

Chapter 5

Advances: Finite-sum Coupled Compositional Optimization

Abstract In this chapter, we study a novel family of stochastic compositional optimization problems namely **finite-sum coupled compositional optimization (FCCO)**, and introduce algorithms for solving them. These algorithms have direct applications in addressing the empirical X-risk minimization challenges discussed in Chapter 2. To ensure broad applicability, we examine various settings of this problem, characterized by different properties of outer and inner functions, including smooth and non-smooth cases, as well as convex, weakly convex, and non-convex scenarios. The results presented here also significantly extend and complement those discussed in Chapter 4. We also discuss how to efficiently optimize compositional optimized certainty equivalent risks, especially compositional entropic risk.

Coupling reveals depth where composition meets reality!

Contents

5.1	Finite-sum Coupled Compositional Optimization	189
5.2	Smooth Functions	190
5.2.1	The SOX Algorithm	191
5.2.2	Multi-block Single-Probe Variance Reduction	199
5.3	Non-Smooth Weakly Convex Functions	208
5.3.1	SONX for Non-smooth Inner Functions	210
5.3.2	SONEX for Non-smooth Outer functions	217
5.4	Convex inner and outer functions	222
5.4.1	The ALEXR Algorithm	224
5.4.2	Technical Lemmas	226
5.4.3	Strongly convex objectives	237
5.4.4	Convex objectives with non-smooth outer functions	242
5.4.5	Double-loop ALEXR for weakly convex inner functions	247
5.4.6	Lower Bounds	249
5.5	Stochastic Optimization of Compositional OCE	255
5.5.1	A Basic Algorithm	256
5.5.2	A Geometry-aware Algorithm for Entropic Risk	264
5.6	History and Notes	294

5.1 Finite-sum Coupled Compositional Optimization

Specifically, we focus on the following optimization problem:

$$\min_{\mathbf{w} \in \mathbb{R}^d} F(\mathbf{w}) := \frac{1}{n} \sum_{i=1}^n f_i(\mathbb{E}_{\zeta \sim \mathbb{P}_i} g_i(\mathbf{w}; \zeta)), \quad (5.1)$$

where $g_i(\cdot; \zeta) : \mathbb{R}^d \rightarrow \mathbb{R}^{d'}$ is a stochastic mapping, $f_i(\cdot) : \mathbb{R}^{d'} \rightarrow \mathbb{R}$ is a deterministic function, and \mathbb{P}_i denotes the distribution of the random variable ζ .

We refer to this problem as **finite-sum coupled compositional optimization (FCCO)**. If we interpret i as an outer random variable, a distinctive feature that sets FCCO apart from standard stochastic compositional optimization (SCO) is that each inner stochastic function $g_i(\mathbf{w}; \zeta)$ depends on both an inner random variable ζ and an outer index i , giving rise to the term *coupled*. While this problem can be cast as a special case of SCO by defining $f(\mathbf{g}) = \frac{1}{n} \sum_{i=1}^n f_i(g_i)$ and $\mathbf{g}(\mathbf{w}) = [g_1(\mathbf{w}), \dots, g_n(\mathbf{w})]$, the high dimensionality of \mathbf{g} due to large n , along with its stochastic components, significantly complicates the construction of unbiased estimators and theoretical analysis. Therefore, FCCO warrants the development of specialized optimization methods.

Below, we revisit several applications of FCCO in ML and discuss the properties of f_i and g_i .

Group DRO

In Section 2.2.3, we have formulated the CVaR divergence regularized group DRO as

$$\min_{\mathbf{w}, \nu} \frac{1}{K\alpha} \sum_{i=1}^K [L_i(\mathbf{w}) - \nu]_+ + \nu, \quad (5.2)$$

where $\alpha \in (0, 1)$, $L_i(\mathbf{w}) = \frac{1}{n_k} \sum_{j=1}^{n_k} \ell(\mathbf{w}; \mathbf{x}_j^i, y_j^i)$ denotes the average loss over data from the i -th group. The first term above is an instance of the FCCO objective, where the outer function $f(g) = ([g]_1 - [g]_2)_+$ is a **convex but non-smooth** function of g , and each inner function $g_i(\mathbf{w}, \nu) = [L_i(\mathbf{w}), \nu]^\top$ could be convex or non-convex, smooth or non-smooth depending on applications.

AP Maximization

In Section 2.3.2, the AP maximization has been formulated as the following problem:

$$\min_{\mathbf{w}} \frac{1}{n_+} \sum_{\mathbf{x}_i \in \mathcal{S}_+} f(g_i(\mathbf{w})), \quad (5.3)$$

where \mathcal{S}_+ is the set of n_+ positive examples, $g_i(\mathbf{w}) = [g_1(\mathbf{w}; \mathbf{x}_i, \mathcal{S}), g_2(\mathbf{w}; \mathbf{x}_i, \mathcal{S})]^\top$ is a vector mapping with two components:

$$g_1(\mathbf{w}; \mathbf{x}_i, \mathcal{S}) = \frac{1}{|\mathcal{S}|} \sum_{\mathbf{x}_j \in \mathcal{S}} \mathbb{I}(y_j = 1) \ell(h(\mathbf{w}; \mathbf{x}_j) - h(\mathbf{w}; \mathbf{x}_i))$$

$$g_2(\mathbf{w}; \mathbf{x}_i, \mathcal{S}) = \frac{1}{|\mathcal{S}|} \sum_{\mathbf{x}_j \in \mathcal{S}} \ell(h(\mathbf{w}; \mathbf{x}_j) - h(\mathbf{w}; \mathbf{x}_i)),$$

and $f(\mathbf{g}) = -\frac{[\mathbf{g}]_1}{[\mathbf{g}]_2}$ is simple function. We can see that f is **non-convex and smooth** if the loss value is upper bounded and $\ell(0)$ is lower bounded. The inner mapping $g_i(\mathbf{w})$ could be convex (e.g., a linear model) or non-convex (e.g., a deep model), smooth or non-smooth depending on applications.

Contrastive Representation Learning

The contrastive objective of self-supervised representation learning presented in (2.50), is the following:

$$\min_{\mathbf{w}} \frac{1}{n} \sum_{i=1}^n \tau \log \left(\varepsilon + \frac{1}{|\mathcal{S}_i^-|} \sum_{\mathbf{y} \in \mathcal{S}_i^-} \exp((s(\mathbf{w}; \mathbf{x}_i, \mathbf{y}) - s(\mathbf{w}; \mathbf{x}_i, \mathbf{x}_i^+))/\tau) \right).$$

The outer function $f(g) = \tau \log(\varepsilon + g)$ is a non-convex function and smooth when ε is lower bounded. Each inner function g_i is a non-convex function of \mathbf{w} in general.

5.2 Smooth Functions

In this section, we consider a non-convex but smooth objective function $F(\mathbf{w})$ with smooth outer functions. In addition, we assume the inner stochastic functions satisfy the following conditions throughout this section.

Assumption 5.1. *We assume that*

- (i) $\mathbb{E}_{\zeta \sim \mathbb{P}_i} [\|g_i(\mathbf{w}; \zeta) - g_i(\mathbf{w})\|_2^2] \leq \sigma_0^2.$
- (ii) $\mathbb{E}_{\zeta \sim \mathbb{P}_i} [\|\nabla g_i(\mathbf{w}; \zeta) - \nabla g_i(\mathbf{w})\|_2^2] \leq \sigma_2^2.$
- (iii) $\mathbb{E}_{\zeta \sim \mathbb{P}_i} [\|\nabla g_i(\mathbf{w}; \zeta)\|_2^2] \leq G_2^2.$

5.2.1 The SOX Algorithm

The first algorithm for solving FCCO is called SOX, named by **S**tochastic **O**ptimization of **X**-risks. Owing to its ease of implementation and favorable practical performance, this algorithm is commonly adopted for addressing FCCO. Below, we outline the assumptions necessary for its analysis.

Assumption 5.2. *There exist $G_1, L_1, L_F > 0$ such that*

- (i) $f_i : \mathbb{R}^{d'} \mapsto \mathbb{R}$ is G_1 -Lipschitz continuous and L_1 -smooth;
- (ii) $F : \mathbb{R}^d \mapsto \mathbb{R}$ is L_F -smooth;
- (iii) $F_* = \min_{\mathbf{w}} F(\mathbf{w}) \geq -\infty$.

Similar to that for SCO, we also need to track and estimate the inner functions. However, the difference is that we need to maintain and update n estimators for the n inner functions $g_i(\mathbf{w}), i \in [n]$.

To this end, we maintain n sequence of estimators $\{\mathbf{u}_{i,t}, t \in [T]\}_{i=1}^n$. At the t -th iteration, we draw a set of B random indices $\mathcal{B}_t \subset [n]$ with $|\mathcal{B}_t| = B$. We update $\mathbf{u}_{i,t}, i \in [n]$ by the following:

$$\mathbf{u}_{i,t} = \begin{cases} (1 - \gamma_t)\mathbf{u}_{i,t-1} + \gamma_t g_i(\mathbf{w}_t; \zeta_{i,t}), & i \in \mathcal{B}_t \\ \mathbf{u}_{i,t-1}, & \text{o.w.} \end{cases}, t = 1, \dots, T, \quad (5.4)$$

where $\zeta_{i,t} \sim \mathbb{P}_i$ is a random variable. We refer to the above estimator as coordinate moving average estimator. Then, similar to SCMA, a moving average estimator of the gradient is computed by:

$$\mathbf{v}_t = (1 - \beta_t)\mathbf{v}_{t-1} + \beta_t \mathbf{z}_t, \\ \text{where } \mathbf{z}_t = \frac{1}{|\mathcal{B}_t|} \sum_{i \in \mathcal{B}_t} \nabla g_i(\mathbf{w}_t; \zeta'_{i,t}) \nabla f_i(\mathbf{u}_{i,t}).$$

Then, the model parameters are updated by:

$$\mathbf{w}_{t+1} = \mathbf{w}_t - \eta_t \mathbf{v}_t.$$

The detailed steps are presented in Algorithm 14.

Convergence Analysis

Let us first define two notations:

$$\Delta_t = \|\mathbf{v}_t - \nabla F(\mathbf{w}_t)\|_2^2, \quad (5.5)$$

$$\delta_t = \frac{1}{n} \sum_{i=1}^n \|\mathbf{u}_{i,t} - g_i(\mathbf{w}_t)\|_2^2. \quad (5.6)$$

Algorithm 14 SOX

```

1: Input: learning rate schedules  $\{\eta_t\}_{t=1}^T, \{\gamma_t\}_{t=1}^T, \{\beta_t\}_{t=1}^T$ ; starting points  $\mathbf{w}_0, \mathbf{u}_0, \mathbf{v}_0$ 
2: Let  $\mathbf{w}_1 = \mathbf{w}_0 - \eta_0 \mathbf{v}_0$ 
3: for  $t = 1, \dots, T$  do
4:   Draw a batch of samples  $\mathcal{B}_t \subset [n]$ 
5:   for  $i \in \mathcal{B}_t$  do
6:     Draw two samples  $\zeta_{i,t}, \zeta'_{i,t} \sim \mathbb{P}_i$ 
7:     Update the inner function value estimators
           
$$\mathbf{u}_{i,t} = (1 - \gamma_t)\mathbf{u}_{i,t-1} + \gamma_t g_i(\mathbf{w}_t; \zeta_{i,t}),$$

8:   end for
9:   Set  $\mathbf{u}_{i,t} = \mathbf{u}_{i,t-1}, i \notin \mathcal{B}_t$ 
10:  Compute the vanilla gradient estimator  $\mathbf{z}_t = \frac{1}{|\mathcal{B}_t|} \sum_{i \in \mathcal{B}_t} \nabla g_i(\mathbf{w}_t; \zeta'_{i,t}) \nabla f_i(\mathbf{u}_{i,t})$ 
11:  Update the MA gradient estimator  $\mathbf{v}_t = (1 - \beta_t)\mathbf{v}_{t-1} + \beta_t \mathbf{z}_t$ 
12:  Update the model  $\mathbf{w}$  by  $\mathbf{w}_{t+1} = \mathbf{w}_t - \eta_t \mathbf{v}_t$ 
13: end for

```

The descent lemma (Lemma 4.9) remains valid. Next, we analyze the recursion of Δ_t and δ_t . One point of deviation is that only some randomly selected coordinates of \mathbf{u} are updated and used for computing the gradient estimator \mathbf{z}_t . To facilitate the proof, we introduce a virtual sequence:

$$\bar{\mathbf{u}}_{i,t} = (1 - \gamma_t)\mathbf{u}_{i,t-1} + \gamma_t g_i(\mathbf{w}_t; \zeta_{i,t}), \forall i = 1, \dots, n. \quad (5.7)$$

This is similar to that is done in the analysis of stochastic coordinate descent method in Section 3.3. Then, we have

$$\mathcal{M}_t = \mathbb{E}_{\mathcal{B}_t, \zeta'_t}[\mathbf{z}_t] = \frac{1}{n} \sum_{i=1}^n \nabla g_i(\mathbf{w}_t) \nabla f_i(\bar{\mathbf{u}}_{i,t}).$$

Critical: Since \mathbf{u}_t is a random variable that depends on \mathcal{B}_t , hence

$$\mathbb{E}_{\mathcal{B}_t, \zeta'_t}[\mathbf{z}_t] \neq \frac{1}{n} \sum_{i=1}^n \nabla g_i(\mathbf{w}_t) \nabla f_i(\mathbf{u}_{i,t}).$$

We first bound the error recursion of δ_t .

Lemma 5.1 *Consider the \mathbf{u}_t updates in Algorithm 14. Under Assumption 5.1, if $\gamma_t \leq 1$, then*

$$\mathbb{E}[\delta_t] \leq \left(1 - \frac{B\gamma_t}{2n}\right) \mathbb{E}[\delta_{t-1}] + \frac{2nG_2^2}{B\gamma_t} \mathbb{E}[\|\mathbf{w}_{t-1} - \mathbf{w}_t\|_2^2] + \frac{B\gamma_t^2 \sigma_0^2}{n}.$$

Proof. Since $\bar{\mathbf{u}}_{i,t}$ is updated using MA, then similar to (4.6), for all $i \in [n]$ we have

$$\mathbb{E}_{\zeta_{i,t}} [\|\bar{\mathbf{u}}_{i,t} - g_i(\mathbf{w}_t)\|_2^2] \leq (1 - \gamma_t)^2 \|\mathbf{u}_{i,t-1} - g_i(\mathbf{w}_t)\|_2^2 + \gamma_t^2 \sigma_0^2.$$

Given $i \in [n]$, with a probability of B/n that $i \in \mathcal{B}_t$, we have $\mathbf{u}_{i,t} = \bar{\mathbf{u}}_{i,t}$; otherwise, $\mathbf{u}_{i,t} = \mathbf{u}_{i,t-1}$. Hence,

$$\begin{aligned} & \mathbb{E}_{\zeta_{i,t}} \mathbb{E}_{\mathcal{B}_t} [\|\mathbf{u}_{i,t} - g_i(\mathbf{w}_t)\|_2^2] \\ &= \frac{B}{n} \mathbb{E}_{\zeta_{i,t}} [\|\bar{\mathbf{u}}_{i,t} - g_i(\mathbf{w}_t)\|_2^2] + \left(1 - \frac{B}{n}\right) \|\mathbf{u}_{i,t-1} - g_i(\mathbf{w}_t)\|_2^2 \\ &\leq \frac{B}{n} (1 - \gamma_t)^2 \|\mathbf{u}_{i,t-1} - g_i(\mathbf{w}_t)\|_2^2 + \frac{B\gamma_t^2 \sigma_0^2}{n} + \left(1 - \frac{B}{n}\right) \|\mathbf{u}_{i,t-1} - g_i(\mathbf{w}_t)\|_2^2 \\ &\leq \left(1 - \frac{B\gamma_t}{2n}\right)^2 \|\mathbf{u}_{i,t-1} - g_i(\mathbf{w}_t)\|_2^2 + \frac{B\gamma_t^2 \sigma_0^2}{n}, \end{aligned}$$

where we use the fact $\frac{B}{n}(1 - \gamma_t)^2 + \left(1 - \frac{B}{n}\right) \leq \left(1 - \frac{\gamma_t B}{2n}\right)^2$. Then, taking expectation over all randomness on both sides yields

$$\mathbb{E} [\|\mathbf{u}_{i,t} - g_i(\mathbf{w}_t)\|_2^2] \leq \left(1 - \frac{B\gamma_t}{2n}\right)^2 \mathbb{E} [\|\mathbf{u}_{i,t-1} - g_i(\mathbf{w}_t)\|_2^2] + \frac{B\gamma_t^2 \sigma_0^2}{n}.$$

Then using the Young's inequality similar to the proof of Lemma 4.1, we have

$$\begin{aligned} & \mathbb{E} [\|\mathbf{u}_{i,t} - g_i(\mathbf{w}_t)\|_2^2] \leq \left(1 + \frac{B\gamma_t}{2n}\right) \left(1 - \frac{B\gamma_t}{2n}\right)^2 \mathbb{E} [\|\mathbf{u}_{i,t-1} - g_i(\mathbf{w}_{t-1})\|_2^2] \\ & \quad + \left(1 + \frac{2n}{B\gamma_t}\right) \left(1 - \frac{B\gamma_t}{2n}\right)^2 \mathbb{E} [\|g_i(\mathbf{w}_{t-1}) - g_i(\mathbf{w}_t)\|_2^2] + \frac{B\gamma_t^2 \sigma_0^2}{n} \\ & \leq \left(1 - \frac{B\gamma_t}{2n}\right) \mathbb{E} [\|\mathbf{u}_{i,t-1} - g_i(\mathbf{w}_{t-1})\|_2^2] + \frac{2nG_2^2}{B\gamma_t} \mathbb{E} [\|\mathbf{w}_{t-1} - \mathbf{w}_t\|_2^2] + \frac{B\gamma_t^2 \sigma_0^2}{n}, \end{aligned}$$

where we use $\gamma_t \leq 1 < \frac{2n}{B}$. The desired result follows by taking average over $i = 1, \dots, n$ on both sides. \square

Lemma 5.2 (Variance of \mathbf{z}_t) Let $\sigma^2 = \frac{G_1^2 \sigma_2^2}{B} + \frac{G_1^2 G_2^2}{B} \frac{n-B}{n-1}$. We have

$$\mathbb{E}_t [\|\mathbf{z}_t - \mathcal{M}_t\|_2^2] \leq \sigma^2.$$

Proof. First, using the variance bound of the average of B independent zero-mean random variables gives

$$A_1 = \mathbb{E}_t \left[\left\| \frac{1}{B} \sum_{i \in \mathcal{B}_t} \nabla g_i(\mathbf{w}_t; \zeta'_{i,t}) \nabla f_i(\bar{\mathbf{u}}_{i,t}) - \frac{1}{B} \sum_{i \in \mathcal{B}_t} \nabla g_i(\mathbf{w}_t) \nabla f_i(\bar{\mathbf{u}}_{i,t}) \right\|_2^2 \right] \leq \frac{G_1^2 \sigma_2^2}{B},$$

and using the variance bound of B random variables without replacement yields

$$A_2 = \mathbb{E}_t \left[\left\| \frac{1}{B} \sum_{i \in \mathcal{B}_t} \nabla g_i(\mathbf{w}_t) \nabla f_i(\bar{\mathbf{u}}_{i,t}) - \frac{1}{n} \sum_{i=1}^n \nabla g_i(\mathbf{w}_t) \nabla f_i(\bar{\mathbf{u}}_{i,t}) \right\|_2^2 \right] \leq \frac{G_1^2 G_2^2}{B} \frac{n-B}{n-1}.$$

As a result,

$$\begin{aligned} & \mathbb{E}_t [\|\mathbf{z}_t - \mathcal{M}_t\|_2^2] \\ &= \mathbb{E}_t \left[\left\| \frac{1}{B} \sum_{i \in \mathcal{B}_t} \nabla g_i(\mathbf{w}_t; \zeta'_{i,t}) \nabla f_i(\bar{\mathbf{u}}_{i,t}) - \frac{1}{n} \sum_{i=1}^n \nabla g_i(\mathbf{w}_t) \nabla f_i(\bar{\mathbf{u}}_{i,t}) \right\|_2^2 \right] \\ &= A_1 + A_2 \leq \frac{G_1^2 \sigma_2^2}{B} + \frac{G_1^2 G_2^2}{B} \frac{n-B}{n-1} := \sigma^2. \end{aligned}$$

□

Lemma 5.3 *Under Assumptions 5.1 and 5.2, if $\beta_t \leq 1$, the gradient estimation error Δ_t can be bounded as*

$$\begin{aligned} \mathbb{E}[\Delta_t] &\leq (1 - \beta_t) \mathbb{E}[\Delta_{t-1}] + \frac{2L_F^2 + 8\beta_t^2 G_2^4 L_1^2}{\beta_t} \mathbb{E}[\|\mathbf{w}_t - \mathbf{w}_{t-1}\|_2^2] + 8\beta_t L_1^2 G_2^2 \mathbb{E}[\delta_{t-1}] \\ &\quad + \beta_t^2 \sigma^2 + 4G_2^2 L_1^2 \beta_t \gamma_t^2 \sigma_0^2, \end{aligned}$$

$$\text{where } \sigma^2 = \frac{G_1^2 \sigma_2^2}{B} + \frac{G_1^2 G_2^2}{B} \frac{n-B}{n-1}.$$

Proof. Since \mathbf{v}_t is updated using MA, we apply Lemma 4.7 in light of Lemma 5.2, yielding

$$\begin{aligned} \mathbb{E}[\|\mathbf{v}_t - \nabla F(\mathbf{w}_t)\|_2^2] &\leq (1 - \beta_t) \mathbb{E}[\|\mathbf{v}_{t-1} - \nabla F(\mathbf{w}_{t-1})\|_2^2] \\ &\quad + \frac{2L_F^2}{\beta_t} \mathbb{E}[\|\mathbf{w}_{t-1} - \mathbf{w}_t\|_2^2] + 4\beta_t \mathbb{E}[\|\mathcal{M}_t - \nabla F(\mathbf{w}_t)\|_2^2] + \beta_t^2 \sigma^2. \end{aligned} \tag{5.8}$$

Next, we bound $\mathbb{E}[\|\mathcal{M}_t - \nabla F(\mathbf{w}_t)\|_2^2]$.

$$\begin{aligned} \|\mathcal{M}_t - \nabla F(\mathbf{w}_t)\|_2^2 &= \left\| \frac{1}{n} \sum_{i=1}^n \nabla g_i(\mathbf{w}_t) \nabla f(\bar{\mathbf{u}}_{i,t}) - \frac{1}{n} \sum_{i=1}^n \nabla g_i(\mathbf{w}_t) \nabla f(g_i(\mathbf{w}_t)) \right\|_2^2 \\ &\leq G_2^2 L_1^2 \frac{1}{n} \sum_{i=1}^n \|\bar{\mathbf{u}}_{i,t} - g_i(\mathbf{w}_t)\|_2^2. \end{aligned}$$

From Lemma 5.1, we have

$$\mathbb{E}_{\zeta_{i,t}} [\|\bar{\mathbf{u}}_{i,t} - g_i(\mathbf{w}_t)\|_2^2] \leq (1 - \gamma_t)^2 \|\mathbf{u}_{i,t-1} - g_i(\mathbf{w}_t)\|_2^2 + \gamma_t^2 \sigma_0^2, \forall i.$$

Hence

$$\begin{aligned}
 \mathbb{E} \left[\frac{1}{n} \sum_{i=1}^n \|\bar{\mathbf{u}}_{i,t} - g_i(\mathbf{w}_t)\|_2^2 \right] &\leq (1 - \gamma_t)^2 \mathbb{E} \left[\frac{1}{n} \sum_{i=1}^n \|\mathbf{u}_{i,t-1} - g_i(\mathbf{w}_t)\|_2^2 \right] + \gamma_t^2 \sigma_0^2 \\
 &\leq (1 - \gamma_t)^2 \mathbb{E} \left[\frac{1}{n} \sum_{i=1}^n (2\|g_i(\mathbf{w}_t) - g_i(\mathbf{w}_{t-1})\|_2^2 + 2\|\mathbf{u}_{i,t-1} - g_i(\mathbf{w}_{t-1})\|_2^2) \right] + \gamma_t^2 \sigma_0^2 \\
 &\leq 2\mathbb{E}[\delta_{t-1}] + \gamma_t^2 \sigma_0^2 + \mathbb{E} \left[2G_2^2 \|\mathbf{w}_{t-1} - \mathbf{w}_t\|_2^2 \right].
 \end{aligned}$$

As a result,

$$\mathbb{E}[\|\mathcal{M}_t - \nabla F(\mathbf{w}_t)\|_2^2] \leq 2G_2^2 L_1^2 \mathbb{E}[\delta_{t-1}] + G_2^2 L_1^2 \gamma_t^2 \sigma_0^2 + \mathbb{E} \left[2G_2^4 L_1^2 \|\mathbf{w}_{t-1} - \mathbf{w}_t\|_2^2 \right].$$

Plugging the above results into (5.8) we finish the proof. \square

For combining the descent lemma and the above lemmas, we present a result similar to Lemma 4.10, with differences highlighted in boxes.

Lemma 5.4 *If $\eta_t \leq 1/L$, assume that there exist non-negative sequences $A_t, B_t, \Gamma_t, \Delta_t, \delta_t, t \geq 0$ satisfying:*

$$\begin{aligned}
 (*) A_{t+1} &\leq A_t + \eta_t \Delta_t - \eta_t B_t - \eta_t \Gamma_t \\
 (\#) \Delta_{t+1} &\leq (1 - \beta_{t+1}) \Delta_t + C_1 \beta_{t+1} \boxed{\delta_t} + \frac{C_2 \eta_t^2}{\beta_{t+1}} \Gamma_t + \beta_{t+1}^2 \sigma^2 + \boxed{\beta_{t+1} \gamma_{t+1}^2 \sigma'^2}, \\
 (\diamond) \delta_{t+1} &\leq (1 - \gamma_{t+1}) \delta_t + \frac{C_3 \eta_t^2}{\gamma_{t+1}} \Gamma_t + \gamma_{t+1}^2 \sigma'^2.
 \end{aligned}$$

If $\beta = \frac{\epsilon^2}{4\sigma^2}, \gamma = \min(\frac{\epsilon^2}{8C_1\sigma'^2}, \frac{\epsilon}{2\sigma''}), \eta = \min(\frac{1}{L}, \frac{\beta}{\sqrt{4C_2}}, \frac{\gamma}{\sqrt{8C_1C_3}})$, then in order to guarantee

$$\sum_{t=0}^{T-1} \frac{1}{T} (B_t + \frac{1}{2} \Gamma_t) \leq \epsilon^2.$$

the iteration complexity is in the order of

$$T = O \left(\max \left\{ \frac{C_Y L_F}{\epsilon^2}, \frac{C_Y \sqrt{C_1 C_3} \sigma''}{\epsilon^3}, \frac{C_Y \sigma^2 \sqrt{C_2}}{\epsilon^4}, \frac{C_Y \sqrt{C_1 C_3} C_1 \sigma'^2}{\epsilon^4} \right\} \right).$$

where $C_Y \leq A_0 - \min_t A_t + \frac{1}{2\sqrt{C_2}} \Delta_0 + \sqrt{\frac{C_1}{2C_3}} \delta_0$.

Proof. Following similar analysis to Lemma 4.10, we have

$$\begin{aligned}
A_{t+1} + \frac{\eta_t}{\beta_{t+1}} \Delta_{t+1} + (C_1 \eta_t \frac{1 + \gamma_{t+1}}{\gamma_{t+1}} - C_1 \eta_t) \delta_{t+1} &\leq A_t - \eta_t B_t - \eta_t \Gamma_t \\
+ \left(\eta_t + \frac{\eta_t}{\beta_{t+1}} (1 - \beta_{t+1}) \right) \Delta_t + \frac{C_2 \eta_t^3}{\beta_{t+1}^2} \Gamma_t + \eta_t (\beta_{t+1} \sigma^2 + \boxed{\gamma_{t+1}^2 \sigma'^2}) &+ \boxed{C_1 \eta_t (\delta_t - \delta_{t+1})} \\
+ C_1 \eta_t \frac{1 + \gamma_{t+1}}{\gamma_{t+1}} (1 - \gamma_{t+1}) \delta_t + \frac{C_3 C_1 \eta_t^3 (1 + \gamma_{t+1})}{\gamma_{t+1}^2} \Gamma_t &+ C_1 \eta_t (1 + \gamma_{t+1}) \gamma_{t+1} \sigma'^2.
\end{aligned}$$

where the terms in the box highlight the difference due to the slight difference in the recursion of Δ_t . Under similar conditions of $\beta_{t+1}, \gamma_{t+1}, \eta_t$ and similar analysis, we get

$$\begin{aligned}
A_{t+1} + \frac{\eta_t}{\beta_{t+1}} \Delta_{t+1} + \frac{C_1 \eta_t}{\gamma_{t+1}} \delta_{t+1} &\leq A_t + \frac{\eta_{t-1}}{\beta_t} \Delta_t + \frac{C_1 \eta_{t-1}}{\gamma_t} \delta_t + \boxed{C_1 \eta_t (\delta_t - \delta_{t+1})} \\
- \eta_t B_t - \frac{1}{2} \eta_t \Gamma_t + \eta_t (\beta_{t+1} \sigma^2 + \boxed{\gamma_{t+1}^2 \sigma'^2}) &+ 2C_1 \eta_t \gamma_{t+1} \sigma'^2.
\end{aligned}$$

Since $\eta_{t+1} \leq \eta_t$, we have

$$\begin{aligned}
A_{t+1} + \frac{\eta_t}{\beta_{t+1}} \Delta_{t+1} + \frac{C_1 \eta_t}{\gamma_{t+1}} \delta_{t+1} + \boxed{C_1 \eta_{t+1} \delta_{t+1}} &\leq A_t + \frac{\eta_{t-1}}{\beta_t} \Delta_t + \frac{C_1 \eta_{t-1}}{\gamma_t} \delta_t + \boxed{C_1 \eta_t \delta_t} \\
- \eta_t B_t - \frac{1}{2} \eta_t \Gamma_t + \eta_t (\beta_{t+1} \sigma^2 + \boxed{\gamma_{t+1}^2 \sigma'^2}) &+ 2C_1 \eta_t \gamma_{t+1} \sigma'^2.
\end{aligned}$$

Define $Y_{t+1} = A_{t+1} + \frac{\eta_t}{\beta_{t+1}} \Delta_{t+1} + \frac{C_1 \eta_t}{\gamma_{t+1}} \delta_{t+1} + \boxed{C_1 \eta_{t+1} \delta_{t+1}}$, we have

$$\eta_t B_t + \frac{1}{2} \eta_t \Gamma_t \leq Y_t - Y_{t+1} + \eta_t (\beta_{t+1} \sigma^2 + \boxed{\gamma_{t+1}^2 \sigma'^2}) + 2C_1 \eta_t \gamma_{t+1} \sigma'^2.$$

Hence

$$\sum_{t=0}^{T-1} (\eta_t B_t + \frac{1}{2} \eta_t \Gamma_t) \leq Y_0 - A_* + \sum_{t=0}^{T-1} (\eta_t \beta_{t+1} \sigma^2 + 2C_1 \eta_t \gamma_{t+1} \sigma'^2 + \eta_t \gamma_{t+1}^2 \sigma'^2).$$

Next, let us consider $\eta_t = \eta, \beta_t = \beta, \gamma_t = \gamma$. Then we have

$$\sum_{t=0}^{T-1} \frac{1}{T} (B_t + \frac{1}{2} \Gamma_t) \leq \frac{C_Y}{T} + (\beta \sigma^2 + 2\gamma C_1 \sigma'^2 + \gamma^2 \sigma'^2).$$

Since $\eta_t = \eta, \gamma_t = \gamma, \beta_t = \beta$, in order to ensure the RHS is less than ϵ^2 , it suffices to have

$$\beta = \frac{\epsilon^2}{4\sigma^2}, \quad \gamma = \min(\frac{\epsilon^2}{8C_1 \sigma'^2}, \frac{\epsilon}{2\sigma''}), \quad T \geq \frac{C_Y}{4\epsilon^2 \eta}.$$

Since

$$\eta = \min\left(\frac{1}{L}, \frac{\beta}{\sqrt{4C_2}}, \frac{\gamma}{\sqrt{8C_1C_3}}\right).$$

Thus the order of T becomes

$$\begin{aligned} T &= O\left(\max\left\{\frac{C_Y L_F}{\epsilon^2}, \frac{C_Y \sqrt{C_2}}{\epsilon^2 \beta}, \frac{C_Y \sqrt{C_1 C_3}}{\gamma \epsilon^2}\right\}\right) \\ &= O\left(\max\left\{\frac{C_Y L_F}{\epsilon^2}, \frac{C_Y \sqrt{C_1 C_3} \sigma''}{\epsilon^3}, \frac{C_Y \sigma^2 \sqrt{C_2}}{\epsilon^4}, \frac{C_Y \sqrt{C_1 C_3} C_1 \sigma'^2}{\epsilon^4}\right\}\right) \end{aligned}$$

where

$$C_Y = A_0 - A_* + \frac{\eta}{\beta} \Delta_0 + \frac{C_1 \eta}{\gamma} \delta_0 + C_1 \eta \delta_0 \leq A_0 - A_* + \frac{1}{2\sqrt{C_2}} \Delta_0 + 2 \frac{\sqrt{C_1}}{\sqrt{8C_3}} \delta_0.$$

□

Finally, we state the convergence of SOX.

Theorem 5.1 *Under Assumption 5.1 and 5.2, SOX with $\beta = \frac{\epsilon^2}{4\sigma^2} < \frac{1}{4L_1 G_2}$, $\gamma = \min(\frac{\epsilon^2}{64G_2^2 L_1^2 \sigma_0^2}, \frac{n}{2BG_1 L_1 \sigma_0})$, $\eta = \min(\frac{1}{2L_F}, \frac{\beta}{2\sqrt{C_2}}, \frac{B\gamma}{n\sqrt{32C_1 C_3}})$, can find \mathbf{w}_τ with τ randomly sampled from $\{1, \dots, T\}$ so that $\mathbb{E}[\|\mathbf{v}_\tau\|_2^2 + \|\nabla F(\mathbf{w}_\tau)\|_2^2] \leq \epsilon^2$ with an iteration complexity of*

$$T = O\left(\max\left\{\frac{C_Y L_F}{\epsilon^2}, \frac{C_Y L_1^2 \sigma_0}{\epsilon^3}, \frac{C_Y L_F \sigma^2}{\epsilon^4}, \frac{C_Y L_1^3 n \sigma_0^2}{\epsilon^4 B}\right\}\right),$$

where $C_1 = 8G_2^2 L_1$, $C_2 = 4L_F^2 + 2$, $C_3 = 2G_2^2$, $\sigma^2 = \frac{G_1^2 \sigma_2^2}{B} + \frac{G_1^2 G_2^2}{B} \frac{n-B}{n-1}$, and $C_Y = O(F(\mathbf{w}_0) - F_* + \frac{1}{L_F} \|\mathbf{v}_0 - \nabla F(\mathbf{w}_0)\|_2^2 + L_1 \frac{1}{n} \|\mathbf{u}_0 - g(\mathbf{w}_0)\|_2^2)$.

💡 Why it matters

Theorem 5.1 shows that SOX achieves a complexity dominated by $O\left(\frac{C_Y L_1^3 n \sigma_0^2}{\epsilon^4 B}\right)$, which is comparable to that of SCMA for finding an ϵ -stationary solution. The key difference is that the complexity of SOX is scaled by a factor of n/B , since it must track and estimate n inner functions.

Proof. Assume that ϵ is sufficiently small such that $8\beta^2 G_2^2 L_1^2 \leq 1$. We have established the following three inequalities:

$$\begin{aligned}
(*) \mathbb{E} [F(\mathbf{w}_{t+1})] &\leq \mathbb{E} [F(\mathbf{w}_t)] + \frac{\eta}{2} \mathbb{E} [\Delta_t] - \frac{\eta}{2} \mathbb{E} [\|\nabla F(\mathbf{w}_t)\|_2^2] - \frac{\eta}{4} \mathbb{E} [\|\mathbf{v}_t\|_2^2], \\
(\sharp) \mathbb{E} [\Delta_{t+1}] &\leq (1 - \beta) \mathbb{E} [\Delta_t] + \frac{2L_F^2 + 1}{\beta} \eta^2 \mathbb{E} [\|\mathbf{v}_t\|_2^2] + 8\beta L_1^2 G_2^2 \mathbb{E} [\delta_t] \\
&\quad + \beta^2 \sigma^2 + 4G_2^2 L_1^2 \beta \gamma^2 \sigma_0^2, \\
(\diamond) \mathbb{E} [\delta_{t+1}] &\leq \left(1 - \frac{B\gamma}{2n}\right) \mathbb{E} [\delta_t] + \frac{2nG_2^2 \eta^2}{B\gamma} \mathbb{E} [\|\mathbf{v}_t\|_2^2] + \frac{B\gamma^2 \sigma_0^2}{n}.
\end{aligned}$$

Let us define $\bar{\gamma} = \frac{B\gamma}{2n}$, the last inequality becomes

$$(\diamond) \mathbb{E} [\delta_{t+1}] \leq (1 - \bar{\gamma}) \mathbb{E} [\delta_t] + \frac{G_2^2 \eta^2}{\bar{\gamma}} \mathbb{E} [\|\mathbf{v}_t\|_2^2] + \frac{4n\bar{\gamma}^2 \sigma_0^2}{B}.$$

Define $A_t = 2(F(\mathbf{w}_t) - F(\mathbf{w}_*))$ and $B_t = \|\nabla F(\mathbf{w}_t)\|_2^2$, $\Gamma_t = \|\mathbf{v}_t\|_2^2/2$, $\Delta_t = \|\nabla F(\mathbf{w}_t) - \mathbf{v}_t\|_2^2$, $\delta_t = \frac{1}{n} \|\mathbf{u}_t - g(\mathbf{w}_t)\|_2^2$, and $\Upsilon_t = A_t + \frac{\eta_{t-1}}{\beta_t} \Delta_t + \frac{C_1 \eta_{t-1}}{\bar{\gamma}_t} \delta_t$.

Then the three inequalities satisfy that in Lemma 4.10 with $C_1 = 8G_2^2 L_1^2$, $C_2 = 2(2L_F^2 + 1)$, $C_3 = 2G_2^2$, $\sigma^2 = \frac{G_1^2 \sigma_z^2}{B} + \frac{G_1^2 G_2^2}{B} \frac{n-B}{n-1}$, $\sigma'^2 = \frac{4n\sigma_0^2}{B}$, $\sigma''^2 = 4G_2^2 L_1^2 \sigma_0^2$. Then $\eta, \beta, \bar{\gamma}$ satisfy

$$\begin{aligned}
\beta &= \frac{\epsilon^2}{4\sigma^2}, \quad \bar{\gamma} = \min \left(\frac{\epsilon^2}{8C_1 \sigma'^2}, \frac{\epsilon}{2\sigma''} \right) = \min \left(\frac{\epsilon^2 B}{128G_2^2 L_1^2 n \sigma_0^2}, \frac{\epsilon}{4G_2 L_1 \sigma_0} \right), \\
\eta &= \min \left(\frac{1}{2L_F}, \frac{\beta}{\sqrt{4C_2}}, \frac{\bar{\gamma}}{\sqrt{8C_1 C_3}} \right).
\end{aligned}$$

Thus the order of T becomes

$$\begin{aligned}
T &= O \left(\max \left\{ \frac{C_Y L_F}{\epsilon^2}, \frac{C_Y \sqrt{C_1 C_3} \sigma''}{\epsilon^3}, \frac{C_Y \sigma^2 \sqrt{C_2}}{\epsilon^4}, \frac{C_Y \sqrt{C_1 C_3} C_1 \sigma'^2}{\epsilon^4} \right\} \right) \\
&= O \left(\max \left\{ \frac{C_Y L_F}{\epsilon^2}, \frac{C_Y L_1^2 \sigma_0}{\epsilon^3}, \frac{C_Y L_F \sigma^2}{\epsilon^4}, \frac{C_Y L_1^3 n \sigma_0^2}{\epsilon^4 B} \right\} \right),
\end{aligned}$$

where

$$\begin{aligned}
C_Y &\leq 2(F(\mathbf{w}_0) - F(\mathbf{w}_*)) + \frac{1}{2\sqrt{C_2}} \|\mathbf{v}_0 - \nabla F(\mathbf{w}_0)\|_2^2 + \frac{\sqrt{C_1}}{\sqrt{2C_3}} \frac{1}{n} \|\mathbf{u}_0 - g(\mathbf{w}_0)\|_2^2 \\
&= 2(F(\mathbf{w}_0) - F(\mathbf{w}_*)) + O \left(\frac{1}{L_F} \right) \|\mathbf{v}_0 - \nabla F(\mathbf{w}_0)\|_2^2 + O(L_1) \frac{1}{n} \|\mathbf{u}_0 - g(\mathbf{w}_0)\|_2^2.
\end{aligned}$$

□

5.2.2 Multi-block Single-Probe Variance Reduction

In this subsection, we present a second algorithm for solving FCCO with an improved complexity than that of SOX under a stronger condition on g_i . We replace Assumption 5.2 by the following:

Assumption 5.3. *There exist $G_1, L_1, L_2 > 0$ such that*

- (i) $f_i : \mathbb{R}^{d'} \mapsto \mathbb{R}$ is G_1 -Lipschitz continuous and L_1 -smooth;
- (ii) $\nabla g_i(\cdot, \zeta) : \mathbb{R}^d \mapsto \mathbb{R}^{d'}$ is mean-squared Lipschitz continuous, i.e.,

$$\mathbb{E}_{\zeta} [\|\nabla g_i(\mathbf{w}, \zeta) - \nabla g_i(\mathbf{w}', \zeta)\|_2^2] \leq L_2^2 \|\mathbf{w} - \mathbf{w}'\|_2^2, \forall \mathbf{w}, \mathbf{w}';$$

- (iii) $F_* = \min_{\mathbf{w}} F(\mathbf{w}) \geq -\infty$.

The idea is to leverage advanced variance reduction for tracking both the inner functions and the gradient. A straightforward approach is to change the update of $\mathbf{u}_{i,t-1}$ by using the STORM estimator and do similarly for the gradient estimator. In particular, one may change the update for $\mathbf{u}_{i,t}$ according to STORM:

$$\mathbf{u}_{i,t} = \begin{cases} (1 - \gamma_t)\mathbf{u}_{i,t-1} + \gamma_t g_i(\mathbf{w}_t; \zeta_{i,t}) + \underbrace{(1 - \gamma_t)(g_i(\mathbf{w}_t; \zeta_{i,t}) - g_i(\mathbf{w}_{t-1}; \zeta_{i,t}))}_{\text{error correction}} & i \in \mathcal{B}_t \\ \mathbf{u}_{i,t-1} & i \notin \mathcal{B}_t \end{cases} \quad (5.9)$$

However, this naive approach does not work as the standard error correction term marked above only accounts for the randomness in $g_i(\mathbf{w}_t; \zeta_{i,t})$ but not in the randomness caused by sampling $i \in \mathcal{B}_t$.

In order to tackle this challenge, we introduce the following estimator termed multi-block single-probe variance reduction estimator (MSVR):

$$\mathbf{u}_{i,t} = \begin{cases} (1 - \gamma_t)\mathbf{u}_{i,t-1} + \gamma_t g_i(\mathbf{w}_t; \zeta_{i,t}) + \gamma'_t (g_i(\mathbf{w}_t; \zeta_{i,t}) - g_i(\mathbf{w}_{t-1}; \zeta_{i,t})) & i \in \mathcal{B}_t \\ \mathbf{u}_{i,t-1} & i \notin \mathcal{B}_t \end{cases} \quad (5.10)$$

The difference from (5.9) lies at the value of γ'_t , which is set as $\frac{n-B}{B(1-\gamma_t)} + (1 - \gamma_t)$ with $B = |\mathcal{B}_t|$. The MSVR estimator can track multiple functional mappings (g_1, g_2, \dots, g_n) , simultaneously, while the number of sampled blocks B_1 for probing can be as small as one. It is notable that when $B = n$, i.e., all blocks are probed at each iteration, $\gamma'_t = 1 - \gamma_t$ and MSVR reduces to STORM applied to $\mathbf{g}(\mathbf{w})$. The additional factor in γ'_t , i.e., $\alpha_t = \frac{n-B}{B(1-\gamma_t)}$ is to account for the randomness in the sampled blocks and noise in those blocks that are not updated.

With \mathbf{u}_t , we compute a vanilla gradient estimator by

$$\mathbf{z}_t = \frac{1}{B} \sum_{i \in \mathcal{B}'_t} \nabla g_i(\mathbf{w}_t; \zeta'_{i,t}) \nabla f_i(\mathbf{u}_{i,t}),$$

where $\mathcal{B}'_t \subset [n]$ is a mini-batch of B indices independent of \mathcal{B}_t .

Similar to SCST, we apply another STORM estimator to estimate

$$\mathcal{M}_t = \mathbb{E}_{\mathcal{B}'_t, \zeta'_t}[\mathbf{z}_t] = \frac{1}{n} \sum_{i=1}^n \nabla g_i(\mathbf{w}_t) \nabla f_i(\mathbf{u}_{i,t}),$$

with an extra vanilla gradient estimator at previous iteration:

$$\tilde{\mathbf{z}}_{t-1} = \frac{1}{B} \sum_{i \in \mathcal{B}'_t} \nabla g_i(\mathbf{w}_{t-1}; \zeta'_{i,t}) \nabla f_i(\mathbf{u}_{i,t-1}).$$

This is computed by the following sequence:

$$\mathbf{v}_t = (1 - \beta_t) \mathbf{v}_{t-1} + \beta_t \mathbf{z}_t + (1 - \beta_t)(\mathbf{z}_t - \tilde{\mathbf{z}}_{t-1}). \quad (5.11)$$

Then we use \mathbf{v}_t to update the model parameter. The full steps are presented in Algorithm 15.

Critical: We use an independent batch \mathcal{B}'_t because \mathbf{z}_t depends on \mathbf{u}_t , which depends on \mathcal{B}_t . If we use the same batch \mathcal{B}_t to compute \mathbf{z}_t , then

$$\begin{aligned} \mathcal{M}_t &= \mathbb{E}_{\mathcal{B}_t, \zeta'_t} \left[\frac{1}{B} \sum_{i \in \mathcal{B}_t} \nabla g_i(\mathbf{w}_t; \zeta'_{i,t}) \nabla f_i(\mathbf{u}_{i,t}) \right] \\ &= \mathbb{E}_{\mathcal{B}_t, \zeta'_t} \left[\frac{1}{B} \sum_{i \in \mathcal{B}_t} \nabla g_i(\mathbf{w}_t; \zeta'_{i,t}) \nabla f_i(\bar{\mathbf{u}}_{i,t}) \right] = \frac{1}{n} \sum_{i=1}^n \nabla g_i(\mathbf{w}_t) \nabla f_i(\bar{\mathbf{u}}_{i,t}). \end{aligned}$$

where $\bar{\mathbf{u}}_t$ independent of \mathcal{B}_t is defined in (5.12). However, we cannot construct an unbiased estimator of \mathcal{M}_{t-1} since $\bar{\mathbf{u}}_{t-1}$ is not available in the algorithm.

An alternative approach is that we use \mathbf{u}_{t-1} and \mathbf{u}_{t-2} to compute \mathbf{z}_t and $\tilde{\mathbf{z}}_{t-1}$, respectively, with \mathcal{B}_t , i.e.,

$$\begin{aligned} \mathbf{z}_t &= \frac{1}{B} \sum_{i \in \mathcal{B}_t} \nabla g_i(\mathbf{w}_t; \zeta'_{i,t}) \nabla f_i(\mathbf{u}_{i,t-1}) \\ \tilde{\mathbf{z}}_{t-1} &= \frac{1}{B} \sum_{i \in \mathcal{B}_t} \nabla g_i(\mathbf{w}_{t-1}; \zeta'_{i,t}) \nabla f_i(\mathbf{u}_{i,t-2}), \end{aligned}$$

and compute \mathbf{v}_t by

$$\mathbf{v}_t = (1 - \beta_t) \mathbf{v}_t + \beta_t \mathbf{z}_t + (1 - \beta_t)(\mathbf{z}_t - \tilde{\mathbf{z}}_{t-1}).$$

Algorithm 15 MSVR

```

1: Input: learning rate schedules  $\{\eta_t\}_{t=1}^T, \{\gamma_t\}_{t=1}^T, \{\beta_t\}_{t=1}^T$ ; starting points  $\mathbf{w}_0, \mathbf{u}_0, \mathbf{v}_0$ 
2: Let  $\mathbf{w}_1 = \mathbf{w}_0 - \eta_0 \mathbf{v}_0$ 
3: for  $t = 1, \dots, T$  do
4:   Draw two batches of samples  $\mathcal{B}_t, \mathcal{B}'_t \subset [n]$ 
5:   for  $i \in \mathcal{B}_t$  do
6:     Draw two samples  $\zeta_{i,t}, \zeta'_{i,t} \sim \mathbb{P}_i$ 
7:     Update the inner function value estimators
           
$$\mathbf{u}_{i,t} = (1 - \gamma_t) \mathbf{u}_{i,t-1} + \gamma_t g_i(\mathbf{w}_t; \zeta_{i,t}) + \gamma'_t (g_i(\mathbf{w}_t; \zeta_{i,t}) - g_i(\mathbf{w}_{t-1}; \zeta_{i,t}))$$

8:   end for
9:   Set  $\mathbf{u}_{i,t} = \mathbf{u}_{i,t-1}, i \notin \mathcal{B}_t$ 
10:  Compute the vanilla gradient estimator  $\mathbf{z}_t = \frac{1}{B} \sum_{i \in \mathcal{B}'_t} \nabla g_i(\mathbf{w}_t; \zeta'_{i,t}) \nabla f_i(\mathbf{u}_{i,t})$ 
11:  Compute the extra vanilla gradient estimator  $\tilde{\mathbf{z}}_{t-1} = \frac{1}{B} \sum_{i \in \mathcal{B}'_t} \nabla g_i(\mathbf{w}_{t-1}; \zeta'_{i,t}) \nabla f_i(\mathbf{u}_{i,t-1})$ 
12:  Update the STORM gradient estimator  $\mathbf{v}_t = (1 - \beta_t) \mathbf{v}_{t-1} + \beta_t \mathbf{z}_t + (1 - \beta_t)(\mathbf{z}_t - \tilde{\mathbf{z}}_{t-1})$ 
13:  Update the model  $\mathbf{w}$  by  $\mathbf{w}_{t+1} = \mathbf{w}_t - \eta_t \mathbf{v}_t$ 
14: end for

```

The converge analysis can be performed similarly with slight modifications.

