

---

# Conditional Gradient Methods via Stochastic Path-Integrated Differential Estimator

---

Alp Yurtsever<sup>1</sup> Suvrit Sra<sup>2</sup> Volkan Cevher<sup>1</sup>

## Abstract

We propose a class of novel variance-reduced stochastic conditional gradient methods. By adopting the recent stochastic path-integrated differential estimator technique (SPIDER) of Fang et al. (2018) for the classical Frank-Wolfe (FW) method, we introduce SPIDER-FW for finite-sum minimization as well as the more general expectation minimization problems. SPIDER-FW enjoys superior complexity guarantees in the non-convex setting, while matching the best known FW variants in the convex case. We also extend our framework à la conditional gradient sliding (CGS) of Lan & Zhou (2016), and propose SPIDER-CGS.

## 1. Introduction

We study two different problem settings in this paper, *finite-sum* and the more general *expectation minimization*:

$$\underset{x \in \Omega}{\text{minimize}} \quad F(x) := \begin{cases} \mathbb{E}_{\xi} f(x, \xi) & \text{(expectation)} \\ \frac{1}{n} \sum_{i=1}^n f_i(x) & \text{(finite-sum)} \end{cases} \quad (1)$$

- ▷  $\Omega \subset \mathbb{R}^d$  is the convex and compact domain;
- ▷  $F$ ,  $f$  and  $f_i$  are differentiable and *possibly non-convex*;
- ▷  $\xi \sim \mathcal{P}$  is a random variable, supported on  $\Xi \subset \mathbb{R}^p$ .

The expectation objective template covers a large number of applications in machine learning and statistics. The finite-sum template frequently arises in M-estimation and empirical risk minimization problems. Accordingly, there are many applications for stochastic conditional gradient methods both in convex and non-convex settings. This includes low-rank matrix and tensor factorizations, structured sparse matrix estimation, dictionary learning applications, multi-class classification (considered as a motivating example in

---

<sup>1</sup>Ecole Polytechnique Fédérale de Lausanne, Switzerland

<sup>2</sup>Massachusetts Institute of Technology, USA. Correspondence to: Alp Yurtsever <alp.yurtsever@epfl.ch>.

a related work by Hazan & Luo (2016)), constrained deep learning problems (e.g., Ravi et al. (2018) present an application in computer vision) and many more.

Template (1) can be solved by using the well-known projected stochastic gradient descent method (SGD). At each iteration, SGD takes a stochastic gradient step followed by a projection to ensure the feasibility of the new point. However, in many applications, projection onto  $\Omega$  can impose a computational bottleneck (e.g., projection onto the nuclear norm-ball may require a full singular value decomposition), or it can be even intractable (e.g., dual structural SVMs (Lacoste-Julien et al., 2013)).

As a result, the Frank-Wolfe (FW) algorithm (*aka* conditional gradient method) has witnessed tremendous interest in the machine learning community in the last decade. FW avoids projection by leveraging the so-called linear minimization oracle instead:

$$\text{lmo}_{\Omega}(v) = \arg \min_{x \in \Omega} \langle x, v \rangle. \quad (\text{lmo})$$

lmo is significantly cheaper to compute than projection. For instance, *lmo* of nuclear norm-ball requires the computation of the leading singular vectors only (*vs.* the full spectrum for projection), which can be efficiently found by using Krylov subspace methods (Jaggi, 2013).

Our focus in this paper is on the theoretical complexity of stochastic and finite-sum FW, with an aim to identify and present the tightest results known so far. To this end, we also propose a class of novel variance-reduced stochastic optimization algorithms, based on the recent *stochastic path-integrated differential estimator* technique (SPIDER) of Fang et al. (2018).

By combining SPIDER with the classical FW method, we introduce SPIDER-FW for finite-sum and expectation minimization problems. We also extend our framework à la conditional gradient sliding (CGS) of Lan & Zhou (2016), and propose SPIDER-CGS.

From SPIDER, we adopt the variance bounds from Lemma 1 of (Fang et al., 2018), which relates the variance of the current estimator to the error of the previous estimator and the distance between the iterates. Nevertheless, Fang et al. (2018) introduce SPIDER for normalized gradient method

which is fundamentally different than the FW method. Accordingly, the analyses are different.

A natural and widely used measure for the convergence of conditional gradient methods is the so-called FW-gap (*cf.*, Section 3). However, we are not aware of any reported FW-gap convergence of CGS in the non-convex settings. Therefore, we present a new compact proof (and an extension for the stochastic setting) in the supplementary material. Although CGS does not seem to provide any improvement upon FW in this setup, we use the proof technique to extend SPIDER-CGS for the non-convex settings.

Finally, for the majority of the variance reduced FW methods in the literature, the analysis relies on the induction technique with respect to the outer loop counter, along with a sufficient improvement condition for each epoch. Consequently, at the beginning of each epoch parameters are typically reset. Instead, we set our learning-rate parameters with respect to the more natural total iteration counter, and we go over the proof without induction.

**Roadmap.** Section 2 provides an extensive discussion on the related works. Section 3 recalls some basic notions from the optimization theory. Sections 4 and 5 present SPIDER-FW and SPIDER-CGS respectively, along with their theoretical guarantees for various problem settings. Section 6 provides an extensive comparison of the theoretical complexity of FW methods in the literature. Finally, Section 7 draws the conclusions.

**Notation.** We work on the real space with Euclidean norms for simplicity. Throughout,  $\langle \cdot, \cdot \rangle$  represents the standard inner product associated with the Euclidean norm  $\| \cdot \|$ . We use the notation  $[n] = \{1, 2, \dots, n\}$ .  $D$  denotes diameter of  $\Omega$ , *i.e.*,  $D = \max_{(x,y) \in \Omega^2} \|x - y\|$ .

## 2. Related Works

**Frank-Wolfe algorithm.** This classical method is first proposed by Frank & Wolfe (1956) for solving smooth convex minimization problems with a polyhedral domain constraint (polyhedral constraint is relaxed for an arbitrary convex compact set by Jaggi (2013)).

---

### Algorithm 1 Frank-Wolfe algorithm

---

**Input:**  $x^1 \in \Omega$   
**for**  $k = 1, 2, \dots, K$  **do**  
    Compute  $w^k \in \text{lmo}_\Omega(\nabla F(x^k))$   
    Update  $x^{k+1} = x^k + \eta_k(w^k - x^k)$   
**end for**

---

Given an initial guess  $x^1 \in \Omega$ , at each iteration, FW minimizes the linear approximation of  $F$  at the current iterate  $x^k$  over  $\Omega$  (this corresponds to the *lmo* step). Clearly, minimization of a linear function returns an extreme point of the

domain. Since the new estimate is constructed as a convex combination of the current iterate and this extreme point, by definition it is a feasible point, hence the method does not require projections.

FW did not attract much attention in the machine learning community due to its slow convergence rate until Hazan (2008) and Jaggi (2013) emphasize the favorable trade-off between the convergence rate and the per-iteration cost provided by FW in key applications. Following then, there has been a resurgence of interest for FW-type algorithms.

FW literature in the stochastic optimization setting is much younger compared to the projection-based stochastic gradient methods. We can trace it back to a variant for online learning proposed by Hazan & Kale (2012). More recently, Hazan & Luo (2016) introduced stochastic FW methods with and without variance reduction for finite-sum problems. Very recently, Mokhtari et al. (2018) have proposed an alternative scheme for expectation minimization setting.

FW methods for non-convex stochastic learning are relatively understudied, most of the known results are due to Reddi et al. (2016). We discuss more details on the theoretical aspects of all these FW variants in Section 6.

**Conditional gradient sliding.** Lan & Zhou (2016) has recently developed the conditional gradient sliding method (CGS) based on the idea of applying accelerated gradient method (AG) of Nesterov (1983) for solving problems from template (1), but applying FW to the projection subproblems. In other words, CGS establishes the convergence of an inexact version of AG. Surprisingly, CGS has superior first-order oracle complexity compared to FW, although they have the same *lmo* complexity. We discuss more details and variants of CGS in Section 6.

**SPIDER.** There has been extensive research on variance reduced stochastic optimization methods in order to address the needs of machine learning and big data applications. Therefore, various variance reduction techniques are proposed in the last few years such as SAG (Roux et al., 2012), SVRG (Johnson & Zhang, 2013), SAGA (Defazio et al., 2014), and more recently SARAH (Nguyen et al., 2017) and SPIDER (Fang et al., 2018).

SARAH and SPIDER are closely related since they use the same sequential update rule for the gradient estimator  $v^k$ :

$$v^k = \nabla f_S(x^k) - \nabla f_S(x^{k-1}) + v^{k-1}.$$

However, SARAH uses this estimator in the classical gradient descent template, while SPIDER adopts a normalized gradient approach, and their results and analyses differ.

As described by Wang et al. (2018), the original SPIDER framework has a restrictive step-size (proportional with the target accuracy  $\epsilon$ ), which makes the algorithm impractical

though its theoretical appeal. Surprisingly, this problem disappears in the conditional gradient framework analysis.

### 3. Preliminaries

**Solution.** We denote a solution and the optimal value of problem (1) by  $x^*$  and  $F^*$  respectively:

$$x^* \in \arg \min_{x \in \Omega} F(x) \quad \text{and} \quad F^* = F(x^*).$$

**The measure of non-stationarity.** For unconstrained non-convex problems, the typical measure of non-stationarity is the gradient norm, because  $\|\nabla f(x)\| \rightarrow 0$  as  $x$  converges to a stationary point. However, this measure cannot be used for constrained problems, because  $\|\nabla f(x)\|$  might not converge to 0 when we approach to a solution on the boundary.

Instead, we will use the quantity

$$\mathcal{G}(x) := \max_{u \in \Omega} \langle u - x, -\nabla F(x) \rangle,$$

which is widely known as the FW gap (because it naturally appears in the analysis of FW-type methods). FW gap is always non-negative, and it gets 0 if and only if we are looking at a stationary point or a solution. Therefore, FW-gap is a meaningful measure of non-stationarity. It was also used by Lacoste-Julien (2016) and Reddi et al. (2016).

**$\epsilon$ -solution.** Due to the fundamental difference in the measure of non-stationarity, we use different definitions of approximate solutions for convex and non-convex problems:

▷ If  $F$  is *convex*, we say  $x_\epsilon^* \in \Omega$  is an  $\epsilon$ -solution if

$$F(x_\epsilon^*) - F^* \leq \epsilon.$$

▷ If  $F$  is *non-convex*, we say that a random variable  $x_\epsilon^*$  chosen uniformly from a finite set of points  $\{x^1, x^2, \dots, x^k\}$  is an  $\epsilon$ -solution if

$$\mathbb{E}[\mathcal{G}(x_\epsilon^*)] \leq \epsilon.$$

It is common to provide convergence guarantees in expectation for a randomly chosen iterate in the non-convex setting. See (Reddi et al., 2016).

**Oracle models.** We adopt the following black-box oracle model from Reddi et al. (2016), to establish a ground for comparing the convergence speed of different algorithms:

- Stochastic first-order oracle (*sfo*)  
For a stochastic function  $\mathbb{E}_\xi f(\cdot, \xi)$  with  $\xi \sim \mathcal{P}$ , *sfo* returns a pair  $(f(x, \xi'), \nabla f(x, \xi'))$  where  $\xi'$  is an *iid* sample from  $\mathcal{P}$ . (Nemirovski & Yudin, 1983)
- Incremental first-order oracle (*ifo*)  
For a finite-sum, *ifo* takes an index  $i \in [n]$  and returns  $(f_i(x), \nabla f_i(x))$ . (Agarwal & Bottou, 2014)

- Linear minimization oracle (*lmo*)  
Well-known oracle of FW-type methods.

