \mathbf{0}_{n-d \times d} & \mathbf{0}_{n-d \times n-d} \end{bmatrix}, \quad \mathbf{A}^\top \mathbf{A} = \begin{bmatrix} 1 & -1 & & \\ -1 & \ddots & -1 & \\ & -1 & 2 & -1 \\ & & -1 & 2 \end{bmatrix}. \quad (3.113)$$

The operator norm of \mathbf{A} will be the spectral norm $\|\mathbf{A}\|_{2,2} \leq 2$, so the multiplier $\frac{G}{2}$ makes f adhere to Assumption 3.2.2 for any $L \geq 0$. Using these specifications and Lemma 3.7.3, we establish the zero-chain property for the primal-dual sequence.

Corollary 3.7.1. *Let $n \geq d$ and use the notation $\mathbf{e}_{j:p}$ to denote the j -th standard basis vector in \mathbb{R}^p . Using \mathbf{A} as defined in (3.113) and $\mathbf{c} = \beta \mathbf{e}_{1:d}$ for a constant $\beta \in \mathbb{R}$, we have that*

$$\begin{aligned} \mathcal{H}_k^x &\subseteq \text{span} \{ \mathbf{e}_{1:d}, \dots, \mathbf{e}_{\min\{k,d\}:d} \} \\ \mathcal{H}_k^y &\subseteq \text{span} \{ \mathbf{e}_{1:n}, \dots, \mathbf{e}_{\min\{k,d\}:n} \}. \end{aligned}$$

Proof. First, note that

$$\mathbf{A}\mathbf{c} = \beta \mathbf{A}\mathbf{e}_{1:d} = \beta \mathbf{e}_{1:d}.$$

Additionally, because $\mathbf{A}\mathbf{A}^\top$ and $\mathbf{A}^\top \mathbf{A}$ are tridiagonal we have that

$$\begin{aligned} \mathbf{u} \in \text{span}\{\mathbf{e}_{j:p} : j \leq p \leq d\} &\implies \mathbf{A}^\top \mathbf{A}\mathbf{u} \in \text{span}\{\mathbf{e}_{j:p+1} : j \leq p \leq d-1\} \\ \mathbf{v} \in \text{span}\{\mathbf{e}_{j:p} : j \leq p \leq n\} &\implies \mathbf{A}\mathbf{A}^\top \mathbf{v} \in \text{span}\{\mathbf{e}_{j:p+1} : j \leq p \leq d-1\}. \end{aligned}$$

Combining the displays above with Lemma 3.7.3, we have that for $k \geq 2$,

$$\begin{aligned} \mathcal{H}_k^x &\subseteq \text{span}\{\mathbf{e}_{1:d}, \dots, \mathbf{e}_{\min\{k,d\}:d}\} \\ \mathcal{H}_k^y &\subseteq \text{span}\{\mathbf{e}_{1:n}, \dots, \mathbf{e}_{\min\{k,d\}:n}\} \end{aligned}$$

the result as desired. \square

Notice in the result above that d is a limiting factor, in that the zero components of $\mathbf{A}\mathbf{A}^\top$ prevent any further non-zero entries from appearing beyond element d . Thus, in the upcoming proofs, we will increase d to achieve the desired properties.

Dual-Linear Lower Bounds After having established Corollary 3.7.1, the rest of the argument follows very similarly to Zhang et al. [2022, Theorem 3.5]. Noting that $n \geq d$ we first partially maximize the objective over \mathbf{y} (which can be solved by hand), to construct the primal gap $\Phi(\mathbf{x})$. We first recall that \mathcal{Y} may include vectors with negative components without harming convexity, as all functions $(f_j)_{j=1}^n$ are linear, giving

$$\Phi(\mathbf{x}) = \max_{\mathbf{y} \in \mathcal{Y}} \mathcal{L}(\mathbf{x}, \mathbf{y}) = -\frac{1}{2} \mathbf{x}^\top \left(\frac{G^2}{4\nu} \mathbf{A}^\top \mathbf{A} + \mu I \right) \mathbf{x} + \mathbf{c}^\top \mathbf{x}.$$

We now fully specify the objective by setting $\beta = \frac{G^2}{4\nu}$ (recalling that $\mathbf{c} = \beta \mathbf{e}_{1:d}$), yielding

$$\Phi(\mathbf{x}) = \frac{G^2}{4\nu} \left(-\frac{1}{2} \mathbf{x}^\top (\mathbf{A}^\top \mathbf{A} + \alpha I) \mathbf{x} + \mathbf{e}_{1:d}^\top \mathbf{x} \right)$$

for $\alpha = \frac{4\mu\nu}{G^2}$. Consider the following result, which defines an approximation to the maximizer of the expression above.

Lemma 3.7.4. [*Zhang et al., 2022, Lemma 3.3*] Consider the system of linear equations

$$(\mathbf{A}^\top \mathbf{A} + \alpha I) \mathbf{x} = \mathbf{e}_{1:d},$$

for $\alpha > 0$ and $\mathbf{A}^\top \mathbf{A}$ defined in (3.113). Denote the unique solution of this system as \mathbf{x}^* , which is guaranteed to exist by positive definiteness of $\mathbf{A}^\top \mathbf{A} + \alpha I$. Denote by

$$q \equiv q(\alpha) = \frac{1}{2} \left((2 + \alpha) - \sqrt{(2 + \alpha)^2 - 4} \right) \in (0, 1)$$

the smallest root of the quadratic equation $1 - (2 + \alpha)q + q^2 = 0$. Then, by defining

$$\hat{x}_i^* = \frac{q^i}{1 - q}, \forall i \in \{1, \dots, d\} \quad (3.114)$$

and $\hat{\mathbf{x}}^* = (\hat{x}_1^*, \dots, \hat{x}_d^*)$ we have that

$$\|\hat{\mathbf{x}}^* - \mathbf{x}^*\|_2 \leq \frac{q^{d+1}}{\alpha(1 - q)}.$$

We may now make the full claim.

Proposition 3.7.1. Consider any algorithm satisfying Definition 3.7.1 with $\mathbf{x}_0 = \mathbf{0}_d$ and $\mathbf{y}_0 = \mathbf{0}_n$ applied to the problem (3.107), with $n \geq d \geq k$. We have that when $\sqrt{\alpha} = 2\sqrt{\mu\nu}/G \leq 1$,

$$\text{Gap}^{\mathbf{x}^*, \mathbf{y}^*}(\mathbf{x}_k, \mathbf{y}_k) \geq \frac{\mu}{2} \|\mathbf{x}_k - \mathbf{x}^*\|_2^2 \geq \frac{\mu}{8\sqrt{2}} (1 - \sqrt{\alpha})^k \|\mathbf{x}_0 - \mathbf{x}^*\|_2^2 - O((1 - \alpha)^{2d+2})/\alpha$$

Before stating the proof, note that d is a free parameter, so the second term in the lower bound can be made arbitrarily small.

Proof. By Corollary 3.7.1, we have that

$$\mathbf{x}_k \in \mathcal{H}_k^{\mathbf{x}} \subseteq \text{span}\{\mathbf{e}_{1:d}, \dots, \mathbf{e}_{k:d}\},$$

indicating that the last $d - k$ elements of \mathbf{x}_k are necessarily zero. On the other hand, by the definition of $\hat{\mathbf{x}}^*$, it must hold that elements $k + 1, \dots, d$ of $\hat{\mathbf{x}}^*$ are non-zero with formula (3.114). Thus, we have that

$$\begin{aligned}\|\mathbf{x}_{2k} - \hat{\mathbf{x}}^*\|_2^2 &\geq \sum_{i=k+1}^d (\hat{x}_i^*)^2 = \left(\frac{q^k}{1-q}\right)^2 (q^2 + \dots + q^{2(d-k)}) \\ &\geq \frac{q^{2k}}{2} \|\hat{\mathbf{x}}^*\|_2^2 = \frac{q^{2k}}{2} \|\hat{\mathbf{x}}^* - \mathbf{x}_0\|_2^2.\end{aligned}$$

Next, by using the triangle inequality $\sqrt{a+b} \leq \sqrt{a} + \sqrt{b}$ and by the assumption that $\alpha = 4\mu\nu/G \leq 1$, we have that

$$q(\alpha) = \frac{1}{2} \left(2 + \alpha - \sqrt{\alpha^2 + 4\alpha}\right) \geq 1 - \sqrt{\alpha} \geq 0,$$

so that

$$\|\mathbf{x}_{2k} - \hat{\mathbf{x}}^*\|_2^2 \geq \frac{q^{2k}}{2} \|\mathbf{x}_0 - \hat{\mathbf{x}}^*\|_2^2 \geq \frac{1}{2} (1 - \sqrt{\alpha})^{2k} \|\mathbf{x}_0 - \hat{\mathbf{x}}^*\|_2^2.$$

By using that $\alpha/4 \leq 1$, we also have the upper bound

$$q(\alpha) = 1 + \alpha/2 - \sqrt{\alpha^2/4 + \alpha} \leq 1 + \alpha/2 - \sqrt{\alpha^2/2} = 1 - \frac{\sqrt{2}-1}{2}\alpha.$$

Then, applying the upper bound and second claim of Lemma 3.7.4, we have that

$$\begin{aligned}\|\mathbf{x}_k - \mathbf{x}^*\|_2^2 &\geq \frac{1}{4\sqrt{2}} (1 - \sqrt{\alpha})^k \|\mathbf{x}_0 - \mathbf{x}^*\|_2^2 - \frac{3q^{2d+2}}{2\alpha(1-q)} \\ &= \frac{1}{4\sqrt{2}} (1 - \sqrt{\alpha})^k \|\mathbf{x}_0 - \mathbf{x}^*\|_2^2 - O((1 - \alpha)^{2d+2}/\alpha).\end{aligned}$$

□

Primal-Linear Lower Bounds While all upper bounds can be proved with perfect symmetry in the primal-linear setting, this is not obvious in the case of the lower bounds due to the condition $n \geq d$. We modify the argument of Section 3.7 to establish hardness guarantees

for the problem

$$\min_{\mathbf{x} \in \mathcal{X}} \max_{\mathbf{y} \in \mathcal{Y}} [\mu(\mathbf{x}, \mathbf{y}) := \langle g(\mathbf{y}), \mathbf{x} \rangle - \psi(\mathbf{y}) + \phi(\mathbf{x})], \quad (3.115)$$

where $\mathcal{X} \subseteq \{\mathbf{x} \in \mathbb{R}^d : x_l \geq l \text{ if } g_l \text{ is not affine}\}$ and the components of $g : \mathcal{Y} \rightarrow \mathcal{X}^*$ satisfy an analogous form of Assumption 3.2.2.

Assumption 3.7.1. Assume that each component of $g = (g_1, \dots, g_d)$ is convex, and there exist $G > 0$ and $L \geq 0$ such that for all $\mathbf{y} \in \mathcal{Y}$,

$$\|\nabla g(\mathbf{y})\|_{\mathcal{Y} \rightarrow \mathcal{X}^*} := \sup \{\|\nabla g(\mathbf{y})\bar{\mathbf{y}}\|_{\mathcal{X}^*} : \bar{\mathbf{y}} \in \mathcal{Y}, \|\bar{\mathbf{y}}\|_{\mathcal{Y}} = 1\} \leq G.$$

and

$$\left\| \sum_{l=1}^d x_j (\nabla g_l(\mathbf{y}) - \nabla g_l(\mathbf{y}')) \right\|_{\mathcal{Y}^*} \leq L \|\mathbf{y} - \mathbf{y}'\|_{\mathcal{Y}} \text{ for all } \mathbf{x} \in \mathcal{X}. \quad (3.116)$$

Here, we still assume $n \geq d$, so that for this primal-linear setting, the linear variable is of smaller dimension than the nonlinear dimension. We also maintain Assumption 3.2.1. We still pursue a bound of the form

$$\|\mathbf{x}_k - \hat{\mathbf{x}}^*\|_{\mathcal{X}}^2 \geq C(1 - \alpha) \|\mathbf{x}_0 - \hat{\mathbf{x}}^*\|_{\mathcal{X}}^2,$$

where C is an absolute constant and $\alpha = \min \{\sqrt{\mu\nu}/G, \nu/L\}$, where $\hat{\mathbf{x}}^*$ arbitrarily close approximation of the optimal primal solution \mathbf{x}^* . Consequently, the parameter regime for this lower bound to be tight is $\nu G^2 \geq \mu L^2$.

To lay the groundwork, we use an analogous algorithm class to Definition 3.7.1.

Definition 3.7.2 (Primal-Linear Proximal Algorithm Class). We define a (primal-linear) *deterministic proximal gradient algorithm* to be a sequence of primal-dual iterates $(\mathbf{x}_k, \mathbf{y}_k)_{k \geq 0}$ satisfying $(\mathbf{x}_k, \mathbf{y}_k) \in \mathcal{H}_k^{\mathcal{X}} \times \mathcal{H}_k^{\mathcal{Y}}$, where $\mathcal{H}_k^{\mathcal{X}} \subseteq \mathcal{X}$ and $\mathcal{H}_k^{\mathcal{Y}} \subseteq \mathcal{Y}$ are constructed from the following

steps. First, for $k \geq 0$, define the subspaces

$$\mathcal{G}_k^{\mathcal{X}} = \text{span} \{g(\mathbf{y}_i) : 0 \leq i \leq k\} \subseteq \mathcal{X}^*, \quad (3.117)$$

$$\mathcal{G}_k^{\mathcal{Y}} = \text{span} \{\nabla g(\mathbf{y}_i)^\top \mathbf{x}_i : 0 \leq i \leq k+1\} \subseteq \mathcal{Y}^*. \quad (3.118)$$

Define $\mathcal{H}_0^{\mathcal{X}} = \text{span}\{\mathbf{x}_0\}$ and $\mathcal{H}_0^{\mathcal{Y}} = \text{span}\{\mathbf{y}_0\}$, and for $k \geq 1$, we recursively define

$$\begin{aligned} \mathcal{H}_k^{\mathcal{X}} &= \text{span} \left(\{p_{\eta_k \phi}(\tilde{\mathbf{x}}_{k-1}, \tilde{\mathbf{g}}_{k-1}) : \tilde{\mathbf{x}}_{k-1} \in \mathcal{H}_{k-1}^{\mathcal{X}}, \tilde{\mathbf{g}}_{k-1} \in \mathcal{G}_{k-1}^{\mathcal{X}}\} \cup \mathcal{H}_{k-1}^{\mathcal{X}} \right) \\ \mathcal{H}_k^{\mathcal{Y}} &= \text{span} \left(\{p_{\gamma_k \psi}(\tilde{\mathbf{y}}_{k-1}, \tilde{\mathbf{f}}_k) : \tilde{\mathbf{y}}_{k-1} \in \mathcal{H}_{k-1}^{\mathcal{Y}}, \tilde{\mathbf{f}}_{k-1} \in \mathcal{G}_{k-1}^{\mathcal{Y}}\} \cup \mathcal{H}_{k-1}^{\mathcal{Y}} \right). \end{aligned}$$

We also use the same exact hard instance as before, namely

$$\mu(\mathbf{x}, \mathbf{y}) = \underbrace{\left\langle \frac{G}{2} \mathbf{A}^\top \mathbf{y}, \mathbf{x} \right\rangle}_{\langle g(\mathbf{y}), \mathbf{x} \rangle} - \underbrace{\frac{\nu}{2} \|\mathbf{y}\|_2^2}_{\psi(\mathbf{y})} + \underbrace{\mathbf{c}^\top \mathbf{x} + \frac{\mu}{2} \|\mathbf{x}\|_2^2}_{\phi(\mathbf{x})}, \quad (3.119)$$

with Euclidean proximal oracles. In this case, just as before, we have for all $k \geq 1$,

$$\mathbf{x}_k = \frac{\eta_k, \mathbf{x}_{k-1} - \bar{\mathbf{g}}_{k-1}}{\mu + \eta_k} \subseteq \text{span} \{\mathbf{x}_{k-1}, \mathbf{c}, \mathbf{A}^\top \mathbf{y}_{k-1}, \mathbf{A}^\top \mathbf{y}_{k-2}, \dots, \mathbf{A}^\top \mathbf{y}_0\} \quad (3.120)$$

$$\mathbf{y}_k = \frac{\gamma_k, \mathbf{y}_{k-1} - \bar{\mathbf{f}}_{k-1}}{\nu + \gamma_k} \subseteq \text{span} \{\mathbf{y}_{k-1}, \mathbf{A} \mathbf{x}_k, \mathbf{A} \mathbf{x}_{k-1}, \dots, \mathbf{A} \mathbf{x}_0\}. \quad (3.121)$$

As a result, we have that Lemma 3.7.3 and Corollary 3.7.1 follow in this setting as well, as they are fully determined by the equations displayed above. After this, the argument follows *exactly* to Section 3.7.2. This results in the following claim.

Proposition 3.7.2. *Consider any algorithm satisfying Definition 3.7.2 with $\mathbf{x}_0 = \mathbf{0}_d$ and $\mathbf{y}_0 = \mathbf{0}_n$ applied to the problem (3.119), with $n \geq d \geq k$. Consider $\alpha = \frac{4\mu\nu}{G^2} \leq 1$. We have that*

$$\text{Gap}^{\mathbf{x}^*, \mathbf{y}^*}(\mathbf{x}_k, \mathbf{y}_k) \geq \frac{\mu}{2} \|\mathbf{x}_k - \mathbf{x}^*\|_2^2 \geq \frac{\mu}{8\sqrt{2}} (1 - \sqrt{\alpha})^k \|\mathbf{x}_0 - \mathbf{x}^*\|_2^2 - O((1 - \alpha)^{2d+2})/\alpha$$

Proof. Follow an identical proof to Proposition 3.7.1. □

Having established preliminary lower bounds for deterministic algorithms, extending

these constructions to the setting in which L/μ comprises the dominant term in the objective would be a natural direction for future research. We discuss this and other directions in Section 3.8.

3.8 Perspectives & Future Work

This chapter presented a deep dive into the semilinear min-max problem, which interpolates between two problem classes that have been fertile ground for optimization research for decades: the bilinearly and nonbilinearly-coupled min-max problems. We employed a constructive convergence analysis and derived algorithms for the convex-concave, strongly convex-strongly concave, strongly convex-concave, and convex-strongly concave in a unified manner. Key ideas were adaptive sampling, which made use of possibly non-uniform Lipschitz and smoothness constants among the component functions f_1, \dots, f_n , and a historical regularization term on the primal update that can achieve complexity improvements in even special cases such as bilinear problems.

Because the focus of this work was on deriving upper bounds on complexity, there are opportunities to derive tight lower bounds for both deterministic and stochastic settings. The constructions generally differ for stochastic methods (e.g., that of [Woodworth and Srebro \[2016\]](#)) from the Nesterov tri-diagonal approach taken in Section 3.7. Furthermore, our lower bound only applies to particular parameter regimes (when $G/\sqrt{\mu\nu} \gtrsim L/\mu$). To extend this analysis to arbitrary parameter regimes, we hypothesize that our high-dimensional quadratic objective would be replaced by a Lipschitz continuous and smooth function, such as a multivariate analog of the Huber loss, such as

$$\ell_\delta(\mathbf{x}) = \begin{cases} \frac{1}{2} \|\mathbf{x}\|_2^2 & \|\mathbf{x}\|_2 \leq \delta \\ \delta (\|\mathbf{x}\|_2 - \frac{1}{2}\delta) & \|\mathbf{x}\|_2 > \delta \end{cases}.$$

An interesting extension from both theoretical and algorithmic viewpoints is the possibility to solve a generic nonbilinearly-coupled optimization problem by alternating between 1) linearizing the objective in the dual variables, and 2) solving the resulting dual linear

min-max problem. Fixing the strongly convex-strongly concave setting for concreteness, and recalling the formulation (3.1), one may consider $(\bar{\mathbf{x}}_k, \bar{\mathbf{y}}_k)$ a solution of

$$\min_{\mathbf{x} \in \mathcal{X}} \max_{\mathbf{y} \in \mathcal{Y}} \langle \mathbf{y}, \nabla_{\mathbf{y}} c(\mathbf{x}, \bar{\mathbf{y}}_{k-1}) \rangle - \psi(\mathbf{y}) + \phi(\mathbf{x}), \quad (3.122)$$

and inspect the limiting value as $k \rightarrow \infty$.

Returning to the outline of Chapter 1, we devoted significant attention to the topics introduced in Section 1.2.2. While the relationship between distributional robustness and convex optimization (e.g., statistical learning with linear models) allows for elegant mathematical analyses, it is also critical to consider highly expressive models such as neural networks (as noted in Section 2.12). Furthermore, core statistical questions regarding generalization remain, such as the optimal selection of training data or the resulting sample complexity of modern learning methods. Both points are addressed in the upcoming chapter using the backdrop of Section 1.2.3).

Reference	Global Complexity
Nesterov and Scramali [2006, Theorem 3]	$O\left(nd\sqrt{\frac{L^2+G^2}{\mu\wedge\nu}}\ln(1/\varepsilon)\right)$
Wang and Li [2020, Theorem 3]	$\tilde{O}\left(nd\sqrt{\frac{L}{\mu} + \frac{G\cdot\max\{L,G\}}{\mu\nu}}\ln(1/\varepsilon)\right)$
Lin et al. [2020, Theorem 9] Borodich et al. [2024, Theorem 3]	$\tilde{O}\left(nd\sqrt{\frac{L^2+G^2}{\mu\nu}}\ln(1/\varepsilon)\right)$
Carmon et al. [2022, Theorem 3] Lan and Li [2023, Theorem 4.2]	$O\left(nd\sqrt{\frac{L^2+G^2}{\mu\nu}}\ln(1/\varepsilon)\right)$
Kovalev and Gasnikov [2022, Theorem 3]	$O\left(nd\sqrt{\frac{L^2+G^2}{\mu(\mu\wedge\nu)}}\ln(1/\varepsilon)\right)$
Jin et al. [2022, Theorem 1] Li et al. [2023, Corollary 3.4]	$\tilde{O}\left(nd\left(\frac{L}{\mu} + \frac{G}{\sqrt{\mu\nu}}\right)\ln(1/\varepsilon)\right)$
This work (Theorem 3.3.1)	$O\left(nd\left(\frac{L}{\mu} + \frac{G}{\sqrt{\mu\nu}}\right)\ln(1/\varepsilon)\right)$

Table 3.2: **Complexity Bounds for General Nonbilinarily-Coupled Objectives** for $\mu, \nu > 0$. Runtime or global complexity (i.e. the total number of elementary operations required to compute (\mathbf{x}, \mathbf{y}) satisfying $\text{Gap}^{\mathbf{u}, \mathbf{v}}(\mathbf{x}, \mathbf{y}) \leq \varepsilon$ for fixed $(\mathbf{u}, \mathbf{v}) \in \mathcal{X} \times \mathcal{Y}$. The methods considered call the entire list of primal first-order oracles $(f_j, \nabla f_j)$ for $j = 1, \dots, n$ on each iteration. The method of Kovalev and Gasnikov [2022, Corollary 1] achieves its claim by swapping the role of \mathbf{x} and \mathbf{y} , which is not possible for (3.2). Carmon et al. [2022, Theorem 3] requires a bounded diameter assumption on \mathcal{Y} , which is generally required for $\mathbf{L}_1, \dots, \mathbf{L}_N$ to be finite if each of f_1, \dots, f_n is nonlinear.

Reference	Additional Structure	Global Complexity (big- \tilde{O})
Alacaoglu and Malitsky [2022, Corollary 6] Cai et al. [2024, Theorem 4.2]	Constants known	$\frac{n(d+N)}{\sqrt{N}\varepsilon} \ \boldsymbol{\lambda}\ _1$
Alacaoglu and Malitsky [2022, Corollary 6] Cai et al. [2024, Theorem 4.2] Pichugin et al. [2024, Corollary 1]		$\frac{n(d+N)}{\varepsilon} \ \boldsymbol{\lambda}\ _2$
Diakonikolas [2025, Thm. 1 & Eq. (38)]	Constants known + Separable	$\frac{n(d+N)}{N\varepsilon} \ \boldsymbol{\lambda}\ _{1/2}$ $\frac{nd}{N\varepsilon} \ \boldsymbol{\lambda}\ _{1/2}$
Diakonikolas [2025, Thm. 1 & Eq. (36)]	+ Separable	$\frac{n(d+N)}{\varepsilon} \sqrt{N} \ \boldsymbol{\lambda}\ _2$ $\frac{nd}{\varepsilon} \sqrt{N} \ \boldsymbol{\lambda}\ _2$
This work (Theorem 3.4.1)	Constants known	$\frac{n(d+N)}{N\varepsilon} \left(\ \boldsymbol{\lambda}\ _1^{1/2} (\ \boldsymbol{\lambda}\ _1^{1/2} + \ \mathbf{G}\ _1^{1/2}) \right)$
This work (Theorem 3.5.1)	+ Separable	$\frac{nd}{\sqrt{N}\varepsilon} \left(\ \boldsymbol{\lambda}\ _1^{1/2} (\ \mathbf{L}\ _1^{1/2} + \ \mathbf{G}\ _\infty^{1/2}) \right)$
This work (Theorem 3.4.2)	Constants known	$\frac{n(d+N)}{N\varepsilon} \left(\ \boldsymbol{\lambda}\ _{1/2}^{1/2} (\ \boldsymbol{\lambda}\ _{1/2}^{1/4} \ \mathbf{L}\ _{1/2}^{1/4} + \ \mathbf{G}\ _{1/2}^{1/2}) \right)$
This work (Theorem 3.5.2)	+ Separable	$\frac{nd}{\sqrt{N}\varepsilon} \left(\ \boldsymbol{\lambda}\ _{1/2}^{1/4} (\ \mathbf{L}\ _{1/2}^{3/4} + \ \mathbf{G}\ _\infty^{3/4}) \right)$

Table 3.3: **Complexity Bounds for Convex-Concave Finite-Sum Objectives.** Arithmetic or global complexity (i.e., the total number of elementary operations required to compute (\mathbf{x}, \mathbf{y}) satisfying $\mathbb{E}[\text{Gap}^{u,v}(\mathbf{x}, \mathbf{y})] \leq \varepsilon$ for fixed $(\mathbf{u}, \mathbf{v}) \in \mathcal{X} \times \mathcal{Y}$. We use $\boldsymbol{\lambda}$ as defined in Section 3.2. The objective is assumed to be **convex-concave** and have a finite sum structure $\mathcal{L}(\mathbf{x}, \mathbf{y}) = \sum_{j=1}^N \mathcal{L}_j(\mathbf{x}, \mathbf{y})$. The expectation is taken over any randomness incurred by the algorithm.

Reference	Additional Structure	Global Complexity (big- \tilde{O})
Palaniappan and Bach [2016, Theorem 2]	Constants known	$\frac{n(n+d)}{N} \left(N + \frac{L^2+G^2}{(\mu\wedge\nu)^2} \right) \ln(1/\varepsilon)$
Alacaoglu and Malitsky [2022, Corollary 27] Cai et al. [2024, Theorem 4.6]	Constants known	$\frac{n(d+N)}{N} \left(N + \frac{\sqrt{N}\ \lambda\ _1}{\mu\wedge\nu} \right) \ln \left(\frac{1}{\varepsilon} \right)$
Alacaoglu and Malitsky [2022, Corollary 27] Cai et al. [2024, Theorem 4.6]		$\frac{n(d+N)}{N} \left(N + \frac{N\ \lambda\ _2}{\mu\wedge\nu} \right) \ln \left(\frac{1}{\varepsilon} \right)$
Diakonikolas [2025, Thm. 1 & Eq. (38)]	Constants known + Separable	$\frac{n(d+N)}{N} \left(N + \frac{\ \lambda\ _{1/2}}{\mu\wedge\nu} \right) \ln \left(\frac{1}{\varepsilon} \right)$ $\frac{nd}{N} \left(N + \frac{\ \lambda\ _{1/2}}{\mu\wedge\nu} \right) \ln \left(\frac{1}{\varepsilon} \right)$
Diakonikolas [2025, Thm. 1 & Eq. (36)]	+ Separable	$\frac{n(d+N)}{N} \left(N + \frac{N^{3/2}\ \lambda\ _2}{\mu\wedge\nu} \right) \ln \left(\frac{1}{\varepsilon} \right)$ $\frac{nd}{N} \left(N + \frac{N^{3/2}\ \lambda\ _2}{\mu\wedge\nu} \right) \ln \left(\frac{1}{\varepsilon} \right)$
This work (Theorem 3.4.1)	Constants known	$\frac{n(d+N)}{N} \left(N + \frac{\sqrt{N}\ \lambda\ _1}{\mu} + \frac{\sqrt{N}\ \lambda\ _1^{1/2}\ \mathbf{G}\ _1^{1/2}}{\sqrt{\mu\nu}} \right) \ln \left(\frac{1}{\varepsilon} \right)$
This work (Theorem 3.5.1)	+ Separable	$\frac{nd}{N} \left(N + \frac{\sqrt{N}\ \lambda\ _1^{1/2}\ \mathbf{L}\ _1^{1/2}}{\mu} + \frac{\sqrt{N}\ \lambda\ _1^{1/2}\ \mathbf{G}\ _\infty^{1/2}}{\sqrt{\mu\nu}} \right) \ln \left(\frac{1}{\varepsilon} \right)$
This work (Theorem 3.4.2)	Constants known	$\frac{n(d+N)}{N} \left(N + \frac{\ \lambda\ _{1/2}^{3/4}\ \mathbf{L}\ _{1/2}^{1/4}}{\mu} + \frac{\ \lambda\ _{1/2}^{1/2}\ \mathbf{G}\ _{1/2}^{1/2}}{\sqrt{\mu\nu}} \right) \ln \left(\frac{1}{\varepsilon} \right)$
This work (Theorem 3.5.2)	+ Separable	$\frac{nd}{N} \left(N + \frac{\ \lambda\ _{1/2}^{1/4}\ \mathbf{L}\ _{1/2}^{3/4}}{\mu} + \frac{\ \lambda\ _{1/2}^{1/4}\ \mathbf{G}\ _\infty^{3/4}}{\sqrt{\mu\nu}} \right) \ln \left(\frac{1}{\varepsilon} \right)$

Table 3.4: **Complexity Bounds for Strongly Convex-Strongly Concave Finite-Sum Objectives:** Arithmetic or global complexity (i.e., the total number of elementary operations required to compute (\mathbf{x}, \mathbf{y}) satisfying $\mathbb{E}[\text{Gap}^{\mathbf{u}, \mathbf{v}}(\mathbf{x}, \mathbf{y})] \leq \varepsilon$ for fixed $(\mathbf{u}, \mathbf{v}) \in \mathcal{X} \times \mathcal{Y}$. We use λ as defined in Section 3.2. The objective is assumed to be (μ, ν) -**strongly convex-strongly concave** and have a finite sum structure $\mathcal{L}(\mathbf{x}, \mathbf{y}) = \sum_{j=1}^N \mathcal{L}_j(\mathbf{x}, \mathbf{y})$. The expectation is taken over any randomness in the algorithm.

Chapter 4

GENERALIZATION CAPABILITIES OF ZERO-SHOT PREDICTION

4.1 Introduction

In Chapter 2 and Chapter 3, we studied statistical and algorithmic aspects of out-of-distribution generalization from the perspective of linear models and convex optimization. While this setting allows for acute statements of algorithm performance, an undisputed driver of modern generalization capabilities—such as zero-shot prediction (ZSP)—is the expressivity of overparametrized neural network models (see Section 1.2.3). Nonetheless, we note that neural networks (in particular, transformers [Vaswani et al., 2017]) merely provide a canvas; the primary driver of performance in both ZSP and language modeling is the data used to pre-train the model [Fang et al., 2023, Gadre et al., 2023, Xu et al., 2024, Li et al., 2024a]. Furthermore, the prompting mechanism (as used in (1.7)), though relatively less-studied, has been seen to be an essential determinant of downstream performance [Pratt et al., 2023]. We now frame our approach to addressing them.

Let us first describe the pre-training phase. Consider random variables X and Z governed by a joint probability measure $P \equiv P_{X,Z}$ over $\mathcal{X} \times \mathcal{Z}$. We may interpret (X, Z) as an image-caption pair. Given this data-generating distribution, we assume access to $(\mathbf{x}_1, \mathbf{z}_1), \dots, (\mathbf{x}_n, \mathbf{z}_n)$ drawn i.i.d. realization from P with empirical measure $P_n = \frac{1}{n} \sum_{i=1}^n \delta_{(\mathbf{x}_i, \mathbf{z}_i)}$. As mentioned in Section 1.2.1, the pre-training phase involves learning a parameter θ by minimizing a self-supervised learning (SSL) loss of the form (1.5). Specifically, losses such as CLIP/InfoNCE [van den Oord et al., 2019], VICReg [Bardes et al., 2022], BarlowTwins [Zbontar et al., 2021], and spectral embedding methods [HaoChen et al., 2021, Balestriero

and LeCun, 2022, Tan et al., 2024] can be written in the form

$$\mathcal{R}(\boldsymbol{\theta}, P_n) = \mathbb{E}_{P_n^{\otimes b}} [\ell_b(\boldsymbol{\theta}, (X_1, Z_1), \dots, (X_b, Z_b))], \quad (4.1)$$

where $P_n^{\otimes b}$ denotes sampling b points uniformly randomly with replacement from $\{(\mathbf{x}_i, \mathbf{z}_i)\}_{i=1}^n$ and ℓ_b is a loss function that, at a high level, involves comparison between points in the mini-batch. The specific forms of the loss are not relevant for our main investigation, in which we analyze a prominent method used to design the pre-training set of the original CLIP model [Radford et al., 2021] as well as open-source reproductions [Xu et al., 2024]. In the creation of both models, the scientists first resampled the pre-training set in order to achieve particular marginal distributions over metadata (such as word counts in the text modality). We prove that when these marginal distributions coincide with the true data-generating distribution, linear functionals (such as expected loss) can be estimated with a reduced variance via a non-asymptotic mean squared error bound. The bound decomposes into an $O(1/n)$ variance term, an $\tilde{O}(1/n^2)$ bias term, and an $\tilde{O}(1/n^{3/2})$ cross term. The first-order variance term recovers the asymptotic variance previously discovered by Bickel et al. [1991] by the tools of asymptotic semiparametric efficiency theory. Our bound not only explicitly exposes the higher-order terms non-asymptotically, but also displays a new closed-form formula for the variance term, which was previously described in terms of projections. We hypothesize, and confirm with experimentation, that this procedure may stabilize the loss estimate (4.1) and lead to improved downstream performance of the resulting models.

For the downstream task, the practitioner aims to use the parameter $\boldsymbol{\theta}$ and an auxiliary prompting procedure (described in Section 1.2.3). We prove a basic identity that relates (1.7) to a two-stage regression function (from X to Z to an unseen label Y). The same identity allows us to describe the fundamental limit of ZSP, i.e., how close its performance may approximate the Bayes optimal predictor for the downstream task, in terms of a conditional dependence relationship on (X, Z, Y) .

The analyses above are formalized and elaborated upon in Section 4.2. Results regarding (upstream) data curation through marginal matching and the improved estimation of

linear functionals are contained in Section 4.3. The connection between ZSP and two-stage regression is proved, and its implications are given in Section 4.4. Experimental results that validate our viewpoints on both improving downstream performance with marginal-corrected pre-training datasets and optimized prompting are contained in Section 4.5. Extensions and future work are discussed in Section 4.6 and Section 4.7.

4.2 Preliminaries

Recall that we used $P_{X,Z}$ to denote the joint probability measure of the image-caption pairs in the pre-training set (e.g. arbitrary content from the Internet that is not necessarily associated with any task). For a downstream task, we consider an additional random variable Y realized in a (possibly non-discrete) set \mathcal{Y} , indicating a category or class label. We define a joint distribution $P^\mathcal{T} \equiv P_{X,Z,Y}^\mathcal{T}$ governing the data from the downstream task (e.g., image classification on CIFAR-10). Crucially, notice that we consider a random variable Z governed by $P^\mathcal{T}$, representing a hypothetical (but unobserved) caption for the image being classified. Instead, this latent variable will be useful for defining a theoretical counterpart of the ZSP procedure (1.7). The first part of this chapter will only concern the pre-training phase and $P_{X,Z}$, whereas the second part will consider the downstream task and $P_{X,Z,Y}^\mathcal{T}$. We outline the main approaches below.

4.2.1 Curating Pre-Training Data

Firstly, to understand the effects of data curation on pre-training, we formally analyze the marginal matching method of Radford et al. [2021], Xu et al. [2024] introduced in Section 4.1. For this example, first assume that \mathcal{X} and \mathcal{Z} are finite sets with $|\mathcal{X}| = m$ and $|\mathcal{Z}| = l$, as the method is applied to discretized versions of the original data in practice. The general approach is to select two user-defined marginal distributions P_X on \mathcal{X} and P_Z on \mathcal{Z} and “adjust” the empirical measure P_n so that it marginalizes to (P_X, P_Z) in either variable. In our statistical analysis, we assume that P_X and P_Z are in fact the true marginal distributions of the data-generating distribution P . Under a frequentist viewpoint, this captures a

setting in which high-quality, paired examples are relatively scarce or expensive (i.e., sampling or integration under P is “hard”), but unpaired examples are so abundant that one can estimate the marginals with negligible error. This was the historical justification for the setting, studied as early as 1940 in the context of cross-tabulated data for census applications [Deming and Stephan, 1940, Ireland and Kullback, 1968]. From a Bayesian viewpoint, the marginals P_X and P_Z can encode prior information, inductive bias, or side information that a distribution estimator may incorporate for improved statistical accuracy [Miller and Liu, 2002]. The marginal matching done in the works above is equivalent to estimating P using some number of steps of the following recursion: we define $P_n^{(0)} = P_n$ as the empirical measure, and for $k \geq 1$ construct

$$P_n^{(k)}(\mathbf{x}, \mathbf{z}) := \begin{cases} \frac{P_X(\mathbf{x})}{P_{n,X}^{(k-1)}(\mathbf{x})} \cdot P_n^{(k-1)}(\mathbf{x}, \mathbf{z}) & k \text{ odd} \\ \frac{P_Z(\mathbf{z})}{P_{n,Z}^{(k-1)}(\mathbf{z})} \cdot P_n^{(k-1)}(\mathbf{x}, \mathbf{z}) & k \text{ even} \end{cases}. \quad (4.2)$$

The conditions under which the iterations of (4.2) are well-defined are given in Section 4.3. These operations reduce to rescaling the rows of an $(m \times l)$ -matrix by $P_X/P_{n,X}^{(k-1)}$ or its columns by $P_Z/P_{n,Z}^{(k-1)}$. This very algorithm has a decades-old history and is known in various contexts as Sinkhorn-Knopp matrix scaling [Sinkhorn, 1967], iterative proportional or bi-proportional fitting [Johnston and Pattie, 1993], and raking ratio estimation [Thompson, 2000]. We hypothesize that the improvements of the data “balancing” procedure (4.2) in large-scale training are derived from the improvement of $P_n^{(k)}$ over P_n as an estimator for P (at a linear functional representing the expected loss). In Section 4.3, we prove non-asymptotic bounds on the mean squared error of the estimator $\mathbb{E}_{P_n^{(k)}}[h(X, Z)]$ for $\mathbb{E}_P[h(X, Z)]$ for an integrable test function $h : \mathcal{X} \times \mathcal{Z} \rightarrow \mathbb{R}$. The bound decomposes into an $O(n^{-1})$ first-order variance term and an $\tilde{O}(k^6 n^{-3/2})$ higher-order term. The first-order term is shown to have a strict improvement over the empirical measure baseline with a fine-grained dependence on the spectra of two conditional mean operators associated to P . The higher-order term can be used to compute the asymptotic variance, resolving some efficiency questions originally

raised in [Bickel et al. \[1991\]](#). Our proof technique relies on a recursion decomposition for balancing-based estimators, which may be of independent interest.

While we focus on the applications of balancing to data curation, we also comment that for specific foundation models, such as CLIP [[Radford et al., 2021](#)], SwaV [[Caron et al., 2020](#)], and variations on this theme [[Jones et al., 2022](#), [Asano et al., 2020](#)] also contain a balancing operation in the computation of their objectives (see [Liu et al. \[2024, Section 2\]](#)). While not originally derived in this way, in the particular case of CLIP, the objective is computed by applying the iterations of (4.2) to an unnormalized measure defined on the elements of the *mini-batch*. In fact, only $k = 1$ iteration is applied to each marginal in the original CLIP objective, whereas [Liu et al. \[2024\]](#) observe performance improvements by increasing the number of iterations and further committing to the viewpoint that the objective implicitly computes (4.2). We identify this intersection of balancing and learning objectives as a valuable line of future work.

4.2.2 Prompting and Downstream Classification

Secondly, we consider the downstream task of predicting Y from X , and specifically, using a predictor of the form (1.7) which relies on pre-trained encoders. Note that we do not use the finiteness assumption on \mathcal{X} or \mathcal{Z} in these results. We identify a population parameter and describe how each component of the pipeline (including prompting) is composed to serve as an estimator of the parameter. Let $f : \mathcal{Y} \rightarrow \mathbb{R}$ be a function and consider the problem of predicting $f(Y)$ given $X = \mathbf{x}$. The use of f serves only to handle multiple tasks such as classification (e.g., $f(\mathbf{y}) = \mathbb{1}\{\mathbf{y} = 1\}$) and regression ($f(\mathbf{y}) = y$) in a unified manner. Now, define the minimum $P_{X,Y}^\mathcal{T}$ -expected risk (Bayes optimal) predictor as

$$\eta_\star(\mathbf{x}) := \mathbb{E}_{P_{X,Y}^\mathcal{T}} [f(Y)|X](\mathbf{x}). \quad (4.3)$$

Because ZSP operates by “translating” the image classification problem to a text classification problem via prompting, it is natural to consider the entire procedure (1.7) to be a

finite-sample estimator of the two-stage regression function

$$\eta_\rho(\mathbf{x}) := \mathbb{E}_{P_{X,Z}} [\mathbb{E}_{\rho_{Y,Z}} [f(Y)|Z] | X] (\mathbf{x}), \quad (4.4)$$

where $\rho_{Y,Z}$ is the *prompting* distribution. This probability measure represents the pseudo-captions used to define the predictor (1.7), which include template-based prompts as well as class-conditional prompts (those for which the caption text differs for each class label) drawn from large language models. To support the claim of (4.4) as the underlying estimand mathematically, we prove a key relationship in Section 4.4. Under the assumption that $P_{X,Z} \ll P_X \otimes P_Z$, we define the likelihood ratio (or Radon-Nikodym derivative) $R(\mathbf{x}, \mathbf{z}) = \left(\frac{dP_{X,Z}}{d(P_X \otimes P_Z)} \right) (\mathbf{x}, \mathbf{z})$. Then, it holds that

$$\eta_\rho(\mathbf{x}) = \mathbb{E}_{\rho_{Y,Z}} [f(Y) \cdot R(\mathbf{x}, Z)] + \text{err}(P_Z, \rho_Z), \quad P_X\text{-almost surely}, \quad (4.5)$$

where $\text{err}(P_Z, \rho_Z)$ measures the disagreement between the distribution of the pre-training captions and the captions drawn from the prompting distribution determined by the user. To understand the significance of this result, recall the encoders $(\boldsymbol{\alpha}, \boldsymbol{\beta})$ introduced in (1.7). When parameterized via $\boldsymbol{\theta}$ as $(\boldsymbol{\alpha}, \boldsymbol{\beta}) \equiv (\boldsymbol{\alpha}_\theta, \boldsymbol{\beta}_\theta)$, they are produced by optimizing the self-supervised learning (SSL) objective (4.1) over a large pre-training set. Several works (see Gutmann and Hyvärinen [2012], Oko et al. [2025] and references therein) have established that popular SSL objectives such as the CLIP and Noise Contrastive Estimation achieve the exact population minimum when $\langle \boldsymbol{\alpha}(\mathbf{x}), \boldsymbol{\beta}(\mathbf{z}) \rangle = \log R(\mathbf{x}, \mathbf{z})$, assuming it exists. This form (4.5) can be related to (1.7), as is done in Section 4.4, connecting an otherwise mysterious empirical procedure to a well-defined estimand with clear limits. In other words, we are able to use the two-stage prediction form (4.4) to comment on the mathematical limits of prompting as a replacement for downstream training data in Section 4.4.

4.3 Non-Asymptotic Analysis of Variance Reduction

We expand on the problem setup introduced in Section 4.2.1, in which we consider sample spaces $(\mathcal{X}, \mathcal{Z})$, along with true and unknown joint distribution P on $\mathcal{X} \times \mathcal{Z}$ with known marginals (P_X, P_Z) . Recall the ERM notation of Section 1.2.1. For an integrable test function $h : \mathcal{X} \times \mathcal{Z} \rightarrow \mathbb{R}$, we introduce the empirical process notation $P(h) := \mathbb{E}_P[h(X, Z)]$. We then let $\Theta = \mathbb{R}$ and define the squared error risk functional

$$\mathcal{R}(\boldsymbol{\theta}, P) = (\boldsymbol{\theta} - P(h))^2.$$

Thus, considering both the empirical measure P_n and the marginal-rebalanced empirical measure $P_n^{(k)}$ from (4.2), we define

$$\boldsymbol{\theta}_n := P_n(h) = \arg \min_{\boldsymbol{\theta} \in \Theta} \mathcal{R}(\boldsymbol{\theta}, P_n) \text{ and } \boldsymbol{\theta}_n^{(k)} := P_n^{(k)}(h) = \arg \min_{\boldsymbol{\theta} \in \Theta} \mathcal{R}(\boldsymbol{\theta}, P_n^{(k)}) \quad (4.6)$$

We present theoretical guarantees on the mean squared error (MSE) of the data-balanced estimator $\boldsymbol{\theta}_n^{(k)}$ (taken over the randomness in the sample) and highlight relevant points in the proofs. We first give context on the main innovations of the analysis and then outline its high-level steps. These innovations include relating the nonlinear iterations of balancing over probability measures to linear operators on a vector space and using a singular value decomposition of these operators to quantify their effect after a finite number of iterations. Furthermore, by scaling the number of iterations appropriately, we can characterize the estimator using the limit of balancing iterations, which is an object of interest in applications including optimal transport.

We make the following assumption throughout, which is usually satisfied by the desired marginals P_X and P_Z in practice: the target marginals $P_X(\mathbf{x}) > 0$ and $P_Z(\mathbf{z}) > 0$ for all $\mathbf{x} \in \mathcal{X}$ and $\mathbf{z} \in \mathcal{Z}$. The iterative procedure given in (4.2) is visualized in Figure 4.1 (left). The iterations are well-defined for all k under the event that $\text{Supp}(P_{n,X}) = \text{Supp}(P_X)$ and $\text{Supp}(P_{n,Z}) = \text{Supp}(P_Z)$, i.e., all observed row counts and column counts are non-empty.¹

¹Due to this technical consideration, we define $P_n^{(k)}$ to be the empirical measure P_n when this condition

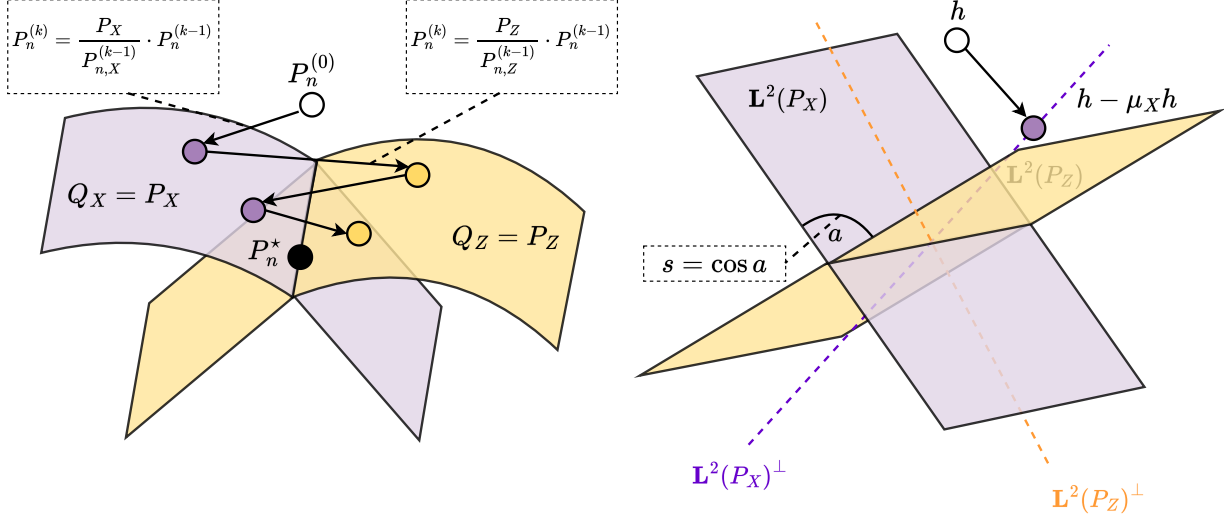


Figure 4.1: **Data Balancing.** Nonlinear and linear operators associated with each iteration of (4.2). **Left:** Visualization of the exact iterations of (4.2) in the space of probability measures. The purple set contains joint distributions with \mathcal{X} -marginal equal to P_X , whereas the golden set contains joint distributions with \mathcal{Z} -marginal equal to P_Z . **Right:** Visualization of $L^2(P)$, the operators defining (4.9), and the singular values given in (4.11).

To provide background, the scheme of alternating the operators (4.2) is often seen as an iterative algorithm to solve the problem

$$\min_{Q \in \Pi(P_X, P_Z)} \text{KL}(Q \| P_n), \quad (4.7)$$

where $\Pi(P_X, P_Z)$ denotes the set of probability measures on $\mathcal{X} \times \mathcal{Z}$ that marginalize to P_X and P_Z in each variable and $\text{KL}(\cdot \| \cdot)$ denotes the Kullback-Leibler divergence. The iterations (4.2) are based on the alternating minimization approach of solving

$$P_n^{(k)}(\mathbf{x}, \mathbf{z}) := \begin{cases} \arg \min_{\{Q: Q_X = P_X\}} \text{KL}(Q \| P_n^{(k-1)}) & k \text{ odd} \\ \arg \min_{\{Q: Q_Z = P_Z\}} \text{KL}(Q \| P_n^{(k-1)}) & k \text{ even} \end{cases},$$

which inspires the viewpoint of balancing as alternating *information projections*. As we

is not satisfied, which we show occurs with low probability. See Appendix B.3.4 for details.

show in Appendix B.2, the iterations of (4.2) can equivalently be defined using the KL, reverse KL, or χ^2 -divergences. This viewpoint is relevant as previously, efforts have been made (e.g., in Bickel et al. [1991]) to analyze the variance reduction afforded by the solution to (4.7) directly. However, quantifying the variance reduction (in terms of properties of P) using this approach is challenging, as there is no closed-form expression for the solution of (4.7). A key mathematical outcome of our analysis is that the closed-form expressions of the projections (4.2) can be used to compute the reduction in mean squared error at each iteration. Thus, by letting $k \equiv k(n) \rightarrow \infty$ (scaled appropriately against n), we can determine the reduction for the solution of (4.7) for large n . This is the subject of Theorem 4.3.1.

From Information Projections to Orthogonal Projections First, we will show that the variance reduction resulting from each nonlinear iteration of (4.2) is associated with a linear operator applied to h . Thus, instead of analyzing the alternating information projections over probability measures, we may use familiar tools to understand alternating orthogonal projections in a vector space. To define them, we first let $\mathbf{L}^2(P)$ to be the set of functions $h : \mathcal{X} \times \mathcal{Z} \rightarrow \mathbb{R}$ satisfying $\mathbb{E}_P[h^2(X, Z)] < \infty$. Even though $\mathcal{X} \times \mathcal{Z}$ is finite, working within $\mathbf{L}^2(P)$ will be analytically convenient. Let $\mathbf{L}^2(P_X)$ be the subspace of $\mathbf{L}^2(P)$ containing functions that only depend on the first argument $\mathbf{x} \in \mathcal{X}$ and define $\mathbf{L}^2(P_Z)$ analogously. These are the solid-colored subspaces in Figure 4.1 (right). Next, let $\mu_X : \mathbf{L}^2(P) \rightarrow \mathbf{L}^2(P_X)$ and $\mu_Z : \mathbf{L}^2(P) \rightarrow \mathbf{L}^2(P_Z)$ be defined as, for any $h \in \mathbf{L}^2(P)$,

$$\mu_X h = \arg \min_{f \in \mathbf{L}^2(P_X)} \mathbb{E}_P [(h(X, Z) - f(X))^2] \implies [\mu_X h](\mathbf{x}, \mathbf{z}) := \mathbb{E}_P [h(X, Z) | X](\mathbf{x})$$

The operator μ_X is an orthogonal projection onto $\mathbf{L}^2(P_X)$. The orthogonal projection operator μ_Z onto $\mathbf{L}^2(P_Z)$ is defined analogously. We may also define the conditional centering/debiasing operators $\mathcal{C}_X = I - \mu_X$ and $\mathcal{C}_Z = I - \mu_Z$, which each project onto the orthogonal complements of $\mathbf{L}^2(P_X)$ and $\mathbf{L}^2(P_Z)$, visualized as subspaces with dotted lines in Figure 4.1 (right). To understand the importance of the conditional mean and debiasing operators, we give a recursive formula that forms the backbone of our analysis. Define $\mu_k = \mu_X$ for k odd

and $\mu_k = \mu_Z$ for k even, and define \mathcal{C}_k similarly. Thus, we have by linearity of expectation that

$$\begin{aligned}
[P_n^{(k)} - P](h) &= [P_n^{(k)} - P](\mathcal{C}_k h) + \overbrace{[P_n^{(k)} - P](\mu_k h)}^{=0} \\
&= [P_n^{(k-1)} - P](\mathcal{C}_k h) + [P_n^{(k)} - P_n^{(k-1)}](\mathcal{C}_k h) \\
&= \underbrace{[P_n^{(0)} - P](\mathcal{C}_1 \dots \mathcal{C}_k h)}_{\text{first-order term}} + \underbrace{\sum_{\ell=1}^k [P_n^{(\ell)} - P_n^{(\ell-1)}](\mathcal{C}_\ell \dots \mathcal{C}_k h)}_{\text{higher-order terms}}. \tag{4.8}
\end{aligned}$$

To justify the first line, we discuss the case when k is odd. Notice that $\mu_X h$ is only a function of X , so its expectation only depends on P_X that is equal to $P_{n,X}^{(k)}$ (the \mathcal{X} -marginal of $P_n^{(k)}$) by (4.2). The last line follows by unrolling the previous step $k - 1$ times. This recursive expansion is proven formally in Proposition B.3.1 in Appendix B.3. Given the expansion, the mean squared error can be computed by taking the expectation of the square of (4.8). We show that the second moment of the first-order term in (4.8) is equal to σ_0^2/n where

$$\sigma_0^2 := \text{Var}(h) \text{ and } \sigma_k^2 := \text{Var}(\mathcal{C}_1 \dots \mathcal{C}_k h) \text{ for } k \geq 1, \tag{4.9}$$

and all other terms are $\tilde{O}(k^6 n^{-3/2})$. Thus, by exactly computing the constant in the dominating term, we may quantify the asymptotic variance reduction. Our first main result concerns the higher-order terms and shows that it is indeed dominated by the first-order term. Note that the empirical mean $\theta_n^{(0)} = \frac{1}{n} \sum_{i=1}^n h(X_i, Z_i)$ is unbiased, and so its MSE is equal to σ_0^2/n . Define in addition

$$p_\star := \min\{\min_{\mathbf{x}} P_X(\mathbf{x}), \min_{\mathbf{z}} P_Z(\mathbf{z})\}$$

which measures the non-uniformity of the target marginals. We have that p_\star is positive because both P_X and P_Z are positive. We now state the first main result.

Theorem 4.3.1. *For a sequence of data balancing estimators $(\theta_n^{(k)})_{k \geq 1}$ as defined in (4.6), there exists an absolute constant $C > 0$ and distribution dependent constants $s \in [0, 1)$ and $\sigma_{gap}^2 \geq 0$, such that the following holds: For $n \geq C[\log_2(2n/p_\star) + m \log(n + 1)]/p_\star^2$ and $k \geq 1$,*

we have

$$\mathbb{E}_P [(\boldsymbol{\theta}_n^{(k)} - \boldsymbol{\theta})^2] \leq \frac{\sigma_0^2 - \sigma_{gap}^2}{n} + O\left(\frac{s^k}{n}\right) + \tilde{O}\left(\frac{k^6}{n^{3/2}}\right). \quad (4.10)$$

Furthermore, $\sigma_{gap}^2 > 0$ except for the pathological case of $\mu_X h$ being a constant function.