Convergence Analysis

We first analyze the error recursion of

$$\delta_t = \frac{1}{n} \|\mathbf{u}_t - \mathbf{g}(\mathbf{w}_t)\|_2^2 := \frac{1}{n} \sum_{i=1}^n \|\mathbf{u}_{i,t} - g_i(\mathbf{w}_t)\|_2^2.$$

Similar to the analysis of SOX, we introduce virtual sequences $\bar{\mathbf{u}}_{i,t}, \forall i$:

$$\bar{\mathbf{u}}_{i,t} = (1 - \gamma_t) \mathbf{u}_{i,t-1} + \gamma_t g_i(\mathbf{w}_t; \zeta_{i,t}) + \gamma'_t (g_i(\mathbf{w}_t; \zeta_{i,t}) - g_i(\mathbf{w}_{t-1}; \zeta_{i,t})), \forall i. \quad (5.12)$$

Lemma 5.5 Consider the \mathbf{u}_t updates in Algorithm 15. Under Assumption 5.1 and 5.3 (ii), by setting $\gamma'_t = \frac{n-B}{B(1-\gamma_t)} + (1 - \gamma_t)$, for $\gamma_t \leq \frac{1}{2}$, we have:

$$\mathbb{E}[\delta_t] \leq \left(1 - \frac{B\gamma_t}{n}\right) \mathbb{E}[\delta_{t-1}] + \frac{2B}{n} \gamma_t^2 \sigma_0^2 + \frac{12nG_2^2}{B} \mathbb{E}[\|\mathbf{w}_t - \mathbf{w}_{t-1}\|^2].$$

Proof. Let us consider a fixed $i \in [n]$. With a probability B/n that $i \in \mathcal{B}_t$, we have $\mathbf{u}_{i,t} = \bar{\mathbf{u}}_{i,t}$; otherwise $\mathbf{u}_{i,t} = \mathbf{u}_{i,t-1}$. Hence,

$$\mathbb{E}[\|\mathbf{u}_{i,t} - g_i(\mathbf{w}_t)\|_2^2] = \underbrace{\frac{B}{n} \mathbb{E}[\|\bar{\mathbf{u}}_{i,t} - g_i(\mathbf{w}_t)\|_2^2]}_{A_1} + \underbrace{\left(1 - \frac{B}{n}\right) \mathbb{E}[\|\mathbf{u}_{i,t-1} - g_i(\mathbf{w}_t)\|_2^2]}_{A_2}.$$

Note that the first term A_1 in the R.H.B. can be bounded similarly as in Lemma 4.12 for using the STORM estimator by building a recursion with $\|\mathbf{u}_{i,t-1} - g_i(\mathbf{w}_{t-1})\|_2^2$. However, there exists the second term due to the randomness of \mathcal{B}_t , which can be decomposed as

$$\begin{aligned} A_2 &= \mathbb{E}[\|\mathbf{u}_{i,t-1} - g_i(\mathbf{w}_{t-1}) + g_i(\mathbf{w}_{t-1}) - g_i(\mathbf{w}_t)\|_2^2] \\ &= \underbrace{\mathbb{E}[\|\mathbf{u}_{i,t-1} - g_i(\mathbf{w}_{t-1})\|_2^2]}_{A_{21}} + \underbrace{\mathbb{E}[\|g_i(\mathbf{w}_{t-1}) - g_i(\mathbf{w}_t)\|_2^2]}_{A_{22}} \\ &\quad + \underbrace{\mathbb{E}[2(\mathbf{u}_{i,t-1} - g_i(\mathbf{w}_{t-1}))^\top (g_i(\mathbf{w}_{t-1}) - g_i(\mathbf{w}_t))]}_{A_{23}}. \end{aligned}$$

The first two terms in RHS (A_{21} and A_{22}) can be easily handled. The difficulty comes from the third term, which cannot be simply bounded by using Young's inequality. If doing so, it will end up with a non-diminishing error of $\mathbf{u}_{i,t}$. To combat this difficulty, we use the additional factor brought by $\gamma'_t(g_i(\mathbf{w}_t; \xi_t^i) - g_i(\mathbf{w}_{t-1}; \xi_t^i))$ in A_1 to cancel A_{23} . This is more clear by the following decomposition of A_1 .

$$\begin{aligned} A_1 &= \underbrace{\mathbb{E}[\|(1 - \gamma_t)(\mathbf{u}_{i,t-1} - g_i(\mathbf{w}_{t-1}))\|_2^2]}_{A_{11}} + \underbrace{\mathbb{E}[\|\alpha_t(g_i(\mathbf{w}_t) - g_i(\mathbf{w}_{t-1}))\|_2^2]}_{A_{12}} \\ &\quad + \underbrace{\mathbb{E}[\|\gamma_t(g_i(\mathbf{w}_t; \zeta_{i,t}) - g_i(\mathbf{w}_t))\|_2^2]}_{A_{13}} \\ &\quad + \underbrace{\mathbb{E}[\|\gamma'_t(g_i(\mathbf{w}_t; \zeta_{i,t}) - g_i(\mathbf{w}_{t-1}; \xi_{i,t}) - g_i(\mathbf{w}_t) + g_i(\mathbf{w}_{t-1}))\|_2^2]}_{A_{14}}, \end{aligned}$$

where $\alpha_t = \gamma'_t + \gamma_t - 1$. Since $\mathbb{E}_t[A_{13}] = 0, \mathbb{E}_t[A_{14}] = 0$, then we have

$$A_1 \leq \mathbb{E}[\|A_{11} + A_{12}\|_2^2] + \mathbb{E}[\|A_{13} + A_{14}\|_2^2].$$

In light of the above decomposition, we can bound $\mathbb{E}[\|A_{11} + A_{12}\|_2^2] \leq \mathbb{E}[\|A_{11}\|_2^2 + \|A_{12}\|_2^2 + 2A_{11}^\top A_{12}]$ and $\mathbb{E}[\|A_{13} + A_{14}\|_2^2] \leq 2\mathbb{E}[\|A_{13}\|_2^2] + 2\mathbb{E}[\|A_{14}\|_2^2]$. The resulting term $\mathbb{E}[2A_{11}^\top A_{12}]$ has a negative sign as A_{23} . Hence, by carefully choosing γ'_t , we can cancel both terms. Specifically, we have

$$\begin{aligned}
\frac{B}{n}A_1 &\leq \frac{B}{n} \left(\mathbb{E}[\|A_{11}\|_2^2 + \|A_{12}\|_2^2 + 2A_{11}^\top A_{12}] + 2\mathbb{E}[\|A_{13}\|_2^2] + 2\mathbb{E}[\|A_{14}\|_2^2] \right) \\
&= \mathbb{E} \left[\frac{B}{n} (1 - \gamma_t)^2 \|\mathbf{u}_{i,t-1} - g_i(\mathbf{w}_{t-1})\|_2^2 \right] + \mathbb{E} \left[\frac{B}{n} \alpha_t^2 \|g_i(\mathbf{w}_t) - g_i(\mathbf{w}_{t-1})\|_2^2 \right] \\
&\quad + \mathbb{E} \left[\frac{B}{n} 2\alpha_t (1 - \gamma_t) (g_i(\mathbf{w}_t) - g_i(\mathbf{w}_{t-1}))^\top (\mathbf{u}_{i,t-1} - g_i(\mathbf{w}_{t-1})) \right] \\
&\quad + \mathbb{E} \left[\frac{B}{n} 2\gamma_t^2 \|g_i(\mathbf{w}_t; \zeta_{i,t}) - g_i(\mathbf{w}_t)\|_2^2 \right] \\
&\quad + \mathbb{E} \left[\frac{B}{n} 2\gamma_t'^2 \|(g_i(\mathbf{w}_t; \zeta_{i,t}) - g_i(\mathbf{w}_{t-1}; \zeta_{i,t}) - g_i(\mathbf{w}_t) + g_i(\mathbf{w}_{t-1}))\|_2^2 \right].
\end{aligned}$$

Combining the upper bounds of A_1 and A_2 , we have

$$\begin{aligned}
&\frac{B}{n}A_1 + \frac{n-B}{n}A_2 \\
&\leq \mathbb{E} \left[\frac{B}{n} (1 - \gamma_t)^2 \|\mathbf{u}_{i,t-1} - g_i(\mathbf{w}_{t-1})\|_2^2 + \frac{B}{n} \alpha_t^2 \|g_i(\mathbf{w}_t) - g_i(\mathbf{w}_{t-1})\|_2^2 \right] \\
&\quad + \mathbb{E} \left[\frac{B}{n} 2\alpha_t (1 - \gamma_t) (g_i(\mathbf{w}_t) - g_i(\mathbf{w}_{t-1}))^\top (\mathbf{u}_{i,t-1} - g_i(\mathbf{w}_{t-1})) \right] \\
&\quad + \mathbb{E} \left[\frac{B}{n} 2\gamma_t^2 \|g_i(\mathbf{w}_t; \zeta_{i,t}) - g_i(\mathbf{w}_t)\|_2^2 \right] \\
&\quad + \mathbb{E} \left[\frac{B}{n} 2(\gamma_t')^2 \|(g_i(\mathbf{w}_t; \zeta_{i,t}) - g_i(\mathbf{w}_{t-1}; \zeta_{i,t}) - g_i(\mathbf{w}_t) + g_i(\mathbf{w}_{t-1}))\|_2^2 \right] \\
&\quad + \mathbb{E} \left[\frac{n-B}{n} \|\mathbf{u}_{i,t-1} - g_i(\mathbf{w}_{t-1})\|_2^2 \right] + \mathbb{E} \left[\frac{n-B}{n} \|g_i(\mathbf{w}_{t-1}) - g_i(\mathbf{w}_t)\|_2^2 \right] \\
&\quad + \mathbb{E} \left[\frac{n-B}{n} 2(\mathbf{u}_{i,t-1} - g_i(\mathbf{w}_{t-1}))^\top (g_i(\mathbf{w}_{t-1}) - g_i(\mathbf{w}_t)) \right].
\end{aligned}$$

Since $\frac{B}{n} 2\alpha_t (1 - \gamma_t) = 2 \frac{B}{n} \frac{(n-B)}{B(1-\gamma_t)} (1 - \gamma_t) = 2 \frac{n-B}{n}$, then cross terms will cancel out. The remaining terms can be merged and handled separately. First,

$$\begin{aligned}
&\mathbb{E} \left[\frac{B}{n} (1 - \gamma_t)^2 \|\mathbf{u}_{i,t-1} - g_i(\mathbf{w}_{t-1})\|_2^2 + \frac{n-B}{n} \|\mathbf{u}_{i,t-1} - g_i(\mathbf{w}_{t-1})\|_2^2 \right] \\
&\leq \left(1 - \frac{B}{n} \gamma_t \right) \mathbb{E} [\|\mathbf{u}_{i,t-1} - g_i(\mathbf{w}_{t-1})\|_2^2],
\end{aligned}$$

where we use $\frac{B}{n} (1 - \gamma_t)^2 + \frac{n-B}{n} \leq 1 - \frac{2B}{n} \gamma_t + \frac{B}{n} \gamma_t^2 \leq 1 - \frac{B}{n} \gamma_t$ due to $\gamma_t < 1$. Second

$$\begin{aligned}
&\frac{B}{n} \alpha_t^2 \|g_i(\mathbf{w}_t) - g_i(\mathbf{w}_{t-1})\|_2^2 + \frac{n-B}{n} \|g_i(\mathbf{w}_{t-1}) - g_i(\mathbf{w}_t)\|_2^2 \\
&\leq \left(\frac{B}{n} \frac{(n-B)^2}{B^2(1-\gamma_t)^2} + \frac{n-B}{n} \right) G_2^2 \|\mathbf{w}_t - \mathbf{w}_{t-1}\|_2^2 \leq \frac{4n-4B}{B} G_2^2 \|\mathbf{w}_t - \mathbf{w}_{t-1}\|_2^2,
\end{aligned}$$

where we use $\frac{B}{n} \frac{(n-B)^2}{B^2(1-\gamma_t)^2} + \frac{n-B}{n} \leq \frac{n-B}{n} \left(\frac{(n-B)}{B(1-\gamma_t)^2} + 1 \right) \leq \frac{n-B}{n} \left(\frac{4(n-B)}{B} + 4 \right) = \frac{4n-4B}{B}$ due to $\gamma_t \leq 1/2$. Third,

$$\begin{aligned} & \mathbb{E} \left[\frac{B}{n} 2\gamma_t'^2 \left\| (g_i(\mathbf{w}_t; \zeta_{i,t}) - g_i(\mathbf{w}_{t-1}; \zeta_{i,t}) - g_i(\mathbf{w}_t) + g_i(\mathbf{w}_{t-1})) \right\|_2^2 \right] \\ & \leq \frac{B}{n} 2\gamma_t'^2 \mathbb{E} \left[\left\| (g_i(\mathbf{w}_t; \zeta_{i,t}) - g_i(\mathbf{w}_{t-1}; \zeta_{i,t})) \right\|_2^2 \right] \\ & \leq \frac{B}{n} 2 \left(\frac{n-B}{B(1-\gamma_t)} + 1 - \gamma_t \right)^2 G_2^2 \|\mathbf{w}_t - \mathbf{w}_{t-1}\|_2^2 \leq \frac{8n-4B}{B} G_2^2 \|\mathbf{w}_t - \mathbf{w}_{t-1}\|_2^2, \end{aligned}$$

where we use $\frac{B}{n} 2 \left(\frac{n-B}{B(1-\gamma_t)} + 1 - \gamma_t \right)^2 \leq \frac{B}{n} 2 \left(\frac{2(n-B)}{B} + 1 \right)^2 \leq \frac{B}{n} 2 \left(\frac{2n-B}{B} \right)^2 = \frac{2(2n-B)(2n-B)}{nB} \leq \frac{8n-4B}{B}$.

Combining the above results, we have

$$\begin{aligned} & \mathbb{E} [\|\mathbf{u}_{i,t} - g_i(\mathbf{w}_t)\|_2^2] \leq \left(1 - \frac{B}{n} \gamma_t\right) \mathbb{E} [\|\mathbf{u}_{i,t-1} - g_i(\mathbf{w}_{t-1})\|_2^2] \\ & + \frac{12n-8B}{B} G_2^2 \|\mathbf{w}_t - \mathbf{w}_{t-1}\|_2^2 + \frac{B}{n} 2\gamma_t'^2 \sigma_0^2. \end{aligned}$$

Averaging over $i = 1, \dots, n$ concludes the proof. \square

Lemma 5.6 Consider the \mathbf{u}_t updates in Algorithm 15. Suppose that Assumption 5.1 and 5.3 hold. With $\gamma_t \leq \frac{1}{2}$ and $\gamma_t' = \frac{n-B}{B(1-\gamma_t)} + (1 - \gamma_t)$, we have

$$\mathbb{E} [\|\mathbf{u}_t - \mathbf{u}_{t-1}\|_2^2] \leq 6B\gamma_t^2 \sigma_0^2 + 6B\gamma_t'^2 \mathbb{E} [\delta_{t-1}] + \frac{10n^2 G_2^2}{B} \mathbb{E} [\|\mathbf{w}_t - \mathbf{w}_{t-1}\|_2^2].$$

Proof. Since $\|\mathbf{u}_t - \mathbf{u}_{t-1}\|_2^2 = \sum_{i=1}^n \|\mathbf{u}_{i,t} - \mathbf{u}_{i,t-1}\|_2^2$, with a probability B/n we have $\mathbf{u}_{i,t} = \bar{\mathbf{u}}_{i,t}$ and a probability $1 - B/n$ we have $\mathbf{u}_{i,t} = \mathbf{u}_{i,t-1}$, then

$$\begin{aligned} & \mathbb{E} [\|\mathbf{u}_t - \mathbf{u}_{t-1}\|_2^2] \\ & = \frac{B}{n} \sum_{i=1}^n \mathbb{E} \left[\left\| -\gamma_t \mathbf{u}_{i,t-1} + \gamma_t g_i(\mathbf{w}_t; \zeta_{i,t}) + \gamma_t' (g_i(\mathbf{w}_t; \zeta_{i,t}) - g_i(\mathbf{w}_{t-1}; \zeta_{i,t})) \right\|_2^2 \right] \\ & \leq \frac{B}{n} \sum_{i=1}^n \mathbb{E} \left[2\gamma_t^2 \|g_i(\mathbf{w}_t; \zeta_{i,t}) - \mathbf{u}_{i,t-1}\|_2^2 + 2(\gamma_t')^2 \|g_i(\mathbf{w}_t; \zeta_{i,t}) - g_i(\mathbf{w}_{t-1}; \zeta_{i,t})\|_2^2 \right] \\ & \leq \frac{B}{n} \sum_{i=1}^n \mathbb{E} \left[2\gamma_t^2 \|g_i(\mathbf{w}_t; \zeta_{i,t}) - \mathbf{u}_{i,t-1}\|_2^2 \right] + 2B(\gamma_t')^2 G_2^2 \|\mathbf{w}_t - \mathbf{w}_{t-1}\|_2^2. \end{aligned}$$

To the first term on the RHS, we use the Young's inequality and Lipschitz continuity of g_i :

$$\begin{aligned}
\mathbb{E} \left[\|g_i(\mathbf{w}_t; \zeta_{i,t}) - \mathbf{u}_{i,t-1}\|_2^2 \right] &\leq 3\mathbb{E} \left[\|g_i(\mathbf{w}_t; \zeta_{i,t}) - g_i(\mathbf{w}_t)\|_2^2 \right] \\
&+ 3\mathbb{E} \left[\|g_i(\mathbf{w}_t) - g_i(\mathbf{w}_{t-1})\|_2^2 \right] + 3\mathbb{E} \left[\|g_i(\mathbf{w}_{t-1}) - \mathbf{u}_{i,t-1}\|_2^2 \right] \\
&\leq 3\sigma_0^2 + 3G_2^2\mathbb{E} \left[\|\mathbf{w}_t - \mathbf{w}_{t-1}\|_2^2 \right] + 3\mathbb{E} \left[\|g_i(\mathbf{w}_{t-1}) - \mathbf{u}_{i,t-1}\|_2^2 \right].
\end{aligned}$$

Combining the above results, we have

$$\begin{aligned}
&\mathbb{E} \left[\|\mathbf{u}_t - \mathbf{u}_{t-1}\|_2^2 \right] \\
&\leq 6B\gamma_t^2\sigma_0^2 + 6B\gamma_t^2\mathbb{E}[\delta_{t-1}] + 2BG_2^2(3\gamma_t^2 + (\gamma'_t)^2)\mathbb{E} \left[\|\mathbf{w}_t - \mathbf{w}_{t-1}\|_2^2 \right].
\end{aligned}$$

With $\gamma_t \leq \frac{1}{2}$, we have $\gamma'_t \leq \frac{2n}{B}$, which yields $(3\gamma_t^2 + (\gamma'_t)^2) \leq \frac{5n^2}{B^2}$. \square

Next, we analyze error recursion of $\Delta_t := \|\mathbf{v}_t - \mathcal{M}_t\|_2^2$.

Lemma 5.7 Consider the \mathbf{v}_t updates in Algorithm 15 and suppose that Assumption 5.1 and 5.3 hold. Then we have

$$\begin{aligned}
\mathbb{E}[\Delta_t] &\leq (1 - \beta_t)\mathbb{E}[\Delta_{t-1}] + \frac{24G_2^2L_1^2B\gamma_t^2}{n}\mathbb{E}[\delta_{t-1}] \\
&+ \left(4L_2^2G_1^2 + \frac{40G_2^4L_1^2n}{B} \right) \mathbb{E} \left[\|\mathbf{w}_{t-1} - \mathbf{w}_t\|_2^2 \right] + 2\beta_t^2\sigma^2 + \frac{24G_2^2L_1^2B}{n}\gamma_t^2\sigma_0^2,
\end{aligned}$$

where $\sigma^2 = \frac{G_1^2\sigma_2^2}{B} + \frac{G_1^2G_2^2}{B} \frac{n-B}{n-1}$.

Proof. Similar to Lemma 5.2, we have $\mathbb{E}_t[\|\mathbf{z}_t - \mathcal{M}_t\|_2^2] \leq \sigma^2$. Since $\mathbf{v}_t = (1 - \beta_t)\mathbf{v}_{t-1} + \beta_t\mathbf{z}_t + (1 - \beta_t)(\mathbf{z}_t - \tilde{\mathbf{z}}_{t-1})$, applying Lemma 4.11, we have

$$\mathbb{E}_t \left[\|\mathbf{v}_t - \mathcal{M}_t\|_2^2 \right] \leq (1 - \beta_t) \|\mathbf{v}_{t-1} - \mathcal{M}_{t-1}\|_2^2 + \mathbb{E}_t[2\|\mathbf{z}_t - \tilde{\mathbf{z}}_{t-1}\|_2^2] + 2\beta_t^2\sigma^2.$$

To bound $\mathbb{E}_t[\|\mathbf{z}_t - \tilde{\mathbf{z}}_{t-1}\|_2^2]$, we have

$$\begin{aligned}
&\mathbb{E}_t[\|\mathbf{z}_t - \tilde{\mathbf{z}}_{t-1}\|_2^2] \\
&\leq 2\mathbb{E}_t \left[\frac{1}{B} \sum_{i \in \mathcal{B}'_t} \|\nabla g_i(\mathbf{w}_t; \zeta'_{i,t}) \nabla f_i(\mathbf{u}_{i,t}) - \nabla g_i(\mathbf{w}_t; \zeta'_{i,t}) \nabla f_i(\mathbf{u}_{i,t-1})\|_2^2 \right] \\
&+ 2\mathbb{E}_t \left[\frac{1}{B} \sum_{i \in \mathcal{B}'_t} \|\nabla g_i(\mathbf{w}_t; \zeta'_{i,t}) \nabla f_i(\mathbf{u}_{i,t-1}) - \nabla g_i(\mathbf{w}_{t-1}; \zeta'_{i,t}) \nabla f_i(\mathbf{u}_{i,t-1})\|_2^2 \right] \\
&\leq 2G_2^2L_1^2\mathbb{E}_t \left[\frac{1}{B} \sum_{i \in \mathcal{B}'_t} \|\mathbf{u}_{i,t} - \mathbf{u}_{i,t-1}\|_2^2 \right] + 2L_2^2G_1^2\|\mathbf{w}_t - \mathbf{w}_{t-1}\|_2^2 \\
&= 2G_2^2L_1^2\mathbb{E}_t \left[\frac{1}{n} \sum_{i=1}^n \|\mathbf{u}_{i,t} - \mathbf{u}_{i,t-1}\|_2^2 \right] + 2L_2^2G_1^2\|\mathbf{w}_t - \mathbf{w}_{t-1}\|_2^2,
\end{aligned}$$

where the last inequality follows the Assumption 5.3.

As a result, we have

$$\begin{aligned}\mathbb{E}[\Delta_t] &\leq (1 - \beta_t)\mathbb{E}[\Delta_{t-1}] + \frac{4G_2^2L_1^2}{n}\mathbb{E}[\|\mathbf{u}_t - \mathbf{u}_{t-1}\|_2^2] + 4L_2^2G_1^2\mathbb{E}[\|\mathbf{w}_{t-1} - \mathbf{w}_t\|_2^2] \\ &\quad + 2\beta_t^2\sigma^2.\end{aligned}$$

Combining with the result in Lemma 5.6, i.e.,

$$\mathbb{E}[\|\mathbf{u}_t - \mathbf{u}_{t-1}\|_2^2] \leq 6B\gamma_t^2\sigma_0^2 + 6B\gamma_t^2\mathbb{E}[\delta_{t-1}] + \frac{10n^2G_2^2}{B}\mathbb{E}[\|\mathbf{w}_t - \mathbf{w}_{t-1}\|_2^2].$$

we have

$$\begin{aligned}\mathbb{E}[\Delta_t] &\leq (1 - \beta_t)\mathbb{E}[\Delta_{t-1}] + \frac{24BG_2^2L_1^2}{n}\gamma_t^2\mathbb{E}[\delta_{t-1}] \\ &\quad + \left(4L_2^2G_1^2 + \frac{40G_2^4L_1^2n}{B}\right)\mathbb{E}[\|\mathbf{w}_{t-1} - \mathbf{w}_t\|_2^2] + 2\beta_t^2\sigma^2 + \frac{24BG_2^2L_1^2\gamma_t^2\sigma_0^2}{n},\end{aligned}$$

which completes the proof. \square

Lemma 5.8 For the update $\mathbf{w}_{t+1} = \mathbf{w}_t - \eta_t \mathbf{v}_t$, $t \geq 0$, if $\eta_t \leq 1/(2L_F)$ we have

$$F(\mathbf{w}_{t+1}) \leq F(\mathbf{w}_t) + G_2^2L_1^2\eta_t\delta_t + \eta_t\Delta_t - \frac{\eta_t}{2}\|\nabla F(\mathbf{w}_t)\|_2^2 - \frac{1}{4\eta_t}\|\mathbf{w}_{t+1} - \mathbf{w}_t\|_2^2. \quad (5.13)$$

Proof. It follows directly from Lemma 4.9 by noting that

$$\begin{aligned}\|\mathbf{v}_t - \nabla F(\mathbf{w}_t)\|_2^2 &= \|\mathbf{v}_t - \mathcal{M}_t + \mathcal{M}_t - \nabla F(\mathbf{w}_t)\|_2^2 \\ &\leq 2\Delta_t + 2\left\|\frac{1}{n}\sum_{i=1}^n \nabla g_i(\mathbf{w}_t)\nabla f_i(\mathbf{u}_{i,t}) - \frac{1}{n}\sum_{i=1}^n \nabla g_i(\mathbf{w}_t)\nabla f_i(g_i(\mathbf{w}_t))\right\|_2^2 \\ &\leq 2\Delta_t + \frac{2G_2^2L_1^2}{n}\sum_{i=1}^n \|\mathbf{u}_{i,t} - g_i(\mathbf{w}_t)\|_2^2.\end{aligned}$$

Taking expectation over all randomness on both sides yields the desired result. \square

Now we state the convergence theorem for MSVR.

Theorem 5.2 Suppose that Assumption 5.1 and 5.3 hold. Let $\beta = O(\frac{\epsilon\eta L_1\sqrt{n}}{\sigma\sqrt{B}})$, $\gamma = \min\left(\frac{\epsilon\eta L_1n}{\sigma_0B}, 1\right)$, $\eta = \min\left(\frac{1}{2L_F}, O(\frac{\epsilon\sqrt{B}}{L_1\sigma\sqrt{n}}), O(\frac{\epsilon B}{L_1^2\sigma_0n}), O(\frac{B}{nL_1})\right)$. Then MSVR can find \mathbf{w}_τ that is sampled randomly from $\{0, \dots, T-1\}$ satisfying

$$\mathbb{E} [\|\mathbf{v}_\tau\|_2^2 + \|\nabla F(\mathbf{w}_\tau)\|_2^2] \leq O(\epsilon).$$

with an iteration complexity of

$$T = O\left(\max\left\{\frac{C_Y L_F}{\epsilon^2}, \frac{C_Y L_1 n}{\epsilon^2 B}, \frac{C_Y L_1 \sigma \sqrt{n}}{\epsilon^3 \sqrt{B}}, \frac{C_Y L_1^2 \sigma_0 n}{\epsilon^3 B}\right\}\right).$$

where $\sigma^2 = \frac{G_1^2 \sigma_z^2}{B} + \frac{G_1^2 G_2^2 (n-B)}{B(n-1)}$, $C_Y = O(F(\mathbf{w}_0) - F_* + \frac{B}{nL_1^2 \eta} \Delta_0 + \frac{B}{nL_1^2 \eta} \delta_0)$.

Why it matters

Theorem 5.2 indicates that when the initial estimators \mathbf{u}_0 and \mathbf{v}_0 have an estimation error in the order of $O(\epsilon)$ such that C_Y is $O(1)$, MSVR attains a better complexity than SOX for finding an ϵ -stationary solution under stronger assumptions of the mean-Lipschitz continuity of g and ∇g . Its complexity is comparable to that of SCST in Theorem 4.4, up to a factor of n/B .

Proof. We have established the following:

$$\begin{aligned} (*) \quad & F(\mathbf{w}_{t+1}) \leq F(\mathbf{w}_t) + G_2^2 L_1^2 \eta_t \delta_t + \eta_t \Delta_t - \frac{\eta_t}{2} \|\nabla F(\mathbf{w}_t)\|_2^2 - \frac{1}{4\eta_t} \|\mathbf{w}_{t+1} - \mathbf{w}_t\|_2^2 \\ (\#) \quad & \mathbb{E}[\Delta_t] \leq (1 - \beta_t) \mathbb{E}[\Delta_{t-1}] + \frac{24BG_2^2 L_1^2}{n} \gamma_t^2 \mathbb{E}[\delta_{t-1}] \\ & + \left(4L_2^2 G_1^2 + \frac{40G_2^4 L_1^2 n}{B}\right) \mathbb{E}[\|\mathbf{w}_{t-1} - \mathbf{w}_t\|_2^2] + 2\beta_t^2 \sigma^2 + \frac{24BG_2^2 L_1^2 \sigma_0^2}{n} \gamma_t^2, \\ (\diamond) \quad & \mathbb{E}[\delta_t] \leq \left(1 - \frac{B\gamma_t}{n}\right) \mathbb{E}[\delta_{t-1}] + \frac{2B}{n} \gamma_t^2 \sigma_0^2 + \frac{12nG_2^2}{B} \mathbb{E}[\|\mathbf{w}_t - \mathbf{w}_{t-1}\|_2^2]. \end{aligned}$$

In order to apply Lemma 4.15, we let $A_t = F(\mathbf{w}_t) - F_*$, $B_t = \|\nabla F(\mathbf{w}_t)\|_2^2/2$, $\Gamma_t = \|\mathbf{v}_t\|_2^2/4$, $\bar{\delta}_t = L_1^2 G_2^2 \delta_t$, $\bar{\gamma}_t = \frac{B\gamma_t}{n}$. Then the following three inequalities

$$\begin{aligned} (*) \quad & \mathbb{E}[A_{t+1}] \leq \mathbb{E}[A_t + \eta_t \Delta_t + \eta_t \bar{\delta}_t - \eta_t B_t - \eta_t \Gamma_t] \\ (\#) \quad & \mathbb{E}[\Delta_{t+1}] \leq \mathbb{E}[(1 - \beta_{t+1})\Delta_t + C_1 \bar{\gamma}_{t+1}^2 \bar{\delta}_t + C_2 \eta_t^2 \Gamma_t + \beta_{t+1}^2 \sigma^2 + \bar{\gamma}_{t+1}^2 \sigma'^2], \\ (\diamond) \quad & \mathbb{E}[\bar{\delta}_{t+1}] \leq \mathbb{E}[(1 - \bar{\gamma}_{t+1})\bar{\delta}_t + C_3 \eta_t^2 \Gamma_t + \bar{\gamma}_{t+1}^2 \sigma''^2]. \end{aligned}$$

hold with $C_1 = O(n/B)$, $C_2 = O(L_1^2 n/B + L_2^2)$, $C_3 = O(L_1^2 n/B)$, $\sigma^2 = \frac{G_1^2 \sigma_z^2}{B} + \frac{G_1^2 G_2^2 (n-B)}{B(n-1)}$, $\sigma'^2 = O(L_1^2 \sigma_0^2 n/B)$, $\sigma''^2 = O(L_1^2 \sigma_0^2 n/B)$. Following the settings in Lemma 4.15, we can finish the proof with

	f_i			g_i			F
	Lipschitz continuity	Weak convexity	Monotonicity	Lipschitz continuity	Weak convexity	Smoothness	Weak convexity (ρ)
5.5(i)	G_1	ρ_1	$\partial f \geq 0$	G_2	ρ_2	-	$G_1\rho_2\sqrt{d'} + \rho_1 G_2^2$
5.5(ii)	G_1	ρ_1	$\frac{\partial f}{\partial f} \geq 0$ or $\frac{\partial f}{\partial f} \leq 0$	G_2	-	L_2	$G_1 L_2 \sqrt{d'} + \rho_1 G_2^2$

Table 5.1: Conditions of f_i and g_i to make $F(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n f_i(g_i(\mathbf{w}))$ weakly convex, where $g_i : \mathbb{R}^d \rightarrow \mathbb{R}^{d'}$ and $f_i : \mathbb{R}^{d'} \rightarrow \mathbb{R}$.

$$\begin{aligned}
\eta &= \min \left(\frac{1}{L}, \frac{\epsilon}{4\sqrt{C_2}\sigma}, \frac{\epsilon\sqrt{C_2}}{8C_3\sigma'}, \frac{\epsilon}{8\sqrt{C_3}\sigma''}, \frac{\sqrt{C_2}}{4C_3\sqrt{C_1}} \right) \\
&= \min \left(\frac{1}{2L_F}, O\left(\frac{\epsilon}{L_1\sigma}\sqrt{\frac{B}{n}}\right), O\left(\frac{\epsilon B}{L_1^2\sigma_0 n}\right), O\left(\frac{B}{nL_1}\right) \right), \\
\beta &= \frac{\epsilon\eta\sqrt{2C_2}}{2\sigma} = O\left(\frac{\epsilon\eta L_1}{2\sigma}\sqrt{\frac{n}{B}}\right), \\
\bar{\gamma} &= \min \left(\frac{\epsilon\eta\sqrt{C_2}}{\sigma'}, \frac{\epsilon\eta\sqrt{C_3}}{\sigma''}, \frac{C_2}{2C_3C_1} \right) = \min \left(O\left(\frac{\epsilon\eta}{\sigma_0}\right), O\left(\frac{B}{n}\right) \right), \\
T &= O\left(\max \left\{ \frac{C_Y L_F}{\epsilon^2}, \frac{C_Y L_1 n}{\epsilon^2 B}, \frac{C_Y L_1 \sigma \sqrt{n}}{\epsilon^3 \sqrt{B}}, \frac{C_Y L_1^2 \sigma_0 n}{\epsilon^3 B} \right\} \right).
\end{aligned}$$

where $C_Y = F(\mathbf{w}_0) - F_* + \frac{1}{4C_2\eta}\Delta_0 + \frac{1}{4C_3\eta}\delta_0$.

□

5.3 Non-Smooth Weakly Convex Functions

In this section, we consider non-smooth weakly convex functions, where either the outer function or the inner function are non-smooth. The group DRO objective (5.2) falls into this category. Another instance is the two-way partial AUC maximization problem as discussed in Section 6.4.3.

Assumption 5.4. We assume that

- (i) $\mathbb{E}_{\zeta \sim \mathbb{P}_i} [\|g_i(\mathbf{w}; \zeta) - g_i(\mathbf{w})\|_2^2] \leq \sigma_0^2$.
- (ii) $\mathbb{E}_{\zeta \sim \mathbb{P}_i} [\|\mathcal{G}_i(\mathbf{w}; \zeta)\|_2^2] \leq G_2^2$ for any $\mathcal{G}_i(\mathbf{w}; \zeta) \in \partial g_i(\mathbf{w}; \zeta)$.

The second condition above implies that g_i is G_2 -Lipschitz continuous.

Assumption 5.5. We assume either of the following conditions holds:

- (i) f_i is ρ_1 -weakly convex, G_1 -Lipschitz continuous, and $\partial f_i(g) \geq 0 \forall g$; g_i is ρ_2 -weakly convex.

(ii) f_i is ρ_1 -weakly convex, G_1 -Lipschitz continuous, and $\partial f_i(g) \geq 0$ or $\partial f_i(g) \leq 0 \forall g$; and g_i is L_2 -smooth.

We first characterize the conditions on f_i and g_i to induce weak convexity of F .

Lemma 5.9 *Under Assumption 5.4 and 5.5, the objective function F is ρ -weakly convex for some $\rho > 0$. If Assumption 5.5(i) holds, then $\rho = G_1 \rho_2 \sqrt{d'} + \rho_1 G_2^2$ and if Assumption 5.5(ii) holds, then $\rho = G_1 L_2 \sqrt{d'} + \rho_1 G_2^2$.*

Proof. The weak convexity of f_i implies that for any $\mathbf{v}_i \in \partial f_i(g_i(\mathbf{w}))$:

$$\begin{aligned} f_i(g_i(\mathbf{w}')) &\geq f_i(g_i(\mathbf{w})) + \mathbf{v}_i^\top (g_i(\mathbf{w}') - g_i(\mathbf{w})) - \frac{\rho_1}{2} \|g_i(\mathbf{w}') - g_i(\mathbf{w})\|_2^2 \\ &\geq f_i(g_i(\mathbf{w})) + \mathbf{v}_i^\top (g_i(\mathbf{w}') - g_i(\mathbf{w})) - \frac{\rho_1 G_2^2}{2} \|\mathbf{w} - \mathbf{w}'\|_2^2. \end{aligned}$$

Let us first prove the weak convexity under Assumption 5.5(i). Since g_i is ρ_2 -weakly convex, we have for any $U_i \in \partial g_i(\mathbf{w})$

$$g_i(\mathbf{w}') - g_i(\mathbf{w}) \geq U_i^\top (\mathbf{w}' - \mathbf{w}) - \frac{\rho_2}{2} \|\mathbf{w}' - \mathbf{w}\|_2^2 \mathbf{1}. \quad (5.14)$$

where $\mathbf{1}$ denotes a vector of all ones. Since $\mathbf{v}_i \geq 0$, we have

$$\begin{aligned} f_i(g_i(\mathbf{w}')) &\geq f_i(g_i(\mathbf{w})) + \mathbf{v}_i^\top (U_i^\top (\mathbf{w}' - \mathbf{w}) - \frac{\rho_2}{2} \|\mathbf{w} - \mathbf{w}'\|_2^2 \mathbf{1}) - \frac{\rho_1 G_2^2}{2} \|\mathbf{w} - \mathbf{w}'\|_2^2 \\ &\geq f_i(g_i(\mathbf{w})) + (U_i \mathbf{v}_i)^\top (\mathbf{w}' - \mathbf{w}) - \frac{G_1 \sqrt{d'} \rho_2 + \rho_1 G_2^2}{2} \|\mathbf{w} - \mathbf{w}'\|_2^2 \end{aligned}$$

Since $U_i \mathbf{v}_i \in \partial g_i(\mathbf{w}) \partial f_i(g_i(\mathbf{w}))$, the above inequality indicates that $f_i(g_i(\mathbf{w}))$ is ρ -weakly convex, where $\rho = G_1 \sqrt{d'} \rho_2 + \rho_1 G_2^2$. As a result, $F(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n f_i(g_i(\mathbf{w}))$ is ρ -weakly convex.

Next, we prove the weak convexity of F under Assumption 5.5(ii). Due to the smoothness of $g(\cdot)$ we have

$$\begin{aligned} g(\mathbf{w}) - g(\mathbf{w}') &\leq \nabla g(\mathbf{w}')^\top (\mathbf{w} - \mathbf{w}') + \frac{L_2}{2} \|\mathbf{w} - \mathbf{w}'\|_2^2 \mathbf{1}, \\ g(\mathbf{w}) - g(\mathbf{w}') &\geq \nabla g(\mathbf{w}')^\top (\mathbf{w} - \mathbf{w}') - \frac{L_2}{2} \|\mathbf{w} - \mathbf{w}'\|_2^2 \mathbf{1}. \end{aligned} \quad (5.15)$$

If $\partial f_i(g_i(\mathbf{w})) \geq 0$, we use the second inequity above and follow the same steps as before to prove the ρ -weak convexity of F with $\rho = G_1 \sqrt{d'} L_2 + \rho_1 G_2^2$. If $\partial f_i(g_i(\mathbf{w})) \leq 0$, we will use the first inequality above to get:

Algorithm 16 SONX

```

1: Input: learning rate schedules  $\{\eta_t\}_{t=1}^T, \{\gamma_t\}_{t=1}^T, \{\beta_t\}_{t=1}^T$ ; starting points  $\mathbf{w}_0, \mathbf{u}_0$ 
2:  $\mathbf{w}_1 = \mathbf{w}_0$ 
3: for  $t = 1, \dots, T$  do
4:   Draw a batch of samples  $\mathcal{B}_t \subset [n]$ 
5:   for  $i \in \mathcal{B}_t$  do
6:     Draw two samples  $\zeta_{i,t} \sim \mathbb{P}_i$ 
7:     Update the inner function value estimators by
        v1:  $\mathbf{u}_{i,t} = (1 - \gamma_t)\mathbf{u}_{i,t-1} + \gamma_t g_i(\mathbf{w}_t; \zeta_{i,t})$ 
        v2:  $\mathbf{u}_{i,t} = (1 - \gamma_t)\mathbf{u}_{i,t-1} + \gamma_t g_i(\mathbf{w}_t; \zeta_{i,t}) + \gamma'_t (g_i(\mathbf{w}_t; \zeta_{i,t}) - g_i(\mathbf{w}_{t-1}; \zeta_{i,t}))$ 
8:   end for
9:   Set  $\mathbf{u}_{i,t} = \mathbf{u}_{i,t-1}, i \notin \mathcal{B}_t$ 
10:  Compute  $\mathbf{z}_t = \frac{1}{B} \sum_{i \in \mathcal{B}'_t} \partial g_i(\mathbf{w}_t; \zeta'_{i,t}) \partial f_i(\mathbf{u}_{i,t})$  ◊ check text for discussion
11:  Update the model  $\mathbf{w}$  by  $\mathbf{w}_{t+1} = \mathbf{w}_t - \eta_t \mathbf{z}_t$ 
12: end for

```

$$\begin{aligned}
f_i(g_i(\mathbf{w}')) &\geq f_i(g_i(\mathbf{w})) + \mathbf{v}_i^\top (g_i(\mathbf{w}') - g_i(\mathbf{w})) - \frac{\rho_1}{2} \|g_i(\mathbf{w}') - g_i(\mathbf{w})\|_2^2 \\
&\geq f_i(g_i(\mathbf{w})) + \mathbf{v}_i^\top (\nabla g_i(\mathbf{w})^\top (\mathbf{w}' - \mathbf{w}) + \frac{L_2}{2} \|\mathbf{w} - \mathbf{w}'\|_2^2 \mathbf{1}) - \frac{\rho_1 G_2^2}{2} \|\mathbf{w} - \mathbf{w}'\|_2^2 \\
&\geq f_i(g_i(\mathbf{w})) + (\nabla g_i(\mathbf{w}) \mathbf{v}_i)^\top (\mathbf{w}' - \mathbf{w}) - \frac{G_1 \sqrt{d'} L_2 + \rho_1 G_2^2}{2} \|\mathbf{w} - \mathbf{w}'\|_2^2.
\end{aligned}$$

This concludes the proof. □

5.3.1 SONX for Non-smooth Inner Functions

Since we do not assume smoothness for the overall objective function, the key difference from the previous two sections is that we no longer have the descent lemma in Lemma 4.9, hence cannot leverage the MA or STORM gradient estimators. Consequently, we employ the vanilla gradient estimator \mathbf{z}_t to update the model parameter \mathbf{w}_{t+1} . The updating steps are summarized in Algorithm 16, referred to as **SONX**. The two options correspond to different strategies for updating the inner function value estimators: v1 uses a coordinate MA estimator, while v2 adopts the MSVR estimator.

For ease of presentation, we compute the vanilla gradient estimator \mathbf{z}_t using a batch \mathcal{B}'_t independent from \mathcal{B}_t :

$$\mathbf{z}_t = \frac{1}{B} \sum_{i \in \mathcal{B}'_t} \partial g_i(\mathbf{w}_t; \zeta'_{i,t}) \partial f_i(\mathbf{u}_{i,t}).$$

However, for SONX-v1 with MA estimator, we can indeed use the same vanilla gradient estimator \mathbf{z}_t as in SOX:

$$\mathbf{z}_t = \frac{1}{B} \sum_{i \in \mathcal{B}_t} \partial g_i(\mathbf{w}_t; \zeta'_{i,t}) \partial f_i(\mathbf{u}_{i,t}).$$

An alternative method for using both options is to compute \mathbf{z}_t by

$$\mathbf{z}_t = \frac{1}{B} \sum_{i \in \mathcal{B}_t} \partial g_i(\mathbf{w}_t; \zeta'_{i,t}) \partial f_i(\mathbf{u}_{i,t-1}).$$

Convergence Analysis

Similar to Section 3.1.4, we state the convergence using the Moreau envelope of F :

$$F_\lambda(\mathbf{w}) := \min_{\mathbf{u}} F(\mathbf{u}) + \frac{1}{2\lambda} \|\mathbf{u} - \mathbf{w}\|_2^2.$$

Recall the definition:

$$\text{prox}_{\lambda F}(\mathbf{w}) = \arg \min_{\mathbf{u}} F(\mathbf{u}) + \frac{1}{2\lambda} \|\mathbf{u} - \mathbf{w}\|_2^2.$$

We first present a result similar to Lemma 3.5 for standard SGD to account for the bias of \mathbf{z}_t .

Lemma 5.10 *Suppose Assumption 5.4 and 5.5 hold. Let $\bar{\rho} = \rho + \rho_2 G_1 + 2\rho_1 G_2^2$. Consider the step update of SONX, we have*

$$\begin{aligned} \mathbb{E}_{\zeta'_t, \mathcal{B}'_t} [F_{1/\bar{\rho}}(\mathbf{w}_{t+1})] &\leq F_{1/\bar{\rho}}(\mathbf{w}_t) + \frac{\eta_t^2 \bar{\rho} G^2}{2} - \frac{\eta_t}{2} \|\nabla F_{1/\bar{\rho}}(\mathbf{w}_t)\|_2^2 \\ &+ \frac{\bar{\rho} \eta_t}{n} \sum_{i=1}^n \left[2G_1 \|g_i(\mathbf{w}_t) - \mathbf{u}_{i,t}\|_2 + \rho_1 \|g_i(\mathbf{w}_t) - \mathbf{u}_{i,t}\|_2^2 \right]. \end{aligned}$$

If f_i is further L_1 -smooth, then

$$\begin{aligned} \mathbb{E}_{\zeta'_t, \mathcal{B}'_t} [F_{1/\bar{\rho}}(\mathbf{w}_{t+1})] &\leq F_{1/\bar{\rho}}(\mathbf{w}_t) + \frac{\eta_t^2 \bar{\rho} G^2}{2} - \frac{\eta_t}{2} \|\nabla F_{1/\bar{\rho}}(\mathbf{w}_t)\|_2^2 \\ &+ \frac{\bar{\rho} \eta_t}{n} \sum_{i=1}^n \left[\frac{L_1}{2} \|g_i(\mathbf{w}_t) - \mathbf{u}_{i,t}\|_2^2 + \rho_1 \|g_i(\mathbf{w}_t) - \mathbf{u}_{i,t}\|_2^2 \right]. \end{aligned}$$

where $G^2 = G_1^2 G_2^2$.

If $\mathbf{u}_{i,t} = g_i(\mathbf{w}_t)$, i.e., there is no bias in \mathbf{z}_t , then the terms in the square bracket are gone, the above lemma reduces to Lemma 3.4.

Proof. Define $\hat{\mathbf{w}}_t := \text{prox}_{F/\bar{\rho}}(\mathbf{w}_t)$ and $\mathbb{E}_t[\cdot] = \mathbb{E}_{\zeta'_t, \mathcal{B}'_t}[\cdot]$. First,

$$\begin{aligned}\mathbb{E}_t[\|\mathbf{z}_t\|_2^2] &\leq \mathbb{E}_t\left[\frac{1}{B}\sum_{i\in\mathcal{B}_t}\|\partial g_i(\mathbf{w}_t, \zeta'_{i,t})\partial f_i(\mathbf{u}_{i,t})\|_2^2\right] \\ &\leq \mathbb{E}_t\left[\frac{1}{B}\sum_{i\in\mathcal{B}_t}\|\partial g_i(\mathbf{w}_t, \zeta'_{i,t})\|_2^2 G_1^2\right] \leq G_2^2 G_1^2 = G^2.\end{aligned}$$

Following Lemma 3.4, we have

$$\mathbb{E}_t[F_{1/\bar{\rho}}(\mathbf{w}_{t+1})] \leq F_{1/\bar{\rho}}(\mathbf{w}_t) + \bar{\rho}\eta_t(\mathbb{E}_t[\mathbf{z}_t])^\top(\hat{\mathbf{w}}_t - \mathbf{w}_t) + \frac{\eta_t^2 \bar{\rho} G^2}{2}. \quad (5.16)$$

Next we bound the term $\mathbb{E}_t[\mathbf{z}_t]^\top(\hat{\mathbf{w}}_t - \mathbf{w}_t)$ on the RHS of (5.16). Note that $\mathbb{E}_t[\mathbf{z}_t] = \frac{1}{n}\sum_{i=1}^n \partial g_i(\mathbf{w}_t)\partial f_i(\mathbf{u}_{i,t})$. For a given $i \in [n]$, we have

$$\begin{aligned}&f_i(g_i(\hat{\mathbf{w}}_t)) - f_i(\mathbf{u}_{i,t}) \\ &\stackrel{(a)}{\geq} \partial f_i(\mathbf{u}_{i,t})^\top(g_i(\hat{\mathbf{w}}_t) - \mathbf{u}_{i,t}) - \frac{\rho_1}{2}\|g_i(\hat{\mathbf{w}}_t) - \mathbf{u}_{i,t}\|_2^2 \\ &\geq \partial f_i(\mathbf{u}_{i,t})^\top(g_i(\hat{\mathbf{w}}_t) - \mathbf{u}_{i,t}) - \rho_1\|g_i(\hat{\mathbf{w}}_t) - g_i(\mathbf{w}_t)\|_2^2 - \rho_1\|g_i(\mathbf{w}_t) - \mathbf{u}_{i,t}\|_2^2 \\ &\geq \partial f_i(\mathbf{u}_{i,t})^\top(g_i(\hat{\mathbf{w}}_t) - \mathbf{u}_{i,t}) - \rho_1 G_2^2\|\hat{\mathbf{w}}_t - \mathbf{w}_t\|_2^2 - \rho_1\|g_i(\mathbf{w}_t) - \mathbf{u}_{i,t}\|_2^2 \\ &\stackrel{(b)}{\geq} \partial f_i(\mathbf{u}_{i,t})^\top\left[g_i(\mathbf{w}_t) - \mathbf{u}_{i,t} + \partial g_i(\mathbf{w}_t)^\top(\hat{\mathbf{w}}_t - \mathbf{w}_t) - \frac{\rho_2}{2}\|\hat{\mathbf{w}}_t - \mathbf{w}_t\|_2^2\right] \\ &\quad - \rho_1 G_2^2\|\hat{\mathbf{w}}_t - \mathbf{w}_t\|_2^2 - \rho_1\|g_i(\mathbf{w}_t) - \mathbf{u}_{i,t}\|_2^2 \\ &\stackrel{(c)}{\geq} \partial f_i(\mathbf{u}_{i,t})^\top(g_i(\mathbf{w}_t) - \mathbf{u}_{i,t}) + \partial f_i(\mathbf{u}_{i,t})^\top \partial g_i(\mathbf{w}_t)^\top(\hat{\mathbf{w}}_t - \mathbf{w}_t) \\ &\quad - \left(\frac{\rho_2 G_1}{2} + \rho_1 G_2^2\right)\|\hat{\mathbf{w}}_t - \mathbf{w}_t\|_2^2 - \rho_1\|g_i(\mathbf{w}_t) - \mathbf{u}_{i,t}\|_2^2,\end{aligned}$$

where (a) follows from the ρ_1 -weak-convexity of f_i , (b) follows from that $\partial f_i(\cdot) \geq 0$ and the weak convexity of g_i , (c) is due to $\|\partial f_i(\mathbf{u}_{i,t})\|_2 \leq G_1$. When $\partial f_i(\cdot) \leq 0$ and g_i is smooth, we can bound similarly with ρ_2 in the last inequality replaced by L_2 .

Then rearranging the above inequality and averaging over i yields

$$\begin{aligned}\mathbb{E}_t[\mathbf{z}_t]^\top(\hat{\mathbf{w}}_t - \mathbf{w}_t) &= \frac{1}{n}\sum_{i=1}^n \partial f_i(\mathbf{u}_{i,t})^\top \partial g_i(\mathbf{w}_t)^\top(\hat{\mathbf{w}}_t - \mathbf{w}_t) \\ &\leq \frac{1}{n}\sum_{i=1}^n \left[f_i(g_i(\hat{\mathbf{w}}_t)) - f_i(g_i(\mathbf{w}_t)) + f_i(g_i(\mathbf{w}_t)) - f_i(\mathbf{u}_{i,t}) \right. \\ &\quad \left. - \partial f_i(\mathbf{u}_{i,t})^\top(g_i(\mathbf{w}_t) - \mathbf{u}_{i,t}) + \left(\frac{\rho_2 G_1}{2} + \rho_1 G_2^2\right)\|\hat{\mathbf{w}}_t - \mathbf{w}_t\|_2^2 + \rho_1\|g_i(\mathbf{w}_t) - \mathbf{u}_{i,t}\|_2^2 \right].\end{aligned} \quad (5.17)$$

Due to the ρ -weak convexity of $F(\mathbf{w})$, we have that $F(\mathbf{w}) + \frac{\bar{\rho}}{2}\|\mathbf{w}_t - \mathbf{w}\|_2^2$ is $(\bar{\rho} - \rho)$ -strongly convex. Then $\left[F(\mathbf{w}_t) + \frac{\bar{\rho}}{2}\|\mathbf{w}_t - \mathbf{w}_t\|_2^2\right] - \left[F(\hat{\mathbf{w}}_t) + \frac{\bar{\rho}}{2}\|\mathbf{w}_t - \hat{\mathbf{w}}_t\|_2^2\right] \geq \frac{\bar{\rho} - \rho}{2}\|\hat{\mathbf{w}}_t - \mathbf{w}_t\|_2^2$. It follows that:

$$\begin{aligned}
 & \frac{1}{n} \sum_{i=1}^n \left[f_i(g_i(\hat{\mathbf{w}}_t)) - f_i(g_i(\mathbf{w}_t)) \right] = F(\hat{\mathbf{w}}_t) - F(\mathbf{w}_t) \\
 & = \left[F(\hat{\mathbf{w}}_t) + \frac{\bar{\rho}}{2} \|\mathbf{w}_t - \hat{\mathbf{w}}_t\|_2^2 \right] - \left[F(\mathbf{w}_t) + \frac{\bar{\rho}}{2} \|\mathbf{w}_t - \mathbf{w}_t\|_2^2 \right] - \frac{\bar{\rho}}{2} \|\mathbf{w}_t - \hat{\mathbf{w}}_t\|_2^2 \quad (5.18) \\
 & \leq \left(\frac{\rho}{2} - \bar{\rho} \right) \|\mathbf{w}_t - \hat{\mathbf{w}}_t\|_2^2
 \end{aligned}$$

Combining inequality (5.17), (5.16) and (5.18) yields

$$\begin{aligned}
 \mathbb{E}_t[F_{1/\bar{\rho}}(\mathbf{w}_{t+1})] & \leq F_{1/\bar{\rho}}(\mathbf{w}_t) + \frac{\eta_t^2 \bar{\rho} G^2}{2} - \frac{\bar{\rho}^2 \eta_t}{2} \|\mathbf{w}_t - \hat{\mathbf{w}}_t\|_2^2 \\
 & + \frac{\bar{\rho} \eta_t}{n} \sum_{i=1}^n \left[f_i(g_i(\mathbf{w}_t)) - f_i(\mathbf{u}_{i,t}) - \partial f_i(\mathbf{u}_{i,t})^\top (g_i(\mathbf{w}_t) - \mathbf{u}_{i,t}) \right. \\
 & \quad \left. + \rho_1 \|g_i(\mathbf{w}_t) - \mathbf{u}_{i,t}\|_2^2 \right].
 \end{aligned}$$

We finish the proof by noting that $\|\nabla F_{1/\bar{\rho}}(\mathbf{w}_t)\|_2 = \bar{\rho} \|\mathbf{w}_t - \hat{\mathbf{w}}_t\|_2$, using

$$f_i(g_i(\mathbf{w}_t)) - f_i(\mathbf{u}_{i,t}) - \partial f_i(\mathbf{u}_{i,t})^\top (g_i(\mathbf{w}_t) - \mathbf{u}_{i,t}) \leq 2G_1 \|g_i(\mathbf{w}_t) - \mathbf{u}_{i,t}\|_2,$$

if f_i is G_1 -Lipschitz continuous, or using

$$f_i(g_i(\mathbf{w}_t)) - f_i(\mathbf{u}_{i,t}) - \partial f_i(\mathbf{u}_{i,t})^\top (g_i(\mathbf{w}_t) - \mathbf{u}_{i,t}) \leq \frac{L_1}{2} \|g_i(\mathbf{w}_t) - \mathbf{u}_{i,t}\|_2^2,$$

if f_i is L_1 -smooth. \square

Convergence of SONX-v1

Recall the definition:

$$\delta_t = \frac{1}{n} \sum_{i=1}^n \|\mathbf{u}_{i,t} - g_i(\mathbf{w}_t)\|_2^2.$$

Let us also define:

$$\delta'_t = \frac{1}{n} \sum_{i=1}^n \|\mathbf{u}_{i,t} - g_i(\mathbf{w}_t)\|_2.$$

From Lemma 5.10, the key is to bound δ_t and δ'_t .