**Assumptions (finite-sum).** For the finite-sum setting, we assume that  $f_i(x)$  has an averaged  $L$ -Lipschitz gradient:

$$\mathbb{E} \|\nabla f_i(x) - \nabla f_i(y)\|^2 \leq L^2 \|x - y\|^2, \quad \forall (x, y) \in \Omega^2.$$

Note that this implies  $F$  is  $L$ -smooth, since

$$\begin{aligned} \|\nabla F(x) - \nabla F(y)\|^2 &= \|\mathbb{E}(\nabla f_i(x) - \nabla f_i(y))\|^2 \\ &\leq \mathbb{E} \|\nabla f_i(x) - \nabla f_i(y)\|^2 \leq L^2 \|x - y\|^2. \end{aligned}$$

**Assumptions (expectation).** For the expectation minimization, we assume that  $\nabla f(x, \xi)$  is an unbiased estimate of the gradient:

$$\mathbb{E} \nabla f(x, \xi) = \nabla F(x).$$

We also assume that the variance is bounded:

$$\mathbb{E} \|\nabla f(x, \xi) - \nabla F(x)\|^2 \leq \sigma^2 < \infty, \quad \forall \xi \in \Xi, \forall x \in \Omega.$$

And finally, we assume an averaged  $L$ -Lipschitz gradient condition, *i.e.*, the following condition holds  $\forall \xi \in \Xi$ :

$$\mathbb{E} \|\nabla f(x, \xi) - \nabla f(y, \xi)\|^2 \leq L^2 \|x - y\|^2, \quad \forall (x, y) \in \Omega^2.$$

Similar to the finite-sum, this implies the smoothness of  $F$ .

**Assumptions (non-convex).** Let us denote the initial point by  $\bar{x}^1$ . Initial suboptimality  $F(\bar{x}^1) - F^*$  appears in the convergence bounds in the non-convex setting. For notational convenience, we denote an upper bound on this term by  $\mathcal{E}$ :

$$F(\bar{x}^1) - F^* \leq \mathcal{E}$$

Assume that  $F^*$  is finite, then there exists a finite  $\mathcal{E}$  which satisfies this bound. This is a direct consequence of the smoothness of  $F$  and the boundedness of domain.

All these assumptions are mild and frequently used in the analysis of stochastic methods and FW-type algorithms.

## 4. SPIDER Frank-Wolfe

This section presents SPIDER-FW algorithm and its convergence guarantees for various problem settings.

Our methods have a double loop structure, hence the iterates and the parameters have two different iteration counters  $t$  and  $k$ , such as  $x^{t,k}$ . For notational simplicity, we drop the first counter when there is no ambiguity, such as  $x^k$ . Throughout,  $s_{t,k}$  denotes the total number of inner iterations until  $k^{\text{th}}$  iteration of  $t^{\text{th}}$  epoch. In our pseudocodes, *draw samples* means *iid* samples for expectation minimization, and uniform selection with replacement in the finite-sum.

**Algorithm 2** SPIDER Frank-Wolfe

**Input:**  $\bar{x}^1 \in \Omega$   
**for**  $t = 1, 2, \dots, T$  **do**  
 Set  $x^1 = \bar{x}^t$   
 Draw  $Q_t$  samples  $\mathcal{Q}_t$   
 Compute  $v^1 = \nabla f_{\mathcal{Q}_t}(x^1)$   
 Compute  $w^1 \in \text{lmo}_{\Omega}(v^1)$   
 Update  $x^2 = x^1 + \eta_{t,1}(w^1 - x^1)$   
**for**  $k = 2, 3, \dots, K_t$  **do**  
 Draw  $S_{t,k}$  samples  $\mathcal{S}_{t,k}$   
 Compute  $v^k = \nabla f_{\mathcal{S}_{t,k}}(x^k) - \nabla f_{\mathcal{S}_{t,k}}(x^{k-1}) + v^{k-1}$   
 Compute  $w^k \in \text{lmo}_{\Omega}(v^k)$   
 Update  $x^{k+1} = x^k + \eta_{t,k}(w^k - x^k)$   
**end for**  
 Set  $\bar{x}^{t+1} = x^{K_t+1}$   
**end for**

**SPIDER-FW: Convex finite-sum**

We consider SPIDER-FW with

$$K_t = 2^{t-1} \text{ for } t = 1, 2, \dots, T.$$

We choose the sampling parameters

$$S_{t,k} = K_t \quad \mathcal{Q}_t = [n]$$

and the learning rate parameter

$$\eta_{t,k} = \frac{2}{s_{t,k} + 1} \text{ where } s_{t,k} = K_t + k - 1.$$

**Theorem 1.** Consider the convex finite-sum optimization template, and suppose that the assumptions in Section 3 for this template hold. Then, estimate  $x^{t,k}$  of SPIDER-FW with the parameter choices described above satisfies

$$\mathbb{E}[F(x^{t,k})] - F^* = \mathcal{O}\left(\frac{LD^2}{s_{t,k}}\right)$$

**Corollary 1.** The ifo and lmo complexities of SPIDER-FW for achieving  $\epsilon$ -solution in this setting are as follows:

$$\begin{aligned} \#(\text{ifo}) &= \mathcal{O}\left(n \ln\left(\frac{LD^2}{\epsilon}\right) + \frac{L^2 D^4}{\epsilon^2}\right) \\ \#(\text{lmo}) &= \mathcal{O}\left(\frac{LD^2}{\epsilon}\right) \end{aligned}$$

**SPIDER-FW: Convex expectation minimization**

We consider SPIDER-FW with

$$K_t = 2^{t-1} \text{ for } t = 1, 2, \dots, T.$$

We choose the sampling parameters

$$S_{t,k} = K_t \quad \mathcal{Q}_t = \left\lceil \frac{\sigma^2 K_t^2}{5L^2 D^2} \right\rceil$$

and the learning rate parameter

$$\eta_{t,k} = \frac{2}{s_{t,k} + 1} \text{ where } s_{t,k} = K_t + k - 1.$$

**Theorem 2.** Consider the convex expectation minimization template, and suppose that the assumptions in Section 3 for this template hold. Then, estimate  $x^{t,k}$  of SPIDER-FW with the parameter choices described above satisfies

$$\mathbb{E}[F(x^{t,k})] - F^* = \mathcal{O}\left(\frac{LD^2}{s_{t,k}}\right)$$

**Corollary 2.** The sfo and lmo complexities of SPIDER-FW for achieving  $\epsilon$ -solution in this setting are as follows:

$$\begin{aligned} \#(\text{sfo}) &= \mathcal{O}\left(\frac{\sigma^2 D^2 + L^2 D^4}{\epsilon^2}\right) \\ \#(\text{lmo}) &= \mathcal{O}\left(\frac{LD^2}{\epsilon}\right) \end{aligned}$$

SPIDER-FW has the same asymptotic oracle complexities as SCGS (Lan & Zhou, 2016) in this setting. In Section 5, we also present the SPIDER-CGS.

**SPIDER-FW: Non-convex finite-sum**

We consider SPIDER-FW with

$$K_t = K = \lceil \sqrt{n} \rceil.$$

Furthermore, we choose the parameters as

$$S_{t,k} = S = \lceil \sqrt{n} \rceil \quad \mathcal{Q}_t = [n]$$

and the learning rate parameter

$$\eta_{t,k} = \eta = \frac{1}{\sqrt{s_{T,K}}} \text{ where } s_{T,K} = TK.$$

**Theorem 3.** Consider the non-convex finite-sum template, and suppose that the assumptions in Section 3 for this template hold. Denote by  $x^{\text{out}}$  an iterate  $x^{t,k}$  of SPIDER-FW chosen uniformly random over all  $(t, k)$  pairs up to  $(T, K)$ . Then, the following bound on the FW-gap holds:

$$\mathbb{E}[\mathcal{G}(x^{\text{out}})] = \mathcal{O}\left(\frac{\mathcal{E} + LD^2}{\sqrt{s_{T,K}}}\right)$$

Although it is impractical to store all estimates until the final iteration, all stochastic methods for the non-convex setting shown in Table 1 have this type convergence guarantees, see (Reddi et al., 2016). More stringently, Lacoste-Julien (2016) proves convergence of non-convex FW in terms of the running best iterate. However, we cannot keep track of best estimate in the stochastic setting, simply because we cannot measure the FW-gap.

**Corollary 3.** *The ifo and lmo complexities of SPIDER-FW for achieving  $\epsilon$ -solution in the non-convex finite-sum setting are as follows:*

$$\begin{aligned}\#(ifo) &= \mathcal{O}\left(\frac{\sqrt{n}}{\epsilon^2}(\mathcal{E}^2 + L^2 D^4)\right) \\ \#(lmo) &= \mathcal{O}\left(\frac{1}{\epsilon^2}(\mathcal{E}^2 + L^2 D^4)\right)\end{aligned}$$

SPIDER-FW has better *ifo* complexity than state-of-the-art in the non-convex finite-sum setting. It improves the dependence on sample size  $n$ . See Table 1 for comparison.

#### SPIDER-FW: Non-convex expectation minimization

We consider SPIDER-FW with

$$K_t = K = \lceil \sigma/\epsilon \rceil.$$

Furthermore, we choose the parameters as

$$S_{t,k} = S = \lceil \sigma/\epsilon \rceil \quad Q_t = Q = \lceil 4(\sigma/\epsilon)^2 \rceil$$

and the learning rate parameter

$$\eta_{t,k} = \eta = \frac{1}{\sqrt{s_{T,K}}} \text{ where } s_{T,K} = TK.$$

**Theorem 4.** *Consider the non-convex expectation minimization template, and suppose that the assumptions in Section 3 for this template hold. Denote by  $x^{out}$  an iterate  $x^{t,k}$  of SPIDER-FW chosen uniformly random over all  $(t, k)$  pairs up to  $(T, K)$ . Then, the following bound holds:*

$$\mathbb{E}[\mathcal{G}(x^{out})] = \mathcal{O}\left(\frac{\mathcal{E} + LD^2}{\sqrt{s_{T,K}}}\right) + \frac{\epsilon}{2}$$

**Corollary 4.** *The sfo and lmo complexities of SPIDER-FW for achieving  $\epsilon$ -solution in this setting are as follows:*

$$\begin{aligned}\#(sfo) &= \mathcal{O}\left(\frac{\sigma}{\epsilon^3}(\mathcal{E}^2 + L^2 D^4)\right) \\ \#(lmo) &= \mathcal{O}\left(\frac{1}{\epsilon^2}(\mathcal{E}^2 + L^2 D^4)\right)\end{aligned}$$

Once again, SPIDER-FW enjoys superior *sfo* complexity while maintaining the same *lmo* complexity as its competitors. SVRF was the state-of-the-art with  $\mathcal{O}(\epsilon^{-10/3})$ , see Reddi et al. (2016).

## 5. SPIDER Conditional Gradient Sliding

This section presents SPIDER-CGS (as shown in Algorithm 3) and its convergence guarantees for various settings. SPIDER-CGS has the same oracle complexity as the SPIDER-FW.