The quantities σ_{gap}^2 and s are quantified toward the end of this section and are dependent on singular decays of the conditional mean operators for each variable under P . Showing Theorem 4.3.1 boils down to showing that the higher-order term in (4.8) is $O(n^{-1})$ with high probability. Using the expression (4.2) and assuming that $\ell \geq 1$ is odd, we see that

$$[P_n^{(\ell)} - P_n^{(\ell-1)}](\mathcal{C}_\ell \dots \mathcal{C}_k h) = \sum_{\mathbf{x}, \mathbf{z}} \left[\frac{P_X(\mathbf{x})}{P_{n,X}^{(\ell-1)}(\mathbf{x})} - 1 \right] \cdot [\mathcal{C}_\ell \dots \mathcal{C}_k h](\mathbf{x}, \mathbf{z}) P_n^{(\ell-1)}(\mathbf{x}, \mathbf{z}).$$

The first (blue) term in the product quantifies the disagreement between the \mathcal{X} -marginal of $P_n^{(\ell-1)}$ and the true marginal, which can be bounded in terms of $\text{KL}(P_{n,X}^{(0)} \| P_X)$ and is shown to be $O(n^{-1/2})$ with high probability via techniques from information theory. The second term can be unrolled recursively in a similar fashion to (4.8) itself, which will consequently be $\tilde{O}(n^{-1/2})$ as well; this is the most technical part of the analysis (see Appendix B.3.3). We also discuss various extensions such as balancing with misspecified marginals and handling continuous data; see Section 4.6.1 and Section 4.6.2.

Given Theorem 4.3.1, a natural next step is to quantify the gap between σ_0^2 and σ_k^2 , which requires finer-grained properties of \mathcal{C}_X and \mathcal{C}_Z . Notably, we show that as $k \rightarrow \infty$, σ_k^2 approaches the limiting value $\sigma^2 - \sigma_{gap}^2$. Thus, via (4.10), by using $k = o(n^{1/12})$ (excluding logarithmic terms for simplicity), one obtains the asymptotic variance of the solution to (4.7). This contrasts with Albertus and Berthet [2019], in which the dependence of a quantity similar to (4.10) is exponential in k , meaning that $k = o(\log(n))$ is required for convergence under this argument.

From Orthogonal Projections to Variance Reduction We now clarify what is precisely meant by the “spectrum” of the conditional mean operators μ_X and μ_Z . As proven

using a *singular value decomposition* (Proposition B.1.1) in Appendix B.1.1, there exists a basis $\{\alpha_j\}_{j=1}^m$ of $\mathbf{L}^2(P_X)$, a basis $\{\beta_j\}_{j=1}^m$ of $\mathbf{L}^2(P_Z)$, and real values $\{s_j\}_{j=1}^m$, that satisfy

$$\mu_Z \alpha_j = s_j \beta_j \text{ and } \mu_X \beta_j = s_j \alpha_j \text{ for } j \in \{1, \dots, m\}. \quad (4.11)$$

Furthermore, $\alpha_1 = \mathbf{1}_X$ and $\beta_1 = \mathbf{1}_Z$ leading to the equality $\langle f, \alpha_1 \rangle_{\mathbf{L}^2(P_X)} = \mathbb{E}_{P_X}[f(X)]$. Finally, $s_1 = 1$ and s_j is non-negative and non-increasing in j . For a concrete example, consider $m = 2$, in which case P can be written as a matrix in $\mathbb{R}^{2 \times 2}$ and elements of $\mathbf{L}^2(P_X)$ and $\mathbf{L}^2(P_Z)$ are vectors in \mathbb{R}^2 . Then, in the case of uniform marginals, we can verify directly that (4.11) can be satisfied by setting

$$\alpha_1 = \beta_1 = \begin{bmatrix} +1 \\ +1 \end{bmatrix}, \alpha_2 = \beta_2 = \begin{bmatrix} +1 \\ -1 \end{bmatrix}, \text{ and } P = \frac{1}{4} \begin{bmatrix} 1+s & 1-s \\ 1-s & 1+s \end{bmatrix} \quad (4.12)$$

for $s = s_2$ (the second largest singular value). Thus, as $s \rightarrow 1$, the distribution becomes “fully dependent” as Z and X are completely determined by one another. As $s \rightarrow 0$, P approaches the product measure. Geometrically, because $\alpha_1 = \beta_1$, we know that the angle a between the subspaces $\mathbf{L}^2(P_X)$ and $\mathbf{L}^2(P_Z)$ is given by the angle between α_2 and β_2 . By computing their inner product in $\mathbf{L}^2(P)$, we have that $\langle \alpha_2, \beta_2 \rangle_{\mathbf{L}^2(P)} = \langle P, \alpha_2 \beta_2^\top \rangle = s = \cos a$. Thus, $s = 0$ indicates orthogonality of these subspaces, alluding to the independence of X and Z (see the right panel of Figure 4.1).

Returning to $m \geq 2$, we consider the following as a sufficient condition for variance reduction: the operators μ_X and μ_Z have a positive spectral gap, i.e., $s_2 < s_1$. Note that this assumption is satisfied when $P(\mathbf{x}, \mathbf{z}) > 0$ for all $(\mathbf{x}, \mathbf{z}) \in \mathcal{X} \times \mathcal{Z}$ by the Perron–Frobenius Theorem [Horn and Johnson, 2013, Chapter 8]. Using the intuition from Figure 4.1, this rules out pathological cases such as Z being a deterministic function of X . Under the spectral gap condition, the singular values $\{s_j\}_{j=2}^m$ that are strictly less than 1 will determine a geometric rate of decay in variance given in Corollary 4.3.1. The left and right singular functions $\alpha_j : \mathcal{X} \rightarrow \mathbb{R}$ and $\beta_j : \mathcal{Z} \rightarrow \mathbb{R}$ will define a useful coordinate system to represent projections

of h when analyzing $\theta_n^{(k)}$.

Indeed, let $\bar{h} = P(h)$ be the centered test function. Because $\mu_X \bar{h} \in \mathbf{L}^2(P_X)$ and $\mu_Z \bar{h} \in \mathbf{L}^2(P_Z)$, we may decompose this function on the two bases to write

$$\mu_X \bar{h} = \sum_{j=1}^m u_j \alpha_j \quad \text{and} \quad \mu_Z \bar{h} = \sum_{j=1}^m v_j \beta_j. \quad (4.13)$$

Corollary 4.3.1 below relates the (normalized) variance σ_k^2 of the first-order term to the variance of the sample mean $\theta_n^{(0)}$. In fact, it shows that the variance reduction $\sigma_0^2 - \sigma_k^2$ decays geometrically to the quantity

$$\sigma_{\text{gap}}^2 := \sum_{j=2}^m \left[\frac{u_j^2 + v_j^2 - 2s_j u_j v_j}{1 - s_j^2} \right].$$

For simplicity, we only present the result for k even, i.e., σ_{2t}^2 .

Corollary 4.3.1. *The variance reduction achieved by $t + 1$ iterations of the $\mathcal{C}_Z \mathcal{C}_X$ operator can be quantified as*

$$\sigma_0^2 - \sigma_{2(t+1)}^2 = \sigma_{\text{gap}}^2 - \sum_{j=2}^m \frac{s_j^2 (v_j - s_j u_j)^2}{1 - s_j^2} s_j^{4t} = \sum_{j=2}^m \left[u_j^2 + (1 - s_j^{4t+2}) \frac{(v_j - s_j u_j)^2}{1 - s_j^2} \right].$$

Intuitively, the operators \mathcal{C}_X and \mathcal{C}_Z are the main sources of the variance reduction via orthogonality. Since $\alpha_1 = \mathbf{1}_X$, we can see that the reduction will always be strictly positive as long as $\mu_X \bar{h}$ is not a constant function. Finally, using $s := s_2 \geq s_j$ for $j \geq 2$ gives the second term in Theorem 4.3.1.

4.4 Theoretical Limits of Zero-Shot Prediction

We now consider the questions raised in Section 4.2.2. Rather than introducing particular estimators and considering their sample complexity, we focus on deriving a form that relates ZSP to the Bayes optimal predictor on the downstream task, giving an outline to analyze many possible estimation procedures. Before entering into the details, we briefly comment on how our analysis compares to similar theory for few-shot learning (FSL).

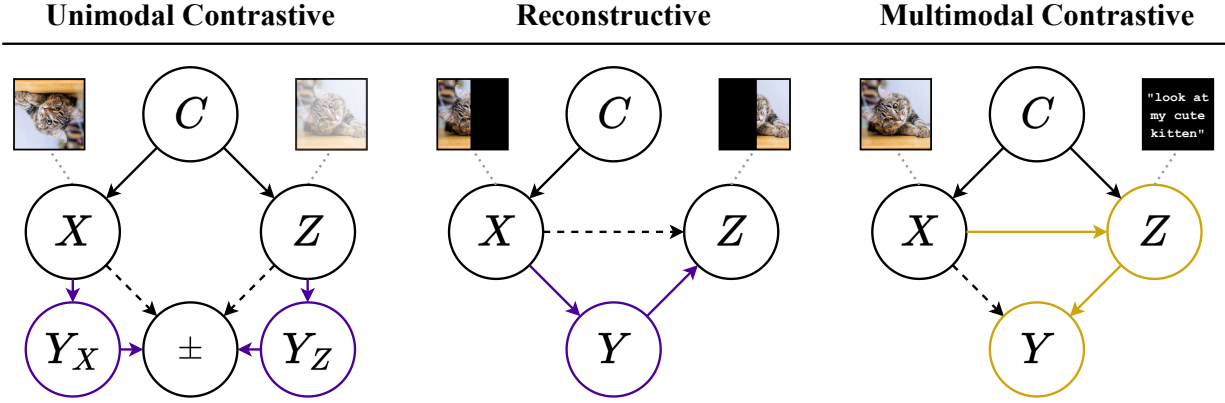


Figure 4.2: **Graphical Models of Self-Supervised Prediction Paths.** Each directed graphical model corresponds to the data types and dependence structures for various SSL pre-training approaches. The variable C represents an unobserved context that determines the observed data-generating distribution. Dotted lines indicate the possibility of presence or absence of the arrow. Methods that are compatible with FSL learn the label Y as a latent variable in the process of solving the pretext task. Methods compatible with ZSP may learn the relationship between X and Z directly, whereas the relationship between Z and Y is estimated via prompting.

The prevailing techniques of recent FSL theory are visualized in Figure 4.2. In unimodal contrastive learning (left), X and Z are augmented/corrupted images, and the pretext task is to identify examples derived from the same (“+”) or different (“−”) underlying image [Chen et al., 2020]. In reconstructive SSL (center), the encoder is pre-trained to predict a hidden portion of the raw/embedded image [Assran et al., 2023]. If the dotted arrows were absent, the only path to solve the pretext task *is* through the label, from which generalization guarantees follow. This motivates another prevalent assumption of exact/approximate conditional independence of X and Z given Y (e.g., as in Lee et al. [2021]). We also avoid this assumption, which is unrealistic in the multimodal context as the dependence between an image and its caption is unlikely to be fully explained by a coarse label such as “cat from CIFAR-10”. In summary, FSL methods learn *through* the label during pre-training. For multimodal contrastive learning (Figure 4.2, right), the ideal dependence structure is fundamentally different.

For illustrative purposes, let us first consider the “compatible” case, that is, when $P_{X,Z} = P_{X,Z}^\mathcal{T}$ and $\rho_{Y,Z} = P_{Y,Z}^\mathcal{T}$, recalling that $P_{X,Y,Z}^\mathcal{T}$ is the evaluation distribution augmented with its latent caption. If X and Y are conditionally independent given Z , then η_ρ (from (4.4)) is in fact identical to η_\star (from (4.3)) because

$$\begin{aligned}\eta_\rho &= \mathbb{E}_{P_{X,Z}} \left[\mathbb{E}_{\rho_{Y,Z}} [f(Y)|Z] | X \right] \\ &= \mathbb{E}_{P_{X,Z}^\mathcal{T}} \left[\mathbb{E}_{P_{Y,Z}^\mathcal{T}} [f(Y)|Z] | X \right] \\ &= \mathbb{E}_{P_{X,Z}^\mathcal{T}} \left[\mathbb{E}_{P_{X,Y,Z}^\mathcal{T}} [f(Y)|Z, X] | X \right] = \eta_\star\end{aligned}$$

by the tower property of conditional expectation under $P_{X,Y,Z}^\mathcal{T}$. Thus, one naturally interprets ZSP as learning *toward* the label using prompting as the essential connection between Z and Y . For this reason, we will quantify the difference between (4.3) and (4.4) based on a conditional dependence measure (see (4.18) below), which formalizes the information-theoretic cost of using natural language as a proxy for image classification. We dub this the *residual dependence* between X and Y that is unexplained by Z . In addition to this, we quantify the incompatibility of the evaluation, pre-training, and prompting distributions, as $P_{X,Z} \neq P_{X,Z}^\mathcal{T}$ and $\rho_{Y,Z} \neq P_{Y,Z}^\mathcal{T}$ in general. We argue that this taxonomy of SSL methods based on their ideal graphical model is useful not only for understanding them conceptually but also for analyzing them mathematically.

We first present the precise description of the identity (4.5) (reproduced below in (4.15)). We then demonstrate the connection between the right-hand side of (4.15) with the empirical ZSP (1.7). Having validated the viewpoint that ZSP is a form of two-stage prediction (i.e., the left-hand side of (4.15)), we use this viewpoint to bound the $\mathbf{L}^2(P_X^\mathcal{T})$ distance between η_ρ and η_\star . Observe the following, which relies on basic measure-theoretic manipulations (see Appendix B.4).

Lemma 4.4.1. *Assume that the pre-training distribution $P_{X,Z}$ satisfies: (i) $P_{X,Z} \ll P_X \otimes P_Z$ with Radon-Nikodym derivative $\mathbf{R} = \left(\frac{dP_{X,Z}}{d(P_X \otimes P_Z)} \right)$, and (ii) that f is integrable under $P_Y^\mathcal{T}$. We*

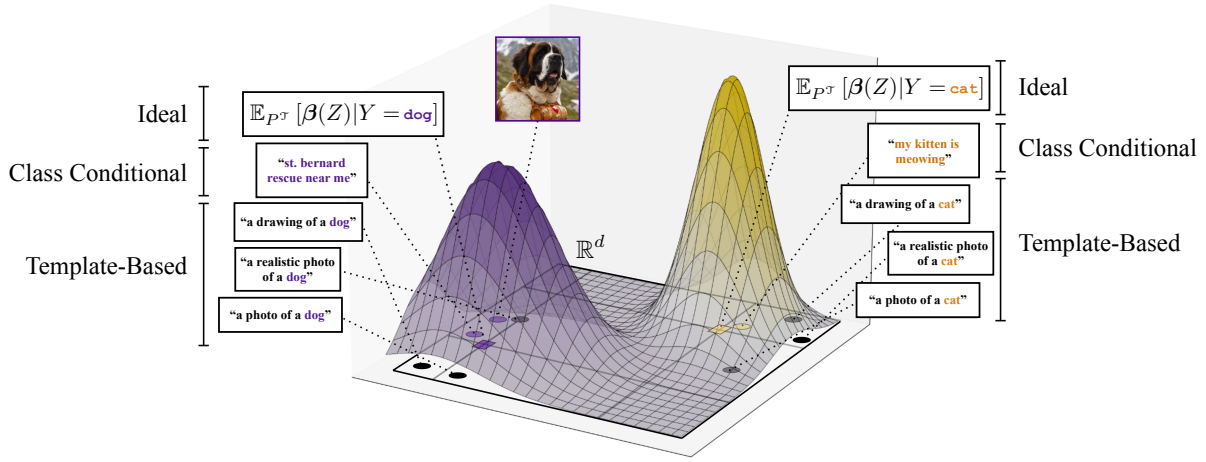


Figure 4.3: **Illustration of Prompting Approaches.** A hypothetical distribution of embeddings $\beta(Z)$ parametrized by two classes (“cat” and “dog”). Three prompting strategies (template-based, class-conditional, and idealized) are shown with example text and resulting embeddings in \mathbb{R}^d . Prompt bias is the distance of the average of the circular points to the square target points.

may then define the function

$$g_\rho(\mathbf{z}) := \mathbb{E}_{\rho_{Y,Z}} [f(Y)|Z](\mathbf{z}). \quad (4.14)$$

Then, P_X -almost surely, we have the equality

$$\eta_\rho(\mathbf{x}) = \mathbb{E}_{\rho_{Y,Z}} [f(Y) \cdot R(\mathbf{x}, Z)] + \underbrace{\int_{\mathcal{Z}} g_\rho(\mathbf{z}) R(\mathbf{x}, \mathbf{z}) (dP_Z(\mathbf{z}) - d\rho_Z(\mathbf{z}))}_{\text{err}(P_Z, \rho_Z)}. \quad (4.15)$$

Proof. By Lemma B.4.1, we already have that for P_X -almost all $\mathbf{x} \in \mathcal{X}$, the identity

$$\begin{aligned} \eta_\rho(\mathbf{x}) &= \mathbb{E}_{P_Z} [g_\rho(Z) R(\mathbf{x}, Z)] \\ &= \mathbb{E}_{\rho_Z} [g_\rho(Z) R(\mathbf{x}, Z)] + \mathbb{E}_{P_Z} [g_\rho(Z) R(\mathbf{x}, Z)] - \mathbb{E}_{\rho_Z} [g_\rho(Z) R(\mathbf{x}, Z)] \\ &= \mathbb{E}_{\rho_Z} [g_\rho(Z) R(\mathbf{x}, Z)] + \int_{\mathcal{Z}} g_\rho(\mathbf{z}) R(\mathbf{x}, \mathbf{z}) (dP_Z(\mathbf{z}) - d\rho_Z(\mathbf{z})). \end{aligned}$$

Now, unpacking the first term on the right-hand side above, we recognize that for fixed

$\mathbf{x} \in \mathcal{X}$, the random variable $R(\mathbf{x}, Z)$ is $\sigma(Z)$ -measurable, so via the properties of conditional expectation [Schilling, 2017, Theorem 27.11 (vii)] in $\mathbf{L}^1(\rho_Z)$, we may write

$$\mathbb{E}_{\rho_Z} [g_\rho(Z)R(\mathbf{x}, Z)] = \mathbb{E}_{\rho_Z} [\mathbb{E}_{\rho_{Y,Z}} [f(Y)|Z] R(\mathbf{x}, Z)] = \mathbb{E}_{\rho_Z} [\mathbb{E}_{\rho_{Y,Z}} [f(Y)R(\mathbf{x}, Z)|Z]] .$$

Using the expression above and the tower property of the conditional expectation, we write

$$\mathbb{E}_{\rho_Z} [g_\rho(Z)R(\mathbf{x}, Z)] = \mathbb{E}_{\rho_Z} [\mathbb{E}_{\rho_{Y,Z}} [f(Y)R(\mathbf{x}, Z)|Z]] = \mathbb{E}_{\rho_{Y,Z}} [f(Y)R(\mathbf{x}, Z)] ,$$

completing the proof. \square

The Radon-Nikodym derivative R appearing in Lemma 4.4.1 is a fundamental quantity known as the *information density* in the information theory literature². It is equal to unity almost surely if and only if X and Z are independent under $P_{X,Z}$. For this reason, its deviation from unity constitutes a dependence measure which we employ later in this section.

Proceeding, we relate the first term of (4.15) to the procedure (1.7). Recall from Section 4.2.2 that the encoders $(\boldsymbol{\alpha}, \boldsymbol{\beta})$ optimize popular loss functions such as CLIP if $\langle \boldsymbol{\alpha}(\mathbf{x}), \boldsymbol{\beta}(\mathbf{z}) \rangle = \log R(\mathbf{x}, \mathbf{z})$, or equivalently, $\exp \langle \boldsymbol{\alpha}(\mathbf{x}), \boldsymbol{\beta}(\mathbf{z}) \rangle = R(\mathbf{x}, \mathbf{z})$. To now make the connection between the right-hand side of (4.15) and (1.7), consider the setting of binary classification ($\mathcal{Y} = \{0, 1\}$) in which ρ_Y places equal weight on each class. Define $\eta_\rho^{(0)}$ and $\eta_\rho^{(1)}$ as instances of (4.4) when setting $f(\mathbf{y}) \equiv f^{(0)}(\mathbf{y}) = \mathbb{1}\{\mathbf{y} = 0\}$ and $f(\mathbf{y}) \equiv f^{(1)}(\mathbf{y}) = \mathbb{1}\{\mathbf{y} = 1\}$, respectively. First, assume that $P_Z = \rho_Z$ (there is no change between the pre-training and prompting distribution, marginally over Z). Then we may then classify \mathbf{x} as the value of \mathbf{y} that maximizes $\eta_\rho^{(\mathbf{y})}(\mathbf{x})$. This classifier is exactly equivalent to

$$\arg \max_{\mathbf{y} \in \{0,1\}} \mathbb{E}_{\rho_{Y,Z}} [\exp \langle \boldsymbol{\alpha}(\mathbf{x}), \boldsymbol{\beta}(Z) \rangle | Y = \mathbf{y}] , \quad (4.16)$$

²This term actually refers to $(\mathbf{x}, \mathbf{z}) \mapsto \log R(\mathbf{x}, \mathbf{z})$, but for simplicity, we use it for R —see Dytso et al. [2023, Eq. (11)].

which can be seen by directly calculating

$$\begin{aligned}
& \mathbb{E}_{\rho_{Y,Z}} [f^{(1)}(Y) \cdot \mathbf{R}(\mathbf{x}, Z)] \\
&= \mathbb{E}_{P_Y^\mathcal{T}} [\mathbb{E}_{\rho_{Y,Z}} [f^{(1)}(Y) \cdot \mathbf{R}(\mathbf{x}, Z) | Y]] \\
&= \frac{1}{2} \mathbb{E}_{\rho_{Y,Z}} [f^{(1)}(Y) \cdot \mathbf{R}(\mathbf{x}, Z) | Y = 1] + \frac{1}{2} \mathbb{E}_{\rho_{Y,Z}} [f^{(1)}(Y) \cdot \mathbf{R}(\mathbf{x}, Z) | Y = 0] \\
&= \frac{1}{2} \mathbb{E}_{\rho_{Y,Z}} [\mathbf{R}(\mathbf{x}, Z) | Y = 1]
\end{aligned}$$

and using $\mathbf{R}(\mathbf{x}, \mathbf{z}) = \exp \langle \boldsymbol{\alpha}(\mathbf{x}), \boldsymbol{\beta}(\mathbf{z}) \rangle$. The expression (4.16) not only resembles (1.7) down to an exponentiation of the inner product but also, by relaxing the assumption that $P_Z = \rho_Z$, we may observe exactly where the effect of distribution shift arises, that is, the introduction of error in the coverage of the prompts.

We now formalize the manner in which we measure the conditional dependence and “incompatibility” of the evaluation distribution $P_{X,Y}^\mathcal{T}$, pre-training distribution $P_{X,Z}$, and prompting distribution $\rho_{Y,Z}$. To do so, we will need to make several mild regularity conditions on $P_{X,Y,Z}$. We use the notion of regular conditional distribution, or r.c.d. (Definition B.4.1), introduced in Appendix B.4.

Assumption 4.4.1. The joint probability $P_{X,Y,Z}^\mathcal{T}$ on $\mathcal{X} \times \mathcal{Y} \times \mathcal{Z}$ satisfies the following constraints.

- **Agrees jointly with the evaluation distribution:** For all measurable sets $A \subseteq \mathcal{X} \times \mathcal{Y}$, we have that $P_{X,Y,Z}^\mathcal{T}(A \times \mathcal{Z}) = P_{X,Y}^\mathcal{T}(A)$ (i.e., $P_{X,Y,Z}^\mathcal{T}$ agrees with the given marginal $P_{X,Y}^\mathcal{T}$).
- **Agrees conditionally with the pre-training distribution:** There exists a measurable set $\mathcal{X}_1 \subseteq \mathcal{X}$ with $P_X^\mathcal{T}(\mathcal{X}_1) = 1$ such that the regular conditional distributions $P_{Z|X=\mathbf{x}}$ and $P_{Z|X=\mathbf{x}}^\mathcal{T}$ on \mathcal{Z} exist. Furthermore, these satisfy $P_{Z|X=\mathbf{x}} = P_{Z|X=\mathbf{x}}^\mathcal{T}$ for all $\mathbf{x} \in \mathcal{X}_1$.
- **Regularity of conditional distributions:** There exists a measurable set $\mathcal{Z}_1 \subseteq \mathcal{Z}$

with $P_Z^\mathcal{T}(\mathcal{Z}_1) = 1$ such that the regular conditional distributions $P_{X,Y|Z=z}^\mathcal{T}$ on $\mathcal{X} \times \mathcal{Y}$ exists for all $\mathbf{z} \in \mathcal{Z}_1$. Furthermore, we have the absolute continuity relation $P_{X,Y|Z=z}^\mathcal{T} \ll P_{X|Z=z}^\mathcal{T} \otimes P_{Y|Z=z}^\mathcal{T}$ with Radon-Nikodym derivative

$$\mathbf{S}_z := \frac{dP_{X,Y|Z=z}^\mathcal{T}}{d(P_{X|Z=z}^\mathcal{T} \otimes P_{Y|Z=z}^\mathcal{T})}, \quad (4.17)$$

that satisfies $\mathbb{E}_{P_{X,Y|Z=z}^\mathcal{T}}[\mathbf{S}_z(X, Y)] < \infty$ for each $\mathbf{z} \in \mathcal{Z}_1$ and $\mathbb{E}_{P_{X,Y,Z}^\mathcal{T}}[\mathbf{S}_Z(X, Y)] < \infty$.

That $P_{X,Y,Z}^\mathcal{T}$ marginalizes to $P_{X,Y}^\mathcal{T}$ is more of an axiomatic property than an assumption, but we phrase it as so to emphasize that $P_{X,Y,Z}^\mathcal{T}$ is meant to describe the evaluation distribution. The assumption that the conditionals $P_{Z|X}$ and $P_{Z|X}^\mathcal{T}$ match almost surely represents the viewpoint that, after fixing an image \mathbf{x} , the latent caption $Z|X = \mathbf{x}$ follows the same relationship to \mathbf{x} as seen during pre-training. Importantly, this does not require or imply that $P_X^\mathcal{T} = P_X$ or that $P_Z^\mathcal{T} = P_Z$. The marginal distribution $P_X^\mathcal{T}$ is supplied entirely by the evaluation distribution $P_{X,Y}^\mathcal{T}$, as for any measurable set $A \subseteq \mathcal{X}$, we have by definition that $P_X^\mathcal{T}(A) = P_{X,Y}^\mathcal{T}(A \times \mathcal{Y})$. On the other hand, the marginal $P_Z^\mathcal{T}$ can be defined using the Markov transition kernel $P_{Z|X=\mathbf{x}}^\mathcal{T}$, in that for any measurable $B \subseteq \mathcal{Z}$, it holds that

$$P_Z^\mathcal{T}(B) := \int_{\mathcal{X}_1} P_{Z|X=\mathbf{x}}^\mathcal{T}(B) dP_X^\mathcal{T}(\mathbf{x}) = \int_{\mathcal{X}_1} P_{Z|X=\mathbf{x}}(B) dP_X(\mathbf{x}).$$

Finally, the absolute continuity condition, i.e., the existence of (4.17), rules out degeneracies such as Y being a deterministic function of X given $Z = \mathbf{z}$ (outside of a set of $P_Z^\mathcal{T}$ -measure zero). This allows us to define the conditional dependence measure

$$I(X; Y|Z = \mathbf{z}) = \mathbb{E}_{P_{X|Z=\mathbf{z}}^\mathcal{T} \otimes P_{Y|Z=\mathbf{z}}^\mathcal{T}}[(1 - \mathbf{S}_z(Y, X))^2], \quad (4.18)$$

which may be known in other contexts as the (square of the) conditional *mean square contingency* [Rényi, 1959, Eq. (13)]. This is the residual dependence measure that appears in the upcoming Theorem 4.4.1.

It is also worth pointing out that the first two conditions Assumption 4.4.1 do not contradict one another. For example, one can consider $P_{X,Y,Z}^\mathcal{T}$ that satisfies the Markov chain

$Y \rightarrow X \rightarrow Z$, where (X, Y) is drawn according to $P_{X,Y}^\mathcal{T}$, and Z and Y are conditionally independent given X . Then, informally, we have that $P_{Z|X,Y}^\mathcal{T} = P_{Z|X}^\mathcal{T} = P_{Z|X}$, so $P_{X,Y,Z}^\mathcal{T}$ is uniquely determined. While this example implies the existence of a valid joint probability measure $P_{X,Y,Z}^\mathcal{T}$, it is also, in a sense, showcasing the “least desirable” distribution for zero-shot prediction, as the dependence between X and Z does not provide any additional information about Y .

Next, we describe how to measure the incompatibility of the evaluation, pre-training, and prompting distributions, which will lead to a *prompt bias* term in the bound. Define the function

$$g_{P_{Y,Z}^\mathcal{T}}(\mathbf{z}) = \mathbb{E}_{P_{Y,Z}^\mathcal{T}}[f(Y)|Z](\mathbf{z}).$$

The prompt bias compares $g_{P_{Y,Z}^\mathcal{T}}$ to g_ρ from (4.14). This reflects the notion that if $P_{X,Y,Z}^\mathcal{T}$ agrees with two of the three fundamental distributions governing the problem, it may not be able to agree with the prompt distribution $\rho_{Y,Z}$ in general. Observe the following.

Theorem 4.4.1. *Assume that f is bounded in absolute value by B_f . Under Assumption 4.4.1, it holds that Then, it holds that*

$$\|\eta_\rho - \eta_\star\|_{\mathbf{L}^2(P_X^\mathcal{T})}^2 \leq \underbrace{2\|g_\rho - g_{P_{Y,Z}^\mathcal{T}}\|_{\mathbf{L}^2(P_Z^\mathcal{T})}^2}_{\text{prompt bias}} + \underbrace{2B_f^2 \mathbb{E}_{P_Z^\mathcal{T}}[I(X; Y|Z)]}_{\text{residual dependence}}. \quad (4.19)$$

Proof. We first establish a useful representation of the conditional mean of $f(Y)$ given $X = \mathbf{x}$, in terms of the (conditional) information density from Lemma B.4.1. Fix $\mathbf{x} \in \mathcal{X}_1$ and $\mathbf{z} \in \mathcal{Z}_1$, the sets on which the regular conditional distributions $P_{Z|X=\mathbf{x}}^\mathcal{T}$ and $P_{X,Y|Z=\mathbf{z}}^\mathcal{T}$ are defined (see Assumption 4.4.1). Because of the existence the Radon-Nikodym derivative S_z from (4.17), we may apply Lemma B.4.1 with $U = Y$, $V = X$, and $h = f$ to write

$$\mathbb{E}_{P_{X,Y|Z=\mathbf{z}}^\mathcal{T}}[f(Y)|X](\mathbf{x}) = \underbrace{\mathbb{E}_{P_{Y|Z=\mathbf{z}}^\mathcal{T}}[f(Y)S_z(Y, \mathbf{x})]}_{=:r(\mathbf{x}, \mathbf{z})} \text{ for all } (\mathbf{x}, \mathbf{z}) \in \mathcal{X}_1 \times \mathcal{Z}_1.$$

We have denoted the right-hand side by the function $r(\mathbf{x}, \mathbf{z})$. Integrate both sides over

$P_{Z|X=\mathbf{x}}^\mathcal{T}$, then use the tower property of conditional expectation to achieve

$$\begin{aligned}\eta_\star(\mathbf{x}) &= \mathbb{E}_{P_{X,Y}^\mathcal{T}} [f(Y)|X](\mathbf{x}) = \int_{\mathcal{Z}} \mathbb{E}_{P_{X,Y|Z=\mathbf{z}}^\mathcal{T}} [f(Y)|X](\mathbf{x}) \, dP_{Z|X=\mathbf{x}}^\mathcal{T}(\mathbf{z}) \\ &= \int_{\mathcal{Z}} r(\mathbf{x}, \mathbf{z}) \, dP_{Z|X=\mathbf{x}}^\mathcal{T}(\mathbf{z}) \\ &= \mathbb{E}_{P_{Z|X=\mathbf{x}}^\mathcal{T}} [r(\mathbf{x}, Z)].\end{aligned}\tag{4.20}$$

Using the identity (4.20) and $P_{Z|X=\mathbf{x}} = P_{Z|X=\mathbf{x}}^\mathcal{T}$ on $\mathbf{x} \in \mathcal{X}_1$ (Assumption 4.4.1), we write

$$\begin{aligned}\eta_\rho(\mathbf{x}) - \eta_\star(\mathbf{x}) &= \mathbb{E}_{P_{Z|X=\mathbf{x}}} [g_\rho(Z)] - \mathbb{E}_{P_{Z|X=\mathbf{x}}^\mathcal{T}} [r(\mathbf{x}, Z)] \\ &= \mathbb{E}_{P_{Z|X=\mathbf{x}}^\mathcal{T}} [g_\rho(Z)] - \mathbb{E}_{P_{Z|X=\mathbf{x}}^\mathcal{T}} [r(\mathbf{x}, Z)] \\ &= \mathbb{E}_{P_{Z|X=\mathbf{x}}^\mathcal{T}} [g_\rho(Z) - g_{P_{Y,Z}^\mathcal{T}}(Z)] + \mathbb{E}_{P_{Z|X=\mathbf{x}}^\mathcal{T}} [g_{P_{Y,Z}^\mathcal{T}}(Z) - r(\mathbf{x}, Z)].\end{aligned}$$

Taking the integral over $P_X^\mathcal{T}$, we have by the decomposition above that

$$\begin{aligned}\|\eta_\rho - \eta_\star\|_{\mathbf{L}^2(P_X^\mathcal{T})}^2 &= \int_{\mathcal{X}} (\eta_\rho(\mathbf{x}) - \eta_\star(\mathbf{x}))^2 \, dP_X^\mathcal{T}(\mathbf{x}) \\ &\leq 2 \int_{\mathcal{X}_1} \left(\mathbb{E}_{P_{Z|X=\mathbf{x}}^\mathcal{T}} [g_\rho(Z) - g_{P_{Y,Z}^\mathcal{T}}(Z)] \right)^2 \, dP_X^\mathcal{T}(\mathbf{x})\end{aligned}\tag{4.21}$$

$$+ 2 \int_{\mathcal{X}_1} \left(\mathbb{E}_{P_{Z|X=\mathbf{x}}^\mathcal{T}} [g_{P_{Y,Z}^\mathcal{T}}(Z) - r(\mathbf{x}, Z)] \right)^2 \, dP_X^\mathcal{T}(\mathbf{x}).\tag{4.22}$$

To handle (4.21), we apply Jensen's inequality for each r.c.d. $P_{Z|X=\mathbf{x}}^\mathcal{T}$ to achieve

$$\begin{aligned}&\int_{\mathcal{X}_1} \left(\mathbb{E}_{P_{Z|X=\mathbf{x}}^\mathcal{T}} [g_\rho(Z) - g_{P_{Y,Z}^\mathcal{T}}(Z)] \right)^2 \, dP_X^\mathcal{T}(\mathbf{x}) \\ &\leq \int_{\mathcal{X}_1} \mathbb{E}_{P_{Z|X=\mathbf{x}}^\mathcal{T}} \left[(g_\rho(Z) - g_{P_{Y,Z}^\mathcal{T}}(Z))^2 \right] \, dP_X^\mathcal{T}(\mathbf{x}) \\ &= \mathbb{E}_{P_Z^\mathcal{T}} \left[(g_\rho(Z) - g_{P_{Y,Z}^\mathcal{T}}(Z))^2 \right] \\ &= \left\| g_\rho - g_{P_{Y,Z}^\mathcal{T}} \right\|_{\mathbf{L}^2(P_Z^\mathcal{T})}^2.\end{aligned}$$

It remains to control (4.22). Applying Jensen's inequality for each r.c.d. $P_{Z|X=\mathbf{x}}^\mathcal{T}$ once again,

we have that

$$\int_{\mathcal{X}_1} \left(\mathbb{E}_{P_{Z|X=\mathbf{x}}^\mathcal{T}} \left[\left(g_{P_{Y,Z}^\mathcal{T}}(Z) - r(\mathbf{x}, Z) \right) \right] \right)^2 dP_X^\mathcal{T}(\mathbf{x}) \quad (4.23)$$

$$\begin{aligned} &\leq \int_{\mathcal{X}_1} \left(\mathbb{E}_{P_{Z|X=\mathbf{x}}^\mathcal{T}} \left[\left(g_{P_{Y,Z}^\mathcal{T}}(Z) - r(\mathbf{x}, Z) \right)^2 \right] \right) dP_X^\mathcal{T}(\mathbf{x}) \\ &= \mathbb{E}_{P_{X,Z}^\mathcal{T}} \left[\left(g_{P_{Y,Z}^\mathcal{T}}(Z) - r(X, Z) \right)^2 \right] \\ &= \int_{\mathcal{Z}_1} \mathbb{E}_{P_{X|Z=\mathbf{z}}^\mathcal{T}} \left[\left(g_{P_{Y,Z}^\mathcal{T}}(\mathbf{z}) - r(X, \mathbf{z}) \right)^2 \right] dP_Z^\mathcal{T}(\mathbf{z}), \end{aligned} \quad (4.24)$$

where the last step follows due to the existence of the r.c.d. $P_{X|Z=\mathbf{z}}^\mathcal{T}$ for $\mathbf{z} \in \mathcal{Z}_1$, as $P_{X|Z=\mathbf{z}}^\mathcal{T}(A) := P_{X,Y|Z=\mathbf{z}}^\mathcal{T}(A \times \mathcal{Y})$ for every measurable $A \subseteq \mathcal{X}$, and the latter exists by assumption. Using the definition of $g_{P_{Y,Z}^\mathcal{T}}$, write

$$g_{P_{Y,Z}^\mathcal{T}}(\mathbf{z}) - r(\mathbf{x}, \mathbf{z}) = \mathbb{E}_{P_{Y|Z=\mathbf{z}}^\mathcal{T}} [f(Y)(1 - \mathbf{S}_\mathbf{z}(Y, \mathbf{x}))].$$

We may substitute this expression into the integrand of (4.24) and apply Jensen's inequality to $P_{Y|Z=\mathbf{z}}^\mathcal{T}$ to achieve

$$\begin{aligned} \mathbb{E}_{P_{X|Z=\mathbf{z}}^\mathcal{T}} \left[\left(g_{P_{Y,Z}^\mathcal{T}}(\mathbf{z}) - r(X, \mathbf{z}) \right)^2 \right] &= \mathbb{E}_{P_{X|Z=\mathbf{z}}^\mathcal{T}} \left[\left(\mathbb{E}_{P_{Y|Z=\mathbf{z}}^\mathcal{T}} [f(Y)(1 - \mathbf{S}_\mathbf{z}(Y, X))] \right)^2 \right] \\ &\leq \mathbb{E}_{P_{X|Z=\mathbf{z}}^\mathcal{T}} \left[\mathbb{E}_{P_{Y|Z=\mathbf{z}}^\mathcal{T}} [(f(Y)(1 - \mathbf{S}_\mathbf{z}(Y, X)))^2] \right] \\ &\leq \|f\|_\infty^2 \mathbb{E}_{P_{X|Z=\mathbf{z}}^\mathcal{T}} \left[\mathbb{E}_{P_{Y|Z=\mathbf{z}}^\mathcal{T}} [(1 - \mathbf{S}_\mathbf{z}(Y, X))^2] \right] \\ &= \|f\|_\infty^2 \mathbb{E}_{P_{X|Z=\mathbf{z}}^\mathcal{T} \otimes P_{Y|Z=\mathbf{z}}^\mathcal{T}} [(1 - \mathbf{S}_\mathbf{z}(Y, X))^2], \end{aligned}$$

where the final step follows by applying Fubini's theorem [Schilling, 2017, Corollary 14.9] to the product measure $P_{X|Z=\mathbf{z}}^\mathcal{T} \otimes P_{Y|Z=\mathbf{z}}^\mathcal{T}$ for fixed $\mathbf{z} \in \mathcal{Z}_1$. By the definition (4.18), it holds that

$$\mathbb{E}_{P_{X|Z=\mathbf{z}}^\mathcal{T} \otimes P_{Y|Z=\mathbf{z}}^\mathcal{T}} [(1 - \mathbf{S}_\mathbf{z}(Y, X))^2] = I(X; Y|Z = \mathbf{z}). \quad (4.25)$$

After confirming that (4.25) is $P_Z^\mathcal{T}$ -integrable, substituting this expression back into (4.24) achieves the desired result. Expand the quadratic term and apply the Radon-Nikodym

theorem [Schilling, 2017, Theorem 20.1] to achieve

$$\begin{aligned}
I(X; Y|Z = \mathbf{z}) &= 1 - 2\mathbb{E}_{P_{X|Z=\mathbf{z}}^\mathcal{T} \otimes P_{Y|Z=\mathbf{z}}^\mathcal{T}} [\mathbf{S}_\mathbf{z}(Y, X)] + \mathbb{E}_{P_{X|Z=\mathbf{z}}^\mathcal{T} \otimes P_{Y|Z=\mathbf{z}}^\mathcal{T}} [\mathbf{S}_\mathbf{z}^2(Y, X)] \\
&= 1 - 2\mathbb{E}_{P_{X,Y|Z=\mathbf{z}}^\mathcal{T}} [1] + \mathbb{E}_{P_{X,Y|Z=\mathbf{z}}^\mathcal{T}} [\mathbf{S}_\mathbf{z}(Y, X)] \\
&= \mathbb{E}_{P_{X,Y|Z=\mathbf{z}}^\mathcal{T}} [\mathbf{S}_\mathbf{z}(Y, X)] - 1.
\end{aligned}$$

Thus, by integrating against $P_Z^\mathcal{T}$, we see that

$$\mathbb{E}_{P_Z^\mathcal{T}} [I(X; Y|Z)] = \mathbb{E}_{P_{X,Y,Z}^\mathcal{T}} [\mathbf{S}_Z(Y, X)] - 1,$$

where the expectation term is finite by Assumption 4.4.1. This completes the proof. \square

To give context to Theorem 4.4.1, conditional independence relations have previously been used to describe the performance of multimodal contrastive SSL for FSL. We are particularly inspired by the *multi-view redundancy* theory of Tosh et al. [2021], which states informally that the population FSL predictor can approach the performance of the idealized direct predictor that is given *both* X and Z at test time, if $X \perp\!\!\!\perp Y|Z$ and $Z \perp\!\!\!\perp Y|X$ approximately hold. However, the theory of graphical models [Lauritzen, 1996, Proposition 3.1] asserts that both conditional independence relations hold only if $(X, Z) \perp\!\!\!\perp Y$, that is, neither view has information predictive of the label. This can be seen intuitively from Figure 4.2 by breaking the arrows $X \rightarrow Y$ and $Z \rightarrow Y$. Notice that we compare only to the Bayes optimal predictor (4.3) given X , so that we need only that $X \perp\!\!\!\perp Y|Z$ (i.e., X depends on Y through Z) to close the gap.

We now present several numerical illustrations to validate the hypotheses that inspired the results in both Section 4.3 and Section 4.4. Namely, we wish to see the effects of 1) balancing procedures during pre-training and 2) inference-time prompting on downstream model performance. This includes an example of *unbiased prompting*, or a choice of $\rho_{Y,Z}$ that renders the prompt bias term of (4.19) as zero.

4.5 Experiments

In our experiments, we illustrate the two subjects of this chapter in practice, that is, how data balancing and various prompt strategies (i.e., estimates of $P_{Y,Z}^{\mathcal{T}}$) affect the empirical performance of CLIP models. We focus on image classification tasks as a canonical testbed for ZSP. See Appendix B.5 for experiment details. Code to reproduce the data and experiments can be found at <https://github.com/ronakdm/balancing> and <https://github.com/ronakdm/zeroshot>.

Models, Datasets, and Evaluation. Throughout, we consider training variants of the CLIP model, which require a dataset of image-caption pairs. In the experiments that simulate pre-training, we train models on top of embedded representations of images and text. The pre-training set chosen is the ImageNet-Captions dataset [Fang et al., 2023], which pairs images from ImageNet [Deng et al., 2009] that were taken from Flickr with their original captions. For prompting-based experiments, we use three publicly available CLIP models from the OpenCLIP repository [Ilharco et al., 2022]: ResNet50 pre-trained on YFCC15M [Thomee et al., 2016], NLLB-CLIP pre-trained on a subset of LAION COCO [Visheratin, 2023], and ViT-B/32 pre-trained on the DataComp medium pool [Gadre et al., 2023].

We evaluate models based on zero-shot classification performance via top- k accuracy, in which a test example is considered to be classified correctly if the true class is contained within the elements of \mathcal{T} with the k largest scores as computed by (1.7). Our evaluation datasets include five standard benchmarks: the Describable Textures Dataset or DTD [Cimpoi et al., 2014], Flowers 102 [Nilsback and Zisserman, 2008], FGVC Aircraft [Maji et al., 2013], SUN397 [Xiao et al., 2010], and ImageNet-1k [Deng et al., 2009]. For some prompting experiments, we also make use of the ImageNet-Captions dataset as a way to estimate prompt embeddings. Evaluation is done using tools from the [CLIP Benchmark repository](#). In Figure 4.5 and Figure 4.6, “templates” refers to using all [default community-curated prompts](#) available in CLIP Benchmark. See Appendix B.5 for specific model tags and full experimental details.

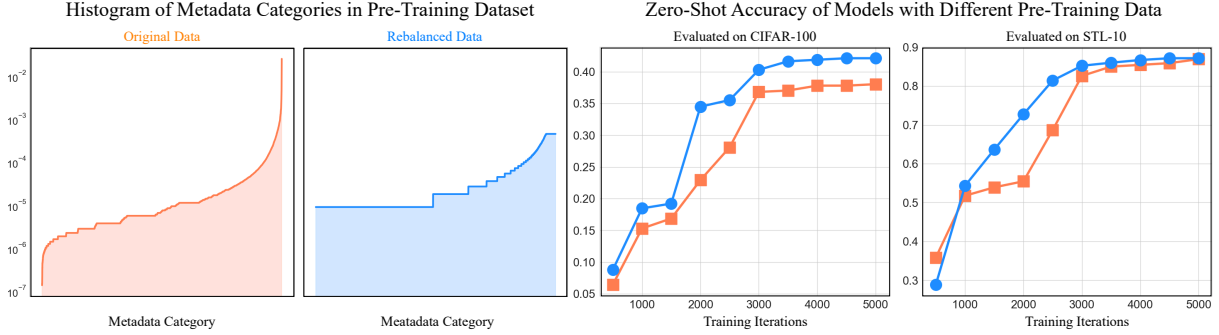


Figure 4.4: **Balancing as Data Curation.** Depiction of balancing and data curation on ImageNet-Captions dataset, in which \mathcal{X} represents image-caption pairs and \mathcal{Y} represents keywords. **Left:** Observed marginal $P_{n,Z}$ (orange) and P_Z (blue), which are sorted by order of increasing probability. **Right:** Zero-shot evaluation of an embedding model trained using the standard CLIP loss on the original versus the balanced training set.

4.5.1 Balanced Pre-Training Effects

We perform a data curation experiment exploring the use of balancing to adjust the entire pre-training set, in the spirit of Xu et al. [2024]. The target marginal P_Z is selected by choosing a threshold for which frequent keywords have their probability mass truncated, and the probability measure is normalized to sum to one. In Figure 4.4, we show the observed marginal $P_{n,Z}$ and the target marginal P_Z sorted in increasing order (left). The original marginal on \mathcal{Y} has approximately 5 orders of magnitude of difference between the most and least probable keyword. After balancing, the target marginal has less than 2 orders of difference between the most frequent and least frequent keywords.. To see how this affects downstream performance, we plot the zero-shot classification accuracy over training iterations in Figure 4.4 (right) when using the original dataset (orange) and using the metadata-balanced dataset (blue). We observe moderate improvement, especially in the small batch regime ($m = 512$) when curating the dataset.

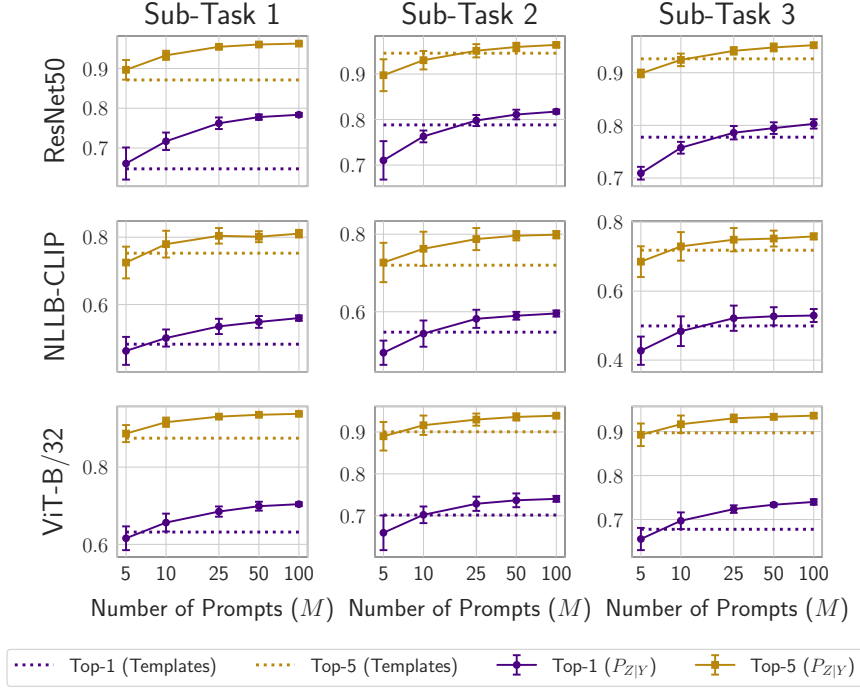


Figure 4.5: **Results: Unbiased Prompting.** Pre-trained models are varied along the rows, and sub-tasks (subsets of 50 ImageNet-1k classes) are varied along columns. In all plots, the x -axis denotes the number of prompts sampled for each class embedding, and the y -axis denotes top- k zero-shot classification accuracy. Error bars indicate standard deviations across 10 seeds for prompt sampling. In this setting $P_{Y,Z} = P_{Y,Z}^\mathcal{T}$.

4.5.2 Unbiased Prompting with Observations from $P_{Z,Y}^\mathcal{T}$

This is an illustrative experiment in which it is possible to use the unbiased prompting mechanism, that is, directly drawing samples from the distribution $P_{Y,Z}^\mathcal{T}$ to approximate the conditional mean $\mathbb{E}_{P_{Y,Z}^\mathcal{T}}[\beta(Z)|Y = \mathbf{y}]$ (see Figure 4.3). Indeed, because the ImageNet-Captions dataset includes ImageNet images, their Flickr captions, and their label (i.e., the joint observation of (X, Y, Z)), we may directly observe the scaling of the variance with respect to M without prompt bias. We design three in-distribution sub-tasks by randomly selecting collections of 50 classes ($\mathcal{Y}_1, \mathcal{Y}_2, \mathcal{Y}_3$) from each of 998 classes, reserving held-out *prompting examples* $(Z_1, Y_1), \dots, (Z_{15,000}, Y_{15,000})$, 100 for each of 150 classes. Then, for task i , using M examples $j_1(\mathbf{y}), \dots, j_M(\mathbf{y})$ selected randomly without replacement for $\mathbf{y} \in$

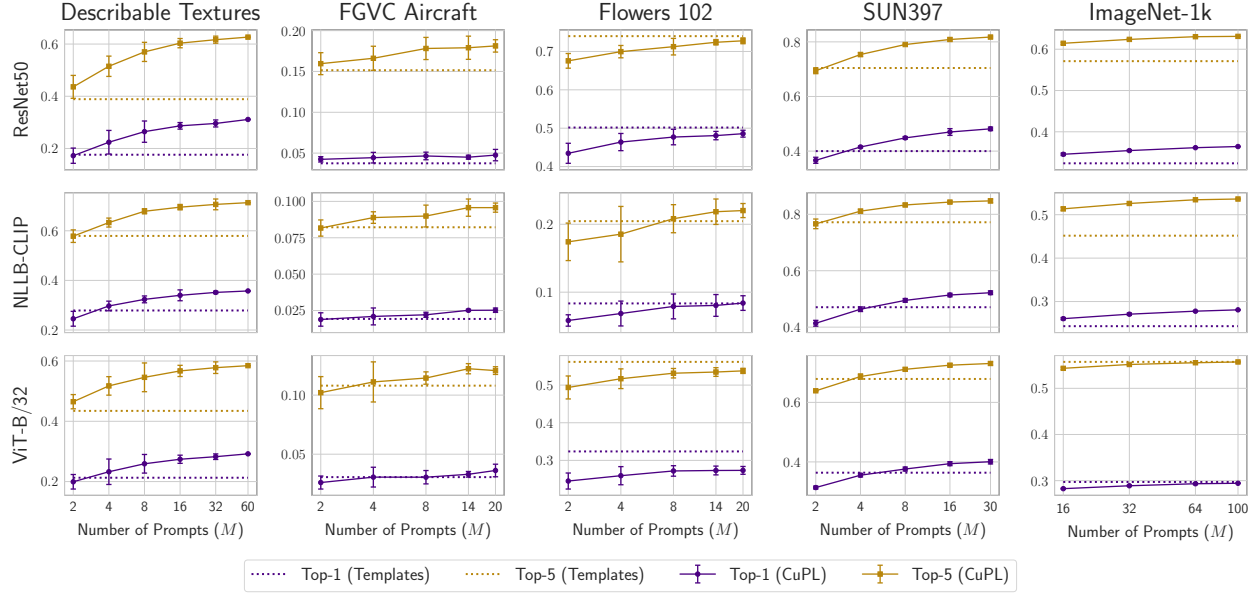


Figure 4.6: **Results: Class-Conditional Prompting.** Pre-trained models are varied along the rows, and evaluation datasets are varied along columns. In all plots, the x -axis denotes the number of prompts sampled for each class embedding, and the y -axis denotes top- k zero-shot classification accuracy. Error bars indicate standard deviations across 10 seeds for prompt sampling.

\mathcal{Y}_i , we use the vector $\frac{1}{M} \sum_{m=1}^M \beta(Z_{j_m(y)})$ as the class embedding (projected to unit norm). Using an evaluation set of approximately 25,000 examples from each sub-task, we compute the classification accuracy of this approach and plot the results in Figure 4.5. Observe that the threshold at which unbiased prompting outperforms the 18 default templates is approximately $M = 10$ across tasks. However, the performance of the theoretically unbiased approach only saturates at $M = 100$ and can have enormous benefits (almost 15% absolute increase in top-1 accuracy for the ResNet50 on Sub-Task 1) in performance. Thus, for models that have not yet been saturated from pre-training, prompting can close surprisingly wide gaps in zero-shot classification accuracy.

4.5.3 Class-Conditional Prompting with Language Models

From both theoretical and empirical investigations, two facts are clear: class-agnostic templates are outperformed by captions that are attuned to the class being evaluated, and a potentially large number of prompts (50-100) is needed to saturate performance, as opposed to the 5-15 defaults used in OpenCLIP. Framed as so, this is a superb use case for modern large language models (LLMs), as they can generate class-aware captions with high throughput. As mentioned in Section 4.1, we investigate CuPL as a means to implement class-conditional prompting with LLMs. Our experimental setup and scientific goals differ from those used in Pratt et al. [2023]: 1) we use lightweight encoders that have not saturated their performance during pre-training, as opposed to their use of the large-scale ViT-L/14 architecture, 2) we quantify the variability of classification accuracy with respect to prompting by generating up to fifty times as many prompts per experiment, and 3) we employ LLaMA 3, which is free and accessible to other, as opposed to GPT-3 [Brown et al., 2020]. The results are shown in Figure 4.6, where we order the datasets in increasing number of classes per task: 47, 100, 102, 397, and 998. Similar phenomena as in the unbiased case are observed, although the approximate saturation threshold varies per dataset from 20 for Flowers 102 and FGVC Aircraft up to 60 for DTD. Note that the choice of defaults heavily influences the baseline performance. Shockingly, the Flowers 102 dataset uses a single default: “a photo of a __, a type of flower”, and is often able to outperform the class-conditional LLM approach on average. On the other hand, the DTD templates of the form “a photo of a __ {texture, pattern, thing, object}” are dramatically outperformed by our LLM-generated captions such as “a gauzy surface is characterized by a thin, translucent, and wispy appearance that is soft and delicate in texture”, with a nearly 20% increase in top-5 accuracy on the ResNet50 and ViT-B/32 architectures. Understanding, both theoretically and experimentally, the properties of the distribution $P_{Z|Y=\mathbf{y}}$ that explain such performance differences between hand-designed and model-generated captions is still a ripe and exciting topic for further research.

4.6 Possible Extensions

4.6.1 Misspecified Marginal Distributions

Throughout the chapter, we have assumed that the marginals (P_X, P_Z) supplied by the user are accurate. A natural question to consider is how the analysis and final result in terms of mean squared error change when exposed to some degree of marginal error. Indeed, we consider the balancing method given marginals (\hat{P}_X, \hat{P}_Z) which satisfy the following structure.