Lemma 5.11 *Consider the update of SONX-v1, under Assumptions 5.4 and 5.5, with constant parameters $\gamma_t = \gamma \leq 1$ and $\eta_t = \eta$, we have*

$$\begin{aligned}\mathbb{E} [\delta_t] &\leq \left(1 - \frac{B\gamma}{4n}\right)^{2t} \delta_0 + \frac{8n^2 G_2^4 G_1^2 \eta^2}{B^2 \gamma^2} + 4\gamma \sigma_0^2. \\ \mathbb{E} [\delta'_t] &\leq \left(1 - \frac{B\gamma}{4n}\right)^t \delta'_0 + \frac{4n G_2^2 G_1 \eta}{B\gamma} + 2\sqrt{\gamma} \sigma_0.\end{aligned}$$

Proof. From the proof of Lemma 5.1, we have

$$\begin{aligned}&\mathbb{E} \left[\left\| \mathbf{u}_{i,t} - g_i(\mathbf{w}_t) \right\|_2^2 \right] \\ &\leq \left(1 - \frac{B\gamma_t}{2n}\right) \mathbb{E} \left[\left\| \mathbf{u}_{i,t-1} - g_i(\mathbf{w}_{t-1}) \right\|_2^2 \right] + \frac{2n G_2^2}{B\gamma_t} \mathbb{E} \left[\left\| \mathbf{w}_{t-1} - \mathbf{w}_t \right\|_2^2 \right] + \frac{B\gamma_t^2 \sigma_0^2}{n} \\ &\leq \left(1 - \frac{B\gamma_t}{2n}\right) \mathbb{E} \left[\left\| \mathbf{u}_{i,t-1} - g_i(\mathbf{w}_{t-1}) \right\|_2^2 \right] + \frac{2n G_2^2 \eta_{t-1}^2}{B\gamma_t} \mathbb{E} \left[\left\| \mathbf{z}_{t-1} \right\|_2^2 \right] + \frac{B\gamma_t^2 \sigma_0^2}{n} \\ &\leq \left(1 - \frac{B\gamma_t}{4n}\right)^2 \mathbb{E} \left[\left\| \mathbf{u}_{i,t-1} - g_i(\mathbf{w}_{t-1}) \right\|_2^2 \right] + \frac{2n G_2^4 G_1^2 \eta_{t-1}^2}{B\gamma_t} + \frac{B\gamma_t^2 \sigma_0^2}{n}.\end{aligned}$$

Applying the above inequality recursively for $\gamma_t = \gamma$ and $\eta_t = \eta$, we obtain

$$\begin{aligned}&\mathbb{E} \left[\left\| \mathbf{u}_{i,t} - g_i(\mathbf{w}_t) \right\|_2^2 \right] \\ &\leq \left(1 - \frac{B\gamma}{4n}\right)^{2t} \left\| \mathbf{u}_{i,0} - g_i(\mathbf{w}_0) \right\|_2^2 + \sum_{j=0}^{t-1} \left(1 - \frac{B\gamma}{4n}\right)^{2j} \left(\frac{2n G_2^4 G_1^2 \eta^2}{B\gamma} + \frac{B\gamma^2 \sigma_0^2}{n} \right) \\ &\leq \left(1 - \frac{B\gamma}{4n}\right)^{2t} \left\| \mathbf{u}_{i,0} - g_i(\mathbf{w}_0) \right\|_2^2 + \frac{8n^2 G_2^4 G_1^2 \eta^2}{B^2 \gamma^2} + 4\gamma \sigma_0^2,\end{aligned}$$

where we use

$$\sum_{j=0}^{t-1} (1 - \alpha)^{2j} \leq \sum_{j=0}^{\infty} (1 - \alpha)^{2j} = \frac{1}{1 - (1 - \alpha)^2} = \frac{1}{\alpha(2 - \alpha)} \leq \frac{1}{\alpha}, \forall \alpha \in (0, 1).$$

Averaging the above inequality over i , we prove the first result in the lemma.

It follows

$$\begin{aligned}\mathbb{E} \left[\left\| \mathbf{u}_{i,t} - g_i(\mathbf{w}_t) \right\|_2 \right] &\leq \sqrt{\mathbb{E} \left[\left\| \mathbf{u}_{i,t} - g_i(\mathbf{w}_t) \right\|_2^2 \right]} \\ &\leq \sqrt{\left(1 - \frac{B\gamma}{4n}\right)^{2t} \left\| \mathbf{u}_{i,0} - g_i(\mathbf{w}_0) \right\|_2^2 + \frac{8n^2 G_2^4 G_1^2 \eta^2}{B^2 \gamma^2} + 4\gamma \sigma_0^2} \\ &\leq \left(1 - \frac{B\gamma}{4n}\right)^t \left\| \mathbf{u}_{i,0} - g_i(\mathbf{w}_0) \right\|_2 + \frac{4n G_2^2 G_1 \eta}{B\gamma} + 2\gamma^{1/2} \sigma_0.\end{aligned}$$

Averaging the above result, we prove the second result. \square

Theorem 5.3 (Convergence of SONX-v1 with Lipschitz f_i) Consider [SONX-v1](#), and suppose Assumption [5.4](#) and [5.5](#) hold and f_i is G_1 -Lipschitz continuous. Let $\eta_t = \eta = O(\frac{B\epsilon^6}{n\sigma_0^2})$, $\gamma_t = \gamma = O(\frac{\epsilon^4}{\sigma_0^2})$. Then after $T = O(\frac{n\sigma_0^2}{B\epsilon^8})$ iterations, we have $\mathbb{E}[\|\nabla F_{1/\bar{\rho}}(\mathbf{w}_t)\|_2^2] \leq O(\epsilon^2)$.

Proof. From Lemma [5.10](#), we have

$$\begin{aligned} \mathbb{E}\left[\frac{1}{T} \sum_{t=1}^T \|\nabla F_{1/\bar{\rho}}(\mathbf{w}_t)\|_2^2\right] &\leq \mathbb{E}\left[\frac{2 \sum_{t=1}^T (F_{1/\bar{\rho}}(\mathbf{w}_t) - F_{1/\bar{\rho}}(\mathbf{w}_{t+1}))}{\eta T}\right] + \eta \bar{\rho} G^2 \\ &\quad + 4\bar{\rho} G_1 \mathbb{E}\left[\frac{1}{T} \sum_{t=1}^T \delta'_t\right] + 2\bar{\rho} \rho_1 \mathbb{E}\left[\frac{1}{T} \sum_{t=1}^T \delta_t\right]. \end{aligned}$$

Next, we bound the last two terms. From Lemma [5.11](#), we have

$$\begin{aligned} \mathbb{E}\left[\frac{1}{T} \sum_{t=1}^T \delta_t\right] &\leq \frac{1}{T} \sum_{t=1}^T \left(1 - \frac{B\gamma}{4n}\right)^{2t} \delta_0 + \frac{8n^2 G_2^4 G_1^2 \eta^2}{B^2 \gamma^2} + 4\gamma \sigma_0^2. \\ \mathbb{E}\left[\frac{1}{T} \sum_{t=1}^T \delta'_t\right] &\leq \frac{1}{T} \sum_{t=1}^T \left(1 - \frac{B\gamma}{4n}\right)^t \delta'_0 + \frac{4n G_2^2 G_1 \eta}{B\gamma} + 2\sqrt{\gamma} \sigma_0. \end{aligned}$$

Since $\sum_{t=1}^T (1 - \mu)^t \leq \frac{1}{\mu}$ for $\mu \in (0, 1)$, we have

$$\begin{aligned} \mathbb{E}\left[\frac{1}{T} \sum_{t=1}^T \delta_t\right] &\leq \frac{4n\delta_0}{B\gamma T} + \frac{8n^2 G_2^4 G_1^2 \eta^2}{B^2 \gamma^2} + 4\gamma \sigma_0^2. \\ \mathbb{E}\left[\frac{1}{T} \sum_{t=1}^T \delta'_t\right] &\leq \frac{4n\delta'_0}{B\gamma T} + \frac{4n G_2^2 G_1 \eta}{B\gamma} + 2\sqrt{\gamma} \sigma_0. \end{aligned}$$

From Proposition [3.2](#), we have

$$\sum_{t=1}^T (F_{1/\bar{\rho}}(\mathbf{w}_t) - F_{1/\bar{\rho}}(\mathbf{w}_{t+1})) = F_{1/\bar{\rho}}(\mathbf{w}_1) - F_{1/\bar{\rho}}(\mathbf{w}_{T+1}) \leq F(\mathbf{w}_1) - F(\mathbf{w}_*).$$

Combining the above results, we have

$$\begin{aligned} \mathbb{E}\left[\frac{1}{T} \sum_{t=1}^T \|\nabla F_{1/\bar{\rho}}(\mathbf{w}_t)\|_2^2\right] &\leq \mathbb{E}\left[\frac{2(F(\mathbf{w}_1) - F_*)}{\eta T}\right] + \eta \bar{\rho} G^2 \\ &\quad + 4\bar{\rho} G_1 \left(\frac{4n\delta'_0}{B\gamma T} + \frac{4n G_2^2 G_1 \eta}{B\gamma} + 2\sqrt{\gamma} \sigma_0\right) + 2\bar{\rho} \rho_1 \left(\frac{4n\delta_0}{B\gamma T} + \frac{8n^2 G_2^4 G_1^2 \eta^2}{B^2 \gamma^2} + 4\gamma \sigma_0^2\right). \end{aligned}$$

Plugging the order of η, γ , we finish the proof. \square

Theorem 5.4 (Convergence of SONX-v1 with smooth f_i) Consider [SONX-v1](#), and suppose Assumption [5.1](#) and [5.5](#) hold and f_i is L_1 -smooth. Let $\eta_t = \eta = O(\frac{B\epsilon^3}{n\sigma_0^2})$, $\gamma_t = \gamma = O(\frac{\epsilon^2}{\sigma_0^2})$, then after $T = O(\frac{n\sigma_0^2}{B\epsilon^3})$ iterations, we have $\mathbb{E}[\|\nabla F_{1/\bar{\rho}}(\mathbf{w}_t)\|_2^2] \leq O(\epsilon^2)$.

Proof. By using the result for smooth f_i in Lemma [5.10](#), we have

$$\begin{aligned} \mathbb{E}\left[\frac{1}{T} \sum_{t=1}^T \|\nabla F_{1/\bar{\rho}}(\mathbf{w}_t)\|_2^2\right] &\leq \mathbb{E}\left[\frac{2 \sum_{t=1}^T (F_{1/\bar{\rho}}(\mathbf{w}_t) - F_{1/\bar{\rho}}(\mathbf{w}_{t+1}))}{\eta T}\right] + \eta \bar{\rho} G^2 \\ &\quad + \bar{\rho}(L_1 + 2\rho_1) \mathbb{E}\left[\frac{1}{T} \sum_{t=1}^T \delta_t\right]. \end{aligned}$$

Plugging the bounds for the first and last term in the RHS, we have

$$\begin{aligned} \mathbb{E}\left[\frac{1}{T} \sum_{t=1}^T \|\nabla F_{1/\bar{\rho}}(\mathbf{w}_t)\|_2^2\right] &\leq \mathbb{E}\left[\frac{2(F(\mathbf{w}_1) - F_*)}{\eta T}\right] + \eta \bar{\rho} G^2 \\ &\quad + \bar{\rho}(L_1 + 2\rho_1) \left(\frac{4n\delta_0}{B\gamma T} + \frac{8n^2 G_2^4 G_1^2 \eta^2}{B^2 \gamma^2} + 4\gamma \sigma_0^2\right). \end{aligned}$$

Plugging the order of η, γ , we finish the proof. \square

Convergence of SONX-v2

Similar to the first option, we need to bound δ_t, δ'_t first.

Lemma 5.12 Under Assumption [5.4](#), [5.5](#), by setting $\gamma_t = \gamma \leq \frac{1}{2}$, $\eta_t = \eta$, $\gamma'_t = \frac{n-B}{B(1-\gamma)} + (1-\gamma)$, we have:

$$\begin{aligned} \mathbb{E}[\delta_t] &\leq \left(1 - \frac{B\gamma}{2n}\right)^{2t} \delta_0 + 4\gamma \sigma_0^2 + \frac{24n^2 G_2^4 G_1^2 \eta^2}{B^2 \gamma}, \\ \mathbb{E}[\delta'_t] &\leq \left(1 - \frac{B\gamma}{2n}\right)^t \delta'_0 + 2\gamma^{1/2} \sigma_0 + \frac{5n G_2^2 G_1 \eta}{B\gamma^{1/2}}. \end{aligned}$$

Proof is omitted as it is similar to that of Lemma [5.11](#) but based on Lemma [5.5](#).

Theorem 5.5 (Convergence of SONX-v2) Consider [SONX-v2](#), and suppose Assumption [5.4](#), and [5.5](#) hold.

- If f_i is G_1 -Lipschitz continuous, by setting $\eta = O(\frac{B\epsilon^4}{n\sigma_0^4})$, $\gamma = O(\frac{\epsilon^4}{\sigma_0^2})$, then after $T = O(\frac{n\sigma_0}{B\epsilon^6})$ iterations, we have $\mathbb{E}[\|\nabla F_{1/\bar{\rho}}(\mathbf{x}_t)\|_2^2] \leq \epsilon^2$.

- If f_i is further L_1 -smooth, by setting $\eta = O(\frac{B\epsilon^2}{n\sigma_0})$, $\gamma = O(\frac{\epsilon^2}{\sigma_0^2})$, then the complexity reduces to $T = O(\frac{n\sigma_0}{B\epsilon^4})$.

The proof follows similarly to that of Theorem 5.3 and Theorem 5.4 and is left as an exercise for interested readers.

5.3.2 SONEX for Non-smooth Outer functions

When f_i is Lipschitz continuous and non-smooth, the best complexity derived in last subsection is $O(n/(B\epsilon^6))$. Can we further improve the complexity when the inner functions are smooth? We present a method and its analysis in this section.

Let us make the following assumptions.

Assumption 5.6. We assume that

- (i) $\mathbb{E}_{\zeta \sim \mathbb{P}_i} [\|g_i(\mathbf{w}; \zeta) - g_i(\mathbf{w})\|_2^2] \leq \sigma_0^2$
- (ii) $\mathbb{E}_{\zeta \sim \mathbb{P}_i} [\|\nabla g_i(\mathbf{w}; \zeta) - \nabla g_i(\mathbf{w})\|_2^2] \leq \sigma_2^2$
- (iii) $\mathbb{E}_{\zeta \sim \mathbb{P}_i} [\|\nabla g_i(\mathbf{w}; \zeta)\|_2^2] \leq G_2^2$.

Assumption 5.7. The following conditions hold:

- (i) f_i is ρ_1 -weakly convex, G_1 -Lipschitz continuous,
- (ii) g_i is L_2 -smooth and G_2 -Lipschitz continuous.

Moreau Envelope Smoothing of the outer function

A classical approach of improving the convergence for non-smooth functions in convex optimization is smoothing, i.e., first smoothing the function and then using an optimizer for solving the resulting smoothed function. We define the Moreau envelope smoothing of f_i as follows:

$$\bar{f}_i(g) = \min_{\mathbf{u} \in \mathbb{R}^{d'}} \frac{\bar{\rho}_1}{2} \|\mathbf{u} - g\|_2^2 + f_i(\mathbf{u}), \quad (5.19)$$

where $\bar{\rho}_1 > \rho_1$. We present a lemma below regarding \bar{f}_i .

Lemma 5.13 If f_i is G_1 -Lipschitz continuous and ρ_1 -weakly convex, then \bar{f}_i is \bar{L}_1 -smooth and G_1 Lipschitz continuous, where $\bar{L}_1 = \frac{\bar{\rho}_1(2\bar{\rho}_1 - \rho_1)}{(\bar{\rho}_1 - \rho_1)}$.

Proof. Define $\text{prox}_{f_i/\bar{\rho}_1}(g) = \arg \min_{\mathbf{u} \in \mathbb{R}^{d'}} \frac{\bar{\rho}_1}{2} \|\mathbf{u} - g\|_2^2 + f_i(\mathbf{u})$. We have

$$\nabla \bar{f}_i(g) = \bar{\rho}_1(g - \text{prox}_{f_i/\bar{\rho}_1}(g)).$$

Due to the optimality condition of $\text{prox}_{f_i/\bar{\rho}_1}(g)$, we have

$$\bar{\rho}_1(g - \text{prox}_{f_i/\bar{\rho}_1}(g)) \in \partial f_i(\text{prox}_{f_i/\bar{\rho}_1}(g)).$$

Hence, $\nabla \bar{f}_i(g) \in \partial f_i(\text{prox}_{f_i/\bar{\rho}_1}(g))$, which implies $\|\nabla \bar{f}_i(g)\| \leq G_1$. The smoothness of \bar{f}_i follows from Proposition 3.1. \square

Relationship with Nesterov Smoothing

When f_i is a convex function, its Moreau envelope smoothing is also equivalent to the well-known Nesterov smoothing. To see this, let f_i^* denote the convex conjugate of f_i , i.e., $f_i^*(\mathbf{u}) = \max_{g \in \mathbb{R}^{d'}} \mathbf{u}^\top g - f_i(g)$. Since f_i is convex, we have $f_i(g) = \max_{\mathbf{u} \in \mathcal{U}} \mathbf{u}^\top g - f_i^*(\mathbf{u})$, where $\mathcal{U} = \text{dom}(f_i^*)$ is bounded as $\|\partial f_i(g)\| \leq G_1$. As a result,

$$\bar{f}_i(g) = \min_{\mathbf{u} \in \mathbb{R}^{d'}} \frac{\bar{\rho}_1}{2} \|\mathbf{u} - g\|_2^2 + f_i(\mathbf{u}) = \min_{\mathbf{u} \in \mathbb{R}^{d'}} \frac{\bar{\rho}_1}{2} \|\mathbf{u} - g\|_2^2 + \max_{\mathbf{u}' \in \mathcal{U}} \mathbf{u}'^\top \mathbf{u} - f_i^*(\mathbf{u}').$$

By Sion's minimax theorem, we can switch the min and max. Hence,

$$\bar{f}_i(g) = \max_{\mathbf{u}' \in \mathcal{U}} \min_{\mathbf{u} \in \mathbb{R}^{d'}} \frac{\bar{\rho}_1}{2} \|\mathbf{u} - g\|_2^2 + \mathbf{u}'^\top \mathbf{u} - f_i^*(\mathbf{u}').$$

By solving the minimization over \mathbf{u} and plugging the optimal solution into the expression, we get

$$\bar{f}_i(g) = \max_{\mathbf{u}' \in \mathcal{U}} g^\top \mathbf{u}' - f_i^*(\mathbf{u}') - \frac{1}{2\bar{\rho}_1} \|\mathbf{u}'\|_2^2. \quad (5.20)$$

This is known as Nesterov smoothing of the function $f_i(g)$. When $\bar{\rho}_1$ is sufficiently large, we can prove that \bar{f}_i is sufficiently close to f_i .

Example

Example 5.1. Let us consider the Nesterov smoothing of the hinge function $f(x) = [x]_+$. Let $\bar{\rho}_1 = 1/\varepsilon$ for some small $\varepsilon \ll 1$. Then, the Nesterov smoothing of the hinge function is

$$\bar{f}(x) = \max_{u \in [0,1]} ux - \frac{\varepsilon}{2} u^2 = \begin{cases} x - \frac{\varepsilon}{2} & \text{if } x \geq \varepsilon \\ \frac{x^2}{2\varepsilon} & \text{if } 0 < x < \varepsilon \\ 0 & \text{o.w.} \end{cases}.$$

This is also known as the smoothed hinge function.

Solving the smoothed problem

With a smoothed outer function \tilde{f}_i , we consider optimizing the following problem with some proper value of $\bar{\rho}_1$:

$$\min_{\mathbf{w}} \bar{F}(\mathbf{w}) := \frac{1}{n} \sum_{i=1}^n \tilde{f}_i(g_i(\mathbf{w})). \quad (5.21)$$

Following Lemma 4.3, $\bar{F}(\cdot)$ is smooth with a smoothness parameter $\bar{L}_F = G_1 L_2 + G_2^2 \bar{L}_1$.

The key concern is how the convergence of solving the above problem translates to the convergence of solving the original problem (5.1). To address this question, we introduce a new convergence measure, named approximate ϵ -stationarity.

Definition 5.1 (Approximate ϵ -stationary solution) A point \mathbf{w} is an approximate ϵ -stationary solution to the original problem (5.1), if there exists $(\mathbf{u}_1, \dots, \mathbf{u}_n)$ and $\lambda_i \in \partial f(\mathbf{u}_i)$, $\forall i$ such that

$$\left\| \frac{1}{n} \sum_{i=1}^n \nabla g_i(\mathbf{w}) \lambda_i \right\|_2 \leq \epsilon, \quad (5.22)$$

$$\|\mathbf{u}_i - g_i(\mathbf{w})\|_2 \leq O(\epsilon), \forall i. \quad (5.23)$$

We note that this condition is closely related to the KKT condition of the following equivalent constrained formulation of the original problem (5.1):

$$\min_{\mathbf{w}, \mathbf{u}} \quad \frac{1}{n} \sum_{i=1}^n f_i(\mathbf{u}_i) \quad (5.24)$$

$$\text{s.t.} \quad g_i(\mathbf{w}) = \mathbf{u}_i, \forall i. \quad (5.25)$$

The Lagrangian function of this constrained formulation is given by

$$F(\mathbf{w}, \mathbf{u}, \lambda) = \frac{1}{n} \sum_{i=1}^n f_i(\mathbf{u}_i) + \sum_{i=1}^n \lambda_i^\top (g_i(\mathbf{w}) - \mathbf{u}_i).$$

A solution $(\mathbf{w}, \mathbf{u}, \lambda)$ satisfies the KKT condition, if

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n \nabla g_i(\mathbf{w}) \lambda_i &= 0, \quad \lambda_i \in \partial f_i(\mathbf{u}_i) \\ \mathbf{u}_i &= g_i(\mathbf{w}). \end{aligned}$$

Hence, an approximate ϵ -stationary solution satisfies the KKT condition approximately when $\epsilon \ll 1$.

If f_i is L_1 -smooth, an approximate ϵ -stationary solution is also a standard $O(\epsilon)$ -stationary solution. To see this, we have

Algorithm 17 SONE

```

1: Input: learning rate schedules  $\{\eta_t\}_{t=1}^T, \{\gamma_t\}_{t=1}^T, \{\beta_t\}_{t=1}^T$ ; starting points  $\mathbf{w}_0, \mathbf{v}_0, \mathbf{u}_0$ 
2:  $\mathbf{w}_1 = \mathbf{w}_0 - \eta_0 \mathbf{v}_0$ 
3: for  $t = 1, \dots, T$  do
4:   Draw a batch of samples  $\mathcal{B}_t \subset [n]$ 
5:   for  $i \in \mathcal{B}_t$  do
6:     Draw two samples  $\zeta_{i,t} \sim \mathbb{P}_i$ 
7:     Update the inner function value estimators by
        v1:  $\mathbf{u}_{i,t} = (1 - \gamma_t) \mathbf{u}_{i,t-1} + \gamma_t g_i(\mathbf{w}_t; \zeta_{i,t})$ 
        v2:  $\mathbf{u}_{i,t} = (1 - \gamma_t) \mathbf{u}_{i,t-1} + \gamma_t g_i(\mathbf{w}_t; \zeta_{i,t}) + \gamma'_t (g_i(\mathbf{w}_t; \zeta_{i,t}) - g_i(\mathbf{w}_{t-1}; \zeta_{i,t}))$ 
8:   end for
9:   Set  $\mathbf{u}_{i,t} = \mathbf{u}_{i,t-1}, i \notin \mathcal{B}_t$ 
10:  Compute the vanilla gradient estimator  $\mathbf{z}_t = \frac{1}{B} \sum_{i \in \mathcal{B}_t} \nabla g_i(\mathbf{w}_t; \zeta'_{i,t}) \nabla \bar{f}_i(\mathbf{u}_{i,t})$ 
11:  Update the MA gradient estimator  $\mathbf{v}_t = (1 - \beta_t) \mathbf{v}_{t-1} + \beta_t \mathbf{z}_t$ 
12:  Update the model  $\mathbf{w}$  by  $\mathbf{w}_{t+1} = \mathbf{w}_t - \eta_t \mathbf{v}_t$ 
13: end for

```

$$\begin{aligned}
& \left\| \frac{1}{n} \sum_{i=1}^n \nabla g_i(\mathbf{w}) \nabla f_i(g_i(\mathbf{w})) \right\|_2 \\
&= \left\| \frac{1}{n} \sum_{i=1}^n \nabla g_i(\mathbf{w}) \nabla f_i(g_i(\mathbf{w})) - \frac{1}{n} \sum_{i=1}^n \nabla g_i(\mathbf{w}) \nabla f_i(\mathbf{u}_i) + \frac{1}{n} \sum_{i=1}^n \nabla g_i(\mathbf{w}) \nabla f_i(\mathbf{u}_i) \right\|_2 \\
&\leq \frac{1}{n} \sum_{i=1}^n G_2 L_1 \|\mathbf{u}_i - g_i(\mathbf{w})\|_2 + \epsilon \leq O(\epsilon).
\end{aligned}$$

The following proposition states that an ϵ -stationary solution to the smoothed problem (5.21) is an approximate ϵ -stationary solution to the original problem when $\bar{\rho}_1$ is sufficiently large.

Proposition 5.1 *Let \mathbf{w} be an ϵ -stationary solution to (5.21), when $\bar{\rho}_1 = 1/\epsilon$, then \mathbf{w} is also an approximate ϵ -stationary solution to (5.1).*

Proof. Given that \mathbf{w} be an ϵ -stationary solution to (5.21), we have

$$\left\| \frac{1}{n} \sum_{i=1}^n \nabla g_i(\mathbf{w}) \nabla \bar{f}_i(g_i(\mathbf{w})) \right\|_2 \leq \epsilon.$$

We define $\mathbf{u}_i = \text{prox}_{f_i/\bar{\rho}_1}(g_i(\mathbf{w})) = \arg \min_{\mathbf{u}} f_i(\mathbf{u}) + \frac{\bar{\rho}_1}{2} \|\mathbf{u} - g_i(\mathbf{w})\|_2^2$ and $\lambda_i = \nabla \bar{f}_i(g_i(\mathbf{w}))$. Since $\nabla \bar{f}_i(g_i(\mathbf{w})) \in \partial f_i(\text{prox}_{f_i/\bar{\rho}_1}(g_i(\mathbf{w}))) = \partial f_i(\mathbf{u}_i)$. As a result, we have $\lambda_i \in \partial f_i(\mathbf{u}_i)$ and $\left\| \frac{1}{n} \sum_{i=1}^n \nabla g_i(\mathbf{w}) \lambda_i \right\|_2 \leq \epsilon$.

Due to the optimality condition of \mathbf{u}_i , we have $g_i(\mathbf{w}) - \mathbf{u}_i \in \partial f_i(\mathbf{u}_i)/\bar{\rho}_1$. Since f_i is G_1 -Lipschitz continuous and $\bar{\rho}_1 \geq 1/\epsilon$, hence, $\|\mathbf{u}_i - g_i(\mathbf{w})\|_2 \leq O(\epsilon)$. \square

Next, we discuss algorithms and complexities for solving the smoothed problem when $\bar{\rho}_1 = 1/\epsilon$. Since both inner and outer functions of the smoothed problem are smooth, we can leverage the moving average gradient estimators. We present detailed steps for solving the smoothed problem in Algorithm 17, which is referred to as SONEX.

A step in implementing SONEX for solving the smoothed problem (5.21) is the calculation of $\nabla \tilde{f}_i(\mathbf{u}_{i,t})$, which amounts to solving a proximal mapping of f_i , i.e.,

$$\text{prox}_{f_i/\bar{\rho}_1}(\mathbf{u}_{i,t}) = \arg \min_{\mathbf{u} \in \mathbb{R}^{d'}} \frac{\bar{\rho}_1}{2} \|\mathbf{u} - \mathbf{u}_{i,t}\|_2^2 + f_i(\mathbf{u}).$$

In fact, $\nabla \tilde{f}_i(\mathbf{u}_{i,t}) = \bar{\rho}_1(\mathbf{u}_{i,t} - \text{prox}_{f_i/\bar{\rho}_1}(\mathbf{u}_{i,t}))$.

Convergence of SONEX-v1

Finally, we present the complexity of SONEX-v1 for finding an ϵ -stationary solution to the smoothed problem when $\bar{\rho}_1 = 1/\epsilon$.

Corollary 5.1 (Convergence of SONEX-v1) *Under Assumptions 5.6 and 5.7, if we set \mathbf{u}_0 such that $\frac{1}{n}\mathbb{E}[\sum_{i=1}^n \|\mathbf{u}_{i,0} - g_i(\mathbf{w}_0)\|_2^2] \leq O(\epsilon)$, $\beta = O(\frac{\epsilon^2}{\sigma^2})$, $\gamma = O(\frac{\epsilon^4}{\sigma_0^2})$, $\eta = \min(\epsilon, O(\beta\epsilon), O(\frac{B\epsilon\gamma}{n}))$, $\bar{\rho}_1 = 1/\epsilon > \rho_1$, then SONEX-v1 finds an approximate $O(\epsilon)$ -stationary solution to the original problem (5.1) with a complexity of $O(\frac{n\sigma_0^2}{B\epsilon^7})$.*

Proof. The proof can be completed by using the convergence result of SOX with noting the order of $\bar{L}_1 = O(\bar{\rho}_1) = O(1/\epsilon)$ and $L_F = O(\bar{L}_1) = O(1/\epsilon)$. \square

Convergence of SONEX-v2

SONEX-v2 is a combination of SOX and MSVR, i.e., with \mathbf{u}_t sequence from MSVR and \mathbf{v}_t from SOX.

Theorem 5.6 (Convergence of SONEX-v2) *Under Assumptions 5.6 and 5.7, if we set \mathbf{u}_1 such that $\frac{1}{n}\mathbb{E}[\sum_{i=1}^n \|\mathbf{u}_{i,0} - g_i(\mathbf{w}_0)\|_2^2] \leq O(\epsilon^3/\sigma_0)$, $\beta = O(\frac{\epsilon^2}{\sigma^2})$, $\gamma = O(\frac{\epsilon^2}{\sigma_0^2})$, $\eta = \min(O(\epsilon), O(\beta\epsilon), O(\frac{B\sqrt{\gamma}\epsilon}{n}))$ and $\bar{\rho}_1 = \frac{1}{\epsilon} > \rho_1$, then SONEX-v2 finds an approximate ϵ -stationary solution to the original problem (5.1) with a complexity of*

$$T = O\left(\max\left\{\frac{1}{\epsilon^3}, \frac{\sigma^2}{\epsilon^5}, \frac{n\sigma_0}{B\epsilon^5}\right\}\right),$$

$$\text{where } \sigma^2 = \frac{G_1^2\sigma_2^2}{B} + \frac{G_1^2G_2^2(n-B)}{B(n-1)}.$$

Proof. The proof is similar to that of Theorem 4.3 except that the \diamond inequality in Lemma 4.10 is replaced by the following for using MSVR estimators (see Lemma 5.5):

$$(\diamond) \quad \mathbb{E}[\delta_{t+1}] \leq \mathbb{E}[(1 - \bar{\gamma})\delta_t + C_3\eta^2\Gamma_t + \bar{\gamma}^2\sigma'^2],$$

where $\bar{\gamma} = \frac{B\gamma}{n}$, $\sigma'^2 = \frac{2n\sigma_0^2}{B}$, $C_3 = O(n/B)$

We only highlight the changes below and leave details as an exercise. First, the condition on η in Lemma 4.10 is changed to

$$\eta \leq O\left(\frac{1}{L}, \frac{\beta}{\sqrt{C_2}}, \sqrt{\frac{\bar{\gamma}}{C_1 C_3}}\right).$$

The settings on $\beta, \bar{\gamma}$ remain the same as $\beta = O(\frac{\epsilon^2}{\sigma^2})$, $\bar{\gamma} = O(\frac{\epsilon^2}{C_1 \sigma'^2})$. The iteration complexity becomes:

$$\begin{aligned} T &= O\left(\max\left\{\frac{C_Y L_F}{\epsilon^2}, \frac{C_Y \sqrt{C_2}}{\epsilon^2 \beta}, \frac{C_Y \sqrt{C_1 C_3}}{\sqrt{\bar{\gamma}} \epsilon^2}\right\}\right) \\ &= O\left(\max\left\{\frac{C_Y L_F}{\epsilon^2}, \frac{C_Y \sigma^2 \sqrt{C_2}}{\epsilon^4}, \frac{C_Y \sqrt{C_3} C_1 \sigma'}{\epsilon^3}\right\}\right). \end{aligned}$$

and C_Y is changed to

$$\begin{aligned} C_Y &= A_0 - A_* + \frac{\eta}{\beta} \Delta_0 + \frac{C_1 \eta}{\bar{\gamma}} \delta_0 \leq A_0 - A_* + O\left(\frac{1}{\sqrt{C_2}}\right) \Delta_0 + O\left(\frac{\sqrt{C_1}}{\sqrt{C_3 \bar{\gamma}}}\right) \delta_0 \\ &= A_0 - A_* + O\left(\frac{1}{\sqrt{C_2}}\right) \Delta_0 + O\left(\frac{C_1 \sigma'}{\sqrt{8 C_3} \epsilon}\right) \delta_0. \end{aligned}$$

Then, as in the proof of Theorem 5.1, we substitute $C_1 = O(\bar{L}_1^2)$, $C_2 = O(\bar{L}_F^2)$, $C_3 = O(n/B)$, $\sigma^2 = \frac{G_1^2 \sigma_0^2}{B} + \frac{G_1^2 G_2^2 (n-B)}{B(n-1)}$, and $\sigma'^2 = O(n\sigma_0^2/B)$ into the above complexity expression and C_Y , and obtain

$$\begin{aligned} T &= O\left(\max\left\{\frac{C_Y \bar{L}_F}{\epsilon^2}, \frac{C_Y \bar{L}_F \sigma^2}{\epsilon^4}, \frac{C_Y n \bar{L}_1^2 \sigma_0}{B \epsilon^3}\right\}\right), \\ C_Y &\leq O(F(\mathbf{w}_0) - \bar{F}_*) + O\left(\frac{1}{\bar{L}_F}\right) \Delta_0 + O\left(\frac{\bar{L}_1^2 \sigma_0}{\epsilon}\right) \delta_0. \end{aligned}$$

We finish the proof by noting that $\bar{L}_1 = O(1/\epsilon)$ and $\bar{L}_F = O(1/\epsilon)$ and $C_Y = O(1)$ if $\delta_0 \leq O(\epsilon^3/\sigma_0)$.

□

5.4 Convex inner and outer functions

In Chapter 3, we discussed standard stochastic convex optimization and established the iteration complexities of various algorithms. For general convex problems,

Algorithm 18 ALEXR

```

1: Input: learning rate schedules  $\{\eta_t\}_{t=1}^T, \{\alpha_t\}_{t=1}^T, \theta \in [0, 1]$ ; starting points  $\mathbf{w}_0, \mathbf{y}_1 \in \mathcal{Y}_1 \times \dots \times \mathcal{Y}_n$ 
2: Let  $\mathbf{w}_1 = \mathbf{w}_0$ 
3: for  $t = 1, \dots, T - 1$  do
4:   Sample a batch  $\mathcal{B}_t \subset \{1, \dots, n\}, |\mathcal{B}_t| = B$ 
5:   for each  $i \in \mathcal{S}_t$  do
6:     Draw a sample  $\zeta_{i,t}, \zeta'_{i,t} \sim \mathbb{P}_i$ 
7:     Compute  $\tilde{g}_{i,t} = g_i(\mathbf{w}_t; \zeta_{i,t}) + \theta(g_i(\mathbf{w}_t; \zeta_{i,t}) - g_i(\mathbf{w}_{t-1}; \zeta_{i,t}))$ 
8:     Update  $y_{i,t+1} = \arg \max_{y_i \in \mathcal{Y}_i} \left\{ y_i^\top \tilde{g}_{i,t} - f_i^*(y_i) - \frac{1}{\alpha_t} D_{\psi_i}(y_i, y_{i,t}) \right\}$ 
9:   end for
10:  For each  $i \notin \mathcal{B}_t, y_{i,t+1} = y_{i,t}$ 
11:  Compute the vanilla gradient estimator  $\mathbf{z}_t = \frac{1}{B} \sum_{i \in \mathcal{B}_t} [\partial g_i(\mathbf{w}_t; \zeta'_{i,t})]^\top y_{i,t+1}$ 
12:  Update  $\mathbf{w}_{t+1} = \arg \min_{\mathbf{w}} \left\{ \mathbf{z}_t^\top \mathbf{w} + \frac{1}{2\eta_t} \|\mathbf{w} - \mathbf{w}_t\|_2^2 + r(\mathbf{w}) \right\}$ 
13: end for
    
```

stochastic gradient descent (SGD) achieves a complexity of $O(1/\epsilon^2)$, while for μ -strongly convex problems, its complexity improves to $O(1/(\mu\epsilon))$. These analyses rely on the assumption of unbiased stochastic gradient estimators, which does not hold for convex compositional optimization problems.

In this section, we introduce stochastic algorithms for a family of convex FCCO problems, where both the inner and outer functions are convex. We establish that these algorithms attain the same order of iteration complexities as SGD in standard stochastic convex optimization. In particular, let us consider a regularized convex FCCO:

$$\min_{\mathbf{w} \in \mathbb{R}^d} F(\mathbf{w}) := \frac{1}{n} \sum_{i=1}^n f_i(g_i(\mathbf{w})) + r(\mathbf{w}), \quad (5.26)$$

where $g_i(\mathbf{w}) = \mathbb{E}_{\zeta \sim \mathbb{P}_i} [g_i(\mathbf{w}; \zeta)]$, the outer and inner functions satisfy the following assumption.

Assumption 5.8. *The following conditions hold:*

- (i) f_i is convex, G_1 -Lipschitz continuous, and $\partial f_i(\cdot) \geq 0$.
- (ii) g_i is convex and G_2 -Lipschitz continuous.
- (iii) r is μ -strongly convex for some $\mu \geq 0$.

Group DRO (5.2) could satisfy the above assumption when the individual loss function is convex and Lipschitz with respect to the model parameter. Two-way partial AUC maximization considered in Section 6.4.3 is another example satisfying the above assumption when the loss function is convex and Lipschitz continuous.

Let f_i^* denote the convex conjugate of f_i . We can write $f_i(g_i(\mathbf{w}))$ as

$$f_i(g_i(\mathbf{w})) = \max_{y_i \in \mathcal{Y}_i} (y_i^\top g_i(\mathbf{w}) - f_i^*(y_i)),$$

where $\mathcal{Y}_i = \text{dom}(f_i^*)$. Since $0 \leq \partial f_i(\cdot)$ and $\|\partial f_i(\cdot)\| \leq G_1$, hence \mathcal{Y}_i is a compact set following from Lemma 1.8.

Then, we can convert (5.26) into an equivalent minimax optimization problem:

$$\min_{\mathbf{w} \in \mathbb{R}^d} \max_{\mathbf{y} \in \mathcal{Y}} \frac{1}{n} \sum_{i=1}^n (y_i^\top g_i(\mathbf{w}) - f_i^*(y_i)) + r(\mathbf{w}), \quad (5.27)$$

where $\mathbf{y} = (y_1, \dots, y_n)^\top$, $\mathcal{Y} = \mathcal{Y}_1 \times \dots \times \mathcal{Y}_n$. Thus, the above problem is convex-concave problem under Assumption 5.8.

We introduce a method to optimize the above minimax problem. However, there are several unique challenges: (i) updating all coordinates of \mathbf{y} is difficult because it is computationally prohibitive to traverse all data points $i = 1, \dots, n$ at each iteration; (ii) we only have access to stochastic evaluations of the functions $g_i(\mathbf{w}; \zeta)$, which limits our ability to update both the corresponding coordinate of \mathbf{y} and the parameter \mathbf{w} .

5.4.1 The ALEXR Algorithm

To present the algorithm, we assume a strongly convex prox-function ψ_i for the i -th coordinate and impose the following conditions.

Assumption 5.9. Suppose ψ_i is differentiable and obeys the following conditions

- (i) ψ_i is μ_ψ -strongly convex with respect to $\|\cdot\|_2$, i.e., $\psi_i(y) \geq \psi_i(y') + \nabla \psi_i(y')^\top (y - y') + \frac{\mu_\psi}{2} \|y - y'\|_2^2$.
- (ii) $D_{f_i^*}(y, y') \geq \rho D_{\psi_i}(y, y')$ for some $\rho \geq 0$.
- (iii) The following proximal mapping can be easily computed:

$$y_{i,t+1} = \arg \max_{y_i \in \mathcal{Y}_i} \left\{ y_i^\top \tilde{g}_{i,t} - f_i^*(y_i) - \frac{1}{\alpha_t} D_{\psi_i}(y_i, y_{i,t}) \right\}.$$

A meta-algorithm, termed ALEXR, is presented in Algorithm 18. ALEXR employs stochastic block-coordinate proximal mirror ascent to update \mathbf{y} , using the prox-function ψ_i for the i -th coordinate, and applies stochastic proximal gradient descent to update \mathbf{w} . Below, we consider different choices of the prox-functions ψ_i and the corresponding updates for $y_{i,t+1}$.

ALEXR-v1 for smooth f_i : using $\psi_i = f_i^*$

When f_i is L_1 -smooth, its convex conjugate f_i^* is $1/L_1$ -strongly convex. We can use $\psi_i = f_i^*$ to define a Bregman divergence $D_{\psi_i}(y, y') = D_{f_i^*}(y, y')$.

Critical: In this case, Assumption 5.9 (i) and (ii) hold with $\mu_\psi = 1/L_1$, and $\rho = 1$.

Let us consider the update of $y_{i,t+1}$, which becomes:

$$y_{i,t+1} = \arg \max_{y_i \in \mathcal{Y}_i} \left\{ y_i^\top \tilde{g}_{i,t} - f_i^*(y_i) - \frac{1}{\alpha_t} D_{f_i^*}(y_i, y_{i,t}) \right\}, \forall i \in \mathcal{B}_t. \quad (5.28)$$

The following lemma provides an efficient way to compute $y_{i,t+1}$, which also builds the connection to the sequence of $\mathbf{u}_{i,t}$ in SOX and MSVR.

Lemma 5.14 *Let $\mathbf{u}_{i,t-1} \in \partial f_i^*(y_{i,t})$. Then for $i \in \mathcal{B}_t$ we have $y_{i,t+1} = \nabla f_i(\mathbf{u}_{i,t})$, where $\mathbf{u}_{i,t} = \frac{1}{1+\alpha_t} \mathbf{u}_{i,t-1} + \frac{\alpha_t}{1+\alpha_t} \tilde{g}_{i,t}$.*

Proof. For the problem (5.28), we have

$$\begin{aligned} & y_i^\top \tilde{g}_{i,t} - f_i^*(y_i) - \frac{1}{\alpha_t} D_{f_i^*}(y_i, y_{i,t}) \\ &= y_i^\top \tilde{g}_{i,t} - f_i^*(y_i) - \frac{1}{\alpha_t} (f_i^*(y_i) - \partial f_i^*(y_{i,t})^\top (y_i - y_{i,t}) - f_i^*(y_{i,t})) \\ &= y_i^\top (\tilde{g}_{i,t} + \frac{1}{\alpha_t} \partial f_i^*(y_{i,t})) - (1 + \frac{1}{\alpha_t}) f_i^*(y_i) - \frac{1}{\alpha_t} \partial f_i^*(y_{i,t})^\top y_{i,t} + \frac{1}{\alpha_t} f_i^*(y_{i,t}). \end{aligned}$$

Hence $y_{i,t+1} \in \arg \max_{y_i \in \mathcal{Y}_i} y_i^\top (\frac{\alpha_t}{1+\alpha_t} \tilde{g}_{i,t} + \frac{1}{1+\alpha_t} \partial f_i^*(y_{i,t})) - f_i^*(y_i)$. If we define $\mathbf{u}_{i,t} = \frac{\alpha_t}{1+\alpha_t} \tilde{g}_{i,t} + \frac{1}{1+\alpha_t} \partial f_i^*(y_{i,t})$, we have

$$f(\mathbf{u}_{i,t}) = \max_{y_i \in \mathcal{Y}_i} y_i^\top \mathbf{u}_{i,t} - f_i^*(y_i) = y_{i,t+1}^\top \mathbf{u}_{i,t} - f_i^*(y_{i,t+1}).$$

Hence, $\mathbf{u}_{i,t} \in \arg \max_{\mathbf{u}} y_{i,t+1}^\top \mathbf{u} - f_i(\mathbf{u})$ and therefore $y_{i,t+1} = \nabla f_i(\mathbf{u}_{i,t})$. \square

If f_i is a Legendre function such that $\nabla f_i^{-1} = \nabla f_i^*$ (see Lemma 1.8). Then, we can derive the following equivalent update of \mathbf{u} sequence such that $y_{i,t} = \nabla f_i(\mathbf{u}_{i,t-1})$.

$$\mathbf{u}_{i,t} = \begin{cases} \frac{1}{1+\alpha_t} \mathbf{u}_{i,t-1} + \frac{\alpha_t}{1+\alpha_t} \tilde{g}_{i,t}, & \text{if } i \in \mathcal{B}_t \\ \mathbf{u}_{i,t-1} & \text{o.w.} \end{cases}. \quad (5.29)$$

When $\theta = 0$, the equivalent \mathbf{u} update (5.64) becomes:

$$\mathbf{u}_{i,t} = (1 - \frac{\alpha_t}{1+\alpha_t}) \mathbf{u}_{i,t-1} + \frac{\alpha_t}{1+\alpha_t} g_i(\mathbf{w}_t; \zeta_{i,t}), \forall i \in \mathcal{B}_t. \quad (5.30)$$

This is the same as the moving average estimator in SOX with $\gamma_t = \alpha_t/(1+\alpha_t)$. Using the equivalent \mathbf{u} sequence, the stochastic gradient estimator becomes $\mathbf{z}_t = \frac{1}{B} \sum_{i \in \mathcal{B}_t} [\partial g_i(x_t; \zeta'_{i,t})]^\top \nabla f_i(\mathbf{u}_{i,t})$. If the regularizer r is not present, the update of the model parameter becomes $\mathbf{w}_{t+1} = \mathbf{w}_t - \eta_t \mathbf{z}_t$. In this case, ALEXR with $\theta = 0$ is the same as SOX with $\beta_t = 1$. We will prove its convergence for convex and strongly convex regularizer r .

When $\theta > 0$, the equivalent \mathbf{u} update (5.64) becomes:

$$\mathbf{u}_{i,t} = \left(1 - \frac{\alpha_t}{1 + \alpha_t}\right) \mathbf{u}_{i,t-1} + \frac{\alpha_t}{1 + \alpha_t} g_i(\mathbf{w}_t; \zeta_{i,t}) + \frac{\theta \alpha_t}{1 + \alpha_t} (g_i(\mathbf{w}_t; \zeta_t) - g_i(\mathbf{w}_{t-1}; \zeta_t)). \quad (5.31)$$

This is similar to the MSVR estimator with $\gamma_t = \frac{\alpha_t}{1 + \alpha_t}$ and $\gamma'_t = \frac{\theta \alpha_t}{1 + \alpha_t}$. However, the key difference is that γ'_t in MSVR is larger than 1, while it is smaller than 1 in ALEXR for convex problems. In practice, setting $\gamma'_t < 1$ is a better choice. We will prove a better convergence of ALEXR with $\theta \in (0, 1)$ for a strongly convex r .

ALEXR-v2 for non-smooth f_i : using a quadratic function $\psi_i(\cdot)$

When f_i is non-smooth, we cannot use f_i^* as the prox function. In this case, we will use a smooth and strongly convex ψ_i , a quadratic function $\psi_i(y) = \frac{1}{2} \|y\|_2^2$.

Critical: In this case, Assumption 5.9 (i) holds with $\mu_\psi = 1$, and Assumption 5.9 (ii) holds with $\rho = 0$.

Example

Example 5.2. For the update of $y_{i,t+1}$, consider the example $f_i(\cdot) = [\cdot]_+$, as used in GDRO and TPAUC maximization. In this case, the conjugate $f_i^*(y)$ is the indicator function of the interval $[0, 1]$. Consequently, $y_{i,t+1}$ can be computed as:

$$y_{i,t+1} = \arg \max_{y_i \in [0,1]} \left\{ y_i^\top \tilde{g}_{i,t} - \frac{1}{2\alpha_t} (y_i - y_{i,t})^2 \right\} = \Pi_{[0,1]}(y_{i,t} - \alpha_t \tilde{g}_{i,t}), \forall i \in \mathcal{B}_t,$$

where $\Pi_{[0,1]}(\cdot)$ projects the input into the range of $[0, 1]$.

5.4.2 Technical Lemmas

To facilitate the analysis, we define $(\mathbf{w}_*, \mathbf{y}_*)$ as the saddle point to the minimax problem and

$$\begin{aligned}
 F(\mathbf{w}, \mathbf{y}) &= \frac{1}{n} \sum_{i=1}^n y_i^\top g_i(\mathbf{w}) - f_i^*(y_i) + r(\mathbf{w}), \\
 \tilde{\mathbf{g}}_t &= (\tilde{g}_{1,t}, \dots, \tilde{g}_{n,t})^\top, \\
 \bar{y}_{i,t+1} &= \arg \max_{y_i \in \mathcal{Y}_i} \left\{ y_i^\top \tilde{\mathbf{g}}_{i,t} - f_i^*(y_i) - \frac{1}{\alpha_t} D_{\psi_i}(y_i, y_{i,t}) \right\}, \forall i \in [n] \\
 D_\psi(\mathbf{y}, \mathbf{y}') &= \sum_{i=1}^n D_{\psi_i}(y_i, y'_i).
 \end{aligned}$$

Note that $\bar{\mathbf{y}}_{t+1}$ is a virtual sequence, which is updated for all coordinates from \mathbf{y}_t making it independent of \mathcal{B}_t .

We make the following assumption regarding the stochastic estimators.

Assumption 5.10. *We assume that*

- (i) $\mathbb{E}_{\zeta \sim \mathbb{P}_i} [\|g_i(\mathbf{w}; \zeta) - g_i(\mathbf{w})\|_2^2] \leq \sigma_0^2$.
- (ii) $\mathbb{E}_{\zeta \sim \mathbb{P}_i} [\|\nabla g_i(\mathbf{w}; \zeta) - \nabla g_i(\mathbf{w})\|_2^2] \leq \sigma_2^2$.
- (iii) $\mathbb{E}_{i \sim \mathbb{U}_n} \left[\left\| y_i \nabla g_i(\mathbf{w}) - \frac{1}{n} \sum_{i=1}^n y_i \nabla g_i(\mathbf{w}) \right\|_2^2 \right] \leq \delta^2$ for any fixed \mathbf{y} , where \mathbb{U}_n denotes a uniform distribution.