#### Algorithm 3 SPIDER Conditional Gradient Sliding

**Input:**  $\bar{x}^1 = \bar{y}^1 \in \Omega$   
**for**  $t = 1, 2, \dots, T$  **do**  
 Set  $x^1 = \bar{x}^t$  and  $y^1 = \bar{y}^t$   
 Update  $z^1 = y^1 + \gamma_{t,1}(x^1 - y^1)$   
 Draw  $Q_t$  samples  $\mathcal{Q}_t$   
 Compute  $v^1 = \nabla f_{\mathcal{Q}_t}(z^1)$   
 $x^2 = \text{CndG}(x^1, v^1, \alpha_{t,1}, \beta_{t,1})$   
 Update  $y^2 = y^1 + \gamma_{t,1}(x^2 - y^1)$   
**for**  $k = 2, 3, \dots, K_t$  **do**  
 Update  $z^k = y^k + \gamma_{t,k}(x^k - y^k)$   
 Draw  $S_{t,k}$  samples  $\mathcal{S}_{t,k}$   
 Compute  $v^k = \nabla f_{\mathcal{S}_{t,k}}(z^k) - \nabla f_{\mathcal{S}_{t,k}}(z^{k-1}) + v^{k-1}$   
 $x^{k+1} = \text{CndG}(x^k, v^k, \alpha_{t,k}, \beta_{t,k})$   
 Update  $y^{k+1} = y^k + \gamma_{t,k}(x^{k+1} - y^k)$   
**end for**  
 Set  $\bar{x}^{t+1} = x^{K_t+1}$  and  $\bar{y}^{t+1} = y^{K_t+1}$   
**end for**  
**function**  $u^+ = \text{CndG}(u, v, \alpha, \beta)$   
 Set  $u^1 = u$   
**for**  $k = 1, 2, \dots$  **do**  
 Compute  $w^k \in \text{lmo}_\Omega(v + \beta(u^k - u))$   
 Evaluate  $\zeta_k = \langle v + \beta(u^k - u), u^k - w^k \rangle$   
**if**  $\zeta_k \leq \alpha$  **then**  
     **break**  
**end if**  
 Set  $\theta_k = \min\{1, \zeta_k / (\beta \|w^k - u^k\|^2)\}$   
 Update  $u^{k+1} = u^k + \theta_k(w^k - u^k)$   
**end for**  
 Set  $u^+ = u^k$   
**end function**

#### SPIDER-CGS: Convex finite-sum

We consider SPIDER-CGS with

$$K_t = \lceil 2^{t/2} \rceil \text{ for } t = 1, 2, \dots, T.$$

Furthermore, we choose the sampling parameters as

$$S_{t,k} = 9K_t s_{t,K_t}^2 \quad Q_t = [n]$$

CndG subsolver parameters as

$$\beta_{t,k} = \frac{3}{2}L\gamma_{t,k} \quad \alpha_{t,k} = \frac{2LD^2}{(s_{t,k} + 1)^2}$$

and the learning rate parameter as

$$\gamma_{t,k} = \frac{3}{s_{t,k} + 2} \text{ where } s_{t,k} = \sum_{\tau=1}^{t-1} K_\tau + k$$

**Theorem 5.** Consider the convex finite-sum template, and suppose that the assumptions in Section 3 for this template hold. Then, estimate  $y^{t,k}$  of SPIDER-CGS with the parameter choices described above satisfies

$$\mathbb{E}[F(y^{t,k})] - F^* = \mathcal{O}\left(\frac{LD^2}{s_{t,k}^2}\right)$$

**Corollary 5.** The *ifo* and *lmo* complexities of SPIDER-CGS for achieving  $\epsilon$ -solution in this template are as follows:

$$\begin{aligned} \#(\text{ifo}) &= \mathcal{O}\left(n \ln\left(\frac{LD^2}{\epsilon}\right) + \frac{L^2 D^4}{\epsilon^2}\right) \\ \#(\text{lmo}) &= \mathcal{O}\left(\frac{LD^2}{\epsilon}\right) \end{aligned}$$

Remark that the STORC (Hazan & Luo, 2016) has a better *ifo* complexity, but under the additional assumption of Lipschitz continuity of  $F$ .

#### SPIDER-CGS: Convex expectation minimization

We consider SPIDER-CGS with

$$K_t = \lceil 2^{t/2} \rceil \text{ for } t = 1, 2, \dots, T.$$

Furthermore, we choose the sampling parameters as

$$S_{t,k} = 9K_t s_{t,K_t}^2 \quad Q_t = \lceil \frac{\sigma^2 s_{t,K_t}^4}{L^2 D^2} \rceil$$

CndG subsolver parameters as

$$\beta_{t,k} = \frac{3}{2}L\gamma_{t,k} \quad \alpha_{t,k} = \frac{2LD^2}{(s_{t,k} + 1)^2}$$

and the learning rate parameter as

$$\gamma_{t,k} = \frac{2}{s_{t,k} + 1} \text{ where } s_{t,k} = \sum_{\tau=1}^{t-1} K_\tau + k$$

**Theorem 6.** Consider the convex expectation minimization template, and suppose that the assumptions in Section 3 for this template hold. Then, estimate  $y^{t,k}$  of SPIDER-CGS with the parameter choices described above satisfies

$$\mathbb{E}[F(y^{t,k})] - F^* = \mathcal{O}\left(\frac{LD^2}{s_{t,k}^2}\right)$$

**Corollary 6.** The *sfo* and *lmo* complexities of SPIDER-CGS for achieving  $\epsilon$ -solution in convex expectation minimization problems are as follows:

$$\begin{aligned} \#(\text{sfo}) &= \mathcal{O}\left(\frac{\sigma^2 D^2 + L^2 D^4}{\epsilon^2}\right) \\ \#(\text{lmo}) &= \mathcal{O}\left(\frac{LD^2}{\epsilon}\right) \end{aligned}$$

#### SPIDER-CGS: Non-convex finite-sum

We consider SPIDER-CGS with

$$K_t = K = \lceil \sqrt{n} \rceil.$$

Furthermore, we choose the sampling parameters as

$$S_{t,k} = K \quad Q_t = \lceil n \rceil$$

CndG subsolver parameters as

$$\beta_{t,k} = \frac{3}{2}L\gamma \quad \alpha_{t,k} = LD^2\gamma$$

and the learning rate parameter as

$$\gamma_{t,k} = \gamma = \frac{1}{\sqrt{s_{T,K}}} \text{ where } s_{T,K} = TK.$$

**Theorem 7.** Consider the non-convex finite-sum template, and suppose that the assumptions in Section 3 for this template hold. Denote by  $y^{\text{out}}$  an iterate  $y^{t,k}$  of SPIDER-CGS chosen uniformly random over all  $(t, k)$  pairs up to  $(T, K)$ . Then, the following bound on the FW-gap holds:

$$\mathbb{E}[\mathcal{G}(y^{\text{out}})] = \mathcal{O}\left(\frac{\mathcal{E} + LD^2}{\sqrt{s_{T,K}}}\right)$$

**Corollary 7.** The *ifo* and *lmo* complexities of SPIDER-CGS for achieving  $\epsilon$ -solution in non-convex finite-sum are

$$\begin{aligned} \#(\text{ifo}) &= \mathcal{O}\left(\frac{\sqrt{n}}{\epsilon^2}(\mathcal{E}^2 + L^2 D^4)\right) \\ \#(\text{lmo}) &= \mathcal{O}\left(\frac{1}{\epsilon^2}(\mathcal{E}^2 + L^2 D^4)\right) \end{aligned}$$

#### SPIDER-CGS: Non-convex expectation minimization

We consider SPIDER-CGS with

$$K_t = K = \lceil \sigma/\epsilon \rceil.$$

Furthermore, we choose the sampling parameters as

$$S_{t,k} = K \quad Q_t = \lceil 4(\sigma/\epsilon)^2 \rceil$$

CndG subsolver parameters as

$$\beta_{t,k} = \frac{3}{2}L\gamma \quad \alpha_{t,k} = LD^2\gamma$$

and the learning rate parameter as

$$\gamma_{t,k} = \gamma = \frac{1}{\sqrt{s_{T,K}}} \text{ where } s_{T,K} = TK.$$

**Theorem 8.** Consider the non-convex expectation minimization template, and suppose that the assumptions in Section 3 for this template hold. Denote by  $y^{\text{out}}$  an iterate  $y^{t,k}$  of

SPIDER-CGS chosen uniformly random over all  $(t, k)$  pairs up to  $(T, K)$ . Then, the following bound holds:

$$\mathbb{E}[\mathcal{G}(y^{out})] = \mathcal{O}\left(\frac{\mathcal{E} + LD^2}{\sqrt{s_{T,K}}}\right) + \frac{\epsilon}{2}$$

**Corollary 8.** *The sfo and lmo complexities of SPIDER-CGS for achieving  $\epsilon$ -solution in non-convex expectation minimization problems are as follows:*

$$\begin{aligned} \#(sfo) &= \mathcal{O}\left(\frac{\sigma}{\epsilon^3}(\mathcal{E}^2 + L^2 D^4)\right) \\ \#(lmo) &= \mathcal{O}\left(\frac{1}{\epsilon^2}(\mathcal{E}^2 + L^2 D^4)\right) \end{aligned}$$

## 6. Comparison & Discussions

This section presents an extensive comparison of theoretical aspects of FW methods. Table 1 compiles a summary of this comparison.

### 6.1. Convex optimization camp

**Batch setting.** FW achieves an  $\epsilon$ -solution after  $\mathcal{O}(1/\epsilon)$  iterations. This complexity is optimal for a large class of methods that construct the decision variable through convex combination of *lmo* outputs (Lan, 2014). CGS, on the other side, enjoys  $\mathcal{O}(1/\sqrt{\epsilon})$  first order oracle complexity while keeping the same  $\mathcal{O}(1/\epsilon)$  *lmo* complexity, by reusing the same gradients over multiple iterations Lan & Zhou (2016).

**Stochastic setting.** Hazan & Kale (2012) propose Online-FW for an online-learning setting, but as mentioned later by Hazan & Luo (2016), these results can be translated to the stochastic template via standard conversion approaches, and gets  $\mathcal{O}(1/\epsilon^4)$  and  $\mathcal{O}(1/\epsilon^2)$  complexities for *sfo* and *lmo* calls respectively.

A natural extension of FW for stochastic setting is described by Hazan & Luo (2016), as shown in Algorithm 4. This method (SFW) is shown to converge with  $\mathcal{O}(1/k)$  rate when the sample size  $S_k = \Theta(k^2)$ , hence it provides an  $\epsilon$ -solution with  $\mathcal{O}(1/\epsilon^3)$  *sfo* and  $\mathcal{O}(1/\epsilon)$  *lmo* complexities.

---

#### Algorithm 4 Stochastic Frank-Wolfe

---

**Input:**  $x^1 \in \Omega$   
**for**  $k = 1, 2, \dots, K$  **do**  
     Draw  $S_k$  samples  $\mathcal{S}_k$   
     Compute  $w^k \in \text{lmo}_\Omega(\nabla f_{\mathcal{S}_k}(x^k))$   
     Update  $x^{k+1} = x^k + \eta_k(w^k - x^k)$   
**end for**

---

Lan & Zhou (2016) extend their CGS framework to the stochastic setting by introducing SCGS in Section 3 of their original work. While keeping the optimal  $\mathcal{O}(1/\epsilon)$  *lmo* complexity, SCGS achieves  $\mathcal{O}(1/\epsilon^2)$  *sfo* complexity, which even gets  $\mathcal{O}(1/\epsilon)$  under strong convexity assumption.

Hazan & Luo (2016) introduce the stochastic variance reduced Frank-Wolfe method (SVRF) by adopting the variance reduction techniques from Johnson & Zhang (2013) and Mahdavi et al. (2013). SVRF is explicitly designed for the finite-sum setting, and it requires  $\mathcal{O}(n \ln(1/\epsilon))$  full gradients as well as  $\mathcal{O}(1/\epsilon^2)$  *ifo* and  $\mathcal{O}(1/\epsilon)$  *lmo* to get an  $\epsilon$ -solution.

To further improve *ifo* complexity of SVRF, Hazan & Luo (2016) design a variant based on CGS. This variant, stochastic variance reduced condition gradient sliding (STORC), also requires  $\mathcal{O}(n \ln(1/\epsilon))$  full gradients and  $\mathcal{O}(1/\epsilon)$  *lmo*, but it enjoys a reduced number of *ifo* calls at  $\mathcal{O}(1/\epsilon^{1.5})$ . Compared to SVRF, however, STORC additionally assumes that  $F$  is Lipschitz continuous in domain  $\Omega$ . Also remark that STORC gets better rates under additional assumptions such as strong-convexity.

Lu & Freund (2018) propose a stochastic FW variant which requires  $\mathcal{O}(1/\epsilon)$  *lmo* and  $\mathcal{O}(n + 1/\epsilon)$  *ifo* complexity for the convex finite-sum. However, the proposed method relies on a special structure of the objective function, that  $f_i$  are univariate functions of the fitted value  $\langle a_i, x \rangle$  for some given data sample  $a_i$ .