Assumption 4.6.1. There exist fixed probability mass functions \hat{P}_X and \hat{P}_Z for some $\varepsilon \in [0, 1)$,

$$\hat{P}_{X,\varepsilon} = (1 - \varepsilon)P_X + \varepsilon\hat{P}_X \text{ and } \hat{P}_{Z,\varepsilon} = (1 - \varepsilon)P_Z + \varepsilon\hat{P}_Z.$$

We also have the existence of the positive quantity

$$\hat{p}_\star := \min\{\min_{\mathbf{x}} \hat{P}_X(\mathbf{x}), \min_{\mathbf{z}} \hat{P}_Z(\mathbf{z})\} > 0.$$

Given the existence of $\hat{p}_\star > 0$, we may also define

$$\hat{p}_{\star,\varepsilon} = \min\{\min_{\mathbf{x}} \hat{P}_{X,\varepsilon}(\mathbf{x}), \min_{\mathbf{z}} \hat{P}_{Z,\varepsilon}(\mathbf{z})\} \geq \varepsilon\hat{p}_\star + (1 - \varepsilon)p_\star > 0.$$

To be precise, the iterations of balancing are initialized at $\hat{P}_n^{(0)} = P_n$ and

$$\hat{P}_n^{(k)}(\mathbf{x}, \mathbf{z}) := \begin{cases} \frac{\hat{P}_{X,\varepsilon}(\mathbf{x})}{\hat{P}_{n,X}^{(k-1)}(\mathbf{x})} \cdot \hat{P}_n^{(k-1)}(\mathbf{x}, \mathbf{z}) & k \text{ odd} \\ \frac{\hat{P}_{Z,\varepsilon}(\mathbf{z})}{\hat{P}_{n,Z}^{(k-1)}(\mathbf{z})} \cdot \hat{P}_n^{(k-1)}(\mathbf{x}, \mathbf{z}) & k \text{ even} \end{cases}. \quad (4.26)$$

While details are deferred to Appendix B.3.5, we give an overview of the proof and state the final bound in this section. Broadly, the proof will proceed by deriving analogous results to those in Appendix B.3.1, Appendix B.3.3, and Appendix B.3.4. As in Section 4.3, the backbone of the argument will be a recursive formula to relate the estimation error across

iteration counter $k \geq 1$. When doing so similarly to (4.8), we achieve the equality

$$[\hat{P}_n^{(k)} - P](h) = [\hat{P}_n^{(k-1)} - P](\mathcal{C}_k h) + \hat{V}_n^{(k-1)}(\mathcal{C}_k h) + \begin{cases} [\hat{P}_{X,\varepsilon} - P_X](\mu_X h) & \text{if } k \text{ odd} \\ [\hat{P}_{Z,\varepsilon} - P_Z](\mu_Z h) & \text{if } k \text{ even} \end{cases}. \quad (4.27)$$

as shown in (B.48). Because we must bound the error terms containing $(\hat{P}_{X,\varepsilon} - P_X)$ and $(\hat{P}_{Z,\varepsilon} - P_Z)$, unlike (4.8), our recursion will be stated in the form of an inequality. We measure the deviation of the marginals (\hat{P}_X, \hat{P}_Z) from the true (P_X, P_Z) using the constant

$$c^2 = \max \left\{ \chi^2(\hat{P}_X \| P_X), \chi^2(\hat{P}_Z \| P_Z) \right\},$$

and show in Proposition B.3.4 that

$$[\hat{P}_n^{(k)} - P](h) \leq \left| [\hat{P}_n^{(k-1)} - P](\mathcal{C}_k h) \right| + \left| \hat{V}_n^{(k-1)}(\mathcal{C}_k h) \right| + c \|h\|_{\mathbf{L}^2(P)} \sqrt{\varepsilon},$$

where, as one might expect, we defined

$$\hat{V}_n^{(k-1)}(h) := \begin{cases} \sum_{\mathbf{x}, y} \left(\frac{\hat{P}_{X,\varepsilon}}{\hat{P}_{n,X}^{(k-1)}}(\mathbf{x}) - 1 \right) h(\mathbf{x}, \mathbf{z}) \hat{P}_n^{(k-1)}(\mathbf{x}, \mathbf{z}) & k \text{ odd} \\ \sum_{\mathbf{x}, y} \left(\frac{\hat{P}_{Z,\varepsilon}}{\hat{P}_{n,Z}^{(k-1)}}(\mathbf{z}) - 1 \right) h(\mathbf{x}, \mathbf{z}) \hat{P}_n^{(k-1)}(\mathbf{x}, \mathbf{z}) & k \text{ even} \end{cases} \quad (4.28)$$

The inequality (4.28) is a form of the centered recursion formula in Proposition B.3.1, and we can perform the unrolling

$$[\hat{P}_n^{(k)} - P](h) \leq \underbrace{[P_n^{(0)} - P](h)(\mathcal{C}_1 \dots \mathcal{C}_k h)}_{\text{first-order term}} + \underbrace{\sum_{\ell=1}^k \left| \hat{V}_n^{(\ell-1)}(\mathcal{C}_\ell \dots \mathcal{C}_k h) \right|}_{\text{higher-order term}} + \underbrace{O(k\sqrt{\varepsilon})}_{\text{misspecification}}$$

and establish an analogous uncentered recursion formula to upper bound the higher-order terms. The remainder of the argument is purely algebraic; the calculations are contained in Appendix B.3.5, and are written in a way to reuse parts of the analysis from throughout Appendix B.3. For convenience, we state the mean squared error bound in terms of its dependence on $(\varepsilon, n, k, \hat{p}_{\star, \varepsilon})$ and remove absolute constants and other problem parameters.

Theorem 4.6.1. *Let Assumption 4.6.1 be true with error $\varepsilon \in [0, 1)$. For a sequence of rebalanced distributions $(\hat{P}_n^{(k)})_{k \geq 1}$, there exists an absolute constant $C > 0$ such that when $n \geq C[\log_2(2n/\hat{p}_{\star, \varepsilon}) + m \log(n+1)]/\min\{p_{\star}, \hat{p}_{\star, \varepsilon}\}^2$, we have that*

$$\begin{aligned} \mathbb{E}_P \left[\left(\hat{P}_n^{(k)}(h) - P(h) \right)^2 \mathbb{1}_S \right] + \mathbb{E}_P \left[(P_n(h) - P(h))^2 \mathbb{1}_{S^c} \right] &\leq \frac{\sigma_k^2}{n} + \tilde{O} \left(\frac{k^6}{n^{3/2}} \right) \\ &+ \tilde{O} \left(\frac{k^4}{\hat{p}_{\star, \varepsilon}^2} \left(\sqrt{\frac{1}{n} \log \frac{1}{1-\varepsilon}} + \log \frac{1}{1-\varepsilon} \right) \left[\frac{k^2}{\hat{p}_{\star, \varepsilon}^2} \left(\sqrt{\frac{1}{n} \log \frac{1}{1-\varepsilon}} + \log \frac{1}{1-\varepsilon} + \frac{1}{n} \right) + \frac{1}{\sqrt{n}} \right] \right) \\ &+ \tilde{O} \left(k^2 \left[\sqrt{\varepsilon} \left(\frac{\hat{p}_{\star, \varepsilon}^4}{n^4} + \frac{1}{\sqrt{n}} + \frac{\hat{p}_{\star, \varepsilon}^2 k}{n^4} \left(n + \frac{k^2}{\hat{p}_{\star, \varepsilon}^2} \right) + \frac{k^2}{\hat{p}_{\star, \varepsilon}^2} \left[\frac{1}{n} + \sqrt{\frac{1}{n} \log \frac{1}{1-\varepsilon}} + \log \frac{1}{1-\varepsilon} \right] \right) + \varepsilon \right] \right). \end{aligned}$$

Notably, we can no longer take $k \rightarrow \infty$ as $n \rightarrow \infty$ in the expression above unless we also have that $\varepsilon \rightarrow 0$ at the appropriate rate.

4.6.2 Alternative Marginal Rebalancing Approches

In Section 4.3, we considered an estimator which incorporated the marginals by approximating the projection (4.7), taking for granted that it is a reasonable way to use this side information. However, when dealing with continuous real-valued data, there is another rather simple method by using the generalized inverse of the cumulative distribution functions (CDFs) of P_X and P_Z . In the estimator discussed in Bickel et al. [1991, Section 4], we have that $\mathcal{X} = [0, 1]$, $\mathcal{Z} = [0, 1]$, and the respective target marginals are uniform. Thus, the joint distribution P is a copula and the goal is to estimate a particular linear functional of P : the probability of the event $\{X \geq s \cap Z \geq t\}$ for $(s, t) \in [0, 1] \times [0, 1]$.³ By using the inverse CDF trick described below, one can simply transform the two modalities individually to (approximately) fit a particular marginal distribution in each variable. It is shown in the same work that this approach is in fact *inefficient* in the nonparametric model with (P_X, P_Z) known, while the raking/iterative proportional fitting method used in this chapter is efficient.

³We use the survival events $\{X \geq s\}$ and $\{Z \geq t\}$ instead of $\{X \leq s\}$ and $\{Z \leq t\}$ because calculations with inverse CDFs will be simpler.

One notable difference between this work and [Bickel et al. \[1991\]](#) is that we consider discrete data directly, whereas they consider continuous data that is discretized via partitions. Letting \mathcal{A}_n be a partition of measurable subsets of \mathcal{X} and \mathcal{B}_n defined similarly for \mathcal{Z} , the condition [Bickel et al. \[1991, F1\]](#) requires that

$$P_X(A) \geq \frac{\lambda_n}{\sqrt{n}} \quad \forall A \in \mathcal{A}_n \quad \text{and} \quad P_Z(B) \geq \frac{\lambda_n}{\sqrt{n}} \quad \forall B \in \mathcal{B}_n,$$

where the sequence $(\lambda_n)_{n \geq 1}$ satisfies $\lambda_n^2 / \log(n) \rightarrow \infty$ and $\lambda_n / \sqrt{n} \rightarrow 0$ as $n \rightarrow \infty$. Moreover, they directly analyze the solution (4.7), so that the efficient influence function and asymptotic variance are expressed in terms of projections. As mentioned in Section 4.3, a key outcome of our analysis is a formula for the asymptotic variance reduction by constructing a sequence of probability measures and Markov operators associated to each iteration of (4.2), where every probability measure of the sequence can be expressed in closed form. Thus, the asymptotic variance can be stated in terms of the spectra of these operators.

Despite the inefficiency of the copula-based estimator for continuous data, it is conceptually interesting to construct such an estimator for discrete data, providing another possibly efficient estimator. Conceptually, the raking approach can be described as one in which the empirical measure P_n is computed from data, and then transformed using the auxiliary marginal information (P_X, P_Z) . We may also consider a reversed approach, that is, to transform the data $(X_1, Z_1), \dots, (X_n, Z_n)$ and then compute the empirical measure afterward, in a way that ensures that the resulting measure adheres to the marginal constraints. To be precise, “adhering” to the marginal constraints will be interpreted differently for continuous and discrete data, as explained in the rest of this section. We now describe a generalization of the copula-based estimator that preserves the essential aspects but can be used in our setting.

Generalizing the Copula Estimator Having observed univariate data $(X_1, Z_1), \dots, (X_n, Z_n)$, let $F_{n,X}$ denote the empirical CDF of $\{X_i\}_{i=1}^n$, $F_{n,Z}$ of $\{Z_i\}_{i=1}^n$, and let $P_n = \frac{1}{n} \sum_{i=1}^n \delta_{(X_i, Z_i)}$ denote the empirical measure. The target marginals can be represented by their respective

CDFs F_X and F_Z . Finally, for any CDF F on \mathbb{R} , define the quantile function $F^{-1}(s) := \inf \{\mathbf{x} \in \mathbb{R} : F(\mathbf{x}) \geq s\}$. The estimator in this case is

$$\begin{aligned} P_n(X \geq F_{n,X}^{-1}(s), Z \geq F_{n,Z}^{-1}(t)) &= \frac{1}{n} \sum_{i=1}^n \mathbb{1} \{X_i \geq F_{n,X}^{-1}(s), Z_i \geq F_{n,Z}^{-1}(t)\} \\ &\stackrel{(1)}{=} \frac{1}{n} \sum_{i=1}^n \mathbb{1} \{F_{n,X}(X_i) \geq s, F_{n,Z}(Z_i) \geq t\} \\ &\stackrel{(2)}{=} \frac{1}{n} \sum_{i=1}^n \mathbb{1} \{F_X^{-1}(F_{n,X}(X_i)) \geq s, F_Z^{-1}(F_{n,Z}(Z_i)) \geq t\}, \end{aligned}$$

where (1) follows by [Bobkov and Ledoux \[2019, Lemma A.3\]](#) and (2) follows because the inverse CDF of the uniform distribution $[0, 1]$ is equal to the identity on $[0, 1]$. Notably, we applied the maps $T_X = F_X^{-1} \circ F_{n,X}$ and $T_Z = F_Z^{-1} \circ F_{n,Z}$ to each of our data sources, and then estimated the empirical measure. As a result, we have that

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n \mathbb{1} \{T_X(X_i) \geq s\} &= 1 - F_X(s) \text{ for } s \in \left\{ \frac{0}{n}, \frac{1}{n}, \dots, \frac{n}{n} \right\} \\ \frac{1}{n} \sum_{i=1}^n \mathbb{1} \{T_Z(Z_i) \geq t\} &= 1 - F_Z(t) \text{ for } t \in \left\{ \frac{0}{n}, \frac{1}{n}, \dots, \frac{n}{n} \right\}. \end{aligned}$$

The restriction that s and t are selected at increments of $1/n$ is the reason why the transformed data only partially adheres to the marginal constraints. Indeed, it is unreasonable to require an empirical CDF to agree with the CDF of a continuous random variable on its entire domain. To complete the story, we notice that T_X and T_Z are the optimal transportation maps from $F_{n,X}$ to F_X and $F_{n,Z}$ to F_Z , respectively, with respect to the squared distance on \mathbb{R} (see [\[Bobkov and Ledoux, 2019, Theorems 2.10 and 2.11\]](#)). This key observation will motivate the discrete version of this estimator.

The Transport Plan Estimator We return to the setting in which the sample spaces $\mathcal{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_m\}$ and $\mathcal{Z} = \{\mathbf{z}_1, \dots, \mathbf{z}_m\}$ are finite. We adopt in this discussion the following assumption. In a full treatment, we argue that this assumption holds with high probability.

Assumption 4.6.2. The empirical marginals $(P_{n,X}, P_{n,Z})$ and the target marginals (P_X, P_Z) are positive on their domains.

Following suit from the previous section, we are interested in defining transformations $T_X : \mathcal{X} \mapsto \mathcal{X}$ and $T_Z : \mathcal{Z} \mapsto \mathcal{Z}$ such that

$$\frac{1}{n} \sum_{i=1}^n \mathbb{1} \{T_X(X_i) = \mathbf{x}\} = P_X(\mathbf{x}) \text{ and } \frac{1}{n} \sum_{i=1}^n \mathbb{1} \{T_Z(Z_i) = \mathbf{z}\} = P_Z(\mathbf{z}), \quad (4.29)$$

without altering the data too much. We first specify a ground cost by endowing \mathcal{X} and \mathcal{Z} with the discrete metrics $d_X(\mathbf{z}, \mathbf{x}') = \mathbb{1} \{\mathbf{x} \neq \mathbf{x}'\}$ and $d_Z(\mathbf{z}, \mathbf{z}') = \mathbb{1} \{\mathbf{z} \neq \mathbf{z}'\}$. This is a natural choice of metric, as we have not made any additional assumptions (such as ordinality) on \mathcal{X} and \mathcal{Z} ; they are purely categorical (although other domain-specific distances may be designed). Next, searching among the maps that satisfy (4.29) is equivalent to solving the Monge assignment problem, which may not be feasible. Thus, we consider the Kantorovich relaxation [Peyré and Cuturi, 2019, Section 2.3] and replace T_X and T_Z with probability measures $\pi_n^* : \mathcal{X} \times \mathcal{X} \rightarrow [0, 1]$ and $\gamma_n^* : \mathcal{Z} \times \mathcal{Z} \rightarrow [0, 1]$ that solve the respective optimal transportation problems

$$\pi_n^* = \arg \min_{\pi \in \Pi(P_{n,X}, P_X)} \sum_{i \neq j} \pi(\mathbf{x}_i, \mathbf{x}_j) \text{ and } \gamma_n^* = \arg \min_{\gamma \in \Pi(P_{n,Z}, P_Z)} \sum_{i \neq j} \gamma(\mathbf{z}_i, \mathbf{z}_j), \quad (4.30)$$

which are always feasible and for which minimizers are guaranteed to exist. Here, Π denotes the set of couplings with the given marginals. For any solution pair (π_n^*, γ_n^*) (called *transport plans*), we will design the estimator P_n^* of P in a way that is semantically similar to the transportation map interpretation of the previous section. Intuitively, we will define the estimator to be the joint distribution of random variables (U, V) on $\mathcal{X} \times \mathcal{Z}$, which are components of a quadruple (X, Z, U, V) satisfying the following properties:

- $(X, Z) \sim P_n$.
- $(X, U) \sim \pi_n^*$ and $(Z, V) \sim \gamma_n^*$.
- U and (Z, V) are conditionally independent given X .

- V and (X, U) are conditionally independent given Z .

These properties are enough to specify the expectation of a function h . Accordingly, we may retrieve probabilities by taking the expectation of the indicator. First, by the law of total expectation:

$$\mathbb{E}_{P_n^*}[h(U, V)] = \sum_{\mathbf{x} \in \mathcal{X}} \sum_{\mathbf{z} \in \mathcal{Z}} \mathbb{E}_{P_n^*}[h(U, V) | (X, Z) = (\mathbf{x}, \mathbf{z})] P_n(\mathbf{x}, \mathbf{z}).$$

For an indicator random variable $h(\mathbf{x}, \mathbf{z}) = \mathbb{1}\{\mathbf{x} = \mathbf{u}\} \mathbb{1}\{\mathbf{z} = \mathbf{v}\}$, we may then compute

$$\begin{aligned} P_n^*(\mathbf{u}, \mathbf{v}) &= \sum_{\mathbf{x} \in \mathcal{X}} \sum_{\mathbf{z} \in \mathcal{Z}} \mathbb{E}_{P_n^*}[\mathbb{1}\{U = \mathbf{u}\} \mathbb{1}\{V = \mathbf{v}\} | (X, Z) = (\mathbf{x}, \mathbf{z})] P_n(\mathbf{x}, \mathbf{z}) \\ &= \sum_{\mathbf{x} \in \mathcal{X}} \sum_{\mathbf{z} \in \mathcal{Z}} \mathbb{E}_{P_n^*}[\mathbb{1}\{U = \mathbf{u}\} | X = \mathbf{x}] \mathbb{E}_{P_n^*}[\mathbb{1}\{V = \mathbf{v}\} | Z = \mathbf{z}] P_n(\mathbf{x}, \mathbf{z}) \\ &= \sum_{\mathbf{x} \in \mathcal{X}} \sum_{\mathbf{z} \in \mathcal{Z}} \frac{\pi_n^*(\mathbf{x}, \mathbf{u})}{P_{n,X}(\mathbf{x})} \cdot \frac{\gamma_n^*(\mathbf{z}, \mathbf{v})}{P_{n,Z}(\mathbf{z})} P_n(\mathbf{x}, \mathbf{z}). \end{aligned}$$

It is easy to verify that P_n^* satisfies the marginal constraints. However, given that this estimator relies on the optimal transport plans π_n^* and γ_n^* , we find ourselves in a similar dilemma as [Bickel et al. \[1991\]](#): it is challenging to compute the variance of estimators defined in terms of projections. Furthermore, while a solution (π_n^*, γ_n^*) exists, it may not be unique. We hypothesize that a similar approach to the recursion from (4.8) can be taken by replacing (π_n^*, γ_n^*) by the iterates of a procedure such as the Sinkhorn algorithm for solving entropy-regularized OT. The resulting couplings would also depend on a particular regularization parameter $\varepsilon \equiv \varepsilon_n > 0$. Computing the asymptotic variance of the estimator $P_n^*(h)$ based on scaling the number of steps of this procedure and the regularization parameter offers an interesting path for future work, especially if the estimator is efficient in the discrete case.

4.7 Perspectives & Future Work

This chapter considers the statistical analysis of foundation models, from pre-training to zero-shot prediction. Inspired by data curation procedures used in practice, we study the effect of known marginals for estimating linear functionals of distributions over multiple modalities.

This yielded a surprising connection to the classical work of [Bickel et al. \[1991\]](#), quantifying the variance reduction afforded by this additional information in the form of non-asymptotic mean squared error bounds. We also provided a formal discussion of prompting-based zero-shot prediction, phrased in terms of the dependence relations between random variables.

An appealing theoretical follow-up to the variance reduction results in [Section 4.3](#) would be to investigate the connections between the information projection viewpoint of the Sinkhorn iterations to an instance of mirror descent [[Léger, 2021](#), [Aubin-Frankowski et al., 2022](#), [Deb et al., 2023](#)]. In particular, the mirror descent viewpoint has yielded a sublinear $O(1/k)$ decay rate of the marginal violation in terms of Kullback-Leibler (KL) divergence (as opposed to the use of monotonicity of Sinkhorn iterations in [Proposition B.2.4](#)). The main benefit of such an analysis would be to meaningfully improve the dependence of the higher-order term of [\(4.10\)](#) on the iteration count k (as the current requirement for asymptotic efficiency is $k = o(n^{1/12})$). Part of the challenge of an approach like this is 1) that errors from each iteration accrue additively, as seen in the unrolled recursion [\(4.8\)](#), and 2) that the mirror descent approaches rely heavily on the use of the KL divergence. Noticing that for iteration k , the higher-order components include terms of the form $\left(\frac{P_X}{P_{n,X}^{(k-1)}}(\mathbf{x}) - 1\right)$, which may suggest the use of χ^2 -divergence as an alternative to KL. Owing to this, an exciting contribution would be to extend the mirror descent viewpoint to a broader class of information divergences, which would serve other communities such as statistical optimal transport as well.

Chapter 5

CONCLUSION

In this dissertation, we address several problems regarding statistical learning methods that exhibit out-of-distribution generalization. The resulting predictors may be applied to data from distributions other than the training distribution without a catastrophic increase in prediction error. As a methodological theme, we focus on the development and analysis of large-scale optimization procedures over the model parameters or even the distribution of the training data.

To summarize Chapter 2, we considered the distributionally robust optimization (DRO) problem, wherein the parameter of interest is estimated by minimizing the worst-case empirical risk achievable by reweighting the training examples. Using ideas such as progressive bias and variance reduction, we constructed stochastic algorithms with linear convergence guarantees for smoothed maximum-type problems with convex loss functions. We focused on several practical aspects of the problem, such as solving the dual maximization problem, selecting the uncertainty set/smoothing parameter, and extensions to changing supports and neural networks.

A question that can spawn many studies is the relationship between distributionally robust optimization and highly expressive, data-interpolating models such as neural networks. While linear models are amenable to fitting via convex optimization, one limitation to their use in distributionally robust optimization is the fact that the optimal primal solution might be close to (or even coincide with) the solution to the empirical risk minimization problem [Zhai et al., 2021]. On the other hand, while more expressive models may have fewer theoretical guarantees for the runtime of optimization procedures, it is of statistical interest whether the perturbations to the training distribution may increase the “effective” sample size of the

training set (akin to data augmentations). In fact, if a network is trained to zero loss on all training examples, then it will also minimize any distributionally robust objective. Because there are possibly many zero-loss parameter settings for overparametrized models, one may investigate whether distributionally robust objectives lead to such solutions with better generalization properties than the optimizers of the empirical risk minimization objective.

In Chapter 3, we augmented the algorithmic ideas of Chapter 2 with notions of adaptive sampling and historical regularization, which led to state-of-the-art algorithms for a broad class of “semilinear” min-max problems. This problem class is simultaneously a generalization of bilinearly coupled min-max problems and a special case of general nonbilinearly coupled min-max problems. When the dual feasible set is decomposable into separate feasible sets for each coordinate of the dual variable, we also apply coordinate-wise stochastic updates for improved complexity guarantees. While originally motivated by applications in distributionally robust optimization, a thorough experimental study of the proposed methods in the fully composite optimization and minimization with functional constraints (in the spirit of Chapter 2) would be a future direction of great practical interest.

Chapter 4 studies the effect of projecting an empirical measure onto the set of probability measures that satisfy particular marginal constraints. These marginal constraints come from prior knowledge of the data-generating distribution. We find that this procedure, which reflects data curation methods used in modern foundation modeling, provides a variance reduction for linear plug-in estimators based on this adjusted empirical measure. Other aspects of foundation modeling, such as downstream zero-shot prediction, are also discussed from the lens of conditional dependence relations between (unlabeled) pre-training and (labeled) task-specific data.

On downstream tasks, we focused on zero-shot prediction (often classification) as this is a canonical task for evaluating foundation models in terms of their pre-training data or architectures [Gadre et al., 2023]. However, the universal representation principle is used much more generally. Most commonly, the image and text embeddings can be applied for tasks such as retrieval or computing similarity/distance/kernel values between data points.

Deriving performance guarantees of zero-shot procedures against optimal predictors for these tasks would be an interesting direction for future work as well.

While many foundational questions remain, this dissertation aims to be a step toward out-of-distribution generalization being a common standard for statistical theory and methods across scientific disciplines.

Software

Links to all of the repositories below can be found at <https://ronakdm.github.io/software>.

Paper Reproducibility:

- `lrm` [Mehta et al., 2023].
- `prospect` [Mehta et al., 2024b].
- `drago` [Mehta et al., 2024a].
- `balancing` [Liu et al., 2024].
- `zeroshot` [Mehta and Harchaoui, 2025].

Standalone Packages:

- `deshift`: Instance-level and group-level distributionally robust optimization for CPU/GPU PyTorch workflows, with support for data distributed computing.
- `drlearn`: Distributionally robust linear predictors in the scikit-learn interface.

BIBLIOGRAPHY

- Carlo Acerbi and Dirk Tasche. On the coherence of expected shortfall. *Journal of Banking & Finance*, 2002.
- Alekh Agarwal, Alina Beygelzimer, Miroslav Dudik, John Langford, and Hanna Wallach. A Reductions Approach to Fair Classification. In *ICML*, 2018.
- Alekh Agarwal, Miroslav Dudik, and Zhiwei Steven Wu. Fair Regression: Quantitative Definitions and Reduction-Based Algorithms. In *ICML*, 2019.
- Medha Agarwal, Kasim Rafiq, Ronak Mehta, Briana Abrahms, and Zaid Harchaoui. Leveraging machine learning and accelerometry to classify animal behaviours with uncertainty. *Methods in Ecology and Evolution*, 2024.
- Uchenna Akujuobi and Xiangliang Zhang. Delve: A Dataset-Driven Scholarly Search and Analysis System. *SIGKDD Explor. Newsl.*, 19, 2017.
- Ahmet Alacaoglu and Yura Malitsky. Stochastic Variance Reduction for Variational Inequality Methods. In *COLT*, 2022.
- Ahmet Alacaoglu, Quoc Tran Dinh, Olivier Fercoq, and Volkan Cevher. Smooth Primal-Dual Coordinate Descent Algorithms for Nonsmooth Convex Optimization. In *NeurIPS*, 2017.
- Ahmet Alacaoglu, Olivier Fercoq, and Volkan Cevher. Random extrapolation for primal-dual coordinate descent. In *ICML*, 2020.
- Ahmet Alacaoglu, Volkan Cevher, and Stephen J. Wright. On the Complexity of a Practical Primal-Dual Coordinate Method, 2022. arXiv Technical Report.

- Mickael Albertus and Philippe Berthet. Auxiliary information: the raking-ratio empirical process. *Electronic Journal of Statistics*, 2019.
- Zeinab Alizadeh, Erfan Yazdandoost Hamedani, and Afrooz Jalilzadeh. Variance-reduction for Variational Inequality Problems with Bregman Distance Function, 2024.
- Philippe Artzner, Freddy Delbaen, Eber Jean-Marc, and David Heath. Coherent Measures of Risk. *Mathematical Finance*, 1999.
- YM. Asano, C. Rupprecht, and A. Vedaldi. Self-labelling via simultaneous clustering and representation learning. In *ICLR*, 2020.
- Mahmoud Assran, Quentin Duval, Ishan Misra, Piotr Bojanowski, Pascal Vincent, Michael Rabbat, Yann LeCun, and Nicolas Ballas. Self-Supervised Learning from Images with a Joint-Embedding Predictive Architecture. In *CVPR*, 2023.
- Pierre-Cyril Aubin-Frankowski, Anna Korba, and Flavien Léger. Mirror Descent with Relative Smoothness in Measure Spaces, with application to Sinkhorn and EM. In *NeurIPS*, 2022.
- Francis Bach. *Learning Theory from First Principles*. MIT Press, 2024.
- Michel Baes, Michael Burgisser, and Arkadi Nemirovski. A randomized mirror-prox method for solving structured large-scale matrix saddle-point problems. *SIAM Journal on Optimization*, 2013.
- Randall Balestriero and Yann LeCun. Contrastive and Non-Contrastive Self-Supervised Learning Recover Global and Local Spectral Embedding Methods. In *NeurIPS*, 2022.
- Adrien Bardes, Jean Ponce, and Yann LeCun. VICReg: Variance-Invariance-Covariance Regularization for Self-Supervised Learning. In *ICLR*, 2022.

- Sandi Baressi Segota, Nikola Andelic, Jan Kudlacek, and Robert Cep. Artificial Neural Network for Predicting Values of Residuary Resistance per Unit Weight of Displacement. *Journal of Maritime & Transportation Science*, 57, 2020.
- Sara Beery, Elijah Cole, and Arvi Gjoka. The iWildCam 2020 Competition Dataset. *arXiv preprint arXiv:2004.10340*, 2020.
- Aharon Ben-Tal and Marc Teboulle. An Old-New Concept of Convex Risk Measures: The Optimized Certainty Equivalent. *Mathematical Finance*, 2007.
- Dimitri P Bertsekas. Nonlinear Programming. *Journal of the Operational Research Society*, 1997.
- Sanjay P. Bhat and L. A. Prashanth. Concentration of risk measures: A Wasserstein distance approach. In *NeurIPS*, 2019.
- Peter J. Bickel, Ya'Acov Ritov, and Jon A. Wellner. Efficient Estimation of Linear Functionals of a Probability Measure P with Known Marginal Distributions. *The Annals of Statistics*, 1991.
- David A. Binder. On the Variances of Asymptotically Normal Estimators from Complex Surveys. *International Statistical Review*, 1983.
- Sergey G. Bobkov and Michel Ledoux. One-Dimensional Empirical Measures, Order Statistics, and Kantorovich Transport Distances. *Memoirs of the American Mathematical Society*, 2019.
- Digvijay Boob and Mohammad Khalafi. Optimal Primal-Dual Algorithm with Last iterate Convergence Guarantees for Stochastic Convex Optimization Problems, 2024. arXiv Technical Report.
- Ekaterina Borodich, Georgiy Kormakov, Dmitry Kovalev, Aleksandr Beznosikov, and

- Alexander Gasnikov. Near-Optimal Algorithm with Complexity Separation for Strongly Convex-Strongly Concave Composite Saddle Point Problems. In *ICOMP*, 2024.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language Models are Few-Shot Learners. In *NeurIPS*, 2020.
- Andreas Buja. Remarks on Functional Canonical Variates, Alternating Least Squares Methods and ACE. *The Annals of Statistics*, 1990.
- Joy Buolamwini and Timnit Gebru. Gender Shades: Intersectional Accuracy Disparities in Commercial Gender Classification. In *ACM Conference on Fairness, Accountability and Transparency*, 2018.
- Jonathon Byrd and Zachary C. Lipton. What is the Effect of Importance Weighting in Deep Learning? In *ICML*, 2019.
- Xufeng Cai, Ahmet Alacaoglu, and Jelena Diakonikolas. Variance Reduced Halpern Iteration for Finite-Sum Monotone Inclusions. In *ICLR*, 2024.
- Yair Carmon, Yujia Jin, Aaron Sidford, and Kevin Tian. Variance Reduction for Matrix Games. In *NeurIPS*, 2019.
- Yair Carmon, Arun Jambulapati, Yujia Jin, and Aaron Sidford. RECAPP: Crafting a More Efficient Catalyst for Convex Optimization. In *ICML*, 2022.
- Mathilde Caron, Ishan Misra, Julien Mairal, Priya Goyal, Piotr Bojanowski, and Armand Joulin. Unsupervised learning of visual features by contrasting cluster assignments. In *NeurIPS*, 2020.

- Antonin Chambolle and Thomas Pock. A First-Order Primal-Dual Algorithm for Convex Problems with Applications to Imaging. *Journal of Mathematical Imaging and Vision*, 2011.
- Antonin Chambolle, Matthias J. Ehrhardt, Peter Richtárik, and Carola-Bibiane Schönlieb. Stochastic Primal-Dual Hybrid Gradient Algorithm with Arbitrary Sampling and Imaging Applications. *SIAM Journal on Optimization*, 2018.
- Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *ICML*, 2020.
- M. Cimpoi, S. Maji, I. Kokkinos, S. Mohamed, , and A. Vedaldi. Describing Textures in the Wild. In *CVPR*, 2014.
- Laurent Condat. Fast projection onto the simplex and the ℓ_1 -ball. *Mathematical Programming*, 2016.
- Peter Coppens and Panagiotis Patrinos. Ordered Risk Minimization: Learning More from Less Data. *IEEE Conference on Decision and Control (CDC)*, 2023.
- Corinna Cortes and Vladimir Vapnik. Support-vector networks. *Machine Learning*, 1995.
- John Cotter and Kevin Dowd. Extreme Spectral Risk Measures: an Application to Futures Clearinghouse Margin Requirements. *Journal of Banking & Finance*, 2006.
- Thomas M Cover. *Elements of Information Theory*. John Wiley & Sons, 1999.
- Ying Cui, Jong-Shi Pang, and Bodhisattva Sen. Composite Difference-Max Programs for Modern Statistical Estimation Problems. *SIAM Journal on Optimization*, 2018.
- Sebastian Curi, Kfir Y. Levy, Stefanie Jegelka, and Andreas Krause. Adaptive Sampling for Stochastic Risk-Averse Learning. In *NeurIPS*, 2020.

- Joshua Cutler, Dmitriy Drusvyatskiy, and Zaid Harchaoui. Stochastic Optimization under Distributional Drift. *Journal of Machine Learning Research*, 2023.
- Abdelaati Daouia, Irène Gijbels, and Gilles Stupfler. Extremiles: A New Perspective on Asymmetric Least Squares. *Journal of the American Statistical Association*, 2019.
- H. A. David and H. N. Nagaraja. *Order Statistics*. Wiley, 3rd edition, 2003.
- Nabarun Deb, Young-Heon Kim, Soumik Pal, and Geoffrey Schiebinger. Wasserstein Mirror Gradient Flow as the limit of the Sinkhorn Algorithm. arXiv Technical Report, 2023.
- Aaron Defazio, Francis Bach, and Simon Lacoste-Julien. SAGA: A Fast Incremental Gradient Method With Support for Non-Strongly Convex Composite Objectives. *NeurIPS*, 2014.
- W. Edwards Deming and Frederick F. Stephan. On a Least Squares Adjustment of a Sampled Frequency Table When the Expected Marginal Totals are Known. *The Annals of Mathematical Statistics*, 1940.
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. ImageNet: A Large-scale Hierarchical Image Database. In *CVPR*, 2009.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *NAACL*, 2019a.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina” Toutanova. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *ACL*, 2019b.
- Jelena Diakonikolas. A Block Coordinate and Variance-Reduced Method for Generalized Variational Inequalities of Minty Type. *JDS*, 2025.
- Jelena Diakonikolas and Lorenzo Orecchia. The approximate duality gap technique: A unified theory of first-order methods. *SIAM Journal on Optimization*, 2019.

- Frances Ding, Moritz Hardt, John Miller, and Ludwig Schmidt. Retiring Adult: New Datasets for Fair Machine Learning. In *NeurIPS*, volume 34. Curran Associates, Inc., 2021.
- Nikita Doikov and Yurii Nesterov. High-Order Optimization Methods for Fully Composite Problems. *SIAM Journal on Optimization*, 2022.
- D. Drusvyatskiy and C. Paquette. Efficiency of minimizing compositions of convex functions and smooth maps. *Mathematical Programming*, 2019.
- John C. Duchi and Hongseok Namkoong. Variance-based Regularization with Convex Objectives. *Journal of Machine Learning Research*, 2019.
- John C. Duchi and Hongseok Namkoong. Learning models with uniform performance via distributionally robust optimization. *The Annals of Statistics*, 2021.
- John C. Duchi, Peter W. Glynn, and Hongseok Namkoong. Statistics of Robust Optimization: A Generalized Empirical Likelihood Approach. *Mathematics of Operations Research*, 2021.
- Alex Dytso, Martina Cardone, and Ian Zieder. Meta Derivative Identity for the Conditional Expectation. *IEEE Transactions on Information Theory*, 2023.
- Yanbo Fan, Siwei Lyu, Yiming Ying, and Baogang Hu. Learning with Average Top- k Loss. In *NeurIPS*, 2017.
- Alex Fang, Gabriel Ilharco, Mitchell Wortsman, Yuhao Wan, Vaishaal Shankar, Achal Dave, and Ludwig Schmidt. Data determines distributional robustness in contrastive language-image pre-training (CLIP). In *ICML*, 2023.
- Tongtong Fang, Nan Lu, Gang Niu, and Masashi Sugiyama. Rethinking Importance Weighting for Deep Learning under Distribution Shift. In *NeurIPS*, 2020.

- Olivier Fercoq. Quadratic error bound of the smoothed gap and the restarted averaged primal-dual hybrid gradient. *Open Journal of Mathematical Optimization*, 2023.
- Hans Föllmer and Alexander Schied. Convex measures of risk and trading constraints. *Finance Stochastics*, 2002.
- Samir Yitzhak Gadre, Gabriel Ilharco, Alex Fang, Jonathan Hayase, Georgios Smyrnis, Thao Nguyen, Ryan Marten, Mitchell Wortsman, Dhruva Ghosh, Jieyu Zhang, Eyal Or-gad, Rahim Entezari, Giannis Daras, Sarah M Pratt, Vivek Ramanujan, Yonatan Bit-ton, Kalyani Marathe, Stephen Mussmann, Richard Vencu, Mehdi Cherti, Ranjay Kr-ishna, Pang Wei Koh, Olga Saukh, Alexander Ratner, Shuran Song, Hannaneh Hajishirzi, Ali Farhadi, Romain Beaumont, Sewoong Oh, Alex Dimakis, Jenia Jitsev, Yair Carmon, Vaishaal Shankar, and Ludwig Schmidt. DataComp: In search of the next generation of multimodal datasets. In *NeurIPS*, 2023.
- I. Gohberg, S. Goldberg, and M.A. Kaashoek. *Classes of Linear Operators Vol. 1*. Springer, 1990.
- Mingming Gong, Kun Zhang, Tongliang Liu, Dacheng Tao, Clark Glymour, and Bernhard Schölkopf. Domain Adaptation with Conditional Transferable Components. In *ICML*, 2016.
- Robert M. Gower, Mark Schmidt, Francis Bach, and Peter Richtárik. Variance-Reduced Methods for Machine Learning. *Proceedings of the IEEE*, 2020.
- Arthur Gretton, Alex Smola, Jiayuan Huang, Marcel Schmittfull, Karsten Borgwardt, and Bernhard Schölkopf. Covariate Shift by Kernel Mean Matching. In *Dataset Shift in Machine Learning*. The MIT Press, 2008.
- Vincent Guigues and Claudia A. Sagastizábal. Risk-averse feasible policies for large-scale multistage stochastic linear programs. *Mathematical Programming*, 2013.

- Michael U. Gutmann and Aapo Hyvärinen. Noise-Contrastive Estimation of Unnormalized Statistical Models, with Applications to Natural Image Statistics. *JMLR*, 2012.
- Jeff Z. HaoChen, Colin Wei, Adrien Gaidon, and Tengyu Ma. Provable Guarantees for Self-Supervised Deep Learning with Spectral Contrastive Loss. In *NeurIPS*, 2021.
- Tatsunori Hashimoto, Megha Srivastava, Hongseok Namkoong, and Percy Liang. Fairness without Demographics in Repeated Loss Minimization. In *ICML*, 2018.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep Residual Learning for Image Recognition. In *CVPR*, 2016.
- Xue Dong He, Steven Kou, and Xianhua Peng. Risk Measures: Robustness, Elicitability, and Backtesting. *Annual Review of Statistics and Its Application*, 2022.
- Jean-Baptiste Hiriart-Urruty and Claude Lemaréchal. *Fundamentals of Convex Analysis*. Springer, 2004.
- Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, Tom Hennigan, Eric Noland, Katie Millican, George van den Driessche, Bogdan Damoc, Aurelia Guy, Simon Osindero, Karen Simonyan, Erich Elsen, Jack W. Rae, Oriol Vinyals, and Laurent Sifre. Training Compute-Optimal Large Language Models. arXiv Technical Report, 2022.
- Matthew J. Holland and El Mehdi Haress. Spectral risk-based learning using unbounded losses. In *AISTATS*, 2022.
- Roger A. Horn and Charles R. Johnson. *Matrix Analysis*. Cambridge University Press, 2013.
- Jiayuan Huang, Arthur Gretton, Karsten Borgwardt, Bernhard Schölkopf, and Alex Smola. Correcting Sample Selection Bias by Unlabeled Data. In *NeurIPS*, 2006.

- Gabriel Ilharco, Mitchell Wortsman, Ross Wightman, Cade Gordon, Nicholas Carlini, Rohan Taori, Achal Dave, Vaishaal Shankar, Hongseok Namkoong, John Miller, Hannaneh Hajishirzi, Ali Farhadi, and Ludwig Schmidt. OpenCLIP. GitHub Repository, 2022.
- C. T. Ireland and S. Kullback. Contingency tables with given marginals. *Biometrika*, 1968.
- Yujia Jin, Aaron Sidford, and Kevin Tian. Sharper Rates for Separable Minimax and Finite Sum Optimization via Primal-Dual Extragradient Methods. In *COLT*, 2022.
- Rie Johnson and Tong Zhang. Accelerating Stochastic Gradient Descent using Predictive Variance Reduction. In *NeurIPS*, 2013.
- R. J. Johnston and C. J. Pattie. Entropy-maximizing and the iterative proportional fitting procedure. *The Professional Geographer*, 1993.
- Corinne Jones, Vincent Roulet, and Zaid Harchaoui. Discriminative clustering with representation learning with any ratio of labeled to unlabeled data. *Statistics and Computing*, 2022.
- Anatoli Juditsky, Arkadi Nemirovski, et al. First-Order Methods for Nonsmooth Convex Large-Scale Optimization, II: Utilizing Problem’s Structure. *Optimization for Machine Learning*, 2011.
- Anatoli Juditsky, Fatma Kılınç Karzan, and Arkadi Nemirovski. Randomized first order algorithm with applications to ℓ_1 -minimization. *Mathematical Programming*, 2013.
- H. Kahn and A. W. Marshall. Methods of Reducing Sample Size in Monte Carlo Computations. *Journal of the Operations Research Society of America*, 1953.
- Takafumi Kanamori, Shohei Hido, and Masashi Sugiyama. A least-squares approach to direct importance estimation. *Journal of Machine Learning Research*, 2009.

- Takafumi Kanamori, Taiji Suzuki, and Masashi Sugiyama. Computational complexity of kernel-based density-ratio estimation: a condition number analysis. *Machine Learning*, 2013.
- Masahiro Kato and Takeshi Teshima. Non-Negative Bregman Divergence Minimization for Deep Direct Density Ratio Estimation. In *ICML*, 2021.
- Kenji Kawaguchi and Haihao Lu. Ordered SGD: A New Stochastic Optimization Framework for Empirical Risk Minimization. In *AISTATS*, 2020.
- Justin Khim, Liu Leqi, Adarsh Prasad, and Pradeep Ravikumar. Uniform Convergence of Rank-weighted Learning. In *ICML*, 2020.
- Diederik Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *ICLR*, 2015.
- Pang Wei Koh, Shiori Sagawa, Henrik Marklund, Sang Michael Xie, Marvin Zhang, Akshay Balsubramani, Weihua Hu, Michihiro Yasunaga, Richard Lanus Phillips, Irena Gao, Tony Lee, Etienne David, Ian Stavness, Wei Guo, Berton A. Earnshaw, Imran S. Haque, Sara Beery, Jure Leskovec, Anshul Kundaje, Emma Pierson, Sergey Levine, Chelsea Finn, and Percy Liang. WILDS: A Benchmark of in-the-Wild Distribution Shifts, 2021.
- G. M. Korpelevich. An Extragradient Method for Finding Saddle Points and for Other Problems. *Ekonomika i Matematicheskie Metody*, 1976.
- Dmitry Kovalev and Alexander Gasnikov. The First Optimal Algorithm for Smooth and Strongly-Convex-Strongly-Concave Minimax Optimization. In *NeurIPS*, 2022.
- Alex Krizhevsky. Learning Multiple Layers of Features from Tiny Images. University of Toronto Technical Report, 2009.
- Arun Kumar Kuchibhotla and Abhishek Chakraborty. Moving beyond sub-Gaussianity in

- high-dimensional statistics: applications in covariance estimation and linear regression. *Information and Inference: A Journal of the IMA*, 2022.
- Daniel Kuhn, Peyman Mohajerin Esfahani, Viet Anh Nguyen, and Soroosh Shafieezadeh-Abadeh. Wasserstein Distributionally Robust Optimization: Theory and Applications in Machine Learning. *Operations research & management science in the age of analytics*, 2019.
- Ramnath Kumar, Kushal Alpesh Majmundar, Dheeraj Mysore Nagaraj, and Arun Suggala. Stochastic re-weighted gradient descent via distributionally robust optimization. *Transactions on Machine Learning Research*, 2024.
- Yassine Laguel, Jerome Malick, and Zaid Harchaoui. First-Order Optimization for Superquantile-Based Supervised Learning. In *IEEE MLSP*, 2020.
- Yassine Laguel, Krishna Pillutla, Jérôme Malick, and Zaid Harchaoui. Superquantiles at Work: Machine Learning Applications and Efficient Subgradient Computation. *Set-Valued and Variational Analysis*, 2021.
- Henry Lam and Enlu Zhou. The empirical likelihood approach to quantifying uncertainty in sample average approximation. *Operations Research Letters*, 2017.
- Guanghui Lan and Yan Li. A Novel Catalyst Scheme for Stochastic Minimax Optimization, 2023. arXiv Technical Report.
- Steffen L. Lauritzen. *Graphical Models*. Oxford University Press, 1996.
- Nicolas Le Roux, Mark Schmidt, and Francis Bach. A Stochastic Gradient Method with an Exponential Convergence Rate for Finite Training Sets. In *NeurIPS*, 2012.
- Thanh Quoc Trinh Le Thi Thanh Hai and Phan Tu Vuong. A refined convergence analysis of Popov’s algorithm for pseudo-monotone variational inequalities. *Optimization*, 2025.

- Jaeho Lee, Sejun Park, and Jinwoo Shin. Learning Bounds for Risk-sensitive Learning. In *NeurIPS*, 2020.
- Jason D Lee, Qi Lei, Nikunj Saunshi, and Jiacheng Zhuo. Predicting What You Already Know Helps: Provable Self-Supervised Learning. In *NeurIPS*, 2021.
- Flavien Léger. A Gradient Descent Perspective on Sinkhorn. *Applied Mathematics & Optimization*, 2021.
- Liu Leqi, Adarsh Prasad, and Pradeep K Ravikumar. On Human-Aligned Risk Minimization. In *NeurIPS*, 2019.
- Daniel Levy, Yair Carmon, John Duchi, and Aaron Sidford. Large-Scale Methods for Distributionally Robust Optimization. In *NeurIPS*, 2020.
- Chris Junchi Li, Angela Yuan, Gauthier Gidel, Quanquan Gu, and Michael I. Jordan. Nesterov Meets Optimism: Rate-Optimal Separable Minimax Optimization. In *ICML*, 2023.
- Jeffrey Li, Alex Fang, Georgios Smyrnis, Maor Ivgi, Matt Jordan, Samir Yitzhak Gadre, Hritik Bansal, Etash Kumar Guha, Sedrick Keh, Kushal Arora, Saurabh Garg, Rui Xin, Niklas Muennighoff, Reinhard Heckel, Jean Mercat, Mayee F Chen, Suchin Gururangan, Mitchell Wortsman, Alon Albalak, Yonatan Bitton, Marianna Nezhurina, Amro Kamal Mohamed Abbas, Cheng-Yu Hsieh, Dhruba Ghosh, Joshua P Gardner, Maciej Kilian, Hanlin Zhang, Rulin Shao, Sarah M Pratt, Sunny Sanyal, Gabriel Ilharco, Giannis Daras, Kalyani Marathe, Aaron Gokaslan, Jieyu Zhang, Khyathi Chandu, Thao Nguyen, Igor Vasiljevic, Sham M. Kakade, Shuran Song, Sujay Sanghavi, Fartash Faghri, Sewoong Oh, Luke Zettlemoyer, Kyle Lo, Alaaeldin El-Nouby, Hadi Pouransari, Alexander T Toshev, Stephanie Wang, Dirk Groeneveld, Luca Soldaini, Pang Wei Koh, Jenia Jitsev, Thomas Kollar, Alex Dimakis, Yair Carmon, Achal Dave, Ludwig Schmidt, and Vaishaal Shankar. DataComp-LM: In search of the next generation of training sets for language models. In *NeurIPS*, 2024a.

- Shuyao Li, Sushrut Karmalkar, Ilias Diakonikolas, and Jelena Diakonikolas. Learning a single neuron robustly to distributional shifts and adversarial label noise. In *NeurIPS*, 2024b.
- Hongzhou Lin, Julien Mairal, and Zaid Harchaoui. Catalyst acceleration for first-order convex optimization: from theory to practice. *Journal of Machine Learning Research*, 2018.
- Tianyi Lin, Chi Jin, and Michael I. Jordan. Near-Optimal Algorithms for Minimax Optimization. In *COLT*, 2020.
- Lang Liu, Ronak Mehta, Soumik Pal, and Zaid Harchaoui. The Benefits of Balance: From Information Projections to Variance Reduction. In *NeurIPS*, 2024.
- Cong Ma, Reese Pathak, and Martin J. Wainwright. Optimally tackling covariate shift in RKHS-based nonparametric regression. *The Annals of Statistics*, 2023.
- A. Rupam Mahmood, Hado P van Hasselt, and Richard S Sutton. Weighted importance sampling for off-policy learning with linear function approximation. In *NeurIPS*, 2014.
- Julien Mairal. Incremental Majorization-Minimization Optimization with Application to Large-Scale Machine Learning. *SIAM Journal on Optimization*, 2014.
- Subhransu Maji, Esa Rahtu, Juho Kannala, Matthew Blaschko, and Andrea Vedaldi. Fine-Grained Visual Classification of Aircraft. arXiv Technical Report, 2013.
- Patrice Marcotte. Application Of Khobotov’s Algorithm To Variational Inequalities And Network Equilibrium Problems. *INFOR: Information Systems and Operational Research*, 1991.
- Andreas Maurer, Daniela Angela Parletta, Andrea Paudice, and Massimiliano Pontil. Robust Unsupervised Learning via L-statistic Minimization. In *ICML*, 2021.
- Ronak Mehta and Zaid Harchaoui. A Generalization Theory for Zero-Shot Prediction. In *ICML*, 2025.

- Ronak Mehta, Vincent Roulet, Krishna Pillutla, Lang Liu, and Zaid Harchaoui. Stochastic Optimization for Spectral Risk Measures. In *AISTATS*, 2023.
- Ronak Mehta, Jelena Diakonikolas, and Zaid Harchaoui. Drago: Primal-Dual Coupled Variance Reduction for Faster Distributionally Robust Optimization. In *NeurIPS*, 2024a.
- Ronak Mehta, Vincent Roulet, Krishna Pillutla, and Zaid Harchaoui. Distributionally Robust Optimization with Bias and Variance Reduction. In *ICLR*, 2024b.
- Ronak Mehta, Jelena Diakonikolas, and Zaid Harchaoui. Min-Max Optimization with Dual-Linear Coupling, 2025. arXiv Technical Report.
- Thomas Mikosch. Regular variation, subexponentiality and their applications in probability theory. *International Journal of Production Economics*, 1999.
- Douglas J. Miller and Wei han Liu. On the recovery of joint distributions from limited information. *Journal of Econometrics*, 2002.
- Aryan Mokhtari, Asuman Ozdaglar, and Sarath Pattathil. A Unified Analysis of Extra-gradient and Optimistic Gradient Methods for Saddle Point Problems: Proximal Point Approach. In *AISTATS*, 2020.
- Debarghya Mukherjee, Felix Petersen, Mikhail Yurochkin, and Yuekai Sun. Domain Adaptation meets Individual Fairness. And they get along. In *NeurIPS*, 2022.
- Hongseok Namkoong and John C Duchi. Stochastic Gradient Methods for Distributionally Robust Optimization with f-divergences. In *NeurIPS*, 2016.
- Deanna Needell, Nathan Srebro, and Rachel Ward. Stochastic Gradient Descent, Weighted Sampling, and the Randomized Kaczmarz algorithm. In *NeurIPS*, 2014.
- Arkadi Nemirovski. Prox-Method with Rate of Convergence $O(1/t)$ for Variational Inequalities with Lipschitz Continuous Monotone Operators and Smooth Convex-Concave Saddle Point Problems. *SIAM Journal on Optimization*, 2004.

- Yurii Nesterov. Dual extrapolation and its applications to solving variational inequalities and related problems. *Mathematical Programming*, 2007a.
- Yurii Nesterov. Smoothing Technique and its Applications in Semidefinite Optimization. *Mathematical Programming*, 2007b.
- Yurii Nesterov. *Lectures on Convex Optimization*. Springer, 2nd edition, 2018.
- Yurii Nesterov and Laura Scramali. Solving strongly monotone variational and quasi-variational inequalities. *Core Discussion Paper*, 2006.
- Jianmo Ni, Jiacheng Li, and Julian McAuley. Justifying Recommendations using Distantly-Labeled Reviews and Fine-grained Aspects. In *EMNLP*, 2019.
- Maria-Elena Nilsback and Andrew Zisserman. Automated Flower Classification over a Large Number of Classes. In *Indian Conference on Computer Vision, Graphics and Image Processing*, 2008.
- Jorge Nocedal and Stephen J Wright. *Numerical optimization*. Springer, 1999.
- Marcel Nutz. Introduction to Entropic Optimal Transport. *Lecture notes, Columbia University*, 2021.
- Ilsang Ohn and Yongdai Kim. Toward fast rates: A review of localization analysis for statistical learning, 2025. Technical report.
- Kazusato Oko, Licong Lin, Yuhang Cai, and Song Mei. A Statistical Theory of Contrastive Pre-training and Multimodal Generative AI. arXiv Technical Report, 2025.
- Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy V. Vo, Marc Szafraniec, Vasil Khalidov, Pierre Fernandez, Daniel HAZIZA, Francisco Massa, Alaaeldin El-Nouby, Mido Assran, Nicolas Ballas, Wojciech Galuba, Russell Howes, Po-Yao Huang, Shang-Wen Li, Ishan Misra, Michael Rabbat, Vasu Sharma, Gabriel Synnaeve, Hu Xu, Herve Jegou,

- Julien Mairal, Patrick Labatut, Armand Joulin, and Piotr Bojanowski. DINOv2: Learning Robust Visual Features without Supervision. *TMLR*, 2024.
- Art Owen. Empirical Likelihood Ratio Confidence Regions. *The Annals of Statistics*, 1990.
- Art Owen. *Empirical Likelihood*. CRC Press, 1st edition, 2001.
- Balamurugan Palaniappan and Francis Bach. Stochastic Variance Reduction Methods for Saddle-Point Problems. *NeurIPS*, 2016.
- Ajay Kumar Pandey, L. A. Prashanth, and Sanjay P. Bhat. Estimation of Spectral Risk Measures. In *AAAI Conference on Artificial Intelligence*, 2019.
- Giuditta Parolini. The Emergence of Modern Statistics in Agricultural Science: Analysis of Variance, Experimental Design and the Reshaping of Research at Rothamsted Experimental Station, 1919-1933. *Journal of the History of Biology*, 2015.
- Gabriel Peyré and Marco Cuturi. Computational Optimal Transport: With Applications to Data Science. *Foundations and Trends in Machine Learning*, 2019.
- Georg Ch Pflug and Andrzej Ruszczyński. Measuring Risk for Income Streams. *Computational Optimization and Applications*, 2005.
- Alexander Pichugin, Maksim Pechin, Aleksandr Beznosikov, Vasilii Novitskii, and Alexander Gasnikov. Method with batching for stochastic finite-sum variational inequalities in non-Euclidean setting. *Chaos, Solitons & Fractals*, 2024.
- L. D. Popov. A modification of the Arrow-Hurwicz method for search of saddle points. *Mathematical notes of the Academy of Sciences of the USSR*, 1980.
- L. A. Prashanth and Sanjay P. Bhat. A Wasserstein Distance Approach for Concentration of Empirical Risk Estimates. *Journal of Machine Learning Research*, 2022.