Lemma 5.15 *The following holds for any $\mathbf{w}, \mathbf{y} \in \mathcal{Y}$ after the t -th iteration of Algorithm 18.*

$$\begin{aligned}
 &F(\mathbf{w}_{t+1}, \mathbf{y}) - F(\mathbf{w}, \bar{\mathbf{y}}_{t+1}) \\
 &\leq \frac{1}{2\eta_t} \|\mathbf{w} - \mathbf{w}_t\|_2^2 - \left(\frac{1}{2\eta_t} + \frac{\mu}{2} \right) \|\mathbf{w} - \mathbf{w}_{t+1}\|_2^2 - \frac{1}{2\eta_t} \|\mathbf{w}_{t+1} - \mathbf{w}_t\|_2^2 \\
 &+ A_t(\mathbf{y}) + B_t(\mathbf{y}) + C_t(\mathbf{w}),
 \end{aligned} \tag{5.32}$$

where

$$\begin{aligned}
 A_t(\mathbf{y}) &= \frac{1}{n\alpha_t} D_\psi(\mathbf{y}, \mathbf{y}_t) - \left(\frac{1}{n\alpha_t} + \frac{\rho}{n} \right) D_\psi(\mathbf{y}, \bar{\mathbf{y}}_{t+1}) - \frac{1}{n\alpha_t} D_\psi(\bar{\mathbf{y}}_{t+1}, \mathbf{y}_t) \\
 B_t(\mathbf{y}) &= \frac{1}{n} \sum_{i=1}^n (g_i(\mathbf{w}_{t+1}) - \tilde{\mathbf{g}}_{i,t})^\top (y_i - \bar{y}_{i,t+1}) \\
 C_t(\mathbf{w}) &= \frac{1}{n} \sum_{i=1}^n (g_i(\mathbf{w}_{t+1}) - g_i(\mathbf{w}))^\top \bar{\mathbf{y}}_{i,t+1} - \mathbf{z}_t^\top (\mathbf{w}_{t+1} - \mathbf{w}).
 \end{aligned}$$

Proof. Following Lemma 3.10, for all $i \in [n]$ the dual update rule implies that for any $y \in \mathcal{Y}$ it holds

$$\begin{aligned}
 &\tilde{g}_{i,t}^\top (y_i - \bar{y}_{i,t+1}) + f_i^*(\bar{y}_{i,t+1}) - f_i^*(y_i) \\
 &\leq \frac{1}{\alpha_t} D_{\psi_i}(y_i, y_{i,t}) - \left(\frac{1}{\alpha_t} + \rho \right) D_{\psi_i}(y_i, \bar{y}_{i,t+1}) - \frac{1}{\alpha_t} D_{\psi_i}(\bar{y}_{i,t+1}, y_{i,t}).
 \end{aligned}$$

Averaging this inequality over $i = 1, \dots, n$.

$$\begin{aligned}
& \frac{1}{n} \sum_{i=1}^n \tilde{g}_{i,t}^\top (y_{i,t} - \bar{y}_{i,t+1}) + \frac{1}{n} \sum_{i=1}^n f_i^*(\bar{y}_{i,t+1}) - \frac{1}{n} \sum_{i=1}^n f_i^*(y_i) \\
& \leq \frac{1}{n\alpha_t} D_\psi(\mathbf{y}, \mathbf{y}_t) - \left(\frac{1}{n\alpha_t} + \frac{\rho}{n} \right) D_\psi(\mathbf{y}, \bar{\mathbf{y}}_{t+1}) - \frac{1}{n\alpha_t} D_\psi(\bar{\mathbf{y}}_{t+1}, \mathbf{y}_t).
\end{aligned} \tag{5.33}$$

According to Lemma 3.6, the primal update rule implies that

$$\begin{aligned}
& \mathbf{z}_t^\top (\mathbf{w}_{t+1} - \mathbf{w}) + r(\mathbf{w}_{t+1}) - r(\mathbf{w}) \\
& \leq \frac{1}{2\eta_t} \|\mathbf{w} - \mathbf{w}_t\|_2^2 - \left(\frac{1}{2\eta_t} + \frac{\mu}{2} \right) \|\mathbf{w} - \mathbf{w}_{t+1}\|_2^2 - \frac{1}{2\eta_t} \|\mathbf{w}_{t+1} - \mathbf{w}_t\|_2^2.
\end{aligned} \tag{5.34}$$

By the definition of $F(\mathbf{w}, \mathbf{y})$, we have

$$\begin{aligned}
& F(\mathbf{w}_{t+1}, \mathbf{y}) - F(\mathbf{w}, \bar{\mathbf{y}}_{t+1}) \\
& = \frac{1}{n} \sum_{i=1}^n y_i^\top g_i(\mathbf{w}_{t+1}) - \frac{1}{n} \sum_{i=1}^n f_i^*(y_i) + r(\mathbf{w}_{t+1}) - \frac{1}{n} \sum_{i=1}^n \bar{y}_{i,t+1}^\top g_i(\mathbf{w}) \\
& \quad + \frac{1}{n} \sum_{i=1}^n f_i^*(\bar{y}_{i,t+1}) - r(\mathbf{w}) \\
& = \frac{1}{n} \sum_{i=1}^n g_i(\mathbf{w}_{t+1})^\top (y_i - \bar{y}_{i,t+1}) + \frac{1}{n} \sum_{i=1}^n f_i^*(\bar{y}_{i,t+1}) - \frac{1}{n} \sum_{i=1}^n f_i^*(y_i) \\
& \quad + \frac{1}{n} \sum_{i=1}^n (g_i(\mathbf{w}_{t+1}) - g_i(\mathbf{w}))^\top \bar{y}_{i,t+1} + r(\mathbf{w}_{t+1}) - r(\mathbf{w}).
\end{aligned}$$

Combining the equation above with (5.34) and (5.33), we can finish the proof. \square

Next, we bound the three terms $A_t(\mathbf{y})$, $B_t(\mathbf{y})$, $C_t(\mathbf{w})$ separately.

Lemma 5.16 *Let $\tau_t = 1/\alpha_t$. For \mathbf{y} that possibly depends on all randomness in the algorithm and any $\lambda_0 > 0$, we have*

$$\begin{aligned}
\mathbb{E}[A_t(\mathbf{y})] &= \mathbb{E} \left[\frac{\tau_t}{n} D_\psi(\mathbf{y}, \mathbf{y}_t) - \frac{\tau_t + \rho}{n} D_\psi(\mathbf{y}, \bar{\mathbf{y}}_{t+1}) - \frac{\tau_t}{n} D_\psi(\bar{\mathbf{y}}_{t+1}, \mathbf{y}_t) \right] \\
&\leq \mathbb{E} \left[\frac{\tau_t + \rho \left(1 - \frac{B}{n}\right)}{B} D_\psi(\mathbf{y}, \mathbf{y}_t) - \frac{\tau_t + \rho}{B} D_\psi(\mathbf{y}, \mathbf{y}_{t+1}) \right] - \frac{\tau_t}{n} \mathbb{E} [D_\psi(\bar{\mathbf{y}}_{t+1}, \mathbf{y}_t)] \\
&\quad + \mathbb{E} \left[\frac{\lambda_0(\tau_t + \rho)}{n} (D_\psi(\mathbf{y}, \hat{\mathbf{y}}_t) - D_\psi(\mathbf{y}, \hat{\mathbf{y}}_{t+1})) \right] \\
&\quad + \frac{(n-B)(\tau_t + \rho)}{2\mu_\psi \lambda_0 n B} \mathbb{E} \left[\sum_{i=1}^n \|\nabla \psi_i(\bar{y}_{i,t+1}) - \nabla \psi_i(y_{i,t})\|_2^2 \right],
\end{aligned} \tag{5.35}$$

where the sequence $\{\hat{\mathbf{y}}_t\}_t$, $\hat{\mathbf{y}}_t \in \mathcal{Y}$ is virtual. In addition, for \mathbf{y}_* , we have

$$\begin{aligned} \mathbb{E}[A_t(\mathbf{y}_*)] &= \mathbb{E} \left[\frac{\tau_t}{n} D_\psi(\mathbf{y}_*, \mathbf{y}_t) - \frac{\tau_t + \rho}{n} D_\psi(\mathbf{y}_*, \bar{\mathbf{y}}_{t+1}) - \frac{\tau_t}{n} D_\psi(\bar{\mathbf{y}}_{t+1}, \mathbf{y}_t) \right] \quad (5.36) \\ &\leq \mathbb{E} \left[\frac{\tau_t + \rho \left(1 - \frac{B}{n}\right)}{B} D_\psi(\mathbf{y}_*, \mathbf{y}_t) - \frac{\tau_t + \rho}{B} D_\psi(\mathbf{y}_*, \mathbf{y}_{t+1}) \right] - \frac{\tau_t}{n} \mathbb{E} [D_\psi(\bar{\mathbf{y}}_{t+1}, \mathbf{y}_t)]. \end{aligned}$$

Proof.

$$\begin{aligned} &\frac{\tau_t}{n} D_\psi(\mathbf{y}, \mathbf{y}_t) - \frac{\tau_t + \rho}{n} D_\psi(\mathbf{y}, \bar{\mathbf{y}}_{t+1}) - \frac{\tau_t}{n} D_\psi(\bar{\mathbf{y}}_{t+1}, \mathbf{y}_t) \quad (5.37) \\ &= \frac{\tau_t + \rho \left(1 - \frac{B}{n}\right)}{B} D_\psi(\mathbf{y}, \mathbf{y}_t) - \frac{\tau_t + \rho}{B} D_\psi(\mathbf{y}, \mathbf{y}_{t+1}) - \frac{\tau_t}{n} D_\psi(\bar{\mathbf{y}}_{t+1}, \mathbf{y}_t) \\ &+ \left(\frac{\tau_t + \rho}{B} D_\psi(\mathbf{y}, \mathbf{y}_{t+1}) - \frac{\tau_t + \rho}{n} D_\psi(\mathbf{y}, \bar{\mathbf{y}}_{t+1}) + \frac{(B-n)(\tau_t + \rho)}{nB} D_\psi(\mathbf{y}, \mathbf{y}_t) \right). \end{aligned}$$

For bounding the last three terms, we consider the following:

$$\begin{aligned} &\frac{1}{B} D_{\psi_i}(y_i, y_{i,t+1}) - \frac{1}{n} D_{\psi_i}(y_i, \bar{y}_{i,t+1}) + \frac{(B-n)}{nB} D_{\psi_i}(y_i, y_{i,t}) \quad (5.38) \\ &= \frac{1}{B} (\psi_i(y_i) - \psi_i(y_{i,t+1}) - \nabla \psi_i(y_{i,t+1})^\top (y_i - y_{i,t+1})) \\ &- \frac{1}{n} (\psi_i(y_i) - \psi_i(\bar{y}_{i,t+1}) - \nabla \psi_i(\bar{y}_{i,t+1})^\top (y_i - \bar{y}_{i,t+1})) \\ &+ \frac{(B-n)}{nB} (\psi_i(y_i) - \psi_i(y_{i,t}) - \nabla \psi_i(y_{i,t})^\top (y_i - y_{i,t})) \\ &= \left[\frac{1}{n} \left(\psi_i(\bar{y}_{i,t+1}) - \frac{n}{B} \psi_i(y_{i,t+1}) + \frac{n-B}{B} \psi_i(y_{i,t}) \right) \right] \\ &+ \left[\frac{1}{B} \nabla \psi_i(y_{i,t+1})^\top y_{i,t+1} - \frac{1}{n} \nabla \psi_i(\bar{y}_{i,t+1})^\top \bar{y}_{i,t+1} + \frac{(B-n)}{nB} \nabla \psi_i(y_{i,t})^\top y_{i,t} \right] \\ &+ \underbrace{\frac{1}{n} \left(-\frac{n}{B} \nabla \psi_i(y_{i,t+1}) + \nabla \psi_i(\bar{y}_{i,t+1}) + \frac{n-B}{B} \nabla \psi_i(y_{i,t}) \right)^\top y_i}_{\#}. \end{aligned}$$

Taking expectation over \mathcal{B}_t for the first two terms in the brackets of the above bound will give zeros. This is because that both $\bar{y}_{i,t+1}$ and $y_{i,t}$ are independent of \mathcal{B}_t such that

$$\begin{aligned} \mathbb{E}_{\mathcal{B}_t} [\psi_i(y_{i,t+1})] &= \frac{B}{n} \psi_i(\bar{y}_{i,t+1}) + \frac{n-B}{n} \psi_i(y_{i,t}), \\ \mathbb{E}_{\mathcal{B}_t} [\nabla \psi_i(y_{i,t+1})^\top y_{i,t+1}] &= \frac{B}{n} \nabla \psi_i(\bar{y}_{i,t+1})^\top \bar{y}_{i,t+1} + \frac{n-B}{n} \nabla \psi_i(y_{i,t})^\top y_{i,t}, \\ \mathbb{E}_{\mathcal{B}_t} [\nabla \psi_i(y_{i,t+1})] &= \frac{B}{n} \nabla \psi_i(\bar{y}_{i,t+1}) + \frac{n-B}{n} \nabla \psi_i(y_{i,t}). \end{aligned}$$

Next, we bound the $\#$ term. When $\mathbf{y} = \mathbf{y}_*$, expectation of $\#$ is also zero which proves (5.36).

When \mathbf{y} is possibly random, let us apply Lemma 3.13 to the update $\hat{y}_{i,t+1} = \arg \min_v -\Delta_{i,t}^\top v + \lambda_0 D_{\psi_i}(v, \hat{y}_{i,t}), \forall i$ (λ_0 to be determined), where

$$\Delta_{i,t} := -\frac{n}{B} \nabla \psi_i(y_{i,t+1}) + \nabla \psi_i(\bar{y}_{i,t+1}) + \frac{n-B}{B} \nabla \psi_i(y_{i,t})$$

is a martingale sequence due to

$$\mathbb{E}_{\mathcal{B}_t}[(\nabla \psi_i(y_{i,t+1}) - \nabla \psi_i(y_{i,t}))] = \frac{B}{n}(\nabla \psi_i(\bar{y}_{i,t+1}) - \nabla \psi_i(y_{i,t})).$$

We have

$$\mathbb{E}[\#] \leq \mathbb{E} \left[\frac{\lambda_0}{n} (D_{\psi_i}(y_i, \hat{y}_{i,t}) - D_{\psi_i}(y_i, \hat{y}_{i,t+1})) \right] + \frac{1}{2n\mu_\psi \lambda_0} \mathbb{E} \left[\|\Delta_{i,t}\|_2^2 \right].$$

Note that $\mathbb{E}_{\mathcal{B}_t}[(\nabla \psi_i(y_{i,t+1}) - \nabla \psi_i(y_{i,t}))] = \frac{B}{n}(\nabla \psi_i(\bar{y}_{i,t+1}) - \nabla \psi_i(y_{i,t}))$ such that

$$\begin{aligned} \mathbb{E}_{\mathcal{B}_t} \left[\|\Delta_{i,t}\|_2^2 \right] &= \mathbb{E}_{\mathcal{B}_t} \left\| (\nabla \psi_i(\bar{y}_{i,t+1}) - \nabla \psi_i(y_{i,t})) - \frac{n}{B}(\nabla \psi_i(y_{i,t+1}) - \nabla \psi_i(y_{i,t})) \right\|_2^2 \\ &\leq \frac{n^2}{B^2} \mathbb{E}_{\mathcal{B}_t} \left\| \nabla \psi_i(y_{i,t+1}) - \nabla \psi_i(y_{i,t}) \right\|_2^2 - \left\| (\nabla \psi_i(\bar{y}_{i,t+1}) - \nabla \psi_i(y_{i,t})) \right\|_2^2 \\ &\leq \frac{n}{B} \left\| \nabla \psi_i(\bar{y}_{i,t+1}) - \nabla \psi_i(y_{i,t}) \right\|_2^2 - \left\| (\nabla \psi_i(\bar{y}_{i,t+1}) - \nabla \psi_i(y_{i,t})) \right\|_2^2. \end{aligned}$$

Thus, we have

$$\begin{aligned} \mathbb{E}[\#] &\leq \mathbb{E} \left[\frac{\lambda_0}{n} (D_{\psi_i}(y_i, \hat{y}_{i,t}) - D_{\psi_i}(y_i, \hat{y}_{i,t+1})) \right] \\ &\quad + \frac{n-B}{2\mu_\psi \lambda_0 n B} \mathbb{E} \left[\left\| \nabla \psi_i(\bar{y}_{i,t+1}) - \nabla \psi_i(y_{i,t}) \right\|_2^2 \right]. \end{aligned}$$

Averaging (5.38) multiplied by $\tau_t + \rho$ and combining (5.37) finishes the proof. \square

Lemma 5.17 Suppose ψ_i is μ_ψ -strongly convex. For any $\lambda_2, \lambda_3, \lambda_4, \lambda_5 > 0$ and \mathbf{y} that possibly depends on all randomness in the algorithm, we have

$$\begin{aligned} \mathbb{E}[B_t(\mathbf{y})] &= \frac{1}{n} \sum_{i=1}^n \mathbb{E}(g_i(\mathbf{w}_{t+1}) - \tilde{g}_{i,t})^\top (y_i - \bar{y}_{i,t+1}) \leq \mathbb{E}[\Gamma_{t+1} - \theta \Gamma_t] \quad (5.39) \\ &\quad + \frac{(\lambda_3 + \lambda_4 \theta) \mathbb{E}[D_\psi(\bar{\mathbf{y}}_{t+1}, \mathbf{y}_t)]}{\mu_\psi n} + \frac{G_2^2 \mathbb{E} \|\mathbf{w}_{t+1} - \mathbf{w}_t\|_2^2}{2\lambda_3} + \frac{G_2^2 \theta \mathbb{E} \|\mathbf{w}_t - \mathbf{w}_{t-1}\|_2^2}{2\lambda_4} \\ &\quad + \frac{(1 + 3.5\theta + 3.5\theta^2) \sigma_0^2 \alpha_t}{\mu_\psi} + \frac{(1 + \theta) \lambda_2 \sigma_0^2}{2\mu_\psi} + \frac{\theta \sigma_0^2 \lambda_5}{2\mu_\psi} \\ &\quad + \frac{1 + \theta}{n\lambda_2} \mathbb{E}[D_\psi(\mathbf{y}, \tilde{\mathbf{y}}_t) - D_\psi(\mathbf{y}, \tilde{\mathbf{y}}_{t+1})] + \frac{\theta}{n\lambda_5} \mathbb{E}[D_\psi(\mathbf{y}, \check{\mathbf{y}}_t) - D_\psi(\mathbf{y}, \check{\mathbf{y}}_{t+1})], \end{aligned}$$

where $\Gamma_t := \frac{1}{n} \sum_{i=1}^n (g_i(\mathbf{w}_t) - g_i(\mathbf{w}_{t-1}))^\top (y_i - y_{i,t})$ and $\check{\mathbf{y}}_t, \bar{\mathbf{y}}_t$ are some virtual sequences. In addition, we have

$$\begin{aligned} \mathbb{E}[B_t(\mathbf{y}_*)] &= \frac{1}{n} \sum_{i=1}^n \mathbb{E}(g_i(\mathbf{w}_{t+1}) - \tilde{g}_{i,t})^\top (y_{i,*} - \bar{y}_{i,t+1}) \leq \mathbb{E}[\Gamma_{t+1}^* - \theta \Gamma_t^*] \quad (5.40) \\ &+ \frac{(\lambda_3 + \lambda_4 \theta) \mathbb{E}[D_\psi(\bar{\mathbf{y}}_{t+1}, \mathbf{y}_t)]}{\mu_\psi n} + \frac{G_2^2 \mathbb{E} \|\mathbf{w}_{t+1} - \mathbf{w}_t\|_2^2}{2\lambda_3} + \frac{G_2^2 \theta \mathbb{E} \|\mathbf{w}_t - \mathbf{w}_{t-1}\|_2^2}{2\lambda_4} \\ &+ \frac{(1 + 3.5\theta + 3.5\theta^2) \sigma_0^2 \alpha_t}{\mu_\psi}, \end{aligned}$$

where $\Gamma_t^* := \frac{1}{n} \sum_{i=1}^n (g_i(\mathbf{w}_t) - g_i(\mathbf{w}_{t-1}))^\top (y_{i,*} - y_{i,t})$.

Proof. Since

$$\tilde{g}_{i,t} = g_i(\mathbf{w}_t; \zeta_{i,t}) + \theta(g_i(\mathbf{w}_t; \zeta_{i,t}) - g_i(\mathbf{w}_{t-1}; \zeta_{i,t})),$$

we have

$$\begin{aligned} &\frac{1}{n} \sum_{i=1}^n (g_i(\mathbf{w}_{t+1}) - \tilde{g}_{i,t})^\top (y_i - \bar{y}_{i,t+1}) \quad (5.41) \\ &= \underbrace{\frac{1+\theta}{n} \sum_{i=1}^n (g_i(\mathbf{w}_t) - g_i(\mathbf{w}_t; \zeta_{i,t}))^\top (y_i - \bar{y}_{i,t+1})}_\text{I} \\ &+ \underbrace{\frac{\theta}{n} \sum_{i=1}^n (g_i(\mathbf{w}_{t-1}; \zeta_{i,t}) - g_i(\mathbf{w}_{t-1}))^\top (y_i - \bar{y}_{i,t+1})}_\text{II} \\ &+ \underbrace{\frac{1}{n} \sum_{i=1}^n (g_i(\mathbf{w}_{t+1}) - g_i(\mathbf{w}_t)) + \theta(g_i(\mathbf{w}_{t-1}) - g_i(\mathbf{w}_t))^\top (y_i - \bar{y}_{i,t+1})}_\text{III}. \end{aligned}$$

Define

$$\dot{y}_{i,t+1} := \arg \max_{v \in \mathcal{Y}_i} \{v^\top ((1+\theta)g_i(\mathbf{w}_t) - \theta g_i(\mathbf{w}_{t-1})) - f_i^*(v) - \frac{1}{\alpha_t} D_{\psi_i}(v, y_{i,t})\}, \forall i \in [n].$$

This update differs from that of $\bar{y}_{i,t+1}$ in that it uses full gradients instead of stochastic gradients. We decompose the I term in (5.41) as

$$\begin{aligned}
\mathbf{I} &= \frac{1+\theta}{n} \sum_{i=1}^n (g_i(\mathbf{w}_t) - g_i(\mathbf{w}_t; \zeta_{i,t}))^\top (y_i - \bar{y}_{i,t+1}) \\
&= \underbrace{\frac{1+\theta}{n} \sum_{i=1}^n (g_i(\mathbf{w}_t) - g_i(\mathbf{w}_t; \zeta_{i,t}))^\top (\dot{y}_{i,t+1} - \bar{y}_{i,t+1})}_{\mathbf{I}_1} \\
&\quad + \underbrace{\frac{1+\theta}{n} \sum_{i=1}^n (g_i(\mathbf{w}_t) - g_i(\mathbf{w}_t; \zeta_{i,t}))^\top y_i}_{\mathbf{I}_2} - \underbrace{\frac{1+\theta}{n} \sum_{i=1}^n (g_i(\mathbf{w}_t) - g_i(\mathbf{w}_t; \zeta_{i,t}))^\top \dot{y}_{i,t+1}}_{\mathbf{I}_3}.
\end{aligned}$$

Taking expectation over $\zeta_{i,t}$, $\forall i$ will make $\mathbb{E}_{\zeta_t}[\mathbf{I}_3] = 0$. Below, we will bound \mathbf{I}_1 and \mathbf{I}_2 .

$$\mathbf{I}_1 \leq \frac{1+\theta}{n} \sum_{i=1}^n \|g_i(\mathbf{w}_t) - g_i(\mathbf{w}_t; \zeta_{i,t})\|_2 \|\dot{y}_{i,t+1} - \bar{y}_{i,t+1}\|_2.$$

Since $D_{\psi_i}(y_i, y_{i,t})$ is μ_ψ -strongly convex, Lemma 3.8 implies that

$$\begin{aligned}
&\|\dot{y}_{i,t+1} - \bar{y}_{i,t+1}\|_2 \\
&\leq \frac{\alpha_t}{\mu_\psi} \left((1+\theta) \|g_i(\mathbf{w}_t) - g_i(\mathbf{w}_t; \zeta_{i,t})\|_2 + \theta \|g_i(\mathbf{w}_{t-1}) - g_i(\mathbf{w}_{t-1}; \zeta_{i,t})\|_2 \right)
\end{aligned}$$

Hence

$$\begin{aligned}
\mathbb{E}_{\zeta_t}[\mathbf{I}_1] &\leq \frac{(1+\theta)\alpha_t}{n\mu_\psi} \\
&\sum_{i=1}^n \mathbb{E} \left[(1+1.5\theta) \|g_i(\mathbf{w}_t) - g_i(\mathbf{w}_t; \zeta_{i,t})\|_2^2 + 0.5\theta \|g_i(\mathbf{w}_{t-1}) - g_i(\mathbf{w}_{t-1}; \zeta_{i,t})\|_2^2 \right] \\
&\leq \frac{(1+\theta)(1+2\theta)\sigma_0^2\alpha_t}{\mu_\psi}. \tag{5.42}
\end{aligned}$$

Next, let us handle \mathbf{I}_2 . Let us define an auxiliary sequence $\{\tilde{\mathbf{y}}_t\}_{t \geq 1}$,

$$\tilde{y}_{i,t+1} = \arg \min_{v \in \mathcal{Y}_i} \{ (g_i(\mathbf{w}_t; \zeta_{i,t}) - g_i(\mathbf{w}_t))^\top v + \frac{1}{\lambda_2} D_{\psi_i}(v, \tilde{y}_{i,t}) \},$$

where $\lambda_2 > 0$. Lemma 3.13 implies that

$$(g_i(\mathbf{w}_t) - g_i(\mathbf{w}_t; \zeta_{i,t}))^\top y_i \leq \frac{1}{\lambda_2} \mathbb{E}[D_{\psi_i}(y_i, \tilde{y}_{i,t}) - D_{\psi_i}(y_i, \tilde{y}_{i,t+1})] + \frac{\lambda_2 \sigma_0^2}{2\mu_\psi}.$$

Averaging over $i = 1, \dots, n$ and multiplying $(1+\theta)$ yields a bound of \mathbf{I}_2 :

$$\mathbb{E}[\text{I}_2] \leq \frac{1+\theta}{n\lambda_2} \mathbb{E}[D_\psi(\mathbf{y}, \tilde{\mathbf{y}}_t) - D_\psi(\mathbf{y}, \tilde{\mathbf{y}}_{t+1})] + \frac{(1+\theta)\lambda_2\sigma_0^2}{2\mu_\psi}.$$

As a result, the I term in (5.41) can be bounded as

$$\mathbb{E}[\text{I}] \leq \frac{1+\theta}{n\lambda_2} \mathbb{E}[D_\psi(\mathbf{y}, \tilde{\mathbf{y}}_t) - D_\psi(\mathbf{y}, \tilde{\mathbf{y}}_{t+1})] + \frac{(1+\theta)\lambda_2\sigma_0^2}{2\mu_\psi} + \frac{(1+\theta)(1+2\theta)\sigma_0^2\alpha_t}{\mu_\psi}. \quad (5.43)$$

Similarly, the II term in (5.41) can be bounded as

$$\mathbb{E}[\text{II}] \leq \frac{\theta}{n\lambda_5} \mathbb{E}[D_\psi(\mathbf{y}, \check{\mathbf{y}}_t) - D_\psi(\mathbf{y}, \check{\mathbf{y}}_{t+1})] + \frac{\theta\lambda_5\sigma_0^2}{2\mu_\psi} + \frac{\theta(0.5+1.5\theta)\sigma_0^2\alpha_t}{\mu_\psi}. \quad (5.44)$$

where

$$\check{\mathbf{y}}_{i,t+1} = \arg \min_{v \in \mathcal{Y}_i} \{(g_i(\mathbf{w}_{t-1}) - g_i(\mathbf{w}_{t-1}; \zeta_{i,t}))^\top v + \lambda_5 D_{\psi_i}(v, \check{\mathbf{y}}_{i,t})\}, \forall i.$$

Next, let us bound III. Recall $\Gamma_t := \frac{1}{n} \sum_{i=1}^n (g_i(\mathbf{w}_t) - g_i(\mathbf{w}_{t-1}))^\top (y_i - y_{i,t})$. For any $\lambda_3, \lambda_4 > 0$, III can be rewritten as

$$\begin{aligned} \text{III} &= \Gamma_{t+1} - \theta\Gamma_t + \frac{1}{n} \sum_{i=1}^n (g_i(\mathbf{w}_{t+1}) - g_i(\mathbf{w}_t))^\top (y_{i,t+1} - \bar{y}_{i,t+1}) \\ &\quad - \frac{\theta}{n} \sum_{i=1}^n (g_i(\mathbf{w}_t) - g_i(\mathbf{w}_{t-1}))^\top (y_{i,t} - \bar{y}_{i,t+1}) \\ &\leq \Gamma_{t+1} - \theta\Gamma_t + \frac{G_2^2 \|\mathbf{w}_{t+1} - \mathbf{w}_t\|_2^2}{2\lambda_3} + \frac{\lambda_3 \|\mathbf{y}_{t+1} - \bar{\mathbf{y}}_{t+1}\|_2^2}{2n} \\ &\quad + \frac{G_2^2 \theta \|\mathbf{w}_t - \mathbf{w}_{t-1}\|_2^2}{2\lambda_4} + \frac{\lambda_4 \theta \|\mathbf{y}_t - \bar{\mathbf{y}}_{t+1}\|_2^2}{2n}. \end{aligned}$$

Note that $y_{i,t+1} = \bar{y}_{i,t+1}$ if $i \in \mathcal{B}_t$ and $y_{i,t+1} = y_{i,t}$ otherwise. Then, $\|\mathbf{y}_{t+1} - \bar{\mathbf{y}}_{t+1}\|_2^2 \leq \|\mathbf{y}_t - \bar{\mathbf{y}}_{t+1}\|_2^2$ such that

$$\begin{aligned} \text{III} &\leq \Gamma_{t+1} - \theta\Gamma_t + \frac{G_2^2 \|\mathbf{w}_{t+1} - \mathbf{w}_t\|_2^2}{2\lambda_3} + \frac{G_2^2 \theta \|\mathbf{w}_t - \mathbf{w}_{t-1}\|_2^2}{2\lambda_4} \\ &\quad + \frac{(\lambda_3 + \lambda_4 \theta) D_\psi(\bar{\mathbf{y}}_{t+1}, \mathbf{y}_t)}{\mu_\psi n}. \end{aligned} \quad (5.45)$$

Combining (5.43), (5.45), (5.44), we have

$$\begin{aligned}
\mathbb{E}[B_t(\mathbf{y})] &\leq \mathbb{E}[\Gamma_{t+1} - \theta\Gamma_t] \\
&+ \frac{1+\theta}{n\lambda_2} \mathbb{E}[D_\psi(\mathbf{y}, \tilde{\mathbf{y}}_t) - D_\psi(\mathbf{y}, \tilde{\mathbf{y}}_{t+1})] + \frac{\theta}{n\lambda_5} \mathbb{E}[D_\psi(\mathbf{y}, \tilde{\mathbf{y}}_t) - D_\psi(\mathbf{y}, \tilde{\mathbf{y}}_{t+1})] \\
&+ \frac{(\lambda_3 + \lambda_4\theta) \mathbb{E}[D_\psi(\tilde{\mathbf{y}}_{t+1}, \mathbf{y}_t)]}{\mu_\psi n} + \frac{G_2^2 \mathbb{E}[\|\mathbf{w}_{t+1} - \mathbf{w}_t\|_2^2]}{2\lambda_3} + \frac{G_2^2 \theta \mathbb{E}[\|\mathbf{w}_t - \mathbf{w}_{t-1}\|_2^2]}{2\lambda_4} \\
&+ \frac{(1+\theta)\lambda_2\sigma_0^2}{2\mu_\psi} + \frac{\theta\sigma_0^2\lambda_5}{2\mu_\psi} + \frac{(1+3.5\theta+3.5\theta^2)\sigma_0^2\alpha_t}{\mu_\psi}.
\end{aligned}$$

□

Lemma 5.18 *When $\theta = 0$, for any $\lambda_2, \lambda_4 \geq 0$ and \mathbf{y} that possibly depends on all randomness in the algorithm, we have*

$$\begin{aligned}
\mathbb{E}[B_t(\mathbf{y})] &= \frac{1}{n} \sum_{i=1}^n \mathbb{E}(g_i(\mathbf{w}_{t+1}) - \tilde{g}_t)^\top (y_i - \bar{y}_{i,t+1}) \leq \frac{G_2^2 \mathbb{E}[\|\mathbf{w}_{t+1} - \mathbf{w}_t\|_2^2]}{4\lambda_4} + 4\lambda_4 G_1^2 \\
&+ \frac{1}{n\lambda_2} \mathbb{E}[D_\psi(\mathbf{y}, \tilde{\mathbf{y}}_t) - D_\psi(\mathbf{y}, \tilde{\mathbf{y}}_{t+1})] + \frac{\lambda_2\sigma_0^2}{2\mu_\psi} + \frac{\sigma_0^2\alpha_t}{\mu_\psi}. \tag{5.46}
\end{aligned}$$

Proof. For ALEXR with $\theta = 0$, we have $\tilde{g}_{i,t} = g_i(\mathbf{w}_t; \zeta_{i,t})$. Then, for any $\lambda_4 > 0$ we have

$$\begin{aligned}
\frac{1}{n} \sum_{i=1}^n \mathbb{E}(g_i(\mathbf{w}_{t+1}) - \tilde{g}_t)^\top (y_i - \bar{y}_{i,t+1}) &= \frac{1}{n} \sum_{i=1}^n \mathbb{E}[(g_i(\mathbf{w}_{t+1}) - g_i(\mathbf{w}_t))^\top (y_i - \bar{y}_{i,t+1})] \\
&+ \frac{1}{n} \sum_{i=1}^n \mathbb{E}[(g_i(\mathbf{w}_t) - g_i(\mathbf{w}_t; \zeta_{i,t}))^\top (y_i - \bar{y}_{i,t+1})]. \tag{5.47}
\end{aligned}$$

We bound the first term on the RHS by Young's inequality:

$$\begin{aligned}
&\frac{1}{n} \sum_{i=1}^n \mathbb{E}[(g_i(\mathbf{w}_{t+1}) - g_i(\mathbf{w}_t))^\top (y_i - \bar{y}_{i,t+1})] \\
&\leq \frac{1}{n} \sum_{i=1}^n \left(\frac{G_2^2 \|\mathbf{w}_{t+1} - \mathbf{w}_t\|_2^2}{4\lambda_4} + \lambda_4 \|y_i - \bar{y}_{i,t+1}\|_2^2 \right) \leq \frac{G_2^2 \|\mathbf{w}_{t+1} - \mathbf{w}_t\|_2^2}{4\lambda_4} + 4G_1^2.
\end{aligned}$$

The second term in (5.47) can be bounded similarly as (5.43) by:

$$\begin{aligned}
&\frac{1}{n} \sum_{i=1}^n \mathbb{E}[(g_i(\mathbf{w}_t) - g_i(\mathbf{w}_t; \zeta_{i,t}))^\top (y_i - \bar{y}_{i,t+1})] \\
&\leq \frac{1}{n\lambda_2} \mathbb{E}[D_\psi(\mathbf{y}, \tilde{\mathbf{y}}_t) - D_\psi(\mathbf{y}, \tilde{\mathbf{y}}_{t+1})] + \frac{\lambda_2\sigma_0^2}{2\mu_\psi} + \frac{\sigma_0^2\alpha_t}{\mu_\psi}.
\end{aligned}$$

Combining the above inequalities together, we have

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n \mathbb{E}(g_i(\mathbf{w}_{t+1}) - \tilde{g}_t)^\top (y_i - \bar{y}_{i,t+1}) &\leq \frac{G_2^2 \mathbb{E}[\|\mathbf{w}_{t+1} - \mathbf{w}_t\|_2^2]}{4\lambda_4} + 4\lambda_4 G_1^2 \\ &+ \frac{1}{n\lambda_2} \mathbb{E}[D_\psi(\mathbf{y}, \tilde{\mathbf{y}}_t) - D_\psi(\mathbf{y}, \tilde{\mathbf{y}}_{t+1})] + \frac{\lambda_2 \sigma_0^2}{2\mu_\psi} + \frac{\sigma_0^2 \alpha_t}{\mu_\psi}. \end{aligned}$$

□

Lemma 5.19 *If g_i is L_2 -smooth and $\eta \leq \frac{1}{2G_1 L_2}$, then*

$$\begin{aligned} \mathbb{E}[C_t(\mathbf{w}_*)] &= \mathbb{E}\left[\frac{1}{n} \sum_{i=1}^n (g_i(\mathbf{w}_{t+1}) - g_i(\mathbf{w}_*))^\top \bar{y}_{i,t+1} - \mathbf{z}_t^\top (\mathbf{w}_{t+1} - \mathbf{w}_*)\right] \\ &\leq \eta \sigma^2 + \frac{1}{4\eta} \|\mathbf{w}_{t+1} - \mathbf{w}_t\|_2^2. \end{aligned} \quad (5.48)$$

If g_i is G_2 -Lipschitz continuous, then

$$\begin{aligned} \mathbb{E}[C_t(\mathbf{w}_*)] &= \mathbb{E}\left[\frac{1}{n} \sum_{i=1}^n (g_i(\mathbf{w}_{t+1}) - g_i(\mathbf{w}_*))^\top \bar{y}_{i,t+1} - \mathbf{z}_t^\top (\mathbf{w}_{t+1} - \mathbf{w}_*)\right] \\ &\leq \eta(\sigma^2 + 4G_1^2 G_2^2) + \frac{1}{4\eta} \|\mathbf{w}_{t+1} - \mathbf{w}_t\|_2^2. \end{aligned} \quad (5.49)$$

where $\sigma^2 = \frac{G_1^2 \sigma_z^2}{B} + \frac{G_1^2 G_2^2 (n-B)}{B(n-1)}$.

Proof. We define $\Delta_t := \frac{1}{B} \sum_{i \in \mathcal{B}_t} [\partial g_i(\mathbf{w}_t; \zeta'_{i,t})]^\top y_{i,t+1} - \frac{1}{n} \sum_{i=1}^n [\partial g_i(\mathbf{w}_t)]^\top \bar{y}_{i,t+1}$. Similar to Lemma 5.2, we have $\mathbb{E}_t[\|\Delta_t\|_2^2] \leq \sigma^2$. To proceed, we have

$$\begin{aligned} &\frac{1}{n} \sum_{i=1}^n (g_i(\mathbf{w}_{t+1}) - g_i(\mathbf{w}_*))^\top \bar{y}_{i,t+1} - \mathbf{z}_t^\top (\mathbf{w}_{t+1} - \mathbf{w}_*) \\ &= \frac{1}{n} \sum_{i=1}^n (g_i(\mathbf{w}_{t+1}) - g_i(\mathbf{w}_t))^\top \bar{y}_{i,t+1} + \frac{1}{n} \sum_{i=1}^n (g_i(\mathbf{w}_t) - g_i(\mathbf{w}_*))^\top \bar{y}_{i,t+1} \\ &+ \frac{1}{n} \sum_{i=1}^n ([\partial g_i(\mathbf{w}_t)]^\top \bar{y}_{i,t+1} + \Delta_t)^\top (\mathbf{w}_* - \mathbf{w}_{t+1}). \end{aligned}$$

Since g_i is convex and $\mathcal{Y}_t \subset \mathbb{R}_+^n$ as $\partial f_i \geq 0$, we have

$$\frac{1}{n} \sum_{i=1}^n (g_i(\mathbf{w}_t) - g_i(\mathbf{w}_*))^\top \bar{y}_{i,t+1} \leq \frac{1}{n} \sum_{i=1}^n [\nabla g_i(\mathbf{w}_t)]^\top (\mathbf{w}_t - \mathbf{w}_*)^\top \bar{y}_{i,t+1}.$$

Adding the above two inequalities, we have

$$\begin{aligned}
& \frac{1}{n} \sum_{i=1}^n (g_i(\mathbf{w}_{t+1}) - g_i(\mathbf{w}_*))^\top \bar{y}_{i,t+1} - \mathbf{z}_t^\top (\mathbf{w}_{t+1} - \mathbf{w}_*) \\
& \leq \frac{1}{n} \sum_{i=1}^n (g_i(\mathbf{w}_{t+1}) - g_i(\mathbf{w}_t) - \nabla g_i(\mathbf{w}_t)^\top (\mathbf{w}_{t+1} - \mathbf{w}_t))^\top \bar{y}_{i,t+1} + \frac{1}{n} \sum_{i=1}^n \Delta_t^\top (\mathbf{w}_* - \mathbf{w}_{t+1}).
\end{aligned} \tag{5.50}$$

If g_i is L_2 -smooth, the first term in (5.50) can be bounded by

$$\begin{aligned}
& \frac{1}{n} \sum_{i=1}^n (g_i(\mathbf{w}_{t+1}) - g_i(\mathbf{w}_t) - \nabla g_i(\mathbf{w}_t)^\top (\mathbf{w}_{t+1} - \mathbf{w}_t))^\top \bar{y}_{i,t+1} \\
& \leq \frac{G_1}{n} \sum_{i=1}^n \|g_i(\mathbf{w}_{t+1}) - g_i(\mathbf{w}_t) - \nabla g_i(\mathbf{w}_t)^\top (\mathbf{w}_{t+1} - \mathbf{w}_t)\|_2 \leq \frac{G_1 L_2}{2} \|\mathbf{w}_{t+1} - \mathbf{w}_t\|_2^2.
\end{aligned} \tag{5.51}$$

To bound the second term in (5.50), we note that $\mathbb{E}_{\mathcal{B}_t, \xi_t} [\Delta_t] = 0$. Let us define $\hat{\mathbf{w}}_{t+1} = \arg \min_{\mathbf{w}} \mathbf{w}^\top \frac{1}{n} \sum_{i=1}^n [\nabla g_i(\mathbf{w}_t)]^\top \bar{y}_{i,t+1} + \frac{1}{2\eta} \|\mathbf{w} - \mathbf{w}_t\|_2^2 + r(\mathbf{w})$. Then we have

$$\begin{aligned}
& \mathbb{E} \left[\frac{1}{n} \sum_{i=1}^n \Delta_t^\top (\mathbf{w}_* - \mathbf{w}_{t+1}) \right] = \mathbb{E} \left[\frac{1}{n} \sum_{i=1}^n \Delta_t^\top (\mathbf{w}_* - \hat{\mathbf{w}}_{t+1} + \hat{\mathbf{w}}_{t+1} - \mathbf{w}_{t+1}) \right] \\
& = \mathbb{E} \left[\frac{1}{n} \sum_{i=1}^n \Delta_t^\top (\hat{\mathbf{w}}_{t+1} - \mathbf{w}_{t+1}) \right],
\end{aligned}$$

where we use the fact that

$$\mathbb{E} \left[\frac{1}{n} \sum_{i=1}^n \Delta_t^\top (\mathbf{w}_* - \hat{\mathbf{w}}_{t+1}) \right] = \mathbb{E} \left[\frac{1}{n} \sum_{i=1}^n \mathbb{E}_{\mathcal{B}_t, \xi_t'} [\Delta_t]^\top (\mathbf{w}_* - \hat{\mathbf{w}}_{t+1}) \right] = 0.$$

According to Lemma 1.7 we have

$$\mathbb{E} [\Delta_t^\top (\hat{\mathbf{w}}_{t+1} - \mathbf{w}_{t+1})] \leq \frac{\eta}{1 + \mu\eta} \mathbb{E} \|\Delta_t\|_2^2 \leq \frac{\eta\sigma^2}{1 + \mu\eta}. \tag{5.52}$$

Then, combining (5.50), (5.51) and (5.52) leads to

$$\begin{aligned}
& \frac{1}{n} \sum_{i=1}^n \mathbb{E} [(g_i(\mathbf{w}_{t+1}) - g_i(\mathbf{w}_*))^\top \bar{y}_{i,t+1}] - \mathbb{E} \mathbf{z}_t^\top (\mathbf{w}_{t+1} - \mathbf{w}_*) \\
& \leq \frac{\eta\sigma^2}{1 + \mu\eta} + \frac{L_2 G_1}{2} \|\mathbf{w}_{t+1} - \mathbf{w}_t\|_2^2,
\end{aligned}$$

which finishes the first part by noting the condition on η .

If g_i is G_2 -Lipschitz continuous, we have

$$\begin{aligned}
 & \frac{G_1}{n} \sum_{i=1}^n \|g_i(\mathbf{w}_{t+1}) - g_i(\mathbf{w}_t) - \partial g_i(\mathbf{w}_t)(\mathbf{w}_{t+1} - \mathbf{w}_t)\|_2 \\
 & \leq 2G_1G_2 \|\mathbf{w}_{t+1} - \mathbf{w}_t\|_2 \leq \eta 4G_1^2G_2^2 + \frac{1}{4\eta} \|\mathbf{w}_{t+1} - \mathbf{w}_t\|_2^2.
 \end{aligned} \tag{5.53}$$

Combining (5.50), (5.52), and (5.53), we get

$$\begin{aligned}
 & \frac{1}{n} \sum_{i=1}^n \mathbb{E}[g_i(\mathbf{w}_{t+1}) - g_i(\mathbf{w}_*)^\top \bar{\mathbf{y}}_{i,t+1}] - \mathbb{E}\mathbf{z}_t^\top (\mathbf{w}_{t+1} - \mathbf{w}_*) \\
 & \leq \eta(\sigma^2 + 4G_1^2G_2^2) + \frac{1}{4\eta} \|\mathbf{w}_{t+1} - \mathbf{w}_t\|_2^2.
 \end{aligned}$$

□

5.4.3 Strongly convex objectives

In this section, we derive a complexity of $O(1/\epsilon)$ under the the following condition.

Assumption 5.11. *We assume that the function r is μ -strongly convex ($\mu > 0$) and each f_i is L_1 -smooth, both with respect to the Euclidean norm $\|\cdot\|_2$.*

With this assumption, the minimax problem becomes strongly convex and strongly concave since the dual f_i^* is $1/L_1$ -strongly convex with respect to $\|\cdot\|_2$. In this case, we will establish the convergence of $\mu\|\mathbf{w} - \mathbf{w}_*\|_2^2$.

Critical: Under Assumption 5.11, parts (i) and (ii) of Assumption 5.9 hold for both variants of ALEXR. For ALEXR-v1, we have $\mu_\psi = 1/L_1$ and $\rho = 1$, whereas for ALEXR-v2, we have $\mu_\psi = 1$ and $\rho = 1/L_1$. Hence, the following theorem holds for both variants of ALEXR.

Let us introduce a few notations:

$$a = \frac{\epsilon\mu_\psi\rho}{24\sigma_0^2}, \quad b_1 = 3(\sigma^2 + 4G_1^2G_2^2), \quad b_2 = 3\sigma^2.$$

Theorem 5.7 *Suppose Assumptions 5.8, 5.10 and 5.11 hold.*

- *If g_i is G_2 -Lipschitz continuous, by setting $\alpha = \frac{1-\theta}{\rho(\theta-(1-B/n))}$, $\eta = \frac{1-\theta}{\theta\mu}$ and*

$$\theta = \max \left\{ 1 - \frac{a \frac{B}{n}}{1+a}, 1 - \frac{\mu\epsilon}{b_1 + \mu\epsilon} \right\}.$$

ALEXR finds a solution \mathbf{w}_{T+1} such that $\mathbb{E}[\mu\|\mathbf{w}_{T+1} - \mathbf{w}_\|_2^2] \leq \epsilon$ with an iteration complexity of*

$$T = O\left(\frac{1}{1-\theta} \log(3Y/\epsilon)\right) = \tilde{O}\left(\max\left(\frac{n}{B}, \frac{(\sigma^2 + G_1^2 G_2^2)}{\mu\epsilon}, \frac{n\sigma_0^2}{B\epsilon\mu_\psi\rho}\right)\right).$$

- If g_i is further L_2 -smooth, by setting $\alpha = \frac{1-\theta}{\rho(\theta-(1-B/n))}$, $\eta = \frac{1-\theta}{\theta\mu}$ and

$$\theta = \max\left\{1 - \frac{a\frac{B}{n}}{1+a}, 1 - \frac{\mu\epsilon}{b_2 + \mu\epsilon}, 1 - \frac{\mu}{2G_1 L_2 + \mu}\right\},$$

for sufficiently small ϵ , ALEXR finds a solution \mathbf{w}_{T+1} such that $\mathbb{E}[\mu\|\mathbf{w}_{T+1} - \mathbf{w}_*\|_2^2] \leq \epsilon$ with an iteration complexity of

$$T = O\left(\frac{1}{1-\theta} \log(3Y/\epsilon)\right) = \tilde{O}\left(\max\left(\frac{G_1 L_2}{\mu}, \frac{n}{B}, \frac{\sigma^2}{\mu\epsilon}, \frac{n\sigma_0^2}{B\epsilon\mu_\psi\rho}\right)\right).$$

where $Y = \frac{\mu}{2}\|\mathbf{w}_1 - \mathbf{w}_*\|_2^2 + \frac{2\rho}{B}D_\psi(\mathbf{y}_*, \mathbf{y}_1)$ and $\sigma^2 = \frac{G_1^2\sigma_1^2}{B} + \frac{G_2^2 G_2^2(n-B)}{B(n-1)}$.

💡 Why it matters

For smooth functions g_i , the iteration complexity is improved in the sense that the $O(1/\epsilon)$ dependence is scaled by the variance of the stochastic estimators. In contrast, for non-smooth g_i , the complexity always has a term $\frac{G_1^2 G_2^2}{\mu\epsilon}$ independent of variance.

Proof. We first consider non-smooth g_i . Combining (5.32), (5.36) for $A_t(\mathbf{y}_*)$, (5.40) for $B_t(\mathbf{y}_*)$, (5.49) for $C_t(\mathbf{w}_*)$ together we have

$$\begin{aligned} & \mathbb{E}[F(\mathbf{w}_{t+1}, \mathbf{y}_*) - F(\mathbf{w}_*, \bar{\mathbf{y}}_{t+1})] \\ & \leq \frac{1}{2\eta_t} \|\mathbf{w}_* - \mathbf{w}_t\|_2^2 - \left(\frac{1}{2\eta_t} + \frac{\mu}{2}\right) \|\mathbf{w}_* - \mathbf{w}_{t+1}\|_2^2 - \frac{1}{2\eta_t} \|\mathbf{w}_{t+1} - \mathbf{w}_t\|_2^2 \\ & + \mathbb{E}\left[\frac{\tau_t + \rho\left(1 - \frac{B}{n}\right)}{B} D_\psi(\mathbf{y}_*, \mathbf{y}_t) - \frac{\tau_t + \rho}{B} D_\psi(\mathbf{y}_*, \mathbf{y}_{t+1})\right] - \frac{\tau_t}{n} \mathbb{E}[D_\psi(\bar{\mathbf{y}}_{t+1}, \mathbf{y}_t)] \\ & + \mathbb{E}[\Gamma_{t+1}^* - \theta\Gamma_t^*] + \frac{(\lambda_3 + \lambda_4\theta)\mathbb{E}[D_\psi(\bar{\mathbf{y}}_{t+1}, \mathbf{y}_t)]}{\mu_\psi n} \\ & + \frac{G_2^2\mathbb{E}\|\mathbf{w}_{t+1} - \mathbf{w}_t\|_2^2}{2\lambda_3} + \frac{G_2^2\theta\mathbb{E}\|\mathbf{w}_t - \mathbf{w}_{t-1}\|_2^2}{2\lambda_4} + \frac{(1 + 3.5\theta + 3.5\theta^2)\sigma_0^2\alpha_t}{\mu_\psi} \\ & + \eta_t(\sigma^2 + 4G_1^2 G_2^2) + \frac{1}{4\eta_t} \|\mathbf{w}_{t+1} - \mathbf{w}_t\|_2^2. \end{aligned}$$

Define $Y_{1,t} := \frac{1}{2}\|\mathbf{w}_* - \mathbf{w}_t\|_2^2$ and $Y_{2,t} = \frac{1}{B}D_\psi(\mathbf{y}_*, \mathbf{y}_t)$. Since

$$F(\mathbf{w}_{t+1}, \mathbf{y}_*) - F(\mathbf{w}_*, \bar{\mathbf{y}}_{t+1}) \geq F(\mathbf{w}_{t+1}, \mathbf{y}_*) - F(\mathbf{w}_*, \mathbf{y}_*) + F(\mathbf{w}_*, \mathbf{y}_*) - F(\mathbf{w}_*, \bar{\mathbf{y}}_{t+1}) \geq 0,$$

multiplying the above inequality by θ^{-t} on both sides, we have

$$\begin{aligned}
 0 \leq & \theta^{-t} \mathbb{E} \left[\frac{1}{\eta_t} Y_{1,t} + (\tau_t + \rho(1 - \frac{B}{n})) Y_{2,t} - \theta \Gamma_t^* \right] \\
 & - \theta^{-t} \mathbb{E} \left[\left(\frac{1}{\eta_t} + \mu \right) Y_{1,t+1} + (\tau_t + \rho) Y_{2,t+1} - \Gamma_{t+1}^* \right] \\
 & - \theta^{-t} \left(\frac{1}{2\eta_t} \mathbb{E} \|\mathbf{w}_{t+1} - \mathbf{w}_t\|_2^2 + \frac{\tau_t}{n} \mathbb{E} [D_\psi(\bar{\mathbf{y}}_{t+1}, \mathbf{y}_t)] \right) + \theta^{-t} \frac{(\lambda_3 + \lambda_4 \theta) \mathbb{E} [D_\psi(\bar{\mathbf{y}}_{t+1}, \mathbf{y}_t)]}{\mu_\psi n} \\
 & + \theta^{-t} \frac{G_2^2 \mathbb{E} \|\mathbf{w}_{t+1} - \mathbf{w}_t\|_2^2}{2\lambda_3} + \theta^{-t} \frac{G_2^2 \theta \mathbb{E} \|\mathbf{w}_t - \mathbf{w}_{t-1}\|_2^2}{2\lambda_4} + \theta^{-t} \frac{1}{4\eta_t} \mathbb{E} \|\mathbf{w}_{t+1} - \mathbf{w}_t\|_2^2 \\
 & + \theta^{-t} \left(\frac{(1 + 3.5\theta + 3.5\theta^2) \sigma_0^2 \alpha_t}{\mu_\psi} + \eta_t (\sigma^2 + 4G_1^2 G_2^2) \right).
 \end{aligned} \tag{5.54}$$

Let

$$\frac{1}{\eta_{t-1}} + \mu = \frac{1}{\eta_t \theta}, \quad (\tau_{t-1} + \rho) = \frac{1}{\theta} (\tau_t + \rho(1 - \frac{B}{n})). \tag{5.55}$$

Hence,

$$\begin{aligned}
 & \sum_{t=1}^T \left\{ \theta^{-t} \left[\frac{1}{\eta_t} Y_{1,t} + (\tau_t + \rho(1 - \frac{B}{n})) Y_{2,t} - \theta \Gamma_t^* \right] \right. \\
 & \quad \left. - \theta^{-t} \left[\left(\frac{1}{\eta_t} + \mu \right) Y_{1,t+1} + (\tau_t + \rho) Y_{2,t+1} - \Gamma_{t+1}^* \right] \right\} \\
 & \leq \sum_{t=1}^T \left\{ \theta^{-(t-1)} \left[\left(\frac{1}{\eta_{t-1}} + \mu \right) Y_{1,t} + (\tau_{t-1} + \rho) Y_{2,t} - \Gamma_t^* \right] \right. \\
 & \quad \left. - \theta^{-t} \left[\left(\frac{1}{\eta_t} + \mu \right) Y_{1,t+1} + (\tau_t + \rho) Y_{2,t+1} - \Gamma_{t+1}^* \right] \right\} \\
 & = \left[\left(\frac{1}{\eta_0} + \mu \right) Y_{1,1} + (\tau_0 + \rho) Y_{2,1} - \Gamma_1 \right] \\
 & \quad - \theta^{-T} \left[\left(\frac{1}{\eta_T} + \mu \right) Y_{1,T+1} + (\tau_T + \rho) Y_{2,T+1} - \Gamma_{T+1} \right].
 \end{aligned}$$

Since

$$\begin{aligned}
 -\Gamma_{T+1} & \geq -\frac{1}{n} \sum_{i=1}^n G_2 \|\mathbf{w}_{T+1} - \mathbf{w}_T\|_2 \|y_{i,*} - y_{i,T+1}\|_2 \\
 & \geq -\frac{1}{n} \sum_{i=1}^n \left(\frac{G_2^2 B}{n(\rho + \tau_T) \mu_\psi} \|\mathbf{w}_{T+1} - \mathbf{w}_T\|_2^2 + \frac{n \mu_\psi (\rho + \tau_T)}{4B} \|y_{i,*} - y_{i,T+1}\|_2^2 \right) \\
 & \geq -\left(\frac{G_2^2 B}{2n(\rho + \tau_T) \mu_\psi} \|\mathbf{w}_{T+1} - \mathbf{w}_T\|_2^2 + \frac{\rho + \tau_T}{2} Y_{2,T+1} \right).
 \end{aligned} \tag{5.56}$$

Summing (5.54) over $t = 1, \dots, T$ and utilizing the above two inequalities, we have

$$\begin{aligned}
& \theta^{-T} \mathbb{E} \left[\left(\frac{1}{\eta_T} + \mu \right) \Upsilon_{1,T+1} + \frac{\rho + \tau_T}{2} \Upsilon_{2,T+1} \right] \\
& \leq \left[\left(\frac{1}{\eta_0} + \mu \right) \Upsilon_{1,1} + (\tau_0 + \rho) \Upsilon_{2,1} - \Gamma_1 \right] + \\
& \quad \frac{\theta^{-T} G_2^2 B}{2n(\rho + \tau_T) \mu_\psi} \mathbb{E} \|\mathbf{w}_{T+1} - \mathbf{w}_T\|_2^2 - \mathbb{E} \left[\sum_{t=1}^T \frac{\theta^{-t}}{2} \left(\frac{1}{\eta_t} - \frac{G_2^2}{\lambda_3} - \frac{G_2^2}{\lambda_4} - \frac{1}{2\eta_t} \right) \|\mathbf{w}_{t+1} - \mathbf{w}_t\|_2^2 \right] \\
& \quad - \mathbb{E} \left[\sum_{t=1}^T \frac{\theta^{-t}}{n} \left(\frac{1}{\alpha_t} - \frac{\lambda_3 + \lambda_4 \theta}{\mu_\psi} \right) D_\psi(\bar{\mathbf{y}}_{t+1}, \mathbf{y}_t) \right] \\
& \quad + \sum_{t=1}^T \theta^{-t} \left(\frac{(1 + 3.5\theta + 3.5\theta^2) \sigma_0^2 \alpha_t}{\mu_\psi} + \eta_t (\sigma^2 + 4G_1^2 G_2^2) \right),
\end{aligned}$$

where we use the fact $\sum_{t=1}^T \|\mathbf{w}_t - \mathbf{w}_{t-1}\|_2^2 \leq \sum_{t=1}^{T+1} \|\mathbf{w}_t - \mathbf{w}_{t-1}\|_2^2 = \sum_{t=1}^T \|\mathbf{w}_{t+1} - \mathbf{w}_t\|_2^2$.