All stochastic FW variants we discussed up to know are based on an increasing mini-batch size. Very recently, Mokhtari et al. (2018) have proposed an alternative scheme (SFW-1) for expectation minimization setting, which requires a single *sfo* at each iteration. Nevertheless, SFW-1 has an arguably worse computational complexity compared to SFW, with its  $\mathcal{O}(1/\epsilon^3)$  calls of *sfo* and *lmo*. We emphasize the applications of SFW-1 in submodular maximization, but this is beyond the scope of our work.

For the convex finite-sum setting, SPIDER-FW and SPIDER-CGS share the same complexities as SVRF.

### 6.2. Non-convex optimization camp

**Batch setting.** FW converges asymptotically to a stationary point, see Section 2.2 in (Bertsekas, 1999). To our knowledge, Yu et al. (2014) shows the first convergence rates for FW in non-convex setting, and Lacoste-Julien (2016) proves a non-asymptotic  $\mathcal{O}(1/\sqrt{k})$  rate in FW-gap for a FW variant with line-search.

**Stochastic setting.** As shown by Reddi et al. (2016), SFW achieves to an  $\epsilon$ -solution with  $\mathcal{O}(1/\epsilon^4)$  *sfo* and  $\mathcal{O}(1/\epsilon^2)$  *lmo* complexities. Moreover, they also analyze SVRF in the non-convex setting (but they call it SVFW), and prove that it takes  $\mathcal{O}(1/\epsilon^{10/3})$  *sfo* and  $\mathcal{O}(1/\epsilon^2)$  *lmo* complexity for this method to get an  $\epsilon$ -solution. In the finite-sum setting, the former is replaced by  $\mathcal{O}(n + n^{2/3}/\epsilon^2)$  *ifo* calls.

Reddi et al. (2016) also propose a variant (SAGAFW) based on the SAGA variance reduction technique described by

## Conditional Gradient Methods via Stochastic Path-Integrated Differential Estimator

	convex				non-convex			
	finite-sum		expectation		finite-sum		expectation	
	(ifo)	(lmo)	(sfo)	(lmo)	(ifo)	(lmo)	(sfo)	(lmo)
FW	$\mathcal{O}(n\epsilon^{-1})$	$\mathcal{O}(\epsilon^{-1})$	-	-	$\mathcal{O}(n\epsilon^{-2})$	$\mathcal{O}(\epsilon^{-2})$	-	-
CGS	$\mathcal{O}(n\epsilon^{-1/2})$	$\mathcal{O}(\epsilon^{-1})$	-	-	$\mathcal{O}(n\epsilon^{-2})$	$\mathcal{O}(\epsilon^{-2})$	-	-
SFW	$\mathcal{O}(\epsilon^{-3})$	$\mathcal{O}(\epsilon^{-1})$	$\mathcal{O}(\epsilon^{-3})$	$\mathcal{O}(\epsilon^{-1})$	$\mathcal{O}(\epsilon^{-4})$	$\mathcal{O}(\epsilon^{-2})$	$\mathcal{O}(\epsilon^{-4})$	$\mathcal{O}(\epsilon^{-2})$
SFW-1	$\mathcal{O}(\epsilon^{-3})$	$\mathcal{O}(\epsilon^{-3})$	$\mathcal{O}(\epsilon^{-3})$	$\mathcal{O}(\epsilon^{-3})$	-	-	-	-
Online-FW	$\mathcal{O}(\epsilon^{-4})$	$\mathcal{O}(\epsilon^{-2})$	$\mathcal{O}(\epsilon^{-4})$	$\mathcal{O}(\epsilon^{-2})$	-	-	-	-
SCGS	$\mathcal{O}(\epsilon^{-2})$	$\mathcal{O}(\epsilon^{-1})$	$\mathcal{O}(\epsilon^{-2})$	$\mathcal{O}(\epsilon^{-1})$	$\mathcal{O}(\epsilon^{-4})$	$\mathcal{O}(\epsilon^{-2})$	$\mathcal{O}(\epsilon^{-4})$	$\mathcal{O}(\epsilon^{-2})$
SVRF / SVFW	$\mathcal{O}(n \ln(\epsilon^{-1}) + \epsilon^{-2})$	$\mathcal{O}(\epsilon^{-1})$	-	-	$\mathcal{O}(n + n^{2/3}\epsilon^{-2})$	$\mathcal{O}(\epsilon^{-2})$	$\mathcal{O}(\epsilon^{-10/3})$	$\mathcal{O}(\epsilon^{-2})$
STORC <sup>†</sup>	$\mathcal{O}(n \ln(\epsilon^{-1}) + \epsilon^{-3/2})$	$\mathcal{O}(\epsilon^{-1})$	-	-	-	-	-	-
<i>SPIDER-FW</i>	$\mathcal{O}(n \ln(\epsilon^{-1}) + \epsilon^{-2})$	$\mathcal{O}(\epsilon^{-1})$	$\mathcal{O}(\epsilon^{-2})$	$\mathcal{O}(\epsilon^{-1})$	$\mathcal{O}(n^{1/2}\epsilon^{-2})$	$\mathcal{O}(\epsilon^{-2})$	$\mathcal{O}(\epsilon^{-3})$	$\mathcal{O}(\epsilon^{-2})$
<i>SPIDER-CGS</i>	$\mathcal{O}(n \ln(\epsilon^{-1}) + \epsilon^{-2})$	$\mathcal{O}(\epsilon^{-1})$	$\mathcal{O}(\epsilon^{-2})$	$\mathcal{O}(\epsilon^{-1})$	$\mathcal{O}(n^{1/2}\epsilon^{-2})$	$\mathcal{O}(\epsilon^{-2})$	$\mathcal{O}(\epsilon^{-3})$	$\mathcal{O}(\epsilon^{-2})$

Table 1: Comparison of conditional gradient methods for stochastic optimization. Contribution of *this work* is highlighted with blue font. See Section 6 for more details.

FW (Frank & Wolfe, 1956; Jaggi, 2013), CGS (Lan & Zhou, 2016), SFW (Hazan & Luo, 2016; Reddi et al., 2016), SFW-1 (Mokhtari et al., 2018), Online-FW (Hazan & Kale, 2012), SCGS (Lan & Zhou, 2016), SVRF / SVFW (Hazan & Luo, 2016; Reddi et al., 2016), STORC (Hazan & Luo, 2016)

Defazio et al. (2014). However, we omit SAGAFW because there is an issue in the analysis of this method (while telescoping Eq.(14), in page 1249).

Qu et al. (2018) show the convergence rate for special instances of CGS and SCGS in the non-convex setting. However, they consider a different convergence criterion based on a proximal gradient mapping rather than the conventional FW-gap. Consequently, their results are incomparable with the rest of the literature. For the fact that we are running a *projection-free* method, the FW-gap is a more natural choice than the projection/proximal gradient norm.

We provide a parameter setting and a compact proof for CGS and SCGS in the supplemental material. Note however this setting simply gets the same guarantees as FW and SFW respectively. Whether or not CGS can provide improved oracle complexities compared to FW in the non-convex setting, is an open problem.

For the non-convex setting, SPIDER-FW and SPIDER-CGS have the same oracle complexities, superior to SVRF (which is the state-of-the-art to our knowledge) for finite-sum and expectation minimization problems.

### 6.3. Results from Concurrent Works

By the time we prepared this manuscript, the idea of combining SPIDER with the FW analysis was not explored yet. However, stochastic variance reduction methods and FW-type algorithms are both very active research fields. In this part, we discuss some results from a few concurrent works that appeared after we submitted our paper for review.

The recent work by Shen et al. (2019) is very closely related to our approach. They propose a class of methods based on the CGM and various variance reduction techniques for the non-convex finite-sum setting, including the SPIDER-FW. Besides, they also propose extensions that use second-order approximations. Finally, they provide simulation studies to compare empirical performance of different variants. We refer to this paper for a numerical comparison.

Hassani et al. (2019) introduce a novel variance reduced CGM method, but their work focuses primarily on the submodular maximization. Accordingly, they consider a more general expectation minimization template (the so-called non-oblivious setting) where the probability distribution depends on the decision variable  $x$  and may change during the optimization procedure. Therefore, the proposed method requires some further assumptions and modifications involving computations with the Hessian approximation. Finally, Zhang et al. (2019) consider a stochastic CGM approach with SPIDER in the distributed and quantized settings.

## 7. Concluding Remarks

We have proposed two novel FW-type methods based on the idea of blending the recent variance reduction technique SPIDER into FW and CGS frameworks. We have shown that the resulting methods enjoy superior oracle complexities in various convex and non-convex optimization templates. Extension of our framework for the strongly convex case is left open. Developing a well-tuned implementation, including one that incorporates parallel optimization, is an important piece of future work.



## Acknowledgements

VC was supported by the Swiss National Science Foundation (SNSF) under grant number 200021\_178865/1. VC has received funding for this project from the European Research Council (ERC) under the European Union’s Horizon 2020 research and innovation programme (grant agreement no 725594 - time-data). SS acknowledges support from an Amazon Research Award and the NSF-CAREER award (id 1846088). The authors thank Maria Vladarean for pointing out an error in the initial version of this work.

## References

- Agarwal, A. and Bottou, L. A lower bound for the optimization of finite sums. arXiv:1410.0723, 2014.
- Bertsekas, D. *Nonlinear Programming*. Athena Scientific, 2 edition, 1999.
- Defazio, A., Bach, F., and Lacoste-Julien, S. SAGA: A fast incremental gradient method with support for non-strongly convex composite objectives. In *Advances in Neural Information Processing Systems 28*, 2014.
- Fang, C., Li, C. J., Lin, Z., and Zhang, T. Spider: Near-optimal non-convex optimization via stochastic path-integrated differential estimator. In *Advances in Neural Information Processing Systems 31*, 2018.
- Frank, M. and Wolfe, P. An algorithm for quadratic programming. *Naval Research Logistics Quarterly*, 3:95–110, 1956.
- Hassani, H., Karbasi, A., Mokhtari, A., and Shen, Z. Stochastic conditional gradient++. arXiv:1902.06992, 2019.
- Hazan, E. Sparse approximate solutions to semidefinite programs. In *Proc. 8th Latin American Conf. Theoretical Informatics*, pp. 306–316, 2008.
- Hazan, E. and Kale, S. Projection-free online learning. In *Proc. 29th Int. Conf. Machine Learning*, 2012.
- Hazan, E. and Luo, H. Variance-reduced and projection-free stochastic optimization. In *Proc. 33rd Int. Conf. Machine Learning*, 2016.
- Jaggi, M. Revisiting Frank-Wolfe: Projection-free sparse convex optimization. In *Proc. 30th Int. Conf. Machine Learning*, 2013.
- Johnson, R. and Zhang, T. Accelerating stochastic gradient descent using predictive variance reduction. In *Advances in Neural Information Processing Systems 26*, 2013.
- Lacoste-Julien, S. Convergence rate of frank-wolfe for non-convex objectives. arXiv:1607.00345, 2016.
- Lacoste-Julien, S., Jaggi, M., Schmidt, M., and Pletscher, P. Block-coordinate Frank-Wolfe optimization for structural SVMs. In *Proc. 30th Int. Conf. Machine Learning*, 2013.
- Lan, G. The complexity of large-scale convex programming under a linear optimization oracle. arXiv:1309.5550v2, 2014.
- Lan, G. and Zhou, Y. Conditional gradient sliding for convex optimization. *SIAM J. Optim.*, 26(2):1379–1409, 2016.
- Lu, H. and Freund, R. M. Generalized stochastic frank-wolfe algorithm with stochastic” substitute” gradient for structured convex optimization. arXiv:1807.07680, 2018.
- Mahdavi, M., Zhang, L., and Jin, R. Mixed optimization for smooth functions. In *Advances in Neural Information Processing Systems 26*, 2013.
- Mokhtari, A., Hassani, H., and Karbasi, A. Stochastic conditional gradient methods: From convex minimization to submodular maximization. arXiv:1804.09554, 2018.
- Nemirovski, A. and Yudin, D. *Problem Complexity and Method Efficiency in Optimization*. John Wiley & Sons, 1983.
- Nesterov, Y. A method of solving a convex programming problem with convergence rate  $O(1/k^2)$ . *Soviet Mathematics Doklady*, 27(2):372–376, 1983.
- Nguyen, L. M., Liu, J., Scheinberg, K., and Takác, M. SARAH: A novel method for machine learning problems using stochastic recursive gradient. In *Proc. 34th Int. Conf. Machine Learning*, 2017.
- Qu, C., Li, Y., and Xu, H. Non-convex conditional gradient sliding. In *Proc. 35th Int. Conf. Machine Learning*, 2018.
- Ravi, S. N., Dinh, T., Lokhande, V. S. R., and Singh, V. Constrained deep learning using conditional gradient and applications in computer vision. arXiv:1803.06453v1, 2018.
- Reddi, S. J., Sra, S., Póczos, B., and Smola, A. Stochastic Frank-Wolfe methods for nonconvex optimization. In *54th Annual Allerton Conf. Communication, Control, and Computing*, pp. 1244–1251, 2016.
- Roux, N. L., Schmidt, M., and Bach, F. A stochastic gradient method with an exponential convergence rate for finite training sets. In *Advances in Neural Information Processing Systems 26*, 2012.
- Shen, Z., Fang, C., Zhao, P., Huang, J., and Qian, H. Complexities in projection-free stochastic non-convex minimization. In *Proc. 22nd Int. Conf. Artificial Intelligence and Statistics*, 2019.