- Sarah Pratt, Ian Covert, Rosanne Liu, and Ali Farhadi. What does a platypus look like? generating customized prompts for zero-shot image classification. In *ICCV*, 2023.
- Joaquin Quiñonero-Candela, Masashi Sugiyama, Anton Schwaighofer, and Neil D. Lawrence. *Dataset Shift in Machine Learning*. MIT Press, 2022.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *ICML*, 2021.
- A. Rényi. On measures of dependence. *Acta Mathematica Academiae Scientiarum Hungarica*, 1959.
- Ali Rizvi, Beata Strack, Jonathan P. DeShazo, Chris Gennings, Juan L. Olmo, Sebastian Ventura, Krzysztof J. Cios, and John N. Clore. Impact of HbA1c Measurement on Hospital Readmission Rates: Analysis of 70,000 Clinical Database Patient Records. *BioMed Research International*, 2014, 2014.
- R. T. Rockafellar and J. O. Royset. Random Variables, Monotone Relations, and Convex Analysis. *Mathematical Programming*, 2014.
- R Tyrrell Rockafellar and Johannes O Royset. Superquantiles and Their Applications to Risk, Random Variables, and Regression. In *Theory Driven by Influential Applications*. INFORMS, 2013.
- R. Tyrrell Rockafellar and Stanislav Uryasev. The fundamental risk quadrangle in risk management, optimization and statistical estimation. *Surveys in Operations Research and Management Science*, 18:33–53, 2013.
- Paul R. Rosenbaum and Donald B. Rubin. The Central Role of the Propensity Score in Observational Studies for Causal Effects. *Biometrika*, 1983.

- Shiori Sagawa, Pang Wei Koh, Tatsunori B. Hashimoto, and Percy Liang. Distributionally Robust Neural Networks. In *ICLR*, 2020.
- René Schilling. *Measures, Integrals, and Martingales*. Springer, 2nd edition, 2017.
- Shai Shalev-Shwartz and Tong Zhang. Stochastic Dual Coordinate Ascent Methods for Regularized Loss. *Journal of Machine Learning Research*, 2013.
- Jun Shao. Functional calculus and asymptotic theory for statistical analysis. *Statistics & Probability Letters*, 1989.
- Alexander Shapiro. Distributionally Robust Stochastic Programming. *SIAM Journal on Optimization*, 2017.
- Hidetoshi Shimodaira. Improving predictive inference under covariate shift by weighting the log-likelihood function. *Journal of Statistical Planning and Inference*, 2000.
- Galen Shorack. *Probability for Statisticians*. Springer, 2017.
- Ravindra Singh and Naurang Singh Mangat. *Elements of Survey Sampling*. Springer, 1996.
- Richard Sinkhorn. Diagonal equivalence to matrices with prescribed row and column sums. *American Mathematical Monthly*, 74(4), 1967.
- Chaobing Song, Stephen J Wright, and Jelena Diakonikolas. Variance Reduction via Primal-Dual Accelerated Dual Averaging for Nonsmooth Convex Finite-Sums. In *ICML*, 2021.
- Masashi Sugiyama, Matthias Krauledat, and Klaus-Robert Müller. Covariate Shift Adaptation by Importance Weighted Cross Validation. *Journal of Machine Learning Research*, 2007.
- Masashi Sugiyama, Shinichi Nakajima, Hisashi Kashima, Paul Buenau, and Motoaki Kawanabe. Direct Importance Estimation with Model Selection and Its Application to Covariate Shift Adaptation. In *NeurIPS*, 2008.

- Adith Swaminathan and Thorsten Joachims. Batch Learning from Logged Bandit Feedback through Counterfactual Risk Minimization. *Journal of Machine Learning Research*, 2015.
- Zhiquan Tan, Yifan Zhang, Jingqin Yang, and Yang Yuan. Contrastive Learning is Spectral Clustering on Similarity Graph. In *ICLR*, 2024.
- Bart Thomee, David A. Shamma, Gerald Friedland, Benjamin Elizalde, Karl Ni, Douglas Poland, Damian Borth, and Li-Jia Li. YFCC100M: the New Data in Multimedia Research. *Communications of the ACM*, 2016.
- M. E. Thompson. *Theory of Sample Surveys*. Chapman & Hall, 2000.
- Christopher Tosh, Akshay Krishnamurthy, and Daniel Hsu. Contrastive learning, multi-view redundancy, and linear models. In *ALT*, 2021.
- Athanasios Tsanas and Angeliki Xifara. Accurate Quantitative Estimation of Energy Performance of Residential Buildings Using Statistical Machine Learning Tools. *Energy and Buildings*, 49, 2012.
- Paul Tseng. On linear convergence of iterative methods for the variational inequality problem. *Journal of Computational and Applied Mathematics*, 1995.
- Pınar Tüfekci. Prediction of Full Load Electrical Power Output of a Base Load Operated Combined Cycle Power Plant using Machine Learning Methods. *International Journal of Electrical Power & Energy Systems*, 60, 2014.
- Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation Learning with Contrastive Predictive Coding. arXiv Technical Report, 2019.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is All you Need. In *NeurIPS*, 2017.
- Alexander Visheratin. NLLB-CLIP - train performant multilingual image retrieval model on a budget. In *NeurIPS Workshop: ENLSP-III*, 2023.

- Maria-Luiza Vladarean, Nikita Doikov, Martin Jaggi, and Nicolas Flammarion. Linearization Algorithms for Fully Composite Optimization. In *COLT*, 2023.
- Mariia Vladimirova, Stéphane Girard, Hien Nguyen, and Julyan Arbel. Sub-Weibull distributions: Generalizing sub-Gaussian and sub-Exponential properties to heavier tailed distributions. *Stat*, 2020.
- Hieu Vu, Toan Tran, Man-Chung Yue, and Viet Anh Nguyen. Distributionally Robust Fair Principal Components via Geodesic Descents. In *ICLR*, 2022.
- Martin J. Wainwright. *High-Dimensional Statistics: A Non-Asymptotic Viewpoint*. Cambridge University Press, 2019.
- Iyad Walwil and Olivier Fercoq. The smoothed duality gap as a stopping criterion. *Mathematical Programming Computation*, 2025.
- Yuanhao Wang and Jian Li. Improved Algorithms for Convex-Concave Minimax Optimization. In *NeurIPS*, 2020.
- Yunke Wang, Chang Xu, Bo Du, and Honglak Lee. Learning to Weight Imperfect Demonstrations. In *ICML*, 2021.
- Junfeng Wen, Chun-Nam Yu, and Russell Greiner. Robust Learning under Uncertain Test Distributions: Relating Covariate Shift to Model Misspecification. In *ICML*, 2014.
- Robert Williamson and Aditya Menon. Fairness Risk Measures. In *ICML*, 2019.
- Blake Woodworth and Nathan Srebro. Tight complexity bounds for optimizing composite objectives. In *NeurIPS*, 2016.
- Xiaojing Xiang. A note on the bias of L -estimators and a bias reduction procedure. *Statistics & Probability Letters*, 1995.

- J. Xiao, J. Hays, K. A. Ehinger, A. Oliva, and A. Torralba. SUN database: Large-scale scene recognition from abbey to zoo. In *CVPR*, 2010.
- Lin Xiao. Dual Averaging Method for Regularized Stochastic Learning and Online Optimization. In *NeurIPS*, 2009.
- Da Xu, Yuting Ye, and Chuanwei Ruan. Understanding the role of importance weighting for deep learning. In *ICLR*, 2021.
- Hu Xu, Saining Xie, Xiaoqing Tan, Po-Yao Huang, Russell Howes, Vasu Sharma, Shang-Wen Li, Gargi Ghosh, Luke Zettlemoyer, and Christoph Feichtenhofer. Demystifying CLIP data. In *ICLR*, 2024.
- Junchi Yang, Siqi Zhang, Negar Kiyavash, and Niao He. A Catalyst Framework for Minimax Optimization. In *NeurIPS*, 2020.
- Sitan Yang, Malcolm Wolff, Shankar Ramasubramanian, Vincent Quenneville-Belair, Ronak Mehta, and Michael Mahoney. GEANN: Scalable Graph Augmentations for Multi-Horizon Time Series Forecasting. In *KDD 2023 Workshop on Mining and Learning with Graphs*, 2023.
- Ivy Yeh. Analysis of Strength of Concrete Using Design of Experiments and Neural Networks. *Journal of Materials in Civil Engineering*, 18, 2006.
- Jure Zbontar, Li Jing, Ishan Misra, Yann LeCun, and Stephane Deny. Barlow Twins: Self-Supervised Learning via Redundancy Reduction. In *ICML*, 2021.
- Runtian Zhai, Chen Dan, Arun Suggala, J. Zico Kolter, and Pradeep Ravikumar. Boosted CVaR Classification. In *NeurIPS*, 2021.
- Huiming Zhang and Haoyu Wei. Sharper Sub-Weibull Concentrations. *Mathematics*, 2022.
- Junyu Zhang, Minyi Hong, and Shuzhong Zhang. On lower iteration complexity bounds for the convex concave saddle point problems. *Mathematical Programming*, 2022.

Kun Zhang, Bernhard Schölkopf, Krikamol Muandet, and Zhikun Wang. Domain Adaptation under Target and Conditional Shift. In *ICML*, 2013.

Tianyi Zhang, Ikko Yamane, Nan Lu, and Masashi Sugiyama. A One-step Approach to Covariate Shift Adaptation. In *ACML*, 2020.

Tong Zhang. *Mathematical Analysis of Machine Learning Algorithms*. Cambridge University Press, 2023.

Yu-Jie Zhang, Zhen-Yu Zhang, Peng Zhao, and Masashi Sugiyama. Adapting to Continuous Covariate Shift via Online Density Ratio Estimation. In *NeurIPS*, 2023.

Appendix A

APPENDIX TO CHAPTER 2

A.1 Technical Background

In this section, we collect several results from convex analysis used throughout the thesis. In the following, let $\|\cdot\|$ denote an arbitrary norm on \mathbb{R}^d and let $\|\cdot\|_*$ denote its associated dual norm.

A.1.1 Smooth and Strongly Convex Functions

The first set of results concerns L -smooth functions, or those with L -Lipschitz continuous gradient.

Theorem A.1.1. [*Nesterov, 2018, Theorem 2.1.5*] *The conditions below are considered for any $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$ and $\alpha \in [0, 1]$. The following are equivalent for a differentiable function $f : \mathbb{R}^d \rightarrow \mathbb{R}$.*

1. f is convex and L -smooth with respect to $\|\cdot\|$.
2. $0 \leq f(\mathbf{y}) - f(\mathbf{x}) - \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle \leq \frac{L}{2} \|\mathbf{x} - \mathbf{y}\|^2$.
3. $f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \frac{1}{2L} \|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|_*^2 \leq f(\mathbf{y})$.
4. $\frac{1}{L} \|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|_*^2 \leq \langle \nabla f(\mathbf{x}) - \nabla f(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle$.
5. $0 \leq \langle \nabla f(\mathbf{x}) - \nabla f(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle \leq L \|\mathbf{x} - \mathbf{y}\|^2$.

Next, we detail the properties of strongly convex functions.

Theorem A.1.2. [[Nesterov, 2018](#), Theorem 2.1.10] If $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is μ -strongly convex and differentiable, then for any $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$,

- $f(\mathbf{y}) \leq f(\mathbf{x}) + \langle \nabla f(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \frac{1}{2\mu} \|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|_*^2.$
- $\langle \nabla f(\mathbf{x}) - \nabla f(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle \leq \frac{1}{\mu} \|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|_*^2.$
- $\mu \|\mathbf{x} - \mathbf{y}\| \leq \|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|_*.$

Finally, functions that are both smooth and strongly convex enjoy a number of relevant primal-dual properties.

Theorem A.1.3. [[Nesterov, 2018](#), Theorem 2.1.12] If f is both L -smooth and μ -strongly convex, then for any $\mathbf{x}, \mathbf{y} \in \mathbb{R}^d$,

$$-\langle \nabla f(\mathbf{x}), \mathbf{x} - \mathbf{y} \rangle = -\frac{\mu L}{\mu + L} \|\mathbf{x} - \mathbf{y}\|^2 - \frac{1}{\mu + L} \|\nabla f(\mathbf{x}) - \nabla f(\mathbf{y})\|^2 - \langle \nabla f(\mathbf{y}), \mathbf{x} - \mathbf{y} \rangle. \quad (\text{A.1})$$

Lemma A.1.1. Let $f : \mathbb{R}^d \rightarrow \mathbb{R}$ be μ -strongly convex and M -smooth. Then, we have for any $\mathbf{w}, \mathbf{v} \in \mathbb{R}^d$,

$$f(\mathbf{v}) \geq f(\mathbf{w}) + \nabla f(\mathbf{w})^\top (\mathbf{v} - \mathbf{w}) + \frac{1}{2(M + \mu)} \|\nabla f(\mathbf{w}) - \nabla f(\mathbf{v})\|_2^2 + \frac{\mu}{4} \|\mathbf{w} - \mathbf{v}\|_2^2.$$

Proof. The function $g = f - \mu \|\cdot\|_2^2/2$ is convex and $M - \mu$ smooth. Hence, we have by line 3 of Theorem [A.1.1](#) for any $\mathbf{w}, \mathbf{v} \in \mathbb{R}^d$,

$$g(\mathbf{v}) \geq g(\mathbf{w}) + \nabla g(\mathbf{w})^\top (\mathbf{v} - \mathbf{w}) + \frac{1}{2(M - \mu)} \|\nabla g(\mathbf{v}) - \nabla g(\mathbf{w})\|_2^2.$$

Expanding g and ∇g , we get

$$\begin{aligned} f(\mathbf{v}) &\geq f(\mathbf{w}) + \nabla f(\mathbf{w})^\top (\mathbf{v} - \mathbf{w}) + \frac{1}{2(M - \mu)} \|\nabla f(\mathbf{w}) - \nabla f(\mathbf{v})\|_2^2 \\ &\quad + \frac{\mu M}{2(M - \mu)} \|\mathbf{w} - \mathbf{v}\|_2^2 - \frac{\mu}{M - \mu} (\nabla f(\mathbf{w}) - \nabla f(\mathbf{v}))^\top (\mathbf{w} - \mathbf{v}). \end{aligned}$$

Using Young's inequality, that is, $\mathbf{a}^\top \mathbf{b} \leq \frac{\alpha}{2} \|\mathbf{a}\|_2^2 + \frac{\alpha^{-1}}{2} \|\mathbf{b}\|_2^2$, we have

$$\begin{aligned} f(\mathbf{v}) &\geq f(\mathbf{w}) + \nabla f(\mathbf{w})^\top (\mathbf{v} - \mathbf{w}) + \frac{1 - \alpha\mu}{2(M - \mu)} \|\nabla f(\mathbf{w}) - \nabla f(\mathbf{v})\|_2^2 \\ &\quad + \frac{\mu(M - \alpha^{-1})}{2(M - \mu)} \|\mathbf{w} - \mathbf{v}\|_2^2. \end{aligned}$$

Taking $\alpha = \frac{2}{\mu + M}$ gives the claim. \square

A.1.2 f -Divergences

Let Q and P be two probability measures over \mathcal{Z} . Consider a convex function $f : [0, \infty) \rightarrow \mathbb{R} \cup \{+\infty\}$ such that $f(1) = 0$, $f(\mathbf{x})$ is finite for $x > 0$, and $\lim_{t \rightarrow 0^+} f(t) = 0$. The f -divergence from Q to P [Zhang, 2023, Appendix B] generated by this function f is

$$D_f(Q \| P) := \begin{cases} \int_{\mathcal{Z}} f\left(\frac{dQ}{dP}(z)\right) dP(z) & \text{if } Q \ll P \\ +\infty & \text{otherwise} \end{cases}.$$

In the special case that \mathbf{q} and \mathbf{p} are two probability mass functions defined on atoms $\{1, \dots, n\}$, we may use the abuse of notation

$$D_f(\mathbf{q} \| \mathbf{p}) := \sum_{i=1}^n f\left(\frac{q_i}{p_i}\right) p_i,$$

where we define $0f(0/0) := 0$ in the formula above. If there is an i such that $p_i = 0$ but $q_i > 0$, we say $D_f(\mathbf{q} \| \mathbf{p}) = \infty$. The χ^2 -divergence is generated by $f_{\chi^2}(t) = t^2 - 1$ and the KL divergence is generated by $f_{\text{KL}}(t) = t \ln t + \iota_+(t)$ where ι_+ denotes the convex indicator that is zero for $t \geq 0$ and $+\infty$ otherwise, and we define $t \ln t = 0$ for all $t < 0$.

The convexity properties of the f -divergence in its first argument can be derived from similar properties of the function f on \mathbb{R} .

Proposition A.1.1. *Assume that $f : \mathbb{R} \rightarrow \mathbb{R}$ is α_n -strongly convex on $[0, n]$. Then, $\mathbf{q} \mapsto D_f(\mathbf{q} \| \mathbf{1}_n/n)$ is $(n\alpha_n)$ -strongly convex with respect to $\|\cdot\|_2$.*

Proof. Due to the α_n -strong convexity of f , for any $\mathbf{q}, \mathbf{p} \in [0, 1]^n$ and any $\lambda \in (0, 1)$ and any

$i \in [n]$,

$$f(\lambda nq_i + (1 - \lambda)np_i) \leq \lambda f(nq_i) + (1 - \lambda)f(np_i) - \frac{\alpha_n}{2}\lambda(1 - \lambda)(nq_i - np_i)^2.$$

We average this inequality over i , yielding

$$\frac{1}{n} \sum_{i=1}^n f(n(\lambda q_i + (1 - \lambda)p_i)) \leq \frac{\lambda}{n} \sum_{i=1}^n f(nq_i) + \frac{1 - \lambda}{n} \sum_{i=1}^n f(np_i) - \frac{\alpha_n}{2}\lambda(1 - \lambda)\|nq_i - np_i\|_2^2.$$

Defining $\text{Reg}(\mathbf{q}) := D_f(\mathbf{q} \parallel \mathbf{1}_n/n)$, the statement above can be succinctly written as

$$\text{Reg}(\lambda \mathbf{q} + (1 - \lambda)\mathbf{p}) \leq \lambda \text{Reg}(\mathbf{q}) + (1 - \lambda) \text{Reg}(\mathbf{p}) - \frac{\alpha_n n}{2}\lambda(1 - \lambda)\|\mathbf{q} - \mathbf{p}\|_2^2.$$

Therefore, Reg is $(\alpha_n n)$ -strongly convex with respect to $\|\cdot\|_2$ on $[0, 1]^n$, completing the proof. \square

A.1.3 Miscellaneous Results

Lemma A.1.2. *For a convex function $f : \mathbb{R} \rightarrow \mathbb{R} \cup \{+\infty\}$, if $x_1 \geq x_2$ and $y_2 \geq y_1$, then*

$$f(y_1 - x_1) + f(y_2 - x_2) \geq f(y_2 - x_1) + f(y_1 - x_2).$$

Proof. First, observe that

$$y_2 - x_2 \geq y_2 - x_1 \geq y_1 - x_1 \text{ and } y_2 - x_2 \geq y_1 - x_2 \geq y_1 - x_1.$$

Thus, $y_2 - x_1$ and $y_1 - x_2$ both lie between $y_2 - x_2$ and $y_1 - x_1$ and can be expressed as a convex combination of the two endpoints, that is

$$y_2 - x_1 = \alpha(y_2 - x_2) + (1 - \alpha)(y_1 - x_1)$$

$$y_1 - x_2 = \beta(y_2 - x_2) + (1 - \beta)(y_1 - x_1)$$

for some $\alpha, \beta \in [0, 1]$. By solving for α we get $\alpha = 1 - \beta$. Apply the definition of convexity

to get

$$\begin{aligned} f(y_2 - x_1) &\leq \alpha f(y_2 - x_2) + (1 - \alpha)f(y_1 - x_1) \\ f(y_1 - x_2) &\leq (1 - \alpha)f(y_2 - x_2) + \alpha f(y_1 - x_1). \end{aligned}$$

Sum both inequalities to achieve the desired result. \square

Lemma A.1.1. *Consider a random variable X with c.d.f. F . If X satisfies $\mathbb{E}[|X|^p] < \infty$ for $p > 2$, then*

$$\int_{-\infty}^{+\infty} \sqrt{F(x)(1 - F(x))} \, dx \leq \sqrt{2} \left(\frac{p}{p-2} \right) \mathbb{E}[|X|^p]^{\frac{1}{p}}.$$

Proof. By definition, $\int_{-\infty}^{\infty} \sqrt{F(x)(1 - F(x))} \, dx = \lim_{a \rightarrow +\infty} \int_{-a}^a \sqrt{F(x)(1 - F(x))} \, dx$. Denote $c = \mathbb{E}[|X|^p]^{1/p}$. For any constant $a \geq c > 0$, we have

$$\begin{aligned} \int_{-a}^a \sqrt{F(x)(1 - F(x))} \, dx &= \int_{-a}^0 \sqrt{F(x)(1 - F(x))} \, dx + \int_0^a \sqrt{F(x)(1 - F(x))} \, dx \\ &\leq \int_{-a}^0 \sqrt{F(x)} \, dx + \int_0^a \sqrt{(1 - F(x))} \, dx \\ &= \int_{-a}^0 \sqrt{\mathbb{P}(X \leq x)} \, dx + \int_0^a \sqrt{\mathbb{P}(X > x)} \, dx \\ &= \int_0^a \sqrt{\mathbb{P}(X \leq -x)} + \sqrt{\mathbb{P}(X > x)} \, dx. \end{aligned}$$

Then, use that for any $a, b \geq 0$,

$$(\sqrt{a} + \sqrt{b})^2 = a + b + 2\sqrt{ab} \leq 2(a + b) \implies \sqrt{a} + \sqrt{b} \leq \sqrt{2(a + b)}.$$

Using this, and that $x \geq 0$, we have

$$\begin{aligned} \sqrt{\mathbb{P}(X \leq -x)} + \sqrt{\mathbb{P}(X > x)} &\leq \sqrt{2(\mathbb{P}(X \leq -x) + \mathbb{P}(X > x))} \\ &= \sqrt{2(\mathbb{P}(|X| > x) + \mathbb{P}(X = -x))} \\ &\leq \sqrt{2(\mathbb{P}(|X| > x) + \mathbb{P}(|X| = x))} \\ &= \sqrt{2\mathbb{P}(|X| \geq x)}. \end{aligned}$$

Combining with the first display, we have that

$$\begin{aligned}
\int_{-a}^a \sqrt{F(x)(1-F(x))} \, dx &\leq \int_0^a \sqrt{\mathbb{P}(X \leq -x)} + \sqrt{\mathbb{P}(X > x)} \, dx \\
&\leq \sqrt{2} \int_0^a \sqrt{\mathbb{P}(|X| \geq x)} \, dx \\
&\leq \sqrt{2} \int_0^a \sqrt{\min \left\{ 1, \frac{c^p}{z^p} \right\}} \, dx && \text{Markov's inequality} \\
&= \sqrt{2} \left(c + c^{p/2} \int_c^a z^{-p/2} \, dz \right).
\end{aligned}$$

Computing the integral yields

$$\int_c^a z^{-p/2} \, dz = \frac{a^{1-p/2} - c^{1-p/2}}{1-p/2}.$$

Because $1 - p/2 < 0$, we have that $\lim_{a \rightarrow \infty} \int_c^a z^{-p/2} \, dz = \frac{c^{1-p/2}}{p/2-1}$. Combining the steps above, we obtain

$$\begin{aligned}
\int_{-\infty}^{\infty} \sqrt{F(x)(1-F(x))} \, dx &= \lim_{a \rightarrow \infty} \int_{-a}^a \sqrt{F(x)(1-F(x))} \, dx \\
&\leq \lim_{a \rightarrow \infty} \sqrt{2} \left(c + c^{p/2} \int_c^a z^{-p/2} \, dz \right) \\
&= \sqrt{2} c \left(1 + \frac{1}{p/2-1} \right) \\
&= \sqrt{2} \frac{pc}{p-2}.
\end{aligned}$$

Resubstituting $c = \mathbb{E}[|X|^p]^{1/p}$ completes the proof. \square

A.2 Convergence Analysis

The results of this section accompany the analysis in Section 2.6.

A.2.1 Intermediate Results

We first prove the generalized descent lemma, which forms the backbone of the argument in both the large and small shift cost settings.

Proposition 2.6.1 (Bias Bound). *Consider any $\boldsymbol{\theta} \in \mathbb{R}^d$, $\mathbf{l} \in \mathbb{R}^n$, and $\bar{\mathbf{q}} \in \mathcal{Q}$. Define*

$$\mathbf{q} := q^{\text{opt}}(\mathbf{l}) = \arg \max_{\mathbf{p} \in \mathcal{Q}} \langle \mathbf{p}, \mathbf{l} \rangle - \nu \text{Reg}(\mathbf{p}).$$

For any $\alpha_1 \in [0, 1]$,

$$\begin{aligned} & -(\nabla r(\boldsymbol{\theta})^\top \mathbf{q} - \nabla r(\boldsymbol{\theta}^\star)^\top \bar{\mathbf{q}})^\top (\boldsymbol{\theta} - \boldsymbol{\theta}^\star) \\ & \leq -(\mathbf{q} - \bar{\mathbf{q}})^\top (\ell(\boldsymbol{\theta}) - \ell(\boldsymbol{\theta}^\star)) - \frac{\mu}{2} \|\boldsymbol{\theta} - \boldsymbol{\theta}^\star\|_2^2 \\ & - \frac{\alpha_1}{4(M + \mu)\kappa_{\mathcal{Q}}} \frac{1}{n} \sum_{i=1}^n \|nq_i \nabla r_i(\boldsymbol{\theta}) - nq_i^\star \nabla r_i(\boldsymbol{\theta}^\star)\|_2^2 + \frac{2\alpha_1 G^2}{\nu(M + \mu)\kappa_{\mathcal{Q}}} n(\mathbf{q} - \mathbf{q}^\star)^\top (\mathbf{l} - \mathbf{l}^\star). \end{aligned}$$

Proof. First, for any $q_i > 0$, we have that $w \mapsto q_i r_i(\boldsymbol{\theta})$ is $(q_i M)$ -smooth and $(q_i \mu)$ -strongly convex. Define the notation $\sigma_n = \kappa_{\mathcal{Q}}/n$. By applying standard convex inequalities (Lemma A.1.1) we have that

$$\begin{aligned} q_i r_i(\boldsymbol{\theta}^\star) & \geq q_i r_i(\boldsymbol{\theta}) + q_i \nabla r_i(\boldsymbol{\theta})^\top (\boldsymbol{\theta}^\star - \boldsymbol{\theta}) \\ & + \frac{1}{2q_i(M + \mu)} \|q_i \nabla r_i(\boldsymbol{\theta}) - q_i \nabla r_i(\boldsymbol{\theta}^\star)\|_2^2 + \frac{q_i \mu}{4} \|\boldsymbol{\theta} - \boldsymbol{\theta}^\star\|_2^2 \\ & \geq q_i r_i(\boldsymbol{\theta}) + q_i \nabla r_i(\boldsymbol{\theta})^\top (\boldsymbol{\theta}^\star - \boldsymbol{\theta}) \\ & + \frac{1}{2\sigma_n(M + \mu)} \|q_i \nabla r_i(\boldsymbol{\theta}) - q_i \nabla r_i(\boldsymbol{\theta}^\star)\|_2^2 + \frac{q_i \mu}{4} \|\boldsymbol{\theta} - \boldsymbol{\theta}^\star\|_2^2 \end{aligned}$$

as $q_i \leq \sigma_n$. The second inequality holds for $q_i = 0$ as well, so by summing the inequality over i and using that $\sum_i q_i = 1$, we have that

$$\begin{aligned} \mathbf{q}^\top r(\boldsymbol{\theta}^\star) & \geq \mathbf{q}^\top r(\boldsymbol{\theta}) + \mathbf{q}^\top \nabla r(\boldsymbol{\theta})(\boldsymbol{\theta}^\star - \boldsymbol{\theta}) \\ & + \frac{1}{2\sigma_n(M + \mu)} \sum_{i=1}^n \|q_i \nabla r_i(\boldsymbol{\theta}) - q_i \nabla r_i(\boldsymbol{\theta}^\star)\|_2^2 + \frac{\mu}{4} \|\boldsymbol{\theta} - \boldsymbol{\theta}^\star\|_2^2. \end{aligned}$$

Applying the same argument replacing \mathbf{q} by $\bar{\mathbf{q}}$ and swapping w and $\boldsymbol{\theta}^\star$ yields

$$\begin{aligned} \bar{\mathbf{q}}^\top r(\boldsymbol{\theta}) & \geq \bar{\mathbf{q}}^\top r(\boldsymbol{\theta}^\star) + \bar{\mathbf{q}}^\top \nabla r(\boldsymbol{\theta}^\star)(\boldsymbol{\theta} - \boldsymbol{\theta}^\star) \\ & + \frac{1}{2\sigma_n(M + \mu)} \sum_{i=1}^n \|\bar{q}_i \nabla r_i(\boldsymbol{\theta}) - \bar{q}_i \nabla r_i(\boldsymbol{\theta}^\star)\|_2^2 + \frac{\mu}{4} \|\boldsymbol{\theta} - \boldsymbol{\theta}^\star\|_2^2. \end{aligned}$$

Summing the two inequalities yields

$$\begin{aligned}
& -(\mathbf{q} - \bar{\mathbf{q}})^\top (r(\boldsymbol{\theta}) - r(\boldsymbol{\theta}^*)) \\
& \geq -(\nabla r(\boldsymbol{\theta})^\top \mathbf{q} - \nabla r(\boldsymbol{\theta}^*)^\top \bar{\mathbf{q}})^\top (\boldsymbol{\theta} - \boldsymbol{\theta}^*) + \frac{\mu}{2} \|\boldsymbol{\theta} - \boldsymbol{\theta}^*\|_2^2 \\
& + \frac{1}{2\sigma_n(M + \mu)} \left[\sum_{i=1}^n \|q_i \nabla r_i(\boldsymbol{\theta}) - q_i \nabla r_i(\boldsymbol{\theta}^*)\|_2^2 + \sum_{i=1}^n \|\bar{q}_i \nabla r_i(\boldsymbol{\theta}) - \bar{q}_i \nabla r_i(\boldsymbol{\theta}^*)\|_2^2 \right].
\end{aligned}$$

Dropping the $\sum_{i=1}^n \|\bar{q}_i \nabla r_i(\boldsymbol{\theta}) - \bar{q}_i \nabla r_i(\boldsymbol{\theta}^*)\|_2^2$ term and applying a weight of $\alpha_1 \in [0, 1]$ to $\sum_{i=1}^n \|q_i \nabla r_i(\boldsymbol{\theta}) - q_i \nabla r_i(\boldsymbol{\theta}^*)\|_2^2$ still satisfies the inequality, which can equivalently be written as

$$\begin{aligned}
& -(\nabla r(\boldsymbol{\theta})^\top \mathbf{q} - \nabla r(\boldsymbol{\theta}^*)^\top \bar{\mathbf{q}})^\top (\boldsymbol{\theta} - \boldsymbol{\theta}^*) \leq -(\mathbf{q} - \bar{\mathbf{q}})^\top (r(\boldsymbol{\theta}) - r(\boldsymbol{\theta}^*)) - \frac{\mu}{2} \|\boldsymbol{\theta} - \boldsymbol{\theta}^*\|_2^2 \\
& - \frac{\alpha_1}{2\sigma_n(M + \mu)} \sum_{i=1}^n \|q_i \nabla r_i(\boldsymbol{\theta}) - q_i \nabla r_i(\boldsymbol{\theta}^*)\|_2^2. \quad (\text{A.2})
\end{aligned}$$

Next, because

$$\|q_i \nabla r_i(\boldsymbol{\theta}) - q_i^* \nabla r_i(\boldsymbol{\theta}^*)\|_2^2 \leq 2 \|q_i \nabla r_i(\boldsymbol{\theta}) - q_i \nabla r_i(\boldsymbol{\theta}^*)\|_2^2 + 2(q_i - q_i^*)^2 \|\nabla r_i(\boldsymbol{\theta}^*)\|_2^2,$$

we have that (by summing over i) that

$$-\sum_{i=1}^n \|q_i \nabla r_i(\boldsymbol{\theta}) - q_i \nabla r_i(\boldsymbol{\theta}^*)\|_2^2 \leq -\frac{1}{2} \sum_{i=1}^n \|q_i \nabla r_i(\boldsymbol{\theta}) - q_i^* \nabla r_i(\boldsymbol{\theta}^*)\|_2^2 + 4G^2 \|\mathbf{q} - \mathbf{q}^*\|_2^2, \quad (\text{A.3})$$

where we used that each $\|\nabla r_i(\boldsymbol{\theta}^*)\|_2 \leq 2G$. To see this, use that $\nabla r(\boldsymbol{\theta}^*)^\top \mathbf{q}^* = 0$ and $\nabla r(\boldsymbol{\theta}^*) = \nabla \ell(\boldsymbol{\theta}^*) + \mu \boldsymbol{\theta}^*$, so

$$\|\nabla r_i(\boldsymbol{\theta}^*)\|_2 = \|\nabla \ell_i(\boldsymbol{\theta}^*) + \mu \boldsymbol{\theta}^*\|_2 = \left\| \nabla \ell_i(\boldsymbol{\theta}^*) - \sum_{j=1}^n q_i^* \nabla \ell_j(\boldsymbol{\theta}^*) \right\|_2 \leq 2G.$$

Because the map q^{opt} is the gradient of a convex and $(1/\nu)$ -smooth map, we also have that

$$\|\mathbf{q} - \mathbf{q}^*\|_2^2 = \|q^{\text{opt}}(\mathbf{l}) - q^{\text{opt}}(\ell(\boldsymbol{\theta}^*))\|_2^2 \leq \frac{1}{\nu} (\mathbf{q} - \mathbf{q}^*)^\top (\mathbf{l} - \ell(\boldsymbol{\theta}^*)), \quad (\text{A.4})$$

so we apply the above to (A.3) to achieve

$$\begin{aligned} & - \sum_{i=1}^n \|q_i \nabla r_i(\boldsymbol{\theta}) - q_i \nabla r_i(\boldsymbol{\theta}^*)\|_2^2 \\ & \leq -\frac{1}{2} \sum_{i=1}^n \|q_i \nabla r_i(\boldsymbol{\theta}) - q_i^* \nabla r_i(\boldsymbol{\theta}^*)\|_2^2 + \frac{4G^2}{\nu} (\mathbf{q} - \mathbf{q}^*)^\top (\mathbf{l} - \ell(\boldsymbol{\theta}^*)), \end{aligned} \quad (\text{A.5})$$

We also use (A.4) to claim non-negativity of $(\mathbf{q} - \mathbf{q}^*)^\top (\mathbf{l} - \ell(\boldsymbol{\theta}^*))$. Finally, because $\sum_i q_i = \sum_i q_i^* = 1$, we have that

$$\begin{aligned} (\mathbf{q} - \bar{\mathbf{q}})^\top (r(\boldsymbol{\theta}) - r(\boldsymbol{\theta}^*)) &= (\mathbf{q} - \bar{\mathbf{q}})^\top \left(\ell(\boldsymbol{\theta}) + \frac{\mu}{2} \|\boldsymbol{\theta}\|_2^2 \mathbf{1} - \ell(\boldsymbol{\theta}^*) - \frac{\mu}{2} \|\boldsymbol{\theta}^*\|_2^2 \mathbf{1} \right) \\ &= (\mathbf{q} - \bar{\mathbf{q}})^\top (\ell(\boldsymbol{\theta}) - \ell(\boldsymbol{\theta}^*)) + (\mathbf{q} - \bar{\mathbf{q}})^\top \mathbf{1} (\|\boldsymbol{\theta}\|_2^2 - \|\boldsymbol{\theta}^*\|_2^2) \\ &= (\mathbf{q} - \bar{\mathbf{q}})^\top (\ell(\boldsymbol{\theta}) - \ell(\boldsymbol{\theta}^*)). \end{aligned} \quad (\text{A.6})$$

Combine (A.2), (A.5), and (A.6) along with $\kappa_Q = n\sigma_n$ to achieve the claim. \square

The upcoming results provide the upper bounds for the Lyapunov function terms introduced in (2.26).

Lemma 2.6.3. *For any value of $\alpha_2 > 0$, we have that*

$$\begin{aligned} \mathbb{E}_k [U^{(k+1)}] &\leq \eta^2(1 + \alpha_2) \textcolor{violet}{Q}^{(k)} + \eta^2(1 + \alpha_2^{-1}) \textcolor{green}{S}^{(k)} \\ &\quad + \frac{\eta M^2}{\mu n} \left(1 - \frac{1}{n}\right) \textcolor{red}{T}^{(k)} + \left(1 - \frac{1}{n}\right) \frac{G^2}{2\nu\mu n} \textcolor{brown}{R}^{(k)} + \left(1 - \frac{1}{n}\right) U^{(k)}. \end{aligned}$$

Proof. First, we apply Condition 2.5.1 to the functions

$$h_i(\mathbf{u}, \mathbf{x}) = \frac{1}{n} \|\mathbf{u} - \mathbf{x}\|_2^2$$

to achieve the equality

$$\begin{aligned} \mathbb{E}_k [U^{(k+1)}] &= \frac{1}{n} \mathbb{E}_k \left[\|\boldsymbol{\theta}^{(k+1)} - \boldsymbol{\theta}^{(k)}\|_2^2 \right] + \left(1 - \frac{1}{n}\right) \mathbb{E}_k \left[\frac{1}{n} \sum_{j=1}^n \|\boldsymbol{\theta}^{(k+1)} - \hat{\boldsymbol{\theta}}_j^{(k)}\|_2^2 \right] \\ &= \frac{\eta^2}{n} \mathbb{E}_k \left[\|\mathbf{v}^{(k)}\|_2^2 \right] + \left(1 - \frac{1}{n}\right) \mathbb{E}_k \left[\frac{1}{n} \sum_{j=1}^n \|\boldsymbol{\theta}^{(k+1)} - \hat{\boldsymbol{\theta}}_j^{(k)}\|_2^2 \right]. \end{aligned}$$

Next, we expand the second term.

$$\begin{aligned}
& \frac{1}{n} \mathbb{E}_k \left[\sum_{j=1}^n \|\boldsymbol{\theta}^{(k+1)} - \hat{\boldsymbol{\theta}}_j^{(k)}\|_2^2 \right] \\
&= \frac{1}{n} \mathbb{E}_k \left[\sum_{j=1}^n \|\boldsymbol{\theta}^{(k+1)} - \boldsymbol{\theta}^{(k)}\|_2^2 \right] + \frac{2}{n} \mathbb{E}_k \left[\sum_{j=1}^n (\boldsymbol{\theta}^{(k+1)} - \boldsymbol{\theta}^{(k)})^\top (\boldsymbol{\theta}^{(k)} - \hat{\boldsymbol{\theta}}_j^{(k)}) \right] + \frac{1}{n} \mathbb{E}_k \left[\sum_{j=1}^n \|\hat{\boldsymbol{\theta}}_j^{(k)} - \boldsymbol{\theta}^{(k)}\|_2^2 \right] \\
&= \eta^2 \mathbb{E}_k [\|\mathbf{v}^{(k)}\|_2^2] - \frac{2\eta}{n} \sum_{j=1}^n \nabla(\mathbf{q}^{(k)\top} r)(\boldsymbol{\theta}^{(k)})^\top (\boldsymbol{\theta}^{(k)} - \hat{\boldsymbol{\theta}}_j^{(k)}) + \frac{1}{n} \sum_{j=1}^n \|\hat{\boldsymbol{\theta}}_j^{(k)} - \boldsymbol{\theta}^{(k)}\|_2^2.
\end{aligned}$$

The first term is simply the noise term that appears in Lemma 2.6.1, whereas the last term is $U^{(k)}$. Next, we have

$$\begin{aligned}
-2\nabla(\mathbf{q}^{(k)\top} r)(\boldsymbol{\theta}^{(k)})^\top (\boldsymbol{\theta}^{(k)} - \hat{\boldsymbol{\theta}}_j^{(k)}) &= -2(\nabla(\mathbf{q}^{(k)\top} r)(\boldsymbol{\theta}^{(k)}) - \nabla(\mathbf{q}^{(k)\top} r)(\hat{\boldsymbol{\theta}}_j^{(k)}))^\top (\boldsymbol{\theta}^{(k)} - \hat{\boldsymbol{\theta}}_j^{(k)}) \\
&\quad - 2(\nabla(\mathbf{q}^{(k)\top} r)(\hat{\boldsymbol{\theta}}_j^{(k)}) - \nabla(\mathbf{q}^{(k)\top} r)(\boldsymbol{\theta}^*))^\top (\boldsymbol{\theta}^{(k)} - \hat{\boldsymbol{\theta}}_j^{(k)}) \\
&\quad - 2(\nabla(\mathbf{q}^{(k)\top} r)(\boldsymbol{\theta}^*) - \nabla(\mathbf{q}^{*\top} r)(\boldsymbol{\theta}^*))^\top (\boldsymbol{\theta}^{(k)} - \hat{\boldsymbol{\theta}}_j^{(k)}),
\end{aligned}$$

where the last term is introduced because $\nabla(\mathbf{q}^{*\top} r)(\boldsymbol{\theta}^*) = 0$. We bound each of the three terms. First,

$$-2(\nabla(\mathbf{q}^{(k)\top} r)(\boldsymbol{\theta}^{(k)}) - \nabla(\mathbf{q}^{(k)\top} r)(\hat{\boldsymbol{\theta}}_j^{(k)}))^\top (\boldsymbol{\theta}^{(k)} - \hat{\boldsymbol{\theta}}_j^{(k)}) \leq -2\mu \left\| \boldsymbol{\theta}^{(k)} - \hat{\boldsymbol{\theta}}_j^{(k)} \right\|_2^2$$

because $\mathbf{q}^{(k)\top} r$ is μ -strongly convex [Nesterov, 2018, Theorem 2.1.9]. Second,

$$\begin{aligned}
& -2(\nabla(\mathbf{q}^{(k)\top} r)(\hat{\boldsymbol{\theta}}_j^{(k)}) - \nabla(\mathbf{q}^{(k)\top} r)(\boldsymbol{\theta}^*))^\top (\boldsymbol{\theta}^{(k)} - \hat{\boldsymbol{\theta}}_j^{(k)}) \\
& \leq \alpha_4 \|\nabla(\mathbf{q}^{(k)\top} r)(\hat{\boldsymbol{\theta}}_j^{(k)}) - \nabla(\mathbf{q}^{(k)\top} r)(\boldsymbol{\theta}^*)\|_2^2 + \alpha_4^{-1} \|\hat{\boldsymbol{\theta}}_j^{(k)} - \boldsymbol{\theta}^{(k)}\|_2^2 \\
& \leq \alpha_4 M^2 \|\hat{\boldsymbol{\theta}}_j^{(k)} - \boldsymbol{\theta}^*\|_2^2 + \alpha_4^{-1} \|\hat{\boldsymbol{\theta}}_j^{(k)} - \boldsymbol{\theta}^{(k)}\|_2^2
\end{aligned}$$

by Young's inequality with parameter α_4 and the M -Lipschitz continuity of $\nabla(\mathbf{q}^{(k)\top} r)$. Third,

$$\begin{aligned}
& -2(\nabla(\mathbf{q}^{(k)\top} r)(\boldsymbol{\theta}^*) - \nabla(\mathbf{q}^{*\top} r)(\boldsymbol{\theta}^*))^\top (\boldsymbol{\theta}^{(k)} - \hat{\boldsymbol{\theta}}_j^{(k)}) \\
& = -2(\nabla((\mathbf{q}^{(k)} - \mathbf{q}^*)^\top \ell)(\boldsymbol{\theta}^*))^\top (\boldsymbol{\theta}^{(k)} - \hat{\boldsymbol{\theta}}_j^{(k)}) \\
& \leq \alpha_5 \|\nabla((\mathbf{q}^{(k)} - \mathbf{q}^*)^\top \ell)(\boldsymbol{\theta}^*)\|_2^2 + \alpha_5^{-1} \|\hat{\boldsymbol{\theta}}_j^{(k)} - \boldsymbol{\theta}^{(k)}\|_2^2 \\
& \leq \alpha_5 G^2 \|\mathbf{q}^{(k)} - \mathbf{q}^*\|_2^2 + \alpha_5^{-1} \|\hat{\boldsymbol{\theta}}_j^{(k)} - \boldsymbol{\theta}^{(k)}\|_2^2,
\end{aligned}$$

by Young's inequality with parameter α_5 and the G -Lipschitz continuity of each ℓ_i . Combining with the above, we have

$$\begin{aligned}
-2 \sum_{j=1}^n \nabla(\mathbf{q}^{(k)\top} r)(\boldsymbol{\theta}^{(k)})^\top (\boldsymbol{\theta}^{(k)} - \hat{\boldsymbol{\theta}}_j^{(k)}) & \leq \alpha_4 M^2 \mathbf{T}^{(k)} + (\alpha_4^{-1} + \alpha_5^{-1} - 2\mu) \mathbf{U}^{(k)} + \alpha_5 G^2 n \|\mathbf{q}^{(k)} - \mathbf{q}^*\|_2^2 \\
& \leq \mu^{-1} M^2 \mathbf{T}^{(k)} + \mu^{-1} G^2 n \|\mathbf{q}^{(k)} - \mathbf{q}^*\|_2^2
\end{aligned}$$

when we set $\alpha_4 = \alpha_5 = \mu^{-1}$. Hence, we get

$$\begin{aligned}
& \mathbb{E}_k [U^{(k+1)}] \\
& = \frac{\eta^2}{n} \mathbb{E}_k [\|\mathbf{v}^{(k)}\|_2^2] + \left(1 - \frac{1}{n}\right) \mathbb{E}_k \left[\frac{1}{n} \sum_{j=1}^n \|\boldsymbol{\theta}^{(k+1)} - \hat{\boldsymbol{\theta}}_j^{(k)}\|_2^2 \right] \\
& \leq \eta^2 \mathbb{E}_k [\|\mathbf{v}^{(k)}\|_2^2] - \frac{\eta}{n} \left(1 - \frac{1}{n}\right) 2 \sum_{j=1}^n \nabla(\mathbf{q}^{(k)\top} r)(\boldsymbol{\theta}^{(k)})^\top (\boldsymbol{\theta}^{(k)} - \hat{\boldsymbol{\theta}}_j^{(k)}) + \left(1 - \frac{1}{n}\right) \mathbf{U}^{(k)} \\
& \leq \eta^2 \mathbb{E}_k [\|\mathbf{v}^{(k)}\|_2^2] + \frac{\eta}{n} \left(1 - \frac{1}{n}\right) [\mu^{-1} M^2 \mathbf{T}^{(k)} + \mu^{-1} G^2 n \|\mathbf{q}^{(k)} - \mathbf{q}^*\|_2^2] + \left(1 - \frac{1}{n}\right) \mathbf{U}^{(k)} \\
& = \eta^2 \mathbb{E}_k [\|\mathbf{v}^{(k)}\|_2^2] + \left(1 - \frac{1}{n}\right) \frac{\eta M^2}{\mu n} \mathbf{T}^{(k)} + \left(1 - \frac{1}{n}\right) \frac{G^2}{2n\mu} 2n\eta \|\mathbf{q}^{(k)} - \mathbf{q}^*\|_2^2 + \left(1 - \frac{1}{n}\right) \mathbf{U}^{(k)} \\
& = \eta^2 \mathbb{E}_k [\|\mathbf{v}^{(k)}\|_2^2] + \left(1 - \frac{1}{n}\right) \frac{\eta M^2}{\mu n} \mathbf{T}^{(k)} + \left(1 - \frac{1}{n}\right) \frac{G^2}{2n\mu\nu} \mathbf{R}^{(k)} + \left(1 - \frac{1}{n}\right) \mathbf{U}^{(k)} \\
& \leq \eta^2 (1 + \alpha_2) \mathbf{Q}^{(k)} + \eta^2 (1 + \alpha_2^{-1}) \mathbf{S}^{(k)} \\
& + \frac{\eta M^2}{\mu n} \left(1 - \frac{1}{n}\right) \mathbf{T}^{(k)} + \left(1 - \frac{1}{n}\right) \frac{G^2}{2\nu\mu n} \mathbf{R}^{(k)} + \left(1 - \frac{1}{n}\right) \mathbf{U}^{(k)},
\end{aligned}$$

where the two last steps follows Lemma 2.4.1 and Theorem A.1.1 to claim $\|\mathbf{q}^{(k)} - \mathbf{q}^*\|_2^2 \leq$

$$\frac{1}{\bar{\nu}}(\mathbf{q}^{(k)} - \mathbf{q}^*)(\mathbf{l}^{(k)} - \mathbf{l}^*).$$

□

Lemma 2.6.4. *For any $\alpha_3 > 0$, it holds that*

$$\begin{aligned} \mathbb{E}_k [\textcolor{brown}{R}^{(k+1)}] &\leq 2\eta(\mathbf{q}^{(k)} - \mathbf{q}^*)^\top (\ell(\boldsymbol{\theta}^{(k)}) - \mathbf{l}^*) + \left(1 - \frac{1}{n}\right) \textcolor{brown}{R}^{(k)} \\ &\quad + \frac{\eta G^2 n}{2\nu} \alpha_3^{-1} \textcolor{red}{T}^{(k)} + \frac{2\eta G^2 n}{\nu} (1 + \alpha_3) \textcolor{blue}{U}^{(k)}. \end{aligned}$$

Proof. First, decompose

$$(\mathbf{q}^{(k+1)} - \mathbf{q}^*)^\top (\mathbf{l}^{(k+1)} - \mathbf{l}^*) = (\mathbf{q}^{(k)} - \mathbf{q}^*)^\top (\mathbf{l}^{(k+1)} - \mathbf{l}^*) + (\mathbf{q}^{(k+1)} - \mathbf{q}^{(k)})^\top (\mathbf{l}^{(k+1)} - \mathbf{l}^{(k)}) \quad (\text{A.7})$$

$$+ (\mathbf{q}^{(k+1)} - \mathbf{q}^{(k)})^\top (\mathbf{l}^{(k)} - \mathbf{l}^*). \quad (\text{A.8})$$

Because $\mathbf{q}^{(k)} = \mathbf{q}^{\text{opt}}(\mathbf{l}^{(k)})$ for all t , and $\mathbf{q}^{\text{opt}}(\cdot)$ is the gradient of a convex and $(1/\nu)$ -smooth function, we have for the second term of (A.8) that

$$(\mathbf{q}^{(k+1)} - \mathbf{q}^{(k)})^\top (\mathbf{l}^{(k+1)} - \mathbf{l}^{(k)}) \leq \frac{1}{\bar{\nu}} \|\mathbf{l}^{(k+1)} - \mathbf{l}^{(k)}\|_2^2.$$

Next, using Young's inequality, that is, $a^\top b \leq \frac{\alpha_3}{2} \|a\|_2^2 + \frac{\alpha_3^{-1}}{2} \|b\|_2^2$ for any $\alpha_3 > 0$, we have for the third term of (A.8) that

$$\begin{aligned} (\mathbf{q}^{(k+1)} - \mathbf{q}^{(k)})^\top (\mathbf{l}^{(k)} - \mathbf{l}^*) &\leq \frac{\alpha_3}{2} \|\mathbf{q}^{(k+1)} - \mathbf{q}^{(k)}\|_2^2 + \frac{\alpha_3^{-1}}{2} \|\mathbf{l}^{(k)} - \mathbf{l}^*\|_2^2 \\ &\leq \frac{\alpha_3}{2\bar{\nu}^2} \|\mathbf{l}^{(k+1)} - \mathbf{l}^{(k)}\|_2^2 + \frac{\alpha_3^{-1}}{2} \|\mathbf{l}^{(k)} - \mathbf{l}^*\|_2^2. \end{aligned}$$

Note that we have

$$\mathbb{E}_k [\mathbf{l}^{(k+1)}] = \frac{1}{n} \ell(\boldsymbol{\theta}^{(k)}) + \left(1 - \frac{1}{n}\right) \mathbf{l}^{(k)}.$$

Hence, we get,

$$\begin{aligned}
\frac{1}{2\eta n} \mathbb{E}_k [R^{(k+1)}] &= \frac{1}{n} (\mathbf{q}^{(k)} - \mathbf{q}^*)^\top (\ell(\boldsymbol{\theta}^{(k)}) - \mathbf{l}^*) + \left(1 - \frac{1}{n}\right) (\mathbf{q}^{(k)} - \mathbf{q}^*)^\top (\mathbf{l}^{(k)} - \mathbf{l}^*) \\
&\quad + \mathbb{E}_k [(\mathbf{q}^{(k+1)} - \mathbf{q}^{(k)})^\top (\mathbf{l}^{(k+1)} - \mathbf{l}^{(k)})] + \mathbb{E}_k [(\mathbf{q}^{(k+1)} - \mathbf{q}^{(k)})^\top (\mathbf{l}^{(k)} - \mathbf{l}^*)] \\
&\leq \frac{1}{n} (\mathbf{q}^{(k)} - \mathbf{q}^*)^\top (\ell(\boldsymbol{\theta}^{(k)}) - \mathbf{l}^*) + \left(1 - \frac{1}{n}\right) (\mathbf{q}^{(k)} - \mathbf{q}^*)^\top (\mathbf{l}^{(k)} - \mathbf{l}^*) \\
&\quad + \left(\frac{1}{\bar{\nu}} + \frac{\alpha_3}{2\bar{\nu}^2}\right) \mathbb{E}_k [\|\mathbf{l}^{(k+1)} - \mathbf{l}^{(k)}\|_2^2] + \frac{\alpha_3^{-1}}{2} \|\mathbf{l}^{(k)} - \mathbf{l}^*\|_2^2 \\
&= \frac{1}{n} (\mathbf{q}^{(k)} - \mathbf{q}^*)^\top (\ell(\boldsymbol{\theta}^{(k)}) - \mathbf{l}^*) + \left(1 - \frac{1}{n}\right) (\mathbf{q}^{(k)} - \mathbf{q}^*)^\top (\mathbf{l}^{(k)} - \mathbf{l}^*) \\
&\quad + \frac{1}{n\bar{\nu}} \left(1 + \frac{\alpha_3}{2\bar{\nu}}\right) \sum_{j=1}^n (\ell_j(\boldsymbol{\theta}^{(k)}) - \ell_j(\hat{\boldsymbol{\theta}}_j))^2 \\
&\quad + \frac{\alpha_3^{-1}}{2} \sum_{j=1}^n (\ell_j(\hat{\boldsymbol{\theta}}_j) - \ell_j(\boldsymbol{\theta}^*))^2.
\end{aligned}$$

Then, apply the G -Lipschitz continuity of each ℓ_i to achieve

$$\begin{aligned}
\frac{1}{2\eta n} \mathbb{E}_k [R^{(k+1)}] &\leq \frac{1}{n} (\mathbf{q}^{(k)} - \mathbf{q}^*)^\top (\ell(\boldsymbol{\theta}^{(k)}) - \mathbf{l}^*) + \left(1 - \frac{1}{n}\right) (\mathbf{q}^{(k)} - \mathbf{q}^*)^\top (\mathbf{l}^{(k)} - \mathbf{l}^*) \\
&\quad + \frac{G^2}{n\bar{\nu}} \left(1 + \frac{\alpha_3}{2\bar{\nu}}\right) \sum_{j=1}^n \|\boldsymbol{\theta}^{(k)} - \hat{\boldsymbol{\theta}}_j\|_2^2 \\
&\quad + \frac{G^2 \alpha_3^{-1}}{2} \sum_{j=1}^n \|\hat{\boldsymbol{\theta}}_j - \boldsymbol{\theta}^*\|_2^2.
\end{aligned}$$

Replacing α_3 by $2\bar{\nu}\alpha_3$ gives the claim. \square

A.2.2 Proof of Main Results

The forthcoming theorems provide the complete convergence analyses for the settings in which 1) the shift cost ν satisfies a particular lower bound and 2) when that bound is violated, respectively.