Let $\eta_t = \eta$, $\alpha_t = \frac{1}{\tau_t} = \alpha$, $\lambda_3 = \lambda_4 = 8\eta G_2^2$. If $\alpha \leq \frac{\mu_\psi}{16\eta G_2^2}$ (to be verified later), we have $\tau \geq \frac{4\eta G_2^2 B}{n\mu_\psi}$. As a result, $\frac{G_2^2 B}{2n(\rho + \tau_t) \mu_\psi} \leq \frac{1}{8\eta}$ and $\frac{1}{\alpha_t} \geq \frac{16\eta G_2^2}{\mu_\psi}$. Then the terms related to $\|\mathbf{w}_{t+1} - \mathbf{w}_t\|$ and $D_\psi(\bar{\mathbf{y}}_{t+1}, \mathbf{y}_t)$ is less than zero. As a result,

$$\begin{aligned}
& \left[\left(\frac{1}{\eta} + \mu \right) \Upsilon_{1,T+1} + \left(\frac{\rho}{2} + \frac{1}{2\alpha} \right) \Upsilon_{2,T+1} \right] \\
& \leq \theta^T \left[\left(\frac{1}{\eta} + \mu \right) \Upsilon_{1,1} + \left(\frac{1}{\alpha} + \rho \right) \Upsilon_{2,1} \right] + \sum_{t=1}^T \theta^{T-t} \left(\frac{8\sigma_0^2 \alpha}{\mu_\psi} + \eta (\sigma^2 + 4G_1^2 G_2^2) \right) \\
& \leq \theta^T \left[\left(\frac{1}{\eta} + \mu \right) \Upsilon_{1,1} + \left(\frac{1}{\alpha} + \rho \right) \Upsilon_{2,1} \right] + \frac{1}{1-\theta} \left(\frac{8\sigma_0^2 \alpha}{\mu_\psi} + \eta (\sigma^2 + 4G_1^2 G_2^2) \right).
\end{aligned}$$

Due to the relationship between η , α and θ in (5.55), we have

$$\begin{aligned}
\theta &= \frac{1}{1 + \mu\eta} = \frac{1 + \alpha\rho(1 - B/n)}{1 + \alpha\rho} \geq \frac{1}{1 + \alpha\rho} \\
\alpha &= \frac{1 - \theta}{\rho(\theta - (1 - B/n))}, \quad \eta = \frac{1 - \theta}{\theta\mu}.
\end{aligned}$$

Then, we have

$$\begin{aligned}
 [\mu Y_{1,T+1}] &= \frac{\mu\eta}{1+\eta\mu} \left(\frac{1}{\eta} + \mu \right) Y_{1,T+1} \\
 &\leq \theta^T \mu \left[Y_{1,1} + \frac{(1+\alpha\rho)\eta}{\alpha(1+\eta\mu)} Y_{2,1} \right] + \frac{1}{1-\theta} \frac{\eta\mu}{1+\eta\mu} \left(\frac{8\sigma_0^2\alpha}{\mu_\psi} + \eta(\sigma^2 + 4G_1^2 G_2^2) \right) \\
 &= \theta^T Y + \frac{8\sigma_0^2\alpha}{\mu_\psi} + \eta(\sigma^2 + 4G_1^2 G_2^2) \\
 &\leq \theta^T Y + \frac{1-\theta}{\rho(\theta - (1-B/n))} \frac{8\sigma_0^2}{\mu_\psi} + \frac{1-\theta}{\theta\mu} (\sigma^2 + 4G_1^2 G_2^2),
 \end{aligned}$$

where $Y = \mu Y_{1,1} + \mu \frac{(1+\alpha\rho)\eta}{\alpha(1+\eta\mu)} Y_{2,1}$.

To let the RHS be less than ϵ , it is sufficient to have

$$\begin{aligned}
 T &\geq \frac{1}{1-\theta} \log(3Y/\epsilon) \geq \frac{-1}{\log(\theta)} \log(3Y/\epsilon) \Rightarrow \theta^T Y \leq \epsilon/3, \\
 \theta &\geq 1 - \frac{\epsilon\mu_\psi\rho B/(24\sigma_0^2 n)}{1 + \epsilon\mu_\psi\rho/(24\sigma_0^2)} \Rightarrow \frac{1-\theta}{\rho(\theta - (1-B/n))} \frac{8\sigma_0^2}{\mu_\psi} \leq \epsilon/3, \\
 \theta &\geq \frac{1}{1 + \mu\epsilon/(3(\sigma^2 + 4G_1^2 G_2^2))} \Rightarrow \frac{1-\theta}{\theta\mu} (\sigma^2 + 4G_1^2 G_2^2) \leq \frac{\epsilon}{3}.
 \end{aligned}$$

As a result,

$$T = O\left(\frac{1}{1-\theta} \log(3Y/\epsilon)\right) = \tilde{O}\left(\max\left(\frac{(\sigma^2 + G_1^2 G_2^2)}{\mu\epsilon}, \frac{n}{B}, \frac{n\sigma_0^2}{B\epsilon\mu_\psi\rho}\right)\right).$$

Finally, we verify that if $\epsilon^2 \leq \frac{9(\sigma^2 + 4G_1^2 G_2^2)\sigma_0^2}{2G_2^2}$, then it holds that

$$\alpha \leq \frac{\mu_\psi\epsilon}{24\sigma_0^2} = \frac{\mu_\psi\epsilon^2}{24\sigma_0^2\epsilon} \leq \frac{\mu_\psi 3(\sigma^2 + 4G_1^2 G_2^2)}{16G_2^2\epsilon} \leq \frac{\mu_\psi\theta\mu}{16G_2^2(1-\theta)} = \frac{\mu_\psi}{16\eta G_2^2}.$$

Since $\alpha\rho \leq O(1)$, we have

$$\frac{(1+\alpha\rho)\eta}{(1+\mu\eta)\alpha} \leq 2\frac{\eta}{\alpha} \leq 2\frac{\rho}{\mu},$$

thus $Y_{1,1} + \frac{(1+\alpha\rho)\eta}{\alpha(1+\mu\eta)} Y_{2,1} \leq Y_{1,1} + \frac{2\rho}{\mu} Y_{2,1}$. Thus, $Y \leq \mu Y_{1,1} + 2\rho Y_{2,1}$.

For smooth g_i , the proof is similar by using (5.48) instead of using (5.49). Hence, $\eta_t(\sigma^2 + 4G_1^2 G_2^2)$ becomes $\eta_t(\sigma^2)$ and there is additional condition $\eta_t \leq \frac{1}{2G_1 L_2}$, which transfers to a condition on θ . \square

5.4.4 Convex objectives with non-smooth outer functions

In this section, we only consider ALEXR-v2 for solving convex objectives with non-smooth f_i . For ALEXR-v2, we have that ψ is 1-smooth and 1-strongly convex. Hence, we have

$$\begin{aligned} & \frac{(n-B)(\tau+\rho)}{2\mu_\psi\lambda_0nB} \mathbb{E} \left[\sum_{i=1}^n \left\| \nabla\psi_i(\bar{\mathbf{y}}_{i,t+1}) - \nabla\psi_i(\mathbf{y}_{i,t}) \right\|_2^2 \right] \\ & \leq \frac{(n-B)(\tau+\rho)}{2\lambda_0nB} \mathbb{E} \left[\sum_{i=1}^n \left\| \bar{\mathbf{y}}_{i,t+1} - \mathbf{y}_{i,t} \right\|_2^2 \right] \leq \frac{(n-B)(\tau+\rho)}{\lambda_0nB} \mathbb{E} [D_\psi(\bar{\mathbf{y}}_{t+1}, \mathbf{y}_t)]. \end{aligned} \quad (5.57)$$

Theorem 5.8 Suppose Assumption 5.9 holds with $\rho = 0, \mu_\psi = 1$, and Assumptions 5.8, 5.10 hold. If g_i is G_2 -Lipschitz continuous, setting $\theta = 0$ and

$$\alpha = \frac{\epsilon}{6\sigma_0^2}, \quad \eta = \frac{\epsilon}{6(\sigma^2 + 8G_1^2G_2^2)},$$

ALEXR-v2 returns an ϵ -optimal solution $\bar{\mathbf{w}}_T = \sum_{t=1}^T \mathbf{w}_t / T$ in expectation with a complexity of

$$T = O \left(\frac{\sigma^2 + G_1^2G_2^2}{\epsilon^2}, \frac{\Omega\sigma_0^2}{B\epsilon^2}, \frac{\Omega\sigma_0^2}{n\epsilon^2} \right).$$

where Ω is a constant such that $\mathbb{E}[D_\psi(\mathbf{y}_T^*, \mathbf{y}_1)] \leq \Omega \leq O(G_1^2n)$, and $\mathbf{y}_T^* = \arg \max_{\mathbf{y} \in \mathcal{Y}_1 \times \dots \times \mathcal{Y}_n} F(\bar{\mathbf{w}}_T, \mathbf{y})$.

💡 Why it matters

In the worst case, the complexity is $O \left(\frac{G_1^2G_2^2}{\epsilon^2} + \frac{nG_1^2\sigma_0^2}{B\epsilon^2} \right)$. This will match the lower bounds established in next section.

Proof. Combining (5.35) with (5.57) yields

$$\mathbb{E} \left[\sum_{t=1}^T A_t(\mathbf{y}) \right] \leq \frac{\tau}{B} \mathbb{E}[D_\psi(\mathbf{y}, \mathbf{y}_1)] - \frac{\tau}{n} \mathbb{E} \left[\sum_{t=1}^T D_\psi(\bar{\mathbf{y}}_{t+1}, \mathbf{y}_t) \right] + \frac{\lambda_0\tau}{n} \mathbb{E} D_\psi(\mathbf{y}, \hat{\mathbf{y}}_1) \quad (5.58)$$

$$+ \frac{(n-B)\tau}{\lambda_0nB} \mathbb{E} \left[\sum_{t=1}^T D_\psi(\bar{\mathbf{y}}_{t+1}, \mathbf{y}_t) \right] \quad (5.59)$$

Adding this inequality with (5.32), (5.46), and (5.49) over $t = 1, \dots, T$, we have

$$\begin{aligned}
 \mathbb{E} \left[\sum_{t=1}^T F(\mathbf{w}_{t+1}, \mathbf{y}) - F(\mathbf{w}_*, \bar{\mathbf{y}}_{t+1}) \right] &\leq \frac{1}{2\eta} \|\mathbf{w}_* - \mathbf{w}_1\|_2^2 - \frac{1}{2\eta} \sum_{t=1}^T \mathbb{E}[\|\mathbf{w}_{t+1} - \mathbf{w}_t\|_2^2] \\
 &+ \frac{\tau}{B} \mathbb{E}[D_\psi(\mathbf{y}, \mathbf{y}_1)] - \frac{\tau}{n} \mathbb{E} \left[\sum_{t=1}^T D_\psi(\bar{\mathbf{y}}_{t+1}, \mathbf{y}_t) \right] + \frac{\lambda_0 \tau}{n} \mathbb{E} D_\psi(\mathbf{y}, \hat{\mathbf{y}}_1) \\
 &+ \frac{(n-B)\tau}{\lambda_0 n B} \mathbb{E} \left[\sum_{t=1}^T D_\psi(\bar{\mathbf{y}}_{t+1}, \mathbf{y}_t) \right], \\
 &+ \frac{G_2^2}{4\lambda_4} \mathbb{E} \left[\sum_{t=1}^T \|\mathbf{w}_{t+1} - \mathbf{w}_t\|_2^2 \right] + 4\lambda_4 T G_1^2 + \frac{1}{n\lambda_2} \mathbb{E}[D_\psi(\mathbf{y}, \bar{\mathbf{y}}_1)] \\
 &+ \frac{\lambda_2 \sigma_0^2}{2} T + \sigma_0^2 \alpha T, \\
 &+ \eta T (\sigma^2 + 4G_1^2 G_2^2) + \frac{1}{4\eta} \sum_{t=1}^T \mathbb{E} \|\mathbf{w}_{t+1} - \mathbf{w}_t\|_2^2.
 \end{aligned}$$

If we set $\lambda_0 = \frac{n-B}{B}$ and $\frac{G_2^2}{4\lambda_4} = \frac{1}{4\eta}$, we observe that the terms involving $\mathbb{E}[\sum_{t=1}^T \|\mathbf{w}_{t+1} - \mathbf{w}_t\|_2^2]$ and $\mathbb{E}[\sum_{t=1}^T D_\psi(\bar{\mathbf{y}}_{t+1}, \mathbf{y}_t)]$ cancel out, leaving us with the following:

$$\begin{aligned}
 &\mathbb{E} \left[\sum_{t=1}^T F(\mathbf{w}_{t+1}, \mathbf{y}) - F(\mathbf{w}_*, \bar{\mathbf{y}}_{t+1}) \right] \\
 &\leq \frac{1}{2\eta} \|\mathbf{w}_* - \mathbf{w}_1\|_2^2 + \left(\frac{\tau(1 + \lambda_0 B/n)}{B} + \frac{1}{n\lambda_2} \right) \mathbb{E} D_\psi(\mathbf{y}, \mathbf{y}_1) \\
 &+ \eta T (\sigma^2 + 8G_1^2 G_2^2) + \frac{\lambda_2 \sigma_0^2}{2} T + \sigma_0^2 \alpha T.
 \end{aligned}$$

Let $\mathbf{y} = \mathbf{y}_T^* = \arg \max F(\bar{\mathbf{w}}_T, \mathbf{y})$. Since $\frac{1}{T} \sum_{t=1}^T F(\mathbf{w}_{t+1}, \mathbf{y}) \geq F(\bar{\mathbf{w}}_T, \mathbf{y}_T^*) = F(\bar{\mathbf{w}}_T)$ and $F(\mathbf{w}_*, \bar{\mathbf{y}}_{t+1}) \leq F(\mathbf{w}_*, \mathbf{y}_*)$, we have

$$\begin{aligned}
 \mathbb{E} \left[F(\bar{\mathbf{w}}_T) - F(\mathbf{w}_*) \right] &\leq \frac{1}{2\eta T} \|\mathbf{w}_* - \mathbf{w}_1\|_2^2 + \frac{1}{T} \left(\frac{\tau(1 + \lambda_0 B/n)}{B} + \frac{1}{n\lambda_2} \right) \Omega \\
 &+ \eta T (\sigma^2 + 8G_1^2 G_2^2) + \frac{\lambda_2 \sigma_0^2}{2} + \sigma_0^2 \alpha.
 \end{aligned} \tag{5.60}$$

Let

$$\begin{aligned}
 \alpha &= \frac{\epsilon}{6\sigma_0^2}, \quad \lambda_2 = \frac{\epsilon}{3\sigma_0^2}, \quad \eta = \frac{\epsilon}{6(\sigma^2 + 8G_1^2 G_2^2)}, \\
 T &\geq O \left(\max \left(\frac{\|\mathbf{w}_1 - \mathbf{w}_*\|_2^2}{12\eta\epsilon}, \frac{\Omega(1 + \lambda_0 B/n)}{6B\epsilon\alpha}, \frac{\Omega}{6n\lambda_2\epsilon} \right) \right).
 \end{aligned}$$

Then, the RHS of (5.60) is less than ϵ . As a result, the complexity is in the order of

$$O\left(\max\left(\frac{\sigma^2 + G_1^2 G_2^2}{\epsilon^2}, \frac{\Omega \sigma_0^2}{B \epsilon^2}, \frac{\Omega \sigma_0^2}{n \epsilon^2}\right)\right).$$

□

Theorem 5.9 Suppose Assumption 5.9 holds with $\rho = 0, \mu_\psi = 1$, Assumptions 5.8, 5.10 hold. If g_i is G_2 -Lipschitz continuous and L_2 -smooth, for sufficiently small ϵ , setting $\theta = 1$ and

$$\alpha = \frac{\epsilon}{64\sigma_0^2}, \quad \eta = \min\left(\frac{\epsilon}{8\sigma^2}, \frac{1}{2G_1 L_2}\right)$$

ALEXR-v2 returns an ϵ -optimal solution $\bar{\mathbf{w}}_T = \sum_{t=1}^T \mathbf{w}_t / T$ in expectation with a complexity of

$$T = O\left(\frac{G_1 L_2}{\epsilon}, \frac{\sigma^2}{\epsilon^2}, \frac{\Omega \sigma_0^2}{B \epsilon^2}, \frac{\Omega \sigma_0^2}{n \epsilon^2}\right).$$

where Ω and \mathbf{y}_T^* are defined similarly as in last theorem.

💡 Why it matters

For smooth functions g_i , the iteration complexity is improved in the sense that the $O(1/\epsilon^2)$ dependence is scaled by the variance of the stochastic estimators. In contrast, for non-smooth g_i , the complexity always includes a term $\frac{G_1^2 G_2^2}{\epsilon^2}$, regardless of the variance.

Proof. The proof is similar to that of previous theorem except that we use (5.39) instead of (5.46), and using (5.48) instead of using (5.49). Additionally, we use

$$\begin{aligned} \sum_{t=1}^T (\Gamma_{t+1} - \Gamma_t) &= \Gamma_{T+1} - \Gamma_1 \leq \frac{1}{n} \sum_{i=1}^n G_2 \|\mathbf{w}_{T+1} - \mathbf{w}_T\|_2 \|y_i - y_{i,T+1}\|_2 \\ &\leq \frac{G_2^2 B}{2n\tau} \|\mathbf{w}_{T+1} - \mathbf{w}_T\|_2^2 + \frac{\tau n/B}{n} D_\psi(\mathbf{y}, \mathbf{y}_{T+1}). \end{aligned} \quad (5.61)$$

Combining this with (5.39), we have

$$\begin{aligned} \mathbb{E}\left[\sum_{t=1}^T B_t(\mathbf{y})\right] &\leq \frac{G_2^2 B}{2n\tau} \mathbb{E}[\|\mathbf{w}_{T+1} - \mathbf{w}_T\|_2^2] + \frac{\tau}{B} \mathbb{E}[D_\psi(\mathbf{y}, \mathbf{y}_{T+1})] \\ &+ \frac{2}{n\lambda_2} \mathbb{E}[D_\psi(\mathbf{y}, \tilde{\mathbf{y}}_1)] + \frac{(\lambda_3 + \lambda_4)}{n} \sum_{t=1}^T \mathbb{E}\left[D_\psi(\tilde{\mathbf{y}}_{t+1}, \mathbf{y}_t)\right] + \frac{G_2^2}{2\lambda_3} \sum_{t=1}^T \mathbb{E}[\|\mathbf{w}_{t+1} - \mathbf{w}_t\|_2^2] \\ &+ \frac{G_2^2}{2\lambda_4} \sum_{t=1}^T \mathbb{E}[\|\mathbf{w}_t - \mathbf{w}_{t-1}\|_2^2] + 8\sigma_0^2 \alpha T + \lambda_2 \sigma_0^2 T + \frac{\sigma_0^2 \lambda_5}{2} T + \frac{1}{n\lambda_5} \mathbb{E}[D_\psi(\mathbf{y}, \tilde{\mathbf{y}}_1)]. \end{aligned} \quad (5.62)$$

Summing the inequalities in (5.32), (5.58), (5.62), and (5.49) over $t = 1, \dots, T$, we have

$$\begin{aligned}
 & \mathbb{E} \left[\sum_{t=1}^T F(\mathbf{w}_{t+1}, \mathbf{y}) - F(\mathbf{w}_*, \bar{\mathbf{y}}_{t+1}) \right] \leq \frac{1}{2\eta} \|\mathbf{w}_* - \mathbf{w}_1\|_2^2 - \frac{1}{2\eta} \sum_{t=1}^T \mathbb{E}[\|\mathbf{w}_{t+1} - \mathbf{w}_t\|_2^2] \\
 & + \frac{\tau}{B} \mathbb{E}[D_\psi(\mathbf{y}, \mathbf{y}_1) - D_\psi(\mathbf{y}, \mathbf{y}_{T+1})] - \frac{\tau}{n} \sum_{t=1}^T \mathbb{E}[D_\psi(\bar{\mathbf{y}}_{t+1}, \mathbf{y}_t)] \\
 & + \frac{\lambda_0 \tau}{n} \mathbb{E} D_\psi(\mathbf{y}, \hat{\mathbf{y}}_1) + \frac{(n-B)\tau}{\lambda_0 n B} \sum_{t=1}^T \mathbb{E}[D_\psi(\bar{\mathbf{y}}_{t+1}, \mathbf{y}_t)], \\
 & + \frac{G_2^2 B}{2n\tau} \mathbb{E}[\|\mathbf{w}_{T+1} - \mathbf{w}_T\|_2^2] + \frac{\tau}{B} \mathbb{E}[D_\psi(\mathbf{y}, \mathbf{y}_{T+1})] \\
 & + \frac{2}{n\lambda_2} \mathbb{E}[D_\psi(\mathbf{y}, \bar{\mathbf{y}}_1)] + \frac{(\lambda_3 + \lambda_4)}{n} \sum_{t=1}^T \mathbb{E}[D_\psi(\bar{\mathbf{y}}_{t+1}, \mathbf{y}_t)] + \frac{G_2^2}{2\lambda_3} \sum_{t=1}^T \mathbb{E}[\|\mathbf{w}_{t+1} - \mathbf{w}_t\|_2^2] \\
 & + \frac{G_2^2}{2\lambda_4} \sum_{t=1}^T \mathbb{E}[\|\mathbf{w}_t - \mathbf{w}_{t-1}\|_2^2] + 8\sigma_0^2 \alpha T + \lambda_2 \sigma_0^2 T + \frac{\sigma_0^2 \lambda_5}{2} T + \frac{1}{n\lambda_5} \mathbb{E}[D_\psi(\mathbf{y}, \check{\mathbf{y}}_1)], \\
 & + \frac{1}{4\eta} \sum_{t=1}^T \mathbb{E} \|\mathbf{w}_{t+1} - \mathbf{w}_t\|_2^2 + \eta T \sigma^2.
 \end{aligned}$$

Similarly as before, if we let $\lambda_0 = \frac{2(n-B)}{B}$, $\frac{G_2^2}{2\lambda_3} = \frac{G_2^2}{2\lambda_4} = \frac{1}{16\eta}$, $\lambda_3 + \lambda_4 = 16\eta G_2^2 \leq \tau/2$, and $\frac{G_2^2 B}{2n\tau} \leq \frac{1}{8\eta}$, we observe that all the cumulated terms cancel out, leaving us the following:

$$\begin{aligned}
 & \mathbb{E} \left[\sum_{t=1}^T F(\mathbf{w}_{t+1}, \mathbf{y}) - F(\mathbf{w}_*, \bar{\mathbf{y}}_{t+1}) \right] \leq \\
 & \frac{1}{2\eta} \|\mathbf{w}_* - \mathbf{w}_1\|_2^2 + \left(\frac{\tau(1 + \lambda_0 B/n)}{B} + \frac{2}{n\lambda_2} + \frac{1}{n\lambda_5} \right) \mathbb{E} D_\psi(\mathbf{y}, \mathbf{y}_1) \\
 & + \eta T \sigma^2 + 8\sigma_0^2 \alpha T + \lambda_2 \sigma_0^2 T + \frac{\sigma_0^2 \lambda_5}{2} T.
 \end{aligned}$$

Let $\mathbf{y} = \mathbf{y}_T^* = \arg \max F(\bar{\mathbf{w}}_T, \mathbf{y})$. Since $\frac{1}{T} \sum_{t=1}^T F(\mathbf{w}_{t+1}, \mathbf{y}) \geq F(\bar{\mathbf{w}}_T, \mathbf{y}_T^*) = F(\bar{\mathbf{w}}_T)$ and $F(\mathbf{w}_*, \bar{\mathbf{y}}_{t+1}) \leq F(\mathbf{w}_*, \mathbf{y}_*)$, we have

$$\begin{aligned}
 & \mathbb{E} \left[F(\bar{\mathbf{w}}_T) - F(\mathbf{w}_*) \right] \leq \frac{1}{2\eta T} \|\mathbf{w}_* - \mathbf{w}_1\|_2^2 + \frac{1}{T} \left(\frac{\tau(1 + \lambda_0 B/n)}{B} + \frac{2}{n\lambda_2} + \frac{1}{n\lambda_5} \right) \Omega \\
 & + \eta(\sigma^2) + 8\sigma_0^2 \alpha + \lambda_2 \sigma_0^2 + \frac{\sigma_0^2 \lambda_5}{2}. \tag{5.63}
 \end{aligned}$$

Let

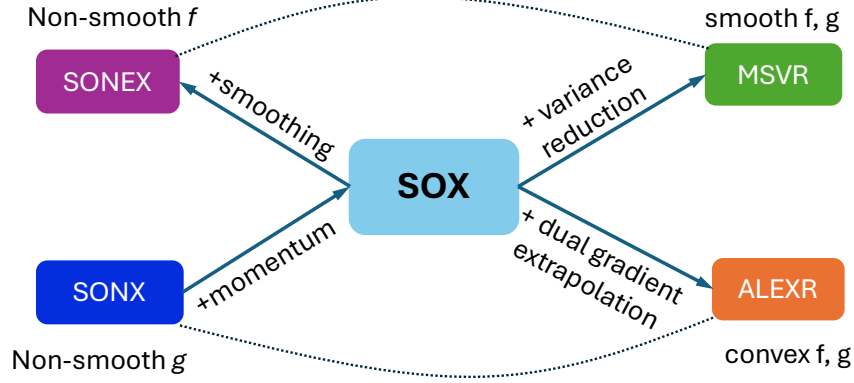


Fig. 5.1: Relationship between different algorithms for FCCO.

$$\alpha = \frac{\epsilon}{64\sigma_0^2}, \quad \lambda_2 = \frac{\epsilon}{8\sigma_0^2}, \quad \lambda_5 = \frac{\epsilon}{4\sigma_0^2}, \quad \eta = \min\left(\frac{\epsilon}{8\sigma^2}, \frac{1}{2G_1L_2}\right)$$

$$T \geq O\left(\max\left(\frac{\|\mathbf{w}_1 - \mathbf{w}_*\|_2^2}{32\eta\epsilon}, \frac{\Omega(1 + \lambda_0 B/n)}{8B\epsilon\alpha}, \frac{\Omega}{4n\lambda_2\epsilon}, \frac{\Omega}{8n\lambda_5\epsilon}\right)\right).$$

Then the conditions $16\eta G_2^2 \leq \tau/2$, $\frac{G_2^2 B}{2n\tau} \leq \frac{1}{8\eta}$ hold for sufficiently small ϵ , and the RHS of (5.63) is less than ϵ . As a result, the complexity is in the order of

$$O\left(\max\left(\frac{G_1L_2}{\epsilon}, \frac{\sigma^2}{\epsilon^2}, \frac{\Omega\sigma_0^2}{B\epsilon^2}, \frac{\sigma_0^2\Omega}{n\epsilon^2}\right)\right).$$

□

Critical: The convergence results above remain valid for ALEXR-v2 even when the outer functions f_i are smooth. If f_i is a smooth Legendre function, ALEXR-v1 can also be applied and its convergence can be established. The key is to note that

$$\|\nabla\psi_i(\bar{y}_{i,t+1}) - \nabla\psi_i(y_{i,t})\|_2^2 = \|\nabla f_i^*(\bar{y}_{i,t+1}) - \nabla f_i^*(y_{i,t})\|_2^2 = \|\bar{\mathbf{u}}_{i,t} - \mathbf{u}_{i,t-1}\|_2^2,$$

where $\mathbf{u}_{i,t-1}$ is defined in Lemma 5.14 and $\bar{\mathbf{u}}_{i,t}$ is a virtual sequence similar to $\mathbf{u}_{i,t}$ (5.64) except that all coordinates are updated by:

$$\bar{\mathbf{u}}_{i,t} = \frac{1}{1 + \alpha_t} \mathbf{u}_{i,t-1} + \frac{\alpha_t}{1 + \alpha_t} \tilde{g}_{i,t}, \forall i. \quad (5.64)$$

Then, similar to the analysis of SOX, we can establish a bound of $\sum_{t=1}^T \sum_{i=1}^n \mathbb{E}[\|\bar{\mathbf{u}}_{i,t} - \mathbf{u}_{i,t-1}\|_2^2]$ and use it to prove the convergence of ALEXR-v1. However, it remains unclear whether ALEXR-v1 provides any convergence advantage over ALEXR-v2 when f_i are smooth.

5.4.5 Double-loop ALEXR for weakly convex inner functions

ALEXR can be also useful for solving non-convex FCCO with convex outer functions and weakly convex inner functions. In particular, we consider the following non-convex problem:

$$\min_{\mathbf{w}} \frac{1}{n} \sum_{i=1}^n f_i(g_i(\mathbf{w})) + r(\mathbf{w}),$$

where $g_i : \mathbb{R}^d \rightarrow \mathbb{R}^{d'}$ and $f_i : \mathbb{R}^{d'} \rightarrow \mathbb{R}$ satisfy the following conditions:

Assumption 5.12. Assume

- (i) f_i is convex, G_1 -Lipschitz continuous and $\partial f(g) \geq 0$.
- (ii) each dimension of g_i is ρ_2 -weakly convex and G_2 -Lipschitz continuous.
- (iii) $r(\mathbf{w})$ is a convex function.

The key idea is to solve the following quadratic problem sequentially:

$$\mathbf{w}_{t+1} \approx \arg \min \bar{F}(\mathbf{w}, \mathbf{w}_t) := \frac{1}{n} \sum_{i=1}^n f_i(g_i(\mathbf{w})) + \frac{\bar{\rho}}{2} \|\mathbf{w} - \mathbf{w}_t\|_2^2,$$

where $\bar{\rho} > \rho$, with ρ being the weak-convexity parameter of $F(\mathbf{w})$. We can employ ALEXR to solve $\min_{\mathbf{w}} \bar{F}(\mathbf{w}, \mathbf{w}_t)$ approximately up to an ϵ -level. This yields a double-loop scheme.

f_i	g_i	r	F	Algorithm	Convergence Measure	Complexity	Theorem
sm	-	0	ncx, sm	SOX	Stationary	$O\left(\frac{n\sigma_0^2}{B\epsilon^4}\right)$	Thm. 5.1
sm	mss	0	ncx	MSVR	Stationary	$O\left(\frac{n\sigma_0}{B\epsilon^3}\right)$	Thm. 5.2
sm	-	pm	ncx, sm	SOX	Stationary	$O\left(\frac{n\sigma_0^2}{B\epsilon^4}\right)$	Thm. 5.1
sm	mss	pm	ncx	MSVR	Stationary	$O\left(\frac{n\sigma_0}{B\epsilon^3}\right)$	Thm. 5.2
wc, nd	wc	0	ncx	SONX (v1)	Nearly Stationary	$O\left(\frac{n\sigma_0^2}{B\epsilon^8}\right)$	Thm. 5.3
sm, nd	wc	0	ncx	SONX (v1)	Nearly Stationary	$O\left(\frac{n\sigma_0^2}{B\epsilon^6}\right)$	Thm. 5.4
wc, nd	wc	0	ncx	SONX (v2)	Nearly Stationary	$O\left(\frac{n\sigma_0}{B\epsilon^6}\right)$	Thm. 5.5
sm, nd	wc	0	ncx	SONX (v2)	Nearly Stationary	$O\left(\frac{n\sigma_0}{B\epsilon^4}\right)$	Thm. 5.5
wc, pm	sm	0	ncx	SONEX (v1)	Approx. Stationary	$O\left(\frac{n\sigma_0^2}{B\epsilon^7}\right)$	Cor. 5.1
wc, pm	sm	0	ncx	SONEX (v2)	Approx. Stationary	$O\left(\frac{n\sigma_0}{B\epsilon^5}\right)$	Thm. 5.6
nd, cvx, f_i^* pm	sm, cvx	cvx, pm	cx	ALEXR (v2)	Obj. Gap	$O\left(\max\left(\frac{\sigma^2}{\epsilon^2}, \frac{n\sigma_0^2}{B\epsilon^2}\right)\right)$	Thm. 5.9
nd, cvx, f_i^* pm	cvx	cvx, pm	cx	ALEXR (v2)	Obj. Gap	$O\left(\max\left(\frac{1}{\epsilon^2}, \frac{n\sigma_0^2}{B\epsilon^2}\right)\right)$	Thm. 5.8
sm, nd, cvx	cvx	scx, pm	cx	ALEXR	Dist. Gap	$O\left(\max\left(\frac{1}{\mu\epsilon}, \frac{n\sigma_0^2}{B\epsilon^4}\right)\right)$	Thm. 5.7
sm, nd, cvx	sm, cvx	scx, pm	cx	ALEXR	Dist. Gap	$O\left(\max\left(\frac{\sigma^2}{\mu\epsilon}, \frac{n\sigma_0^2}{B\epsilon^4}\right)\right)$	Thm. 5.7
sm, nd, cvx, f_i^* pm	wc	cx, pm	ncx	ALEXR-DL	Nearly Stationary	$O\left(\max\left(\frac{1}{\epsilon^4}, \frac{n\sigma_0^2}{B\epsilon^4}\right)\right)$	-
nd, cvx, f_i^* pm	wc	cx, pm	ncx	ALEXR-DL	Approx. Stationary	$O\left(\max\left(\frac{1}{\epsilon^5}, \frac{n\sigma_0^2}{B\epsilon^5}\right)\right)$	-

Table 5.2: Complexity of solving FCCO $F(\mathbf{w}) = \frac{1}{n} \sum_{i=1}^n f_i(g_i(\mathbf{w})) + r(\mathbf{w})$ under different conditions of f_i and g_i , where f_i is a deterministic Lipschitz continuous and g_i is mean Lipschitz continuous. pms means "proximal mapping is simple to compute", mss mean "mean squared smoothness", and ALEXR-DL denotes a double-loop method that employs ALEXR in the inner loop.

We highlight the key results as follows. If each f_i is non-smooth, the double loop method achieves a sample complexity of $O\left(\frac{n\sigma_0^2}{B\epsilon^6}\right)$ for finding a nearly ϵ -stationary solution. The analysis can be found in (Zhou et al., 2025).

If each f_i is L_1 -smooth, the sample complexity improves to $O\left(\frac{nL_1\sigma_0^2}{B\epsilon^4}\right)$ for obtaining a nearly ϵ -stationary solution. This result further implies that, for non-smooth f_i , we may apply the Nesterov smoothing \tilde{f}_i in (5.20) with $\bar{\rho}_1 = 1/\epsilon$, so that \tilde{f}_i becomes $L_1 = \bar{\rho}_1$ -smooth. Hence, Proposition 5.1 implies that the double-loop ALEXR algorithm can find an approximate ϵ -stationary stationary solution of $F(\mathbf{w})$ with a sample complexity $O\left(\frac{nL_1\sigma_0^2}{B\epsilon^4}\right) = O\left(\frac{n\sigma_0^2}{B\epsilon^5}\right)$. The analysis can be found in (Chen et al., 2025b).

Finally, we summarize the sample complexities of all methods introduced in this chapter in Table 5.2, and illustrate the relationship between different methods in Figure 5.1.

Algorithm 19 Abstract Stochastic Update Scheme for Convex FCCO

```

1: Initialize affine subspaces  $\mathfrak{X}_0, \mathfrak{Y}_0, \mathfrak{g}_0, \mathfrak{G}_0$ 
2: for  $t = 0, 1, \dots, T - 1$  do
3:   Sample a batch  $\mathcal{B}_t \subset \{1, \dots, n\}, |\mathcal{B}_t| = B$ 
4:   for each  $i \in \mathcal{B}_t$  do
5:     Sample  $\zeta_{i,t}, \tilde{\zeta}_{i,t}$  from  $\mathbb{P}_i$ 
6:      $\mathfrak{g}_{t+1}^{(i)} = \mathfrak{g}_t^{(i)} + \text{span}\{g_i(\hat{x}; \zeta_{i,t}) \mid \hat{x} \in \mathfrak{X}_t\}$ 
7:

$$\mathfrak{Y}_{t+1}^{(i)} = \mathfrak{Y}_t^{(i)} + \text{span}\left\{\arg \max_{y_i} \left\{y_i \hat{g}^{(i)} - f_i^*(y_i) - \frac{1}{\alpha} D_{\psi_i}(y_i, \hat{y}^{(i)})\right\} \mid \hat{g}^{(i)} \in \mathfrak{g}_{t+1}^{(i)}, \hat{y}^{(i)} \in \mathfrak{Y}_t^{(i)}\right\}$$

8:   end for
9:   For each  $i \notin \mathcal{B}_t$ ,  $\mathfrak{g}_{t+1}^{(i)} = \mathfrak{g}_t^{(i)}, \mathfrak{Y}_{t+1}^{(i)} = \mathfrak{Y}_t^{(i)}$ 
10:   $\mathfrak{G}_{t+1} = \mathfrak{G}_t + \text{span}\left\{\frac{1}{B} \sum_{i \in \mathcal{B}_t} \hat{y}^{(i)} \nabla g_i(\hat{x}; \tilde{\zeta}_{i,t}) \mid \hat{x} \in \mathfrak{X}_t, \hat{y} \in \mathfrak{Y}_{t+1}\right\}$ 
11:   $\mathfrak{X}_{t+1} = \mathfrak{X}_t + \text{span}\left\{\hat{G}^\top x + r(x) + \frac{1}{2\eta} \|x - \hat{x}\|_2^2 \mid \hat{x} \in \mathfrak{X}_t, \hat{G} \in \mathfrak{G}_{t+1}\right\}$ 
12: end for

```

5.4.6 Lower Bounds

In this section, we prove that the complexities of ALEXR for strongly convex and convex FCCO problems are nearly optimal by establishing the matching *lower* bounds.

What is a lower bound?

A lower bound states: for any algorithm in a certain class, there exists a “hard” optimization problem such that the algorithm cannot converge faster than a specified rate.

Lower bounds for convex optimization are typically derived under the standard oracle model, where the algorithm has access only to first-order information—either exact gradients in the deterministic setting or unbiased stochastic gradients in the stochastic setting. In the latter case, a classical result by Nemirovski and Yudin establishes that no stochastic algorithm using unbiased gradient oracles can achieve a convergence rate faster than $O(1/\sqrt{T})$ in terms of the objective gap after T iterations. For strongly convex problems, this lower bound improves to $O(1/T)$. Nevertheless, these lower bounds do not apply to convex FCCO problems or to ALEXR, because the algorithm does not have access to unbiased stochastic gradients.

Below, we establish lower bounds for an abstract stochastic update scheme described in Algorithm 19, where the symbol “+” denotes Minkowski addition. We consider an oracle model that, upon receiving a query point, returns unbiased stochastic function values and stochastic gradients of the inner functions g_i , as well as the solution to the proximal mirror-descent update of f_i^* with respect to a proximal function ψ_i . Since there are n inner functions in total, we assume that at each iteration the algorithm is allowed to access information from only B randomly selected in-

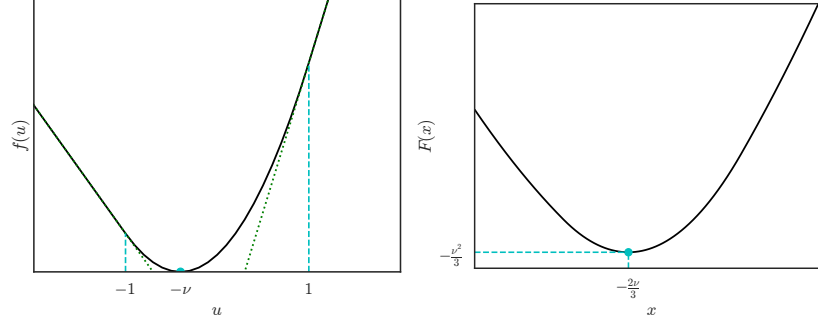


Fig. 5.2: Visualization of f (left) and F (right) in (5.65).

ner functions. Algorithm 19 is sufficiently general to encompass ALEXR, as well as SOX and MSVR.

Theorem 5.10 Consider the abstract scheme (Algorithm 19) with an initialization $\mathfrak{X}_0^{(i)} = \{0\}$, $\mathfrak{Y}_0^{(i)} = \{0\}$, $\mathfrak{g}_0^{(i)} = \emptyset$, $\mathfrak{G}_0^{(i)} = \emptyset$.

- There exists a convex FCCO problem (5.26) with smooth f_i and μ -strongly convex r such that any algorithm in the abstract scheme requires at least $T = \Omega\left(\frac{n\sigma_0^2}{B\epsilon}\right)$ iterations to find an \bar{x} that satisfies $\mathbb{E}\left[\frac{\mu}{2} \|\bar{x} - x_*\|_2^2\right] \leq \epsilon$ or $\mathbb{E}[F(\bar{x}) - F(x_*)] \leq \epsilon$.

- There exists a convex FCCO problem (5.26) with non-smooth f_i such that any algorithm in the abstract scheme requires at least $T = \Omega\left(\frac{n\sigma_0^2}{B\epsilon^2}\right)$ iterations to find an \bar{x} that satisfies $\mathbb{E}[F(\bar{x}) - F(x_*)] \leq \epsilon$.

💡 Why it matters

In light of this theorem, we see that ALEXR (v1/v2) attains a nearly optimal complexity up to a logarithmic factor for solving strongly convex FCCO problems, as its upper bounds are $\tilde{O}\left(\max\left(\frac{1}{\mu\epsilon}, \frac{n\sigma_0^2}{B\epsilon}\right)\right)$. Moreover, ALEXR-v2 achieves the optimal complexity for solving convex FCCO problems with non-smooth outer functions.

Proof. We construct the hard problems for (i) smooth f_i ; and (ii) non-smooth f_i separately.

(i) Smooth f_i and strongly convex r : Consider the following strongly convex FCCO problem

$$\min_{x \in \mathcal{X}} F(x) = \frac{1}{n} \sum_{i=1}^n f(g_i(x)) + r(x),$$

$$f(u) = \begin{cases} (\nu - 1)u + \frac{1}{2}(\nu - 1)^2 + \nu - 1 - \frac{\nu^2}{2}, & u \in (-\infty, -1) \\ \frac{1}{2}(u + \nu)^2 - \frac{\nu^2}{2}, & u \in [-1, 1] \\ (1 + \nu)u + \frac{1}{2}(1 + \nu)^2 - 1 - \nu - \frac{\nu^2}{2}, & u \in (1, \infty) \end{cases}, \quad r(x) = \frac{1}{4n} \|x\|_2^2, \quad (5.65)$$

where $\mathcal{X} = [-1, 1]^n$, the outer function $f : \mathbb{R} \rightarrow \mathbb{R}$ is smooth and Lipschitz continuous for some $\nu \in (0, 1/2)$. Besides, the inner function $g_i : \mathbb{R}^n \rightarrow \mathbb{R}$ is $g_i(x) = \mathbb{E}_{\zeta \sim \mathbb{P}}[g_i(x; \zeta)]$ and $g_i(x; \zeta) = x_i + \zeta$, where ζ follows a distribution \mathbb{P} defined below:

$$\mathbb{P} : \begin{cases} \Pr(\zeta = -\nu) = 1 - p, \\ \Pr(\zeta = \nu(1 - p)/p) = p \end{cases}, \quad \text{where } p := \frac{\nu^2}{\sigma_0^2} < 1.$$

We will determine the values of ν later. We can verify that

$$\mathbb{E}_{\zeta}[|g_i(x; \zeta) - g_i(x)|^2] = \mathbb{E}_{\zeta}[\zeta^2] = \nu^2(1 - p) + \frac{\nu^2(1 - p)^2}{p} = \frac{\nu^2(1 - p)}{p} \leq \sigma_0^2.$$

By the definition of convex conjugate, for any $y_i \in \mathbb{R}$ we have

$$f^*(y_i) = \max \left\{ \sup_{u < -1} \left\{ uy_i - \left((\nu - 1)u + \frac{1}{2}(\nu - 1)^2 + \nu - 1 - \frac{\nu^2}{2} \right) \right\}, \right. \\ \left. \sup_{-1 \leq u \leq 1} \left\{ uy_i - \frac{1}{2}(u + \nu)^2 + \frac{\nu^2}{2} \right\}, \right. \quad (5.66)$$

$$\left. \sup_{u > 1} \left\{ uy_i - \left((1 + \nu)u + \frac{1}{2}(1 + \nu)^2 - 1 - \nu - \frac{\nu^2}{2} \right) \right\} \right\} \\ = \begin{cases} +\infty, & y_i \in (-\infty, \nu - 1) \cup (\nu + 1, \infty) \\ \frac{1}{2}(y_i - \nu)^2, & y_i \in [\nu - 1, \nu + 1]. \end{cases} \quad (5.67)$$

We define $F_i(x_i) := f(g_i(x)) + \frac{1}{4}[x_i]^2$ such that $F(x) = \frac{1}{n} \sum_{i=1}^n F_i(x_i)$. Let $x_* = \arg \min_{x \in \mathcal{X}} F(x)$. Since the problem is separable over the coordinates, we have $x_{i,*} = \arg \min_{x_i \in [-1, 1]} F_i(x_i)$. Thus, we have $x_{i,*} = -\frac{2\nu}{3}$ and $F_i(x_{i,*}) = -\frac{\nu^2}{3}$.

Since $\mathbb{P}_i = \mathbb{P}$ in the “hard” problem (5.65), the abstract scheme (Algorithm 19) only needs to sample shared $\zeta_t, \tilde{\zeta}_t \sim \mathbb{P}$ for all coordinates $i \in \mathcal{S}_t$ in the t -th iteration. For any $i \in [n]$, suppose that $\mathbf{g}_{\tau}^{(i)} = \emptyset$ or $\{-\nu\}$, $\mathbf{y}_{\tau}^{(i)} = \{0\}$, $\mathbf{x}_{\tau}^{(i)} = \{0\}$ for all $\tau \leq t$. Note that when $\mathbf{g}_{\tau}^{(i)} = \emptyset$, it means that the corresponding $y^{(i)}$ will not be updated. Then,

- If $i \notin \mathcal{B}_t$, the abstract scheme (Algorithm 19) leads to

$$\mathbf{g}_{t+1}^{(i)} = \emptyset \text{ or } \{-\nu\}, \quad \mathbf{y}_{t+1}^{(i)} = \{0\}, \quad \mathbf{x}_{t+1}^{(i)} = \{0\}.$$

- If $i \in \mathcal{B}_t$ and $\zeta_t = -\nu$, the abstract scheme (Algorithm 19) proceeds as

$$\begin{aligned}
\mathbf{g}_{t+1}^{(i)} &= \mathbf{g}_t^{(i)} + \text{span} \left\{ \hat{x}_i + \zeta_t \mid \hat{x}_i \in \mathbf{x}_t^{(i)} \right\}, \\
\mathfrak{Y}_{t+1}^{(i)} &= \mathfrak{Y}_t^{(i)} \\
&+ \text{span} \left\{ \arg \max_{y_i \in [\nu-1, \nu+1]} \left\{ y_i \hat{g}_i - \frac{1}{2} (y_i - \nu)^2 - \frac{1}{\alpha} D_{\psi_i}(y_i, \hat{y}_i) \right\} \mid \hat{g}_i \in \mathbf{g}_{t+1}^{(i)}, \hat{y}_i \in \mathfrak{Y}_t^{(i)} \right\}, \\
\mathbf{x}_{t+1}^{(i)} &= \mathbf{x}_t^{(i)} \\
&+ \text{span} \left\{ \arg \min_{x_i \in [-1, 1]} \left\{ \frac{1}{B} \hat{y}_i x_i + \frac{1}{4n} [x_i]^2 + \frac{1}{2\eta} (x_i - \hat{x}_i)^2 \right\} \mid \hat{y}_i \in \mathfrak{Y}_{t+1}^{(i)}, \hat{x}_i \in \mathbf{x}_t^{(i)} \right\}.
\end{aligned}$$

Then, we can derive that $\mathbf{g}_{t+1}^{(i)} = \emptyset$ or $\{-\nu\}$, $\mathfrak{Y}_{t+1}^{(i)} = \{0\}$, and $\mathbf{x}_{t+1}^{(i)} = \{0\}$.

To sum up, given the event $\bigcap_{\tau=1}^t \{\mathbf{g}_\tau^{(i)} = \emptyset \text{ or } \{-\nu\}, \mathfrak{Y}_\tau^{(i)} = \{0\}, \mathbf{x}_\tau^{(i)} = \{0\}\}$, we can make sure that $\{\mathbf{g}_{t+1}^{(i)} = \emptyset \text{ or } \{-\nu\} \wedge \mathfrak{Y}_{t+1}^{(i)} = \{0\} \wedge \mathbf{x}_{t+1}^{(i)} = \{0\}\}$ for the abstract scheme in Algorithm 19 when one of the following mutually exclusive events happens:

- Event I: $i \notin \mathcal{B}_t$;
- Event II: $i \in \mathcal{B}_t$ and $\zeta_t = -\nu$.

Note that the random variable ζ_t is independent of \mathcal{B}_t . Thus, the probability of the event $E_{t+1}^{(i)} := \{\mathbf{g}_{t+1}^{(i)} = \emptyset \text{ or } \{-\nu\} \wedge \mathfrak{Y}_{t+1}^{(i)} = \{0\} \wedge \mathbf{x}_{t+1}^{(i)} = \{0\}\}$ conditioned on $\bigcap_{\tau=1}^t E_\tau^{(i)}$ can be bounded as

$$\begin{aligned}
\Pr \left[E_{t+1}^{(i)} \mid \bigcap_{\tau=1}^t E_\tau^{(i)} \right] &= \mathbb{P} \left[\left\{ \mathbf{g}_{t+1}^{(i)} = \emptyset \text{ or } \{-\nu\} \wedge \mathfrak{Y}_{t+1}^{(i)} = \{0\} \wedge \mathbf{x}_{t+1}^{(i)} = \{0\} \right\} \mid \bigcap_{\tau=1}^t E_\tau^{(i)} \right] \\
&\geq \mathbb{P} [\{i \notin \mathcal{B}_t\}] + \mathbb{P} [\{i \in \mathcal{B}_t\} \wedge \{\zeta_t = -\nu\}] \\
&= \mathbb{P} [\{i \notin \mathcal{B}_t\}] + \mathbb{P} [\{i \in \mathcal{B}_t\}] \mathbb{P} [\{\zeta_t = -\nu\}] \\
&= \left(1 - \frac{B}{n} \right) + \frac{B}{n} (1 - p) = 1 - \frac{Bp}{n}.
\end{aligned}$$

Since \mathcal{B}_t and ζ_t in different iterations t are mutually independent, we have

$$\Pr \left[E_T^{(i)} \right] \geq \mathbb{P} \left[\bigcap_{t=0}^{T-1} E_{t+1}^{(i)} \right] = \prod_{t=0}^{T-1} \mathbb{P} \left[E_{t+1}^{(i)} \mid \bigcap_{t=1}^t E_t^{(i)} \right] = \left(1 - \frac{Bp}{n} \right)^T > 3/4 - \frac{TBp}{n},$$

where the last inequality is due to the Bernoulli inequality $(1+x)^r \geq 1+rx$ for every integer $r \geq 1$ and $x \geq -1$.