Wang, K., Ji, K., Zhou, Y., Liang, Y., and Tarokh, V. Spider-Boost: A class of faster variance-reduced algorithms for nonconvex optimization. arXiv:1810.10690, 2018.

Yu, Y., Zhang, X., and Schuurmans, D. Generalized conditional gradient for sparse estimation. arXiv:1410.4828v1, 2014.

Zhang, M., Chen, L., Mokhtari, A., Hassani, H., and Karbasi, A. Quantized frank-wolfe: Communication-efficient distributed optimization. arXiv:1902.06332, 2019.

## A. Preliminaries

This section presents some known results from the existing literature, key to our analysis, for the sake of completeness.

The following Lemma from Fang et al. (2018) provides an error bound of the estimator  $v^k$  obtained by the SPIDER approach.

**Lemma 1** (Lemma 1 from (Fang et al., 2018), more specifically Eqn.(A.3) in its supplemental). *Suppose that  $\mathcal{S}_{t,k}$  is a subset that samples  $S_{t,k}$  iid realizations from the distribution  $\mathcal{P}$ . Let the stochastic estimator  $\nabla f_{\mathcal{S}_{t,k}}$  satisfy the averaged  $L$ -Lipschitz gradients condition from Section 3. Set the estimator  $v^k$  as*

$$v^k = \nabla f_{\mathcal{S}_{t,k}}(x^k) - \nabla f_{\mathcal{S}_{t,k}}(x^{k-1}) + v^{k-1}. \quad (\text{A.1})$$

Then, the following bound holds:

$$\mathbb{E}\|\nabla F(x^k) - v^k\|^2 \leq \frac{L^2}{S_{t,k}}\|x^k - x^{k-1}\|^2 + \|\nabla F(x^{k-1}) - v^{k-1}\|^2. \quad (\text{A.2})$$

The following Lemma draws a well-known bound on the variance in terms of the mini-batch size.

**Lemma 2** (Eqn.(3.5) from (Lan & Zhou, 2016), or Lemma 2 from (Reddi et al., 2016)). *Suppose that  $\mathcal{S}_{t,k}$  is a subset that samples  $S_{t,k}$  iid realizations from the distribution  $\mathcal{P}$ . Let the stochastic estimator  $\nabla f_{\mathcal{S}_{t,k}}$  satisfy the bounded variance condition from Section 3. Then, the following bound holds:*

$$\mathbb{E}\|\nabla f_{\mathcal{S}_{t,k}}(x) - \nabla F(x)\|^2 \leq \frac{\sigma^2}{S_{t,k}} \quad \forall x \in \Omega. \quad (\text{A.3})$$

Finally, we recall the convergence guarantees for the CndG procedure of CGS-type methods.

**Lemma 3** (Similar to Theorem 2.2 part (c) from (Lan & Zhou, 2016), or more generally Theorem 2 from (Jaggi, 2013)). *Remark that CndG procedure simply applies FW (with exact line-search) for the following projection subproblem:*

$$\min_{x \in \Omega} \frac{\beta}{2} \|x - u + \frac{1}{\beta} v\|^2. \quad (\text{A.4})$$

This problem is  $\beta$ -smooth, hence FW requires at most  $\mathcal{O}(\frac{4\beta D^2}{\alpha})$  iterations to satisfy the convergence criterion.

## B. Non-convex Conditional Gradient Sliding

In this section, we prove convergence of a CGS instance, and derive its oracle complexities in the non-convex settings. We also extend our results for SCGS.

### Proof of convergence for non-convex CGS

---

#### Algorithm 5 Conditional Gradient Sliding

---

**Input:**  $x^1 \in \Omega$   
**Set:**  $\alpha = \gamma LD^2$ ,  $\beta = \gamma L/2$ ,  $\gamma = 1/\sqrt{K}$   
**for**  $k = 1, 2, \dots, K$  **do**  
     Update  $z^k = y^k + \gamma(x^k - y^k)$   
      $x^{k+1} = \text{CndG}(x^k, \nabla F(z^k), \alpha, \beta)$   
     Update  $y^{k+1} = y^k + \gamma(x^{k+1} - y^k)$   
**end for**

---

**Theorem 9.** *Consider CGS algorithm with the parameters as described in Algorithm 5 (in the batch setting). Denote by  $y^{\text{out}}$  a random iterate  $y^k$  drawn uniformly random over all iterates of CGS. Then, the following bound holds:*

$$\mathbb{E}[\mathcal{G}(y^{\text{out}})] = \frac{\mathcal{E} + 3LD^2}{\sqrt{K}} \quad (\text{B.1})$$

**Corollary 9.** *The ifo and lmo complexities of CGS for achieving an  $\epsilon$ -solution in the non-convex minimization setting are*

$$\begin{aligned}\#(\text{ifo}) &= \mathcal{O}\left((\mathcal{E}^2 + L^2 D^4) \frac{n}{\epsilon^2}\right) \\ \#(\text{lmo}) &= \mathcal{O}\left((\mathcal{E}^2 + L^2 D^4) \frac{1}{\epsilon^2}\right)\end{aligned}\tag{B.2}$$

*Proof.* We start by the Taylor expansion and smoothness:

$$\begin{aligned}F(y^{k+1}) &\leq F(y^k) + \langle \nabla F(y^k), y^{k+1} - y^k \rangle + \frac{L}{2} \|y^{k+1} - y^k\|^2 \\ &= F(y^k) + \gamma \langle \nabla F(y^k), x^{k+1} - y^k \rangle + \gamma^2 \frac{L}{2} \|x^{k+1} - y^k\|^2 \\ &\leq F(y^k) + \gamma \langle \nabla F(y^k), x^{k+1} - y^k \rangle + \gamma^2 \frac{L}{2} D^2 \\ &= F(y^k) + \gamma \langle \nabla F(y^k), w_*^k - y^k \rangle + \gamma \langle \nabla F(y^k), x^{k+1} - w_*^k \rangle + \gamma^2 \frac{L}{2} D^2 \\ &= F(y^k) - \gamma \mathcal{G}(y^k) + \gamma \langle \nabla F(y^k), x^{k+1} - w_*^k \rangle + \gamma^2 \frac{L}{2} D^2\end{aligned}\tag{B.3}$$

where we define  $w_*^k = \arg \max_{x \in \Omega} \langle x, -\nabla F(y^k) \rangle$ .

We can equivalently write this inequality as

$$F(y^{k+1}) \leq F(y^k) - \gamma \mathcal{G}(y^k) + \gamma \langle \nabla F(y^k) - \nabla F(z^k), x^{k+1} - w_*^k \rangle + \gamma \langle \nabla F(z^k), x^{k+1} - w_*^k \rangle + \gamma^2 \frac{L}{2} D^2\tag{B.4}$$

Focus on the last inner-product term

$$\begin{aligned}\gamma \langle \nabla F(z^k), x^{k+1} - w_*^k \rangle &= \gamma \langle \nabla F(z^k) + \beta(x^{k+1} - x^k), x^{k+1} - w_*^k \rangle - \gamma \beta \langle x^{k+1} - x^k, x^{k+1} - w_*^k \rangle \\ &\leq \gamma \alpha - \gamma \beta \langle x^{k+1} - x^k, x^{k+1} - w_*^k \rangle \\ &\leq \gamma \alpha + \gamma \beta D^2\end{aligned}\tag{B.5}$$

where the first inequality follows from the role of  $\alpha$  in CndG, and the second one from Cauchy-Schwarz.

Combining these two inequalities, we obtain

$$\begin{aligned}F(y^{k+1}) &\leq F(y^k) - \gamma \mathcal{G}(y^k) + \gamma \langle \nabla F(y^k) - \nabla F(z^k), x^{k+1} - w_*^k \rangle + \gamma \alpha + \gamma \beta D^2 + \gamma^2 \frac{L}{2} D^2 \\ &\leq F(y^k) - \gamma \mathcal{G}(y^k) + \gamma D \|\nabla F(y^k) - \nabla F(z^k)\| + \gamma \alpha + \gamma \beta D^2 + \gamma^2 \frac{L}{2} D^2 \\ &\leq F(y^k) - \gamma \mathcal{G}(y^k) + \gamma^2 L D \|x^k - y^k\| + \gamma \alpha + \gamma \beta D^2 + \gamma^2 \frac{L}{2} D^2 \\ &\leq F(y^k) - \frac{1}{\sqrt{K}} \mathcal{G}(y^k) + \frac{3LD^2}{K}\end{aligned}\tag{B.6}$$

Taking expectation of both sides, rearranging, and summing over all iterations, we obtain

$$\frac{1}{\sqrt{K}} \sum_{k=1}^K \mathbb{E}[\mathcal{G}(y^k)] \leq F(\bar{x}^1) - \mathbb{E}[F(y^K)] + \sum_{k=1}^K \frac{3LD^2}{K} \leq F(\bar{x}^1) - F(x^*) + 3LD^2\tag{B.7}$$

Hence, by definition of  $y^{\text{out}}$ , we get

$$\mathbb{E}[\mathcal{G}(y^{\text{out}})] \leq \frac{F(\bar{x}^1) - F(x^*)}{\sqrt{K}} + \frac{3LD^2}{\sqrt{K}} = \frac{\mathcal{E} + 3LD^2}{\sqrt{K}}\tag{B.8}$$

This completes the convergence rate proof.