Theorem 2.5.2. *Assume that $n \geq 2$ and that the shift cost $\nu \leq 8nG^2/\mu$. The sequence of*

iterates produced by Algorithm 1 with

$$\eta = \frac{1}{16n\mu} \min \left\{ \frac{1}{6[8\delta + (\kappa + 1)\kappa_{\mathcal{P}}]}, \frac{1}{4\delta^2 \max \{2n\kappa^2, \delta\}} \right\}$$

achieves

$$\mathbb{E}_0 \|\boldsymbol{\theta}^{(t)} - \boldsymbol{\theta}^*\|_2^2 \leq \left(5 + 16\delta + \frac{6\kappa^2}{\sigma_n} \right) \exp(-t/\tau) \|\boldsymbol{\theta}^{(0)} - \boldsymbol{\theta}^*\|_2^2$$

for

$$\tau = 32n \max \{ 6[8\delta + (\kappa + 1)\kappa_{\mathcal{P}}], 4\delta^2 \max \{ 2n\kappa^2, \delta \}, 1/16 \}.$$

Proof. First, invoke Lemma 2.6.1 with $\mathbf{q}' = \mathbf{q}^{(k)}$ and $\alpha_1 = 1$ to obtain

$$\mathbb{E}_k \|\boldsymbol{\theta}^{(k+1)} - \boldsymbol{\theta}^*\|_2^2 \leq (1 - \eta\mu) \|\boldsymbol{\theta}^{(k)} - \boldsymbol{\theta}^*\|_2^2 \quad (\text{A.9})$$

$$- 2\eta(\boldsymbol{\theta}^{(k)} - \boldsymbol{\theta}^*)^\top \nabla r(\boldsymbol{\theta}^*) \mathbf{q}^{(k)} + \frac{2G^2}{\bar{\nu}(M + \mu)\kappa_\sigma} \mathbf{R}^{(k)} \quad (\text{A.10})$$

$$- \eta \left(\frac{1}{2(M + \mu)\kappa_\sigma} - \eta(1 + \alpha_2) \right) \mathbf{Q}^{(k)} + \eta^2(1 + \alpha_2^{-1}) \mathbf{S}^{(k)}. \quad (\text{A.11})$$

We will first bound (A.10), by using that $\nabla r(\boldsymbol{\theta}^*) \mathbf{q}^* = 0$ and Young's inequality with parameter $a > 0$ to write

$$\begin{aligned} |(\boldsymbol{\theta}^{(k)} - \boldsymbol{\theta}^*)^\top \nabla r(\boldsymbol{\theta}^*) \mathbf{q}^{(k)}| &= |(\boldsymbol{\theta}^{(k)} - \boldsymbol{\theta}^*)^\top \nabla r(\boldsymbol{\theta}^*) (\mathbf{q}^{(k)} - \mathbf{q}^*)| \\ &\leq \frac{a}{2} \|\nabla r(\boldsymbol{\theta}^*)^\top (\mathbf{q}^{(k)} - \mathbf{q}^*)\|_2^2 + \frac{1}{2a} \|\boldsymbol{\theta}^{(k)} - \boldsymbol{\theta}^*\|_2^2 \\ &\leq \frac{aG^2\gamma_\star^2}{2\nu^2} \mathbf{T}^{(k)} + \frac{1}{2a} \|\boldsymbol{\theta}^{(k)} - \boldsymbol{\theta}^*\|_2^2, \end{aligned}$$

where we used in the second inequality that:

$$\begin{aligned} \|\nabla r(\boldsymbol{\theta}^*)^\top (\mathbf{q}^{(k)} - \mathbf{q}^*)\|_2^2 &= \|\nabla \ell(\boldsymbol{\theta}^*)^\top (\mathbf{q}^{(k)} - \mathbf{q}^*)\|_2^2 \leq \gamma_\star^2 \|\mathbf{q}^{(k)} - \mathbf{q}^*\|_2^2 \leq \frac{\gamma_\star^2}{\nu^2} \|\mathbf{l}^{(k)} - \mathbf{l}^*\|_2^2 \\ &\leq \frac{G^2\gamma_\star^2}{\nu^2} \sum_{i=1}^n \|\hat{\boldsymbol{\theta}}_i^{(k)} - \boldsymbol{\theta}^*\|_2^2 = \frac{G^2\gamma_\star^2}{\nu^2} \mathbf{T}^{(k)}. \end{aligned}$$

We also have by Cauchy-Schwartz and Lipschitz continuity that

$$\textcolor{brown}{R}^{(k)} = 2\eta n(\mathbf{q}^{(k)} - \mathbf{q}^\star)^\top (\mathbf{l}^{(k)} - \mathbf{l}^\star) \leq \frac{2\eta n}{\nu} \|\mathbf{l}^{(k)} - \mathbf{l}^\star\|_2^2 \leq \frac{2\eta n G^2}{\nu} \textcolor{violet}{T}^{(k)}.$$

Combining the above displays yields

$$\begin{aligned} & -2\eta(\boldsymbol{\theta}^{(k)} - \boldsymbol{\theta}^\star)^\top \nabla r(\boldsymbol{\theta}^\star) \mathbf{q}^{(k)} + \frac{2G^2}{\bar{\nu}(M + \mu)\kappa_\sigma} \textcolor{brown}{R}^{(k)} \\ & \leq \frac{\eta G^2}{\nu^2} \left[a\gamma_\star^2 + \frac{4nG^2}{(M + \mu)\kappa_\sigma} \right] \textcolor{violet}{T}^{(k)} + \frac{\eta}{a} \|\boldsymbol{\theta}^{(k)} - \boldsymbol{\theta}^\star\|_2^2. \end{aligned}$$

We take $\alpha_2 = 2$, $c_3 = c_4 = 0$, and apply Lemma 2.6.2 to achieve

$$\begin{aligned} \mathbb{E}_k[V^{(k+1)}] - (1 - \tau^{-1})V^{(k)} & \leq [\tau^{-1} - \eta\mu + \eta a^{-1} + c_2] \|\boldsymbol{\theta}^{(k)} - \boldsymbol{\theta}^\star\|_2^2 \\ & \quad + \left[\tau^{-1} + \frac{3\eta^2}{2c_1} - \frac{1}{n} \right] c_1 \textcolor{green}{S}^{(k)} \\ & \quad + \left[\tau^{-1} + \frac{\eta G^2}{\nu^2 c_2} \left(a\gamma_\star^2 + \frac{4nG^2}{(M + \mu)\kappa_\sigma} \right) - \frac{1}{n} \right] c_2 \textcolor{violet}{T}^{(k)} \\ & \quad + \left[-\frac{\eta}{2(M + \mu)\kappa_\sigma} + 3\eta^2 + \frac{c_1}{n} \right] \textcolor{violet}{Q}^{(k)}, \end{aligned}$$

where $\tau > 0$ is a to-be-specified rate constant. We now need to set the various free parameters a , c_1 , c_2 , and η to make each of the squared bracketed terms be non-positive. We enforce $\tau \geq 2n$ throughout. By setting

$$\eta = \frac{1}{12(\mu + M)\kappa_\sigma} \text{ and } c_1 = \frac{n\eta}{4(\mu + M)\kappa_\sigma},$$

we have that the bracketed constants before $c_1 \textcolor{green}{S}^{(k)}$ and $\textcolor{violet}{Q}^{(k)}$ vanish. Then, setting

$$a^{-1} = \frac{\mu}{2} \text{ and } c_2 = \frac{1}{48(\kappa + 1)\kappa_\sigma}$$

make the bracketed constant before $\|\boldsymbol{\theta}^{(k)} - \boldsymbol{\theta}^\star\|_2^2$, assuming that we enforce

$$\tau \geq 48(\kappa + 1)\kappa_\sigma.$$

We turn to the final constant after substituting the values of a , c_2 , and η . We need that

$$\frac{\eta G^2}{\nu^2 c_2} \left(a \gamma_\star^2 + \frac{8nG^2}{(M + \mu)\kappa_\sigma} \right) = \frac{8G^2}{\nu^2 \mu^2} \left(\gamma_\star^2 + \frac{2nG^2}{(\kappa + 1)\kappa_\sigma} \right) \leq \frac{1}{2n},$$

which occurs when

$$\nu^2 \geq \frac{16nG^2}{\mu^2} \left[\gamma_\star^2 + \frac{2nG^2}{(\kappa + 1)\kappa_\sigma} \right].$$

Because $\gamma_\star^2 \leq nG^2 \leq 2nG^2$, this is achieved when

$$\nu \geq \frac{8nG^2}{\mu},$$

completing the proof of the claim

$$\mathbb{E}_k [V^{(k+1)}] \leq (1 - \tau^{-1})V^{(k)}.$$

To complete the proof, we bound the initial terms. Because $c_3 = c_4 = 0$, we need only to bound $S^{(0)}$ and $T^{(0)}$.

$$\begin{aligned} S^{(0)} &= \frac{1}{n} \sum_{i=1}^n \|n\rho_i^{(0)} \nabla r_i(\tilde{\boldsymbol{\theta}}_i^{(0)}) - nq_i^\star \nabla r_i(\boldsymbol{\theta}^\star)\|_2^2 \\ &= \frac{1}{n} \sum_{i=1}^n \|nq_i^{(0)} \nabla r_i(\boldsymbol{\theta}^{(0)}) - nq_i^\star \nabla r_i(\boldsymbol{\theta}^\star)\|_2^2 \\ &\leq \frac{2}{n} \sum_{i=1}^n \|nq_i^{(0)} \nabla(r_i(\boldsymbol{\theta}^{(0)}) - \nabla r_i(\boldsymbol{\theta}^\star))\|_2^2 + \frac{2}{n} \sum_{i=1}^n \|n(q_i^{(0)} - q_i^\star) \nabla r_i(\boldsymbol{\theta}^\star)\|_2^2 \\ &\leq 2n \sum_{i=1}^n (q_i^{(0)})^2 M^2 \|\boldsymbol{\theta}^{(0)} - \boldsymbol{\theta}^\star\|_2^2 + 8nG^2 \|\mathbf{q}^{(0)} - \mathbf{q}^\star\|_2^2 \\ &\leq \left[2n \|\sigma\|_2^2 M^2 + \frac{8n^2 G^4}{\nu^2} \right] \|\boldsymbol{\theta}^{(0)} - \boldsymbol{\theta}^\star\|_2^2 \\ &\leq [2n \|\sigma\|_2^2 M^2 + \mu^2/8] \|\boldsymbol{\theta}^{(0)} - \boldsymbol{\theta}^\star\|_2^2 \leq 3nM^2 \|\boldsymbol{\theta}^{(0)} - \boldsymbol{\theta}^\star\|_2^2. \end{aligned}$$

This means ultimately that

$$c_1 S^{(0)} \leq \frac{n^2}{16(1 + \kappa^{-1})^2 \kappa_\sigma^2} \|\boldsymbol{\theta}^{(0)} - \boldsymbol{\theta}^\star\|_2^2.$$

Next, we have

$$c_2 T^{(0)} = \frac{n}{48(\kappa + 1)\kappa_\sigma} \|\boldsymbol{\theta}^{(0)} - \boldsymbol{\theta}^\star\|_2^2.$$

Thus, we can write

$$\begin{aligned} V^{(0)} &\leq \left[1 + \frac{n^2}{16(1 + \kappa^{-1})^2 \kappa_\sigma^2} + \frac{n}{48(\kappa + 1)\kappa_\sigma} \right] \|\boldsymbol{\theta}^{(0)} - \boldsymbol{\theta}^\star\|_2^2 \\ &\leq (1 + \sigma_n^{-1} + \sigma_n^{-2}) \|\boldsymbol{\theta}^{(0)} - \boldsymbol{\theta}^\star\|_2^2, \end{aligned}$$

completing the proof. \square

Theorem 2.5.1. *Suppose the shift cost satisfies*

$$\nu \geq 8nG^2/\mu.$$

Then, the sequence of iterates produced by Algorithm 1 with $\eta = 1/(12(\mu + M)\kappa_{\mathcal{P}})$ achieves

$$\mathbb{E}_0 \|\boldsymbol{\theta}^{(t)} - \boldsymbol{\theta}^\star\|_2^2 \leq (1 + \sigma_n^{-1} + \sigma_n^{-2}) \exp(-t/\tau) \|\boldsymbol{\theta}^{(0)} - \boldsymbol{\theta}^\star\|_2^2.$$

with

$$\tau = 2 \max\{n, 24\kappa_{\mathcal{P}}(\kappa + 1)\}.$$

Proof. First, we apply Lemma 2.6.1 with $\mathbf{q}' = \mathbf{q}^\star$, as well as Lemma 2.6.4, Lemma 2.6.2, and

Lemma 2.6.3, set $c_4 = 1$, and consolidate all constants to write

$$\mathbb{E}_k[V^{(k+1)}] - (1 - \tau^{-1})V^{(k)} \leq (\tau^{-1} - \eta\mu + c_2) \|\boldsymbol{\theta}^{(k)} - \boldsymbol{\theta}^*\|_2^2 \quad (\text{A.12})$$

$$+ \left[\tau^{-1} - \frac{1}{n} + \frac{2\alpha_1 G^2}{\nu(M + \mu)\kappa_\sigma} + \left(1 - \frac{1}{n}\right) \frac{G^2 c_3}{2\nu\mu n} \right] R^{(k)} \quad (\text{A.13})$$

$$+ \left[\tau^{-1} + \frac{1 + c_3}{c_1} \eta^2 (1 + \alpha_2^{-1}) - \frac{1}{n} \right] c_1 S^{(k)} \quad (\text{A.14})$$

$$+ \left[\tau^{-1} + \frac{\eta G^2 n}{2c_2 \nu} \alpha_3^{-1} + \frac{c_3 \eta M^2}{c_2 \mu n} \left(1 - \frac{1}{n}\right) - \frac{1}{n} \right] c_2 T^{(k)} \quad (\text{A.15})$$

$$+ \left[\tau^{-1} + \frac{2\eta G^2 n}{c_3 \nu} (1 + \alpha_3) - \frac{1}{n} \right] c_3 U^{(k)} \quad (\text{A.16})$$

$$+ \left[-\frac{\eta\alpha_1}{2(M + \mu)\kappa_\sigma} + \eta^2 (1 + c_3)(1 + \alpha_2) + \frac{c_1}{n} \right] Q^{(k)}. \quad (\text{A.17})$$

We first set $c_1 = \frac{n\eta\alpha_1}{4(M + \mu)\kappa_\sigma}$ and $c_2 = \eta\mu/2$ to clean up (A.12) and (A.17). We also drop the terms $(1 - 1/n) \leq 1$. Then, we notice in (A.13) that to achieve

$$\frac{2\alpha_1 G^2}{\nu(M + \mu)\kappa_\sigma} \leq \frac{1}{4n},$$

we need that $\alpha_1 \leq ((M + \mu)\kappa_\sigma)/(8nG^2/\nu)$. Combined with the requirement that $\alpha_1 \in [0, 1]$, we set $\alpha_1 = ((M + \mu)\kappa_\sigma)/(8nG^2/\nu + (M + \mu)\kappa_\sigma)$. We set $\alpha_2 = 2$, and can rewrite the expression above.

$$\begin{aligned} \mathbb{E}_k[V^{(k+1)}] - (1 - \tau^{-1})V^{(k)} &\leq \left(\tau^{-1} - \frac{\eta\mu}{2} \right) \|\boldsymbol{\theta}^{(k)} - \boldsymbol{\theta}^*\|_2^2 \\ &+ \left[\tau^{-1} - \frac{3}{4n} + \frac{G^2 c_3}{2\nu\mu n} \right] R^{(k)} \\ &+ \left[\tau^{-1} + \frac{6(1 + c_3)(M + \mu)\kappa_\sigma}{n\alpha_1} \eta - \frac{1}{n} \right] c_1 S^{(k)} \\ &+ \left[\tau^{-1} + \frac{G^2 n}{\mu\nu} \alpha_3^{-1} + \frac{c_3 M^2}{\mu^2 n} - \frac{1}{n} \right] c_2 T^{(k)} \\ &+ \left[\tau^{-1} + \frac{2\eta G^2 n}{c_3 \nu} (1 + \alpha_3) - \frac{1}{n} \right] c_3 U^{(k)} \\ &+ \left[-\frac{\eta\alpha_1}{4(M + \mu)\kappa_\sigma} + 3\eta^2 (1 + c_3) \right] Q^{(k)}. \end{aligned}$$

Next, set the learning rate to be

$$\eta \leq \frac{\alpha_1}{12(1+c_3)(M+\mu)\kappa_\sigma} \quad (\text{A.18})$$

to cancel out $Q^{(k)}$ and achieve

$$\begin{aligned} \mathbb{E}_k [V^{(k+1)}] - (1 - \tau^{-1})V^{(k)} &\leq \left(\tau^{-1} - \frac{\eta\mu}{2} \right) \|\boldsymbol{\theta}^{(k)} - \boldsymbol{\theta}^\star\|_2^2 \\ &\quad + \left[\tau^{-1} - \frac{3}{4n} + \frac{G^2 c_3}{2\nu\mu n} \right] R^{(k)} \\ &\quad + \left[\tau^{-1} - \frac{1}{2n} \right] c_1 S^{(k)} \\ &\quad + \left[\tau^{-1} + \frac{G^2 n}{\mu\nu} \alpha_3^{-1} + \frac{c_3 M^2}{\mu^2 n} - \frac{1}{n} \right] c_2 T^{(k)} \\ &\quad + \left[\tau^{-1} + \frac{2\eta G^2 n}{c_3 \nu} (1 + \alpha_3) - \frac{1}{n} \right] c_3 U^{(k)}. \end{aligned}$$

Requiring now that $\tau \geq 2n$, we may also cancel the $S^{(k)}$ term. We substitute $\delta = nG^2/(\mu\nu)$ to achieve

$$\begin{aligned} \mathbb{E}_k [V^{(k+1)}] - (1 - \tau^{-1})V^{(k)} &\leq \left(\tau^{-1} - \frac{\eta\mu}{2} \right) \|\boldsymbol{\theta}^{(k)} - \boldsymbol{\theta}^\star\|_2^2 \\ &\quad + \left[-\frac{1}{4n} + \frac{c_3 \delta}{2n^2} \right] R^{(k)} \\ &\quad + \left[-\frac{1}{2n} + \frac{\delta}{\alpha_3} + \frac{c_3 M^2}{\mu^2 n} \right] c_2 T^{(k)} \\ &\quad + \left[-\frac{1}{2n} + \frac{2\mu\eta\delta}{c_3} (1 + \alpha_3) \right] c_3 U^{(k)}. \end{aligned}$$

It remains to select c_3 and α_3 . As such, we set $\alpha_3 = 4n\delta$ and use that $1 + 4n\delta \leq 8n\delta$ when

$n \geq 2$ and $\delta \geq 1/8$ as assumed, and so

$$\begin{aligned} \mathbb{E}_k [V^{(k+1)}] - (1 - \tau^{-1})V^{(k)} &\leq \left(\tau^{-1} - \frac{\eta\mu}{2} \right) \|\boldsymbol{\theta}^{(k)} - \boldsymbol{\theta}^*\|_2^2 \\ &\quad + \left[-\frac{1}{4n} + \frac{c_3\delta}{2n^2} \right] \textcolor{brown}{R}^{(k)} \\ &\quad + \left[-\frac{1}{4n} + \frac{c_3\kappa^2}{n} \right] c_2 \textcolor{violet}{T}^{(k)} \\ &\quad + \left[-\frac{1}{2n} + \frac{16n\mu\eta\delta^2}{c_3} \right] c_3 \textcolor{blue}{U}^{(k)}. \end{aligned}$$

We require now that

$$c_3 = \frac{1}{2} \min \left\{ \frac{1}{2\kappa^2}, \frac{n}{\delta} \right\},$$

which cancels $\textcolor{violet}{T}^{(k)}$ and $\textcolor{brown}{R}^{(k)}$, leaving

$$\begin{aligned} \mathbb{E}_k [V^{(k+1)}] - (1 - \tau^{-1})V^{(k)} &\leq \left(\tau^{-1} - \frac{\eta\mu}{2} \right) \|\boldsymbol{\theta}^{(k)} - \boldsymbol{\theta}^*\|_2^2 \\ &\quad + \left[-\frac{1}{2n} + 32\mu\eta\delta^2 \max \{2n\kappa^2, \delta\} \right] c_3 \textcolor{blue}{U}^{(k)}. \end{aligned}$$

From the above, we retrieve the requirement that

$$\eta \leq \frac{1}{64n\mu\delta^2 \max \{2n\kappa^2, \delta\}}. \quad (\text{A.19})$$

It now remains to set η . By substituting in the values for α_1 and c_3 into (A.18), we have that

$$\begin{aligned} \eta &\stackrel{\text{want}}{\leq} \frac{\alpha_1}{12(1+c_3)(M+\mu)\kappa_\sigma} = \frac{1}{12(1+c_3)[8\mu\delta + (M+\mu)\kappa_\sigma]} \\ &\geq \frac{1}{(12+6n/\delta)[8\mu\delta + (M+\mu)\kappa_\sigma]} \\ &\geq \frac{1}{(12+48n)[8\mu\delta + (M+\mu)\kappa_\sigma]} \\ &\geq \frac{1}{96n[8\mu\delta + (M+\mu)\kappa_\sigma]}. \end{aligned}$$

The combination of (A.19) and the above display yields

$$\begin{aligned}\eta &= \min \left\{ \frac{1}{96n[8\mu\delta + (M + \mu)\kappa_\sigma]}, \frac{1}{64n\mu\delta^2 \max \{2n\kappa^2, \delta\}} \right\} \\ &= \frac{1}{16n\mu} \min \left\{ \frac{1}{6[8\delta + (\kappa + 1)\kappa_\sigma]}, \frac{1}{4\delta^2 \max \{2n\kappa^2, \delta\}} \right\}.\end{aligned}$$

We need finally that $\tau \geq 2/(\mu\eta)$, resulting in the requirement

$$\tau \geq 32n \max \left\{ 6[8\delta + (\kappa + 1)\kappa_\sigma], 4\delta^2 \max \{2n\kappa^2, \delta\} \right\}.$$

This is achieved by setting

$$\tau = 32n \max \left\{ 6[8\delta + (\kappa + 1)\kappa_\sigma], 4\delta^2 \max \{2n\kappa^2, \delta\}, 1/16 \right\}.$$

completing the proof of the claim

$$\mathbb{E}_k [V^{(k+1)}] \leq (1 - \tau^{-1})V^{(k)}.$$

Next, we bound the initial terms to achieve the final rate. First, we bound η which is used in all of the terms. Because $\delta \geq 1/8$,

$$\eta \leq \frac{1}{16n\mu} \cdot \frac{1}{4\delta^2 \max \{2n\kappa^2, \delta\}} \leq \frac{1}{64n\mu\delta^3} \leq \frac{8}{n\mu}. \quad (\text{A.20})$$

Then,

$$\begin{aligned}
S^{(0)} &= \frac{1}{n} \sum_{i=1}^n \|n\rho_i^{(0)} \nabla r_i(\tilde{\boldsymbol{\theta}}_i^{(0)}) - nq_i^* \nabla r_i(\boldsymbol{\theta}^*)\|_2^2 \\
&= \frac{1}{n} \sum_{i=1}^n \|nq_i^{(0)} \nabla r_i(\boldsymbol{\theta}^{(0)}) - nq_i^* \nabla r_i(\boldsymbol{\theta}^*)\|_2^2 \\
&\leq \frac{2}{n} \sum_{i=1}^n \|nq_i^{(0)} \nabla(r_i(\boldsymbol{\theta}^{(0)}) - \nabla r_i(\boldsymbol{\theta}^*))\|_2^2 + \frac{2}{n} \sum_{i=1}^n \|n(q_i^{(0)} - q_i^*) \nabla r_i(\boldsymbol{\theta}^*)\|_2^2 \\
&\leq 2n \sum_{i=1}^n (q_i^{(0)})^2 M^2 \|\boldsymbol{\theta}^{(0)} - \boldsymbol{\theta}^*\|_2^2 + 8nG^2 \|\mathbf{q}^{(0)} - \mathbf{q}^*\|_2^2 \\
&\leq \left[2n \|\sigma\|_2^2 M^2 + \frac{8n^2 G^2}{\nu^2} \right] \|\boldsymbol{\theta}^{(0)} - \boldsymbol{\theta}^*\|_2^2 \\
&\leq [2n \|\sigma\|_2^2 M^2 + \mu^2/8] \|\boldsymbol{\theta}^{(0)} - \boldsymbol{\theta}^*\|_2^2 \leq 3nM^2 \|\boldsymbol{\theta}^{(0)} - \boldsymbol{\theta}^*\|_2^2.
\end{aligned}$$

Continuing with $\alpha_1 \leq 1$ and (A.20),

$$\begin{aligned}
c_1 S^{(0)} &= \frac{n\eta\alpha_1}{4(M+\mu)\kappa_\sigma} S^{(0)} \\
&\leq \frac{2}{\mu(M+\mu)\kappa_\sigma} \cdot 3nM^2 \|\boldsymbol{\theta}^{(0)} - \boldsymbol{\theta}^*\|_2^2 \\
&\leq \frac{6n\kappa^2}{(1+\kappa)\kappa_\sigma} \|\boldsymbol{\theta}^{(0)} - \boldsymbol{\theta}^*\|_2^2 \\
&\leq \frac{6\kappa^2}{\sigma_n} \|\boldsymbol{\theta}^{(0)} - \boldsymbol{\theta}^*\|_2^2.
\end{aligned}$$

Next, we have $T^{(0)} = n \|\boldsymbol{\theta}^{(0)} - \boldsymbol{\theta}^*\|_2^2$ and by (A.20),

$$\begin{aligned}
c_2 T^{(0)} &= \frac{\eta\mu}{2} \cdot n \|\boldsymbol{\theta}^{(0)} - \boldsymbol{\theta}^*\|_2^2 \\
&\leq 4 \|\boldsymbol{\theta}^{(0)} - \boldsymbol{\theta}^*\|_2^2.
\end{aligned}$$

Because $U^{(0)} = 0$, it is bounded trivially. For $R^{(0)}$, with $c_4 = 1$ we have

$$\begin{aligned}
R^{(0)} &= 2n\eta(q^{\text{opt}}\ell(\boldsymbol{\theta}^{(0)}) - q^{\text{opt}}\ell(\boldsymbol{\theta}^*))^\top (\ell(\boldsymbol{\theta}^{(0)}) - \ell(\boldsymbol{\theta}^*)) \\
&\leq \frac{2n\eta}{\nu} \|\ell(\boldsymbol{\theta}^{(0)}) - \ell(\boldsymbol{\theta}^*)\|_2^2 \\
&\leq \frac{2n^2\eta G^2}{\nu} \|\boldsymbol{\theta}^{(0)} - \boldsymbol{\theta}^*\|_2^2 \\
&\leq \frac{16nG^2}{\mu\nu} \|\boldsymbol{\theta}^{(0)} - \boldsymbol{\theta}^*\|_2^2 \\
&= 16\delta \|\boldsymbol{\theta}^{(0)} - \boldsymbol{\theta}^*\|_2^2.
\end{aligned}$$

Combining each of these terms together, we have that

$$V^{(0)} \leq \left(5 + 16\delta + \frac{6\kappa^2}{\sigma_n}\right) \|\boldsymbol{\theta}^{(0)} - \boldsymbol{\theta}^*\|_2^2,$$

completing the proof. \square

A.3 Implementation Details

In this section, we describe Prospect including computational details, in a way that is amenable to implementation. Particular attention is given to the case when $\mathcal{Q} \equiv \mathcal{Q}(\sigma)$ is the spectral risk measure uncertainty set and the penalty is an f -divergence.

Efficient Implementation We exactly solve the maximization problem

$$\mathbf{q} = q^{\text{opt}}(\mathbf{l}) = \arg \max_{\mathbf{q} \in \mathcal{Q}(\sigma)} \left\{ \langle \mathbf{q}, \mathbf{l} \rangle - (\nu/n) \sum_{i=1}^n f(nq_i) \right\}. \quad (\text{A.21})$$

by a sequence of three steps:

- **Sorting:** Find π such that $l_{\pi(1)} \leq \dots \leq l_{\pi(n)}$.
- **Isotonic regression:** Apply Pool Adjacent Violators (PAV) (Subroutine 1) to solve the isotonic regression minimization problem (2.28), yielding solution $\mathbf{z} = z^{\text{opt}}(\mathbf{l})$.
- **Conversion:** Use (2.29) to convert \mathbf{z} back to $\mathbf{q} = q^{\text{opt}}(\mathbf{l})$.

The sorting step runs in $O(n \ln n)$ elementary operations, whereas the isotonic regression and conversion steps run in $O(n)$ operations. Crucially, retrieving \mathbf{q} from the output $\mathbf{z} = \mathbf{z}^{\text{opt}}(\mathbf{l})$ in the third step can be done by a single $O(n)$ -time pass by setting

$$q_{\pi(i)} = \frac{1}{n} [f^*]' \left(\frac{1}{\nu} (l_{\pi(i)} - z_i) \right)$$

for $i = 1, \dots, n$, as opposed to computing the inverse π^{-1} and using (2.29) directly, which in fact requires another sorting operation and can be avoided. Because only one element of \mathbf{l} changes on every iteration, we may sort it by simply bubbling the value of the index that changed into its correct position to generate the newly sorted version. The full algorithm is given Algorithm 6. We give a brief explanation on the PAV algorithm for general f -divergences below.

Pool Adjacent Violators (PAV) Algorithm First, recall the optimization problem we wish to solve:

$$\min_{\substack{\mathbf{z} \in \mathbb{R}^n \\ \mathbf{z}_1 \leq \dots \leq \mathbf{z}_n}} \sum_{i=1}^n g_i(z_i; \mathbf{l}), \quad \text{where} \quad g_i(z_i; \mathbf{l}) := \sigma_i z_i + \frac{\nu}{n} f^* \left(\frac{l_{\pi(i)} - z_i}{\nu} \right). \quad (\text{A.22})$$

The objective can be thought of as fitting a real-valued monotonic function to the points $(1, l_{\pi(1)}), \dots, (n, l_{\pi(n)})$, which would require specifying its values (c_1, \dots, c_n) on $(1, \dots, n)$ and defining the function as any $x \in [c_j, c_{j+1}]$ on $(j, j+1)$. Because $l_{\pi(1)} \leq \dots \leq l_{\pi(n)}$, if we evaluated our function (c_1, \dots, c_n) on a loss such as $\sum_{i=1}^n (l_{\pi(i)} - z_i)^2$, we may easily solve the problem by returning $c_1 = l_{\pi(1)}, \dots, c_n = l_{\pi(n)}$. However, by specifying functions g_1, \dots, g_n we allow our loss function to change in different regions of the input space $\{1, \dots, n\}$. In such cases, the monotonicity constraint $c_1 \leq \dots \leq c_n$ is often violated because individually minimizing $g_i(z_i)$ for each z_i has no guarantee of yielding a function that is monotonic.

The idea behind the PAV algorithm is to attempt a pass at minimizing each g_i individually, and correcting *violations* as they appear. To provide intuition, define $c_i^* \in \arg \min_{z_i \in \mathbb{R}} g_i(z_i)$, and consider $i < j$ such that $c_i^* > c_j^*$. If f^* is strictly convex, then $g_i(x) > g_i(c_i^*)$ for any $x < c_i^*$ and similarly $g_j(x) > g_j(c_j^*)$ for any $x > c_j^*$. Thus, to correct

the violation, we decrease c_i^* to \bar{c}_i and increase c_j^* to \bar{c}_j until $\bar{c}_i = \bar{c}_j$. We determine this midpoint precisely by

$$\bar{c}_i = \bar{c}_j = \arg \min_{x \in \mathbb{R}} g_i(x) + g_j(x)$$

as these are exactly the contributions made by these terms in the overall objective. The computation above is called *pooling* the indices i and j . We may generalize this viewpoint to *violating chains*, that is collections of contiguous indices $(i, i+1, \dots, i+m)$ such that $c_j^* < c_i^*$ for all $j < i$ and $c_j^* > c_{i+m}^*$ for all $j > i+m$, but $c_i^* > c_{i+m}^*$. One approach is to use dynamic programming to identify such chains and then compute the pooled quantities

$$\bar{c}_i = \arg \min_{x \in \mathbb{R}} \sum_{j=1}^m g_{i+j}(x).$$

This requires two passes through the vector: one for identifying violators and the other for pooling. The Pool Adjacent Violators algorithm, on the other hand, is able to perform both operations in one pass by greedily pooling violators as they appear. This can be viewed as a meta-algorithm, as it hinges on the notion that the solution of “larger” pooling problems can be easily computed from solutions of “smaller” pooling problems. Precisely, for indices $S \subseteq [n] = \{1, \dots, n\}$ define

$$\text{Sol}(S) = \arg \min_{x \in \mathbb{R}} \sum_{i \in S} g_i(x).$$

We rely on the existence of an operation Pool , such that for any $S, T \subseteq [n]$ such that $S \cap T = \emptyset$, we have that

$$\text{Sol}(S \cup T) = \text{Pool}(\text{Sol}(S), \text{Met}(S), \text{Sol}(T), \text{Met}(T)), \quad (\text{A.23})$$

where $\text{Met}(S)$ denotes “metadata” associated to S , and that the number of elementary operations in the Pool function is $O(1)$ with respect to $|S| + |T|$. We review our running examples.

For the χ^2 -divergence, we have that $f_{\chi^2}(x) = x^2 - 1$ and $f_{\chi^2}^*(y) = y^2/4 + 1$, so

$$\begin{aligned} \text{Sol}(S) &= \arg \min_{x \in \mathbb{R}} \left\{ x \left(\sum_{i \in S} \sigma_i \right) + |S| + \frac{1}{4n\nu} \sum_{i \in S} (l_{\pi(i)} - x)^2 \right\} \\ &= \frac{1}{|S|} \left[\sum_{i \in S} l_{\pi(i)} - 2n\nu \sum_{i \in S} \sigma_i \right] \\ \text{Sol}(S \cup T) &= \frac{1}{|S| + |T|} \left[\sum_{i \in S \cup T} l_{\pi(i)} - 2n\nu \sum_{i \in S \cup T} \sigma_i \right] \\ &= \frac{|S| \text{Sol}(S) + |T| \text{Sol}(T)}{|S| + |T|}. \end{aligned}$$

Thus, the metadata $\text{Met}(S) = |S|$ used in the pooling step (A.23) is the size of each subset.

For the KL divergence, $f_{\text{KL}}(x) = x \ln x$ and $f_{\text{KL}}^*(y) = e^{-1} \exp(y)$, so so

$$\begin{aligned} \text{Sol}(S) &= \arg \min_{x \in \mathbb{R}} \left\{ x \left(\sum_{i \in S} \sigma_i \right) + \frac{\nu}{ne} \sum_{i \in S} \exp(l_{\pi(i)}/\nu) \exp(-x/\nu) \right\} \\ &= \nu \left[\ln \sum_{i \in S} \exp(l_{\pi(i)}/\nu) - \ln \sum_{i \in S} \sigma_i - \ln n - 1 \right] \\ \text{Sol}(S \cup T) &= \nu \left[\ln \sum_{i \in S \cup T} \exp(l_{\pi(i)}/\nu) - \ln \sum_{i \in S \cup T} \sigma_i - \ln n - 1 \right] \\ &= \nu \left[\ln \left(\sum_{i \in S} \exp(l_{\pi(i)}/\nu) + \sum_{i \in T} \exp(l_{\pi(i)}/\nu) \right) - \ln \left(\sum_{i \in S} \sigma_i + \sum_{i \in T} \sigma_i \right) - \ln n - 1 \right]. \end{aligned}$$

Here, we carry the metadata $\text{Met}(S) = (\ln \sum_{i \in S} \exp(l_{\pi(i)}/\nu), \ln \sum_{i \in S} \sigma_i)$, which can easily be combined and plugged into the function

$$(m_1, m_2), (m'_1, m'_2) \mapsto \nu [\ln(\exp m_1 + \exp m'_1) - \ln(\exp m_2 + \exp m'_2) - \ln n - 1]. \quad (\text{A.24})$$

for two instances of metadata (m_1, m_2) and (m'_1, m'_2) . We carry the “logsumexp” instead of just the sum of exponential quantities for numerical stability, and Equation (A.24) applies this operation as well. It might be that $\sum_{i \in S} \sigma_i = 0$, e.g. for the superquantile. In this case, we may interpret $\text{Sol}(S) = -\infty$ and evaluate $\exp(-\infty) = 0$ in the conversion formula (A.22). Two examples of the PAV algorithm are given in Subroutine 1 and Subroutine 2,

Algorithm 6 Prospect (with exact implementation details)

Inputs: Initial points $\theta_0(= \mathbf{0})$, spectrum $\sigma(= 2\text{-extremile})$, stepsize $\eta(= 0.01)$, number of iterations $t(= 1000)$, tolerance $\varepsilon(= 10^{-4})$, regularization parameter $\mu(= 1/n)$, shift cost $\nu(= 0.05)$, loss/gradient oracles ℓ_1, \dots, ℓ_n and $\nabla \ell_1, \dots, \nabla \ell_n$.

- 1: $\mathbf{l} \leftarrow \ell(\theta_0) \in \mathbb{R}^n$.
- 2: $\mathbf{g} \leftarrow (\nabla \ell_i(\theta_0) + \mu \theta_0)_{i=1}^n \in \mathbb{R}^{n \times d}$.
- 3: $\pi \leftarrow \text{argsort}(\mathbf{l})$.
- 4: $\mathbf{c} \leftarrow \text{PAV}(\mathbf{l}, \pi, \sigma)$. ▷ Subroutine 1 or Subroutine 2
- 5: $\mathbf{q} \leftarrow \text{Convert}(\mathbf{c}, \mathbf{l}, \pi, \nu, \mathbf{0}_n)$. ▷ Subroutine 3
- 6: $\rho \leftarrow \mathbf{q}$.
- 7: $\bar{\mathbf{g}} \leftarrow \sum_{i=1}^n \rho_i \mathbf{g}_i \in \mathbb{R}^d$.
- 8: **for** $k = 1, \dots, t$ **do**
- 9: If n divides k , then check if certificate $\leq \varepsilon$ (Section 3.7.1). If so, terminate.
- 10: Sample $i, j \sim \text{Unif}[n]$.
- 11: $\mathbf{v} \leftarrow n q_i (\nabla \ell_i(\theta) + \mu \theta) - n \rho_i \mathbf{g}_i - \bar{\mathbf{g}}$. ▷ Iterate Update
- 12: $\theta \leftarrow \theta - \eta \mathbf{v}$.
- 13: $\mathbf{l}_j \leftarrow \ell_j(\theta)$. ▷ Bias Reducing Update
- 14: $\pi \leftarrow \text{BubbleArgSort}(\pi, \mathbf{l})$. ▷ Subroutine 4
- 15: $\mathbf{c} \leftarrow \text{PAV}(\mathbf{l}, \pi, \sigma)$.
- 16: $\mathbf{q} \leftarrow \text{Convert}(\mathbf{c}, \mathbf{l}, \pi, \nu, \mathbf{q})$.
- 17: $\bar{\mathbf{g}} \leftarrow \bar{\mathbf{g}} - \rho_i \mathbf{g}_i + q_i (\nabla \ell_i(\theta) + \mu \theta)$. ▷ Variance Reducing Update
- 18: $\mathbf{g}_i \leftarrow \nabla \ell_i(\theta) + \mu \theta$.
- 19: $\rho_i \leftarrow q_i$.

Output: Final point θ .

respectively. These operate by selecting the unique values of the optimizer and partitions of indices that achieve that value.

Hardware Acceleration Finally, note that all of the subroutines in Algorithm 6 (Subroutine 1/Subroutine 2, Subroutine 3, and Subroutine 4) all require primitive operations such as control flow and linear scans through vectors. Because these steps are outside of the purview of oracle calls or matrix multiplications, they benefit from just-in-time compilation on the CPU. We accelerate these subroutines using the Numba package in Python and are able to achieve an approximate 50%-60% decrease in runtime across benchmarks.

Subroutine 1 Pool Adjacent Violators (PAV) Algorithm for χ^2 divergence

Inputs: Losses $(\ell_i)_{i \in [n]}$, argsort π , and spectrum $(\sigma_i)_{i \in [n]}$.

- 1: Initialize partition endpoints $(b_0, b_1) = (0, 1)$, partition value $v_1 = l_{\pi(1)} - 2n\nu\sigma_1$, number of parts $K = 1$.
- 2: **for** $i = 2, \dots, n$ **do**
- 3: Add part $K = K + 1$.
- 4: Compute $v_K = l_{\pi(i)} - 2n\nu\sigma_i$.
- 5: **while** $K \geq 2$ and $v_{K-1} \geq v_K$ **do**
- 6: $v_{K-1} = \frac{(b_K - b_{K-1})v_{K-1} + (i - b_K)v_K}{i - b_{K-1}}$.
- 7: Set $K = K - 1$.
- 8: $b_K = i$.

Output: Vector c containing $z_i = v_K$ for $b_{K-1} < i \leq b_K$.

Subroutine 2 Pool Adjacent Violators (PAV) Algorithm for KL divergence

Inputs: Losses $(\ell_i)_{i \in [n]}$, argsort π , and spectrum $(\sigma_i)_{i \in [n]}$.

- 1: Initialize partition endpoints $(b_0, b_1) = (0, 1)$, number of parts $K = 1$.
- 2: Initialize partition value $v_1 = \nu (l_{\pi(1)}/\nu - \ln \sigma_1 - \ln n - 1)$.
- 3: Initialize metadata $m_1 = \ell_{\pi(1)}/\nu$ and $t_1 = \ln \sigma_1$.
- 4: **for** $i = 2, \dots, n$ **do**
- 5: Add part $K = K + 1$.
- 6: Compute $v_K = \nu (l_{\pi(i)}/\nu - \ln \sigma_i - \ln n - 1)$.
- 7: Compute $m_K = \ell_{\pi(i)}/\nu$ and $t_K = \ln \sigma_i$
- 8: **while** $k \geq 2$ and $v_{K-1} \geq v_K$ **do**
- 9: $m_{K-1} = \text{logsumexp}(m_{K-1}, m_K)$ and $t_{K-1} = \text{logsumexp}(t_{K-1}, t_K)$.
- 10: $v_{K-1} = \nu (m_{K-1} - t_{K-1} - \ln n - 1)$.
- 11: Set $K = K - 1$.
- 12: $b_K = i$.

Output: Vector c containing $z_i = v_K$ for $b_{K-1} < i \leq b_K$.

Subroutine 3 Convert

Require: Sorted vector $c \in \mathbb{R}$, vector $\mathbf{l} \in \mathbb{R}^n$, argsort π of \mathbf{l} , shift cost $\nu \geq 0$, vector $\mathbf{q} \in \mathbb{R}^n$.

- 1: **for** $i = 1, \dots, n$ **do**
 - 2: Set $q_{\pi(i)} = (1/n)[f^*]'((l_{\pi(i)} - z_i)/\nu)$.
 - 3: **return** \mathbf{q} .
-

Subroutine 4 BubbleArgSort

Require: Index j_{init} , sorting permutation π , loss table \mathbf{l} .

- 1: $j = j_{\text{init}}$. ▷ If $l_{\pi(j_{\text{init}})}$ too small, bubble left.
 - 2: **while** $j > 1$ and $l_{\pi(j)} < l_{\pi(j-1)}$ **do**
 - 3: Swap $\pi(j)$ and $\pi(j-1)$.
 - 4: $j = j_{\text{init}}$. ▷ If $l_{\pi(j_{\text{init}})}$ too large, bubble right.
 - 5: **while** $j < n$ and $l_{\pi(j)} > l_{\pi(j+1)}$ **do**
 - 6: Swap $\pi(j)$ and $\pi(j+1)$.
 - 7: **return** π
-

Dataset	d	n_{train}	n_{test}	Task	Source
yacht	6	244	62	Regression	UCI
energy	8	614	154	Regression	UCI
concrete	8	824	206	Regression	UCI
kin8nm	8	6,553	1,639	Regression	OpenML
power	4	7,654	1,914	Regression	UCI
diabetes	33	4,000	1,000	Binary Classification	Fairlearn
acsincome	202	4,000	1,000	Regression	Fairlearn
amazon	535	10,000	10,000	Multiclass Classification	WILDS
iwildcam	9420	20,000	5,000	Multiclass Classification	WILDS

Table A.1: Dataset attributes and dimensionality d , train sample size n_{train} , and test sample size n_{test} .

A.4 Experimental Details

A.4.1 Tasks & Objectives

In all settings, we consider supervised learning tasks specified by losses of the form

$$\ell_i(\boldsymbol{\theta}) = h(y_i, \langle \boldsymbol{\theta}, \varphi(\mathbf{x}_i) \rangle),$$

where we consider an input $\mathbf{x}_i \in \mathcal{X}$, a feature map $\varphi : \mathcal{X} \rightarrow \mathbb{R}^d$, and a label $y_i \in \mathcal{Y}$. The function $h : \mathcal{Y} \times \mathbb{R} \rightarrow \mathbb{R}$ measures the error between the true label and another value which is the prediction in regression and the logit probabilities of the associated classes in classification. In the regression tasks, $\mathcal{Y} = \mathbb{R}$ and we used the squared loss

$$\ell_i(\boldsymbol{\theta}) = \frac{1}{2}(y_i - \langle \boldsymbol{\theta}, \varphi(\mathbf{x}_i) \rangle)^2.$$

For binary classification, we have $\mathcal{Y} = \{-1, 1\}$, denoting a negative and positive class. We used the binary logistic loss

$$\ell_i(\boldsymbol{\theta}) = -y_i \langle \boldsymbol{\theta}, \varphi(\mathbf{x}_i) \rangle + \ln(1 + e^{\langle \boldsymbol{\theta}, \varphi(\mathbf{x}_i) \rangle}).$$

For multiclass classification, $\mathcal{Y} = \{1, \dots, C\}$ where C is the number of classes. We used the multinomial logistic loss:

$$\ell_i(\boldsymbol{\theta}) = -\ln p_{y_i}(\mathbf{x}_i; \boldsymbol{\theta}), \text{ where } p_{y_i}(\mathbf{x}_i; \boldsymbol{\theta}) := \frac{\exp(\langle \boldsymbol{\theta}_{\cdot y_i}, \varphi(\mathbf{x}_i) \rangle)}{\sum_{y'=1}^C \exp(\langle \boldsymbol{\theta}_{\cdot y'}, \varphi(\mathbf{x}_i) \rangle)}, \quad \boldsymbol{\theta} \in \mathbb{R}^{d \times C}$$

The design matrix $(\varphi(\mathbf{x}_1), \dots, \varphi(\mathbf{x}_n)) \in \mathbb{R}^{n \times d}$ is standardized to have columns with zero mean and unit variance, and the estimated mean and variance from the training set is used to standardize the test sets as well. Our final objectives are of the form

$$\mathcal{L}(\boldsymbol{\theta}) = \max_{\mathbf{q} \in \mathcal{Q}(\sigma)} \sum_{i=1}^n q_i \ell_i(\boldsymbol{\theta}) - \nu n \|\mathbf{q} - \mathbf{1}/n\|_2^2 + \frac{\mu}{2} \|\boldsymbol{\theta}\|_2^2$$

for shift cost $\nu \geq 0$ and regularization constant $\mu \geq 0$.

A.4.2 Datasets

We detail the datasets used in the experiments. If not specified below, the input space $\mathcal{X} = \mathbb{R}^d$ and φ is the identity map. The sample sizes, dimensions, and source of the datasets are summarized in Table A.1, where d refers to the dimension of each $\varphi(\mathbf{x}_i)$.

- (a) **yacht**: prediction of the residuary resistance of a sailing yacht based on its physical attributes [Tsanas and Xifara, 2012].
- (b) **energy**: prediction of the cooling load of a building based on its physical attributes Baressi Segota et al. [2020].
- (c) **concrete**: prediction of the compressive strength of a concrete type based on its physical and chemical attributes [Yeh, 2006].
- (d) **kin8nm**: prediction of the distance of an 8-link all-revolute robot arm to a spatial endpoint [Akujuobi and Zhang, 2017].

- (e) **power**: prediction of net hourly electrical energy output of a power plant given environmental factors [Tüfekci, 2014].
- (f) **diabetes**: prediction of readmission for diabetes patients based on 10 years' worth of clinical care data at 130 US hospitals [Rizvi et al., 2014].
- (g) **acsincome**: prediction of income of US adults given features compiled from the American Community Survey (ACS) Public Use Microdata Sample (PUMS) [Ding et al., 2021].
- (h) **amazon**: prediction of the review score of a sentence taken from Amazon products. Each input $x \in \mathcal{X}$ is a sentence in natural language and the feature map $\varphi(\mathbf{x}) \in \mathbb{R}^d$ is generated by the following steps:

- A BERT neural network [Devlin et al., 2019a] (fine-tuned on 10,000 held-out examples) is applied to the text \mathbf{x}_i , resulting in vector \mathbf{x}'_i .
- The $\mathbf{x}'_1, \dots, \mathbf{x}'_n$ are normalized to have unit norm.
- Principle Components Analysis (PCA) is applied, resulting in 105 components that explain 99% of the variance, resulting in vectors $\mathbf{x}''_i \in \mathbb{R}^{105}$. The d in Table A.1 refers to the total dimension of the parameter vectors for all 5 classes.

- (i) **iwildcam**: prediction of an animal or flora in an image from wilderness camera traps, with heterogeneity in illumination, camera angle, background, vegetation, color, and relative animal frequencies [Beery et al., 2020]. Each input $x \in \mathcal{X}$ is an image the feature map $\varphi(\mathbf{x}) \in \mathbb{R}^d$ is generated by the following steps:

- A ResNet50 neural network [He et al., 2016] that is pretrained on ImageNet [Deng et al., 2009] is applied to the image \mathbf{x}_i , resulting in vector \mathbf{x}'_i .
- The $\mathbf{x}'_1, \dots, \mathbf{x}'_n$ are normalized to have unit norm.
- Principle Components Analysis (PCA) is applied, resulting in $d = 157$ components that explain 99% of the variance. The d in Table A.1 refers to the total dimension of the parameter vectors for all 60 classes.

A.4.3 Hyperparameter Selection

We fix a minibatch size of 64 SGD and SRDA and an epoch length of $N = n$ for LSVRG. For SaddleSAGA we consider three schemes for selecting the primal and dual learning rates that reduce the search to a single parameter $\eta > 0$ by tuning a scaling of the primal and dual learning rates that performs well across experiments. In practice, the regularization parameter μ and shift cost ν are tuned by a statistical metric, i.e., generalization error as measured on a validation set. We study the optimization performance of the methods for multiple values of each in Appendix A.4.5.

For the tuned hyperparameters, we use the following method. Let $k \in \{1, \dots, K\}$ be a seed that determines algorithmic randomness. This corresponds to sampling a minibatch without replacement for SGD and SRDA and a single sampled index for SaddleSAGA, LSVRG, and Prospect. Letting $\mathcal{L}_k(\eta)$ denote the average value of the training loss of the last ten passes using learning rate η and seed k , the quantity $\mathcal{L}(\eta) = \frac{1}{K} \sum_{k=1}^K \mathcal{L}_k(\eta)$ was minimized to select η . The learning rate η is chosen in the set $\{1 \times 10^{-4}, 3 \times 10^{-4}, 1 \times 10^{-3}, 3 \times 10^{-3}, 1 \times 10^{-2}, 3 \times 10^{-2}, 1 \times 10^{-1}, 3 \times 10^{-1}, 1 \times 10^0, 3 \times 10^0\}$, with two orders of magnitude lower numbers used in `acsincome` due to its sparsity. We discard any learning rates that cause the optimizer to diverge for any seed.

A.4.4 Compute Environment

No GPUs were used in the study; Experiments were run on a CPU workstation with an Intel i9 processor, a clock speed of 2.80GHz, 32 virtual cores, and 126G of memory. The code used in this project was written in Python 3 using the PyTorch and Numba packages for automatic differentiation and just-in-time compilation, respectively.

A.4.5 Additional Experiments

Varying Risk Parameters We study the effect of varying the risk parameters, that is (τ, r, γ) for the τ -CVaR (Equation (2.13)), r -extremile (Equation (2.14)), γ -ESRM (Equa-

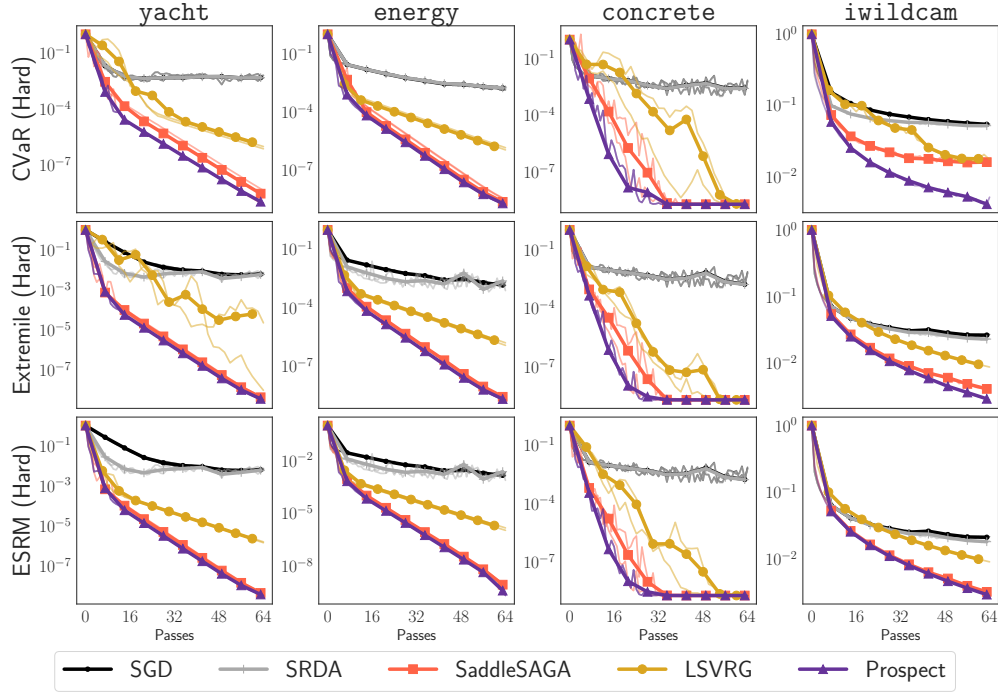


Figure A.1: **Harder Risk Parameter Settings.** Each row represents a different “hard” variant of the superquantile, extremile, and ESRM spectra. Columns represent different datasets. Suboptimality (2.49) is measured on the y -axis while the x -axis measures the total number of gradient evaluations made divided by n , i.e., the number of passes through the training set.

tion (2.15)), choosing the spectrum to increase the condition number $\kappa_\sigma = n\sigma_n$ compared to the experiments in the main text. We use $\tau = 0.75$, $r = 2.5$, and $\gamma = 1/e^{-2}$ to generate “hard” version of the superquantile, extremile, and ESRM. Figure A.1 plots the corresponding training curves for four datasets of varying sample sizes: **yacht**, **energy**, **concrete**, and **iwildcam**. We see that the comparison of methods is the same as the original methods, that is that Prospect performs the best or close to best in terms of optimization trajectories. Except on **concrete**, SaddleSAGA generally matches the performance of Prospect. The trajectory of LSVRG is noticeably noisier than on the original settings; we hypothesize that the bias accrued by this epoch-based algorithm is exacerbated by the skewness in the spectrum, as mentioned in Mehta et al. [2023, Proposition 1].

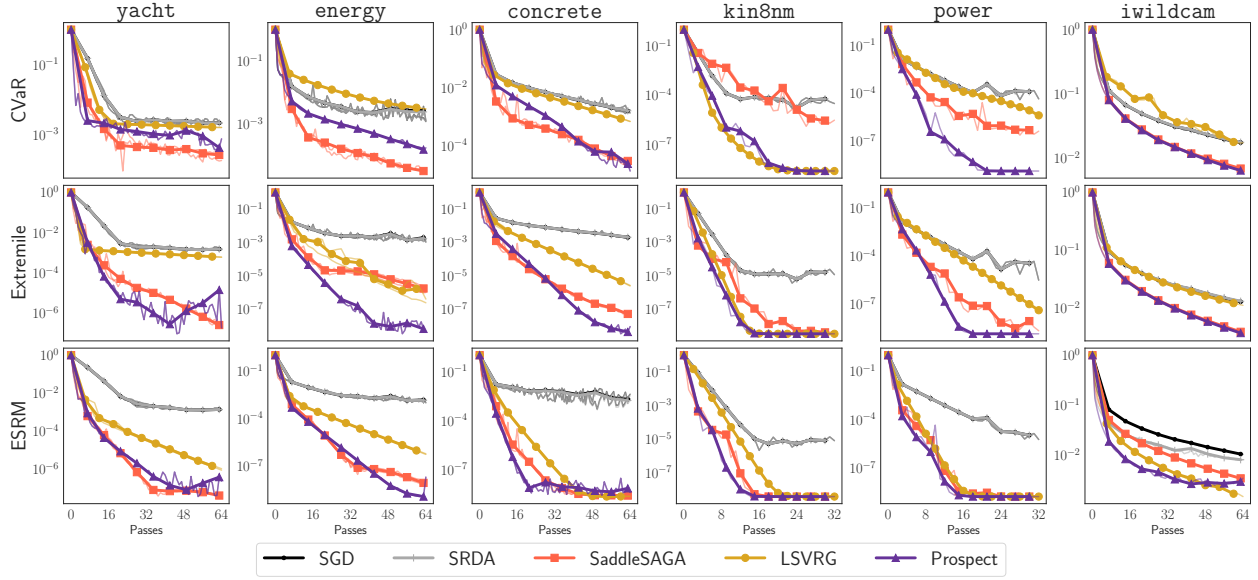


Figure A.2: **No Shift Cost Settings.** Each row represents a different spectral risk objective with $\nu = 0$ (instead of $\nu = 1$) while each column represents a different datasets. Suboptimality (2.49) is measured on the y -axis while the x -axis measures the total number of gradient evaluations made divided by n , i.e., the number of passes through the training set.