Thus, if $T < \frac{n}{4Bp}$ we have $\Pr \left[E_T^{(i)} \right] > \frac{1}{2}$. Let us set $\nu = 3\sqrt{2\epsilon}$ such that $p = \frac{\nu^2}{\sigma_0^2} = \frac{18\epsilon}{\sigma_0^2}$. For any $i \in [n]$ and any output $\bar{x}_i \in \mathbf{x}_T^{(i)}$, we have

$$\begin{aligned}
 \mathbb{E} \left[(\bar{x}_i - x_{i,*})^2 \right] &= \mathbb{E} \left[\mathbb{I}_{E_T^{(i)}} (\bar{x}_i - x_{i,*})^2 + \mathbb{I}_{\overline{E_T^{(i)}}} (\bar{x}^{(i)} - x_{i,*})^2 \right] \\
 &\geq \mathbb{E} \left[\mathbb{I}_{E_T^{(i)}} (\bar{x}_i - x_{i,*})^2 \right] \\
 &= \mathbb{E} \left[\mathbb{I}_{E_T^{(i)}} (x_{i,*})^2 \right] = \Pr \left[E_T^{(i)} \right] (x_{i,*})^2 > \frac{2\nu^2}{9} = 4\epsilon,
 \end{aligned}$$

where \mathbb{I}_E denotes the indicator function of an event E . Moreover, we have

$$\begin{aligned}
 \mathbb{E}[F_i(\bar{x}_i) - F_i(x_{i,*})] &= \mathbb{E} \left[\mathbb{I}_{E_T^{(i)}} (F_i(\bar{x}_i) - F_i(x_{i,*})) + \mathbb{I}_{\overline{E_T^{(i)}}} (F_i(\bar{x}^{(i)}) - F_i(x_{i,*})) \right] \\
 &\geq \mathbb{E} \left[\mathbb{I}_{E_T^{(i)}} (F_i(\bar{x}_i) - F_i(x_{i,*})) \right] \\
 &= \mathbb{E} \left[\mathbb{I}_{E_T^{(i)}} (F_i(0) - F_i(x_{i,*})) \right] = \Pr[E_T^{(i)}] (F_i(0) - F_i(x_{i,*})) \\
 &> \frac{\nu^2}{6} > \epsilon.
 \end{aligned}$$

Since the derivations above hold for arbitrary $i \in [n]$ and the $r(x)$ in (5.65) is $\frac{1}{2n}$ -strongly convex ($\mu = \frac{1}{2n}$), we can derive that

$$\begin{aligned}
 \mathbb{E} \left[\frac{\mu}{2} \|\bar{x} - x_*\|_2^2 \right] &= \mathbb{E} \left[\frac{1}{4n} \|\bar{x} - x_*\|_2^2 \right] = \frac{1}{4n} \sum_{i=1}^n \mathbb{E} \left[(\bar{x}_i - x_{i,*})^2 \right] > \epsilon, \\
 \mathbb{E} [F(\bar{x}) - F(x_*)] &= \frac{1}{n} \sum_{i=1}^n \mathbb{E} [F_i(\bar{x}_i) - F_i(x_{i,*})] > \epsilon.
 \end{aligned}$$

Thus, to find an output \bar{x} that satisfies $\mathbb{E} \left[\frac{\mu}{2} \|\bar{x} - x_*\|_2^2 \right] \leq \epsilon$ or $\mathbb{E} [F(\bar{x}) - F(x_*)] \leq \epsilon$, the abstract scheme requires at least $T \geq \frac{n}{4B\rho} = \frac{n\sigma_0^2}{72B\epsilon}$ iterations.

(ii) Non-smooth f_i : Let $g_i(x) = \mathbb{E}_\zeta [x_i + \zeta] = x_i$ be defined the same as in the smooth case. Let $F_i(x_i) := f(g_i(x)) + \frac{\alpha}{2} [x_i]^2 = \beta \max\{x_i, -\nu\} + \frac{\alpha}{2} [x_i]^2$ such that $F(x) = \frac{1}{n} \sum_{i=1}^n F_i(x_i)$, where $\alpha, \beta > 0$. Let the domain \mathcal{X} be $[-2\nu, 2\nu]^n$. Hence, f is β -Lipschitz continuous and F is α -strongly convex. By the definition of convex conjugate, we have $f(\hat{g}_i) = \max_{y_i \in [0, \beta]} \{y_i \hat{g}_i - \nu(\beta - y_i)\}$.

Since the problem is separable over the coordinates, we have

$$x_{i,*} = \arg \min_{x \in [-2\nu, 2\nu]} F_i(x_i) = \arg \min_{x_i \in [-2\nu, 2\nu]} \left\{ \beta \max\{x_i, -\nu\} + \frac{\alpha}{2} [x_i]^2 \right\}.$$

Considering

$$F_i(x_i) = \begin{cases} \beta x_i + \frac{\alpha}{2} [x_i]^2 & x_i \geq -\nu \\ -\beta\nu + \frac{\alpha}{2} [x_i]^2 & x_i < -\nu \end{cases},$$

we have

$$x_{i,*} = \begin{cases} -\beta/\alpha & \text{if } \alpha > \beta/\nu \\ -\nu & \text{if } \alpha \in \frac{\beta}{\nu}[0, 1] \end{cases}, \quad F_i(x_{i,*}) \leq \begin{cases} -\beta^2/(2\alpha) & \text{if } \alpha > \beta/\nu \\ -\beta\nu/2 & \text{if } \alpha \in \frac{\beta}{\nu}[0, 1] \end{cases}.$$

Since $F_i(0) = 0$, we can derive that $F_i(0) - F_i(x_{i,*}) \geq \frac{1}{2} \min\{\beta\nu, \beta^2/\alpha\}$. Consider an arbitrary $i \in [n]$. Suppose that $\mathfrak{g}_\tau^{(i)} = \emptyset$ or $\{-\nu\}$, $\mathfrak{X}_\tau^{(i)} = \{0\}$, $\mathfrak{Y}_\tau^{(i)} = \{0\}$ for all $\tau \leq t$.

- If $i \notin \mathcal{B}_t$, the abstract scheme (Algorithm 19) leads to

$$\mathfrak{g}_{t+1}^{(i)} = \emptyset \text{ or } \{-\nu\}, \quad \mathfrak{Y}_{t+1}^{(i)} = \{0\}, \quad \mathfrak{X}_{t+1}^{(i)} = \{0\}.$$

- If $i \in \mathcal{B}_t$, the abstract scheme (Algorithm 19) proceeds as

$$\begin{aligned} \mathfrak{g}_{t+1}^{(i)} &= \mathfrak{g}_t^{(i)} + \text{span} \left\{ \hat{x}_i + \zeta_t \mid \hat{x}_i \in \mathfrak{X}_t^{(i)} \right\}, \\ \mathfrak{Y}_{t+1}^{(i)} &= \mathfrak{Y}_t^{(i)} \\ &+ \text{span} \left\{ \arg \max_{y_i \in [0, \beta]} \left\{ y_i \hat{g}_i - \nu(\beta - y_i) - \frac{1}{\alpha} D_\psi(y_i, \hat{y}_i) \right\} \mid \hat{g}_i \in \mathfrak{g}_{t+1}^{(i)}, \hat{y}_i \in \mathfrak{Y}_t^{(i)} \right\}, \\ \mathfrak{X}_{t+1}^{(i)} &= \mathfrak{X}_t^{(i)} \\ &+ \text{span} \left\{ \arg \min_{x_i \in [-2\nu, 2\nu]} \left\{ \frac{1}{B} \hat{y}_i x_i + \frac{1}{n} [x_i]^2 + \frac{1}{2\eta} (x_i - \hat{x}_i)^2 \right\} \mid \hat{y}_i \in \mathfrak{Y}_{t+1}^{(i)}, \hat{x}_i \in \mathfrak{X}_t^{(i)} \right\}. \end{aligned}$$

Due to the same reason as in the smooth f_i case, the probability of the event $E_T^{(i)} := \{\mathfrak{g}_T^{(i)} = \emptyset \text{ or } \{-\nu\} \wedge \mathfrak{Y}_T^{(i)} = \{0\} \wedge \mathfrak{X}_T^{(i)} = \{0\}\}$ can be bounded as

$$\Pr \left[E_T^{(i)} \right] \geq \mathbb{P} \left[\bigcap_{t=0}^{T-1} E_{t+1}^{(i)} \right] = \prod_{t=0}^{T-1} \mathbb{P} \left[E_{t+1}^{(i)} \mid \bigcap_{t=1}^t E_t^{(i)} \right] = \left(1 - \frac{Bp}{n} \right)^T > 3/4 - \frac{TBp}{n}.$$

Thus, if $T < \frac{n}{4Bp}$ we have $\mathbb{P} \left[E_T^{(i)} \right] > \frac{1}{2}$. Let us set $\beta = G_1$, $\nu = \frac{4\epsilon}{G_1}$ such that $p := \frac{\nu^2}{\sigma_0^2} = \frac{16\epsilon^2}{G_1^2 \sigma_0^2}$. For any $i \in [n]$ and any output $\bar{x}_i \in \mathfrak{X}_T^{(i)}$, we have

$$\begin{aligned} \mathbb{E}[F_i(\bar{x}_i) - F_i(x_{i,*})] &= \mathbb{E} \left[\mathbb{I}_{E_T^{(i)}} (F_i(\bar{x}_i) - F_i(x_{i,*})) + \mathbb{I}_{\overline{E_T^{(i)}}} (F_i(\bar{x}_i) - F_i(x_{i,*})) \right] \\ &\geq \mathbb{E} \left[\mathbb{I}_{E_T^{(i)}} (F_i(\bar{x}_i) - F_i(x_{i,*})) \right] \\ &= \mathbb{E} \left[\mathbb{I}_{E_T^{(i)}} (F_i(0) - F_i(x_{i,*})) \right] \\ &= \Pr[E_T^{(i)}] (F_i(0) - F_i(x_{i,*})) > \min\{\beta\nu, \beta^2/\alpha\}/4 = \epsilon. \end{aligned}$$

Since the derivations above hold for arbitrary $i \in [n]$, we can derive that

$$\mathbb{E}[F(\bar{x}) - F(x_*)] = \frac{1}{n} \sum_{i=1}^n \mathbb{E}[F_i(\bar{x}_i) - F_i(x_{i,*})] > \epsilon.$$

Thus, to find an output \bar{x} that satisfies $\mathbb{E}[F(\bar{x}) - F(x_*)] \leq \epsilon$, the abstract scheme requires at least $T \geq \frac{n}{4B\rho} = \frac{nG_1^2\sigma_0^2}{64B\epsilon^2}$ iterations. \square

Critical: From the proof of the non-smooth case, we can see that even when the overall objective is strongly convex, the lower bound complexity is still $T = \Omega\left(\frac{n\sigma_0^2}{B\epsilon^2}\right)$ as long as f_i is non-smooth. This behavior contrasts with standard strongly stochastic optimization with an optimal complexity of $O(1/\epsilon)$ and highlights a fundamental challenge in solving compositional problems.

5.5 Stochastic Optimization of Compositional OCE

The goal of this section is to present and analyze stochastic algorithms for solving compositional OCE (COCE) risk minimization as introduced in Chapter 3. In particular, we consider the following abstract problem:

$$\min_{\mathbf{w} \in \mathbb{R}^d, \mathbf{v} \in \mathbb{R}^n} F(\mathbf{w}, \mathbf{v}) := \frac{1}{n} \sum_{i=1}^n F_i(\mathbf{w}, v_i), \quad (5.68)$$

where

$$F_i(\mathbf{w}, v_i) = \mathbb{E}_{\zeta \sim \mathbb{P}_i} [\Phi_i(\mathbf{w}, v_i; \zeta)], \quad \Phi_i(\mathbf{w}, v_i; \zeta) = \tau \phi^* \left(\frac{s_i(\mathbf{w}; \zeta) - v_i}{\tau} \right) + v_i,$$

where $\tau > 0$ is a constant.

In the special case when $\phi^*(\cdot) = [\cdot]_+/\alpha$ for some $\alpha \in (0, 1)$, the general COCE minimization problem reduces to

$$\min_{\mathbf{w}, \mathbf{v}} F(\mathbf{w}, \mathbf{v}) := \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{\zeta \sim \mathbb{P}_i} \frac{[s_i(\mathbf{w}; \zeta) - v_i]_+}{\alpha} + v_i. \quad (5.69)$$

We refer to this problem as the **compositional CVaR minimization (CCVaR)** problem. The direct one-way partial AUC optimization problem (2.39) can be reformulated as an instance of CCVaR minimization as shown in (6.26).

In the special case when $\phi^*(\cdot) = \exp(\cdot) - 1$, the problem (5.68) reduces to

$$\min_{\mathbf{w}} F(\mathbf{w}) := \frac{1}{n} \sum_{i=1}^n \tau \log \left(\mathbb{E}_{\zeta \sim \mathbb{P}_i} \exp \left(\frac{s_i(\mathbf{w}; \zeta)}{\tau} \right) \right). \quad (5.70)$$

Algorithm 20 The ASGD Algorithm for solving (5.68)

```
1: Initialize  $\mathbf{w}_0, \mathbf{v}_0$ , step sizes  $\eta_t$  and  $\gamma_t$ 
2: for  $t = 0, 1, \dots, T - 1$  do
3:   Sample  $\mathcal{B}_t \subset \{1, \dots, n\}$  and  $|\mathcal{B}_t| = B$ 
4:   for each  $i \in \mathcal{B}_t$  do
5:     Update  $\mathbf{v}_{i,t+1} = \mathbf{v}_{i,t} - \gamma_t \partial_2 \Phi_i(\mathbf{w}_t, \mathbf{v}_{i,t}; \zeta_{i,t})$ 
6:   end for
7:   Compute  $\mathbf{z}_t = \frac{1}{B} \sum_{i \in \mathcal{B}_t} \partial_1 \Phi_i(\mathbf{w}_t, \mathbf{v}_{i,t}; \zeta_{i,t})$ 
8:   Update  $\mathbf{w}_{t+1} = \mathbf{w}_t - \eta_t \mathbf{z}_t$ 
9: end for
```

We refer to this problem as the **compositional entropic risk minimization (CERM)** problem. The cross-entropy loss for multi-class classification, the listwise cross-entropy loss for ranking, the indirect one-way partial AUC loss for imbalanced classification, and the contrastive losses for representation learning discussed in Chapter 2 are all instances of the CERM problem. In particular, for cross-entropy loss minimization, the proposed framework becomes especially relevant when the number of classes is very large, so that the normalization term in the loss cannot be computed efficiently. This setting naturally motivates the stochastic algorithms developed in this section.

Although we can cast the CERM problem into a special instance of FCCO, there remain some gaps to be filled. (i) For the convex CERM problem with a convex loss function $s_i(\cdot; \zeta)$, the ALEXR algorithm and its convergence analysis are not directly applicable, since the outer function $f(\cdot) = \tau \log(\cdot)$ is *not* convex, as required by ALEXR. Consequently, a convergence rate of $O(1/\epsilon^2)$ for solving convex CERM remains to be developed. (ii) For the CCVaR problem, the optimal solution of \mathbf{v} given \mathbf{w} is typically difficult to derive in closed form, and hence the problem cannot be reduced to an instance of FCCO. As a result, previous analyses for FCCO do not directly apply. We address these gaps in this section.

5.5.1 A Basic Algorithm

For optimizing the general COCE minimization problem, we present a basic stochastic algorithm in Algorithm 20. It alternates the stochastic block-coordinate update for \mathbf{v} and a SGD update for \mathbf{w} , which is referred to as ASGD. Below, we present its convergence analysis for both convex and non-convex loss functions.

5.5.1.1 Convex loss

For notational simplicity, we set $\tau = 1$ throughout the analysis.

Assumption 5.13. $s_i(\cdot, \zeta)$ is a convex function.

Lemma 5.20 $F(\mathbf{w}, \mathbf{v})$ is jointly convex in terms of $(\mathbf{w}^\top, \mathbf{v}^\top)^\top$ if $s_i(\cdot; \zeta)$ is convex.

Proof. We prove that $\Phi_i(\mathbf{w}, v_i; \zeta)$ is jointly convex in terms of $(\mathbf{w}^\top, v_i)^\top$. Then the convexity of $F(\mathbf{w}, \mathbf{v})$ follows. Let $\mathbf{u} = (\mathbf{w}^\top, v)^\top$. Consider $\mathbf{u}_1, \mathbf{u}_2$, $\alpha \in [0, 1]$, and $\bar{\mathbf{u}} = \alpha \mathbf{u}_1 + (1 - \alpha) \mathbf{u}_2$. Then

$$\Phi_i(\bar{\mathbf{u}}; \zeta) = \phi^*(s_i(\bar{\mathbf{w}}; \zeta) - \bar{v}) + \bar{v}.$$

If $s_i(\cdot; \zeta)$ is convex, we have $s_i(\bar{\mathbf{w}}; \zeta) \leq \alpha s_i(\mathbf{w}_1; \zeta) + (1 - \alpha) s_i(\mathbf{w}_2; \zeta)$. Since $\phi^*(\cdot)$ is non-decreasing (cf. Lemma 2.3), we have

$$\phi^*(s_i(\bar{\mathbf{w}}; \zeta) - \bar{v}) \leq \phi^*(\alpha(s_i(\mathbf{w}_1; \zeta) - v_1) + (1 - \alpha)(s_i(\mathbf{w}_2; \zeta) - v_2)).$$

Since $\phi^*(\cdot)$ is convex, we further have

$$\begin{aligned} & \phi^*(\alpha(s_i(\mathbf{w}_1; \zeta) - v_1) + (1 - \alpha)(s_i(\mathbf{w}_2; \zeta) - v_2)) \\ & \leq \alpha \phi^*(s_i(\mathbf{w}_1; \zeta) - v_1) + (1 - \alpha) \phi^*(s_i(\mathbf{w}_2; \zeta) - v_2). \end{aligned}$$

As a result,

$$\Phi_i(\bar{\mathbf{u}}; \zeta) \leq \alpha \Phi_i(\mathbf{u}_1; \zeta) + (1 - \alpha) \Phi_i(\mathbf{u}_2; \zeta),$$

which proves the convexity of $\Phi_i(\mathbf{u}; \zeta)$. □

Assumption 5.14. Assume that either of the following conditions hold:

- (i) $F(\mathbf{w}, \mathbf{v})$ is smooth satisfying:

$$\|\nabla_1 F(\mathbf{w}, \mathbf{v})\|_2^2 + \|\nabla_2 F(\mathbf{w}, \mathbf{v})\|_2^2 \leq 2L_F(F(\mathbf{w}, \mathbf{v}) - F(\mathbf{w}_*, \mathbf{v}_*)),$$

- (ii) $F(\mathbf{w}, \mathbf{v})$ non-smooth such that for any $\mathbf{v}_1 \in \partial_1 F(\mathbf{w}, \mathbf{v})$, $\mathbf{v}_{2,i} \in \partial_2 F_i(\mathbf{w}, v_i)$ it holds

$$\|\mathbf{v}_1\|_2^2 \leq G_1^2, \quad |\mathbf{v}_{2,i}|^2 \leq G_2^2,$$

where $\mathbf{w}_*, \mathbf{v}_*$ denotes an optimal solution to (5.68), and $\nabla_1 F(\mathbf{w}, \mathbf{v})$ ($\partial_1 F(\mathbf{w}, \mathbf{v})$), and $\nabla_2 F(\mathbf{w}, \mathbf{v})$ ($\partial_2 F(\mathbf{w}, \mathbf{v})$) denote (partial) gradients with respect to \mathbf{w}, \mathbf{v} , respectively.

Critical: For CERM, the smoothness assumption is satisfied when $s_i(\mathbf{w}; \zeta)$ is bounded, Lipschitz, and smooth. For CCVaR, the non-smoothness assumption is satisfied when $s_i(\mathbf{w}; \zeta)$ is bounded and Lipschitz.

Assumption 5.15 (Bounded Variance). There exist $\sigma_1^2, \sigma_2^2, \delta^2$ such that

$$\begin{aligned}
\mathbb{E}_\zeta \|\nabla_1 \Phi_i(\mathbf{w}, \nu_i; \zeta) - \nabla_1 F_i(\mathbf{w}, \nu_i)\|_2^2 &\leq \sigma_1^2, \quad \forall i \in [n], \\
\mathbb{E}_\zeta \|\nabla_2 \Phi_i(\mathbf{w}, \nu_i; \zeta) - \nabla_2 F_i(\mathbf{w}, \nu_i)\|_2^2 &\leq \sigma_2^2, \quad \forall i \in [n], \\
\frac{1}{n} \sum_{i=1}^n \|\nabla_1 F_i(\mathbf{w}, \nu_i) - \nabla_1 F(\mathbf{w}, \nu)\|_2^2 &\leq \delta^2.
\end{aligned}$$

In the non-smooth case, the gradients above are replaced by subgradients. The subsequent analysis proceeds analogously.

Lemma 5.21 *Let $D_{\mathbf{w},0}^2 := \mathbb{E}\|\mathbf{w}_0 - \mathbf{w}_*\|_2^2$ and $\eta_t = \eta$, we have*

$$\frac{1}{T} \sum_{t=0}^{T-1} (2\mathbb{E}[\nabla_1 F(\mathbf{w}_t, \nu_t)^\top (\mathbf{w}_t - \mathbf{w}_*)] - \eta \mathbb{E}\|\nabla_1 F(\mathbf{w}_t, \nu_t)\|_2^2) \leq \frac{D_{\mathbf{w},0}^2}{\eta T} + \eta \sigma^2.$$

where $\nu_t = (\nu_{1,t}, \dots, \nu_{n,t})^\top$ and $\sigma^2 = \frac{\sigma_1^2}{B} + \frac{\delta^2(n-B)}{B(n-1)}$.

Proof. Let \mathbb{E}_t denote the expectation over the random samples in the t -th iteration. First, we note that $\mathbb{E}_t[\mathbf{z}_t] = \nabla_1 F(\mathbf{w}_t, \nu_t)$. Similar to Lemma 5.2, we have

$$\begin{aligned}
&\mathbb{E}_t \|\mathbf{z}_t - \nabla_1 F(\mathbf{w}_t, \nu_t)\|_2^2 \\
&= \mathbb{E}_t \left\| \mathbf{z}_t - \frac{1}{B} \sum_{i \in \mathcal{B}_t} \partial_1 F_i(\mathbf{w}_t, \nu_{i,t}) + \frac{1}{B} \sum_{i \in \mathcal{B}_t} \partial_1 F_i(\mathbf{w}_t, \nu_{i,t}) - \nabla_1 F(\mathbf{w}_t, \nu_t) \right\|_2^2 \\
&= \mathbb{E}_t \left\| \mathbf{z}_t - \frac{1}{B} \sum_{i \in \mathcal{B}_t} \partial_1 F_i(\mathbf{w}_t, \nu_{i,t}) \right\|_2^2 + \mathbb{E}_t \left\| \frac{1}{B} \sum_{i \in \mathcal{B}_t} \partial_1 F_i(\mathbf{w}_t, \nu_{i,t}) - \nabla_1 F(\mathbf{w}_t, \nu_t) \right\|_2^2 \\
&\leq \frac{\sigma_1^2}{B} + \frac{\delta^2(n-B)}{B(n-1)} := \sigma^2.
\end{aligned}$$

Due to the update of \mathbf{w} , we have

$$\|\mathbf{w}_{t+1} - \mathbf{w}_*\|_2^2 = \|\mathbf{w}_t - \mathbf{w}_*\|_2^2 - 2\eta \mathbf{z}_t^\top (\mathbf{w}_t - \mathbf{w}_*) + \eta^2 \|\mathbf{z}_t\|_2^2.$$

Then,

$$\begin{aligned}
\mathbb{E}\|\mathbf{w}_{t+1} - \mathbf{w}_*\|_2^2 &\leq \mathbb{E}\|\mathbf{w}_t - \mathbf{w}_*\|_2^2 - 2\eta \mathbb{E}[\nabla_1 F(\mathbf{w}_t, \nu_t)^\top (\mathbf{w}_t - \mathbf{w}_*)] \\
&\quad + \eta^2 \mathbb{E}\|\nabla_1 F(\mathbf{w}_t, \nu_t)\|_2^2 + \eta^2 \sigma^2.
\end{aligned} \tag{5.71}$$

Summing over $t = 0, \dots, T-1$ and rearranging it finishes the proof. \square

Lemma 5.22 *Let $D_{\nu,0}^2 := \mathbb{E}\|\nu_0 - \nu_*\|_2^2$ and $\gamma_t = \gamma$, we have*

$$\frac{1}{T} \sum_{t=0}^{T-1} (2\mathbb{E}[\nabla_2 F(\mathbf{w}_t, \nu_t)^\top (\nu_t - \nu_*)] - \gamma n \mathbb{E}\|\nabla_2 F(\mathbf{w}_t, \nu_t)\|_2^2) \leq \frac{D_{\nu,0}^2}{\gamma B T} + \gamma \sigma_2^2.$$

Proof. Let \mathbb{E}_t denote the expectation over the random samples in the t -th iteration. Note that $\mathbb{E}_t[\nabla_2\Phi_i(\mathbf{w}_t, v_{i,t}; \zeta_{i,t})] = \nabla_2F_i(\mathbf{w}_t, v_{i,t})$ and $\mathbb{E}_t\|\nabla_2\Phi_i(\mathbf{w}_t, v_{i,t}; \zeta_{i,t}) - \nabla_2F_i(\mathbf{w}_t, v_{i,t})\|_2^2 \leq \sigma_0^2$ for each $i \in [n]$ (For those $i \notin \mathcal{B}_t$, $\nabla_2\Phi_i(\mathbf{w}_t, v_{i,t}; \zeta_{i,t})$ are not explicitly computed). For each $i \in [n]$, we have

$$\begin{aligned} & \mathbb{E}\|v_{i,t+1} - v_{i,*}\|_2^2 \\ &= (1 - \frac{B}{n})\mathbb{E}\|v_{i,t} - v_{i,*}\|_2^2 + \frac{B}{n}\mathbb{E}\|v_{i,t} - \gamma\nabla_2\Phi_i(\mathbf{w}_t, v_{i,t}; \zeta_{i,t}) - v_{i,*}\|_2^2 \\ &\leq \mathbb{E}\|v_{i,t} - v_{i,*}\|_2^2 - \frac{2\gamma B}{n}\mathbb{E}[\nabla_2F_i(\mathbf{w}_t, v_{i,t})^\top (v_{i,t} - v_{i,*})] + \frac{\gamma^2 B}{n}\mathbb{E}\|\nabla_2F_i(\mathbf{w}_t, v_{i,t})\|_2^2 \\ &\quad + \frac{\gamma^2 \sigma_2^2 B}{n}. \end{aligned}$$

Summing over $i \in [n]$ leads to

$$\begin{aligned} \mathbb{E}\|v_{t+1} - v_*\|_2^2 &= \mathbb{E}\|v_t - v_*\|_2^2 - \frac{2\gamma B}{n}\mathbb{E}\left[\sum_{i=1}^n \nabla_2F_i(\mathbf{w}_t, v_{i,t})^\top (v_{i,t} - v_{i,*})\right] \\ &\quad + \frac{\gamma^2 B}{n}\mathbb{E}\left[\sum_{i=1}^n \|\nabla_2F_i(\mathbf{w}_t, v_{i,t})\|_2^2\right] + \gamma^2 \sigma_2^2 B. \end{aligned} \quad (5.72)$$

Since

$$\begin{aligned} \nabla_2F(\mathbf{w}_t, v_t)^\top (v_t - v_*) &= \frac{1}{n} \sum_{i=1}^n \nabla_2F_i(\mathbf{w}_t, v_{i,t})(v_{i,t} - v_{i,*}) \\ \|\nabla_2F(\mathbf{w}_t, v_t)\|_2^2 &= \frac{1}{n^2} \sum_{i=1}^n \|\nabla_2F_i(\mathbf{w}_t, v_{i,t})\|_2^2, \end{aligned}$$

plugging these into (5.73) we have

$$\begin{aligned} \mathbb{E}\|v_{t+1} - v_*\|_2^2 &\leq \mathbb{E}\|v_t - v_*\|_2^2 - 2\gamma B\mathbb{E}\left[\nabla_2F(\mathbf{w}_t, v_t)^\top (v_t - v_*)\right] \\ &\quad + \gamma^2 n B\mathbb{E}\left[\|\nabla_2F(\mathbf{w}_t, v_t)\|_2^2\right] + \gamma^2 \sigma_2^2 B. \end{aligned} \quad (5.73)$$

Summing over $t = 0, \dots, T-1$ and rearranging it finishes the proof. \square

Theorem 5.11 (Smooth case) Suppose Assumption 5.13, 5.14(i) and 5.15 hold. If we set $\gamma = \min\{\frac{1}{2nL_F}, \frac{\epsilon}{2\sigma_2^2}\}$, $\eta = \min\{\frac{1}{2L_F}, \frac{\epsilon}{2\sigma_2^2}\}$ and $T = \max(\frac{2D_{\mathbf{w},0}^2}{\eta\epsilon}, \frac{2D_{v,0}^2}{\gamma B\epsilon})$, then ASGD guarantees that

$$\mathbb{E}\left[\frac{1}{T} \sum_{t=0}^{T-1} (F(\mathbf{w}_t, v_t) - F(\mathbf{w}_*, v_*))\right] \leq \epsilon.$$

The iteration complexity is

$$T = O \left(\max \left\{ \frac{D_{\mathbf{w},0}^2 L_F}{\epsilon}, \frac{n D_{\mathbf{v},0}^2 L_F}{B \epsilon}, \frac{D_{\mathbf{w},0}^2 \sigma_1^2}{\epsilon^2}, \frac{D_{\mathbf{v},0}^2 \sigma_2^2}{B \epsilon^2} \right\} \right),$$

where $\sigma^2 = \frac{\sigma_1^2}{B} + \frac{\delta^2(n-B)}{B(n-1)}$.

Proof. From Lemma 5.21 and Lemma 5.22, we have

$$\begin{aligned} \frac{1}{T} \sum_{t=0}^{T-1} (2\mathbb{E} \nabla_1 F(\mathbf{w}_t, \mathbf{v}_t)^\top (\mathbf{w}_t - \mathbf{w}_*) - \eta \mathbb{E} \|\nabla_1 F(\mathbf{w}_t, \mathbf{v}_t)\|_2^2) &\leq \frac{D_{\mathbf{w},0}^2}{\eta T} + \eta \sigma^2, \\ \frac{1}{T} \sum_{t=0}^{T-1} (2\mathbb{E} \nabla_2 F(\mathbf{w}_t, \mathbf{v}_t)^\top (\mathbf{v}_t - \mathbf{v}_*) - \gamma n \mathbb{E} \|\nabla_2 F(\mathbf{w}_t, \mathbf{v}_t)\|_2^2) &\leq \frac{D_{\mathbf{v},0}^2}{\gamma B T} + \gamma \sigma_2^2. \end{aligned}$$

If F is smooth and $\eta \leq \frac{1}{2L_F}$ and $\gamma n \leq \frac{1}{2L_F}$,

$$\begin{aligned} \eta \|\nabla_1 F(\mathbf{w}_t, \mathbf{v}_t)\|_2^2 + \gamma n \|\nabla_2 F(\mathbf{w}_t, \mathbf{v}_t)\|_2^2 &\leq \frac{1}{2L_F} \left(\|\nabla_1 F(\mathbf{w}, \mathbf{v})\|_2^2 + \|\nabla_2 F(\mathbf{w}, \mathbf{v})\|_2^2 \right) \\ &\leq F(\mathbf{w}_t, \mathbf{v}_t) - F(\mathbf{w}_*, \mathbf{v}_*), \end{aligned}$$

where the last inequality uses the Lemma 1.5(b).

On the other hand, the joint convexity of F implies

$$F(\mathbf{w}_t, \mathbf{v}_t) - F(\mathbf{w}_*, \mathbf{v}_*) \leq \nabla_1 F(\mathbf{w}_t, \mathbf{v}_t)^\top (\mathbf{w}_t - \mathbf{w}_*) + \nabla_2 F(\mathbf{w}_t, \mathbf{v}_t)^\top (\mathbf{v}_t - \mathbf{v}_*).$$

Then combining the above inequalities, we have

$$\mathbb{E} \left[\frac{1}{T} \sum_{t=0}^{T-1} [F(\mathbf{w}_t, \mathbf{v}_t) - F(\mathbf{w}_*, \mathbf{v}_*)] \right] \leq \frac{D_{\mathbf{w},0}^2}{2\eta T} + \frac{\eta \sigma^2}{2} + \frac{D_{\mathbf{v},0}^2}{2\gamma B T} + \frac{\gamma \sigma_2^2}{2}.$$

In order to let the RHS above be less than ϵ , we set $\gamma = \min\{\frac{1}{2nL_F}, \frac{\epsilon}{2\sigma_2^2}\}$ and $\eta = \min\{\frac{1}{2L_F}, \frac{\epsilon}{2\sigma^2}\}$, and $T \geq \max(\frac{2D_{\mathbf{w},0}^2}{\eta\epsilon}, \frac{2D_{\mathbf{v},0}^2}{\gamma B\epsilon})$. As a result, the complexity is the in the order of

$$T = O \left(\max \left\{ \frac{D_{\mathbf{w},0}^2 L_F}{\epsilon}, \frac{n D_{\mathbf{v},0}^2 L_F}{B \epsilon}, \frac{D_{\mathbf{w},0}^2 \sigma^2}{\epsilon^2}, \frac{D_{\mathbf{v},0}^2 \sigma_2^2}{B \epsilon^2} \right\} \right).$$

□

Theorem 5.12 (Non-smooth case) Suppose Assumption 5.13, 5.14(ii) and 5.15 hold. If we set $\gamma = \frac{\epsilon}{2(G_2^2 + \sigma_2^2)}$, $\eta = \frac{\epsilon}{2(G_1^2 + \sigma^2)}$ and $T = \max(\frac{2D_{\mathbf{w},0}^2}{\eta\epsilon}, \frac{2D_{\mathbf{v},0}^2}{\gamma B\epsilon})$, then ASGD guarantees that

$$\mathbb{E} \left[\frac{1}{T} \sum_{t=0}^{T-1} (F(\mathbf{w}_t, \mathbf{v}_t) - F(\mathbf{w}_*, \mathbf{v}_*)) \right] \leq \epsilon.$$

The iteration complexity is

$$T = O \left(\max \left\{ \frac{D_{\mathbf{w},0}^2(G_1^2 + \sigma^2)}{\epsilon^2}, \frac{D_{\mathbf{v},0}^2(G_2^2 + \sigma_2^2)}{B\epsilon^2} \right\} \right).$$

We leave the proof as an exercise for the reader.

💡 Why it matters

Since $F(\mathbf{w}, \mathbf{v})$ is jointly convex in (\mathbf{w}, \mathbf{v}) , the above two theorems imply convergence of the objective with respect to the primary variable \mathbf{w} , i.e., $F_1(\mathbf{w}) = \min_{\mathbf{v}} F(\mathbf{w}, \mathbf{v})$. In particular, if we define the averaged iterate $\bar{\mathbf{w}}_T = \frac{1}{T} \sum_{t=0}^{T-1} \mathbf{w}_t$, we have

$$\begin{aligned} \mathbb{E}[F_1(\bar{\mathbf{w}}_T) - F_1(\mathbf{w}_*)] &\leq \mathbb{E} \left[\frac{1}{T} \sum_{t=0}^{T-1} (F_1(\mathbf{w}_t) - F_1(\mathbf{w}_*)) \right] \\ &\leq \mathbb{E} \left[\frac{1}{T} \sum_{t=0}^{T-1} (F(\mathbf{w}_t, \mathbf{v}_t) - F(\mathbf{w}_*, \mathbf{v}_*)) \right] \leq \epsilon. \end{aligned}$$

5.5.1.2 Non-convex loss

If $s_i(\mathbf{w}, \zeta)$ is non-convex, we consider two different cases: (1) smooth case and (2) non-smooth weakly convex case. If $F(\mathbf{w}, \mathbf{v})$ is smooth in terms of \mathbf{w}, \mathbf{v} and is strongly convex in terms of \mathbf{v} (e.g, compositional entropic risk or COCE with χ^2 divergence for $\phi(\cdot)$), we can follow the analysis in Chapter 4 [Section 4.5] to design an algorithm and an analysis to prove the convergence for finding an ϵ -stationary point of $F_1(\mathbf{w}) = \min_{\mathbf{v}} F(\mathbf{w}, \mathbf{v})$. We leave this as an exercise for the reader.

Below, we analyze the convergence of ASGD for non-smooth weakly convex losses. We also assume ϕ^* is non-smooth such that it covers the CCVaR minimization.

Assumption 5.16. Suppose the following conditions hold:

- $s_i(\mathbf{w}; \zeta)$ is ρ_0 -weakly convex with respect to \mathbf{w} , and $\mathbb{E}_{\zeta} [\|\partial s_i(\mathbf{w}; \zeta)\|_2^2] \leq G_{\ell}^2$;
- Assume $|\frac{\partial \phi^*(q)}{\partial q}| \leq G_0$ for any $q = s_i(\mathbf{w}, \zeta) - v_i$.

Lemma 5.23 $F(\mathbf{w}, \mathbf{v})$ is ρ -weakly convex with respect to (\mathbf{w}, \mathbf{v}) , where $\rho = \rho_0 G_0$.

Proof. We first prove that $\phi^*(s_i(\mathbf{w}; \zeta) - v_i)$ is weakly convex in terms of (\mathbf{w}, v_i) , i.e. there exists $\rho > 0$ such that $\phi^*(s_i(\mathbf{w}; \zeta) - v_i) + \frac{\rho}{2} \|\mathbf{w}\|_2^2 + \frac{\rho}{2} v_i^2$ is jointly convex in terms of \mathbf{w}, v_i .

Since $s_i(\mathbf{w}; \zeta)$ is ρ_0 -weakly convex, we have that $q(\mathbf{w}, v_i, \zeta) = s_i(\mathbf{w}, \zeta) - v_i$ is ρ_0 -weakly convex in terms of $\mathbf{v}_i = (\mathbf{w}, v_i)$:

$$q(\mathbf{v}_i, \zeta) \geq q(\mathbf{v}'_i, \zeta) + \partial q(\mathbf{v}'_i, \zeta)^\top (\mathbf{v}_i - \mathbf{v}'_i) - \frac{\rho_0}{2} \|\mathbf{v}'_i - \mathbf{v}_i\|_2^2, \forall \mathbf{v}_i, \mathbf{v}'_i.$$

For any ζ , we abbreviate $q(\mathbf{v}_i; \zeta)$ as $q(\mathbf{v}_i)$. Since ϕ^* is convex and monotonically non-decreasing, for any $\omega \in \partial\phi^*(q(\mathbf{v}'_i)) \in [0, G_0]$ we have

$$\begin{aligned}\phi^*(q(\mathbf{v}_i)) &\geq \phi^*(q(\mathbf{v}'_i)) + \omega(q(\mathbf{v}_i) - q(\mathbf{v}'_i)) \\ &\geq \phi^*(q(\mathbf{v}'_i)) + \omega(\partial q(\mathbf{v}'_i))^\top (\mathbf{v}_i - \mathbf{v}'_i) - \frac{\rho_0}{2} \|\mathbf{v}_i - \mathbf{v}'_i\|_2^2 \\ &\geq \phi^*(q(\mathbf{v}'_i)) + \partial\phi^*(q(\mathbf{v}'_i))^\top (\mathbf{v}_i - \mathbf{v}'_i) - \frac{G_0\rho_0}{2} \|\mathbf{v}_i - \mathbf{v}'_i\|_2^2.\end{aligned}$$

The above inequality implies that $\phi^*(s_i(\mathbf{w}; \zeta) - \nu_i)$ is $\rho = G_0\rho_0$ -weakly convex in terms of (\mathbf{w}, ν_i) , i.e., $\mathbb{E}_\zeta \phi^*(s_i(\mathbf{w}; \zeta) - \nu_i) + \frac{\rho}{2} (\|\mathbf{w}\|_2^2 + |\nu_i|^2)$ is convex. As a result $F(\mathbf{w}, \nu) + \frac{\rho}{2} \|\mathbf{w}\|_2^2 + \frac{\rho}{2} \|\nu\|_2^2$ is jointly convex in terms of (\mathbf{w}, ν) . \square

Similar to the SGD for weakly convex objectives in Chapter 3[Section 3.1.4], we use the Moreau envelope of $F(\mathbf{w}; \nu)$. In particular, let $\mathbf{v} = (\mathbf{w}^\top, \nu^\top)^\top$ and consider some $\bar{\rho} > \rho$, we define:

$$F_{1/\bar{\rho}}(\mathbf{v}) = \min_{\mathbf{u}} F(\mathbf{u}) + \frac{\bar{\rho}}{2} \|\mathbf{u} - \mathbf{v}\|_2^2, \quad (5.74)$$

$$\text{prox}_{F/\bar{\rho}}(\mathbf{v}) := \arg \min_{\mathbf{u}} F(\mathbf{u}) + \frac{\bar{\rho}}{2} \|\mathbf{u} - \mathbf{v}\|_2^2. \quad (5.75)$$

Convergence Analysis

Lemma 5.24 *Under Assumption 5.16, we have*

$$\mathbb{E}_t [\|\mathbf{z}_t\|_2^2] \leq G_1^2, \quad |\partial_2 F_i(\mathbf{w}, \nu_i)|^2 \leq G_2^2,$$

where $G_1^2 = G_0^2 G_\ell^2$, and $G_2^2 = (1 + G_0)^2$.

Proof. For the first part,

$$\mathbb{E}_t [\|\mathbf{z}_t\|_2^2] = \mathbb{E}_t \left[\left\| \frac{1}{B} \sum_{i \in \mathcal{B}_t} \partial_1 \Phi_i(\mathbf{w}_t, \nu_{i,t}; \zeta_{i,t}) \right\|_2^2 \right] \leq G_0^2 G_\ell^2.$$

For the second part,

$$|\partial_2 F_i(\mathbf{w}, \nu_i)|^2 = \left| \mathbb{E}_\zeta \left[-\frac{\partial \phi^*(q(\mathbf{w}, \nu_i; \zeta))}{\partial q} + 1 \right] \right|^2 \leq (1 + G_0)^2.$$

\square

Lemma 5.25 *Under Assumption (5.16), let $\mathbf{v}_t = (\mathbf{w}_t^\top, \nu_t^\top)^\top$, for one iteration of ASGD, we have*

$$\begin{aligned}\mathbb{E}_t[F_{1/\bar{\rho}}(\mathbf{v}_{t+1})] &\leq F_{1/\bar{\rho}}(\mathbf{v}_t) + \bar{\rho}\eta_t(F(\bar{\mathbf{v}}_t) - F(\mathbf{v}_t) + \frac{\bar{\rho}}{2}\|\mathbf{v}_t - \bar{\mathbf{v}}_t\|_2^2) \\ &\quad + \frac{\bar{\rho}\eta_t^2(G_1^2 + G_2^2/B)}{2},\end{aligned}$$

where $\bar{\mathbf{v}}_t = \text{prox}_{F/\bar{\rho}}(\mathbf{v}_t)$.

Proof. Let \mathbb{E}_t denote the expectation over the random samples at the t -th iteration conditioned on that in all previous iterations.

$$\begin{aligned}\mathbb{E}_t[F_{1/\bar{\rho}}(\mathbf{v}_{t+1})] &\leq \mathbb{E}_t\left[F(\bar{\mathbf{v}}_t) + \frac{\bar{\rho}}{2}\|\mathbf{v}_{t+1} - \bar{\mathbf{v}}_t\|_2^2\right] \\ &\leq F(\bar{\mathbf{v}}_t) + \frac{\bar{\rho}}{2}\mathbb{E}_t[\|\mathbf{w}_t - \eta_t\mathbf{z}_t - \bar{\mathbf{w}}_t\|_2^2 + \|\mathbf{v}_{t+1} - \bar{\mathbf{v}}_t\|_2^2] \\ &\leq F(\bar{\mathbf{v}}_t) + \frac{\bar{\rho}}{2}\mathbb{E}_t[\|\mathbf{w}_t - \eta_t\mathbf{z}_t - \bar{\mathbf{w}}_t\|_2^2] + \frac{\bar{\rho}}{2}\mathbb{E}_t[\|\mathbf{v}_{t+1} - \bar{\mathbf{v}}_t\|_2^2] \\ &\leq F(\bar{\mathbf{v}}_t) + \frac{\bar{\rho}}{2}\|\mathbf{w}_t - \bar{\mathbf{w}}_t\|_2^2 + \bar{\rho}\eta_t\mathbb{E}_t[(\bar{\mathbf{w}}_t - \mathbf{w}_t)^\top \partial_1 F(\mathbf{w}_t, \mathbf{v}_t)] + \frac{\bar{\rho}\eta_t^2 G_1^2}{2} \\ &\quad + \frac{\bar{\rho}}{2}\mathbb{E}_t[\|\mathbf{v}_{t+1} - \bar{\mathbf{v}}_t\|_2^2]\end{aligned}$$

where the last step uses $\mathbb{E}_t[\mathbf{z}_t] = \partial_1 F(\mathbf{w}_t, \mathbf{v}_t)$ and $\mathbb{E}[\|\mathbf{z}_t\|_2^2] \leq G_1^2$.

Similar to (5.73), we can prove that

$$\mathbb{E}_t\|\mathbf{v}_{t+1} - \bar{\mathbf{v}}_t\|_2^2 = \|\mathbf{v}_t - \bar{\mathbf{v}}_t\|_2^2 - 2\gamma_t B \partial_2 F(\mathbf{w}_t, \mathbf{v}_t)^\top (\mathbf{v}_t - \bar{\mathbf{v}}_t) + \gamma_t^2 G_2^2 B.$$

Let $\gamma_t B = \eta_t$, combining the above we have

$$\begin{aligned}\mathbb{E}_t[F_{1/\bar{\rho}}(\mathbf{v}_{t+1})] &\leq F(\bar{\mathbf{v}}_t) + \frac{\bar{\rho}}{2}\|\mathbf{v}_t - \bar{\mathbf{v}}_t\|_2^2 + \bar{\rho}\eta_t\mathbb{E}_t[(\bar{\mathbf{v}}_t - \mathbf{v}_t)^\top \partial F(\mathbf{v}_t)] + \frac{\bar{\rho}\eta_t^2(G_1^2 + G_2^2/B)}{2} \\ &\leq F(\bar{\mathbf{v}}_t) + \frac{\bar{\rho}}{2}\|\mathbf{v}_t - \bar{\mathbf{v}}_t\|_2^2 + \bar{\rho}\eta_t\mathbb{E}_t[(\bar{\mathbf{v}}_t - \mathbf{v}_t)^\top \partial F(\mathbf{v}_t)] + \frac{\bar{\rho}\eta_t^2(G_1^2 + G_2^2/B)}{2} \\ &\leq F_{1/\bar{\rho}}(\mathbf{v}_t) + \bar{\rho}\eta_t(F(\bar{\mathbf{v}}_t) - F(\mathbf{v}_t) + \frac{\bar{\rho}}{2}\|\mathbf{v}_t - \bar{\mathbf{v}}_t\|_2^2) + \frac{\bar{\rho}\eta_t^2(G_1^2 + G_2^2/B)}{2}.\end{aligned}$$

where the last step uses the definition of $F_{1/\bar{\rho}}(\mathbf{v}_t)$ and the ρ -weak convexity of F . Rearranging this inequality finishes the proof. \square

Theorem 5.13 Suppose Assumption (5.16) holds and $F_* = \inf F(\mathbf{w}, \mathbf{v}) \geq \infty$, by setting $\bar{\rho} = 2\rho$, $\eta = \epsilon^2/(2\bar{\rho}(G_1^2 + G_2^2/B))$, $\gamma = \eta/B$ and $T \geq \frac{4(F(\mathbf{w}_0, \mathbf{v}_0) - F_*)}{\epsilon^2\eta}$, ASGD guarantees that

$$\mathbb{E}\left[\frac{1}{T} \sum_{t=1}^T \|\nabla F_{1/\bar{\rho}}(\mathbf{v}_t)\|_2^2\right] \leq \epsilon^2$$

with a complexity of $T = O\left(\frac{\rho(G_1^2 + G_2^2/B)}{\epsilon^4}\right)$.

Proof. Since $F(\mathbf{v}) + \frac{\bar{\rho}}{2}\|\mathbf{v} - \mathbf{v}_t\|_2^2$ is $(\bar{\rho} - \rho)$ -strongly convex and have a minimum solution at $\bar{\mathbf{v}}_t$, then we have

$$\begin{aligned} & F(\mathbf{v}_t) - F(\bar{\mathbf{v}}_t) - \frac{\rho}{2}\|\mathbf{v}_t - \bar{\mathbf{v}}_t\|_2^2 \\ &= (F(\mathbf{v}_t) + \frac{\bar{\rho}}{2}\|\mathbf{v}_t - \mathbf{v}_t\|_2^2) - (F(\bar{\mathbf{v}}_t) + \frac{\bar{\rho}}{2}\|\bar{\mathbf{v}}_t - \mathbf{v}_t\|_2^2) + (\frac{\bar{\rho}}{2} - \frac{\rho}{2})\|\mathbf{v}_t - \bar{\mathbf{v}}_t\|_2^2 \\ &\geq \frac{(\bar{\rho} - \rho)}{2}\|\mathbf{v}_t - \bar{\mathbf{v}}_t\|_2^2 + \frac{(\bar{\rho} - \rho)}{2}\|\mathbf{v}_t - \bar{\mathbf{v}}_t\|_2^2 = (\bar{\rho} - \rho)\|\mathbf{v}_t - \bar{\mathbf{v}}_t\|_2^2 \\ &= \frac{\bar{\rho} - \rho}{\bar{\rho}^2}\|\nabla F_{1/\bar{\rho}}(\mathbf{v}_t)\|_2^2. \end{aligned}$$

Combining this result with that in Lemma 5.25 and noting that $\bar{\rho} = 2\rho, \eta_t = \eta$, we have

$$\begin{aligned} \mathbb{E}\left[\frac{1}{T} \sum_{t=0}^{T-1} \|\nabla F_{1/\bar{\rho}}(\mathbf{v}_t)\|_2^2\right] &\leq \frac{2(F_{1/\bar{\rho}}(\mathbf{v}_0) - F_*)}{\eta T} + \bar{\rho}\eta(G_1^2 + G_2^2/B) \\ &\leq \frac{2(F(\mathbf{v}_0) - F_*)}{\eta T} + \bar{\rho}\eta(G_1^2 + G_2^2/B) \end{aligned}$$

By setting $\eta = \epsilon^2/(2\bar{\rho}(G_1^2 + G_2^2/B))$ and $T \geq \frac{4(F(\mathbf{v}_0) - F_*)}{\epsilon^2\eta}$, we have $\mathbb{E}[\|\nabla F_{1/\bar{\rho}}(\mathbf{v}_\tau)\|_2^2] \leq \epsilon^2$ for a randomly selected $\tau \in \{0, \dots, T-1\}$. \square

5.5.2 A Geometry-aware Algorithm for Entropic Risk

Although last section presents a general algorithm for solving COCE risk minimization, it may exhibits numerical instability issue and slow convergence when solving compositional entropic risk minimization:

$$\begin{aligned} \min_{\mathbf{w}} \min_{\nu} \left[F(\mathbf{w}, \nu) &= \frac{1}{n} \sum_{i=1}^n \{\mathbb{E}_{\zeta} \exp(s_i(\mathbf{w}; \zeta) - \nu_i) - 1 + \nu_i\} \right] \\ &= \min_{\mathbf{w}} \frac{1}{n} \sum_{i=1}^n \log(\mathbb{E}_{\zeta} \exp(s_i(\mathbf{w}; \zeta))). \end{aligned}$$

The numerical instability issue is caused by dealing with exponential functions, e.g., $\exp(s_i(\mathbf{w}; \zeta) - \nu_i)$, in calculation of stochastic gradients of ν_i . The slow convergence arises because the standard SGD update for ν_i fails to exploit the geometric structure of the problem.

5.5.2.1 Stochastic Optimization of Log-E-Exp

We first consider a simplified problem where there is only one component $n = 1$, i.e.,

$$\min_{\mathbf{w}} F_1(\mathbf{w}) := \log(\mathbb{E}_{\zeta} \exp(s(\mathbf{w}; \zeta))) . \quad (5.76)$$

The KL-regularized DRO problem (2.14) is a special case. It is also known as log-E-Exp, a more general form of the log-Sum-Exp function, where the middle “E” denotes an expectation and highlights the associated computational challenges.

Application of SCGD

At the beginning of Section 4.1, we treat this problem as a special case of stochastic compositional optimization (SCO), where the outer function is $f(\cdot) = \log(\cdot)$ and the inner function is $g(\mathbf{w}) = \mathbb{E}_{\zeta}[\exp(s(\mathbf{w}; \zeta))]$. Let us first apply the SCGD algorithm. The key updates are presented below:

$$\begin{aligned} u_t &= (1 - \gamma_t)u_{t-1} + \gamma_t \exp(s(\mathbf{w}_t; \zeta_t)), \\ \mathbf{z}_t &= \frac{1}{u_t} \exp(s(\mathbf{w}_t; \zeta'_t)) \nabla s(\mathbf{w}_t; \zeta'_t), \\ \mathbf{w}_{t+1} &= \mathbf{w}_t - \eta_t \mathbf{z}_t, \end{aligned} \quad (5.77)$$

where u_t is an estimator of the inner function value $g(\mathbf{w}_t)$ and $\mathbf{z}_t = \nabla f(u_t) \nabla g(\mathbf{w}_t; \zeta'_t)$ is a gradient estimator of \mathbf{w}_t .