To get an  $\epsilon$ -solution, we set the number of iterations  $K_\epsilon$  such that

$$\mathbb{E}[\mathcal{G}(y^{\text{out}})] \leq \frac{\mathcal{E} + 3LD^2}{\sqrt{K_\epsilon}} \leq \epsilon. \quad (\text{B.9})$$

Hence, we can calculate the *lmo* complexity using Lemma 3 as

$$\#(\text{lmo}) = \mathcal{O}\left(K_\epsilon \frac{4\beta D^2}{\alpha}\right) = \mathcal{O}\left((\mathcal{E}^2 + L^2 D^4) \frac{1}{\epsilon^2}\right) \quad (\text{B.10})$$

which completes the proof.  $\square$

### Proof of convergence for the non-convex SCGS

---

#### Algorithm 6 Stochastic Conditional Gradient Sliding

---

**Input:**  $x^1 \in \Omega$   
**Set:**  $\alpha = \gamma LD^2$ ,  $\beta = \gamma L/2$ ,  $\gamma = 1/\sqrt{K}$   
**for**  $k = 1, 2, \dots, K$  **do**  
     Update  $z^k = y^k + \gamma(x^k - y^k)$   
     Draw  $K$  samples  $\mathcal{S}_k$   
      $x^{k+1} = \text{CndG}(x^k, \nabla f_{\mathcal{S}_k}(z^k), \alpha, \beta)$   
     Update  $y^{k+1} = y^k + \gamma(x^{k+1} - y^k)$   
**end for**

---

**Theorem 10.** Consider the SCGS algorithm with the parameters as described in the Algorithm 6. Assume that the conditions for the expectation minimization setting from Section 3 hold. Denote by  $y^{\text{out}}$  a random iterate  $y^k$  drawn uniformly random over all iterates of the SCGS. Then, the following bound holds:

$$\mathbb{E}[\mathcal{G}(y^{\text{out}})] = \frac{\mathcal{E} + \sigma D + 3LD^2}{\sqrt{K}} \quad (\text{B.11})$$

**Corollary 10.** The *sfo* and *lmo* complexities of SCGS for achieving an  $\epsilon$ -solution in the non-convex minimization setting are

$$\#(\text{sfo}) = \mathcal{O}\left((\mathcal{E} + \sigma D + LD^2)^4 \frac{1}{\epsilon^4}\right) \quad \text{and} \quad \#(\text{lmo}) = \mathcal{O}\left((\mathcal{E} + \sigma D + LD^2)^2 \frac{1}{\epsilon^2}\right)$$

*Proof.* From (B.3), and similar to (B.4) and (B.5), we can show

$$F(y^{k+1}) \leq F(y^k) - \gamma \mathcal{G}(y^k) + \gamma \langle \nabla F(y^k) - \nabla f_{\mathcal{S}_k}(z^k), x^{k+1} - w_\star^k \rangle + \gamma \alpha + \gamma \beta D^2 + \gamma^2 \frac{L}{2} D^2 \quad (\text{B.12})$$

where  $w_\star^k = \arg \max_{x \in \Omega} \langle x, -\nabla F(y^k) \rangle$ .

Focusing on the inner-product term, we get the following bound:

$$\begin{aligned} \gamma \langle \nabla F(y^k) - \nabla f_{\mathcal{S}_k}(z^k), x^{k+1} - w_\star^k \rangle &\leq \gamma \|\nabla F(y^k) - \nabla f_{\mathcal{S}_k}(z^k)\| \|x^{k+1} - w_\star^k\| \\ &\leq \gamma D \|\nabla F(y^k) - \nabla f_{\mathcal{S}_k}(z^k)\| \\ &\leq \gamma D \|\nabla F(y^k) - \nabla F(z^k)\| + \gamma D \|\nabla F(z^k) - \nabla f_{\mathcal{S}_k}(z^k)\| \\ &\leq \gamma LD \|y^k - z^k\| + \gamma D \|\nabla F(z^k) - \nabla f_{\mathcal{S}_k}(z^k)\| \\ &= \gamma^2 LD \|x^k - y^k\| + \gamma D \|\nabla F(z^k) - \nabla f_{\mathcal{S}_k}(z^k)\| \\ &\leq \gamma^2 LD^2 + \gamma D \|\nabla F(z^k) - \nabla f_{\mathcal{S}_k}(z^k)\| \end{aligned} \quad (\text{B.13})$$

Substituting back and taking expectations, we obtain

$$\begin{aligned} \mathbb{E}[F(y^{k+1})] &\leq \mathbb{E}[F(y^k)] - \frac{1}{\sqrt{K}} \mathbb{E}[\mathcal{G}(y^k)] + \frac{D}{\sqrt{K}} \mathbb{E} \|\nabla F(z^k) - \nabla f_{\mathcal{S}_k}(z^k)\| + 3\gamma^2 LD^2 \\ &= \mathbb{E}[F(y^k)] - \frac{1}{\sqrt{K}} \mathbb{E}[\mathcal{G}(y^k)] + \frac{D}{\sqrt{K}} \mathbb{E} \|\nabla F(z^k) - \nabla f_{\mathcal{S}_k}(z^k)\| + \frac{3LD^2}{K} \end{aligned} \quad (\text{B.14})$$

Now, we use Lemma 2 with the Jensen's inequality to obtain

$$\mathbb{E}[F(y^{k+1})] \leq \mathbb{E}[F(y^k)] - \frac{1}{\sqrt{K}} \mathbb{E}[\mathcal{G}(y^k)] + \frac{\sigma D + 3LD^2}{K} \quad (\text{B.15})$$

From here, we follow the same steps as in the proof of CGS and get (B.11).

Then, to achieve an  $\epsilon$ -solution, we can calculate *sfo* complexity as

$$\#(\text{sfo}) = \sum_{k=1}^{K_\epsilon} K_\epsilon = K_\epsilon^2 = \mathcal{O}\left((\mathcal{E} + \sigma D + LD^2)^4 \frac{1}{\epsilon^4}\right) \quad (\text{B.16})$$

Finally, *lmo* complexity can be found using Lemma 3

$$\#(\text{lmo}) = \mathcal{O}\left(K_\epsilon \frac{4\beta D^2}{\alpha}\right) = \mathcal{O}\left((\mathcal{E} + \sigma D + LD^2)^2 \frac{1}{\epsilon^2}\right) \quad (\text{B.17})$$

This completes the proof.  $\square$

### C. Proofs for SPIDER-FW

**Lemma 4.** Suppose that the assumptions listed in Section 3 hold. Then, for  $k = 1, \dots, K_t$ , we have the following bounds:

$$\text{Convex finite-sum} \quad \mathbb{E}\|\nabla F(x^k) - v^k\| \leq 2LD/K_t \quad (\text{C.1})$$

$$\text{Convex expectation} \quad \mathbb{E}\|\nabla F(x^k) - v^k\| \leq 3LD/K_t \quad (\text{C.2})$$

$$\text{Non-convex finite-sum} \quad \mathbb{E}\|\nabla F(x^k) - v^k\| \leq LD/\sqrt{TK} \quad (\text{C.3})$$

$$\text{Non-convex expectation} \quad \mathbb{E}\|\nabla F(x^k) - v^k\| \leq LD/\sqrt{TK} + \epsilon/2 \quad (\text{C.4})$$

*Proof.* From Lemma 1, we have the following inequality for  $k = 2, 3, \dots, K_t$ :

$$\mathbb{E}\|\nabla F(x^k) - v^k\|^2 \leq \frac{L^2}{S_{t,k}} \|x^k - x^{k-1}\|^2 + \|\nabla F(x^{k-1}) - v^{k-1}\|^2. \quad (\text{C.5})$$

By definition,  $\|x^{k-1} - x^k\|^2 = \|\eta_{t,k-1}(w^{k-1} - x^{k-1})\|^2 \leq \eta_{t,k-1}^2 D^2$ . Hence, we get

$$\mathbb{E}\|\nabla F(x^k) - v^k\|^2 \leq \frac{\eta_{t,k-1}^2 L^2 D^2}{S_{t,k}} + \|\nabla F(x^{k-1}) - v^{k-1}\|^2. \quad (\text{C.6})$$

*Convex finite-sum:*

We take the telescopic sum of (C.6) from  $i = 2$  to  $k$

$$\mathbb{E}\|\nabla F(x^k) - v^k\|^2 \leq \sum_{i=2}^k \frac{\eta_{t,i-1}^2 L^2 D^2}{S_{t,i}} + \underbrace{\|\nabla F(x^1) - v^1\|^2}_0 \leq \frac{4L^2 D^2}{K_t} \sum_{i=2}^k \frac{1}{(s_{t,i-1} + 1)^2} \quad (\text{C.7})$$

By definition of  $s_{t,k}$ , for any  $i \geq 2$  we have

$$s_{t,i-1} + 1 = K_t + i - 1 \geq K_t. \quad (\text{C.8})$$

Hence, we get

$$\mathbb{E}\|\nabla F(x^k) - v^k\|^2 \leq \frac{4L^2 D^2}{K_t^3} \sum_{i=2}^k 1 \leq \frac{4L^2 D^2}{K_t^3} k \leq \frac{4L^2 D^2}{K_t^2}. \quad (\text{C.9})$$

By using Jensen's inequality, we get (C.1).

Convex expectation minimization:

We take the telescopic sum of (C.6) from  $i = 2$  to  $k$

$$\mathbb{E}\|\nabla F(x^k) - v^k\|^2 \leq \frac{4L^2D^2}{K_t^2} + \|\nabla F(x^1) - v^1\|^2 \leq \frac{9L^2D^2}{K_t^2}, \quad (\text{C.10})$$

where the bound on  $\|\nabla F(x^1) - v^1\|^2$  follows from Lemma 2 as

$$\|\nabla F(x^1) - v^1\|^2 \leq \frac{\sigma^2}{Q_t} \leq \frac{5L^2D^2}{K_t^2}. \quad (\text{C.11})$$

We get (C.2) by using Jensen's inequality.

Non-convex finite-sum:

We take the telescopic sum of (C.6) from  $i = 2$  to  $k$

$$\mathbb{E}\|\nabla F(x^k) - v^k\|^2 \leq \sum_{i=2}^k \frac{\eta^2 L^2 D^2}{S} + \underbrace{\|\nabla F(x^1) - v^1\|^2}_0 \leq \frac{L^2 D^2}{TK^2} \sum_{i=2}^k 1 \leq \frac{L^2 D^2}{TK} \quad (\text{C.12})$$

We get (C.3) by using Jensen's inequality.

Non-convex expectation minimization:

We take the telescopic sum of (C.6) from  $i = 2$  to  $k$

$$\mathbb{E}\|\nabla F(x^k) - v^k\|^2 \leq \frac{L^2 D^2}{TK} + \|\nabla F(x^1) - v^1\|^2 \leq \frac{L^2 D^2}{TK} + \frac{\epsilon^2}{4} \quad (\text{C.13})$$

where the bound on  $\|\nabla F(x^1) - v^1\|^2$  follows from Lemma 2 as

$$\|\nabla F(x^1) - v^1\|^2 \leq \frac{\sigma^2}{Q_t} \leq \frac{\epsilon^2}{4}. \quad (\text{C.14})$$

We get (C.4) by using Jensen's inequality. □

### Proof for Theorem 1 and Corollary 1

We start by the Taylor expansion and smoothness:

$$\begin{aligned} F(x^{k+1}) &\leq F(x^k) + \langle \nabla F(x^k), x^{k+1} - x^k \rangle + \frac{L}{2} \|x^{k+1} - x^k\|^2 \\ &\leq F(x^k) + \eta_{t,k} \langle \nabla F(x^k), w^k - x^k \rangle + \eta_{t,k}^2 \frac{L}{2} D^2 \\ &= F(x^k) + \eta_{t,k} \langle v^k, w^k - x^k \rangle + \eta_{t,k} \langle \nabla F(x^k) - v^k, w^k - x^k \rangle + \eta_{t,k}^2 \frac{L}{2} D^2 \end{aligned} \quad (\text{C.15})$$

By definition of  $w^k$ , we have

$$\langle v^k, w^k - x^k \rangle = \min_{x \in \Omega} \langle v^k, x - x^k \rangle \leq \langle v^k, x^* - x^k \rangle \quad (\text{C.16})$$

Substituting this inequality back, and rearranging, we get

$$\begin{aligned} F(x^{k+1}) &\leq F(x^k) + \eta_{t,k} \langle v^k, x^* - x^k \rangle + \eta_{t,k} \langle \nabla F(x^k) - v^k, w^k - x^k \rangle + \eta_{t,k}^2 \frac{L}{2} D^2 \\ &= F(x^k) + \eta_{t,k} \langle \nabla F(x^k), x^* - x^k \rangle + \eta_{t,k} \langle \nabla F(x^k) - v^k, w^k - x^* \rangle + \eta_{t,k}^2 \frac{L}{2} D^2 \end{aligned} \quad (\text{C.17})$$

From the convexity of  $F$ , we know

$$\langle \nabla F(x^k), x^* - x^k \rangle \leq F^* - F(x^k) \quad (\text{C.18})$$

and by using Cauchy-Schwarz, we have

$$\langle \nabla F(x^k) - v^k, w^k - x^* \rangle \leq \|\nabla F(x^k) - v^k\| \|w^k - x^*\| \leq \|\nabla F(x^k) - v^k\| D \quad (\text{C.19})$$