Removing Shift Cost A relevant setting is the low or no shift cost regime ($\nu = 0$), as this allows the adversary to make arbitrary distribution shifts (while still constrained to $\mathcal{Q}(\sigma)$). Figure A.2 displays these curves for this no-cost experiment. When $\nu = 0$, the optimization problem can equivalently be written as

$$\min_{\boldsymbol{\theta} \in \mathbb{R}^d} \left[\max_{\mathbf{q} \in \mathcal{Q}(\sigma)} \langle \mathbf{q}, \ell(\boldsymbol{\theta}) \rangle + \frac{\mu}{2} \|\boldsymbol{\theta}\|_2^2 = \sum_{i=1}^n \sigma_i \ell_{(i)}(\boldsymbol{\theta}) + \frac{\mu}{2} \|\boldsymbol{\theta}\|_2^2 \right].$$

In this case, we always have that $\mathbf{q}(\mathbf{l}) = (\sigma_{\pi^{-1}(1)}, \dots, \sigma_{\pi^{-1}(n)})$, where π sorts \mathbf{l} . Here, $\boldsymbol{\theta}$ is chosen to optimize a linear combination of order statistics of the losses. In the low shift cost settings, performance trends are qualitatively similar to those seen from $\nu = 1$. Interestingly, for the no-cost setting, SaddleSAGA, LSVRG, and Prospect seem to converge linearly empirically even without smoothness of the objective.

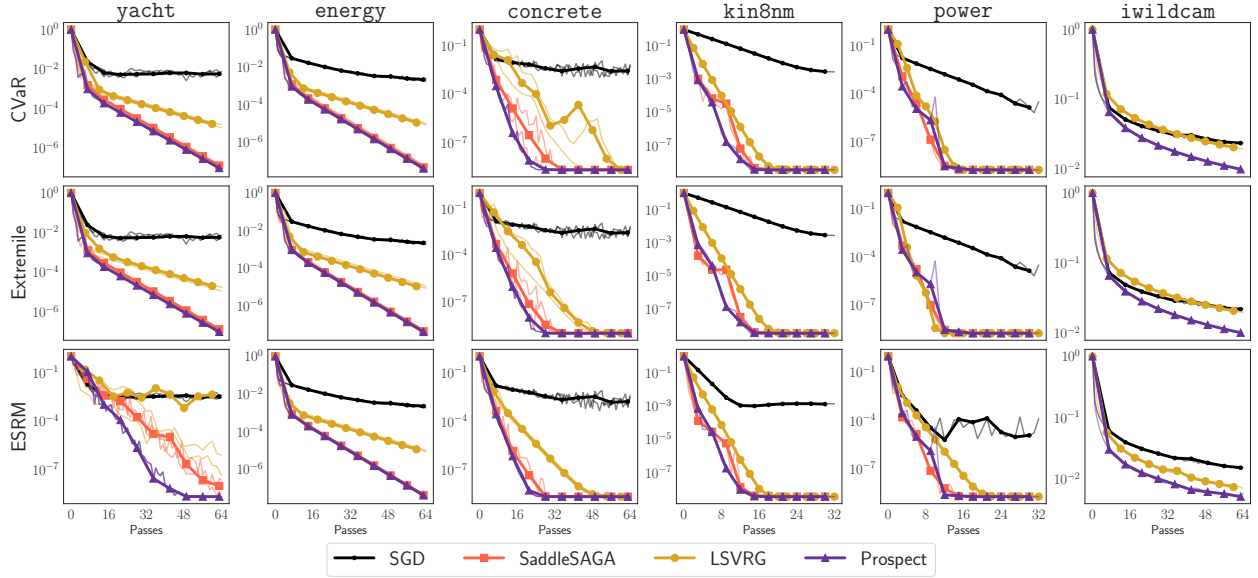


Figure A.3: **Reduced ℓ_2 -regularization settings** ($\mu = 1/(10n)$). Each row represents a different spectral risk objective with $\mu = 1/(10n)$ (instead of $\mu = 1/n$) while each column represents a different dataset. Suboptimality (2.49) is measured on the y -axis while the x -axis measures the total number of gradient evaluations made divided by n , i.e., the number of passes through the training set.

Lowering Regularization Next, we decrease the ℓ_2 -regularization from $\mu = 1/n$ to $\mu = 1/(10n)$. These settings are plotted in Figure A.3. Performance rankings among methods reflect those of the original parameters. For five of the six datasets, that is **yacht**, **energy**, **concrete**, **kin8nm**, and **power**, the regression tasks involve optimizing the squared error. This function is already strongly convex, with a constant depending on the smallest eigenvalue of the empirical second-moment matrix. When assuming that the input data vectors are bounded, this function is also G -Lipschitz. Thus, if the problem is already well-conditioned, we may observe similar behavior even at negligible regularization ($\mu = 5 \cdot 10^{-7}$ for **iwildcam**, for example).

Appendix B

APPENDIX TO CHAPTER 4

B.1 Linear Operators and Variance Reduction

This section is dedicated to establishing the variance reduction result in Corollary 4.3.1 by employing properties of the conditional mean operators introduced in Section 4.3. In the first part, we establish Proposition B.1.1, the singular value decomposition that defines the quantities appearing in Corollary 4.3.1. In the second part, we quantify the difference between σ_0^2 and σ_k^2 for even and odd iterations of k .

B.1.1 Singular Value Decomposition

Recall the conditional mean operators μ_X and μ_Z from Section 4.3,

$$[\mu_X h](\mathbf{x}) := \mathbb{E}[h(X, Z)|X](\mathbf{x}) \text{ and } [\mu_Z h](\mathbf{z}) := \mathbb{E}[h(X, Z)|Z](\mathbf{z}),$$

with the corresponding debiasing (a.k.a. centering) operators defined by $\mathcal{C}_X = I - \mu_X$ and $\mathcal{C}_Z = I - \mu_Z$.

Proposition B.1.1. *There exists a basis $\{\alpha_j\}_{j=1}^m$ of $\mathbf{L}^2(P_X)$, a basis $\{\beta_j\}_{j=1}^m$ of $\mathbf{L}^2(P_Z)$, and real values $\{s_j\}_{j=1}^m$, which satisfy:*

$$\mu_Z \alpha_j = s_j \beta_j \text{ and } \mu_X \beta_j = s_j \alpha_j \text{ for } j \in \{1, \dots, m\}, \quad (\text{B.1})$$

$\alpha_1 = \mathbf{1}_X$, $\beta_1 = \mathbf{1}_Z$, $s_1 = 1$ and s_j is non-negative and non-increasing in j .

Proof. When μ_X is restricted to $\mathbf{L}^2(P_Z)$ and μ_Z is restricted to $\mathbf{L}^2(P_X)$, these operators are

in fact adjoint in $\mathbf{L}^2(P)$, as by the tower property we have the relation

$$\langle f, \mu_X g \rangle_{\mathbf{L}^2(P_X)} = \mathbb{E} [f(X) \mathbb{E} [g(Z)|X]] = \mathbb{E} [\mathbb{E} [f(X)|Z] g(Z)] = \langle \mu_Z f, g \rangle_{\mathbf{L}^2(P_Z)}.$$

Since $\mu_Z : \mathbf{L}^2(P_X) \rightarrow \mathbf{L}^2(P_Z)$ is a compact linear operator, by [Gohberg et al. \[1990, Section IV.1 Theorem 1.1\]](#) and [Gohberg et al. \[1990, Section IV.1 Corollary 1.2\]](#), we have that μ_Z admits a singular value decomposition satisfying (B.1). Next, we show that $s_1 \leq 1$ and that $\mathbf{1}_X$ is an eigenvector of $\mu_X \mu_Z : \mathbf{L}^2(P_X) \rightarrow \mathbf{L}^2(P_X)$ with eigenvalue 1, which confirms that $s_1 = 1$ and $\alpha_1 = \mathbf{1}_X$ by the definition of singular values (arguing symmetrically achieves $\beta_1 = \mathbf{1}_Z$). By the variational representation of singular values [[Gohberg et al., 1990, Section IV.1 Equation \(2\)](#)], we have that

$$\sup_{f: \|f\|_{\mathbf{L}^2(P_X)}=1} \|\mu_Z f\|_{\mathbf{L}^2(P_Z)} = s_1.$$

Consider any $f \in \mathbf{L}^2(P_X)$ such that $\|f\|_{\mathbf{L}^2(P_X)} = 1$. Define the conditional probability $P_{X|Z}(\mathbf{x}|\mathbf{z}) = P(\mathbf{x}, \mathbf{z})/P_Z(\mathbf{z})$ which is well-defined by assumption. Then, by the Cauchy-Schwarz inequality in $\mathbf{L}^2(P_{X|Z})$,

$$\begin{aligned} \|\mu_Z f\|_{\mathbf{L}^2(P_Z)}^2 &= \sum_{\mathbf{z} \in \mathcal{Z}} \left(\sum_{\mathbf{x} \in \mathcal{X}} f(\mathbf{x}) P_{X|Z}(\mathbf{x}|\mathbf{z}) \right)^2 P_Z(\mathbf{z}) \\ &\leq \sum_{\mathbf{z} \in \mathcal{Z}} \sum_{\mathbf{x} \in \mathcal{X}} f^2(\mathbf{x}) P_{X|Z}(\mathbf{x}|\mathbf{z}) P_Z(\mathbf{z}) \\ &= \sum_{\mathbf{x} \in \mathcal{X}} f^2(\mathbf{x}) \sum_{\mathbf{z} \in \mathcal{Z}} P(\mathbf{x}, \mathbf{z}) \\ &= \|f\|_{\mathbf{L}^2(P_X)}^2 = 1. \end{aligned}$$

This proves that $s_1 \leq 1$. For equality, notice that $\mu_X \mu_Z \mathbf{1}_X = \mu_X \mathbf{1}_Z = \mathbf{1}_X$, completing the proof. \square

B.1.2 Proof of Main Results

From Proposition B.1.1, we establish two bases $\{\alpha_j\}_{j=1}^m$ and $\{\beta_j\}_{j=1}^m$ of $\mathbf{L}^2(P_X)$ and $\mathbf{L}^2(P_Z)$, respectively. These bases span the range of the operators μ_X and μ_Z . We will consider the repeated application of the operator $\mathcal{C}_Z\mathcal{C}_X$, a sequence of two centering operations on some function $h \in \mathbf{L}^2(P)$, and compare

$$\mathbb{E} [((\mathcal{C}_Z\mathcal{C}_X)^t \bar{h})^2] \text{ against } \mathbb{E} [\bar{h}^2]$$

for $\bar{h} = h - \mathbb{E}_P[h]$. We establish the main result by measuring the reduction in variance from a single application, in terms of the coordinates of the function of interest on each of the two subspaces. We will then observe how these coordinates change iteration-to-iteration to give the final result.

Lemma B.1.1. *For any $h \in \mathbf{L}^2(P)$ such that $\mathbb{E}_P[h] = 0$, let*

$$\mu_X h = \sum_{j=1}^m u_j \alpha_j \text{ and } \mu_Z h = \sum_{j=1}^m v_j \beta_j.$$

Then, we have that

$$\mathbb{E} [(\mathcal{C}_Z\mathcal{C}_X h)^2] = \mathbb{E} [h^2] - \sum_{j=2}^m u_j^2 - \sum_{j=2}^m (v_j - s_j u_j)^2.$$

Proof. By orthogonality, we have that

$$\begin{aligned} \mathbb{E} [(\mathcal{C}_Z\mathcal{C}_X h)^2] &= \mathbb{E} [((I - \mu_Z)\mathcal{C}_X h)^2] \\ &= \mathbb{E} [(\mathcal{C}_X h)^2] - 2\mathbb{E} [(\mathcal{C}_X h)(\mu_Z \mathcal{C}_X h)] + \mathbb{E} [(\mu_Z \mathcal{C}_X h)^2] \\ &= \mathbb{E} [(\mathcal{C}_X h)^2] - 2P_Z((\mu_Z \mathcal{C}_X h)^2) + P_Z((\mu_Z \mathcal{C}_X h)^2) \\ &= \mathbb{E} [(\mathcal{C}_X h)^2] - P_Z((\mu_Z \mathcal{C}_X h)^2) \\ &= \mathbb{E} [h^2] - P_X((\mu_X h)^2) - P_Z((\mu_Z \mathcal{C}_X h)^2). \end{aligned}$$

Because $P(h) = 0$, it holds by the tower property of conditional expectation that $P_X(\mu_X h) =$

0, which implies that

$$u_1 = \langle \mu_X h, \alpha_1 \rangle_{\mathbf{L}^2(P_X)} = 0 \implies P_X((\mu_X h)^2) = \sum_{j=2}^m u_j^2.$$

For the second term, observe that $P_X(\mathcal{C}_X h) = 0$, so it holds by the tower property that $P_Z(\mu_Z \mathcal{C}_X h) = 0$, so

$$P_Z((\mu_Z \mathcal{C}_X h)^2) = \sum_{j=2}^m \left(\langle \mu_Z \mathcal{C}_X h, \beta_j \rangle_{\mathbf{L}^2(P_Z)} \right)^2.$$

Next, we compute the term in the square by applying Proposition B.1.1:

$$\begin{aligned} \langle \mu_Z \mathcal{C}_X h, \beta_j \rangle_{\mathbf{L}^2(P_Z)} &= \langle \mu_Z h, \beta_j \rangle_{\mathbf{L}^2(P_Z)} - \langle \mu_Z \mu_X h, \beta_j \rangle_{\mathbf{L}^2(P_Z)} \\ &= v_j - \left\langle \mu_Z \sum_{k=1}^m u_k \alpha_k, \beta_j \right\rangle_{\mathbf{L}^2(P_Z)} \\ &= v_j - \left\langle \sum_{k=1}^m u_k s_k \beta_k, \beta_j \right\rangle_{\mathbf{L}^2(P_Z)} \\ &= v_j - s_j u_j, \end{aligned}$$

which completes the proof. \square

Lemma B.1.1 ensures that we have a reduction on each iteration, with a formula that depends on the coordinates of the function on each subspace. Because these coordinates change every iteration, we track them in the next lemma. Define $h_0 = \bar{h}$ and $h_{t+1} = (\mathcal{C}_Z \mathcal{C}_X) h_t$, along with the constants $\{u_{t,j}\}_{j=1}^m$ and $\{v_{t,j}\}_{j=1}^m$ given by

$$\mu_X h_t = \sum_{j=1}^m u_{t,j} \alpha_j \text{ and } \mu_Z h_t = \sum_{j=1}^m v_{t,j} \beta_j.$$

We have the following.

Lemma B.1.2. *For all $t \geq 0$, it holds that*

$$\begin{aligned} u_{t+1,j} &= s_j^2 u_{t,j} - s_j v_{t,j}, \\ v_{t+1,j} &= 0. \end{aligned}$$

Proof. Fix any $j \in [m]$, and use Proposition B.1.1 to write

$$\begin{aligned} u_{t+1,j} &= \langle \mu_X \mathcal{C}_Z \mathcal{C}_X h_t, \alpha_j \rangle_{\mathbf{L}^2(P_X)} \\ &= \langle \mu_X (I - \mu_X - \mu_Z + \mu_Z \mu_X) h_t, \alpha_j \rangle_{\mathbf{L}^2(P_X)} \\ &= \langle \mu_X \mu_Z \mu_X h_t, \alpha_j \rangle_{\mathbf{L}^2(P_X)} - \langle \mu_X \mu_Z h_t, \alpha_j \rangle_{\mathbf{L}^2(P_X)} \\ &= \left\langle \mu_X \mu_Z \sum_{k=1}^m u_{t,k} \alpha_k, \alpha_j \right\rangle_{\mathbf{L}^2(P_X)} - \left\langle \mu_X \sum_{k=1}^m v_{t,k} \beta_k, \alpha_j \right\rangle_{\mathbf{L}^2(P_X)} \\ &= s_j^2 u_{t,j} - s_j v_{t,j}, \end{aligned}$$

which proves the first part of the claim. For the second part, note that $\mu_Z \mathcal{C}_Z = 0$, so $\langle \mu_Z \mathcal{C}_Z \mathcal{C}_X h_t, \alpha_j \rangle_{\mathbf{L}^2(P_Z)} = 0$. \square

Using Lemma B.1.1 and Lemma B.1.2, we can simply accumulate the reduction incurred on every iteration.

Proposition B.1.2. *Define the constants $(u_j)_{j=1}^m$ and $(v_j)_{j=1}^m$ by*

$$\mu_X \bar{h} = \sum_{j=1}^m u_j \alpha_j \text{ and } \mu_Z \bar{h} = \sum_{j=1}^m v_j \beta_j.$$

Then, we may quantify the variance reduction achieved by $t+1$ iterations of the $\mathcal{C}_Z \mathcal{C}_X$ operator as

$$\begin{aligned} \mathbb{E} [\bar{h}^2] - \mathbb{E} [((\mathcal{C}_Z \mathcal{C}_X)^{t+1} \bar{h})^2] &= \sum_{j=2}^m \left\{ u_j^2 + (v_j - s_j u_j)^2 \left[1 + \frac{s_j^2 (1 - s_j^{4t})}{1 - s_j^2} \right] \right\} \\ &\rightarrow \sum_{j=2}^m \left[u_j^2 + \frac{(v_j - s_j u_j)^2}{1 - s_j^2} \right] \end{aligned}$$

as $t \rightarrow \infty$.

Proof. Apply Lemma B.1.1 $(t + 1)$ -times so that

$$\begin{aligned}\mathbb{E} [((\mathcal{C}_Z \mathcal{C}_X)^{t+1} \bar{h})^2] &= \mathbb{E} [\bar{h}^2] - \sum_{j=2}^m \sum_{\tau=0}^t [(1 + s_j^2) u_{\tau,j}^2 + v_{\tau,j}^2 - 2s_j u_{\tau,j} v_{\tau,j}] \\ &= \mathbb{E} [\bar{h}^2] - \sum_{j=2}^m \left[v_{0,j}^2 - 2s_j u_{0,j} v_{0,j} + \sum_{\tau=0}^t (1 + s_j^2) u_{\tau,j}^2 \right]\end{aligned}$$

as by Lemma B.1.2, we have that $v_{\tau,j} = 0$ for $\tau > 0$. Next, we unroll the definition of $u_{\tau,j}$ so that

$$\begin{aligned}u_{\tau,j} &= s_j^2 u_{\tau-1,j} - s_j v_{\tau-1,j} \\ &= s_j^2 (s_j^2 u_{\tau-2,j} - s_j v_{\tau-2,j}) - s_j v_{\tau-1,j} \\ &= s_j^{2\tau-2} (s_j^2 u_{0,j} - s_j v_{0,j})\end{aligned}$$

for $\tau > 0$, yielding

$$\begin{aligned}\mathbb{E} [\bar{h}^2] - \mathbb{E} [((\mathcal{C}_Z \mathcal{C}_X)^{t+1} \bar{h})^2] &= \sum_{j=2}^m \left[u_{0,j}^2 + (v_{0,j} - s_j u_{0,j})^2 + (1 + s_j^2) (s_j^2 u_{0,j} - s_j v_{0,j})^2 \sum_{\tau=1}^t (s_j^4)^{\tau-1} \right] \\ &= \sum_{j=2}^m \left[u_{0,j}^2 + (v_{0,j} - s_j u_{0,j})^2 + (1 + s_j^2) (s_j^2 u_{0,j} - s_j v_{0,j})^2 \sum_{\tau=0}^{t-1} (s_j^4)^{\tau} \right] \\ &= \sum_{j=2}^m \left[u_{0,j}^2 + (v_{0,j} - s_j u_{0,j})^2 + \frac{s_j^2 (1 + s_j^2) (v_{0,j} - s_j u_{0,j})^2 (1 - s_j^{4t})}{1 - s_j^4} \right] \\ &= \sum_{j=2}^m \left[u_{0,j}^2 + (v_{0,j} - s_j u_{0,j})^2 + \frac{s_j^2 (v_{0,j} - s_j u_{0,j})^2 (1 - s_j^{4t})}{1 - s_j^2} \right].\end{aligned}$$

Substitute $u_{0,j} = u_j$ and $v_{0,j} = v_j$ to complete the proof. \square

We also present the corresponding result for k odd. The proof follows similarly by repeated application of the operator $\mathcal{C}_Z \mathcal{C}_X$. However, the iterations will be compared to

$\sigma_1^2 = \mathbb{E}_P [(\mathcal{C}_X \bar{h})^2]$, as we consider $\mathcal{C}_X \bar{h}$ as the “first” iteration to this process.

Proposition B.1.3. *Define the constants $(u_j)_{j=1}^m$ by*

$$\mu_Z \mathcal{C}_X \bar{h} = \sum_{j=1}^m u_j \beta_j.$$

Then, we may quantify the variance reduction achieved by $t+1$ iterations of the $\mathcal{C}_X \mathcal{C}_Z$ operator as

$$\begin{aligned} \mathbb{E} [(\mathcal{C}_X \bar{h})^2] - \mathbb{E} [(\mathcal{C}_X \mathcal{C}_Z)^{t+1} \mathcal{C}_X \bar{h})^2] &= \sum_{j=2}^m \left\{ u_j^2 + (s_j u_j)^2 \left[1 + \frac{s_j^2(1 - s_j^{4t})}{1 - s_j^2} \right] \right\} \\ &\rightarrow \sum_{j=2}^m \left(\frac{1 + s_j^2}{1 - s_j^2} \right) u_j^2 \end{aligned}$$

as $t \rightarrow \infty$.

In order to have full monotonicity, we also need that $\sigma_0^2 \geq \sigma_1^2$. This follows by orthogonality, as

$$\sigma_0^2 = \mathbb{E} [\bar{h}^2] = \mathbb{E} [(\mathcal{C}_X \bar{h})^2] + \mathbb{E} [(\mu_X \bar{h})^2] = \sigma_1^2 + \mathbb{E} [(\mu_X \bar{h})^2] \geq \sigma_1^2. \quad (\text{B.2})$$

Thus, we can combine Proposition B.1.3 and (B.2) to fully quantify the relationship between σ_0^2 and σ_k^2 for k odd.

B.2 Information Projections

This section is dedicated to deriving three representations of the balancing procedure as projections in various statistical divergences, as shown in Figure 4.1.

We consider two sets of probability measures denoted by $\Pi_X = \{Q : Q_X = P_X\}$ and $\Pi_Z = \{Q : Q_Z = P_Z\}$. The marginal matching steps are written as projections in terms of a statistical divergence D (precisely, an f -divergence) in the form

$$\frac{P_X}{P_{n,X}^{(k-1)}} \otimes P_n^{(k-1)} = \arg \min_{Q \in \Pi_X} D(Q \| P_n^{(k-1)}), \quad \frac{P_Z}{P_{n,Z}^{(k-1)}} \otimes R = \arg \min_{Q \in \Pi_Z} D(Q \| P_n^{(k-1)}).$$

We provide the derivations for three common choices of D : Kullback-Leibler (KL), reverse KL, and χ^2 . Using this viewpoint, and simply assuming the positivity of the marginal measures P_X and P_Z , we derive an upper bound in Proposition B.2.5 that is *constant* in k . This is an improvement over the recent work of Albertus and Berthet [2019], in which they show an upper bound that scales *exponentially* in k .

The KL representation will be used in the proof of Proposition B.2.5, which (recalling the sequence $(P_n^{(k)})_{k \geq 1}$ from (4.2)), controls the error between $P_{n,Z}^{(k)}$ and P_Z for k odd and $P_{n,X}^{(k)}$ and P_X for k even.

B.2.1 Balancing as Information Projections

The arguments for three information divergences (KL, reverse KL, and χ^2) are contained in the following propositions.

Proposition B.2.1 (Projection in KL-Divergence). *Assume that $P_X \ll R_X$ and $P_Z \ll R_Z$, and define*

$$Q^* := \arg \min_{Q \in \Pi_X} \text{KL}(Q \| R), \quad P^* := \arg \min_{Q \in \Pi_Z} \text{KL}(Q \| R). \quad (\text{B.3})$$

Then, it holds that

$$Q^*(\mathbf{x}, \mathbf{z}) = \begin{cases} P_X(\mathbf{x}) R_{Z|X}(\mathbf{z}|\mathbf{x}) & \text{if } R_X(\mathbf{x}) > 0 \\ 0 & \text{if } R_X(\mathbf{x}) = 0 \end{cases} \quad (\text{B.4})$$

and

$$P^*(\mathbf{x}, \mathbf{z}) = \begin{cases} P_Z(\mathbf{z}) R_{X|Z}(\mathbf{x}|\mathbf{z}) & \text{if } R_Z(\mathbf{z}) > 0 \\ 0 & \text{if } R_Z(\mathbf{z}) = 0 \end{cases}. \quad (\text{B.5})$$

Proof. In the case that $Q(\mathbf{x}, \mathbf{z}) = 0$, we apply the convention that $0 \log 0 = 0$. Consider the

case Q^* , the projection of R onto Π_X . Write

$$\begin{aligned}
\text{KL}(Q\|R) &= \sum_{\mathbf{x} \in \mathcal{X}} \sum_{\mathbf{z} \in \mathcal{Z}} Q(\mathbf{x}, \mathbf{z}) \log \frac{Q_{Z|X}(\mathbf{z}|\mathbf{x})Q_X(\mathbf{x})}{R_{Z|X}(\mathbf{z}|\mathbf{x})R_X(\mathbf{x})} \\
&= \sum_{\mathbf{x} \in \mathcal{X}} Q_X(\mathbf{x}) \left[\sum_{\mathbf{z} \in \mathcal{Z}} Q_{Z|X}(\mathbf{z}|\mathbf{x}) \log \frac{Q_{Z|X}(\mathbf{z}|\mathbf{x})Q_X(\mathbf{x})}{R_{Z|X}(\mathbf{z}|\mathbf{x})R_X(\mathbf{x})} \right] \\
&= \sum_{\mathbf{x} \in \mathcal{X}} Q_X(\mathbf{x}) \left[\sum_{\mathbf{z} \in \mathcal{Z}} Q_{Z|X}(\mathbf{z}|\mathbf{x}) \log \frac{Q_{Z|X}(\mathbf{z}|\mathbf{x})}{R_{Z|X}(\mathbf{z}|\mathbf{x})} + \sum_{\mathbf{z} \in \mathcal{Z}} Q_{Z|X}(\mathbf{z}|\mathbf{x}) \log \frac{Q_X(\mathbf{x})}{R_X(\mathbf{x})} \right] \\
&= \sum_{\mathbf{x} \in \mathcal{X}} Q_X(\mathbf{x}) \left[\sum_{\mathbf{z} \in \mathcal{Z}} Q_{Z|X}(\mathbf{z}|\mathbf{x}) \log \frac{Q_{Z|X}(\mathbf{z}|\mathbf{x})}{R_{Z|X}(\mathbf{z}|\mathbf{x})} \right] + \sum_{\mathbf{x} \in \mathcal{X}} Q_X(\mathbf{x}) \log \frac{Q_X(\mathbf{x})}{R_X(\mathbf{x})} \\
&= \sum_{\mathbf{x} \in \mathcal{X}} \textcolor{blue}{Q}_X(\mathbf{x}) \text{KL}(\textcolor{blue}{Q}_{Z|X}(\cdot|\mathbf{x})\|R_{Z|X}(\cdot|\mathbf{x})) + \text{KL}(\textcolor{blue}{Q}_X\|R_X) \\
&= \sum_{\mathbf{x} \in \mathcal{X}} \textcolor{red}{P}_X(\mathbf{x}) \text{KL}(\textcolor{blue}{Q}_{Z|X}(\cdot|\mathbf{x})\|R_{Z|X}(\cdot|\mathbf{x})) + \text{KL}(\textcolor{red}{P}_X\|R_X),
\end{aligned}$$

where the last line is due to the marginal constraint $Q \in \Pi_X$. For the above to be well defined, we need that $P_X \ll R_X$ so that $\text{KL}(P_X\|R_X) < +\infty$. The above is minimized when $Q_{Z|X}(\mathbf{z}|\mathbf{x}) = R_{Z|X}(\mathbf{z}|\mathbf{x})$ for all $(\mathbf{x}, \mathbf{z}) \in \mathcal{X} \times \mathcal{Z}$ such that $Q_X(\mathbf{x}) = P_X(\mathbf{x}) > 0$. The case of P^* follows analogously when using that $P_Z \ll R_Z$. \square

Proposition B.2.2 (Projection in Reverse KL-Divergence). *Assume that $P_Z \ll R_X$ and $P_Z \ll R_Z$, and define*

$$Q^* := \arg \min_{Q \in \Pi_X} \text{KL}(R\|Q), \quad P^* := \arg \min_{Q \in \Pi_Z} \text{KL}(R\|Q). \quad (\text{B.6})$$

Then, it holds that

$$Q^*(\mathbf{x}, \mathbf{z}) = \begin{cases} P_X(\mathbf{x})R_{Z|X}(\mathbf{z}|\mathbf{x}) & \text{if } R_X(\mathbf{x}) > 0 \\ 0 & \text{if } R_X(\mathbf{x}) = 0 \end{cases} \quad (\text{B.7})$$

and

$$P^*(\mathbf{x}, \mathbf{z}) = \begin{cases} P_Z(\mathbf{z})R_{X|Z}(\mathbf{x}|\mathbf{z}) & \text{if } R_Z(\mathbf{z}) > 0 \\ 0 & \text{if } R_Z(\mathbf{z}) = 0 \end{cases}. \quad (\text{B.8})$$

Proof. In the case that $R(\mathbf{x}, \mathbf{z}) = 0$, we apply the convention that $0 \log 0 = 0$. Note that minimizing $\text{KL}(R\|Q)$ over Q is equivalent to minimizing $-\sum_{\mathbf{x}, \mathbf{y}} R(\mathbf{x}, \mathbf{z}) \log Q(\mathbf{x}, \mathbf{z})$ (i.e. the cross entropy). Consider the case Q^* , the projection of R onto Π_X . Because $R \ll Q$ for $\text{KL}(R\|Q) < +\infty$ to hold, we have that $R(\mathbf{x}) > 0 \implies Q(\mathbf{x}) > 0$, so that $Q_{Z|X}(\mathbf{z}|\mathbf{x})$ is well-defined. Write

$$\begin{aligned} & - \sum_{\mathbf{x}, \mathbf{y}} R(\mathbf{x}, \mathbf{z}) \log Q(\mathbf{x}, \mathbf{z}) \\ &= - \sum_{\mathbf{x} \in \mathcal{X}} R_X(\mathbf{x}) \log Q_X(\mathbf{x}) - \sum_{\mathbf{x} \in \mathcal{X}} R(\mathbf{x}) \sum_{\mathbf{z} \in \mathcal{Z}} R_{Z|X}(\mathbf{z}|\mathbf{x}) \log Q_{Z|X}(\mathbf{z}|\mathbf{x}) \\ &= - \sum_{\mathbf{x} \in \mathcal{X}} R_X(\mathbf{x}) \log P_X(\mathbf{x}) + \sum_{\mathbf{x} \in \mathcal{X}} R_X(\mathbf{x}) \left[- \sum_{\mathbf{z} \in \mathcal{Z}} R_{Z|X}(\mathbf{z}|\mathbf{x}) \log Q_{Z|X}(\mathbf{z}|\mathbf{x}) \right]. \end{aligned}$$

The first term does not depend on Q due to the marginal constraint $Q \in \Pi_X$. The second term is the expectation of the cross entropy from $R_{Z|X}$ to $Q_{Z|X}$ over R_X , which is minimized if $R_{Z|X} = Q_{Z|X}$. We have specified $Q_{Z|X}$ and Q_X , completing the proof. \square

The projection result for χ^2 -divergence requires a few more intermediate steps. Let $\mathbf{1}$ denote the function that is identically equal to 1. Consider the following optimization problem, which is the subject of the subsequent lemmas:

$$\min_{\zeta \in \mathcal{A}_X} \|\mathbf{1} - \zeta\|_{\mathbf{L}^2(R)}^2, \quad (\text{B.9})$$

where

$$\mathcal{A}_X := \left\{ f : \mathcal{X} \times \mathcal{Z} \rightarrow \mathbb{R} \text{ satisfying } \sum_{\mathbf{z} \in \mathcal{Z}} f(\mathbf{x}, \mathbf{z}) R(\mathbf{x}, \mathbf{z}) = P_X(\mathbf{x}) \text{ for any } \mathbf{x} \in \mathcal{X} \right\}.$$

Lemma B.2.1. *Assume that $P_X \ll R_X$, and define The problem (B.9) is feasible, and its solution can be written as*

$$\zeta^\star = \mathcal{C}_X^R(\mathbf{1} - f) + f$$

for any $f \in \mathbf{L}^2(R)$, where the linear operator \mathcal{C}_X^R is specified by

$$[\mathcal{C}_X^R g](\mathbf{x}, \mathbf{z}) = g(\mathbf{x}, \mathbf{z}) - \sum_{\mathbf{z}' \in \mathcal{Z}} g(\mathbf{x}, \mathbf{z}') R_{\mathcal{Z}|X}(\mathbf{z}'|\mathbf{x}).$$

Proof. First, we establish feasibility by letting

$$f(\mathbf{x}, \mathbf{z}) := \begin{cases} P_X(\mathbf{x})/R_X(\mathbf{x}) & \text{if } R_X(\mathbf{x}) > 0 \\ 1 & \text{otherwise} \end{cases}.$$

This function does not depend on the second input \mathbf{z} . Because we assumed that $P_X \ll R_X$, we have that the terms of $f(\mathbf{x}, \mathbf{z})$ for which $R_X(\mathbf{x}) = 0$ do not affect whether $\sum_{\mathbf{z} \in \mathcal{Z}} f(\mathbf{x}, \mathbf{z}) R(\mathbf{x}, \mathbf{z}) = P_X(\mathbf{x})$, because $P_X(\mathbf{x}) = 0$ in these cases. In the remainder of this proof, we will show that (B.9) is an affine projection problem, and find its solution by converting it to a subspace projection problem. Indeed, consider $f_1, \dots, f_r \in \mathcal{A}_X$, and $\alpha_1, \dots, \alpha_r \in \mathbb{R}$ such that $\sum_{j=1}^r \alpha_j = 1$. Then,

$$\sum_{\mathbf{z} \in \mathcal{Z}} \left[\sum_{j=1}^r \alpha_j f_j(\mathbf{x}, \mathbf{z}) \right] \cdot R(\mathbf{x}, \mathbf{z}) = \sum_{j=1}^r \alpha_j \left[\sum_{\mathbf{z} \in \mathcal{Z}} f_j(\mathbf{x}, \mathbf{z}) R(\mathbf{x}, \mathbf{z}) \right] = P_X(\mathbf{x}),$$

indicating that $\sum_{j=1}^r \alpha_j f_j(\mathbf{x}, \mathbf{z}) \in \mathcal{A}_X$ and \mathcal{A}_X is an affine subset of $\mathbf{L}^2(R)$. Define

$$\mathcal{S}_X := \left\{ g : \mathcal{X} \times \mathcal{Z} \rightarrow \mathbb{R} \text{ satisfying } \sum_{\mathbf{z} \in \mathcal{Z}} g(\mathbf{x}, \mathbf{z}) R(\mathbf{x}, \mathbf{z}) = 0 \text{ for any } \mathbf{x} \in \mathcal{X} \right\}.$$

Then, for any $f \in \mathcal{A}_X$, we have that $g \in \mathcal{S}_X$ if and only if $g + f \in \mathcal{A}_X$. Taking any $f \in \mathcal{A}_X$, letting ϕ^\star be the solution of

$$\min_{\phi \in \mathcal{S}_X} \|\mathbf{1} - f - \phi\|_{\mathbf{L}^2(R)}^2, \quad (\text{B.10})$$

we will have that $\phi^* + f$ will be the solution of (B.9). The remainder of the proof is showing that $\phi^* = \mathcal{C}_X^R(\mathbf{1} - f)$.

First, define the operator μ_X^R by $[\mu_X g](\mathbf{x}, \mathbf{z}) = \sum_{\mathbf{z}' \in \mathcal{Z}} g(\mathbf{x}, \mathbf{z}') R_{Z|X}(\mathbf{z}'|\mathbf{x})$, and note (by factoring out $R_X(\mathbf{x})$) that $g \in \mathcal{S}_X$ if and only if $\mu_X^R g = 0$. In addition, $\mu_X^R g$ is linear and idempotent as $\mu_X^R \mu_X^R g = \mu_X^R g$, so it is a projection operator in $\mathbf{L}^2(R)$. Thus, \mathcal{S}_X is the orthogonal complement of $\text{range}(\mu_X^R)$, and the solution of (B.10) is given by $(I - \mu_X^R)(\mathbf{1} - f) = \mathcal{C}_X^R(\mathbf{1} - f)$, because $\mathcal{C}_X^R = I - \mu_X^R$. The claim is proved. \square

Lemma B.2.2. *Assume that $P_X \ll R_X$. Define*

$$Q^* := \arg \min_{Q \in \Pi_X} \chi^2(Q \| R). \quad (\text{B.11})$$

and let ζ^* be the solution of problem (B.9). Then,

$$Q^*(\mathbf{x}, \mathbf{z}) = \zeta^*(\mathbf{x}, \mathbf{z}) R(\mathbf{x}, \mathbf{z}) = \begin{cases} P_X(\mathbf{x}) R_{Z|X}(\mathbf{z}|\mathbf{x}) & \text{if } R_X(\mathbf{x}) > 0 \\ 0 & \text{if } R_X(\mathbf{x}) = 0 \end{cases}. \quad (\text{B.12})$$

Proof. First, by reparametrizing the problem (B.11) as finding ζ such that $Q(\mathbf{x}, \mathbf{z}) = \zeta(\mathbf{x}, \mathbf{z}) R(\mathbf{x}, \mathbf{z})$, we can compute its solution by solving

$$\min_{\zeta \in \mathcal{A}_X, \zeta \geq 0} \|\mathbf{1} - \zeta\|_{\mathbf{L}^2(R)}^2, \quad (\text{B.13})$$

Notice that we also have a non-negativity constraint, as opposed to (B.9). If ζ^* solves (B.9) and happens to be non-negative, then we have that ζ^* solves (B.13) as well and the first equality of (B.12) is satisfied by definition. We show the second equality of (B.12) by direct computation, which also establishes the non-negativity of ζ^* simultaneously.

Apply Lemma B.2.1 with

$$f(\mathbf{x}, \mathbf{z}) := \begin{cases} P_X(\mathbf{x})/R_X(\mathbf{x}) & \text{if } R_X(\mathbf{x}) > 0 \\ 1 & \text{otherwise} \end{cases}.$$

so that

$$\begin{aligned}
\zeta^*(\mathbf{x}, \mathbf{z}) &= \mathcal{C}_X^R(\mathbf{1} - f)(\mathbf{x}, \mathbf{z}) + f(\mathbf{x}, \mathbf{z}) \\
&= \left[\sum_{z \in \mathcal{Z}} f(\mathbf{x}, z) R_{Z|X}(z|\mathbf{x}) - f(\mathbf{x}, \mathbf{z}) \right] + f(\mathbf{x}, \mathbf{z}) \\
&= f(\mathbf{x}, y')
\end{aligned}$$

for any $\mathbf{z}' \in \mathcal{Z}$. Thus, the likelihood ratio of Q^* with respect to R is a marginal reweighting. Accordingly,

$$Q^*(\mathbf{x}, \mathbf{z}) = \zeta^*(\mathbf{x}, \mathbf{z}) R(\mathbf{x}, \mathbf{z}) = \begin{cases} P_X(\mathbf{x}) R_{Z|X}(\mathbf{z}|\mathbf{x}) & \text{if } R_X(\mathbf{x}) > 0 \\ 0 & \text{if } R_X(\mathbf{x}) = 0 \end{cases},$$

completing the proof. \square

Proposition B.2.3 (Projection in χ^2 -Divergence). *Assume that $P_X \ll R_X$ and $P_Z \ll R_Z$. Define*

$$Q^* := \arg \min_{Q \in \Pi_X} \chi^2(Q \| R), \quad P^* := \arg \min_{Q \in \Pi_Z} \chi^2(Q \| R). \quad (\text{B.14})$$

Then, it holds that

$$\begin{aligned}
Q^*(\mathbf{x}, \mathbf{z}) &= \begin{cases} P_X(\mathbf{x}) R_{Z|X}(\mathbf{z}|\mathbf{x}) & \text{if } R_X(\mathbf{x}) > 0 \\ 0 & \text{if } R_X(\mathbf{x}) = 0 \end{cases} \\
P^*(\mathbf{x}, \mathbf{z}) &= \begin{cases} P_Z(\mathbf{z}) R_{X|Z}(\mathbf{x}|\mathbf{z}) & \text{if } R_Z(\mathbf{z}) > 0 \\ 0 & \text{if } R_Z(\mathbf{z}) = 0 \end{cases}.
\end{aligned} \quad (\text{B.15})$$

Proof. The first equality of (B.15) follows by the claim of Lemma B.2.2. The second equality follows by repeating the argument of Lemma B.2.1 and Lemma B.2.2 with (X, \mathbf{x}) and (Z, \mathbf{z}) swapped. \square

B.2.2 Proof of Main Results

We may now control the errors of the ratio of marginals using the projection interpretation established in the previous sections. Recall the event \mathcal{S} . The following result, the monotonicity of the marginal violation terms in terms of KL, will be useful in the bound.

Proposition B.2.4. [*Nutz, 2021, Proposition 6.10*] *Under the event \mathcal{S} , it holds that*

$$\text{KL}(P_{n,X}^{(0)} \| P_X) \geq \text{KL}(P_Z \| P_{n,Z}^{(1)}) \geq \text{KL}(P_{n,X}^{(2)} \| P_X) \geq \dots$$

We give the following result for \mathcal{X} , and the analogous claim holds on \mathcal{Z} .

Proposition B.2.5. *Assume that $P_{n,X}(\mathbf{x}) > 0$ for all $\mathbf{x} \in \mathcal{X}$. It holds that*

$$\max_{\mathbf{x} \in \mathcal{X}} \left| \frac{P_X(\mathbf{x})}{P_{n,X}^{(k-1)}(\mathbf{x})} - 1 \right| \leq \begin{cases} \max\{n-1, 1\} & \text{if } k = 1 \\ \max\{1/p_\star^2 - 1, 1\} & \text{if } k > 1. \end{cases} \quad (\text{B.16})$$

In addition, we have that

$$\max_{\mathbf{x} \in \mathcal{X}} \left| \frac{P_X(\mathbf{x})}{P_{n,X}^{(k-1)}(\mathbf{x})} - 1 \right| \leq \begin{cases} n \sqrt{\frac{1}{2} \text{KL}(P_{n,X} \| P_X)} & \text{if } k = 1 \\ \frac{1}{p_\star^2} \sqrt{\frac{1}{2} \text{KL}(P_{n,X} \| P_X)} & \text{if } k > 1 \end{cases}.$$

Moreover, when $\text{KL}(P_{n,X} \| P_X) \leq p_\star^2/2$, we have

$$\max_{\mathbf{x} \in \mathcal{X}} \left| \frac{P_X(\mathbf{x})}{P_{n,X}^{(k-1)}(\mathbf{x})} - 1 \right| \leq \frac{2}{p_\star} \sqrt{\frac{1}{2} \text{KL}(P_{n,X} \| P_X)}. \quad (\text{B.17})$$

Proof. We first show that $P_n^{(k-1)}(\mathbf{x}) \geq 1/n$ for $k = 1$ and $P_n^{(k-1)}(\mathbf{x}) \geq p_\star^2$ for $k > 1$. In the case that $k = 1$, the result follows directly from the event \mathcal{S} . For $k > 1$ such that k is odd, we have that for $\mathbf{x} \in \mathcal{X}$,

$$\begin{aligned} P_{n,X}^{(k-1)}(\mathbf{x}) &= \sum_{\mathbf{z} \in \mathcal{Z}} P_n^{(k-1)}(\mathbf{x}, \mathbf{z}) = \sum_{\mathbf{z} \in \mathcal{Z}} \frac{P_Z(\mathbf{z})}{P_{n,Z}^{(k-2)}(\mathbf{z})} P_n^{(k-2)}(\mathbf{x}, \mathbf{z}) \\ &\geq p_\star \sum_{\mathbf{z} \in \mathcal{Z}} P_n^{(k-2)}(\mathbf{x}, \mathbf{z}) = p_\star P_{n,X}^{(k-2)}(\mathbf{x}) = p_\star P_X(\mathbf{x}) \geq p_\star^2. \end{aligned}$$

The result for k even can be proven similarly. We now proceed to prove the inequalities given in the statement, which will rely on the lower bound above.

Proving the first inequality. Then, for any $\mathbf{x} \in \mathcal{X}$,

$$\left| \frac{P_X(\mathbf{x})}{P_{n,X}^{(k-1)}(\mathbf{x})} - 1 \right| = \max \left\{ \frac{P_X(\mathbf{x})}{P_{n,X}^{(k-1)}(\mathbf{x})} - 1, 1 - \frac{P_X(\mathbf{x})}{P_{n,X}^{(k-1)}(\mathbf{x})} \right\} \leq \begin{cases} \max\{n-1, 1\} & \text{if } k = 1 \\ \max\{1/p_\star^2 - 1, 1\} & \text{if } k > 1 \end{cases},$$

which is the desired result for the first inequality.

Proving the second and third inequalities. Consider an odd $k \geq 1$. By the definition of total variation distance, it holds that

$$\max_{\mathbf{x} \in \mathcal{X}} |P_X(\mathbf{x}) - P_{n,X}^{(k-1)}(\mathbf{x})| \leq \text{TV}(P_{n,X}^{(k-1)}, P_X).$$

According to Pinsker's inequality, we have that $\text{TV}(P_{n,X}^{(k-1)}, P_X) \leq \sqrt{\frac{1}{2} \text{KL}(P_{n,X}^{(k-1)} \| P_X)}$, and so we have that

$$\max_{\mathbf{x} \in \mathcal{X}} |P_X(\mathbf{x}) - P_{n,X}^{(k-1)}(\mathbf{x})| \leq \sqrt{\frac{1}{2} \text{KL}(P_{n,X}^{(k-1)} \| P_X)} \leq \sqrt{\frac{1}{2} \text{KL}(P_{n,X}^{(0)} \| P_X)},$$

where the last inequality follows from the monotonicity of Sinkhorn iterations given in Proposition B.2.4. We apply the lower bounds to write

$$\max_{\mathbf{x} \in \mathcal{X}} \left| \frac{P_X(\mathbf{x})}{P_{n,X}^{(k-1)}(\mathbf{x})} - 1 \right| \leq \begin{cases} n \sqrt{\frac{1}{2} \text{KL}(P_{n,X} \| P_X)} & \text{if } k = 1 \\ \frac{1}{p_\star^2} \sqrt{\frac{1}{2} \text{KL}(P_{n,X} \| P_X)} & \text{if } k > 1 \end{cases}.$$

Finally, when $\sqrt{\frac{1}{2} \text{KL}(P_{n,X} \| P_X)} \leq p_\star/2$, we have that $\max_{\mathbf{x} \in \mathcal{X}} |P_X(\mathbf{x}) - P_{n,X}^{(k-1)}(\mathbf{x})| \leq p_\star/2$ and thus

$$\min_{\mathbf{x} \in \mathcal{X}} P_{n,X}^{(k-1)}(\mathbf{x}) \geq \min_{\mathbf{x} \in \mathcal{X}} P_X(\mathbf{x}) - \max_{\mathbf{x} \in \mathcal{X}} |P_{n,X}^{(k-1)}(\mathbf{x}) - P_X(\mathbf{x})| \geq \frac{p_\star}{2}.$$

Hence,

$$\max_{\mathbf{x} \in \mathcal{X}} \left| \frac{P_X(\mathbf{x})}{P_{n,X}^{(k-1)}(\mathbf{x})} - 1 \right| \leq \frac{\max_{\mathbf{x} \in \mathcal{X}} |P_{n,X}^{(k-1)}(\mathbf{x}) - P_X(\mathbf{x})|}{\min_{\mathbf{x} \in \mathcal{X}} P_{n,X}^{(k-1)}(\mathbf{x})} \leq \frac{2}{p_\star} \sqrt{\frac{1}{2} \text{KL}(P_{n,X} \| P_X)}.$$

Now, for k even, set $k = 2t$ for $t \geq 0$. We have that

$$\max_{\mathbf{z} \in \mathcal{Z}} |P_{n,Z}^{(2t-1)}(\mathbf{z}) - P_Z(\mathbf{z})| \leq \text{TV}(P_{n,Z}^{(2t-1)}, P_Z) \leq \sqrt{\frac{1}{2} \text{KL}(P_Z \| P_{n,Z}^{(2t-1)})}.$$

Invoke Proposition B.2.4 once again to achieve

$$\sqrt{\frac{1}{2} \text{KL}(P_Z \| P_{n,Z}^{(2t-1)})} \leq \sqrt{\frac{1}{2} \text{KL}(P_{n,X} \| P_X)},$$

which completes the proof. \square

B.3 Statistical Analysis of Balancing Estimators

This section contains the proof of the main result, namely Theorem 4.3.1. We first consolidate notation and then give a broad outline of the proof for readability. Let the expectation of a function h under a probability measure Q on $\mathcal{X} \times \mathcal{Z}$ by denoted by

$$Q(h) = \sum_{\mathbf{x} \in \mathcal{X}, \mathbf{z} \in \mathcal{Z}} h(\mathbf{x}, \mathbf{z}) Q(\mathbf{x}, \mathbf{z})$$

so that

$$\boldsymbol{\theta}_n^{(k)} = P_n^{(k)}(h), \quad \boldsymbol{\theta} = P(h),$$

and

$$\mathbb{G}_n^{(k)}(h) = \sqrt{n}[P_n^{(k)} - P](h) = \sqrt{n}(P_n^{(k)}(h) - P(h)). \quad (\text{B.18})$$

Recalling in addition that $\mathcal{C}_k = \mathcal{C}_X$ for k odd and $\mathcal{C}_k = \mathcal{C}_Z$ for k even. The event

$$\mathcal{S} := \{\text{Supp}(P_{n,X}) = \text{Supp}(P_X) \text{ and } \text{Supp}(P_{n,Z}) = \text{Supp}(P_Z)\}, \quad (\text{B.19})$$

is used for purely technical reasons in many results.

Proof Outline We first establish that the recursion formula

$$[P_n^{(k)} - P](h) = [P_n^{(k-1)} - P](\mathcal{C}_k h) + V_n^{(k-1)}(\mathcal{C}_k h)$$

holds in Proposition B.3.1, where

$$V_n^{(k-1)}(h) = \begin{cases} \sum_{\mathbf{x}, y} \left(\frac{P_X}{P_{n,X}^{(k-1)}}(\mathbf{x}) - 1 \right) h(\mathbf{x}, \mathbf{z}) P_n^{(k-1)}(\mathbf{x}, \mathbf{z}) & k \text{ odd} \\ \sum_{\mathbf{x}, y} \left(\frac{P_Z}{P_{n,Z}^{(k-1)}}(\mathbf{z}) - 1 \right) h(\mathbf{x}, \mathbf{z}) P_n^{(k-1)}(\mathbf{x}, \mathbf{z}) & k \text{ even} \end{cases}. \quad (\text{B.20})$$

The quantity $V_n^{(k-1)}(\mathcal{C}_k h)$ describes an error term that accumulates for each iteration of balancing, which explains why k must be scaled appropriately against n to ensure the error does not accumulate too fast. Applying the recursion repeatedly to the balanced sequence $(P_n^{(k)})_{k \geq 1}$ and unrolling the recursion, we see that when k is odd,

$$\begin{aligned} [P_n^{(k)} - P](h) &= [P_n^{(k-1)} - P](\mathcal{C}_X h) + V_n^{(k-1)}(\mathcal{C}_X h) \\ &= [P_n^{(k-2)} - P](\mathcal{C}_Z \mathcal{C}_X h) + V_n^{(k-2)}(\mathcal{C}_Z \mathcal{C}_X h) + V_n^{(k-1)}(\mathcal{C}_X h) \\ &= \underbrace{[P_n^{(0)} - P](\mathcal{C}_1 \dots \mathcal{C}_k h)}_{\text{first-order term}} + \underbrace{\sum_{\ell=1}^k V_n^{(\ell-1)}(\mathcal{C}_\ell \dots \mathcal{C}_k h)}_{\text{higher-order term}} \end{aligned} \quad (\text{B.21})$$

Additionally, let $h_{\ell,k} := \mathcal{C}_\ell \dots \mathcal{C}_k h$, so that the first-order term can be written as $P_n^{(0)}(h_{1,k}) - P(h_{1,k})$ higher-order term can also be written as $\sum_{\ell=1}^k V_n^{(\ell-1)}(h_{\ell,k})$. Because our original goal is to upper bound the mean squared error, we use the expansion above to write

$$\begin{aligned} \mathbb{E} |P_n^{(k)}(h) - P(h)|^2 &\leq \mathbb{E} |P_n^{(0)}(h_{1,k}) - P(h_{1,k})|^2 \\ &\quad + 2\mathbb{E} |P_n^{(0)}(h_{1,k}) - P(h_{1,k})| \left| \sum_{\ell=1}^k V_n^{(\ell-1)}(h_{\ell,k}) \right| + \mathbb{E} \left| \sum_{\ell=1}^k V_n^{(\ell-1)}(h_{\ell,k}) \right|^2 \end{aligned}$$

Regarding the first term, we have that $\mathbb{E} |P_n^{(0)}(h_{1,k}) - P(h_{1,k})|^2 = \sigma_k^2/n$, which is the dominant term in Theorem 4.3.1. Thus, the remaining challenge of the proof will be to upper bound the cross term and squared term and show their dependence on n . The dominant term of these

two will be the cross term, as we will essentially show that $|P_n^{(0)}(h_{1,k}) - P(h_{1,k})|$ is $O(n^{-1/2})$ with high probability and that $|\sum_{\ell=1}^k V_n^{(\ell-1)}(h_{\ell,k})|$ is in fact $O(n^{-1})$ with high probability. As stated in Section 4.3, a key intermediate result in controlling the higher-order term is Proposition B.2.5, whose proof is given in Appendix B.2. The remaining subsections walk through these steps in detail.

B.3.1 Recursion of Estimation Error

We first recall that the sequence $(P_n^{(k)})_{k \geq 1}$ can be computed with the following formula:

$$P_n^{(0)}(\mathbf{x}, \mathbf{z}) := P_n(\mathbf{x}, \mathbf{z}) \text{ and } P_n^{(k)}(\mathbf{x}, \mathbf{z}) := \begin{cases} \frac{P_X}{P_n^{(k-1)}}(\mathbf{x}) P_n^{(k-1)}(\mathbf{x}, \mathbf{z}) & k \text{ odd} \\ \frac{P_Z}{P_n^{(k-1)}}(\mathbf{z}) P_n^{(k-1)}(\mathbf{x}, \mathbf{z}) & k \text{ even} \end{cases}. \quad (\text{B.22})$$

Proposition B.3.1 establishes the conditions under which these steps are well-defined (i.e. $P_{n,X}^{(k-1)}(\mathbf{x}) > 0$ and $P_{n,Z}^{(k-1)}(\mathbf{z}) > 0$).

Proposition B.3.1. *Let $(P_n^{(k)})_{k \geq 1}$, be a sequence computed according to (4.2). These iterations are well-defined under the event \mathcal{S} , and for $\mathbb{G}_n^{(k)}$ defined in (B.18) and $V_n^{(k)}$ defined in (B.20), it holds that*

$$\mathbb{G}_n^{(k)}(h) = \mathbb{G}_n^{(k-1)}(h) + \sqrt{n} V_n^{(k-1)}(h). \quad (\text{B.23})$$

and

$$\mathbb{G}_n^{(k)}(h) = \mathbb{G}_n^{(k-1)}(\mathcal{C}_k h) + \sqrt{n} V_n^{(k-1)}(\mathcal{C}_k h). \quad (\text{B.24})$$

Proof. First, assume that $P_{n,X}^{(k-1)}(\mathbf{x}) > 0$ and $P_{n,Z}^{(k-1)}(\mathbf{z}) > 0$ for all $\mathbf{x} \in \mathcal{X}$ and $\mathbf{z} \in \mathcal{Z}$ so that we may establish the recursion, which we will show by induction toward the end of the proof.

Consider the following steps in the case that k is odd:

$$\begin{aligned}
P_n^{(k)}(h) &= \sum_{\mathbf{x}, \mathbf{y}} h(\mathbf{x}, \mathbf{z}) P_n^{(k)}(\mathbf{x}, \mathbf{z}) = \sum_{\mathbf{x}, \mathbf{y}} h(\mathbf{x}, \mathbf{z}) \frac{P_X}{P_{n,X}^{(k-1)}}(\mathbf{x}) P_n^{(k-1)}(\mathbf{x}, \mathbf{z}) && \text{by (B.22) for } k \text{ odd} \\
&= \sum_{\mathbf{x}, \mathbf{y}} \mathbf{1} \cdot h(\mathbf{x}, \mathbf{z}) P_n^{(k-1)}(\mathbf{x}, \mathbf{z}) + \sum_{\mathbf{x}, \mathbf{y}} \left[\frac{P_X}{P_{n,X}^{(k-1)}}(\mathbf{x}) - 1 \right] \cdot h(\mathbf{x}, \mathbf{z}) P_n^{(k-1)}(\mathbf{x}, \mathbf{z}) \\
&= P_n^{(k-1)}(h) + V_n^{(k-1)}(h).
\end{aligned}$$

Arguing analogously for k even and subtracting $P(h)$ on both sides, we have that

$$[P_n^{(k)} - P](h) = [P_n^{(k-1)} - P](h) + V_n^{(k-1)}(h). \quad (\text{B.25})$$

We refer to this as the “uncentered” recursion, which proves (B.23).