From a practitioner’s perspective, the algorithm can be readily implemented and applied to real applications. However, from a theoretical perspective, several open problems remain. In particular: (1) Can we establish an $O(1/\epsilon^2)$ convergence rate for this algorithm to find an ϵ -optimal solution when $s(\mathbf{w}; \zeta)$ is convex? (2) If yes, what are the practical advantages of this algorithm compared with the ASGD method presented in the previous section?

Wait! Shouldn’t we established the convergence rate of SCGD in Chapter 4? It is true that we presented a convergence analysis of the above algorithm for non-convex problems under proper conditions, however, it remains an open problem to establish the complexity of $O(1/\epsilon^2)$ for finding an ϵ -optimal solution under the convexity of $s(\mathbf{w}; \zeta)$. A naive analysis of SCGD for convex problems yields a complexity of $O(1/\epsilon^4)$ (see Wang et al. (2017a)).

A Novel Algorithm

To address these open questions, we present a novel algorithm based on the min-min reformulation of log-E-exp, i.e.,

$$\min_{\mathbf{w}} \min_{\nu} F(\mathbf{w}, \nu) := \mathbb{E}_{\zeta} \exp(s(\mathbf{w}; \zeta) - \nu) + \nu. \quad (5.78)$$

where we ignored the constant -1 in the objective. As proved in Lemma 5.20, $F(\mathbf{w}; \nu)$ is jointly convex in terms of \mathbf{w}, ν when $s(\mathbf{w}; \zeta)$ is convex.

Motivation

The key novelty of our design is a **geometry-aware algorithm** for solving the equivalent min-min optimization (5.78). Let us first discuss the motivation. One challenge for solving the min-min optimization problem is that the objective function $F(\mathbf{w}, \nu)$ could have exponentially large smoothness constant in terms of ν . We will formally analyze this phenomenon in next section. Hence, a vanilla gradient method that uses the first-order approximation of F will inevitably be impacted by the large smoothness parameter.

To mitigate the adverse effects of a large smoothness parameter with respect to ν , we resort to the classical approach of employing a proximal mapping. Proximal mappings have been widely used to handle a non-smooth function in composite objectives consisting of a smooth loss and a non-smooth regularizer. This approach enables optimization algorithms to retain the favorable convergence properties of smooth optimization and often leads to faster convergence despite the presence of non-smooth terms. Analogously, even when a function is smooth but characterized by a very large smoothness parameter, applying its proximal mapping can effectively alleviate the negative impact of this large smoothness constant.

However, there is an important distinction from classical proximal methods, which typically rely on direct access to the function of interest for computing the proximal mapping. In our setting, we cannot directly apply the proximal mapping of $F(\mathbf{w}, \nu)$. Instead, we only have access to a stochastic estimator

$$\Phi(\mathbf{w}, \nu; \zeta) = e^{s(\mathbf{w}; \zeta) - \nu} + \nu,$$

defined for a random sample ζ . As a result, it becomes necessary to explicitly account for the noise introduced by this stochastic approximation.

Algorithm

To account for the stochastic noise, we introduce a Bregman divergence $D_{\varphi}(\cdot, \cdot)$ and update ν_t according to the following scheme:

$$\nu_t = \arg \min_{\nu} \Phi(\mathbf{w}_t, \nu; \zeta_t) + \frac{1}{\alpha_t} D_{\varphi}(\nu, \nu_{t-1}), \quad (5.79)$$

where $\zeta_t \sim \mathbb{P}$ is a random sample and $\alpha_t > 0$ is a step size parameter. We refer to this step as **stochastic proximal mirror descent (SPMD)** update. To respect the geometry of the stochastic objective $\Phi(\mathbf{w}_t, \nu; \zeta_t)$, we construct a tailored Bregman divergence induced by the function $\varphi(\nu) = e^{-\nu}$, namely,

Algorithm 21 The SCENT Algorithm for solving Log-E-Exp (5.76)

```

1: Initialize  $\mathbf{w}_1, v_0$ , step sizes  $\eta_t$  and  $\alpha_t$ ,  $\varphi(v) = e^{-v}$ .
2: for  $t = 1 \dots T - 1$  do
3:   Sample  $\zeta_t, \zeta'_t$ 
4:   Update  $v_t = \arg \min_v \exp(s(\mathbf{w}_t; \zeta_t) - v) + v + \frac{1}{\alpha_t} D_\varphi(v, v_{t-1})$ 
5:   Compute  $\mathbf{z}_t = \exp(s(\mathbf{w}_t; \zeta'_t) - v_t) \nabla s(\mathbf{w}_t; \zeta'_t)$ 
6:   Compute  $\mathbf{v}_t = (1 - \beta_t) \mathbf{v}_{t-1} + \beta_t \mathbf{z}_t$ 
7:   Update  $\mathbf{w}_{t+1} = \mathbf{w}_t - \eta_t \mathbf{v}_t$ 
8: end for
    
```

$$D_\varphi(v, v_{t-1}) = e^{-v} - e^{-v_{t-1}} + e^{-v_{t-1}}(v - v_{t-1}). \quad (5.80)$$

Once we have v_t , we compute a vanilla gradient estimator by

$$\mathbf{z}_t = \exp(s(\mathbf{w}_t; \zeta'_t) - v_t) \nabla s(\mathbf{w}_t; \zeta'_t). \quad (5.81)$$

If the problem is non-convex, we compute a moving-average estimator $\mathbf{v}_t = (1 - \beta_t) \mathbf{v}_{t-1} + \beta_t \mathbf{z}_t$ and then update the model parameter \mathbf{w}_{t+1} . We present the full steps in Algorithm 21, which is referred to SCENT.

SCGD is just a special case of SCENT

To see the connection with SCGD, we present the following lemma.

Lemma 5.26 *The update of v_t defined by (5.79) can be computed by*

$$e^{v_t} = \frac{1}{1 + \alpha_t e^{v_{t-1}}} e^{v_{t-1}} + \frac{\alpha_t e^{v_{t-1}}}{1 + \alpha_t e^{v_{t-1}}} \exp(s(\mathbf{w}_t; \zeta_t)). \quad (5.82)$$

If $y_t = e^{-v_t}$, we have

$$y_t = \frac{y_{t-1} + \alpha_t}{1 + \alpha_t e^{s(\mathbf{w}_t; \zeta_t)}}.$$

Proof. We compute the gradient of the problem (5.79) and set it to zero for computing v_t , i.e.,

$$-\exp(s(\mathbf{w}_t; \zeta_t) - v_t) + 1 + \frac{1}{\alpha_t} (-e^{-v_t} + e^{-v_{t-1}}) = 0.$$

Solving this equation finishes the proof. \square

If we define $u_t = e^{v_t}$ and $\gamma'_t = \frac{\alpha_t e^{v_{t-1}}}{1 + \alpha_t e^{v_{t-1}}}$, then the updates of SCENT ($\beta_t = 1$) are equivalent to

$$\begin{aligned}
u_t &= (1 - \gamma'_t)u_{t-1} + \gamma'_t \exp(s(\mathbf{w}_t; \zeta_t)) \\
\mathbf{z}_t &= \frac{1}{u_t} \exp(s(\mathbf{w}_t; \zeta'_t)) \nabla s(\mathbf{w}_t; \zeta'_t), \\
\mathbf{w}_{t+1} &= \mathbf{w}_t - \eta_t \mathbf{z}_t.
\end{aligned} \tag{5.83}$$

Comparing this update with that of SCGD (5.77), the key difference lies in the choice of the moving-average parameter: SCENT adopts an adaptive parameter $\gamma'_t = \frac{\alpha_t e^{\nu_{t-1}}}{1 + \alpha_t e^{\nu_{t-1}}}$, whereas SCGD uses a non-adaptive γ_t . If we set $\alpha_t = \frac{\gamma_t}{1 - \gamma_t} e^{-\nu_{t-1}}$, then the updates of SCENT reduce to that of SCGD.

Convergence analysis for convex problems

Since $\mathbf{z}_t = \nabla_{\mathbf{w}} \exp(s(\mathbf{w}_t; \zeta'_t) - \nu_t)$, we have

$$\mathbb{E}_{\zeta'_t}[\mathbf{z}_t] = \nabla_{\mathbf{w}} \mathbb{E}_{\zeta'_t}[\exp(s(\mathbf{w}_t; \zeta'_t) - \nu_t)] = \nabla F(\mathbf{w}_t, \nu_t).$$

Let \mathbf{w}_*, ν_* be the optimal solution:

$$(\mathbf{w}_*, \nu_*) = \arg \min_{\mathbf{w}, \nu} F(\mathbf{w}, \nu).$$

It is straightforward to derive $\nu_* = \log[\mathbb{E} \exp(s(\mathbf{w}_*; \zeta))]$.

Assumption 5.17. Assume that the following conditions hold:

- (i) $s(\mathbf{w}; \zeta)$ is convex;
- (ii) the loss function is bounded such that $s(\mathbf{w}; \zeta) \in [c_0, c_1], \forall \mathbf{w}, \zeta$.
- (iii) there exists G such that $\mathbb{E}_{\zeta} \|\nabla s(\mathbf{w}_t, \zeta)\|_2^2 \leq G^2, \forall t$.

Critical: To relax the second assumption, we can assume that \mathbf{w} is restricted to a bounded domain \mathcal{W} and $s(\mathbf{w}; \zeta)$ is regular. In practice, we always enforce the boundness of \mathbf{w}_t through either projection onto \mathcal{W} or using a regularizer $r(\mathbf{w})$. The update of \mathbf{w}_{t+1} can be modified as the SPGD update:

$$\mathbf{w}_{t+1} = \arg \min_{\mathbf{w}} \mathbf{z}_t^\top \mathbf{w} + r(\mathbf{w}) + \frac{1}{\eta_t} \|\mathbf{w} - \mathbf{w}_t\|_2^2.$$

The analysis can be performed similarly.

Lemma 5.27 Under Assumption 5.17(ii), $\nu_* \in [c_0, c_1]$ and if $\nu_0 \in [c_0, c_1]$ then $\nu_t \in [c_0, c_1], \forall t$.

Proof. $\nu_* \in [c_0, c_1]$ can be seen from $\nu_* = \log[\mathbb{E} \exp(s(\mathbf{w}_*; \zeta))]$. The second result can be easily seen from the update of e^{ν_t} as in (5.82) by induction. \square

For the ease of analysis, we define two quantities to capture the variance terms caused by using stochastic estimators.

$$\begin{aligned}\sigma_t^2 &:= \mathbb{E}_{\zeta'_t} \|\exp(s(\mathbf{w}_t; \zeta'_t) - \nu_t) \nabla s(\mathbf{w}_t; \zeta'_t)\|_2^2, \\ \delta_t^2 &:= \mathbb{E}_{\zeta_t} [e^{-\nu_{t-1}} |e^{s(\mathbf{w}_t; \zeta_t)} - \mathbb{E}_{\zeta} [e^{s(\mathbf{w}_t; \zeta)}]|^2].\end{aligned}$$

Under Assumption 5.17 (ii) and (iii), σ_t, δ_t are bounded because $e^{\nu_t}, e^{\nu_{t-1}}$ and $e^{s(\mathbf{w}_t; \zeta_t)}$ is upper and lower bounded.

Critical: These two quantities are related to the variance of stochastic estimators in terms of \mathbf{w}_t and ν_t , respectively. Both quantities have a normalization term $e^{-\nu_t}$ or $e^{-\nu_{t-1}}$.

Lemma 5.28 *Under Assumption 5.17 and $\beta_t = 1$, we have*

$$\mathbb{E}[\eta_t \nabla_1 F(\mathbf{w}_t, \nu_t)^\top (\mathbf{w}_t - \mathbf{w}_*)] \leq \mathbb{E} \left[\frac{1}{2} \|\mathbf{w}_t - \mathbf{w}_*\|_2^2 - \frac{1}{2} \|\mathbf{w}_{t+1} - \mathbf{w}_*\|_2^2 \right] + \frac{\eta_t^2 \sigma_t^2}{2}.$$

Proof. The proof is a simple application of Lemma 3.3. \square

If the SPGD update is used, we can use Lemma 3.6 giving us

$$\begin{aligned}\mathbf{z}_t^\top (\mathbf{w}_{t+1} - \mathbf{w}_*) + r(\mathbf{w}_{t+1}) - r(\mathbf{w}_*) &\leq \frac{1}{2\eta_t} (\|\mathbf{w}_t - \mathbf{w}_*\|_2^2 - \|\mathbf{w}_{t+1} - \mathbf{w}_*\|_2^2) \\ &\quad - \frac{1}{2\eta_t} \|\mathbf{w}_t - \mathbf{w}_{t+1}\|_2^2.\end{aligned}$$

Then,

$$\begin{aligned}\mathbf{z}_t^\top (\mathbf{w}_t - \mathbf{w}_*) + r(\mathbf{w}_t) - r(\mathbf{w}_*) &\leq \frac{1}{2\eta_t} (\|\mathbf{w}_t - \mathbf{w}_*\|_2^2 - \|\mathbf{w}_{t+1} - \mathbf{w}_*\|_2^2) \\ &\quad + \mathbf{z}_t^\top (\mathbf{w}_t - \mathbf{w}_{t+1}) - \frac{1}{2\eta_t} \|\mathbf{w}_t - \mathbf{w}_{t+1}\|_2^2 + r(\mathbf{w}_t) - r(\mathbf{w}_{t+1}) \\ &\leq \frac{1}{2\eta_t} (\|\mathbf{w}_t - \mathbf{w}_*\|_2^2 - \|\mathbf{w}_{t+1} - \mathbf{w}_*\|_2^2) + \frac{\eta_t}{2} \|\mathbf{z}_t\|_2^2 + r(\mathbf{w}_t) - r(\mathbf{w}_{t+1}).\end{aligned}$$

Taking expectation on both sides, we have

$$\begin{aligned}&\mathbb{E}[\eta_t \nabla_1 F(\mathbf{w}_t, \nu_t)^\top (\mathbf{w}_t - \mathbf{w}_*)] + \eta_t (r(\mathbf{w}_t) - r(\mathbf{w}_*)) \\ &\leq \mathbb{E} \left[\left(\eta_t r(\mathbf{w}_t) + \frac{1}{2} \|\mathbf{w}_t - \mathbf{w}_*\|_2^2 \right) - \left(\eta_t r(\mathbf{w}_{t+1}) + \frac{1}{2} \|\mathbf{w}_{t+1} - \mathbf{w}_*\|_2^2 \right) \right] + \frac{\eta_t^2 \sigma_t^2}{2}.\end{aligned}$$

If $\eta_{t+1} \leq \eta_t$ and $r(\mathbf{w}) \geq 0$, then $\eta_t r(\mathbf{w}_{t+1}) \leq \eta_{t+1} r(\mathbf{w}_{t+1})$, then the terms in the square bracket will form a telescoping series over $t = 1, \dots, T$. As a result, the following analysis will proceed similarly.

Lemma 5.29 Under Assumption 5.17 (ii), we have

$$\alpha_t \nabla_2 \Phi(\mathbf{w}_t, v_t; \zeta_t)^\top (v_t - v_*) \leq D_\varphi(v_*, v_{t-1}) - D_\varphi(v_*, v_t) - D_\varphi(v_t, v_{t-1}).$$

Proof. Recall the definition

$$\begin{aligned} \Phi(\mathbf{w}_t, v; \zeta_t) &= \exp(s(\mathbf{w}_t; \zeta_t) - v) + v \\ \varphi(v) &= e^{-v}, \quad D_\varphi(a, b) = \varphi(a) - \varphi(b) - \langle \nabla \varphi(b), a - b \rangle, \end{aligned}$$

and the update of v_t :

$$v_t = \arg \min_v \alpha_t \Phi(\mathbf{w}_t, v; \zeta_t) + D_\varphi(v, v_{t-1}).$$

The first-order optimality gives

$$\alpha_t \nabla_2 \Phi(\mathbf{w}_t, v_t; \zeta_t) + \nabla \varphi(v_t) - \nabla \varphi(v_{t-1}) = 0.$$

Taking inner product with $(v_t - v_*)$ and rearranging gives

$$\begin{aligned} \alpha_t \langle \nabla_2 \Phi(\mathbf{w}_t, v_t; \zeta_t), v_t - v_* \rangle &= \langle \nabla \varphi(v_{t-1}) - \nabla \varphi(v_t), v_t - v_* \rangle \\ &= D_\varphi(v_*, v_{t-1}) - D_\varphi(v_*, v_t) - D_\varphi(v_t, v_{t-1}) \end{aligned}$$

where the last equality holds by three-point identity as in Lemma 3.9. \square

Critical: To proceed the analysis, we need to bound $\mathbb{E}[\alpha_t \nabla_2 F(\mathbf{w}_t, v_t)^\top (v_t - v_*)]$. In light of the above lemma, we will bound the following difference in expectation:

$$\mathbb{E}[(\nabla_2 F(\mathbf{w}_t, v_t)^\top (v_t - v_*) - \nabla_2 \Phi(\mathbf{w}_t, v_t; \zeta_t)^\top (v_t - v_*))].$$

The challenge lies at v_t depends on ζ_t , making the above expectation not equal to zero.

Lemma 5.30 Assume $\alpha_t \leq \rho e^{-v_{t-1}}$ for any constant $\rho > 0$, then we have

$$|\mathbb{E}[(\nabla_2 F(\mathbf{w}_t, v_t)^\top (v_t - v_*) - \nabla_2 \Phi(\mathbf{w}_t, v_t; \zeta_t)^\top (v_t - v_*))]| \leq \alpha_t \delta_t^2 C. \quad (5.84)$$

where $C = (1 + \rho)(1 + c_1 - c_0)$.

Proof. In the following proof, we let \mathcal{F}_{t-1} denote the filtration (the “information available”) up to time $t - 1$.

Let us define $z_t = e^{s(\mathbf{w}_t; \zeta_t)}$, $m_t = \mathbb{E}_\zeta[e^{s(\mathbf{w}_t; \zeta)} | \mathcal{F}_{t-1}]$, and $y_t = e^{-v_t}$. Let z and z' two independent variables so that $\mathbb{E}[z | \mathcal{F}_{t-1}] = \mathbb{E}[z' | \mathcal{F}_{t-1}] = m_t$. Since v_t depends on z_t , let us define random functions:

$$\begin{aligned} y_t(z) &= \frac{y_{t-1} + \alpha_t}{\alpha_t z + 1}, \quad v_t(z) = -\log y_t(z) \\ h_t(z) &= e^{-v_t(z)}(v_t(z) - v_*) = y_t(z)(v_t(z) - v_*). \end{aligned}$$

According to the update of v_t , we have $y_t = y_t(z_t)$, $v_t = v_t(z)$. For the target, we have

$$\begin{aligned} &\mathbb{E}[(\nabla_2 \Phi(\mathbf{w}_t, v_t; \zeta_t) - \nabla_2 F(\mathbf{w}_t, v_t))^\top (v_t - v_*) \mid \mathcal{F}_{t-1}] \\ &= \mathbb{E}[\mathbb{E}_\zeta[e^{s(\mathbf{w}_t; \zeta)}] - e^{s(\mathbf{w}_t; \zeta_t)} e^{-v_t} (v_t - v_*) \mid \mathcal{F}_{t-1}] \\ &= \mathbb{E}[(m_t - z_t)h_t(z_t) \mid \mathcal{F}_{t-1}] = \mathbb{E}_z[(m_t - z)h_t(z) \mid \mathcal{F}_{t-1}]. \end{aligned} \quad (5.85)$$

Since z' is an i.i.d. copy of z and independent of z given \mathcal{F}_{t-1} ,

$$m_t = \mathbb{E}[z \mid \mathcal{F}_{t-1}] = \mathbb{E}[z' \mid \mathcal{F}_{t-1}].$$

Using the conditional independence,

$$\mathbb{E}[(m_t - z)h_t(z) \mid \mathcal{F}_{t-1}] = \mathbb{E}[(z' - z)h_t(z) \mid \mathcal{F}_{t-1}].$$

By exchangeability of (z, z') conditional on \mathcal{F}_{t-1} ,

$$\mathbb{E}[(z' - z)h_t(z') \mid \mathcal{F}_{t-1}] = -\mathbb{E}[(z' - z)h_t(z) \mid \mathcal{F}_{t-1}].$$

Averaging the last two displays gives the standard symmetrization:

$$\mathbb{E}[(m_t - z)h_t(z) \mid \mathcal{F}_{t-1}] = \frac{1}{2} \mathbb{E}[(z' - z)(h_t(z) - h_t(z')) \mid \mathcal{F}_{t-1}]. \quad (5.86)$$

Next, we show that $h(z)$ is Lipschitz continuous. By definition,

$$y_t(z) = \frac{y_{t-1} + \alpha_t}{\alpha_t z + 1}, \quad h_t(z) = y_t(z)(v_t(z) - v_*).$$

Differentiate with respect to z :

$$\frac{dy_t(z)}{dz} = (y_{t-1} + \alpha_t) \frac{d}{dz}((\alpha_t z + 1)^{-1}) = -\frac{\alpha_t(y_{t-1} + \alpha_t)}{(\alpha_t z + 1)^2}.$$

Using $y_t(z)(\alpha_t z + 1) = y_{t-1} + \alpha_t$, we can rewrite this as

$$\frac{dy_t(z)}{dz} = -\frac{\alpha_t y_t(z)}{\alpha_t z + 1}.$$

Since $v_t(z) = -\log y_t(z)$, we have

$$\frac{dv_t(z)}{dz} = -\frac{1}{y_t(z)} \frac{dy_t(z)}{dz} = \frac{\alpha_t}{\alpha_t z + 1}.$$

As a result,

$$\frac{dh_t(z)}{dz} = \frac{dy_t(z)}{dz} (v_t(z) - v_*) + y_t(z) \frac{dv_t(z)}{dz} = \frac{\alpha_t y_t(z)}{\alpha_t z + 1} (1 - (v_t(z) - v_*)).$$

Since $v_t(z), v_* \in [c_0, c_1]$, then

$$|1 - (v_t(z) - v_*)| \leq 1 + c_1 - c_0,$$

and since $y_t(z) = \frac{y_{t-1} + \alpha_t}{\alpha_t z + 1} \leq y_{t-1} + \alpha_t \leq (1 + \rho)y_{t-1}$, we have

$$\left| \frac{dh_t}{dz} \right| \leq \alpha_t y_{t-1} (1 + \rho) (1 + c_1 - c_0),$$

which means i.e. h_t is L_t -Lipschitz with

$$L_t \leq \alpha_t y_{t-1} C.$$

Then, it holds

$$|(z' - z)(h_t(z) - h_t(z'))| \leq L_t (z' - z)^2 \leq C \alpha_t y_{t-1} (z' - z)^2.$$

Thus,

$$\begin{aligned} \mathbb{E} \left[|(z' - z)(h_t(z) - h_t(z'))| \mid \mathcal{F}_{t-1} \right] &\leq C \alpha_t \mathbb{E}[y_{t-1} (z' - z)^2 \mid \mathcal{F}_{t-1}] \\ &= C \alpha_t \cdot 2 \mathbb{E}[y_{t-1} (z - \mathbb{E}[z])^2 \mid \mathcal{F}_{t-1}] \leq 2C \alpha_t \delta_t^2, \end{aligned}$$

where the last step uses the definition of δ_t^2 . Applying this result to (5.86), we have

$$\left| \mathbb{E}[(\mu_t - z)h_t(z) \mid \mathcal{F}_{t-1}] \right| \leq \frac{1}{2} \mathbb{E} \left[|(z' - z)(h_t(z) - h_t(z'))| \mid \mathcal{F}_{t-1} \right] \leq C \alpha_t \delta_t^2.$$

By noting (5.85), we finish the proof. \square

Combining Lemma 5.29 and Lemma 5.30, we have the following lemma for one-step analysis of the v -update.

Lemma 5.31 *Under Assumption (5.17) (ii), we have*

$$\mathbb{E}[\alpha_t \nabla_2 F(\mathbf{w}_t, v_t)^\top (v_t - v_*)] \leq \mathbb{E}[D_\varphi(v_*, v_{t-1}) - D_\varphi(v_*, v_t) + C \alpha_t^2 \delta_t^2]. \quad (5.87)$$

Finally, we state the convergence result of SCENT in the following theorem.

Theorem 5.14 *Suppose Assumption 5.17 holds. Let $\beta_t = 1$, $\eta_t = \eta \alpha_t$, $\alpha_t < \rho e^{-v_{t-1}}$, then SCENT guarantees that*

$$\begin{aligned} & \mathbb{E} \left[\sum_{t=1}^T \alpha_t (F(\mathbf{w}_t, \nu_t) - F(\mathbf{w}_*, \nu_*)) \right] \\ & \leq \frac{1}{2\eta} \|\mathbf{w}_1 - \mathbf{w}\|_2^2 + D_\varphi(\nu_*, \nu_0) + \mathbb{E} \left[\sum_{t=1}^T \frac{\eta \alpha_t^2 \sigma_t^2}{2} + \sum_{t=1}^T C \alpha_t^2 \delta_t^2 \right]. \end{aligned}$$

Proof. Since $\eta_t = \eta \alpha_t$, from Lemma 5.28, we obtain

$$\mathbb{E}[\alpha_t \nabla_1 F(\mathbf{w}_t, \nu_t)^\top (\mathbf{w}_t - \mathbf{w}_*)] \leq \mathbb{E} \left[\frac{1}{2\eta} \|\mathbf{w}_t - \mathbf{w}\|_2^2 - \frac{1}{2\eta} \|\mathbf{w}_{t+1} - \mathbf{w}_*\|_2^2 + \frac{\eta \alpha_t^2 \sigma_t^2}{2} \right].$$

Combining this with Lemma 5.31, we have

$$\begin{aligned} & \mathbb{E}[\alpha_t (\nabla_1 F(\mathbf{w}_t, \nu_t)^\top (\mathbf{w}_t - \mathbf{w}_*) + \nabla_2 F(\mathbf{w}_t, \nu_t)^\top (\nu_t - \nu_*))] \\ & \leq \mathbb{E} \left[\frac{1}{2\eta} \|\mathbf{w}_t - \mathbf{w}\|_2^2 - \frac{1}{2\eta} \|\mathbf{w}_{t+1} - \mathbf{w}_*\|_2^2 + D_\varphi(\nu_*, \nu_{t-1}) - D_\varphi(\nu_*, \nu_t) \right] \\ & \quad + \mathbb{E} \left[\frac{\eta \alpha_t^2 \sigma_t^2}{2} + C \alpha_t^2 \delta_t^2 \right]. \end{aligned}$$

By the joint convexity of $F(\mathbf{w}, \nu)$, we have

$$\alpha_t (F(\mathbf{w}_t, \nu_t) - F(\mathbf{w}_*, \nu_*)) \leq \alpha_t (\nabla_1 F(\mathbf{w}_t, \nu_t)^\top (\mathbf{w}_t - \mathbf{w}_*) + \nabla_2 F(\mathbf{w}_t, \nu_t)^\top (\nu_t - \nu_*)).$$

Combining the last two inequalities and summing over $t = 1, \dots, T$, we have

$$\begin{aligned} & \mathbb{E} \left[\sum_{t=1}^T \alpha_t (F(\mathbf{w}_t, \nu_t) - F(\mathbf{w}_*, \nu_*)) \right] \\ & \leq \frac{1}{2\eta} \|\mathbf{w}_1 - \mathbf{w}\|_2^2 + D_\varphi(\nu_*, \nu_0) + \mathbb{E} \left[\sum_{t=1}^T \frac{\eta \alpha_t^2 \sigma_t^2}{2} + \sum_{t=1}^T C \alpha_t^2 \delta_t^2 \right]. \end{aligned}$$

□

We present two corollaries of the above theorem.

Corollary 5.2 Suppose Assumption 5.17 holds. Let $\beta_t = 1, \eta_t = \eta \alpha_t, \alpha_t = \frac{\alpha}{\sqrt{t}} < \rho e^{-\nu_{t-1}}$ for some constant $\rho > 0$, then SCENT guarantees that

$$\mathbb{E}[(F_1(\bar{\mathbf{w}}_T) - F_1(\mathbf{w}_*))] \leq \frac{D_0}{\alpha \sqrt{T}} + \frac{\alpha V}{\sqrt{T}}.$$

where $\bar{\mathbf{w}}_T = \frac{\sum_{t=1}^T \mathbf{w}_t}{T}$, $D_0 = \frac{1}{2\eta} \|\mathbf{w}_1 - \mathbf{w}_*\|_2^2 + D_\varphi(\nu_*, \nu_0)$ and

$$V = \mathbb{E} \left[\frac{\eta \sum_{t=1}^T \sigma_t^2}{2T} + \frac{\sum_{t=1}^T C \delta_t^2}{T} \right].$$

Proof. Plugging $\alpha_t = \alpha/\sqrt{T}$ into Theorem 5.14, we have

$$\mathbb{E} \left[\frac{1}{T} \sum_{t=1}^T (F(\mathbf{w}_t, \nu_t) - F(\mathbf{w}_*, \nu_*)) \right] \leq \frac{D_0}{\alpha\sqrt{T}} + \frac{\alpha V}{\sqrt{T}}.$$

Using $F_1(\mathbf{w}) = \min_{\nu} F(\mathbf{w}, \nu)$, $F_1(\mathbf{w}_*) = F(\mathbf{w}_*, \nu_*)$ and the Jensen inequality, we can finish the proof. \square

💡 Why it matters

Since δ_t, σ_t are finite, the above result implies a convergence rate of $O(1/\sqrt{T})$ for SCENT.

Corollary 5.3 Suppose Assumption 5.17 holds. Let $\beta_t = 1, \eta_t = \eta\alpha_t, \alpha_t = \frac{\alpha e^{-\nu_{t-1}}}{\sqrt{T}}$, if $\frac{1}{T} \sum_{t=1}^T e^{-\nu_{t-1}} \geq S$ almost surely, then SCENT guarantees that

$$\mathbb{E} [F_1(\hat{\mathbf{w}}_T) - F_1(\mathbf{w}_*)] \leq \frac{D_0}{\alpha\sqrt{T}S} + \frac{\alpha\bar{V}}{\sqrt{T}S}.$$

where $\hat{\mathbf{w}}_T = \frac{\sum_{t=1}^T \alpha_t \mathbf{w}_t}{\sum_{t=1}^T \alpha_t}$ and

$$\bar{V} = \mathbb{E} \left[\frac{\eta \sum_{t=1}^T e^{-2\nu_{t-1}} \sigma_t^2}{2T} + \frac{\sum_{t=1}^T C e^{-2\nu_{t-1}} \delta_t^2}{T} \right].$$

Proof. Let $\hat{\alpha}_t = \frac{\alpha_t}{\sum_{t=1}^T \alpha_t}$. From Theorem 5.14, we have

$$\begin{aligned} & \mathbb{E} \left[\sum_{t=1}^T \alpha_t \sum_{t=1}^T \hat{\alpha}_t (F(\mathbf{w}_t, \nu_t) - F(\mathbf{w}_*, \nu_*)) \right] \\ & \leq \frac{1}{2\eta} \|\mathbf{w}_1 - \mathbf{w}\|_2^2 + D_{\varphi}(\nu_*, \nu_0) + \mathbb{E} \left[\sum_{t=1}^T \frac{\eta \alpha_t^2 \sigma_t^2}{2} + \sum_{t=1}^T C \alpha_t^2 \delta_t^2 \right]. \end{aligned}$$

Since $\sum_{t=1}^T \alpha_t = \sum_{t=1}^T \frac{\alpha e^{-\nu_{t-1}}}{\sqrt{T}} \geq \alpha\sqrt{T}S$, then

$$\mathbb{E} \left[\sum_{t=1}^T \hat{\alpha}_t (F(\mathbf{w}_t, \nu_t) - F(\mathbf{w}_*, \nu_*)) \right] \leq \frac{\frac{1}{2\eta} \|\mathbf{w}_1 - \mathbf{w}\|_2^2 + D_{\varphi}(\nu_*, \nu_0)}{\alpha\sqrt{T}S} + \frac{\alpha\bar{V}}{\sqrt{T}S}.$$

Applying the joint convexity of $F(\mathbf{w}, \nu)$ and $F_1 = \min_{\nu} F(\mathbf{w}, \nu)$, we can finish the proof. \square

💡 Why it matters

Under the stated setting, SCENT reduces to SCGD with $\gamma_t = \frac{\alpha}{\sqrt{T} + \alpha}$. Since S can be lower bounded by a constant, the above corollary implies $O(1/\sqrt{T})$ convergence rate for SCGD to minimize log-E-Exp.

Analysis of the Variance Terms

Since the final convergence bound depends on the variance terms σ_t^2, δ_t^2 , we would like to provide further analysis on them.

Let us introduce some notations:

$$z(\mathbf{w}; \zeta) = e^{s(\mathbf{w}; \zeta)}, \quad \mu(\mathbf{w}) = \log \mathbb{E}_\zeta e^{s(\mathbf{w}; \zeta)}, \quad (5.88)$$

$$m_t = \mathbb{E}_\zeta e^{s(\mathbf{w}_t; \zeta)}, \quad \mu_t = \mu(\mathbf{w}_t) = \log m_t. \quad (5.89)$$

For the analysis, we make two reasonable assumptions.

Assumption 5.18. Assume there exist constants κ, σ'^2 such that (i) $\mathbb{E} \left[\frac{\mathbb{E}[z(\mathbf{w}; \zeta)^2]}{(\mathbb{E}[z(\mathbf{w}; \zeta)]^2)} \right] \leq \kappa$ for all \mathbf{w} ; (ii) $\mathbb{E} \|e^{s(\mathbf{w}_t; \zeta') - \mu_t} \nabla s(\mathbf{w}_t; \zeta')\|^2 \leq \sigma'^2$ for all t ;

Critical: These assumptions are necessary. In next section, we show that the dependence on κ is unavoidable. The second assumption is the standard bounded stochastic gradient assumption for optimizing $F_1(\mathbf{w})$.

Lemma 5.32 (Dual Variance Term) Under Assumption 5.18, we have

$$\delta_t^2 \leq 2(\kappa - 1)m_t \left(F(\mathbf{w}_t, \nu_{t-1}) - F(\mathbf{w}_*, \nu_*) + 1 \right). \quad (5.90)$$

💡 Why it matters

When $F(\mathbf{w}_t, \nu_{t-1}) - F(\mathbf{w}_*, \nu_*) \rightarrow 0$, the variance term in the convergence bound caused by the stochastic update of ν_t will be dominated by $2(\kappa - 1)m_t$. Large m_t can be mitigated by choosing small α_t .

Proof. Recall that

$$\delta_t^2 = \mathbb{E}_{\zeta_t} \left[e^{-\nu_{t-1}} (z(\mathbf{w}_t; \zeta_t) - m_t)^2 \right]$$

By Assumption 5.18(i),

$$\text{Var}(z(\mathbf{w}_t; \zeta)) \leq (\kappa - 1)m_t^2.$$

Hence

$$\delta_t^2 = e^{-\nu_{t-1}} \text{Var}(z(\mathbf{w}_t; \zeta)) \leq (\kappa - 1)e^{-\nu_{t-1}} m_t^2 = (\kappa - 1)m_t \cdot (m_t e^{-\nu_{t-1}}).$$

Let $\tilde{r}_{t-1} := m_t e^{-\nu_{t-1}}$. By the definition:

$$F(\mathbf{w}_t, \nu_{t-1}) = \mathbb{E} e^{s(\mathbf{w}_t; \zeta) - \nu_{t-1}} + \nu_{t-1} = \tilde{r}_{t-1} + \nu_{t-1}.$$

Since $\tilde{r}_{t-1} = e^{\log m_t - \nu_{t-1}}$, we have

$$F(\mathbf{w}_t, \nu_{t-1}) - (1 + \mu_t) = \tilde{r}_{t-1} + \nu_{t-1} - (1 + \log m_t) = \tilde{r}_{t-1} - \log \tilde{r}_{t-1} - 1.$$

Using $r \leq 2(r - \log r)$ for all $r > 0$ yields

$$\tilde{r}_{t-1} \leq 2(F(\mathbf{w}_t, \nu_{t-1}) - (1 + \mu_t) + 1).$$

Since \mathbf{w}_* minimizes $\mu(\mathbf{w})$, we have $\mu_t = \mu(\mathbf{w}_t) \geq \mu(\mathbf{w}_*)$ and thus $(1 + \mu_t) \geq (1 + \mu(\mathbf{w}_*)) = F(\mathbf{w}_*, \nu_*)$, implying

$$F(\mathbf{w}_t, \nu_{t-1}) - (1 + \mu_t) \leq F(\mathbf{w}_t, \nu_{t-1}) - F(\mathbf{w}_*, \nu_*).$$

As a result, we have

$$\tilde{r}_{t-1} \leq 2(F(\mathbf{w}_t, \nu_{t-1}) - F(\mathbf{w}_*, \nu_*) + 1). \quad (5.91)$$

Combining this with the bound of δ_t^2 , we complete the proof. \square

Lemma 5.33 (Primal Variance Term) *Under Assumption 5.18, we have*

$$\sigma_t^2 \leq 4\sigma'^2 (F(\mathbf{w}_t, \nu_t) - F(\mathbf{w}_*, \nu_*) + 1)^2.$$

Why it matters

When $F(\mathbf{w}_t, \nu_t) - F(\mathbf{w}_*, \nu_*) \rightarrow 0$, the variance term in the convergence bound caused by the stochastic update of \mathbf{w}_t will be dominated by $O(\sigma'^2)$.

Proof.

$$\begin{aligned} \sigma_t^2 &= \mathbb{E}_{\zeta'_t} \|\exp(s(\mathbf{w}_t; \zeta'_t) - \nu_t) \nabla s(\mathbf{w}_t; \zeta'_t)\|_2^2, \\ &= \mathbb{E}_{\zeta'_t} [e^{2(\mu_t - \nu_t)} \|\exp(s(\mathbf{w}_t; \zeta'_t) - \mu_t) \nabla s(\mathbf{w}_t; \zeta'_t)\|_2^2] \leq r_t^2 \sigma'^2, \end{aligned}$$

where $r_t = e^{\mu_t - \nu_t}$. Similar to (5.91), we have show that

$$r_t \leq 2(F(\mathbf{w}_t, \nu_t) - F(\mathbf{w}_*, \nu_*) + 1).$$

Hence,

$$\sigma_t^2 \leq 4\sigma'^2 (F(\mathbf{w}_t, \nu_t) - F(\mathbf{w}_*, \nu_*) + 1)^2.$$

\square

Algorithm 22 The SCENT Algorithm for solving CERM

```

1: Initialize  $\mathbf{w}_1, \mathbf{v}_0$ , step sizes  $\eta_t$  and  $\alpha_t$ ,  $\varphi(\mathbf{v}) = e^{-\mathbf{v}}$ .
2: for  $t = 1 \dots T - 1$  do
3:   Sample  $\mathcal{B}_t \subset \{1, \dots, n\}$  with  $|\mathcal{B}_t| = B$ 
4:   for each  $i \in \mathcal{B}_t$  do
5:     Sample  $\zeta_{i,t}, \zeta'_{i,t} \sim \mathbb{P}_i$ 
6:     Update  $\mathbf{v}_{i,t} = \arg \min_{\mathbf{v}} \exp(s_i(\mathbf{w}_t; \zeta_{i,t}) - \mathbf{v}) + \mathbf{v} + \frac{1}{\alpha_t} D_{\varphi}(\mathbf{v}, \mathbf{v}_{i,t-1})$ 
7:   end for
8:   Compute  $\mathbf{z}_t = \frac{1}{B} \sum_{i \in \mathcal{B}_t} \exp(s_i(\mathbf{w}_t; \zeta'_{i,t}) - \mathbf{v}_{i,t}) \nabla s_i(\mathbf{w}_t; \zeta'_{i,t})$ 
9:   Compute  $\mathbf{v}_t = (1 - \beta_t) \mathbf{v}_{t-1} + \beta_t \mathbf{z}_t$ 
10:  Update  $\mathbf{w}_{t+1} = \mathbf{w}_t - \eta_t \mathbf{z}_t$ 
11: end for
    
```

5.5.2.2 Compositional Entropic Risk Minimization

In this section, we extend the results to solving compositional entropic risk minimization (CERM):

$$\min_{\mathbf{w}} F_1(\mathbf{w}) := \frac{1}{n} \sum_{i=1}^n \log(\mathbb{E}_{\zeta \sim \mathbb{P}_i} \exp(s_i(\mathbf{w}; \zeta)))$$

via its equivalent min-min formulation:

$$\min_{\mathbf{w}} \min_{\mathbf{v}} F(\mathbf{w}, \mathbf{v}) := \frac{1}{n} \sum_{i=1}^n \{\mathbb{E}_{\zeta \sim \mathbb{P}_i} \exp(s_i(\mathbf{w}; \zeta) - \mathbf{v}_i) + \mathbf{v}_i\}.$$

The difference from Log-E-Exp is that there are multiple $\mathbf{v}_i, i = 1, \dots, n$, which needs to be updated using stochastic block coordinate method. The technique has been used in algorithms presented in previous sections of this chapter.

We present an extension of SCENT to solving CERM in Algorithm 22. The major change lies at the stochastic block coordinate update of \mathbf{v} in Step 5. This extension is analogous to SOX for FCCO, employing stochastic block-coordinate updates for the inner estimators. Indeed, SOX applied to CERM can be recovered as a special case of SCENT by choosing the coordinate-wise step size $\alpha_{t,i} = \frac{\gamma_t}{1-\gamma_t} e^{-\mathbf{v}_{i,t-1}}$, using an argument similar to (5.83).

Convergence analysis for convex problems

Let us define some notations:

$$\begin{aligned}
 \Phi_i(\mathbf{w}_t, \mathbf{v}_i; \zeta) &= \exp(s_i(\mathbf{w}_t; \zeta) - \mathbf{v}_i) + \mathbf{v}_i \\
 F_i(\mathbf{w}_t, \mathbf{v}_i) &= \mathbb{E}_{\zeta \sim \mathbb{P}_i} [\Phi_i(\mathbf{w}_t, \mathbf{v}_i; \zeta)] \\
 (\mathbf{w}_*, \mathbf{v}_*) &= \arg \min_{\mathbf{w}, \mathbf{v}} F(\mathbf{w}, \mathbf{v}).
 \end{aligned}$$

Similar as before, $\nu_{i,*} = \log[\mathbb{E}_{\zeta \sim \mathbb{P}_i} \exp(s_i(\mathbf{w}_*; \zeta))]$. Since we deal with stochastic block coordinate update, we introduce a virtual sequence $\bar{\nu}_t$, where

$$\bar{\nu}_{i,t} = \arg \min_{\nu} \exp(s_i(\mathbf{w}_t; \zeta_{i,t}) - \nu) + \nu + \frac{1}{\alpha_t} D_{\varphi}(\nu, \nu_{i,t-1}), \forall i$$

Following Lemma 5.26, we have

$$e^{\bar{\nu}_{i,t}} = \frac{1}{1 + \alpha_t e^{\nu_{i,t-1}}} e^{\nu_{i,t-1}} + \frac{\alpha_t e^{\nu_{i,t-1}}}{1 + \alpha_t e^{\nu_{i,t-1}}} \exp(s_i(\mathbf{w}_t; \zeta_t)), \forall i.$$

Assumption 5.19. Assume that the following conditions hold:

- (i) $s_i(\mathbf{w}; \zeta)$ is convex;
- (ii) the loss function is bounded such that $s_i(\mathbf{w}; \zeta) \in [c_0, c_1], \forall \mathbf{w}, \zeta, i$.
- (iii) there exists G such that $\mathbb{E}_{\zeta} \|\nabla s_i(\mathbf{w}_t, \zeta)\|_2^2 \leq G^2, \forall t, i$

Define $\sigma_{i,t}, \delta_{i,t}$ as

$$\begin{aligned} \sigma_{i,t}^2 &:= \mathbb{E}_{\zeta'_{i,t} \sim \mathbb{P}_i} \|\exp(s_i(\mathbf{w}_t; \zeta'_{i,t}) - \nu_{i,t}) \nabla s_i(\mathbf{w}_t; \zeta'_{i,t})\|_2^2, \forall i, t, \\ \delta_{i,t}^2 &:= \mathbb{E}_{\zeta_{i,t} \sim \mathbb{P}_i} [e^{-\nu_{i,t-1}} |e^{s_i(\mathbf{w}_t; \zeta_{i,t})} - \mathbb{E}_{\zeta_{i,t} \sim \mathbb{P}_i} [e^{s_i(\mathbf{w}_t; \zeta_{i,t})}]]^2, \forall i, t. \end{aligned}$$

Similar to Lemma 5.27, the following lemma can be proved.

Lemma 5.34 Under Assumption 5.19, if $\nu_0 \in [c_0, c_1]$ then $\nu_t \in [c_0, c_1], \forall t$.

Similar to Lemma 5.28, we have the following lemma regarding one-step update of \mathbf{w}_t .

Lemma 5.35 Under Assumption (5.19) and $\beta_t = 1$, we have

$$\mathbb{E}[\eta_t \nabla_1 F(\mathbf{w}_t, \bar{\nu}_t)^\top (\mathbf{w}_t - \mathbf{w}_*)] \leq \mathbb{E} \left[\frac{1}{2} \|\mathbf{w}_t - \mathbf{w}_*\|_2^2 - \frac{1}{2} \|\mathbf{w}_{t+1} - \mathbf{w}_*\|_2^2 \right] + \frac{\eta_t^2 \sigma_t^2}{2},$$

where $\sigma_t^2 = \frac{1}{n} \sum_{i=1}^n \sigma_{i,t}^2$.

Proof. We first bound $\mathbb{E}_t[\|\mathbf{z}_t\|_2^2 \mid \mathcal{F}_{t-1}]$, where \mathbb{E}_t denotes the expectation over randomness in t -th iteration given \mathbf{w}_t, ν_{t-1} .

$$\begin{aligned} \mathbb{E}_t[\|\mathbf{z}_t\|_2^2] &= \mathbb{E}_t \left[\left\| \frac{1}{B} \sum_{i \in \mathcal{B}_t} \exp(s_i(\mathbf{w}_t; \zeta'_{i,t}) - \nu_{i,t}) \nabla s_i(\mathbf{w}_t; \zeta'_{i,t}) \right\|_2^2 \right] \\ &= \mathbb{E}_{\mathcal{B}_t, \zeta_t} \mathbb{E}_{\zeta'_t \mid \mathcal{B}_t, \zeta_t} \left[\left\| \frac{1}{B} \sum_{i \in \mathcal{B}_t} \exp(s_i(\mathbf{w}_t; \zeta'_{i,t}) - \nu_{i,t}) \nabla s_i(\mathbf{w}_t; \zeta'_{i,t}) \right\|_2^2 \right] \\ &\leq \mathbb{E}_{\mathcal{B}_t, \zeta_t} \left[\frac{1}{B} \sum_{i \in \mathcal{B}_t} \sigma_{i,t}^2 \right] = \frac{1}{n} \sum_{i=1}^n \sigma_{i,t}^2. \end{aligned}$$

Since $\bar{\nu}_{i,t} = \nu_{i,t}, \forall i \in \mathcal{B}_t$, we have

$$\mathbb{E}_t[\mathbf{z}_t] = \mathbb{E}_{\zeta'_t, \zeta_t, \mathcal{B}_t} \left[\frac{1}{B} \sum_{i \in \mathcal{B}_t} \nabla \Phi_i(\mathbf{w}_t, \bar{v}_{i,t}; \zeta'_{i,t}) \right] = \nabla_1 F(\mathbf{w}_t, \bar{\mathbf{v}}_t).$$

Then following Lemma 3.3, we can finish the proof. \square

Next, we analyze the update of \bar{v}_t .

Lemma 5.36 *Under Assumption (5.19) (ii) and $\alpha_t \leq \min_i \rho e^{-v_{i,t-1}}$, we have*

$$\mathbb{E}[\alpha_t \nabla_2 F(\mathbf{w}_t, \bar{\mathbf{v}}_t)^\top (\bar{\mathbf{v}}_t - \mathbf{v}_*)] \leq \frac{1}{B} \mathbb{E} [D_\varphi(\mathbf{v}_*, \mathbf{v}_{t-1}) - D_\varphi(\mathbf{v}_*, \mathbf{v}_t)] + C \alpha_t^2 \delta_t^2.$$

where $D_\varphi(\mathbf{v}_*, \mathbf{v}_t) = \sum_{i=1}^n D_\varphi(v_{i,*}, v_{i,t})$ and $\delta_t^2 = \frac{1}{n} \sum_{i=1}^n \delta_{i,t}^2$.

Proof. By applying Lemma 5.30 and Lemma 5.29 for each coordinate of $\bar{v}_{i,t}$, we have

$$\mathbb{E}[\alpha_t \nabla_2 F_i(\mathbf{w}_t, \bar{v}_{i,t})^\top (\bar{v}_{i,t} - v_{i,*})] \leq D_\varphi(v_{i,*}, v_{i,t-1}) - D_\varphi(v_{i,*}, \bar{v}_{i,t}) + C \alpha_t^2 \delta_{i,t}^2, \forall i.$$

Averaging the above inequality over $i = 1, \dots, n$, we have

$$\mathbb{E}[\alpha_t \nabla_2 F(\mathbf{w}_t, \bar{\mathbf{v}}_t)^\top (\bar{\mathbf{v}}_t - \mathbf{v}_*)] \leq \frac{1}{n} \sum_{i=1}^n (D_\varphi(v_{i,*}, v_{i,t-1}) - D_\varphi(v_{i,*}, \bar{v}_{i,t})) + C \alpha_t \delta_t^2. \quad (5.92)$$

Due to the randomness of \mathcal{B}_t , we have

$$\mathbb{E}[D_\varphi(v_{i,*}, v_{i,t})] = \mathbb{E} \left[\left(1 - \frac{B}{n}\right) D_\varphi(v_{i,*}, v_{i,t-1}) + \frac{B}{n} D_\varphi(v_{i,*}, \bar{v}_{i,t}) \right], \forall i.$$

Hence

$$\begin{aligned} & \mathbb{E} \left[\frac{1}{n} \sum_{i=1}^n (D_\varphi(v_{i,*}, v_{i,t-1}) - D_\varphi(v_{i,*}, \bar{v}_{i,t})) \right] \\ &= \mathbb{E} \left[\frac{1}{n} \sum_{i=1}^n \left(D_\varphi(v_{i,*}, v_{i,t-1}) - \frac{n}{B} D_\varphi(v_{i,*}, v_{i,t}) + \left(\frac{n}{B} - 1\right) D_\varphi(v_{i,*}, v_{i,t-1}) \right) \right] \\ &= \frac{1}{B} \mathbb{E} \left[\sum_{i=1}^n (D_\varphi(v_{i,*}, v_{i,t-1}) - D_\varphi(v_{i,*}, v_{i,t})) \right]. \end{aligned}$$

Combining this with (5.92), we finish the proof. \square

Finally, we state the convergence result of SCENT in the following theorem.

Theorem 5.15 *Suppose Assumption 5.19 holds. Let $\beta_t = 1$, $\eta_t = \eta \alpha_t$, and $\alpha_t = \frac{\alpha}{\sqrt{t}} < \rho \min_i e^{-v_{i,t-1}}$, then SCENT guarantees that*

$$\mathbb{E}[(F_1(\bar{\mathbf{w}}_T) - F_1(\mathbf{w}_*))] \leq \frac{1}{2\eta\alpha\sqrt{T}} \|\mathbf{w}_1 - \mathbf{w}_*\|_2^2 + \frac{D_\varphi(\mathbf{v}_*, \mathbf{v}_0)}{\alpha B\sqrt{T}} + \frac{\alpha V}{\sqrt{T}}.$$

where $\bar{\mathbf{w}}_T = \frac{\sum_{t=1}^T \mathbf{w}_t}{T}$, and $V = \mathbb{E} \left[\frac{\eta \sum_{t=1}^T \sigma_t^2}{2T} + \frac{\sum_{t=1}^T C \delta_t^2}{T} \right]$.

💡 Why it matters

In order to achieve an ϵ -optimal solution, the above convergence bound implies the following complexity:

$$T = O \left(\frac{\|\mathbf{w}_1 - \mathbf{w}_*\|_2^4}{\eta^2 \alpha^2 \epsilon^2} + \frac{D_\varphi(\mathbf{v}_*, \mathbf{v}_0)^2}{\alpha^2 B^2 \epsilon^2} + \frac{\alpha^2 V^2}{\epsilon^2} \right).$$

For simplicity of discussion, let us consider a setting of η such that the first term matches the second term. As a result, the complexity becomes:

$$T = O \left(\frac{D_\varphi(\mathbf{v}_*, \mathbf{v}_0)^2}{\alpha^2 B^2 \epsilon^2} + \frac{\alpha^2 V^2}{\epsilon^2} \right).$$

Insight 1: Since σ_t, δ_t are finite, and $D_\varphi(\mathbf{v}_*, \mathbf{v}_0) = O(n)$, if $\alpha \propto \sqrt{n/B}$, the above result implies an iteration complexity of $O(\frac{n}{B\epsilon^2})$ for SCENT.