Putting (C.18) and (C.19) back into (C.17), and subtracting  $F^*$  from both sides, we obtain:

$$F(x^{k+1}) - F^* \leq (1 - \eta_{t,k})(F(x^k) - F^*) + \eta_{t,k} D \|\nabla F(x^k) - v^k\| + \eta_{t,k}^2 \frac{L}{2} D^2 \quad (\text{C.20})$$

The, we take expectation of both sides and use (C.1) to get

$$\begin{aligned} \mathbb{E}[F(x^{k+1})] - F^* &\leq (1 - \eta_{t,k})(\mathbb{E}[F(x^k)] - F^*) + \eta_{t,k} D \mathbb{E} \|\nabla F(x^k) - v^k\| + \eta_{t,k}^2 \frac{L}{2} D^2 \\ &\leq (1 - \eta_{t,k})(\mathbb{E}[F(x^k)] - F^*) + \eta_{t,k} \frac{2LD^2}{K_t} + \eta_{t,k}^2 \frac{L}{2} D^2 \end{aligned} \quad (\text{C.21})$$

Telescopic sum of this inequality over  $(t, k)$  pairs gives

$$\mathbb{E}[F(x^{k+1})] - F^* \leq \sum_{(\tau,i)} \left( \eta_{\tau,i} \frac{2LD^2}{K_\tau} + \eta_{\tau,i}^2 \frac{L}{2} D^2 \right) \prod_{(\tau',j)=(\tau,i)}^{(t,k)} (1 - \eta_{\tau',j}) + (\mathbb{E}[F(x^{1,1})] - F^*) \prod_{(\tau,i)} (1 - \eta_{\tau,i}). \quad (\text{C.22})$$

The last term vanishes due to 0 factor ( $\eta_{1,1} = 1$ ). Remark that

$$\prod_{(\tau',j)=(\tau,i)}^{(t,k)} (1 - \eta_{\tau',j}) = \prod_{r=i}^{K_\tau} \frac{s_{\tau,r} - 1}{s_{\tau,r} + 1} \prod_{\tau'=\tau+1}^{t-1} \prod_{j=1}^{K_{\tau'}} \frac{s_{\tau',j} - 1}{s_{\tau',j} + 1} \prod_{\ell=1}^k \frac{s_{t,\ell} - 1}{s_{t,\ell} + 1} = \frac{(s_{\tau,i} - 1)s_{\tau,i}}{s_{t,k}(s_{t,k} + 1)} \quad (\text{C.23})$$

Combining these, we get

$$\mathbb{E}[F(x^{k+1})] - F^* \leq \sum_{(\tau,i)} \left( \eta_{\tau,i} \frac{2LD^2}{K_\tau} + \eta_{\tau,i}^2 \frac{L}{2} D^2 \right) \frac{(s_{\tau,i} - 1)s_{\tau,i}}{s_{t,k}(s_{t,k} + 1)}. \quad (\text{C.24})$$

We focus on the individual terms:

$$\sum_{(\tau,i)} \eta_{\tau,i} \frac{2LD^2}{K_\tau} \frac{(s_{\tau,i} - 1)s_{\tau,i}}{s_{t,k}(s_{t,k} + 1)} \leq \frac{8LD^2}{s_{t,k}(s_{t,k} + 1)} \sum_{(\tau,i)} 1 \leq \frac{8LD^2}{s_{t,k} + 1} \quad (\text{C.25})$$

$$\sum_{(\tau,i)} \eta_{\tau,i}^2 \frac{L}{2} D^2 \frac{(s_{\tau,i} - 1)s_{\tau,i}}{s_{t,k}(s_{t,k} + 1)} \leq \frac{2LD^2}{s_{t,k}(s_{t,k} + 1)} \sum_{(\tau,i)} 1 \leq \frac{2LD^2}{s_{t,k} + 1} \quad (\text{C.26})$$

We proved the convergence rate:

$$\mathbb{E}[F(x^{k+1})] - F^* \leq \frac{10LD^2}{s_{t,k} + 1} \quad (\text{C.27})$$

To get  $\epsilon$ -solution, we set the number of outer iterations  $T_\epsilon$  such that

$$\mathbb{E}[F(\bar{x}^{T_\epsilon})] - F^* \leq \frac{10LD^2}{K_{T_\epsilon}} \leq \epsilon. \quad (\text{C.28})$$

Then, it is sufficient to choose

$$T_\epsilon = \log_2 \left( \frac{10LD^2}{\epsilon} \right) + 1. \quad (\text{C.29})$$

Then, to achieve  $(1 - \epsilon)$  accuracy, we can calculate the *ifo* complexity as

$$\begin{aligned} \#(ifo) &= \sum_{t=1}^{T_\epsilon} \left( Q_t + \sum_{k=2}^{K_t} S_{t,k} \right) = \sum_{t=1}^{T_\epsilon} \left( n + \sum_{k=2}^{K_t} 2^{t-1} \right) \leq \sum_{t=1}^{T_\epsilon} \left( n + 2^{2(t-1)} \right) = \mathcal{O}(nT_\epsilon + 2^{2T_\epsilon}) \\ &= \mathcal{O} \left( n \ln \left( \frac{LD^2}{\epsilon} \right) + \frac{L^2 D^4}{\epsilon^2} \right) \end{aligned} \quad (\text{C.30})$$

and the *lmo* complexity as

$$\#(lmo) = \sum_{t=1}^{T_\epsilon} K_t \leq 2K_{T_\epsilon} = 2^{T_\epsilon} = \mathcal{O} \left( \frac{LD^2}{\epsilon} \right). \quad (\text{C.31})$$



**Proof for Theorem 2 and Corollary 2**

Proof is similar to that for finite-sum setting, but we use (C.2) instead of (C.1) at (C.21), hence the constants change:

$$\mathbb{E}[F(x^{k+1})] - F^* \leq \frac{14LD^2}{s_{t,k} + 1}. \quad (\text{C.32})$$

To get  $\epsilon$ -solution, we set the number of outer iterations  $T_\epsilon$  as

$$T_\epsilon = \log_2 \left( \frac{14LD^2}{\epsilon} \right) + 1. \quad (\text{C.33})$$

Then, to achieve an  $\epsilon$ -solution, we can calculate *sfo* complexity as

$$\#(\text{sfo}) = \sum_{t=1}^{T_\epsilon} \left( Q_t + \sum_{k=2}^{K_t} S_{t,k} \right) \leq \sum_{t=1}^{T_\epsilon} \left( \lceil \frac{\sigma^2 K_t^2}{5L^2 D^2} \rceil + K_t^2 \right) = \mathcal{O} \left( \frac{\sigma^2 D^2}{\epsilon^2} + \frac{L^2 D^4}{\epsilon^2} \right). \quad (\text{C.34})$$

The *lmo* complexity is the same as the finite-sum case.

**Proof for Theorem 3 and Corollary 3**

We start by the Taylor expansion and smoothness:

$$\begin{aligned} F(x^{k+1}) &\leq F(x^k) + \langle \nabla F(x^k), x^{k+1} - x^k \rangle + \frac{L}{2} \|x^{k+1} - x^k\|^2 \\ &\leq F(x^k) + \eta \langle \nabla F(x^k), w^k - x^k \rangle + \eta^2 \frac{L}{2} D^2 \\ &= F(x^k) + \eta \langle v^k, w^k - x^k \rangle + \eta \langle \nabla F(x^k) - v^k, w^k - x^k \rangle + \eta^2 \frac{L}{2} D^2 \\ &\leq F(x^k) + \eta \langle v^k, w_\star^k - x^k \rangle + \eta \langle \nabla F(x^k) - v^k, w^k - x^k \rangle + \eta^2 \frac{L}{2} D^2 \\ &= F(x^k) + \eta \langle \nabla F(x^k), w_\star^k - x^k \rangle + \eta \langle \nabla F(x^k) - v^k, w^k - w_\star^k \rangle + \eta^2 \frac{L}{2} D^2 \\ &\leq F(x^k) - \eta \mathcal{G}(x^k) + \eta D \|\nabla F(x^k) - v^k\| + \eta^2 \frac{L}{2} D^2 \end{aligned} \quad (\text{C.35})$$

where  $w_\star^k = \arg \max_{x \in \Omega} \langle x, -\nabla F(x^k) \rangle$ . Taking expectation of both sides and using (C.3), we get

$$\mathbb{E}[F(x^{k+1})] \leq \mathbb{E}[F(x^k)] - \frac{1}{\sqrt{TK}} \mathbb{E}[\mathcal{G}(x^k)] + \frac{3LD^2}{2TK} \quad (\text{C.36})$$

Rearranging, and summing over all  $(t, k)$  pairs up to  $(T, K)$ , we obtain

$$\frac{1}{\sqrt{TK}} \sum_{(\tau, i)}^{(T, K)} \mathbb{E}[\mathcal{G}(x^{\tau, i})] \leq F(\bar{x}^1) - \mathbb{E}[F(x^{T, K})] + \sum_{(\tau, i)}^{(T, K)} \frac{3LD^2}{2TK} \leq F(\bar{x}^1) - F(x^*) + \frac{3LD^2}{2}. \quad (\text{C.37})$$

Hence, by definition of  $x^{\text{out}}$ , we have

$$\mathbb{E}[\mathcal{G}(x^{\text{out}})] \leq \frac{F(\bar{x}^1) - F(x^*)}{\sqrt{TK}} + \frac{3LD^2}{2\sqrt{TK}} = \frac{2\mathcal{E} + 3LD^2}{2\sqrt{TK}} \quad (\text{C.38})$$

This completes the convergence rate proof.

To get  $\epsilon$ -solution, we set the number of outer iterations  $T_\epsilon$  such that

$$\mathbb{E}[\mathcal{G}(x^{\text{out}})] \leq \frac{2\mathcal{E} + 3LD^2}{2\sqrt{T_\epsilon K}} \leq \epsilon. \quad (\text{C.39})$$

Hence, it suffices to choose

$$T_\epsilon = \left( \frac{2\mathcal{E} + 3LD^2}{2\epsilon\sqrt{K}} \right)^2 = \frac{(2\mathcal{E} + 3LD^2)^2}{4\epsilon^2 K}. \quad (\text{C.40})$$

Then, to achieve an  $\epsilon$ -solution, we can calculate *ifo* complexity as

$$\begin{aligned} \#(\text{ifo}) &= \sum_{t=1}^{T_\epsilon} \left( Q_t + \sum_{k=2}^K S_{t,k} \right) = \sum_{t=1}^{T_\epsilon} \left( n + \sum_{k=2}^K K \right) \leq (n + K^2)T_\epsilon \leq \sqrt{n} \frac{3(2\mathcal{E} + 3LD^2)^2}{4\epsilon^2} \\ &= \mathcal{O} \left( (\mathcal{E}^2 + L^2 D^4) \frac{\sqrt{n}}{\epsilon^2} \right) \end{aligned} \quad (\text{C.41})$$

and the *lmo* complexity as

$$\#(\text{lmo}) = \sum_{t=1}^{T_\epsilon} K = KT_\epsilon = \frac{(2\mathcal{E} + 3LD^2)^2}{4\epsilon^2} = \mathcal{O} \left( (\mathcal{E}^2 + L^2 D^4) \frac{1}{\epsilon^2} \right) \quad (\text{C.42})$$

#### Proof for Theorem 4 and Corollary 4

Proof is similar to that for non-convex finite-sum setting, but we use (C.4) instead of (C.3) at (C.35), and we get

$$\mathbb{E}[F(x^{k+1})] \leq \mathbb{E}[F(x^k)] - \frac{1}{\sqrt{TK}} \mathbb{E}[\mathcal{G}(x^k)] + \frac{3LD^2}{2TK} + \frac{\epsilon}{2\sqrt{TK}} \quad (\text{C.43})$$

Rearranging, and summing over all  $(t, k)$  pairs up to  $(T, K)$ , we get

$$\begin{aligned} \frac{1}{\sqrt{TK}} \sum_{(\tau, i)}^{(T, K)} \mathbb{E}[\mathcal{G}(x^{\tau, i})] &\leq F(\bar{x}^1) - \mathbb{E}[F(x^{T, K})] + \sum_{(\tau, i)}^{(T, K)} \left( \frac{3LD^2}{2TK} + \frac{\epsilon}{2\sqrt{TK}} \right) \\ &\leq F(\bar{x}^1) - F(x^*) + \frac{3LD^2}{2} + \frac{\epsilon\sqrt{TK}}{2}. \end{aligned} \quad (\text{C.44})$$

Hence, by definition of  $x^{\text{out}}$ , we have

$$\mathbb{E}[\mathcal{G}(x^{\text{out}})] \leq \frac{F(\bar{x}^1) - F(x^*)}{\sqrt{TK}} + \frac{3LD^2}{2\sqrt{TK}} + \frac{\epsilon}{2} = \frac{2\mathcal{E} + 3LD^2}{2\sqrt{TK}} + \frac{\epsilon}{2} \quad (\text{C.45})$$

This completes the convergence proof.