We can then establish the following “centered” recursion using the following decomposition in the case of k odd.

$$\begin{aligned}
[P_n^{(k)} - P](h) &= [P_n^{(k)} - P](\mathcal{C}_X h) + [P_n^{(k)} - P](\mu_X h) && h = \mathcal{C}_X h + \mu_X h \\
&= [P_n^{(k-1)} - P](\mathcal{C}_X h) + V_n^{(k-1)}(\mathcal{C}_X h) + [P_n^{(k)} - P](\mu_X h) && \text{apply (B.25) to } \mathcal{C}_X h \\
&= [P_n^{(k-1)} - P](\mathcal{C}_X h) + V_n^{(k-1)}(\mathcal{C}_X h). && P_n^{(k)}(\mu_X h) = P(\mu_X h)
\end{aligned}$$

The last line follows because $\mu_X h$ is only a function on \mathcal{X} , and due to the definition of the marginal rebalancing iterations, $P_{n,X}^{(k)} = P_X$. This gives the desired formula by substituting (B.18).

We proceed to show that the iterations are well-defined. We will in fact show that $P_{n,X}^{(k-1)}(\mathbf{x}) > 0$ and $P_{n,Z}^{(k-1)}(\mathbf{z}) > 0$ for all $\mathbf{x} \in \mathcal{X}$ and $\mathbf{z} \in \mathcal{Z}$. For $k = 1$, $P_{n,X}^{(0)}(\mathbf{x}) = P_{n,X}(\mathbf{x}) > 0$ and $P_{n,Z}^{(0)}(\mathbf{z}) = P_{n,Z}(\mathbf{z}) > 0$ for all $\mathbf{x} \in \mathcal{X}$ and $\mathbf{z} \in \mathcal{Z}$ this holds under the event \mathcal{S} by assumption. We argue by induction that this holds for all $k > 1$. Assume that the claim is

true for $\{1, \dots, k-1\}$, and that k is even. Then,

$$\begin{aligned} P_{n,X}^{(k-1)}(\mathbf{x}) &= P_X(\mathbf{x}) > 0, \\ P_{n,Z}^{(k-1)}(\mathbf{z}) &= \sum_{\mathbf{x} \in \mathcal{X}} P_n^{(k-1)}(\mathbf{x}, \mathbf{z}) = \sum_{\mathbf{x} \in \mathcal{X}} \frac{P_X}{P_{n,X}^{(k-2)}}(\mathbf{x}) P_n^{(k-2)}(\mathbf{x}, \mathbf{z}) \\ &\geq \min_{\mathbf{x} \in \mathcal{X}} \frac{P_X}{P_{n,X}^{(k-2)}}(\mathbf{x}) \cdot P_{n,Z}^{(k-2)}(\mathbf{z}) > 0 \end{aligned}$$

as $P_{n,X}^{(k-2)}(\mathbf{x}) > 0$ and $P_{n,Z}^{(k-2)}(\mathbf{z}) > 0$ by the inductive hypothesis. Arguing analogously for k odd achieves the claim. \square

B.3.2 Technical Tools & Intermediate Results

Having established the backbone of the argument, we collect in this subsection some useful tools that are used in the remainder of the proofs.

The following result follows from the method of types in information theory and will be helpful in deriving the dependence of the higher-order term on n .

Theorem B.3.1. [[Cover, 1999](#), Theorem 11.2.1] *Let ν be a discrete probability measure supported on m atoms. Let $U_1, \dots, U_n \stackrel{i.i.d.}{\sim} \nu$ and ν_n be the associated empirical measure. Then, we have for any $\varepsilon > 0$ that*

$$\mathbb{P}(\text{KL}(\nu_n \| \nu) \geq \varepsilon) \leq 2^{-n(\varepsilon - m \frac{\log(n+1)}{n})}.$$

We then provide a result that counts the number of terms that appear when repeatedly centering via the operators $\mathcal{C}_1, \dots, \mathcal{C}_k$. This formalizes the pattern

$$\begin{aligned} \mathcal{C}_X &= I - \mu_X \\ \mathcal{C}_Z \mathcal{C}_X &= I - \mu_X - \mu_Z + \mu_Z \mu_X \\ \mathcal{C}_X \mathcal{C}_Z \mathcal{C}_X &= I - \mu_X - \mu_Z + \mu_Z \mu_X + \mu_X \mu_Z - \mu_X \mu_Z \mu_X, \end{aligned}$$

and so on. This will be useful when bounding $h_{\ell,k}$ uniformly.

Lemma B.3.1. For any $k \geq 1$ and $\ell \in \{1, \dots, k\}$,

$$\begin{aligned} \mathcal{C}_\ell \dots \mathcal{C}_k &= I - \sum_{\tau=0}^{(k-\ell-1)/2} (\mu_X \mu_Z)^\tau \mu_X - \sum_{\tau=0}^{(k-\ell-1)/2} (\mu_Z \mu_X)^\tau \mu_Z \\ &\quad + \sum_{\tau=1}^{(k-\ell)/2} (\mu_X \mu_Z)^\tau + \sum_{\tau=1}^{(k-\ell)/2} (\mu_Z \mu_X)^\tau + (-1)^{k-\ell+1} \mu_\ell \dots \mu_k, \end{aligned}$$

where the sum $\sum_{\tau=i}^j$ is 0 when $i > j$ and is $\sum_{\tau=i}^{\lfloor j \rfloor}$ when j is not an integer by convention.

Proof. We prove the claim by backward induction on ℓ , for the case that k is odd. In the case $\ell = k$, the claim holds because $\mathcal{C}_k = I - \mu_k$. Next, for any $\ell < k$, assume that the stated result holds for $\{\ell + 1, \dots, k\}$. Then, if ℓ is also odd (so that $\mu_\ell = \mu_X$),

$$\begin{aligned} \mathcal{C}_\ell \dots \mathcal{C}_k &= \mathcal{C}_\ell \mathcal{C}_{\ell+1} \dots \mathcal{C}_k \\ &= I - \sum_{\tau=0}^{(k-\ell-2)/2} (\mu_X \mu_Z)^\tau \mu_X - \sum_{\tau=0}^{(k-\ell-2)/2} (\mu_Z \mu_X)^\tau \mu_Z \\ &\quad + \sum_{\tau=1}^{(k-\ell-1)/2} (\mu_X \mu_Z)^\tau + \sum_{\tau=1}^{(k-\ell-1)/2} (\mu_Z \mu_X)^\tau + \mu_Z \underbrace{\dots}_{k-\ell \text{ terms}} \mu_X \\ &\quad - \mu_X + \sum_{\tau=0}^{(k-\ell-2)/2} (\mu_X \mu_Z)^\tau \mu_X + \sum_{\tau=0}^{(k-\ell-2)/2} \mu_X (\mu_Z \mu_X)^\tau \mu_Z \\ &\quad - \sum_{\tau=1}^{(k-\ell-1)/2} (\mu_X \mu_Z)^\tau - \sum_{\tau=1}^{(k-\ell-1)/2} \mu_X (\mu_Z \mu_X)^\tau - (\mu_X \mu_Z)^{(k-\ell)/2} \mu_X \end{aligned}$$

The red terms and blue terms cancel out to zero. This leaves

$$\begin{aligned} \mathcal{C}_\ell \dots \mathcal{C}_k &= I - \sum_{\tau=0}^{(k-\ell-2)/2} (\mu_X \mu_Z)^\tau \mu_X - \sum_{\tau=0}^{(k-\ell-2)/2} (\mu_Z \mu_X)^\tau \mu_Z \\ &\quad + \sum_{\tau=1}^{(k-\ell-1)/2} (\mu_Z \mu_X)^\tau + (\mu_Z \mu_X)^{(k-\ell)/2} \\ &\quad + \sum_{\tau=0}^{(k-\ell-2)/2} \mu_X (\mu_Z \mu_X)^\tau \mu_Z + (-1)^{k-\ell+1} \mu_\ell \dots \mu_k \end{aligned}$$

wherein we combine the red terms and re-index the blue terms to get

$$\begin{aligned} \mathcal{C}_\ell \dots \mathcal{C}_k &= I - \sum_{\tau=0}^{(k-\ell-2)/2} (\mu_X \mu_Z)^\tau \mu_X - \sum_{\tau=0}^{(k-\ell-2)/2} (\mu_Z \mu_X)^\tau \mu_Z \\ &\quad + \sum_{\tau=1}^{(k-\ell)/2} (\mu_Z \mu_X)^\tau + \sum_{\tau=1}^{(k-\ell)/2} (\mu_X \mu_Z)^\tau + (-1)^{k-\ell+1} \mu_\ell \dots \mu_k. \end{aligned}$$

Finally, because $k - \ell$ is even when k is odd and ℓ is odd, we can set the upper bound of the first two sums to $(k - \ell - 1)/2$ without changing the number of terms. This proves the desired result. The result can be proved similarly when ℓ is even. As a result, we have proved the claim for any odd k and $\ell \leq k$. Similar arguments can be used for the case of k even and $\ell \leq k$. \square

B.3.3 Analysis of Higher-Order Term

Returning to the outline at the start of this section, we may now bound the higher-order remainder term in (B.21), namely

$$\sum_{\ell=1}^k V_n^{(\ell-1)}(h_{\ell,k}) = \sum_{\ell=1}^k V_n^{(\ell-1)}(\mathcal{C}_\ell \dots \mathcal{C}_k h),$$

depends on controlling the quantity $V_n^{(k-1)}$ in the summation, which we recall for convenience:

$$V_n^{(k-1)}(h) = \begin{cases} \sum_{\mathbf{x}, y} \left(\frac{P_X}{P_{n,X}^{(k-1)}}(\mathbf{x}) - 1 \right) h(\mathbf{x}, \mathbf{z}) P_n^{(k-1)}(\mathbf{x}, \mathbf{z}) & k \text{ odd} \\ \sum_{\mathbf{x}, y} \left(\frac{P_Z}{P_{n,Z}^{(k-1)}}(\mathbf{z}) - 1 \right) h(\mathbf{x}, \mathbf{z}) P_n^{(k-1)}(\mathbf{x}, \mathbf{z}) & k \text{ even} \end{cases}. \quad (\text{B.26})$$

Because we have established uniform control over the functions $P_X/P_{n,X}^{(k-1)} - 1$ and $P_Z/P_{n,Z}^{(k-1)} - 1$, via Proposition B.2.5 in Appendix B.2 we can now bound the full remainder in Proposition B.3.2.

We also make use of the following intermediate result, which controls how large the ℓ_∞ -norm of the function h can grow after centering.

Lemma B.3.2. $\|h_{\ell,k}\|_\infty \leq 2(k - \ell + 1) \|h\|_\infty$.

Proof. Apply Lemma B.3.1 and the triangle inequality, so that we only need to count the number of terms that appear in the sums, adding 2 for the first and last term in the expression. We subtract 1 from the total, as one of either $(k - \ell)/2$ or $(k - \ell + 1)/2$ will be a fraction. This yields $2(k - \ell + 1)$ terms total, the desired result. \square

We upper bound the sum in Proposition B.3.2. To do so, we introduce some notation. Consider B_1 and B_2 defined by

$$B_1 := M_1 \quad \text{and} \quad B_2 := \max_{2 \leq \ell \leq k} M_\ell \quad \text{for} \quad M_\ell := \begin{cases} \max_{\mathbf{x} \in \mathcal{X}} \left| \frac{P_X(\mathbf{x})}{P_{n,X}^{(\ell-1)}(\mathbf{x})} - 1 \right| & \ell \text{ odd} \\ \max_{\mathbf{z} \in \mathcal{Z}} \left| \frac{P_Z(\mathbf{z})}{P_{n,Z}^{(\ell-1)}(\mathbf{z})} - 1 \right| & \ell \text{ even} \end{cases}$$

for $k \geq 1$. We also enumerate the sample spaces as $\mathcal{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_m\}$ and $\mathcal{Z} = \{\mathbf{z}_1, \dots, \mathbf{z}_m\}$, and define the function

$$\mathbf{1}_{jk}(\mathbf{x}, \mathbf{z}) := \begin{cases} \mathbb{1}\{\mathbf{x} = \mathbf{x}_j\} & k \text{ odd} \\ \mathbb{1}\{\mathbf{z} = \mathbf{z}_j\} & k \text{ even} \end{cases}.$$

This is an indicator function on the j -th element of either \mathcal{X} or \mathcal{Z} , depending on whether k is odd or even. Finally, for any function h , use (under the event \mathcal{S}) recall the empirical process notation

$$\mathbb{G}_n^{(k)}(h) := \sqrt{n}(P_n^{(k)}(h) - P(h)). \quad (\text{B.27})$$

Using this notation, we can rewrite the recursion in terms of the quantity $\mathbb{G}_n^{(k)}(h)$ itself. This is established in the following lemma.

Lemma B.3.3. *For k odd, it holds that*

$$\mathbb{G}_n^{(k)}(h) = \mathbb{G}_n^{(k-1)}(\mathcal{C}_X h) + \sum_{j=1}^m \left[\frac{P_X(\mathbf{x}_j)}{P_{n,X}^{(k-1)}(\mathbf{x}_j)} - 1 \right] \mathbb{G}_n^{(k-1)}(\mathcal{C}_X h \mathbf{1}_{jk}),$$

whereas for k even, it holds that

$$\mathbb{G}_n^{(k)}(h) = \mathbb{G}_n^{(k-1)}(\mathcal{C}_Z h) + \sum_{j=1}^m \left[\frac{P_Z(\mathbf{z}_j)}{P_{n,Z}^{(k-1)}(\mathbf{z}_j)} - 1 \right] \mathbb{G}_n^{(k-1)}(\mathcal{C}_Z h \mathbf{1}_{jk}),$$

Proof. We give the proof for k odd. By (B.24) from Proposition B.3.1 and by the definition of $\mathbb{G}_n^{(k)}(h)$, we need only show that $P(\mathcal{C}_X h \mathbf{1}_{jk}) = 0$. Indeed,

$$\mathbb{E}[\mathcal{C}_X h \mathbf{1}_{jk} | X](\mathbf{x}) = \begin{cases} \mathbb{E}[\mathcal{C}_X h | X](\mathbf{x}_j) & \text{if } x = \mathbf{x}_j \\ 0 & \text{if } x \neq \mathbf{x}_j \end{cases}.$$

But $\mathbb{E}[\mathcal{C}_X h | X](\mathbf{x}_j) = 0$ by definition of \mathcal{C}_X . Taking an expectation over P_X gives that $P(\mathcal{C}_X h \mathbf{1}_{jk}) = 0$, which implies the desired result. The proof for k even follows symmetrically. \square

The higher-order term in (B.21), can be bounded using Proposition B.3.2.

Proposition B.3.2. *For any $k \geq 1$, the following holds under the event \mathcal{S} :*

$$\begin{aligned} \sqrt{n} \sum_{\ell=1}^k |V_n^{(\ell-1)}(\mathcal{C}_\ell \dots \mathcal{C}_k h)| &\leq \sum_{j=1}^m \left(B_1 |\mathbb{G}_n^{(0)}(h_{1,k} \mathbf{1}_{j\ell})| + B_2 \sum_{\ell=2}^k |\mathbb{G}_n^{(0)}(h_{\ell,k} \mathbf{1}_{j\ell})| \right) \\ &\quad + m B_2 \|h\|_\infty \sqrt{n} k(k-1) [B_1 + B_2(k+1)/3]. \end{aligned}$$

Proof. First, for any $\ell \in \{1, \dots, k\}$, recall the notation $h_{\ell,k} := \mathcal{C}_\ell \dots \mathcal{C}_k h$. By (B.23) from Proposition B.3.1 and by Lemma B.3.3, we have that for ℓ odd,

$$\sqrt{n} V_n^{(\ell-1)}(h_{\ell,k}) = \sum_{j=1}^m \left[\frac{P_X}{P_{n,X}^{(\ell-1)}}(\mathbf{x}_j) - 1 \right] \mathbb{G}_n^{(\ell-1)}(h_{\ell,k} \mathbf{1}_{j\ell}). \quad (\text{B.28})$$

Using the statement above, we have that

$$\sqrt{n} |V_n^{(\ell-1)}(h_{\ell,k})| \leq M_\ell \sum_{j=1}^m |\mathbb{G}_n^{(\ell-1)}(h_{\ell,k} \mathbf{1}_{j\ell})|.$$

The bound above holds for ℓ even as well. Then, using the (B.23) from Proposition B.3.1

again, we have that for $\ell \geq 2$,

$$[P_n^{(\ell-1)} - P](h_{\ell,k} \mathbf{1}_{j\ell}) = [P_n^{(\ell-2)} - P](h_{\ell,k} \mathbf{1}_{j\ell}) + V_n^{(\ell-2)}(h_{\ell,k} \mathbf{1}_{j\ell})$$

which implies that

$$\begin{aligned} |\mathbb{G}_n^{(\ell-1)}(h_{\ell,k} \mathbf{1}_{j\ell})| &\leq |\mathbb{G}_n^{(\ell-2)}(h_{\ell,k} \mathbf{1}_{j\ell})| + \sqrt{n} |V_n^{(\ell-2)}(h_{\ell,k} \mathbf{1}_{j\ell})| \\ &\leq |\mathbb{G}_n^{(0)}(h_{\ell,k} \mathbf{1}_{j\ell})| + \sqrt{n} |V_n^{(0)}(h_{\ell,k} \mathbf{1}_{j\ell})| + \dots + \sqrt{n} |V_n^{(\ell-2)}(h_{\ell,k} \mathbf{1}_{j\ell})| \\ &\leq |\mathbb{G}_n^{(0)}(h_{\ell,k} \mathbf{1}_{j\ell})| + M_1 \sqrt{n} P_n^{(0)}(|h_{\ell,k}| \mathbf{1}_{j\ell}) + \dots + M_\ell \sqrt{n} P_n^{(\ell-2)}(|h_{\ell,k}| \mathbf{1}_{j\ell}) \\ &\leq |\mathbb{G}_n^{(0)}(h_{\ell,k} \mathbf{1}_{j\ell})| + 2 \|h\|_\infty \sqrt{n} [B_1 + B_2(\ell-1)] (k - \ell + 1), \end{aligned} \quad (\text{B.29})$$

by Lemma B.3.2 and $M_1 \leq B_1$ and $M_\ell \leq B_2$ for $\ell \geq 2$. Summing these bounds, we have that

$$\begin{aligned} &\sqrt{n} \sum_{\ell=1}^k |V_n^{(\ell-1)}(h_{\ell,k})| \\ &\leq M_1 \sum_{j=1}^m |\mathbb{G}_n^{(0)}(h_{1,k} \mathbf{1}_{j1})| + \sum_{\ell=2}^k M_\ell \sum_{j=1}^m |\mathbb{G}_n^{(\ell-1)}(h_{\ell,k} \mathbf{1}_{j\ell})| \\ &\leq B_1 \sum_{j=1}^m |\mathbb{G}_n^{(0)}(h_{1,k} \mathbf{1}_{j1})| + B_2 \sum_{\ell=2}^k \sum_{j=1}^m |\mathbb{G}_n^{(\ell-1)}(h_{\ell,k} \mathbf{1}_{j\ell})| \\ &\leq B_1 \sum_{j=1}^m |\mathbb{G}_n^{(0)}(h_{1,k} \mathbf{1}_{j1})| + \\ &\quad B_2 \sum_{\ell=2}^k \sum_{j=1}^m (|\mathbb{G}_n^{(0)}(h_{\ell,k} \mathbf{1}_{j\ell})| + 2 \|h\|_\infty \sqrt{n} [B_1 + B_2(\ell-1)] (k - \ell + 1)) \quad \text{apply (B.29)} \\ &= \sum_{j=1}^m \left(B_1 |\mathbb{G}_n^{(0)}(h_{1,k} \mathbf{1}_{j1})| + B_2 \sum_{\ell=2}^k |\mathbb{G}_n^{(0)}(h_{\ell,k} \mathbf{1}_{j\ell})| \right) + \\ &\quad 2mB_2 \|h\|_\infty \sqrt{n} \sum_{\ell=2}^k [B_1 + B_2(\ell-1)] (k - \ell + 1), \end{aligned}$$

because $|\mathcal{X}| = m$. We sum up the last term:

$$\begin{aligned} \sum_{\ell=2}^k [B_1 + B_2(\ell - 1)] (k - \ell + 1) &= B_1 \sum_{\ell=1}^{k-1} (k - \ell) + B_2 \sum_{\ell=1}^{k-1} \ell(k - \ell) \\ &= \frac{k(k-1)}{2} [B_1 + B_2(k+1)/3]. \end{aligned}$$

completing the proof. \square

B.3.4 Proof of Main Results

We can now show the main result of this section: the bound on the mean squared error of the rebalanced estimator. Recall the event

$$\mathcal{S} := \{\text{Supp}(P_{n,X}) = \text{Supp}(P_X) \text{ and } \text{Supp}(P_{n,Z}) = \text{Supp}(P_Z)\} \quad (\text{B.30})$$

as introduced in (B.19). To remind the reader of the high-level steps of the proof, we may decompose the error on the event \mathcal{S} we used the estimator

$$\tilde{\boldsymbol{\theta}}_n^{(k)} := \boldsymbol{\theta}_n^{(k)} \mathbb{1}_{\mathcal{S}} + \boldsymbol{\theta}_n^{(0)} \mathbb{1}_{\mathcal{S}^c}$$

so we decompose on the event \mathcal{S} to write

$$\mathbb{E}_P \left[\left(\tilde{P}_n^{(k)}(h) - P(h) \right)^2 \right] = \mathbb{E}_P \left[(P_n(h) - P(h))^2 \mathbb{1}_{\mathcal{S}^c} \right] + \mathbb{E}_P \left[(P_n^{(k)}(h) - P(h))^2 \mathbb{1}_{\mathcal{S}} \right]. \quad (\text{B.31})$$

Then, we use the upcoming Proposition B.3.3 to bound the first term, which will in turn require showing that \mathcal{S} occurs with high probability. As for the second term, we will apply Proposition B.3.1 and the derivation (B.21) to write

$$\mathbb{E}_P \left[(P_n^{(k)}(h) - P(h))^2 \mathbb{1}_{\mathcal{S}} \right] = \mathbb{E}_P \left[T_1^2 \mathbb{1}_{\mathcal{S}} \right] + 2\mathbb{E}_P \left[T_1 T_2 \mathbb{1}_{\mathcal{S}} \right] + \mathbb{E}_P \left[T_2^2 \mathbb{1}_{\mathcal{S}} \right] \quad (\text{B.32})$$

for

$$T_1 := [P_n^{(0)} - P](\mathcal{C}_1 \dots \mathcal{C}_k h) \text{ and } T_2 := \sum_{\ell=1}^k V_n^{(\ell-1)}(\mathcal{C}_\ell \dots \mathcal{C}_k h). \quad (\text{B.33})$$

By definition, we have that $\mathbb{E}_P [T_1^2 \mathbb{1}_{\mathcal{S}}] \leq \mathbb{E}_P [T_1^2] = \sigma_k^2/n$. It then remains to bound the cross term $\mathbb{E}_P [T_1 T_2 \mathbb{1}_{\mathcal{S}}]$ and squared term $\mathbb{E}_P [T_2^2 \mathbb{1}_{\mathcal{S}}]$. This is accomplished by Lemma B.3.5 and Lemma B.3.4, respectively.

Proposition B.3.3. *It holds that $P(\mathcal{S}^c) \leq 2m(1 - p_\star)^n$. Moreover, for any $\delta \in (0, 1)$, we have*

$$\mathbb{E}_P [(P_n(h) - P(h))^2 \mathbb{1}_{\mathcal{S}^c}] \leq 4 \|h\|_\infty^2 \min \{2m(1 - p_\star)^n, \delta\} + \frac{2 \log(2/\delta)}{n} \|h\|_\infty^2 2m(1 - p_\star)^n.$$

Proof. Define $\mathcal{F}_X := \{\text{Supp}(P_{n,X}) \neq \text{Supp}(P_X)\}$ and $\mathcal{F}_Z := \{\text{Supp}(P_{n,Z}) \neq \text{Supp}(P_Z)\}$, so that $\mathcal{S}^c = \mathcal{F}_X \cup \mathcal{F}_Z$. We first control the probability of \mathcal{F}_X . Let $F_j := \{P_{n,X}(\mathbf{x}_j) = 0\}$ for $j \in [m]$. We then obtain $\mathcal{F}_X = \cup_{j=1}^m F_j$, which implies by the union bound that

$$P(\mathcal{F}_X) \leq \sum_{j=1}^m P(F_j) = \sum_{j=1}^m (1 - P_X(\mathbf{x}_j))^n \leq m(1 - p_\star)^n.$$

Similarly, we have that $P(\mathcal{F}_Z) \leq m(1 - p_\star)^n$ and thus $P(\mathcal{S}^c) \leq 2m(1 - p_\star)^n$, which gives the first claim.

To control the expectation, consider any $\delta > 0$, and define the event

$$\mathcal{E}_\delta := \left\{ |P_n^{(0)}(h) - P(h)| \leq \sqrt{\frac{2 \log(2/\delta)}{n}} \|h\|_\infty \right\}.$$

By Hoeffding's inequality, it holds that $P(\mathcal{E}_\delta) \geq 1 - \delta$. Furthermore, we get

$$\begin{aligned} \mathbb{E}[\mathbb{1}_{\mathcal{S}^c} (P_n^{(0)}(h) - P(h))^2] &= \mathbb{E}[\mathbb{1}_{\mathcal{S}^c} \mathbb{1}_{\mathcal{E}_\delta^c} (P_n^{(0)}(h) - P(h))^2] + \mathbb{E}[\mathbb{1}_{\mathcal{S}^c} \mathbb{1}_{\mathcal{E}_\delta} (P_n^{(0)}(h) - P(h))^2] \\ &\leq 4 \|h\|_\infty^2 \mathbb{E}[\mathbb{1}_{\mathcal{S}^c} \mathbb{1}_{\mathcal{E}_\delta^c}] + \frac{2 \log(2/\delta)}{n} \|h\|_\infty^2 \mathbb{E}[\mathbb{1}_{\mathcal{S}^c} \mathbb{1}_{\mathcal{E}_\delta}] \\ &\leq 4 \|h\|_\infty^2 \min\{P(\mathcal{S}^c), P(\mathcal{E}_\delta^c)\} + \frac{2 \log(2/\delta)}{n} \|h\|_\infty^2 P(\mathcal{S}^c) \\ &\leq 4 \|h\|_\infty^2 \min\{2m(1 - p_\star)^n, \delta\} + \frac{2 \log(2/\delta)}{n} \|h\|_\infty^2 2m(1 - p_\star)^n. \end{aligned}$$

□

In order to bound the terms appearing in (B.32), we introduce the events \mathcal{E}_1^δ , \mathcal{E}_2^δ , and \mathcal{E}_3^δ ,

defined by

$$\begin{aligned}\mathcal{E}_1^\delta &:= \left\{ \max \{ \text{KL}(P_{n,X} \| P_X), \text{KL}(P_{n,Z} \| P_Z) \} \leq \frac{1}{n} \log_2 \frac{2}{\delta} + m \frac{\log(n+1)}{n} \right\} \\ \mathcal{F}_\ell^\delta &:= \left\{ |\mathbb{G}_n^{(0)}(h_{\ell,k} \mathbf{1}_{j\ell})| \leq \sqrt{2 \log(2mk/\delta)} 2(k-\ell+1) \|h\|_\infty \right\}, \quad \ell = 1, \dots, k, \quad j = 1, \dots, m \\ \mathcal{E}_2^\delta &:= \bigcap_{\ell=1}^k \mathcal{F}_\ell^\delta \\ \mathcal{E}_3^\delta &:= \left\{ |\mathbb{G}_n^{(0)}(h_{1,k})| \leq \sqrt{2 \log(2/\delta)} 2k \|h\|_\infty \right\}.\end{aligned}$$

The events are constructed such that $\mathbb{P}(\mathcal{E}_1^\delta) \geq 1 - \delta$, $\mathbb{P}(\mathcal{E}_2^\delta) \geq 1 - \delta$, and $\mathbb{P}(\mathcal{E}_3^\delta) \geq 1 - \delta$, as we used in the upcoming proofs of Lemma B.3.5, Lemma B.3.4, and Theorem B.3.2.

Lemma B.3.4 (Squared term bound). *Let T_2 be defined as in (B.33). For any $\delta > 0$, assuming that $n \geq 2[\log_2(2/\delta) + m \log(n+1)]/p_\star^2$, we have that*

$$\begin{aligned}\mathbb{E}_P [T_2^2 \mathbb{1}_{\mathcal{S}}] &\leq \frac{2 \|h\|_\infty^2 m^2 k^2}{p_\star^2} [\log_2(2/\delta) + m \log(n+1)]^{2-1\{k=1\}} \times \\ &\left[\left(4n + \frac{k-1}{p_\star^2} \left(n+2 + \frac{k+1}{p_\star^2} \right) \right)^2 \delta + \frac{8}{n^2} \left(\sqrt{2 \log \frac{2mk}{\delta}} (k+1) + \frac{(k-1)(k+4)}{p_\star^2} \right)^2 \right].\end{aligned}$$

Proof. The following computations are done under the event \mathcal{S} . First, apply Proposition B.3.2 to write

$$\begin{aligned}|T_2| &\leq \frac{1}{\sqrt{n}} \sum_{j=1}^m \left(B_1 |\mathbb{G}_n^{(0)}(h_{1,k} \mathbf{1}_{j\ell})| + B_2 \sum_{\ell=2}^k |\mathbb{G}_n^{(0)}(h_{\ell,k} \mathbf{1}_{j\ell})| \right) + \\ &m B_2 \|h\|_\infty k(k-1) [B_1 + B_2(k+1)/3].\end{aligned}\tag{B.34}$$

We decompose on the event $\mathcal{E}_1^\delta \cap \mathcal{E}_2^\delta$. Note that by Theorem B.3.1, we have that $\mathbb{P}(\mathcal{E}_1^\delta) \geq 1 - \delta$. It follows from Hoeffding's inequality, the union bound, and boundedness of $\|h_{\ell,k} \mathbf{1}_{j\ell}\|$ by Lemma B.3.2 that $\mathbb{P}(\mathcal{E}_2^\delta) \geq 1 - \delta$. As a result, $\mathbb{P}(\mathcal{E}_1^\delta \cap \mathcal{E}_2^\delta) \geq 1 - 2\delta$.

Bound $|T_2|$ under the event $\mathcal{S} \setminus (\mathcal{E}_1^\delta \cap \mathcal{E}_2^\delta)$. In this case, we apply (B.16) from Proposition B.2.5 to get $B_1 \leq n$ and $B_2 \leq 1/p_\star^2$, along with the universal bounds from Lemma B.3.2:

$$\begin{aligned} \frac{1}{\sqrt{n}} |\mathbb{G}_n^{(0)}(h_{1,k} \mathbf{1}_{j\ell})| &\leq 2 \|h_{1,k}\|_\infty \leq 4k \|h\|_\infty \\ \frac{1}{\sqrt{n}} \sum_{\ell=2}^k |\mathbb{G}_n^{(0)}(h_{\ell,k} \mathbf{1}_{j\ell})| &\leq 2 \sum_{\ell=2}^k \|h_{\ell,k}\|_\infty \leq \sum_{\ell=2}^k 4(k-\ell+1) \|h\|_\infty = 2k(k-1) \|h\|_\infty \end{aligned}$$

so that by plugging into (B.34),

$$|T_2| \leq \|h\|_\infty mk \left[4n + \frac{k-1}{p_\star^2} \left(n + 2 + \frac{k+1}{3p_\star^2} \right) \right],$$

and in turn,

$$\mathbb{E}_P \left[T_2^2 \mathbb{1}_{\mathcal{S} \setminus (\mathcal{E}_1^\delta \cap \mathcal{E}_2^\delta)} \right] \leq 2 \|h\|_\infty^2 m^2 k^2 \left[4n + \frac{k-1}{p_\star^2} \left(n + 2 + \frac{k+1}{3p_\star^2} \right) \right]^2 \delta. \quad (\text{B.35})$$

Bound $|T_2|$ under the event $\mathcal{S} \cap \mathcal{E}_1^\delta \cap \mathcal{E}_2^\delta$. In this case, we may use that $n \geq 2[\log_2(2/\delta) + m \log(n+1)]/p_\star^2$ apply (B.17) from Proposition B.2.5 to get

$$\max \{B_1, B_2\} \leq \frac{2}{p_\star} \sqrt{\frac{1}{2} \text{KL}(P_{n,X} \| P_X)} \leq \frac{1}{p_\star \sqrt{n}} \sqrt{2 \log_2(2/\delta) + 2m \log(n+1)}$$

and the bounds based on \mathcal{E}_2^δ , which give

$$\begin{aligned} |\mathbb{G}_n^{(0)}(h_{1,k} \mathbf{1}_{j\ell})| &\leq \sqrt{2 \log \frac{2mk}{\delta}} 2k \|h\|_\infty \\ \sum_{\ell=2}^k |\mathbb{G}_n^{(0)}(h_{\ell,k} \mathbf{1}_{j\ell})| &\leq \sum_{\ell=2}^k \sqrt{2 \log \frac{2mk}{\delta}} 2(k-\ell+1) \|h\|_\infty \leq \sqrt{2 \log \frac{2mk}{\delta}} k(k-1) \|h\|_\infty, \end{aligned}$$

By plugging into (B.34),

$$|T_2| \leq \frac{2m \|h\|_\infty \sqrt{2 \log(2mk/\delta) [2 \log_2(2/\delta) + 2m \log(n+1)]}}{np_\star} k(k+1) + \quad (\text{B.36})$$

$$\frac{m \|h\|_\infty [2 \log_2(2/\delta) + 2m \log(n+1)]}{3np_\star^2} k(k-1)(k+4) \quad (\text{B.37})$$

$$\leq \frac{4mk \|h\|_\infty [\log_2(2/\delta) + 2m \log(n+1)]^{1-\mathbb{1}\{k=1\}/2}}{np_\star^2} \times \quad (\text{B.38})$$

$$\left[p_\star \sqrt{2 \log(2mk/\delta)} (k+1) + (k-1)(k+4) \right]. \quad (\text{B.39})$$

In turn,

$$\begin{aligned} \mathbb{E}_P \left[T_2^2 \mathbb{1}_{\mathcal{S} \setminus (\mathcal{E}_1^{\delta} \cap \mathcal{E}_2^{\delta})} \right] &\leq \frac{16 \|h\|_\infty^2 m^2 k^2 [\log_2(2/\delta) + m \log(n+1)]^{2-\mathbb{1}\{k=1\}}}{n^2 p_\star^4} \times \\ &\quad \left[p_\star \sqrt{2 \log(2mk/\delta)} (k+1) + (k-1)(k+4) \right]^2. \end{aligned} \quad (\text{B.40})$$

Combining together both (B.40) and (B.35) and using that $[\log_2(2/\delta) + 2m \log(n+1)] \geq 1$, we have that

$$\begin{aligned} \mathbb{E}_P [T_2^2 \mathbb{1}_{\mathcal{S}}] &\leq \frac{2 \|h\|_\infty^2 m^2 k^2}{p_\star^2} [\log_2(2/\delta) + m \log(n+1)]^{2-\mathbb{1}\{k=1\}} \times \\ &\quad \left[\left(4n + \frac{k-1}{p_\star^2} \left(n + 2 + \frac{k+1}{p_\star^2} \right) \right)^2 \delta + \frac{8}{n^2} \left(\sqrt{2 \log(2mk/\delta)} (k+1) + \frac{(k-1)(k+4)}{p_\star^2} \right)^2 \right], \end{aligned}$$

the result as desired. \square

Lemma B.3.5 (Cross term bound). *Let T_1 and T_2 be defined as in (B.33). For any $\delta > 0$, assuming that $n \geq 2[\log_2(2/\delta) + m \log(n+1)]/p_\star^2$, we have that*

$$\begin{aligned} &\mathbb{E}_P [T_1 T_2 \mathbb{1}_{\mathcal{S}}] \\ &\leq \frac{2mk^2 \|h\|_\infty^2 \sqrt{2 \log(2/\delta)} [\log_2(2/\delta) + 2m \log(n+1)]^{1-\mathbb{1}\{k=1\}/2}}{p_\star^2} \times \\ &\quad \left[\frac{p_\star \sqrt{2 \log(2mk/\delta)} (k+1) + (k-1)(k+4)}{n^{3/2}} + 6 \left(4np_\star^2 + (k-1) \left(n + 2 + \frac{k+1}{p_\star^2} \right) \right) \delta \right], \end{aligned}$$

Proof. The following computations are done under the event \mathcal{S} . First, apply Proposi-

tion B.3.2 to write

$$|T_1 T_2| \leq \frac{1}{\sqrt{n}} |\mathbb{G}_n^{(0)}(h_{1,k})| \left[\frac{1}{\sqrt{n}} \sum_{j=1}^m \left(B_1 |\mathbb{G}_n^{(0)}(h_{1,k} \mathbf{1}_{j\ell})| + B_2 \sum_{\ell=2}^k |\mathbb{G}_n^{(0)}(h_{\ell,k} \mathbf{1}_{j\ell})| \right) + m B_2 \|h\|_\infty k(k-1)[B_1 + B_2(k+1)/3] \right]. \quad (\text{B.41})$$

We decompose on the event $\mathcal{E}_1^\delta \cap \mathcal{E}_2^\delta \cap \mathcal{E}_3^\delta$. Note that by Theorem B.3.1 and that $n \geq \log_2(2/\delta) + m \log(n+1)$, we have that $\mathbb{P}(\mathcal{E}_1^\delta) \geq 1 - \delta$. It follows from Hoeffding's inequality and the union bound that $\mathbb{P}(\mathcal{E}_2^\delta) \geq 1 - \delta$. Similarly, we also have by Hoeffding's inequality that $\mathbb{P}(\mathcal{E}_3^\delta) \geq 1 - \delta$. As a result, $\mathbb{P}(\mathcal{E}_1^\delta \cap \mathcal{E}_2^\delta \cap \mathcal{E}_3^\delta) \geq 1 - 3\delta$.

Bound $|T_1 T_2|$ under the event $\mathcal{S} \setminus (\mathcal{E}_1^\delta \cap \mathcal{E}_2^\delta \cap \mathcal{E}_3^\delta)$. In this case, we apply (B.16) from Proposition B.2.5 to get $B_1 \leq n$ and $B_2 \leq 1/p_\star^2$, along with the universal bounds from Lemma B.3.2:

$$\begin{aligned} \frac{1}{\sqrt{n}} |\mathbb{G}_n^{(0)}(h_{1,k})| &\leq 2 \|h_{1,k}\|_\infty \leq 4k \|h\|_\infty \\ \frac{1}{\sqrt{n}} |\mathbb{G}_n^{(0)}(h_{1,k} \mathbf{1}_{j\ell})| &\leq 2 \|h_{1,k}\|_\infty \leq 4k \|h\|_\infty \\ \frac{1}{\sqrt{n}} \sum_{\ell=2}^k |\mathbb{G}_n^{(0)}(h_{\ell,k} \mathbf{1}_{j\ell})| &\leq 2 \sum_{\ell=2}^k \|h_{\ell,k}\|_\infty \leq \sum_{\ell=2}^k 4(k-\ell+1) \|h\|_\infty = 2k(k-1) \|h\|_\infty, \end{aligned}$$

so that by plugging into (B.41),

$$|T_1 T_2| \leq 4k^2 \|h\|_\infty^2 m \left[4n + \frac{k-1}{p_\star^2} \left(n + 2 + \frac{k+1}{3p_\star^2} \right) \right],$$

and in turn,

$$\mathbb{E}_P \left[T_1 T_2 \mathbb{1}_{\mathcal{S} \setminus (\mathcal{E}_1^\delta \cap \mathcal{E}_2^\delta \cap \mathcal{E}_3^\delta)} \right] \leq \frac{12k^2 \|h\|_\infty^2 m}{p_\star^2} \left[4np_\star^2 + (k-1) \left(n + 2 + \frac{k+1}{3p_\star^2} \right) \right] \delta. \quad (\text{B.42})$$

Bound $|T_1 T_2|$ under the event $\mathcal{S} \cap \mathcal{E}_1^\delta \cap \mathcal{E}_2^\delta \cap \mathcal{E}_3^\delta$. In this case, we may use that $n \geq 2[\log_2(2/\delta) + m \log(n+1)]/p_\star^2$ apply (B.17) from Proposition B.2.5 to get

$$\max\{B_1, B_2\} \leq \frac{2}{p_\star} \sqrt{\frac{1}{2} \text{KL}(P_{n,X} \| P_X)} \leq \frac{1}{\sqrt{n}} \frac{1}{p_\star} \sqrt{2 \log_2(2/\delta) + 2m \log(n+1)}$$

and the bounds based on $\mathcal{E}_2^\delta \cap \mathcal{E}_2^\delta \cap \mathcal{E}_3^\delta$ which give

$$\begin{aligned} |\mathbb{G}_n^{(0)}(h_{1,k})| &\leq \sqrt{2 \log(2/\delta) 2k} \|h\|_\infty \\ |\mathbb{G}_n^{(0)}(h_{1,k} \mathbf{1}_{j\ell})| &\leq \sqrt{2 \log(2mk/\delta) 2k} \|h\|_\infty \\ \sum_{\ell=2}^k |\mathbb{G}_n^{(0)}(h_{\ell,k} \mathbf{1}_{j\ell})| &\leq \sum_{\ell=2}^k \sqrt{2 \log \frac{2mk}{\delta} 2(k-\ell+1)} \|h\|_\infty \leq \sqrt{2 \log \frac{2mk}{\delta} k(k-1)} \|h\|_\infty, \end{aligned}$$

By plugging into (B.41),

$$\begin{aligned} |T_2| &\leq \frac{m \|h\|_\infty \sqrt{2 \log(2mk/\delta) [2 \log_2(2/\delta) + 2m \log(n+1)]}}{np_\star} k(k+1) + \\ &\quad \frac{m \|h\|_\infty [2 \log_2(2/\delta) + 2m \log(n+1)]}{3np_\star^2} k(k-1)(k+4) \\ &\leq \frac{mk \|h\|_\infty [\log_2(2/\delta) + 2m \log(n+1)]^{1-\mathbb{1}\{k=1\}/2}}{np_\star^2} \times \\ &\quad \left[p_\star \sqrt{2 \log(2mk/\delta)} (k+1) + (k-1)(k+4) \right] \\ |T_1 T_2| &\leq \frac{2mk^2 \|h\|_\infty^2 \sqrt{2 \log(2/\delta)} [\log_2(2/\delta) + 2m \log(n+1)]^{1-\mathbb{1}\{k=1\}/2}}{n^{3/2} p_\star^2} \times \\ &\quad \left[p_\star \sqrt{2 \log(2mk/\delta)} (k+1) + (k-1)(k+4) \right], \end{aligned}$$

In turn,

$$\begin{aligned} \mathbb{E}_P \left[T_2^2 \mathbb{1}_{\mathcal{S} \setminus (\mathcal{E}_1^\delta \cap \mathcal{E}_2^\delta \cap \mathcal{E}_3^\delta)} \right] &\leq \frac{2mk^2 \|h\|_\infty^2 \sqrt{2 \log(2/\delta)} [\log_2(2/\delta) + 2m \log(n+1)]^{1-\mathbb{1}\{k=1\}/2}}{n^{3/2} p_\star^2} \times \\ &\quad \left[p_\star \sqrt{2 \log(2mk/\delta)} (k+1) + (k-1)(k+4) \right], \end{aligned} \tag{B.43}$$

Combining together both (B.43) and (B.42) and using that $[\log_2(2/\delta) + 2m \log(n+1)] \geq 1$, we have that

$$\begin{aligned} & \mathbb{E}_P [T_1 T_2 \mathbb{1}_S] \\ & \leq \frac{2mk^2 \|h\|_\infty^2 \sqrt{2 \log(2/\delta)} [\log_2(2/\delta) + 2m \log(n+1)]^{1-\mathbb{1}\{k=1\}/2}}{p_\star^2} \times \\ & \quad \left[\frac{p_\star \sqrt{2 \log(2mk/\delta)} (k+1) + (k-1)(k+4)}{n^{3/2}} + 6 \left(4np_\star^2 + (k-1) \left(n+2 + \frac{k+1}{p_\star^2} \right) \right) \delta \right], \end{aligned}$$

the result as desired. \square

We now combine the previous results to prove Theorem B.3.2.

Theorem B.3.2. *For a sequence of rebalanced distributions $(\tilde{P}_n^{(k)})_{k \geq 1}$, there exists an absolute constant $C > 0$ such that when $n \geq C[\log_2(2n/p_\star) + m \log(n+1)]/p_\star^2$,*

$$\mathbb{E}_P[(\tilde{P}_n^{(k)}(h) - P(h))^2] \leq \frac{\sigma_k^2}{n} + \frac{CB}{n^{3/2}}, \quad (\text{B.44})$$

where

$$B = \frac{\sqrt{\log(2n/p_\star)} m^2 k^4 \|h\|_\infty^2}{p_\star^2} \left(\log_2 \frac{2n}{p_\star} + m \log(n+1) \right)^{2-\mathbb{1}\{k\}} \left(\log \frac{2mkn}{p_\star} + \frac{(k-1)^2}{p_\star^2} \right).$$

Proof. We apply the decomposition (B.31), and subsequently handle the second term using bounds on the terms in (B.32). Set $\delta = p_\star^4/n^4$. We apply Lemma B.3.4 and Lemma B.3.5 with this choice of δ , so that there exists absolute constants \tilde{C} , C_1 , and C_2 such that

$$\begin{aligned} \mathbb{E}_P [T_1 T_2 \mathbb{1}_S] & \leq C_1 \frac{\|h\|_\infty^2 m^2 k^3 \sqrt{\log(2n/p_\star)}}{n^{3/2} p_\star^2} [\log_2(2n/p_\star) + m \log(n+1)]^{1-\mathbb{1}\{k=1\}/2} \times \\ & \quad \left(\log \frac{2mnk}{p_\star} + \frac{k-1}{p_\star^2} \right) \\ \mathbb{E}_P [T_2^2 \mathbb{1}_S] & \leq C_2 \frac{\|h\|_\infty^2 m^2 k^4}{n^2 p_\star^2} [\log_2(2n/p_\star) + m \log(n+1)]^{2-\mathbb{1}\{k=1\}} \times \\ & \quad \left(\log \frac{2mnk}{p_\star} + \frac{(k-1)^2}{p_\star^2} \right), \end{aligned}$$

when $n \geq \tilde{C}[\log_2(2n/p_\star) + m \log(n+1)]/p_\star^2$. This then implies that there is an absolute

constant C_3 such that

$$\begin{aligned} & \mathbb{E}_P \left[\left(\tilde{P}_n^{(k)}(h) - P(h) \right)^2 \right] \\ & \leq \mathbb{E}_P \left[(P_n^{(0)}(h) - P(h))^2 \mathbb{1}_{\mathcal{S}^c} \right] + \frac{\sigma_k^2}{n} + \\ & \quad \frac{C_3 \|h\|_\infty^2 m^2 k^4 \sqrt{\log(2n/p_\star)}}{n^{3/2} p_\star^2} \left[\log_2 \frac{2n}{p_\star} + m \log(n+1) \right]^{2-\mathbb{1}\{k=1\}} \left(\log \frac{2mnk}{p_\star} + \frac{(k-1)^2}{p_\star^2} \right). \end{aligned}$$

Next, we apply Proposition B.3.3 with the same choice of δ . Because $2[\log_2(2/\delta) + m \log(n+1)] \geq \log(m/\delta)$ and $-\log(1-p_\star) \geq p_\star \geq p_\star^2$, we have that $n \geq \log(\delta/m)/\log(1-p_\star)$, which implies that $m(1-p_\star)^n \leq \delta$. Combining with the display above, we have that there exists an absolute constant $C > 0$ such that

$$\begin{aligned} \mathbb{E}_P \left[\left(\tilde{P}_n^{(k)}(h) - P(h) \right)^2 \right] & \leq \frac{\sigma_k^2}{n} + \frac{C \|h\|_\infty^2 m^2 k^4 \sqrt{\log(2n/p_\star)}}{n^{3/2} p_\star^2} \\ & \quad \times [\log_2(2/\delta) + m \log(n+1)]^{2-\mathbb{1}\{k=1\}} \left(\log \frac{2mnk}{p_\star} + \frac{(k-1)^2}{p_\star^2} \right), \end{aligned}$$

which is the claimed result. \square

While not shown in the main text, similar techniques to those used above can also control the bias of $\tilde{P}_n^{(k)}(h)$ as in Theorem B.3.3. Interestingly, this bias is of order $O(n^{-2})$, which confirms the intuition that even though $\tilde{P}_n^{(k)}(h)$ may be biased, the dominant term is the variance.

Theorem B.3.3. *For a sequence of rebalanced distributions $(P^{(k)})_{k \geq 1}$, there exists an absolute constant $C > 0$ such that when $n \geq C[\log_2(2n/p_\star) + m \log(n+1)]/p_\star^2$,*

$$\left| \mathbb{E}_P[\tilde{P}_n^{(k)}(h) - P(h)] \right|^2 \leq \frac{CB}{n^2}, \quad (\text{B.45})$$

where B is as defined in Theorem B.3.2.

Proof. First, apply the decomposition (B.31) so that

$$\left| \mathbb{E}_P \left[\tilde{P}_n^{(k)}(h) - P(h) \right] \right| \leq |\mathbb{E}_P[(P_n(h) - P(h)) \mathbb{1}_{\mathcal{S}^c}]| + |\mathbb{E}_P[(P_n^{(k)}(h) - P(h)) \mathbb{1}_{\mathcal{S}}]|.$$

By using the argument of Proposition B.3.3, we have that

$$|\mathbb{E}_P [P_n(h) - P(h)] \mathbf{1}_{\mathcal{S}^c}| \leq 2 \|h\|_\infty \min \{2m(1 - p_\star)^n, \delta\} + \sqrt{\frac{2 \log(2/\delta)}{n}} \|h\|_\infty 2m(1 - p_\star)^n.$$

Then, by the recursion formula Equation (B.21), we have that

$$\begin{aligned} & \sqrt{n} |\mathbb{E}_P [(P_n^{(k)}(h) - P(h)) \mathbf{1}_{\mathcal{S}}]| \\ &= |\mathbb{E}_P [\mathbb{G}_n^{(k)}(h) \mathbf{1}_{\mathcal{S}}]| = \left| \mathbb{E}_P \left[(1 - \mathbf{1}_{\mathcal{S}^c}) \mathbb{G}_n^{(0)}(\mathcal{C}_1 \dots \mathcal{C}_k h) + \sqrt{n} \mathbf{1}_{\mathcal{S}} \sum_{\ell=1}^k V_n^{(\ell-1)}(\mathcal{C}_\ell \dots \mathcal{C}_k h) \right] \right|. \end{aligned}$$

Because $\mathbb{G}_n^{(0)}(\mathcal{C}_1 \dots \mathcal{C}_k h)$ has zero mean, it follows that

$$\sqrt{n} |\mathbb{E}_P [(P_n^{(k)}(h) - P(h)) \mathbf{1}_{\mathcal{S}}]| \leq |\mathbb{E}_P [\mathbf{1}_{\mathcal{S}^c} \mathbb{G}_n^{(0)}(\mathcal{C}_1 \dots \mathcal{C}_k h)]| + \sqrt{n} |\mathbb{E}_P [\mathbf{1}_{\mathcal{S}} T_2]|$$

We have by Hoeffding's inequality that $\mathbb{P}(\mathcal{E}_3^\delta) \geq 1 - \delta$, and that by Lemma B.3.2 that $\mathbb{G}_n^{(0)}(\mathcal{C}_1 \dots \mathcal{C}_k h) \leq 4k\sqrt{n} \|h\|_\infty$ universally. As a result, applying Proposition B.3.3 once again,

$$\begin{aligned} & |\mathbb{E}_P [\mathbf{1}_{\mathcal{S}^c} \mathbb{G}_n^{(0)}(\mathcal{C}_1 \dots \mathcal{C}_k h)]| \\ & \leq \left| \mathbb{E}_P [\mathbf{1}_{\mathcal{S}^c} \mathbf{1}_{\mathcal{E}_3^\delta} \mathbb{G}_n^{(0)}(\mathcal{C}_1 \dots \mathcal{C}_k h)] \right| + \left| \mathbb{E}_P [\mathbf{1}_{\mathcal{S}^c} \mathbf{1}_{\mathcal{E}_3^\delta} \mathbb{G}_n^{(0)}(\mathcal{C}_1 \dots \mathcal{C}_k h)] \right| \\ & \leq 4k\sqrt{n} \|h\|_\infty \min \{2m(1 - p_\star)^n, \delta\} + \sqrt{2 \log(2/\delta)} 2k \|h\|_\infty 2m(1 - p_\star)^n. \end{aligned}$$

Using a similar argument to Lemma B.3.4, we have that under $\mathcal{S} \setminus (\mathcal{E}_1^\delta \cap \mathcal{E}_2^\delta)$ (which occurs with probability no more than 2δ),

$$|T_2| \leq \|h\|_\infty mk \left[4n + \frac{k-1}{p_\star^2} \left(n + 2 + \frac{k+1}{3p_\star^2} \right) \right],$$

and that under $\mathcal{S} \cap \mathcal{E}_1^\delta \cap \mathcal{E}_2^\delta$ (which occurs with probability at least $1 - 2\delta$),

$$\begin{aligned} |T_2| & \leq \frac{4mk \|h\|_\infty [\log_2(2/\delta) + 2m \log(n+1)]^{1-\mathbf{1}_{\{k=1\}}/2}}{np_\star^2} \\ & \quad \left[p_\star \sqrt{2 \log(2mk/\delta)} (k+1) + (k-1)(k+4) \right]. \end{aligned}$$

Applying the decomposition $|\mathbb{E}_P [\mathbb{1}_{\mathcal{S}} T_2]| \leq \left| \mathbb{E}_P \left[\mathbb{1}_{\mathcal{S} \setminus (\mathcal{E}_1^\delta \cap \mathcal{E}_2^\delta)} T_2 \right] \right| + \left| \mathbb{E}_P \left[\mathbb{1}_{\mathcal{S} \cap \mathcal{E}_1^\delta \cap \mathcal{E}_2^\delta} T_2 \right] \right|$ and setting $\delta = \frac{p_\star^2}{n^2}$ achieves the desired result. \square

B.3.5 Misspecified Marginal Distributions

As described in Section 4.6.1, we now adapt the main results to cases in which the marginal distributions (P_X, P_Y) are misspecified, in that the user is provided marginal distributions $(\hat{P}_{X,\varepsilon}, \hat{P}_{Z,\varepsilon})$ which satisfy Assumption 4.6.1. The sequence $(\hat{P}_n^{(k)})_{k \geq 1}$ is generated via (4.26).

We start by deriving a result similar to Proposition B.3.1. Since $\varepsilon < 1$, the (possibly misspecified) target marginals $\hat{P}_{X,\varepsilon}(\mathbf{x}) > 0$ and $\hat{P}_{Z,\varepsilon}(\mathbf{z}) > 0$ for all $\mathbf{x} \in \mathcal{X}$ and $\mathbf{z} \in \mathcal{Z}$. Define the error term

$$\hat{V}_n^{(k-1)}(h) := \begin{cases} \sum_{\mathbf{x}, \mathbf{z}} \left(\frac{\hat{P}_{X,\varepsilon}}{\hat{P}_{n,X}^{(k-1)}}(\mathbf{x}) - 1 \right) h(\mathbf{x}, \mathbf{z}) \hat{P}_n^{(k-1)}(\mathbf{x}, \mathbf{z}) & k \text{ odd} \\ \sum_{\mathbf{x}, \mathbf{z}} \left(\frac{\hat{P}_{Z,\varepsilon}}{\hat{P}_{n,Z}^{(k-1)}}(\mathbf{z}) - 1 \right) h(\mathbf{x}, \mathbf{z}) \hat{P}_n^{(k-1)}(\mathbf{x}, \mathbf{z}) & k \text{ even} \end{cases} \quad (\text{B.46})$$

as well as the empirical process-style notation

$$\hat{\mathbb{G}}_n^{(k)}(h) := \sqrt{n} \left(\hat{P}_n^{(k)}(h) - P(h) \right).$$

The format of this section will be to derive results analogous to the building blocks of the previous section. From that point, the computations from Appendix B.3.4 will achieve the desired result. For the sake of comparison to Theorem 4.3.1 we consider error terms containing ε only by their dependence on $(\varepsilon, k, n, \hat{p}_{\star, \varepsilon})$.

Intermediate Results The following result provides the necessary recursion formula, although as an inequality rather than an equality.