Insight 2: When the loss $s_i(\mathbf{w}_t; \zeta) \geq 0$ is large, the term $e^{-\nu_{i,t-1}}$ becomes very small, suggesting that the step size parameter α should be chosen small so as to mitigate the large variance term δ_t . In contrast, when the loss $s_i(\mathbf{w}_t; \zeta) < 0$ is small, the term $e^{-\nu_{i,t-1}}$ can become large, allowing α to be set relatively larger, which helps offset the large distance measure $D_\varphi(\mathbf{v}_*, \mathbf{v}_0)$.

Proof. Since $\eta_t = \eta\alpha_t$, from Lemma 5.35, we obtain

$$\mathbb{E}[\alpha_t \nabla_1 F(\mathbf{w}_t, \bar{\mathbf{v}}_t)^\top (\mathbf{w}_t - \mathbf{w}_*)] \leq \mathbb{E} \left[\frac{1}{2\eta} \|\mathbf{w}_t - \mathbf{w}_*\|_2^2 - \frac{1}{2\eta} \|\mathbf{w}_{t+1} - \mathbf{w}_*\|_2^2 + \frac{\eta\alpha_t^2 \sigma_t^2}{2} \right].$$

Adding this to the inequality in Lemma 5.36, we have

$$\begin{aligned} & \mathbb{E}[\alpha_t (\nabla_1 F(\mathbf{w}_t, \bar{\mathbf{v}}_t)^\top (\mathbf{w}_t - \mathbf{w}_*) + \nabla_2 F(\mathbf{w}_t, \bar{\mathbf{v}}_t)^\top (\bar{\mathbf{v}}_t - \mathbf{v}_*))] \\ & \leq \mathbb{E} \left[\frac{1}{2\eta} \|\mathbf{w}_t - \mathbf{w}_*\|_2^2 - \frac{1}{2\eta} \|\mathbf{w}_{t+1} - \mathbf{w}_*\|_2^2 + \frac{1}{B} D_\varphi(\mathbf{v}_*, \mathbf{v}_{t-1}) - \frac{1}{B} D_\varphi(\mathbf{v}_*, \mathbf{v}_t) \right] \\ & \quad + \mathbb{E} \left[\frac{\eta\alpha_t^2 \sigma_t^2}{2} + C\alpha_t^2 \delta_t^2 \right]. \end{aligned}$$

By the joint convexity of $F(\mathbf{w}, \mathbf{v})$, we have

$$\alpha_t (F(\mathbf{w}_t, \bar{\mathbf{v}}_t) - F(\mathbf{w}_*, \mathbf{v}_*)) \leq \alpha_t (\nabla_1 F(\mathbf{w}_t, \bar{\mathbf{v}}_t)^\top (\mathbf{w}_t - \mathbf{w}_*) + \nabla_2 F(\mathbf{w}_t, \bar{\mathbf{v}}_t)^\top (\bar{\mathbf{v}}_t - \mathbf{v}_*)).$$

Combining the last two inequalities and summing over $t = 1, \dots, T$, we have

$$\begin{aligned} & \mathbb{E} \left[\sum_{t=1}^T \alpha_t (F(\mathbf{w}_t, \bar{\nu}_t) - F(\mathbf{w}_*, \nu_*)) \right] \\ & \leq \frac{1}{2\eta} \|\mathbf{w}_1 - \mathbf{w}\|_2^2 + \frac{1}{B} D_\varphi(\nu_*, \nu_0) + \mathbb{E} \left[\sum_{t=1}^T \frac{\eta \alpha_t^2 \sigma_t^2}{2} + \sum_{t=1}^T C \alpha_t^2 \delta_t^2 \right]. \end{aligned}$$

Since $F_1(\mathbf{w}_*) = F(\mathbf{w}_*, \nu_*)$, and $F_1(\mathbf{w}_t) \leq F(\mathbf{w}_t, \bar{\nu}_t)$, we have

$$\begin{aligned} & \mathbb{E} \left[\sum_{t=1}^T \alpha_t (F_1(\mathbf{w}_t) - F_1(\mathbf{w}_*)) \right] \\ & \leq \frac{1}{2\eta} \|\mathbf{w}_1 - \mathbf{w}\|_2^2 + \frac{1}{B} D_\varphi(\nu_*, \nu_0) + \mathbb{E} \left[\sum_{t=1}^T \frac{\eta \alpha_t^2 \sigma_t^2}{2} + \sum_{t=1}^T C \alpha_t^2 \delta_t^2 \right]. \end{aligned}$$

Plugging the value of α_t , we finish the proof. \square

5.5.2.3 Why SCENT is better than ASGD?

In this section, we provide theoretical insight into why SCENT outperforms ASGD for entropic risk minimization. The key distinction between the two methods lies in their updates of the dual variable ν : SCENT employs a stochastic proximal mirror descent (SPMD) update, whereas ASGD relies on a standard SGD update. Accordingly, our analysis focuses exclusively on the ν -update while keeping \mathbf{w} fixed. In particular, we consider the following problem:

$$\min_{\nu} F(\nu) := \mathbb{E}_{\zeta} e^{s(\zeta) - \nu} + \nu, \quad (5.93)$$

where we omit \mathbf{w} in $s(\zeta)$.

Recall the definitions $z := e^{s(\zeta)}$, $m := \mathbb{E}[z]$, $r(\nu) := m e^{-\nu} = e^{\nu_* - \nu}$ as used previously, and the facts $\nu_* = \arg \min_{\nu} F(\nu) = \log m$, $F(\nu_*) = m e^{-\nu_*} + \nu_* = 1 + \nu_*$. Recall the SPMD update:

$$e^{\nu_t} = \frac{1}{1 + \alpha_t e^{\nu_{t-1}}} e^{\nu_{t-1}} + \frac{\alpha_t e^{\nu_{t-1}}}{1 + \alpha_t e^{\nu_{t-1}}} e^{s(\zeta_t)}.$$

Let us define an important quantity to characterize the difficulty of the problem:

$$\kappa = \frac{\mathbb{E}[z^2]}{(\mathbb{E}[z])^2},$$

which is known as second-order moment ratio. Larger κ indicates heavier tails or higher variability relative to the mean.

A Clean Bound of SPMD

The optimality gap can be written as

$$F(v) - F(v_*) = me^{-v} + v - (1 + v_*) = r(v) - \log r(v) - 1. \quad (5.94)$$

We assume $s(\zeta) \in [c_0, c_1]$ and without loss of generality we assume $c_1 \leq 0$. If not, we can define $s'(\zeta) = s(\zeta) - c_1$, $z' = e^{s'(\zeta)}$ and $F'(v') = \mathbb{E}[z' e^{-v'}] + v'$. Then $F(v) - F(v_*) = F'(v') - \min F'(v')$ if $v = v' - c_1$.

Lemma 5.37 (Self-bounding inequality) *For all $r > 0$,*

$$r \leq 2(r - \log r). \quad (5.95)$$

Equivalently, for all $v \in \mathbb{R}$,

$$r(v) \leq 2(F(v) - F(v_*) + 1). \quad (5.96)$$

Proof. If $0 < r \leq 2$, then $r \leq 2 \leq 2(r - \log r)$ since $r - \log r \geq 1$ for all $r > 0$. If $r \geq 2$, then $\log r \leq r/2$, hence $r - \log r \geq r/2$, i.e. $r \leq 2(r - \log r)$. Substituting $r = r(v)$ and using (5.94) yields (5.96). \square

Theorem 5.16 *Suppose $s(\zeta) \in [c_0, c_1] \leq 0$ holds. By setting $\alpha_t = \sqrt{\frac{D_\varphi(v_*, v_0)m}{2CT\text{Var}(z)}} \leq \min(\frac{m}{4C\text{Var}(z)}, \rho)$ for sufficiently large T , SPMD guarantees that*

$$\frac{1}{T} \sum_{t=1}^T \mathbb{E}[F(v_t) - F(v_*)] \leq 4\sqrt{2} \sqrt{\frac{C(\kappa - 1)(1 - r_0 + r_0 \log r_0)}{T}} + \frac{F(v_0) - F(v_*)}{T}. \quad (5.97)$$

where $C = (1 + \rho)(1 + c_1 - c_0)$, and $r_0 = r(v_0) = e^{v_* - v_0}$.

Why it matters

When $v_0 \gg v_*$ (over-estimation), then $1 - r_0 + r_0 \log r_0 = O(1)$, the dominating term becomes $O(\sqrt{\frac{\kappa}{T}})$. This upper bound characterizes the intrinsic complexity of SPMD, which depends on the second-order moment ratio κ . If $s(\zeta) \sim \mathcal{N}(\mu_s, \sigma_s^2)$, then $\kappa = e^{\sigma_s^2}$, which does not depend on the exponential of the mean μ_s but rather $e^{\sigma_s^2}$.

Proof. From Lemma 5.31, we obtain the SPMD averaged bound

$$\bar{G}_T := \frac{1}{T} \sum_{t=1}^T \mathbb{E}[F(v_t) - F(v_*)] \leq \frac{D_\varphi(v_*, v_0)}{\alpha T} + C \alpha V, \quad (5.98)$$

where

$$V := \frac{1}{T} \sum_{t=1}^T \mathbb{E}[\delta_t^2], \quad \delta_t^2 = \mathbb{E}[e^{-\nu_{t-1}}(z_t - m)^2] = e^{-\nu_{t-1}} \text{Var}(z).$$

Since $e^{-\nu_{t-1}} = r(\nu_{t-1})/m$, we can rewrite

$$V = \frac{\text{Var}(z)}{m} \cdot \frac{1}{T} \sum_{t=1}^T \mathbb{E}[r(\nu_{t-1})]. \quad (5.99)$$

By Lemma 5.37,

$$\begin{aligned} \frac{1}{T} \sum_{t=1}^T \mathbb{E}[r(\nu_{t-1})] &\leq \frac{2}{T} \sum_{t=1}^T \mathbb{E}[F(\nu_{t-1}) - F(\nu_*) + 1] \\ &= 2 \left(1 + \frac{1}{T} \sum_{t=1}^T \mathbb{E}[F(\nu_{t-1}) - F(\nu_*)] \right). \end{aligned}$$

Next, observe the index shift:

$$\begin{aligned} \sum_{t=1}^T \mathbb{E}[F(\nu_{t-1}) - F(\nu_*)] &= \mathbb{E}[F(\nu_0) - F(\nu_*)] + \sum_{t=1}^{T-1} \mathbb{E}[F(\nu_t) - F(\nu_*)] \\ &\leq \mathbb{E}[F(\nu_0) - F(\nu_*)] + \sum_{t=1}^T \mathbb{E}[F(\nu_t) - F(\nu_*)]. \end{aligned}$$

Dividing by T yields

$$\frac{1}{T} \sum_{t=1}^T \mathbb{E}[F(\nu_{t-1}) - F(\nu_*)] \leq \frac{\mathbb{E}[F(\nu_0) - F(\nu_*)]}{T} + \bar{G}_T. \quad (5.100)$$

Combining this with (5.99) we have

$$V \leq \frac{2 \text{Var}(z)}{m} \left(1 + \bar{G}_T + \frac{\mathbb{E}[F(\nu_0) - F(\nu_*)]}{T} \right). \quad (5.101)$$

Plugging (5.101) into (5.98) yields

$$\bar{G}_T \leq \frac{D_{\varphi}(\nu_*, \nu_0)}{\alpha T} + \frac{2C\alpha \text{Var}(z)}{m} \left(1 + \bar{G}_T + \frac{\mathbb{E}[F(\nu_0) - F(\nu_*)]}{T} \right).$$

If $\alpha \leq \frac{m}{4C \text{Var}(z)}$, then $\frac{2C\alpha \text{Var}(z)}{m} \leq \frac{1}{2}$, and therefore

$$\begin{aligned}\bar{G}_T &\leq \frac{2D_\varphi(v_*, v_0)}{\alpha T} + \frac{4C\alpha \text{Var}(z)}{m} \left(1 + \frac{\mathbb{E}[F(v_0) - F(v_*)]}{T}\right) \\ &\leq \frac{2D_\varphi(v_*, v_0)}{\alpha T} + \frac{4C\alpha \text{Var}(z)}{m} + \frac{F(v_0) - F(v_*)}{T}.\end{aligned}$$

Optimizing the right-hand side over α (assuming T is large enough) gives:

$$\frac{1}{T} \sum_{t=1}^T \mathbb{E}[F(v_t) - F(v_*)] \leq 4\sqrt{2} \sqrt{\frac{C D_\varphi(v_*, v_0) \text{Var}(z)}{mT}} + \frac{F(v_0) - F(v_*)}{T}.$$

With $r_0 := r(v_0) = e^{v_* - v_0}$,

$$D_\varphi(v_*, v_0) = e^{-v_*} - e^{-v_0} + e^{-v_0}(v_* - v_0) = \frac{1}{m}(1 - r_0 + r_0 \log r_0).$$

Since $\text{Var}(z)/m^2 = \kappa - 1$, thus the convergence upper bound becomes

$$4\sqrt{2} \sqrt{\frac{C(\kappa - 1)(1 - r_0 + r_0 \log r_0)}{T}} + \frac{F(v_0) - F(v_*)}{T}.$$

□

Comparison with SGD.

Benefit under the noise setting

In order to control the variance, we consider projected SGD. Let $\Pi_{[c_0, c_1]}$ denote projection onto $[c_0, c_1]$. The projected SGD update is

$$v_{t+1} = \Pi_{[c_0, c_1]}(v_t - \alpha' g_t), \quad g_t := 1 - z_t e^{-v_t}, \quad (5.102)$$

where $\{z_t\}_{t \geq 0}$ are i.i.d. copies of z and $\alpha' > 0$ is a constant step size. Note that $\mathbb{E}[g_t | v_t] = \nabla F(v_t) = 1 - m e^{-v_t}$.

We present a corollary of Theorem 3.5 for SGD to minimize $F(v)$ below.

Corollary 5.4 *Suppose $s(\zeta) \in [c_0, c_1]$ holds and $F(\cdot)$ is L -smooth in the range of $[c_0, c_1]$. Let $\{v_t\}$ follow (5.102). If $\eta \leq \frac{1}{L}$, Then*

$$\bar{G}_T^{\text{SGD}} := \frac{1}{T} \sum_{t=1}^T \mathbb{E}[F(v_t) - F(v_*)] \leq \frac{(v_0 - v_*)^2}{2\alpha' T} + \alpha' V'.$$

where

$$V' = \frac{\alpha'}{T} \sum_{t=0}^{T-1} (\delta'_t)^2 = \frac{\text{Var}(z)}{T} \sum_{t=0}^{T-1} \mathbb{E}[e^{-2v_t}].$$

We quantify the smoothness on the bounded domain of the objective, which introduces an exponential constant.

Lemma 5.38 *On $[c_0, c_1]$, the function $F(v) = me^{-v} + v$ is L -smooth with*

$$L = \sup_{v \in [c_0, c_1]} F''(v) = \sup_{v \in [c_0, c_1]} me^{-v} = me^{-c_0} = e^{v_* - c_0}.$$

Proof. We have $F''(v) = me^{-v}$, which is decreasing in v , so the maximum over $[c_0, c_1]$ is attained at c_0 . \square

Theorem 5.17 *By choosing the optimal $\alpha' = \frac{|v_0 - v_*|e^{c_0}}{\sqrt{2T\text{Var}(z)}} \leq \frac{1}{L} = \frac{e^{c_0}}{m}$, SGD's upper bound becomes*

$$\bar{G}_T^{\text{SGD}} \leq \sqrt{2}|v_0 - v_*|e^{v_* - c_0} \sqrt{\frac{\kappa - 1}{T}}. \quad (5.103)$$

where $\kappa = \mathbb{E}[z^2]/(\mathbb{E}[z])^2$.

Proof. The proof follows Corollary 5.4 by noting that $V' \leq \text{Var}(z)e^{-2c_0}$ and $\text{Var}(z) = m^2(\kappa - 1) = e^{2v_*}(\kappa - 1)$. \square

💡 Why it matters

By comparing the convergence bound of SPMD with that of SGD, the resulting ratio is:

$$\frac{1}{|v_0 - v_*|e^{v_* - c_0}}.$$

Notably, this ratio becomes exponentially small in regimes where $v_* \gg c_0$, highlighting the superior efficiency of SPMD.

Benefit under the noiseless setting

We further show that, even in the noiseless setting, the dependence of the GD update on $|v_0 - v_*|$ is unavoidable, whereas the PMD update does not exhibit such dependence when $v_0 \gg v_*$.

In the noiseless setting, where $m = \mathbb{E}[e^{s(\zeta)}]$ is known, the gradient descent (GD) iteration becomes:

$$v_{t+1} = v_t - \alpha' \nabla F(v_t) = v_t - \alpha' (1 - me^{-v_t}), \quad t \geq 0, \quad (5.104)$$

where $\alpha' > 0$ is a step size. For deterministic PMD, its update is equivalent to (cf. Lemma 5.26):

$$y_{t+1} = \frac{y_t + \alpha}{1 + \alpha m}, \quad (5.105)$$

where $y_t = e^{-v_t}$.

Lemma 5.39 (GD vs PMD) *Assume $v_0 \gg v_*$. Let $\{v_t\}_{t \geq 0}$ follow (5.104) with $\alpha' \leq 1$. Then in order to have $|\nabla F(v_t)| \leq \epsilon$, then we need at least*

$$t \geq \frac{\nu_0 - \nu_* - \log\left(\frac{1}{1-\epsilon}\right)}{\alpha'}. \quad (5.106)$$

In contrast, for deterministic PMD update (5.105), in order to ensure $|\nabla F(\nu_t)| \leq \epsilon$, it suffices that

$$t = \left\lceil \frac{\log(|1 - r_0|/\epsilon)}{\log(1 + \alpha m)} \right\rceil. \quad (5.107)$$

Proof. Recall the definition $r(\nu) := me^{-\nu} = e^{\nu_* - \nu}$. We have $|\nabla F(\nu)| = |1 - r(\nu)|$. From (5.104),

$$\nu_{t+1} = \nu_t - \alpha'(1 - e^{\nu_* - \nu_t}).$$

If $\nu_t \geq \nu_*$, then $\nu_{t+1} - \nu_* = \nu_t - \nu_* - \alpha'(1 - e^{\nu_* - \nu_t}) \geq 0$ provided $\alpha' \leq 1$. Let $r_t = e^{\nu_* - \nu_t} > 0$. Then, from GD update we have

$$r_{t+1} = r_t e^{\alpha'(1 - r_t)} \leq r_t e^{\alpha'} \leq r_0 e^{\alpha'(t+1)}.$$

In order to have $\|\nabla F(\nu_t)\|_2^2 \leq \epsilon^2$, it is necessary to have $r_t \geq 1 - \epsilon$. Hence, we need at least $t \geq \frac{\log \frac{1-\epsilon}{r_0}}{\alpha'} = \frac{\nu_0 - \nu_* - \log\left(\frac{1}{1-\epsilon}\right)}{\alpha'}$.

For deterministic PMD update (5.105), since $r_t = my_t$ we have

$$r_{t+1} - 1 = \frac{r_t - 1}{1 + \alpha m}.$$

Taking absolute value yields

$$|\nabla F(\nu_{t+1})| = \frac{|\nabla F(\nu_t)|}{(1 + \alpha m)}.$$

Solving $|\nabla F(\nu_t)| \leq |\nabla F(\nu_0)|/(1 + \alpha m)^t \leq \epsilon$ yields (5.107). \square

💡 Why it matters

Deterministic GD needs at least $\Omega((\nu_0 - \nu_*)/\alpha')$ steps to enter a constant-accuracy region, whereas PMD reduces $|\nabla F(\nu_t)|$ geometrically with rate $(1 + \alpha m)^{-1}$, yielding a complexity of order $O\left(\frac{1}{\log(1 + \alpha m)} \log \frac{1}{\epsilon}\right)$, which does not scale with ν_0 due to $|1 - r_0| = |1 - e^{\nu_* - \nu_0}| \leq 1$.

Indeed, in the noiseless setting for PMD, taking the formal limit $\alpha \rightarrow \infty$ yields $y_1 \rightarrow 1/m$ thus $\nu_1 \rightarrow \nu_*$. This highlights that the PMD update is an implicit, geometry-matched step.

5.5.2.4 An Optimal bound for SPMD

In fact, we can improve the convergence rate of SPMD to $O\left(\frac{\kappa-1}{T}\right)$, which matches a lower bound to be established. The key is just to use a specially designed learning

rate scheme α_t . Recall the SPMD update:

$$y_t = \frac{y_{t-1} + \alpha_t}{1 + \alpha_t z_t}, \quad \forall t \geq 1, \quad (5.108)$$

where $y_{t-1} = e^{-v_{t-1}}$, $z_t = e^{s(\zeta_t)}$.

Lemma 5.40 *Let $S_t := \sum_{i=1}^t z_i$ and $\bar{z}_t := S_t/t$. Initialize $y_1 = 1/z_1$ (or equivalently $\alpha_1 = \infty$) and for $t \geq 2$ choose*

$$\alpha_t := \frac{y_{t-1}}{t-1} = \frac{1}{S_{t-1}}. \quad (5.109)$$

Then for all $t \geq 1$,

$$y_t = \frac{t}{S_t}, \quad v_t = -\log y_t = \log\left(\frac{S_t}{t}\right) = \log \bar{z}_t. \quad (5.110)$$

In particular, v_t is the exact minimizer of the empirical objective

$$\widehat{F}_t(v) := \bar{z}_t e^{-v} + v \quad \text{since} \quad \arg \min_v \widehat{F}_t(v) = \log \bar{z}_t.$$

Proof. We prove (5.110) by induction. For $t = 1$, $y_1 = 1/z_1 = 1/S_1$ holds by initialization. Assume $y_{t-1} = (t-1)/S_{t-1}$. Then (5.109) gives $\alpha_t = 1/S_{t-1}$, and the recursion (5.108) yields

$$y_t = \frac{\frac{t-1}{S_{t-1}} + \frac{1}{S_{t-1}}}{1 + \frac{z_t}{S_{t-1}}} = \frac{\frac{t}{S_{t-1}}}{\frac{S_{t-1} + z_t}{S_{t-1}}} = \frac{t}{S_{t-1} + z_t} = \frac{t}{S_t}.$$

Thus $y_t = t/S_t$ and $v_t = -\log y_t = \log(S_t/t) = \log \bar{z}_t$. \square

Assumption 5.20. Assume $s(\zeta)$ is σ^2 -subgaussian, i.e.,

$$\mathbb{E}\left[e^{\lambda(s(\zeta) - \mathbb{E}[s(\zeta)])}\right] \leq e^{\lambda^2 \sigma^2 / 2} \quad \forall \lambda \in \mathbb{R}.$$

This includes Bernoulli distribution (indeed, if $s(\zeta) \in [c_0, c_1]$ a.s., then $s(\zeta) - \mathbb{E}[s(\zeta)]$ is $(c_1 - c_0)^2/4$ -subgaussian by Hoeffding's lemma).

Since $\frac{\text{Var}(z)}{(\mathbb{E}[z])^2} = \kappa - 1$, we have

$$\text{Var}(\bar{z}_T) = \frac{\text{Var}(z)}{T} = \frac{(\kappa - 1)m^2}{T}.$$

Since Lemma 5.40 gives $v_T = \log \bar{z}_T$, in light of (5.94) we can write

$$F(v_T) - F(v_*) = \frac{m}{\bar{z}_T} - 1 + \log\left(\frac{\bar{z}_T}{m}\right) = \frac{1}{Q_T} + \log Q_T - 1, \quad Q_T := \frac{\bar{z}_T}{m}. \quad (5.111)$$

Note that $\mathbb{E}[Q_T] = 1$ and $\text{Var}(Q_T) = (\kappa - 1)/T$.

Let $U_T := Q_T - 1 = (\bar{z}_T - m)/m$. Then $\mathbb{E}[U_T] = 0$ and $\mathbb{E}[U_T^2] = (\kappa - 1)/T$. Define

$$g(u) := \frac{1}{1+u} + \log(1+u) - 1, \forall u > -1$$

so that by (5.111) we have $F(v_T) - F(v_*) = g(U_T)$.

Lemma 5.41 For all $u \geq -\frac{1}{2}$,

$$g(u) \leq 2u^2.$$

Proof. Define $h(u) := 2u^2 - g(u)$ for $u > -1$. Since $g'(u) = \frac{u}{(1+u)^2}$, we have

$$h'(u) = 4u - \frac{u}{(1+u)^2} = u \left(4 - \frac{1}{(1+u)^2} \right).$$

For $u \geq -\frac{1}{2}$, $(1+u)^2 \geq \frac{1}{4}$, hence $\frac{1}{(1+u)^2} \leq 4$. Therefore $h'(u) \leq 0$ for $u \in [-\frac{1}{2}, 0]$ and $h'(u) \geq 0$ for $u \geq 0$. Thus h attains its minimum over $[-\frac{1}{2}, \infty)$ at $u = 0$, where $h(0) = 0$. Hence $h(u) \geq 0$ on $[-\frac{1}{2}, \infty)$, i.e., $g(u) \leq 2u^2$ there. \square

Lemma 5.42 Let $z_i \geq 0$ i.i.d. with finite κ . Then

$$\mathbb{P}(Q_T \leq 1/2) = \mathbb{P}(\bar{z}_T \leq m/2) \leq \exp\left(-\frac{T}{8\kappa}\right).$$

Proof. For any $\lambda > 0$, by Chernoff bound,

$$\mathbb{P}\left(\sum_{i=1}^T z_i \leq \frac{Tm}{2}\right) = \mathbb{P}\left(e^{-\lambda \sum_{i=1}^T z_i} \geq e^{-\lambda Tm/2}\right) \leq e^{\lambda Tm/2} \left(\mathbb{E}[e^{-\lambda z}]\right)^T.$$

Using $e^{-x} \leq 1 - x + x^2/2$ for $x \geq 0$,

$$\mathbb{E}[e^{-\lambda z}] \leq 1 - \lambda m + \frac{\lambda^2}{2} \mathbb{E}[z^2] \leq \exp\left(-\lambda m + \frac{\lambda^2}{2} \mathbb{E}[z^2]\right).$$

Therefore

$$\mathbb{P}(\bar{z}_T \leq m/2) \leq \exp\left(T\left(\lambda m/2 - \lambda m + \frac{\lambda^2}{2} \mathbb{E}[z^2]\right)\right) = \exp\left(-T\left(\frac{\lambda m}{2} - \frac{\lambda^2}{2} \mathbb{E}[z^2]\right)\right).$$

Choose $\lambda = m/(2\mathbb{E}[z^2])$ to get the exponent $-Tm^2/(8\mathbb{E}[z^2]) = -T/(8\kappa)$. \square

Lemma 5.43 If s is σ^2 -subgaussian, then

$$m^2 \mathbb{E}[z^{-2}] = (\mathbb{E}[e^s])^2 \mathbb{E}[e^{-2s}] \leq e^{3\sigma^2}.$$

Proof. Let $\mu = \mathbb{E}[s]$ and $X = s - \mu$. Then $\mathbb{E}[X] = 0$ and $z = e^s = e^\mu e^X$. Thus

$$m^2 \mathbb{E}[z^{-2}] = (e^\mu \mathbb{E}[e^X])^2 \cdot (e^{-2\mu} \mathbb{E}[e^{-2X}]) = (\mathbb{E}[e^X])^2 \mathbb{E}[e^{-2X}].$$

By subgaussianity,

$$\mathbb{E}[e^X] \leq e^{\sigma^2/2}, \quad \mathbb{E}[e^{-2X}] \leq e^{(2^2)\sigma^2/2} = e^{2\sigma^2}.$$

Hence $m^2 \mathbb{E}[z^{-2}] \leq e^{\sigma^2} e^{2\sigma^2} = e^{3\sigma^2}$. \square

Theorem 5.18 *Under Assumption 5.20, the SPMD iterate v_T produced by $\alpha_t = y_{t-1}/(t-1)$ satisfies*

$$\mathbb{E}[F(v_T) - F(v_*)] \leq \frac{2(\kappa-1)}{T} + e^{\frac{3}{2}\sigma^2} \exp\left(-\frac{T}{16\kappa}\right). \quad (5.112)$$

In particular, since the second term is exponentially small in T/κ ,

$$\mathbb{E}[F(v_T) - F(v_*)] = O(\kappa/T),$$

for every σ^2 -subgaussian $s(\zeta)$.

Proof. Since $F(v_T) - F(v_*) = g(U_T)$, we split the expectation on the events $\{U_T \geq -1/2\}$ and $\{U_T < -1/2\}$:

$$\mathbb{E}[g(U_T)] = \mathbb{E}[g(U_T)\mathbf{1}\{U_T \geq -1/2\}] + \mathbb{E}[g(U_T)\mathbf{1}\{U_T < -1/2\}].$$

On $\{U_T \geq -1/2\}$, Lemma 5.41 yields

$$\mathbb{E}[g(U_T)\mathbf{1}\{U_T \geq -1/2\}] \leq 2\mathbb{E}[U_T^2] = 2\text{Var}(Q_T) = 2\frac{\text{Var}(z)}{m^2T} = \frac{2(\kappa-1)}{T}.$$

On $\{U_T < -1/2\}$ we have $Q_T \leq 1/2$, and since $\log Q_T - 1 \leq 0$,

$$g(U_T) = \frac{1}{Q_T} + \log Q_T - 1 \leq \frac{1}{Q_T}.$$

Hence, by Cauchy–Schwarz,

$$\mathbb{E}[g(U_T)\mathbf{1}\{U_T < -1/2\}] \leq \mathbb{E}[Q_T^{-1}\mathbf{1}\{Q_T \leq 1/2\}] \leq (\mathbb{E}[Q_T^{-2}])^{1/2} \mathbb{P}(Q_T \leq 1/2)^{1/2}.$$

By Jensen inequality and Lemma 5.43,

$$\mathbb{E}[Q_T^{-2}] = m^2 \mathbb{E}[\bar{z}_T^{-2}] \leq m^2 \mathbb{E}[z^{-2}] \leq e^{3\sigma^2}.$$

By Lemma 5.42, $\mathbb{P}(Q_T \leq 1/2) \leq \exp(-T/(8\kappa))$. Therefore,

$$\mathbb{E}[g(U_T)\mathbf{1}\{U_T < -1/2\}] \leq e^{\frac{3}{2}\sigma^2} \exp\left(-\frac{T}{16\kappa}\right).$$

Combining the two pieces proves (5.112). \square

A Distribution-free lower bound

Indeed, we can show that $O\left(\frac{\kappa-1}{T}\right)$ is an optimal bound by establishing matching a lower bound for a black-box oracle model where the underlying distribution of z is unknown and for any query v the oracle returns

$$\Phi(v; \zeta) = ze^{-v} + v, \quad g(v; \zeta) = \nabla_v \Phi(v; \zeta) = 1 - ze^{-v}.$$

Since

$$z(\zeta) = e^v (\Phi(v; \zeta) - v) = e^v (1 - g(v; \zeta)),$$

hence, any T -query algorithm can reconstruct T i.i.d. samples z_1, \dots, z_T from P . Thus, it suffices to prove the lower bound in the standard i.i.d. sampling model for z .

Let us define a distribution class. For $\kappa \geq 2$, define

$$\mathcal{P}_\kappa := \left\{ P : z \geq 0, 0 < \mathbb{E}_P[z] < \infty, \frac{\mathbb{E}_P[z^2]}{(\mathbb{E}_P[z])^2} \leq \kappa \right\}.$$

Equivalently, $\text{Var}_P(z)/(\mathbb{E}_P[z])^2 \leq \kappa - 1$. For $P \in \mathcal{P}_\kappa$ let $m(P) = \mathbb{E}_P[z]$ and $v_*(P) = \log m(P)$.

Lemma 5.44 *Let $\phi(u) := e^{-u} + u - 1$. Then $\phi(0) = \phi'(0) = 0$ and $\phi''(u) = e^{-u}$. In particular, for all $|u| \leq 1$,*

$$\phi(u) \geq \frac{e^{-1}}{2} u^2. \quad (5.113)$$

Proof. On the interval $[-1, 1]$, $\phi''(u) = e^{-u} \geq e^{-1}$, so ϕ is e^{-1} -strongly convex on $[-1, 1]$. Since $\phi(0) = \phi'(0) = 0$, strong convexity implies $\phi(u) \geq \frac{e^{-1}}{2} u^2$ for all $|u| \leq 1$. \square

Lemma 5.45 *Let $\phi(u) = e^{-u} + u - 1$. Fix $v_0 < v_1$ and let $\Delta := v_1 - v_0$. Define*

$$H(v) := \phi(v - v_0) + \phi(v - v_1).$$

Then H is strictly convex and its unique minimizer v^\dagger lies in (v_0, v_1) . Moreover, if $\Delta \leq 1$, then

$$\inf_{v \in \mathbb{R}} H(v) \geq \frac{e^{-1}}{4} \Delta^2. \quad (5.114)$$

Proof. We have $\phi'(u) = 1 - e^{-u}$ and $\phi''(u) = e^{-u} > 0$, hence H is strictly convex with

$$H'(v) = \phi'(v - v_0) + \phi'(v - v_1) = 2 - e^{-(v-v_0)} - e^{-(v-v_1)}.$$

At the endpoints,

$$H'(v_0) = 2 - 1 - e^{-(v_0-v_1)} = 1 - e^{-\Delta} < 0, \quad H'(v_1) = 2 - e^{-(v_1-v_0)} - 1 = 1 - e^{-\Delta} > 0.$$

Since H' is strictly increasing (because $H'' > 0$), there is a unique root $v^\dagger \in (v_0, v_1)$ and thus $\inf_{v \in \mathbb{R}} H(v) = \inf_{v \in [v_0, v_1]} H(v)$.

Assume $\Delta \leq 1$. Then for all $v \in [v_0, v_1]$ we have $|v - v_0| \leq \Delta \leq 1$ and $|v - v_1| \leq \Delta \leq 1$. On $[-1, 1]$, $\phi''(u) = e^{-u} \geq e^{-1}$, so $\phi(u) \geq \frac{e^{-1}}{2}u^2$ for all $|u| \leq 1$. Therefore, for all $v \in [v_0, v_1]$,

$$H(v) \geq \frac{e^{-1}}{2}((v - v_0)^2 + (v - v_1)^2).$$

Minimizing the RHS over v yields $\inf_v ((v - v_0)^2 + (v - v_1)^2) = \Delta^2/2$, hence $\inf_{v \in \mathbb{R}} H(v) \geq \frac{e^{-1}}{4}\Delta^2$. \square

Lemma 5.46 (Le Cam's Two-point Method) *Let P_0, P_1 be two distributions and let $L_0(\cdot), L_1(\cdot)$ be nonnegative loss functions. For any estimator \hat{a} measurable w.r.t. the data,*

$$\max\{\mathbb{E}_{P_0}[L_0(\hat{a})], \mathbb{E}_{P_1}[L_1(\hat{a})]\} \geq \frac{1 - \text{TV}(P_0, P_1)}{2} \inf_a (L_0(a) + L_1(a)). \quad (5.115)$$

Proof. Let $M := (P_0 + P_1)/2$ and write $dP_0 = (1 + f) dM$, $dP_1 = (1 - f) dM$ where $|f| \leq 1$ and $\int |f| dM = \text{TV}(P_0, P_1)$. Then for any (possibly random) decision A ,

$$\begin{aligned} \mathbb{E}_{P_0}[L_0(A)] + \mathbb{E}_{P_1}[L_1(A)] &= \int \left(L_0(A)(1 + f) + L_1(A)(1 - f) \right) dM \\ &= \int \left((L_0(A) + L_1(A)) + f(L_0(A) - L_1(A)) \right) dM \\ &\geq \int \left((L_0(A) + L_1(A)) - |f|(L_0(A) + L_1(A)) \right) dM \\ &= \int (L_0(A) + L_1(A))(1 - |f|) dM \\ &\geq \inf_a (L_0(a) + L_1(a)) \int (1 - |f|) dM \\ &= (1 - \text{TV}(P_0, P_1)) \inf_a (L_0(a) + L_1(a)). \end{aligned}$$

Taking half and using $\max\{x, y\} \geq (x + y)/2$ yields (5.115). \square

The final distribution-free suboptimality lower bound is stated in the following theorem.

Theorem 5.19 *Let $z = e^{s(\zeta)} \geq 0$ with $m(P) = \mathbb{E}_P[z]$ and $v_*(P) = \log m(P)$. For $\kappa \geq 2$, define*

$$\mathcal{P}_\kappa := \left\{ P : z \geq 0, 0 < \mathbb{E}_P[z] < \infty, \frac{\mathbb{E}_P[z^2]}{\mathbb{E}_P[z]^2} \leq \kappa \right\}.$$

Let $F_P(v) := m(P)e^{-v} + v$ and $v_(P) = \arg \min_v F_P(v)$. Then there exists an absolute constant $c > 0$ such that for all $T \geq \kappa$, any (possibly adaptive) algorithm using*

T value/gradient oracle calls and outputting \hat{v} satisfies

$$\sup_{P \in \mathcal{P}_\kappa} \mathbb{E}_P[F_P(\hat{v}) - F_P(v_*(P))] \geq c \frac{\kappa - 1}{T}. \quad (5.116)$$

Proof. We construct two strictly positive hard instances in \mathcal{P}_κ . Fix $\varepsilon \in (0, 1]$ and define two distributions supported on $\{\varepsilon, \kappa\}$:

$$P_i^\varepsilon : \quad \mathbb{P}(z = \kappa) = p_i, \quad \mathbb{P}(z = \varepsilon) = 1 - p_i, \quad i \in \{0, 1\},$$

where

$$p_0 := \frac{1}{\kappa}, \quad p_1 := p_0 + h, \quad h := \frac{1}{8\sqrt{\kappa T}}.$$

Since $T \geq \kappa$, we have $h \leq \frac{1}{8\kappa}$ so $p_1 \in (0, 1)$.

Next we show that $P_0^\varepsilon, P_1^\varepsilon \in \mathcal{P}_\kappa$. For a generic $p \in (0, 1)$ and support $\{\varepsilon, \kappa\}$, define

$$R_\varepsilon(p) := \frac{\mathbb{E}[z^2]}{\mathbb{E}[z]^2} = \frac{p\kappa^2 + (1-p)\varepsilon^2}{(p\kappa + (1-p)\varepsilon)^2}.$$

Let $u := \varepsilon/\kappa \in (0, 1/\kappa] \subset (0, 1]$. Then

$$R_\varepsilon(p) = \frac{p + (1-p)u^2}{(p + (1-p)u)^2}.$$

We claim $R_\varepsilon(p) \leq \frac{1}{p}$ for all $u \in [0, 1]$. Indeed,

$$\begin{aligned} & (p + (1-p)u)^2 - p(p + (1-p)u^2) \\ &= p^2 + 2p(1-p)u + (1-p)^2u^2 - p^2 - p(1-p)u^2 \\ &= (1-p)u(2p + (1-2p)u) \geq 0, \end{aligned}$$

because $u \in [0, 1]$ and $2p + (1-2p)u \geq \min\{2p, 1\} \geq 0$. Thus $R_\varepsilon(p) \leq 1/p$. Since $p_0 = 1/\kappa$ and $p_1 \geq p_0$, we have $1/p_i \leq \kappa$, hence $R_\varepsilon(p_i) \leq \kappa$ and therefore $P_0^\varepsilon, P_1^\varepsilon \in \mathcal{P}_\kappa$.

Next, we compute the separation Δ between v_* 's. Let $m_i^\varepsilon = \mathbb{E}_{P_i^\varepsilon}[z] = \varepsilon + p_i(\kappa - \varepsilon)$ and $v_i^\varepsilon = \log m_i^\varepsilon$. Then

$$m_1^\varepsilon - m_0^\varepsilon = h(\kappa - \varepsilon) \geq h(\kappa - 1), \quad m_0^\varepsilon = \varepsilon + p_0(\kappa - \varepsilon) = 1 + \left(1 - \frac{1}{\kappa}\right)\varepsilon \in [1, 2].$$

Hence

$$\Delta := |v_1^\varepsilon - v_0^\varepsilon| = \log\left(1 + \frac{m_1^\varepsilon - m_0^\varepsilon}{m_0^\varepsilon}\right) \geq \frac{1}{2} \cdot \frac{h(\kappa - 1)}{2} = \frac{\kappa - 1}{32\sqrt{\kappa T}},$$

where we used $\log(1+x) \geq x/2$ for $x \in [0, 1/2]$ and the fact that $\frac{h(\kappa-\varepsilon)}{m_0^\varepsilon} \leq h\kappa \leq 1/8$. In particular, $\Delta \leq h\kappa \leq 1/8 < 1$.

Next, we show the lower bound of $\inf_v \left((F_0(v) - F_0(v_0^\varepsilon)) + (F_1(v) - F_1(v_1^\varepsilon)) \right)$. Under P_i^ε the objective is $F_i(v) = m_i^\varepsilon e^{-v} + v$ and the optimal value is $F_i(v_i^\varepsilon) = 1 + v_i^\varepsilon$. Thus the suboptimality can be written as

$$F_i(v) - F_i(v_i^\varepsilon) = e^{v_i^\varepsilon - v} + (v - v_i^\varepsilon) - 1 = \phi(v - v_i^\varepsilon), \quad \phi(u) = e^{-u} + u - 1.$$

Let $v_0^\varepsilon < v_1^\varepsilon$ and set $u = v - v_0^\varepsilon$. Then

$$\phi(v - v_0^\varepsilon) + \phi(v - v_1^\varepsilon) = \phi(u) + \phi(u - \Delta).$$

The function $u \mapsto \phi(u) + \phi(u - \Delta)$ is convex and its minimizer lies in $[0, \Delta]$. Since $\Delta \leq 1$, applying Lemma 5.45 gives

$$\phi(u) + \phi(u - \Delta) \geq \frac{e^{-1}}{4} \Delta^2.$$

Therefore,

$$\inf_v \left((F_0(v) - F_0(v_0^\varepsilon)) + (F_1(v) - F_1(v_1^\varepsilon)) \right) \geq \frac{e^{-1}}{4} \Delta^2. \quad (5.117)$$

Next, we show the total variation between P_0^ε and P_1^ε is bounded. Because the two distributions differ only in the Bernoulli parameter,

$$\text{KL}(P_0^\varepsilon, P_1^\varepsilon) = p_0 \log \frac{p_0}{p_1} + (1 - p_0) \log \frac{1 - p_0}{1 - p_1}.$$

Using the bound $\text{KL}(P, Q) \leq \chi^2(P, Q)$ and the fact that for Bernoulli measures $\chi^2(P_0^\varepsilon, P_1^\varepsilon) = \frac{h^2}{p_1(1-p_1)}$, we get

$$\text{KL}(P_0^\varepsilon, P_1^\varepsilon) \leq \frac{h^2}{p_1(1-p_1)}.$$

Since $h \leq \frac{1}{2\kappa}$, we have $p_1 \leq p_0 + h \leq \frac{3}{2\kappa} \leq \frac{3}{4}$, hence $1 - p_1 \geq 1/4$, and also $p_1 \geq p_0 = 1/\kappa$. Therefore $p_1(1-p_1) \geq \frac{1}{4\kappa}$ and

$$\text{KL}(P_0^\varepsilon, P_1^\varepsilon) \leq 4\kappa h^2.$$

For T i.i.d. samples, this gives

$$\text{KL}((P_0^\varepsilon)^{\otimes T}, (P_1^\varepsilon)^{\otimes T}) = T \text{KL}(P_0^\varepsilon, P_1^\varepsilon) \leq 4\kappa T h^2 = \frac{1}{16}.$$

By Pinsker's inequality,

$$\text{TV}((P_0^\varepsilon)^{\otimes T}, (P_1^\varepsilon)^{\otimes T}) \leq \sqrt{\frac{1}{2} \text{KL}((P_0^\varepsilon)^{\otimes T}, (P_1^\varepsilon)^{\otimes T})} \leq \sqrt{\frac{1}{32}} \leq \frac{1}{4}.$$

Finally, we apply Lemma 5.46 to $P_0 = (P_0^\varepsilon)^{\otimes T}$, $P_1 = (P_1^\varepsilon)^{\otimes T}$ and losses

$$L_i(v) := F_i(v) - F_i(v_i^\varepsilon) \geq 0.$$

Using (5.117) and $\text{TV} \leq 1/4$ yields for any estimator \widehat{v} ,

$$\max_{i \in \{0,1\}} \mathbb{E}_{P_i^\varepsilon} [F_i(\widehat{v}) - F_i(v_i^\varepsilon)] \geq \frac{1 - \text{TV}}{2} \cdot \frac{e^{-1}}{4} \Delta^2 \geq \frac{3}{8} \cdot \frac{e^{-1}}{4} \Delta^2 = \frac{3e^{-1}}{32} \Delta^2.$$

Substituting $\Delta^2 \geq \frac{(\kappa-1)^2}{1024 \kappa T} \geq \frac{\kappa-1}{2048 T}$ (since $\kappa \geq 2$) gives

$$\max_{i \in \{0,1\}} \mathbb{E}_{P_i^\varepsilon} [F_i(\widehat{v}) - F_i(v_i^\varepsilon)] \geq \frac{3}{65536 e} \cdot \frac{\kappa-1}{T}.$$

Since $P_0^\varepsilon, P_1^\varepsilon \in \mathcal{P}_\kappa$, this implies (5.116) with $c = \frac{3}{65536 e}$. \square

5.6 History and Notes

Finite-sum coupled compositional optimization (FCCO) was first formalized in our work (Qi et al., 2021c) for optimizing average precision, an empirical estimator of the area under the precision–recall curve. We proposed the SOAP algorithm for AP maximization and established the first complexity bound of $O\left(\frac{n}{\varepsilon^5}\right)$ for finding an ε -stationary solution. Their algorithm is closely related to SOX, but differs in that it does not employ a moving-average gradient estimator. The framework was demonstrated on applications including image classification and molecular property prediction for drug discovery. The analysis of SOAP draws inspiration from the original SCGD analysis Wang et al. (2017a), while significantly improving upon its $O(1/\varepsilon^8)$ complexity with the a better hyper-parameter setting, leading to Theorem 4.1.

To accelerate convergence, we subsequently adopted the moving average gradient estimator for FCCO (Wang et al., 2022). While this approach achieves a complexity order of $O\left(\frac{n}{B\varepsilon^4}\right)$, it does not benefit from the variance reduction gained by using mini-batches to estimate inner function values. The limitation arises because we treat all inner functions as a single vector variable and compute a sparse unbiased stochastic estimator for this vector; consequently, the estimator does not enjoy the advantages of inner mini-batching. This improved rate and analysis was inspired by the stochastic compositional momentum method (Ghadimi et al., 2020).

Subsequently, we proposed the SOX algorithm—a significant advancement for solving FCCO (Wang and Yang, 2022), encompassing new design, theoretical analysis, and practical applications. In that work, we established a complexity of $O\left(\frac{n\sigma_0^2}{B\varepsilon^4}\right)$

for SOX to find an ϵ -stationary solution in non-convex smooth FCCO problems. It integrates the analysis of stochastic block coordinate update of the \mathbf{u} sequences with that of stochastic compositional momentum method.

Building on this, we developed a double-loop restarted algorithm that utilizes SOX in the inner loop to address non-convex problems under the μ -PL (Polyak-Lojasiewicz) condition, i.e., $\|\nabla F(\mathbf{w})\|_2^2 \geq \mu(F(\mathbf{w}) - \min_{\mathbf{w}} F(\mathbf{w}))$. This approach yields an improved complexity of $O\left(\frac{n\sigma_0^2}{\mu^2 B \epsilon}\right)$ for finding an ϵ -optimal solution. This result further implies a complexity of $O\left(\frac{n\sigma_0^2}{\mu^2 B \epsilon}\right)$ for strongly convex FCCO problems and $O\left(\frac{n\sigma_0^2}{B \epsilon^3}\right)$ for convex FCCO problems, requiring no assumptions on the individual convexity of inner and outer functions beyond the overall convexity of the objective. The improved convergence analysis under the PL condition for the double-loop restarted algorithm was inspired by our prior work on stochastic compositional optimization for distributionally robust learning (Qi et al., 2021b). A comparable complexity bound of $O\left(\frac{1}{\mu^2 \epsilon}\right)$ for a single-loop algorithm in the context of Stochastic Convex Optimization (SCO) under the PL condition was subsequently established in (Jiang et al., 2023), which considers the application of SCO in training energy-based models.

Furthermore, for convex FCCO instances where the outer function is both convex and monotonically non-decreasing and the inner functions are convex, (Wang and Yang, 2022) reformulated the problem as a convex-concave min-max optimization problem and established a complexity of $O\left(\frac{n\sigma_0^2}{B \epsilon^2}\right)$ under a weak duality convergence measure. Finally, when a μ -strongly convex regularizer is present, the complexity is further refined to $O\left(\frac{n\sigma_0^2}{\mu^2 B \epsilon}\right)$ for finding an ϵ -optimal solution in terms of Euclidean distance to the optimum. This analysis was mostly inspired by (Zhang and Lan, 2024), which is the first work that establishes the optimal complexity for solving convex SCO where the outer function is both convex and monotonically non-decreasing and the inner function is convex.

Later, Jiang et al. (2022) proposed the Multi-Block-Single-Probe Variance Reduction (MSVR) algorithm for FCCO, establishing improved complexity bounds over SOX by leveraging the mean squared smoothness of the inner functions. For non-convex smooth FCCO problems, MSVR improves the complexity to $O\left(\frac{n\sigma_0}{B \epsilon^3}\right)$ for identifying an ϵ -stationary solution.

For objectives satisfying the μ -PL condition, a double-loop restarted MSVR algorithm achieves an improved complexity of $O\left(\frac{n\sigma_0}{\mu B \epsilon}\right)$ to find an ϵ -optimal solution. Consequently, this approach yields a complexity of $O\left(\frac{n\sigma_0}{\mu B \epsilon}\right)$ for strongly convex FCCO problems and $O\left(\frac{n\sigma_0}{B \epsilon^2}\right)$ for convex FCCO problems.

The analysis for non-smooth weakly convex FCCO and the SONX (v2) algorithm was studied in our work (Hu et al., 2024b). This work established a complexity of $O\left(\frac{n\sigma_0}{B \epsilon^6}\right)$ for finding a nearly ϵ -stationary solution for weakly convex inner and outer

functions. A similar analysis for a special case of weakly-convex SCO was conducted in (Zhu et al., 2023c). When the outer function is smooth, the complexity is improved in this book to $O\left(\frac{n\sigma_0}{B\epsilon^4}\right)$. The SONEX algorithm for solving weakly convex FCCO with non-smooth outer functions was proposed in our work (Chen et al., 2025b).

The ALEXR algorithm and its analysis for convex FCCO instances appeared in our work (Wang and Yang, 2023), where the outer function is both convex and monotonically non-decreasing and the inner functions are convex. For the first time, we established a complexity of $O\left(\frac{n\sigma_0^2}{B\epsilon^2}\right)$ for finding an ϵ -optimal solution of convex FCCO. Our analysis of the stochastic block coordinate update for the dual variables is primarily informed by the framework in Alacaoglu et al. (2025), which addresses convex-concave minimax problems with bilinear structures. The extrapolation for the gradient of the dual variable is inspired by (Zhang et al., 2021). It is worth mentioning that for strongly convex FCCO with smooth outer functions, we only established the convergence of ALEXR for the Euclidean distance to the optimum. However, it is possible to establish the convergence for the objective gap and even the duality gap following our work on strongly-convex strongly-concave min-max optimization (Yan et al., 2020b).

In (Wang and Yang, 2023), we also established the lower bounds for convex FCCO and strongly convex FCCO, which matches the upper bounds. Our derivation of the lower bound for convex FCCO with non-smooth outer functions builds upon the construction presented in (Zhang and Lan, 2024) for SCO.

The double-loop ALEXR was developed in Chen et al. (2025b), which was mostly inspired by a line of work on weakly-convex concave min-max problems (Rafique et al., 2018; Yan et al., 2020b; Zhang et al., 2022). (Rafique et al., 2018) is the first work that proves the convergence for weakly-convex (strongly)-concave problems. Yan et al. (2020b) simplified the algorithm for weakly-convex strongly-concave problems with μ -strong concavity on the dual variable and established a complexity of $O\left(\frac{1}{\mu^2\epsilon^4}\right)$ for finding an nearly ϵ -stationary point. The later work (Zhang et al., 2022) improved the complexity to $O\left(\frac{1}{\mu\epsilon^4}\right)$ with a simple change on the number of iteration for the inner loop.

The non-convex analysis of ASGD for compositional CVaR minimization first appeared in (Zhu et al., 2022b) for one-way partial AUC optimization. The geometric-aware algorithm SCENT for CERM and its analysis were developed in (Wei et al., 2026). It remains an interesting problem to conduct fine-grained analysis of SCENT for non-convex problems.

A more general framework than FCCO is the so-called conditional stochastic optimization (CSO), defined as:

$$\min_{\mathbf{w}} \mathbb{E}_{\xi} \left[f_{\xi} \left(\mathbb{E}_{\zeta|\xi} [g(\mathbf{w}; \zeta, \xi)] \right) \right].$$

This paradigm was formally introduced by Hu et al. (2020), who analyzed a biased SGD (BSGD) algorithm employing a large inner mini-batch and a constant outer mini-batch. For non-convex smooth problems, using an inner batch size of $O(\epsilon^{-2})$ results in an iteration complexity of $O(\epsilon^{-4})$, which translates to a total sample com-

plexity of $O(\epsilon^{-6})$. This performance is inferior to that of SOX when $n/B < \epsilon^{-2}$. For convex and μ -strongly convex CSO problems, an inner batch size of $O(\epsilon^{-1})$ yields iteration complexities of $O(\epsilon^{-2})$ and $O(\mu^{-2}\epsilon^{-1})$, respectively. Notably, the latter complexity is likewise worse than that of restarted SOX when $n/B < O(\epsilon^{-1})$.