Proposition B.3.4. *Let $(\hat{P}_n^{(k)})_{k \geq 1}$ be a sequence computed according to (4.26). Define*

$$c^2 = \max \left\{ \chi^2(\hat{P}_X \| P_X), \chi^2(\hat{P}_Z \| P_Y) \right\}.$$

These iterations are well-defined under the event \mathcal{S} , and for $\mathbb{G}_n^{(k)}$ defined in (B.27), it holds

that

$$\hat{\mathbb{G}}_n^{(k)}(h) = \hat{\mathbb{G}}_n^{(k-1)}(h) + \sqrt{n}\hat{V}_n^{(k-1)}(h) \quad (\text{B.47})$$

and

$$\hat{\mathbb{G}}_n^{(k)}(h) = \hat{\mathbb{G}}_n^{(k-1)}(\mathcal{C}_k h) + \sqrt{n}\hat{V}_n^{(k-1)}(\mathcal{C}_k h) + \begin{cases} \sqrt{n}[\hat{P}_{X,\varepsilon} - P_X](\mu_X h) & \text{if } k \text{ odd} \\ \sqrt{n}[\hat{P}_{Z,\varepsilon} - P_Y](\mu_Z h) & \text{if } k \text{ even} \end{cases}. \quad (\text{B.48})$$

Furthermore,

$$\begin{aligned} \left| \hat{\mathbb{G}}_n^{(k)}(h) \right| &\leq \left| \hat{\mathbb{G}}_n^{(k-1)}(\mathcal{C}_k h) \right| + \sqrt{n} \left| \hat{V}_n^{(k-1)}(\mathcal{C}_k h) \right| + \textcolor{red}{c} \|h\|_{\mathbf{L}^2(P)} \sqrt{n\varepsilon} \\ &= \left| \hat{\mathbb{G}}_n^{(k-1)}(\mathcal{C}_k h) \right| + \sqrt{n} \left| \hat{V}_n^{(k-1)}(\mathcal{C}_k h) \right| + \textcolor{red}{O}(\sqrt{n\varepsilon}). \end{aligned} \quad (\text{B.49})$$

Proof. The proof that $\hat{P}_{n,X}^{(k-1)}(\mathbf{x}) > 0$ and $\hat{P}_{n,Z}^{(k-1)}(\mathbf{z}) > 0$ for all $\mathbf{x} \in \mathcal{X}$ and $\mathbf{z} \in \mathcal{Z}$ follows the exact same steps as in the proof of Proposition B.3.1. We take this for granted and establish the recursion.

Consider the following steps in the case that k is odd:

$$\begin{aligned} \hat{P}_n^{(k)}(h) &= \sum_{\mathbf{x}, \mathbf{z}} h(\mathbf{x}, \mathbf{z}) \hat{P}_n^{(k)}(\mathbf{x}, \mathbf{z}) = \sum_{\mathbf{x}, \mathbf{z}} h(\mathbf{x}, \mathbf{z}) \frac{\hat{P}_{X,\varepsilon}}{\hat{P}_{n,X}^{(k-1)}}(\mathbf{x}) \hat{P}_n^{(k-1)}(\mathbf{x}, \mathbf{z}) \\ &= \sum_{\mathbf{x}, \mathbf{z}} \textcolor{blue}{1} \cdot h(\mathbf{x}, \mathbf{z}) \hat{P}_n^{(k-1)}(\mathbf{x}, \mathbf{z}) + \sum_{\mathbf{x}, \mathbf{z}} \left[\frac{\hat{P}_{X,\varepsilon}}{\hat{P}_{n,X}^{(k-1)}}(\mathbf{x}) - 1 \right] \cdot h(\mathbf{x}, \mathbf{z}) \hat{P}_n^{(k-1)}(\mathbf{x}, \mathbf{z}) \\ &= \hat{P}_n^{(k-1)}(h) + \hat{V}_n^{(k-1)}(h). \end{aligned}$$

Subtracting $P(h)$ on both sides, we have that

$$[\hat{P}_n^{(k)} - P](h) = [\hat{P}_n^{(k-1)} - P](h) + \hat{V}_n^{(k-1)}(h). \quad (\text{B.50})$$

This proves the uncentered recursion formula given in (B.47). We then show the centered

version.

$$\begin{aligned}
& [\hat{P}_n^{(k)} - P](h) \\
&= [\hat{P}_n^{(k)} - P](\mathcal{C}_X h) + [\hat{P}_n^{(k)} - P](\mu_X h) \\
&= [\hat{P}_n^{(k)} - P](\mathcal{C}_X h) + [\hat{P}_{X,\varepsilon} - P_X](\mu_X h) \\
&= [\hat{P}_n^{(k-1)} - P](\mathcal{C}_X h) + \hat{V}_n^{(k-1)}(\mathcal{C}_X h) + [\hat{P}_{X,\varepsilon} - P_X](\mu_X h).
\end{aligned}$$

Next, we bound the additional error term. By the Cauchy-Schwarz inequality,

$$\begin{aligned}
[\hat{P}_{X,\varepsilon} - P_X](\mu_X h) &\leq \left\| \frac{\hat{P}_{X,\varepsilon}}{P_X} - \mathbf{1} \right\|_{\mathbf{L}^2(P_X)} \cdot \|\mu_X h\|_{\mathbf{L}^2(P_X)} \\
&= \sqrt{\chi^2(\hat{P}_{X,\varepsilon} \| P_X)} \cdot \|\mu_X h\|_{\mathbf{L}^2(P_X)} \\
&\leq \sqrt{\chi^2(\hat{P}_{X,\varepsilon} \| P_X)} \cdot \|h\|_{\mathbf{L}^2(P)},
\end{aligned}$$

as μ_X is an orthogonal projection in $\mathbf{L}^2(P)$. Using convexity of f -divergences, we have that

$$\chi^2(\hat{P}_{X,\varepsilon} \| P_X) \leq \varepsilon \chi^2(\hat{P}_X \| P_X) + (1 - \varepsilon) \chi^2(P_X \| P_X) = \varepsilon \chi^2(\hat{P}_X \| P_X).$$

This achieves the desired result. □

Using similar ideas, we then prove an analog of Lemma B.3.3.

Lemma B.3.6. *For k odd, it holds that*

$$\sqrt{n} \hat{V}_n^{(k-1)}(\mathcal{C}_X h) = \sum_{j=1}^m \left(\frac{\hat{P}_{X,\varepsilon}}{\hat{P}_{n,X}^{(k-1)}}(\mathbf{x}_j) - 1 \right) \hat{\mathbb{G}}_n^{(k-1)}(\mathcal{C}_X h \mathbf{1}_{jk}),$$

whereas for k even, it holds that

$$\sqrt{n} \hat{V}_n^{(k-1)}(\mathcal{C}_Z h) = \sum_{j=1}^m \left(\frac{\hat{P}_{Z,\varepsilon}}{\hat{P}_{n,Z}^{(k-1)}}(\mathbf{x}_j) - 1 \right) \hat{\mathbb{G}}_n^{(k-1)}(\mathcal{C}_Z h \mathbf{1}_{jk}).$$

Proof. We give the proof for k odd. We claim that we need only show that $P(\mathcal{C}_X h \mathbf{1}_{jk}) = 0$.

This would show that

$$\begin{aligned}
\sqrt{n}\hat{V}_n^{(k-1)}(\mathcal{C}_X h) &= \sum_{\mathbf{x}, \mathbf{z}} \left(\frac{\hat{P}_{X,\varepsilon}}{\hat{P}_{n,X}^{(k-1)}}(\mathbf{x}) - 1 \right) [\mathcal{C}_X h](\mathbf{x}, \mathbf{z}) \hat{P}_n^{(k-1)}(\mathbf{x}, \mathbf{z}) \\
&= \sqrt{n} \sum_{j=1}^m \sum_{\mathbf{x}, \mathbf{z}} \left(\frac{\hat{P}_{X,\varepsilon}}{\hat{P}_{n,X}^{(k-1)}}(\mathbf{x}) - 1 \right) [\mathcal{C}_X h \mathbf{1}_{jk}](\mathbf{x}, \mathbf{z}) \hat{P}_n^{(k-1)}(\mathbf{x}, \mathbf{z}) \\
&= \sqrt{n} \sum_{j=1}^m \sum_{\mathbf{x}, \mathbf{z}} \left(\frac{\hat{P}_{X,\varepsilon}}{\hat{P}_{n,X}^{(k-1)}}(\mathbf{x}_j) - 1 \right) [\mathcal{C}_X h \mathbf{1}_{jk}](\mathbf{x}, \mathbf{z}) \hat{P}_n^{(k-1)}(\mathbf{x}, \mathbf{z}) \\
&= \sum_{j=1}^m \left(\frac{\hat{P}_{X,\varepsilon}}{\hat{P}_{n,X}^{(k-1)}}(\mathbf{x}_j) - 1 \right) \sqrt{n} \sum_{\mathbf{x}, \mathbf{z}} [\mathcal{C}_X h \mathbf{1}_{jk}](\mathbf{x}, \mathbf{z}) \hat{P}_n^{(k-1)}(\mathbf{x}, \mathbf{z}) \\
&= \sum_{j=1}^m \left(\frac{\hat{P}_{X,\varepsilon}}{\hat{P}_{n,X}^{(k-1)}}(\mathbf{x}_j) - 1 \right) \hat{\mathbb{G}}_n^{(k-1)}(\mathcal{C}_X h \mathbf{1}_{jk}),
\end{aligned}$$

where $P(\mathcal{C}_X h \mathbf{1}_{jk}) = 0$ is employed in the last step. Now the result follows from (B.48) in Proposition B.3.4 and the definition of $\hat{\mathbb{G}}_n^{(k)}(h)$. To prove the claim, as in Lemma B.3.3, write

$$\mathbb{E}[\mathcal{C}_X h \mathbf{1}_{jk} | X](\mathbf{x}) = \begin{cases} \mathbb{E}[\mathcal{C}_X h | X](\mathbf{x}_j) & \text{if } x = x_j \\ 0 & \text{if } x \neq x_j \end{cases}.$$

But $\mathbb{E}[\mathcal{C}_X h | X](\mathbf{x}_j) = 0$ by definition of \mathcal{C}_X . Taking an expectation over P_X gives that $P(\mathcal{C}_X h \mathbf{1}_{jk}) = 0$, which implies the desired result. The proof for k even follows symmetrically. \square

For the remainder of the argument, we see that (B.49) can be unrolled so that

$$\left| \hat{\mathbb{G}}_n^{(k)}(h) \right| \leq \underbrace{|\mathbb{G}_n^{(0)}(\mathcal{C}_1 \dots \mathcal{C}_k h)|}_{\text{first-order term}} + \underbrace{\sqrt{n} \sum_{\ell=1}^k \left| \hat{V}_n^{(\ell-1)}(\mathcal{C}_\ell \dots \mathcal{C}_k h) \right|}_{\text{higher-order term}} + \underbrace{O(k\sqrt{n\varepsilon})}_{\text{misspecification}}, \quad (\text{B.51})$$

where we use that $\mathbb{G}_n^{(0)} = \hat{\mathbb{G}}_n^{(0)}$.

Next, we need to bound $\left| \hat{V}_n^{(\ell-1)}(\mathcal{C}_\ell \dots \mathcal{C}_k h) \right|$, in particular accounting for the marginal violation term. We follow similar steps as in the analysis of the higher-order term in Appendix B.3.3.

Proposition B.3.5. Assume that $P_{n,X}(\mathbf{x}) > 0$ for all $\mathbf{x} \in \mathcal{X}$. It holds that

$$\max_{\mathbf{x} \in \mathcal{X}} \left| \frac{\hat{P}_{X,\varepsilon}(\mathbf{x})}{\hat{P}_{n,X}^{(k-1)}(\mathbf{x})} - 1 \right| \leq \begin{cases} \max\{n-1, 1\} & \text{if } k = 1 \\ \max\{1/\hat{p}_{\star,\varepsilon}^2 - 1, 1\} & \text{if } k > 1. \end{cases} \quad (\text{B.52})$$

In addition, we have that

$$\max_{\mathbf{x} \in \mathcal{X}} \left| \frac{\hat{P}_{X,\varepsilon}(\mathbf{x})}{\hat{P}_{n,X}^{(k-1)}(\mathbf{x})} - 1 \right| \leq \begin{cases} O\left(n\sqrt{\log \frac{1}{1-\varepsilon}}\right) + n\sqrt{\frac{1}{2}\text{KL}(P_{n,X}\|P_X)} & \text{if } k = 1 \\ O\left(\frac{1}{\hat{p}_{\star,\varepsilon}^2}\sqrt{\log \frac{1}{1-\varepsilon}}\right) + \frac{1}{\hat{p}_{\star,\varepsilon}^2}\sqrt{\frac{1}{2}\text{KL}(P_{n,X}\|P_X)} & \text{if } k > 1 \end{cases},$$

Moreover, when $\text{KL}(P_{n,X}\|P_X) \leq \frac{\hat{p}_{\star,\varepsilon}^2}{8}$ and $\varepsilon \leq 1 - \exp\left(-\frac{\hat{p}_{\star,\varepsilon}^2}{8}\right)$, we have

$$\max_{\mathbf{x} \in \mathcal{X}} \left| \frac{\hat{P}_{X,\varepsilon}(\mathbf{x})}{\hat{P}_{n,X}^{(k-1)}(\mathbf{x})} - 1 \right| \leq O\left(\frac{1}{\hat{p}_{\star,\varepsilon}}\sqrt{\log \frac{1}{1-\varepsilon}}\right) + \frac{2}{\hat{p}_{\star,\varepsilon}}\sqrt{\frac{1}{2}\text{KL}(P_{n,X}\|P_X)}. \quad (\text{B.53})$$

Proof. First, observe that $\hat{P}_{n,X}^{(0)}(\mathbf{x}) = P_{n,X}^{(0)}(\mathbf{x}) \geq 1/n$ under the event \mathcal{S} . For $k > 1$ such that k is odd, we have that for $\mathbf{x} \in \mathcal{X}$,

$$\begin{aligned} \hat{P}_{n,X}^{(k-1)}(\mathbf{x}) &= \sum_{\mathbf{z} \in \mathcal{Z}} \hat{P}_n^{(k-1)}(\mathbf{x}, \mathbf{z}) = \sum_{\mathbf{z} \in \mathcal{Z}} \frac{\hat{P}_{Z,\varepsilon}(\mathbf{z})}{\hat{P}_{n,Z}^{(k-2)}(\mathbf{z})} \hat{P}_n^{(k-2)}(\mathbf{x}, \mathbf{z}) \\ &\geq \hat{p}_{\star,\varepsilon} \sum_{\mathbf{z} \in \mathcal{Z}} \hat{P}_n^{(k-2)}(\mathbf{x}, \mathbf{z}) = \hat{p}_{\star,\varepsilon} \hat{P}_{n,X}^{(k-2)}(\mathbf{x}) = \hat{p}_{\star,\varepsilon} \hat{P}_{X,\varepsilon}(\mathbf{x}) \geq \hat{p}_{\star,\varepsilon}^2. \end{aligned}$$

The result for k even can be proven similarly. We now prove the inequalities listed in the statement using on the lower bounds above.

Proving the first inequality. For any $\mathbf{x} \in \mathcal{X}$,

$$\left| \frac{\hat{P}_{X,\varepsilon}(\mathbf{x})}{\hat{P}_{n,X}^{(k-1)}(\mathbf{x})} - 1 \right| = \max \left\{ \frac{\hat{P}_{X,\varepsilon}(\mathbf{x})}{\hat{P}_{n,X}^{(k-1)}(\mathbf{x})} - 1, 1 - \frac{\hat{P}_{X,\varepsilon}(\mathbf{x})}{\hat{P}_{n,X}^{(k-1)}(\mathbf{x})} \right\} \leq \begin{cases} \max\{n-1, 1\} & \text{if } k = 1 \\ \max\{1/\hat{p}_{\star,\varepsilon}^2 - 1, 1\} & \text{if } k > 1 \end{cases},$$

which is the desired result.

Proving the second and third inequalities. Consider an odd $k \geq 1$. By the definition

of total variation distance, it holds that

$$\max_{\mathbf{x} \in \mathcal{X}} \left| \hat{P}_{X,\varepsilon}(\mathbf{x}) - \hat{P}_{n,X}^{(k-1)}(\mathbf{x}) \right| \leq \text{TV}(\hat{P}_{n,X}^{(k-1)}, \hat{P}_{X,\varepsilon}).$$

According to Pinsker's inequality, we have that $\text{TV}(\hat{P}_{n,X}^{(k-1)}, \hat{P}_{X,\varepsilon}) \leq \sqrt{\frac{1}{2} \text{KL}(\hat{P}_{n,X}^{(k-1)} \parallel \hat{P}_{X,\varepsilon})}$, and so we have that

$$\max_{\mathbf{x} \in \mathcal{X}} \left| \hat{P}_{X,\varepsilon}(\mathbf{x}) - \hat{P}_{n,X}^{(k-1)}(\mathbf{x}) \right| \leq \sqrt{\frac{1}{2} \text{KL}(\hat{P}_{n,X}^{(k-1)} \parallel \hat{P}_{X,\varepsilon})} \leq \sqrt{\frac{1}{2} \text{KL}(P_{n,X}^{(0)} \parallel \hat{P}_{X,\varepsilon})},$$

where the last inequality follows by the monotonicity of Sinkhorn iterations given in Proposition B.2.4. Notice that the remaining term is $\text{KL}(P_{n,X}^{(0)} \parallel \hat{P}_{X,\varepsilon}) = \text{KL}(P_{n,X} \parallel \hat{P}_{X,\varepsilon})$, which may not decay to zero as $n \rightarrow \infty$. Because $\varepsilon < 1$, write

$$\begin{aligned} \text{KL}(P_{n,X} \parallel \hat{P}_{X,\varepsilon}) &= \sum_{\mathbf{x} \in \mathcal{X}} P_{n,X}(\mathbf{x}) \log \frac{P_{n,X}(\mathbf{x})}{(1-\varepsilon)P_X(\mathbf{x}) + \varepsilon \hat{P}_X(\mathbf{x})} \\ &\leq \sum_{\mathbf{x} \in \mathcal{X}} P_{n,X}(\mathbf{x}) \log \frac{P_{n,X}(\mathbf{x})}{(1-\varepsilon)P_X(\mathbf{x})} \\ &= \text{KL}(P_{n,X} \parallel P_X) + \log \frac{1}{1-\varepsilon} \\ \Rightarrow \sqrt{\frac{1}{2} \text{KL}(P_{n,X} \parallel \hat{P}_{X,\varepsilon})} &\leq \sqrt{\frac{1}{2} \text{KL}(P_{n,X} \parallel P_X)} + \sqrt{\frac{1}{2} \log \frac{1}{1-\varepsilon}}. \end{aligned}$$

We can then apply the lower bounds

$$\max_{\mathbf{x} \in \mathcal{X}} \left| \frac{\hat{P}_{X,\varepsilon}(\mathbf{x})}{\hat{P}_{n,X}^{(k-1)}(\mathbf{x})} - 1 \right| \leq \begin{cases} n \left(\sqrt{\frac{1}{2} \text{KL}(P_{n,X} \parallel P_X)} + \sqrt{\frac{1}{2} \log \frac{1}{1-\varepsilon}} \right) & \text{if } k = 1 \\ \frac{1}{\hat{p}_{\star,\varepsilon}^2} \left(\sqrt{\frac{1}{2} \text{KL}(P_{n,X} \parallel P_X)} + \sqrt{\frac{1}{2} \log \frac{1}{1-\varepsilon}} \right) & \text{if } k > 1 \end{cases}.$$

Finally, combining the arguments above, we have that

$$\begin{aligned} \max_{\mathbf{x} \in \mathcal{X}} \left| \hat{P}_{X,\varepsilon}(\mathbf{x}) - \hat{P}_{n,X}^{(k-1)}(\mathbf{x}) \right| &\leq \sqrt{\frac{1}{2} \text{KL}(P_{n,X} \parallel P_X)} + \sqrt{\frac{1}{2} \log \frac{1}{1-\varepsilon}} \\ &\leq \frac{\hat{p}_{\star,\varepsilon}}{4} + \frac{\hat{p}_{\star,\varepsilon}}{4} = \frac{\hat{p}_{\star,\varepsilon}}{2}, \end{aligned}$$

where the last step invoked the assumption that

$$\text{KL}(P_{n,X} \| P_X) \leq \frac{\hat{p}_{\star,\varepsilon}^2}{8} \quad \text{and} \quad \varepsilon \leq 1 - \exp\left(-\frac{\hat{p}_{\star,\varepsilon}^2}{8}\right).$$

This means that

$$\min_{\mathbf{x} \in \mathcal{X}} \hat{P}_{n,X}^{(k-1)}(\mathbf{x}) \geq \min_{\mathbf{x} \in \mathcal{X}} \hat{P}_{X,\varepsilon}(\mathbf{x}) - \max_{\mathbf{x} \in \mathcal{X}} \left| \hat{P}_{n,X}^{(k-1)}(\mathbf{x}) - \hat{P}_{X,\varepsilon}(\mathbf{x}) \right| \geq \frac{\hat{p}_{\star,\varepsilon}}{2}.$$

Hence,

$$\max_{\mathbf{x} \in \mathcal{X}} \left| \frac{\hat{P}_{X,\varepsilon}(\mathbf{x})}{\hat{P}_{n,X}^{(k-1)}(\mathbf{x})} - 1 \right| \leq \frac{\max_{\mathbf{x} \in \mathcal{X}} \left| \hat{P}_{n,X}^{(k-1)}(\mathbf{x}) - \hat{P}_{X,\varepsilon}(\mathbf{x}) \right|}{\min_{\mathbf{x} \in \mathcal{X}} \hat{P}_{n,X}^{(k-1)}(\mathbf{x})} \leq \frac{2}{\hat{p}_{\star,\varepsilon}} \sqrt{\frac{1}{2} \text{KL}(P_{n,X} \| \hat{P}_{X,\varepsilon})}.$$

Now, for k even, set $k = 2t$ for $t \geq 0$. We have that

$$\max_{\mathbf{z} \in \mathcal{Z}} \left| \hat{P}_{n,Z}^{(2t-1)}(\mathbf{z}) - \hat{P}_{Z,\varepsilon}(\mathbf{z}) \right| \leq \text{TV}(\hat{P}_{n,Z}^{(2t-1)}, \hat{P}_{Z,\varepsilon}) \leq \sqrt{\frac{1}{2} \text{KL}(\hat{P}_{Z,\varepsilon} \| \hat{P}_{n,Z}^{(2t-1)})}.$$

Invoke Proposition B.2.4 once again to achieve

$$\sqrt{\frac{1}{2} \text{KL}(\hat{P}_{Z,\varepsilon} \| \hat{P}_{n,Z}^{(2t-1)})} \leq \sqrt{\frac{1}{2} \text{KL}(P_{n,X} \| \hat{P}_{X,\varepsilon})} \leq \sqrt{\frac{1}{2} \text{KL}(P_{n,X} \| P_X)} + \sqrt{\frac{1}{2} \log \frac{1}{1-\varepsilon}},$$

which completes the proof. \square

Proceeding with similar steps, define the quantities

$$\hat{B}_1 := \hat{M}_1 \quad \text{and} \quad \hat{B}_2 := \max_{2 \leq \ell \leq k} \hat{M}_\ell \quad \text{for} \quad \hat{M}_\ell := \begin{cases} \max_{\mathbf{x} \in \mathcal{X}} \left| \frac{\hat{P}_{X,\varepsilon}(\mathbf{x})}{\hat{P}_{n,X}^{(\ell-1)}(\mathbf{x})} - 1 \right| & \ell \text{ odd} \\ \max_{\mathbf{z} \in \mathcal{Z}} \left| \frac{\hat{P}_{Z,\varepsilon}(\mathbf{z})}{\hat{P}_{n,Z}^{(\ell-1)}(\mathbf{z})} - 1 \right| & \ell \text{ even} \end{cases}.$$

We must now establish an analog of Proposition B.3.2.

Proposition B.3.6. *For any $k \geq 1$, the following holds under the event \mathcal{S} :*

$$\begin{aligned} \sqrt{n} \sum_{\ell=1}^k \left| \hat{V}_n^{(\ell-1)}(\mathcal{C}_\ell \dots \mathcal{C}_k h) \right| &\leq \sum_{j=1}^m \left(\hat{B}_1 |\mathbb{G}_n^{(0)}(h_{1,k} \mathbf{1}_{j\ell})| + \hat{B}_2 \sum_{\ell=2}^k |\mathbb{G}_n^{(0)}(h_{\ell,k} \mathbf{1}_{j\ell})| \right) \\ &\quad + m \hat{B}_2 \|h\|_\infty \sqrt{n} k(k-1) [\hat{B}_1 + \hat{B}_2(k+1)/3]. \end{aligned}$$

Proof. This proof largely follows the argument of Proposition B.3.2, while accounting for the misspecified marginal error. Using again the notation $h_{\ell,k} := \mathcal{C}_\ell \dots \mathcal{C}_k h$, it follows from Lemma B.3.6 that, for odd ℓ ,

$$\sqrt{n} \hat{V}_n^{(\ell-1)}(h_{\ell,k}) = \sum_{j=1}^m \left[\frac{\hat{P}_{X,\varepsilon}}{\hat{P}_{n,X}^{(\ell-1)}}(\mathbf{x}_j) - 1 \right] \hat{\mathbb{G}}_n^{(\ell-1)}(h_{\ell,k} \mathbf{1}_{j\ell}) \leq \hat{M}_\ell \sum_{j=1}^m \left| \hat{\mathbb{G}}_n^{(\ell-1)}(h_{\ell,k} \mathbf{1}_{j\ell}) \right|.$$

The bound above holds for ℓ even as well. Then, using (B.47) from Proposition B.3.4 along with the triangle inequality, we have that for $\ell \geq 2$,

$$\left| [\hat{P}_n^{(\ell-1)} - P](h_{\ell,k} \mathbf{1}_{j\ell}) \right| \leq \left| [\hat{P}_n^{(\ell-2)} - P](h_{\ell,k} \mathbf{1}_{j\ell}) \right| + \left| \hat{V}_n^{(\ell-2)}(h_{\ell,k} \mathbf{1}_{j\ell}) \right|$$

which implies that

$$\left| \hat{\mathbb{G}}_n^{(\ell-1)}(h_{\ell,k} \mathbf{1}_{j\ell}) \right| \tag{B.54}$$

$$\begin{aligned} &\leq \left| \hat{\mathbb{G}}_n^{(\ell-2)}(h_{\ell,k} \mathbf{1}_{j\ell}) \right| + \sqrt{n} \left| \hat{V}_n^{(\ell-2)}(h_{\ell,k} \mathbf{1}_{j\ell}) \right| \\ &\leq \left| \hat{\mathbb{G}}_n^{(0)}(h_{\ell,k} \mathbf{1}_{j\ell}) \right| + \sqrt{n} \left| \hat{V}_n^{(0)}(h_{\ell,k} \mathbf{1}_{j\ell}) \right| + \dots + \sqrt{n} \left| \hat{V}_n^{(\ell-2)}(h_{\ell,k} \mathbf{1}_{j\ell}) \right| \\ &\leq \left| \hat{\mathbb{G}}_n^{(0)}(h_{\ell,k} \mathbf{1}_{j\ell}) \right| + \hat{M}_1 \sqrt{n} \hat{P}_n^{(0)}(|h_{\ell,k}| \mathbf{1}_{j\ell}) + \dots + \hat{M}_\ell \sqrt{n} \hat{P}_n^{(\ell-2)}(|h_{\ell,k}| \mathbf{1}_{j\ell}) \\ &\leq \left| \hat{\mathbb{G}}_n^{(0)}(h_{\ell,k} \mathbf{1}_{j\ell}) \right| + 2 \|h\|_\infty \sqrt{n} \left[\hat{B}_1 + \hat{B}_2(\ell-1) \right] (k - \ell + 1), \end{aligned} \tag{B.55}$$

by Lemma B.3.2 and $\hat{M}_1 \leq \hat{B}_1$ and $\hat{M}_\ell \leq \hat{B}_2$ for $\ell \geq 2$. The bound above holds trivially for

$\ell = 1$. Summing these bounds over ℓ and j , we have that

$$\begin{aligned}
& \sqrt{n} \sum_{\ell=1}^k \left| \hat{V}_n^{(\ell-1)}(h_{\ell,k}) \right| \\
& \leq \hat{M}_1 \sum_{j=1}^m |\mathbb{G}_n^{(0)}(h_{1,k} \mathbf{1}_{j\ell})| + \sum_{\ell=2}^k \hat{M}_\ell \sum_{j=1}^m \left| \hat{\mathbb{G}}_n^{(\ell-1)}(h_{\ell,k} \mathbf{1}_{j\ell}) \right| \\
& \leq \hat{B}_1 \sum_{j=1}^m |\mathbb{G}_n^{(0)}(h_{1,k} \mathbf{1}_{j\ell})| + \hat{B}_2 \sum_{\ell=2}^k \sum_{j=1}^m \left| \hat{\mathbb{G}}_n^{(\ell-1)}(h_{\ell,k} \mathbf{1}_{j\ell}) \right| \\
& \leq \hat{B}_1 \sum_{j=1}^m |\mathbb{G}_n^{(0)}(h_{1,k} \mathbf{1}_{j\ell})| \\
& \quad + \hat{B}_2 \sum_{\ell=2}^k \sum_{j=1}^m \left(\left| \hat{\mathbb{G}}_n^{(0)}(h_{\ell,k} \mathbf{1}_{j\ell}) \right| + 2 \|h\|_\infty \sqrt{n} \left[\hat{B}_1 + \hat{B}_2(\ell-1) \right] (k-\ell+1) \right) \quad \text{apply (B.55)} \\
& = \sum_{j=1}^m \left(\hat{B}_1 |\mathbb{G}_n^{(0)}(h_{1,k} \mathbf{1}_{j\ell})| + \hat{B}_2 \sum_{\ell=2}^k |\mathbb{G}_n^{(0)}(h_{\ell,k} \mathbf{1}_{j\ell})| \right) \\
& \quad + 2m\hat{B}_2 \|h\|_\infty \sqrt{n} \sum_{\ell=2}^k \left[\hat{B}_1 + \hat{B}_2(\ell-1) \right] (k-\ell+1),
\end{aligned}$$

because $|\mathcal{X}| = m$. We sum up the last term:

$$\begin{aligned}
\sum_{\ell=2}^k \left[\hat{B}_1 + \hat{B}_2(\ell-1) \right] (k-\ell+1) &= \hat{B}_1 \sum_{\ell=1}^{k-1} (k-\ell) + \hat{B}_2 \sum_{\ell=1}^{k-1} \ell(k-\ell) \\
&= \frac{k(k-1)}{2} \left[\hat{B}_1 + \hat{B}_2(k+1)/3 \right],
\end{aligned}$$

which completes the proof. \square

Mean Squared Error Bound Ultimately, we wish to construct an upper bound for

$$\mathbb{E}_P \left[\left(\hat{P}_n^{(k)}(h) - P(h) \right)^2 \mathbb{1}_{\mathcal{S}} \right] + \mathbb{E}_P \left[(P_n(h) - P(h))^2 \mathbb{1}_{\mathcal{S}^c} \right], \quad (\text{B.56})$$

as the method returns $P_n(h)$ when \mathcal{S} is not satisfied. The first term will be controlled by intermediate tools developed above. The second term that includes \mathcal{S}^c is no different from the one analyzed in Proposition B.3.3. We handle the second term first. Recall from

Proposition B.3.3 that for any $\delta \in (0, 1)$,

$$\begin{aligned} \mathbb{E}_P \left[(P_n(h) - P(h))^2 \mathbf{1}_{S^c} \right] &\leq \\ 4 \|h\|_\infty^2 \min \{2m(1 - p_\star)^n, \delta\} &+ \frac{2 \log(2/\delta)}{n} \|h\|_\infty^2 2m(1 - p_\star)^n. \end{aligned} \quad (\text{B.57})$$

Repeat the argument from the proof of Theorem B.3.2: because $2[\log_2(2/\delta) + m \log(n+1)] \geq \log(m/\delta)$ and $-\log(1 - p_\star) \geq p_\star \geq p_\star^2$, we have that

$$n \geq 2[\log_2(2/\delta) + m \log(n+1)]/p_\star^2 \implies n \geq \log(\delta/m)/\log(1 - p_\star). \quad (\text{B.58})$$

This in turn implies that $m(1 - p_\star)^n \leq \delta$, and gives as a condition on the sample size n . Further in the analysis, we will set $\delta = (\hat{p}_{\star, \varepsilon}/n)^4$, so right-hand side of (B.57) can then be upper bounded further, resulting in

$$\mathbb{E}_P \left[(P_n(h) - P(h))^2 \mathbf{1}_{S^c} \right] \leq 4 \|h\|_\infty^2 \delta \left(2 + \frac{\log(2/\delta)}{n} \right) = \tilde{O} \left(\frac{\hat{p}_{\star, \varepsilon}^4}{n^4} \right),$$

a higher-order term compared to other components of the bound.

Next, we must control the left-hand side of (B.56). We perform the decomposition based on (B.51):

$$\begin{aligned} \mathbb{E}_P \left[\left(\hat{P}_n^{(k)}(h) - P(h) \right)^2 \mathbf{1}_S \right] \\ \leq \mathbb{E}_P \left[T_1^2 \mathbf{1}_S \right] + 2\mathbb{E}_P \left[\left| T_1 \hat{T}_2 \right| \mathbf{1}_S \right] + \mathbb{E}_P \left[\hat{T}_2^2 \mathbf{1}_S \right] \end{aligned} \quad (\text{B.59})$$

$$+ O(k\sqrt{\varepsilon}) \cdot \mathbb{E}_P \left[\left(|T_1| + |\hat{T}_2| \right) \mathbf{1}_S \right] + O(k^2\varepsilon) \quad (\text{B.60})$$

for

$$T_1 := [P_n - P](\mathcal{C}_1 \dots \mathcal{C}_k h) \text{ and } \hat{T}_2 := \sum_{\ell=1}^k \left| \hat{V}_n^{(\ell-1)}(\mathcal{C}_\ell \dots \mathcal{C}_k h) \right|. \quad (\text{B.61})$$

Recall the events \mathcal{E}_1^δ and \mathcal{E}_2^δ and \mathcal{E}_3^δ from Appendix B.3.4. To perform this computation efficiently, we will split the bounds on each term into two components. In particular, we will show that

- Under the event $\mathcal{S} \cap \mathcal{E}_1^\delta \cap \mathcal{E}_2^\delta : |\hat{T}_2| \leq \mathcal{T}_2 + E_2$,
- Under the event $\mathcal{S} \setminus (\mathcal{E}_1^\delta \cap \mathcal{E}_2^\delta) : |\hat{T}_2| \leq \mathcal{T}_2^c + E_2^c$,
- Under the event $\mathcal{S} \cap \mathcal{E}_3^\delta : |T_1| \leq \mathcal{T}_1$,
- Under the event $\mathcal{S} \setminus \mathcal{E}_3^\delta : |T_1| \leq \mathcal{T}_1^c$,

where any term denoted with “ E ” will represent all error terms that include ε and will be written in big- O notation. There are no errors for the bounds on T_1 , as this term does not depend on the misspecified marginals. The idea is that for the “ \mathcal{T}_2 ” terms we may reuse the bounds derived in Appendix B.3.4 by simply replacing p_\star with $\hat{p}_{\star, \varepsilon}$. This is due to the fact that the dependence of the analogous terms from Appendix B.3.4 depend on p_\star only through Proposition B.2.5; similarly, the corresponding terms in this section depend on $\hat{p}_{\star, \varepsilon}$ through Proposition B.3.5. We return to the terms in (B.59) and (B.60).

Decomposing on \mathcal{E}_3^δ will result in a bound of the form

$$O(k\sqrt{\varepsilon}) \cdot \mathbb{E}_P [|T_1| \mathbb{1}_\mathcal{S}] \leq O(k\sqrt{\varepsilon}) \cdot (\delta \mathcal{T}_1^c + \mathcal{T}_1).$$

Decomposing on $\mathcal{E}_1^\delta \cap \mathcal{E}_2^\delta$ will result in a bounds of the form

$$\begin{aligned} \mathbb{E}_P [\hat{T}_2^2 \mathbb{1}_\mathcal{S}] &\leq 2\delta(\mathcal{T}_2^c)^2 + \mathcal{T}_2^2 + \tilde{O}(\delta((E_2^c)^2 + E_2^c \mathcal{T}_2^c) + (E_2^2 + E_2 \mathcal{T}_2)) \\ O(k\sqrt{\varepsilon}) \cdot \mathbb{E}_P [|\hat{T}_2| \mathbb{1}_\mathcal{S}] &\leq O(k\sqrt{\varepsilon}) \cdot (\delta(\mathcal{T}_2^c + E_2^c) + \mathcal{T}_2 + E_2). \end{aligned}$$

Finally, decomposing on $\mathcal{E}_1^\delta \cap \mathcal{E}_2^\delta \cap \mathcal{E}_3^\delta$ will result in a bound of the form

$$\mathbb{E}_P [|T_1 \hat{T}_2| \mathbb{1}_\mathcal{S}] \leq 3\delta \mathcal{T}_1^c \mathcal{T}_2^c + \mathcal{T}_1 \mathcal{T}_2 + \tilde{O}(\delta \mathcal{T}_1^c E_2^c + \mathcal{T}_1 E_2).$$

The leading terms $2\delta(\mathcal{T}_2^c)^2 + \mathcal{T}_2^2$ and $3\delta \mathcal{T}_1^c \mathcal{T}_2^c + \mathcal{T}_1 \mathcal{T}_2$ from both bounds should have the exact same form as the terms in Lemma B.3.4 and Lemma B.3.5, with p_\star replaced by $\hat{p}_{\star, \varepsilon}$, thus retaining the same dependence on (n, k) . By setting $\delta = \hat{p}_{\star, \varepsilon}^4 / n^4$, we will achieve a similar

result to Theorem B.3.2, i.e., that

$$\begin{aligned} & \mathbb{E}_P \left[\left(\hat{P}_n^{(k)}(h) - P(h) \right)^2 \mathbb{1}_{\mathcal{S}} \right] \\ & \leq \frac{\sigma_k^2}{n} + \tilde{O} \left(\frac{k^6}{n^{3/2}} \right) \\ & + \tilde{O} \left((\hat{p}_{*,\varepsilon}/n)^4 (E_2^c(E_2^c + \mathcal{T}_2^c)) + E_2(E_2 + \mathcal{T}_2) + (\hat{p}_{*,\varepsilon}/n)^4 \mathcal{T}_1^c E_2^c + \mathcal{T}_1 E_2 \right). \end{aligned} \quad (\text{B.62})$$

$$+ \tilde{O} \left(k\sqrt{\varepsilon} \left((\hat{p}_{*,\varepsilon}/n)^4 \mathcal{T}_1^c + \mathcal{T}_1 + (\hat{p}_{*,\varepsilon}/n)^4 (\mathcal{T}_2^c + E_2^c) + \mathcal{T}_2 + E_2 \right) + k^2 \varepsilon \right). \quad (\text{B.63})$$

It remains to quantify the \tilde{O} terms by computing the order of the 6 constants $(\mathcal{T}_2, E_2, \mathcal{T}_2^c, E_2^c, \mathcal{T}_1, \mathcal{T}_1^c)$. We follow similar steps to Lemma B.3.4 and Lemma B.3.5 to achieve this.

Lemma B.3.7. *For $\delta = (\hat{p}_{*,\varepsilon}/n)^4$, assume that $n \geq 8[\log_2(2/\delta) + m \log(n+1)]/\hat{p}_{*,\varepsilon}^2$ and $\varepsilon \leq 1 - \exp\left(-\frac{\hat{p}_{*,\varepsilon}^2}{8}\right)$. Then, it holds that*

$$\begin{aligned} \mathcal{T}_2^c &= \tilde{O} \left(\frac{k^2}{\hat{p}_{*,\varepsilon}^2} \left(n + \frac{k}{\hat{p}_{*,\varepsilon}} \right) \right), \quad E_2^c = 0 \\ \mathcal{T}_2 &= \tilde{O} \left(\frac{k^3}{n\hat{p}_{*,\varepsilon}^2} \right), \quad E_2 = \tilde{O} \left(\frac{k^3}{\hat{p}_{*,\varepsilon}^2} \left(\sqrt{\frac{1}{n} \log \frac{1}{1-\varepsilon}} + \log \frac{1}{1-\varepsilon} \right) \right). \end{aligned}$$

Proof. The following computations are done under the event \mathcal{S} . First, apply Proposition B.3.6 to write

$$\begin{aligned} \sqrt{n} |\hat{T}_2| &\leq \sum_{j=1}^m \left(\hat{B}_1 |\mathbb{G}_n^{(0)}(h_{1,k} \mathbf{1}_{j\ell})| + \hat{B}_2 \sum_{\ell=2}^k |\mathbb{G}_n^{(0)}(h_{\ell,k} \mathbf{1}_{j\ell})| \right) \\ &+ m \hat{B}_2 \|h\|_{\infty} k(k-1) [\hat{B}_1 + \hat{B}_2(k+1)/3]. \end{aligned} \quad (\text{B.64})$$

We decompose on the event $\mathcal{E}_1^{\delta} \cap \mathcal{E}_2^{\delta}$.

Bound $|T_2|$ under the event $\mathcal{S} \setminus (\mathcal{E}_1^\delta \cap \mathcal{E}_2^\delta)$. In this case, we apply (B.52) from Proposition B.3.5 to get $\hat{B}_1 \leq n$ and $\hat{B}_2 \leq 1/\hat{p}_{\star,\varepsilon}^2$, along with the universal bounds from Lemma B.3.2:

$$\begin{aligned} \frac{1}{\sqrt{n}} |\mathbb{G}_n^{(0)}(h_{1,k} \mathbf{1}_{j\ell})| &\leq 2 \|h_{1,k}\|_\infty \leq 4k \|h\|_\infty \\ \frac{1}{\sqrt{n}} \sum_{\ell=2}^k |\mathbb{G}_n^{(0)}(h_{\ell,k} \mathbf{1}_{j\ell})| &\leq 2 \sum_{\ell=2}^k \|h_{\ell,k}\|_\infty \leq \sum_{\ell=2}^k 4(k-\ell+1) \|h\|_\infty = 2k(k-1) \|h\|_\infty \end{aligned}$$

so that by plugging into (B.64),

$$|\hat{T}_2| \leq \underbrace{\|h\|_\infty mk \left[4n + \frac{k-1}{\hat{p}_{\star,\varepsilon}^2} \left(n + 2 + \frac{k+1}{3\hat{p}_{\star,\varepsilon}^2} \right) \right]}_{\mathcal{T}_2^c} + \underbrace{0}_{E_2^c}.$$

Bound $|T_2|$ under the event $\mathcal{S} \cap \mathcal{E}_1^\delta \cap \mathcal{E}_2^\delta$. In this case, we may use that $n \geq 8/\hat{p}_{\star,\varepsilon}^2$ (because $[\log_2(2/\delta) + m \log(n+1)] \geq 1$ for $\delta \in (0, 1)$) and apply (B.53) from Proposition B.3.5 to get

$$\max \left\{ \hat{B}_1, \hat{B}_2 \right\} \leq O \left(\frac{1}{\hat{p}_{\star,\varepsilon}} \sqrt{\log \frac{1}{1-\varepsilon}} \right) + \frac{2}{\hat{p}_{\star,\varepsilon}} \sqrt{\frac{2 \log_2(2/\delta) + 2m \log(n+1)}{2n}}$$

The bounds based on \mathcal{E}_2^δ give

$$\begin{aligned} |\mathbb{G}_n^{(0)}(h_{1,k} \mathbf{1}_{j\ell})| &\leq \sqrt{2 \log \frac{2mk}{\delta}} 2k \|h\|_\infty \\ \sum_{\ell=2}^k |\mathbb{G}_n^{(0)}(h_{\ell,k} \mathbf{1}_{j\ell})| &\leq \sum_{\ell=2}^k \sqrt{2 \log \frac{2mk}{\delta}} 2(k-\ell+1) \|h\|_\infty \leq \sqrt{2 \log \frac{2mk}{\delta}} k(k-1) \|h\|_\infty. \end{aligned}$$

By plugging into (B.64), we can reuse the steps in the bound from (B.39) (for all terms without ε) to write

$$\begin{aligned} |\hat{T}_2| &\leq \frac{4mk \|h\|_\infty [\log_2(2/\delta) + 2m \log(n+1)]^{1-\mathbb{1}\{k=1\}/2}}{n\hat{p}_{\star,\varepsilon}^2} \times \\ &\quad \left[\hat{p}_{\star,\varepsilon} \sqrt{2 \log(2mk/\delta)} (k+1) + (k-1)(k+4) \right] + E_2, \end{aligned}$$

so that

$$\begin{aligned} \mathcal{T}_2 &= \frac{4mk \|h\|_\infty [\log_2(2/\delta) + 2m \log(n+1)]^{1-1\{k=1\}/2}}{n\hat{p}_{\star,\varepsilon}^2} \\ &\quad \times \left[\hat{p}_{\star,\varepsilon} \sqrt{2 \log(2mk/\delta)} (k+1) + (k-1)(k+4) \right]. \end{aligned}$$

We compute E_2 by using that

$$\begin{aligned} \max \left\{ \hat{B}_1, \hat{B}_2 \right\} &\leq O \left(\frac{1}{\hat{p}_{\star,\varepsilon}} \sqrt{\log \frac{1}{1-\varepsilon}} \right) + \tilde{O} \left(\frac{1}{\hat{p}_{\star,\varepsilon} \sqrt{n}} \right) \\ |\mathbb{G}_n^{(0)}(h_{1,k} \mathbf{1}_{j\ell})| &\leq \tilde{O}(k) \\ \sum_{\ell=2}^k |\mathbb{G}_n^{(0)}(h_{\ell,k} \mathbf{1}_{j\ell})| &\leq \tilde{O}(k^2), \end{aligned}$$

which gives

$$E_2 = \tilde{O} \left(\frac{k^3}{\hat{p}_{\star,\varepsilon}^2} \left(\sqrt{\frac{1}{n} \log \frac{1}{1-\varepsilon}} + \log \frac{1}{1-\varepsilon} \right) \right).$$

□

We now make the corresponding argument for the term T_1 .

Lemma B.3.8. *For $\delta = (\hat{p}_{\star,\varepsilon}/n)^4$, it holds that*

$$\mathcal{T}_1^c = \tilde{O}(k), \quad \mathcal{T}_1 = \tilde{O} \left(\frac{k}{\sqrt{n}} \right).$$

Proof. The following computations are done under the event \mathcal{S} .

Bound $|T_1|$ under the event $\mathcal{S} \setminus \mathcal{E}_3^\delta$. Here we simply apply a universal bound on the empirical process term:

$$\frac{1}{\sqrt{n}} |\mathbb{G}_n^{(0)}(h_{1,k})| \leq 2 \|h_{1,k}\|_\infty \leq 4k \|h\|_\infty,$$

so that $\mathcal{T}_1^c = 4k \|h\|_\infty$

Bound $|T_1|$ under the event $\mathcal{S} \cap \mathcal{E}_3^\delta$. Now, we may use the definition of the event \mathcal{E}_3^δ to achieve

$$\frac{1}{\sqrt{n}} |\mathbb{G}_n^{(0)}(h_{1,k})| \leq \sqrt{\frac{2 \log(2/\delta)}{n}} 2k \|h\|_\infty = \mathcal{T}_1.$$

□

Knowing that $E_2^c = 0$, we simplify (B.62) and (B.62) to read

$$\begin{aligned} & \tilde{O}\left(E_2(E_2 + \mathcal{T}_2 + \mathcal{T}_1)\right) \\ & \tilde{O}\left(k\sqrt{\varepsilon} \left((\hat{p}_{*,\varepsilon}/n)^4 \mathcal{T}_1^c + \mathcal{T}_1 + (\hat{p}_{*,\varepsilon}/n)^4 \mathcal{T}_2^c + \mathcal{T}_2 + E_2\right) + k^2 \varepsilon\right). \end{aligned}$$

We now combine the bounds from the previous two lemmas to compute (B.62) and (B.63) to achieve the main result, Theorem 4.6.1.

B.4 Zero-Shot Prediction

This appendix section contains basic definitions and identities used in Section 4.4. Let $(\Omega, \mathcal{F}, \mathbb{P})$ denote a probability space. The following definition allows for the discussion of conditional dependence measures based on regular conditional distributions.

Definition B.4.1. Consider random variables $(U, V) : \Omega \rightarrow \mathcal{U} \times \mathcal{V}$. Let $\mathcal{B}(\mathcal{U})$ denote the Borel σ -algebra on \mathcal{U} . A map: $\mu : \mathcal{V} \times \mathcal{B}(\mathcal{U}) \rightarrow [0, 1]$ is called a *regular conditional distribution* (*r.c.d.*) if the following two properties hold:

1. For each $A \in \mathcal{B}(\mathcal{U})$ and $\mathbf{v} \in \mathcal{V}$, it holds that

$$\mu(\mathbf{v}, A) = \mathbb{E}_{P_{U,V}} [\mathbb{1}_A(U) | V](\mathbf{v}).$$

2. For P_V -almost every $\mathbf{v} \in \mathcal{V}$, $\mu(\mathbf{v}, \cdot)$ is a probability measure on $\mathcal{B}(\mathcal{U})$.

For a joint probability measure $P_{X,Z}$, the Radon-Nikodym derivative $R = \frac{dP_{X,Z}}{d(P_X \otimes P_Z)}$ is useful for converting conditional expectation computations into marginal expectation com-

putations. The following identity is referenced by [Buja \[1990, Section 3\]](#) and [Dytso et al. \[2023, Lemma 1, Eq. \(14\)\]](#). We provide a self-contained proof below.

Lemma B.4.1. *Assume that $P_{X,Z} \ll P_X \otimes P_Z$, in which case there exists a Radon-Nikodym derivative $R = \frac{dP_{X,Z}}{d(P_X \otimes P_Z)}$. Then, for all $g \in \mathbf{L}^2(P_Z)$ and $h \in \mathbf{L}^2(P_X)$, it holds that*

$$\begin{aligned}\mathbb{E}_{P_{X,Z}}[g(Z)|X](\mathbf{x}) &= \mathbb{E}_{P_Z}[g(Z)R(\mathbf{x}, Z)] \text{ for } P_X\text{-almost all } \mathbf{x} \in \mathcal{X}, \\ \mathbb{E}_{P_{X,Z}}[h(X)|Z](\mathbf{z}) &= \mathbb{E}_{P_X}[h(X)R(X, \mathbf{z})] \text{ for } P_Z\text{-almost all } \mathbf{z} \in \mathcal{Z}.\end{aligned}$$

Proof. We prove the first identity, whereas the second follows by a symmetric argument. To confirm that the two functions are equal almost surely, it is sufficient to prove that for any measurable set $A \in \sigma(X)$ (the σ -algebra generated by X) the relation

$$\int_A \mathbb{E}_{P_{X,Z}}[g(Z)|X](\mathbf{x}) dP_X(\mathbf{x}) = \int_A \mathbb{E}_{P_Z}[g(Z)R(\mathbf{x}, Z)] dP_X(\mathbf{x}). \quad (\text{B.65})$$

By the definition of conditional expectation, we have that

$$\begin{aligned}\int_A \mathbb{E}_{P_{X,Z}}[g(Z)|X](\mathbf{x}) dP_X(\mathbf{x}) &= \int_{\mathcal{X}} \mathbb{E}_{P_{X,Z}}[g(Z)|X](\mathbf{x}) \mathbb{1}_A(\mathbf{x}) dP_X(\mathbf{x}) \\ &= \mathbb{E}_{P_{X,Z}}[g(Z) \mathbb{1}_A(X)] \\ &= \mathbb{E}_{P_X \otimes P_Z}[g(Z) \mathbb{1}_A(X) R(X, Z)],\end{aligned}$$

where the last step follows from the Radon-Nikodym theorem [[Schilling, 2017, Theorem 20.2](#)]. Next, we compute the expectation, taken under the product measure, using Fubini's theorem [[Schilling, 2017, Corollary 14.9](#)]. That is,

$$\begin{aligned}\int_A \mathbb{E}_{P_{X,Z}}[g(Z)|X](\mathbf{x}) dP_X(\mathbf{x}) &= \mathbb{E}_{P_X \otimes P_Z}[g(Z) \mathbb{1}_A(X) R(X, Z)] \\ &= \int_A \left(\int_{\mathcal{Z}} g(\mathbf{z}) R(\mathbf{x}, \mathbf{z}) dP_Z(\mathbf{z}) \right) dP_X(\mathbf{x}) \\ &= \int_A \mathbb{E}_{P_Z}[g(Z) R(\mathbf{x}, Z)] dP_X(\mathbf{x}).\end{aligned}$$

This achieves (B.65) and completes the proof. \square

B.5 Experimental Details

We provide the full details of the experimental results from Section 4.5.

B.5.1 Datasets

Pre-Training Data For the balancing-based experiments, the pre-training data was taken from the public [ImageNet-Captions](#) dataset [Fang et al., 2023]. We subset the dataset by selecting the 250 classes that were most frequent in the dataset, resulting in 174,594 images and associated Flickr captions. The exact images used and their associated captions are given in the GitHub repo <https://github.com/ronakdm/balancing>.

Evaluation Data We perform zero-shot classification with various image classification and image-caption datasets. For the balanced pre-training experiments, we used the default prompt templates for classification from the [CLIP Benchmark](#) repo. Other, customized prompting strategies used in our experiments are described at the end of this appendix. The datasets (test splits) used were:

- **CIFAR-10:** 10,000 colored natural images labeled with one of 10 classes.
- **CIFAR-100:** 10,000 colored natural images labeled with one of 100 classes.
- **STL-10:** 80,000 colored natural images labeled with one of 10 classes.
- **MS-COCO:** 41,000 colored natural images with associated captions.
- **Flickr8k:** 8,000 colored natural images with associated captions.
- **Rendered SST2:** 1,821 images of typed natural language with sentiment label (2 classes).
- **VOC2007:** 4,952 colored natural images labeled with one of 20 classes.

- **FGVC Aircraft:** 34,000 colored natural images labeled with one of 102 classes.

Evaluation scripts using the various embedding models (described below) are provided.

B.5.2 Model Specification and Hyperparameters

CLIP Architectures First, we specify which OpenCLIP models and pre-training sets were used. These models were chosen due to their range of top-1 zero-shot accuracies on the ImageNet-1k benchmark (as shown below). As opposed to already highly performant models ($\geq 50\%$ on ImageNet-1k), these models benefitted more from optimized prompting techniques in our initial experiments.

Model	OpenCLIP Model Tag	Pre-Training Set Tag	ImageNet-1k Top-1 Acc.
ResNet-50	RN50	yfcc15m	28.11%
NLLB-CLIP	nllb-clip-base	v1	33.51%
ViT-B/32	ViT-B-32	datacomp_m_s128m_b4k	32.81%

Optimizer For optimization, models were trained with stochastic gradient descent (SGD) with the learning rate tuned along the grid $\{1^{-3}, 3^{-3}, 1^{-2}, 3^{-2}, 1^{-1}\}$ and a fixed weight decay parameter of 0.01. Momentum-variants such as Adam [Kingma and Ba, 2015] were not used to isolate the effect of variance reduction as described in the balancing example of Section 4.5.

B.5.3 Compute Environment

Experiments were run on a CPU/GPU workstation with 12 virtual cores, 126G of memory, and four NVIDIA TITAN Xp GPUs with 12G memory each. The code was written in Python 3 and we use PyTorch for automatic differentiation. The [OpenCLIP](#) and [CLIP Benchmark](#) repositories were used for zero-shot evaluation.

Prompt-Generating Model We employed the meta-llama/Llama-3.2-1B-Instruct model publicly available on [HuggingFace](#). For the purpose of generation, we used a **top- p** hyperparameter of **0.9** and **temperature** hyperparameter of **0.99** for more diverse responses. Meta-prompting was based on the following instructions per dataset, which are slight variations of those used in [Pratt et al. \[2023\]](#):

- **Describable Textures Dataset (DTD):**

- “What does __ material look like?”,
- “What does a __ surface look like?”,
- “What does a __ texture look like?”,
- “What does a __ object look like?”,
- “What does a __ pattern look like?”

- **Flowers 102:**

- “Describe how to identify a(n) __, a type of flower.”,
- “What does a(n) __ flower looks like?”

- **FGVC Aircraft:**

- “Describe a(n) __ aircraft.”,
- “Describe the __ aircraft.”

- **SUN397:**

- “Describe what a(n) __ looks like.”,
- “How can you identify a(n) __?”,
- “Describe a photo of a(n) __.”,

- “Describe the scene of $a(n)$ ____.”

- **ImageNet-1k:**

- “Describe what $a(n)$ ____ looks like.”,
- “How can you identify $a(n)$ ____?”,
- “What does $a(n)$ ____ look like?”,
- “Describe an image from the Internet of $a(n)$ ____.”,
- “Write a caption of an image of $a(n)$ ____.”

The following additional instruction was appended for better-formatted responses: *“Please format your response as one that contains only lower case letters and no special characters (including new lines, bold, and any markdown artifacts) other than a period (‘.’) or commas (‘,’). The response should be a single sentence ending in a period that is directed toward the final instruction in this message. Your sentence should be a minimum of three words and a maximum of thirty.”*

Our reproducibility effort includes not only the full list of all 164,400 prompts generated from LLaMA 3, but the subset of prompts used for each class and each seed used to generate the figures in Section 4.5.