To get  $\epsilon$ -solution, we set the number of outer iterations  $T_\epsilon$  such that

$$\mathbb{E}[\mathcal{G}(x^{\text{out}})] \leq \frac{2\mathcal{E} + 3LD^2}{2\sqrt{T_\epsilon K}} + \frac{\epsilon}{2} \leq \epsilon. \quad (\text{C.46})$$

Then, to achieve  $(1 - \epsilon)$  accuracy, we can calculate *sfo* complexity as

$$\begin{aligned} \#(\text{sfo}) &= \sum_{t=1}^{T_\epsilon} \left( Q_t + \sum_{k=2}^K S_{t,k} \right) = \sum_{t=1}^{T_\epsilon} \left( n + \sum_{k=2}^K K \right) \leq ([4(\sigma/\epsilon)^2] + ([\sigma/\epsilon])^2)T_\epsilon \\ &= \mathcal{O} \left( (\mathcal{E}^2 + L^2 D^4) \frac{\sigma}{\epsilon^3} \right) \end{aligned} \quad (\text{C.47})$$

Finally, *lmo* complexity is the same as the non-convex finite-sum case.

## D. Proofs for SPIDER-CGS

**Lemma 5.** Suppose that the assumptions listed in Section 3 hold. Then, for  $k = 1, \dots, K_t$ , we have the following bounds:

$$\text{Convex finite-sum} \quad \mathbb{E} \|\nabla F(z^k) - v^k\| \leq 2\sqrt{2}LD/(s_{t,k} + 1)^2 \quad (\text{D.1})$$

$$\text{Convex expectation} \quad \mathbb{E} \|\nabla F(z^k) - v^k\| \leq 3LD/(s_{t,k} + 1)^2 \quad (\text{D.2})$$

$$\text{Non-convex finite-sum} \quad \mathbb{E} \|\nabla F(z^k) - v^k\| \leq 2LD/\sqrt{TK} \quad (\text{D.3})$$

$$\text{Non-convex expectation} \quad \mathbb{E} \|\nabla F(z^k) - v^k\| \leq 2LD/\sqrt{TK} + \epsilon/2 \quad (\text{D.4})$$

*Proof.* From Lemma 1, we have the following inequality for all  $k = 2, \dots, K_t$ :

$$\mathbb{E}\|\nabla F(z^k) - v^k\|^2 \leq \frac{L^2}{S_{t,k}} \|z^k - z^{k-1}\|^2 + \|\nabla F(z^{k-1}) - v^{k-1}\|^2. \quad (\text{D.5})$$

By definition,

$$\begin{aligned} z^k &= y^k + \gamma_{t,k}(x^k - y^k) = y^{k-1} + \gamma_{t,k-1}(x^k - y^{k-1}) + \gamma_{t,k}(x^k - y^k) \\ \& \quad z^{k-1} &= y^{k-1} + \gamma_{t,k-1}(x^{k-1} - y^{k-1}) \\ \implies \|z^k - z^{k-1}\|^2 &= \|\gamma_{t,k-1}(x^k - x^{k-1}) + \gamma_{t,k}(x^k - y^k)\|^2 \\ &= \gamma_{t,k-1}^2 \|x^k - x^{k-1}\|^2 + \gamma_{t,k}^2 \|x^k - y^k\|^2 + 2\gamma_{t,k-1}\gamma_{t,k} \langle x^k - y^k, x^k - x^{k-1} \rangle \\ &\leq \gamma_{t,k-1}^2 \|x^k - x^{k-1}\|^2 + \gamma_{t,k-1}^2 \|x^k - y^k\|^2 + 2\gamma_{t,k-1}^2 \|x^k - y^k\| \|x^k - x^{k-1}\| \\ &\leq 4\gamma_{t,k-1}^2 D^2 \end{aligned} \quad (\text{D.6})$$

Substituting into (D.5), we get

$$\mathbb{E}\|\nabla F(z^k) - v^k\|^2 \leq \frac{4\gamma_{t,k-1}^2 L^2 D^2}{S_{t,k}} + \|\nabla F(z^{k-1}) - v^{k-1}\|^2. \quad (\text{D.7})$$

Convex finite-sum:

We take the telescopic sum of (D.7) from  $i = 2$  to  $k$

$$\begin{aligned} \mathbb{E}\|\nabla F(z^k) - v^k\|^2 &\leq \sum_{i=2}^k \frac{4\gamma_{t,i-1}^2 L^2 D^2}{S_{t,i}} + \underbrace{\|\nabla F(z^1) - v^1\|^2}_0, \text{ since we take full batch} \\ &= \frac{4L^2 D^2}{9K_t(s_{t,K_t} + 1)^2} \sum_{i=2}^k \frac{9}{(s_{t,i-1} + 2)^2} = \frac{4L^2 D^2}{K_t(s_{t,K_t} + 1)^2} \sum_{i=2}^k \frac{1}{(s_{t,i} + 1)^2} \end{aligned} \quad (\text{D.8})$$

Clearly,  $s_{t,i} + 1 \geq \frac{s_{t,K_t} + 1}{\sqrt{2}}$  for all  $2 \leq i \leq K_t$ . Hence, we get

$$\mathbb{E}\|\nabla F(z^k) - v^k\|^2 \leq \frac{8L^2 D^2}{K_t(s_{t,K_t} + 1)^2} \sum_{i=2}^k \frac{1}{(s_{t,K_t} + 1)^2} \leq \frac{8L^2 D^2}{(s_{t,k} + 1)^4} \quad (\text{D.9})$$

Convex expectation minimization:

We take the telescopic sum of (D.7) from  $i = 2$  to  $k$

$$\mathbb{E}\|\nabla F(z^k) - v^k\|^2 \leq \frac{8L^2 D^2}{(s_{t,k} + 1)^4} + \|\nabla F(z^1) - v^1\|^2 \leq \frac{8L^2 D^2}{(s_{t,k} + 1)^4} + \frac{\sigma^2}{Q_t} \leq \frac{9L^2 D^2}{(s_{t,k} + 1)^4} \quad (\text{D.10})$$

where the bound on  $\|\nabla F(x^1) - v^1\|^2$  follows from Lemma 2.

Non-convex finite-sum:

We take the telescopic sum of (D.7) from  $i = 2$  to  $k$

$$\mathbb{E}\|\nabla F(z^k) - v^k\|^2 \leq \sum_{i=2}^k \frac{4\gamma^2 L^2 D^2}{S} + \underbrace{\|\nabla F(z^1) - v^1\|^2}_0, \text{ since we take full batch} \leq \sum_{i=2}^k \frac{4L^2 D^2}{TK^2} \leq \frac{4L^2 D^2}{TK} \quad (\text{D.11})$$

We get (C.3) by using Jensen's inequality.

*Non-convex expectation minimization:*

We take the telescopic sum of (D.7) from  $i = 2$  to  $k$

$$\mathbb{E}\|\nabla F(z^k) - v^k\|^2 \leq \sum_{i=2}^k \frac{4\gamma^2 L^2 D^2}{S} + \|\nabla F(z^1) - v^1\|^2 \leq \frac{4L^2 D^2}{TK} + \frac{\epsilon^2}{4} \quad (\text{D.12})$$

where the bound on  $\|\nabla F(z^1) - v^1\|^2$  follows from Lemma 2. We get (C.4) by using Jensen's inequality.  $\square$

### Proof for Theorem 5 and Corollary 5

We start by Taylor expansion and smoothness:

$$\begin{aligned} F(y^{k+1}) &\leq F(z^k) + \langle \nabla F(z^k), y^{k+1} - z^k \rangle + \frac{L}{2} \|y^{k+1} - z^k\|^2 \\ &= F(z^k) + \langle \nabla F(z^k), y^k - z^k \rangle + \gamma_{t,k} \langle \nabla F(z^k), x^{k+1} - x^* \rangle \\ &\quad + \gamma_{t,k} \langle \nabla F(z^k), x^* - y^k \rangle + \frac{L\gamma_{t,k}^2}{2} \|x^{k+1} - x^k\|^2 \\ &= (1 - \gamma_{t,k})(F(z^k) + \langle \nabla F(z^k), y^k - z^k \rangle) + \gamma_{t,k}(F(z^k) + \langle \nabla F(z^k), x^* - z^k \rangle) \\ &\quad + \gamma_{t,k} \langle \nabla F(z^k), x^{k+1} - x^* \rangle + \frac{L\gamma_{t,k}^2}{2} \|x^{k+1} - x^k\|^2 \end{aligned} \quad (\text{D.13})$$

From the convexity of  $F$ , we have

$$F(y^k) \geq F(z^k) + \langle \nabla F(z^k), y^k - z^k \rangle \quad \text{and} \quad F(x^*) \geq F(z^k) + \langle \nabla F(z^k), x^* - z^k \rangle \quad (\text{D.14})$$

Hence, we get

$$\begin{aligned} F(y^{k+1}) &\leq (1 - \gamma_{t,k})F(y^k) + \gamma_{t,k}F^* + \gamma_{t,k} \langle \nabla F(z^k), x^{k+1} - x^k \rangle + \frac{L\gamma_{t,k}^2}{2} \|x^{k+1} - x^k\|^2 \\ &= (1 - \gamma_{t,k})F(y^k) + \gamma_{t,k}F^* + \gamma_{t,k} \langle v^k + \beta_{t,k}(x^{k+1} - x^k), x^{k+1} - x^* \rangle - \gamma_{t,k}\beta_{t,k} \langle x^{k+1} - x^k, x^{k+1} - x^* \rangle \\ &\quad + \gamma_{t,k} \langle \nabla F(x^k) - v^k, x^{k+1} - x^* \rangle + \frac{L\gamma_{t,k}^2}{2} \|x^{k+1} - x^k\|^2 \\ &\leq (1 - \gamma_{t,k})F(y^k) + \gamma_{t,k}F^* + \gamma_{t,k}\alpha_{t,k} - \gamma_{t,k}\beta_{t,k} \langle x^{k+1} - x^k, x^{k+1} - x^* \rangle \\ &\quad + \gamma_{t,k} \langle \nabla F(x^k) - v^k, x^{k+1} - x^* \rangle + \frac{L\gamma_{t,k}^2}{2} \|x^{k+1} - x^k\|^2 \\ &= (1 - \gamma_{t,k})F(y^k) + \gamma_{t,k}F^* + \gamma_{t,k}\alpha_{t,k} + \frac{\beta_{t,k}\gamma_{t,k}}{2} (\|x^k - x^*\|^2 - \|x^{k+1} - x^*\|^2) \\ &\quad + \gamma_{t,k} \left( \frac{L\gamma_{t,k} - \beta_{t,k}}{2} \|x^{k+1} - x^k\|^2 + \langle \nabla F(z^k) - v^k, x^{k+1} - x^* \rangle \right) \\ &\leq (1 - \gamma_{t,k})F(y^k) + \gamma_{t,k}F^* + \gamma_{t,k}\alpha_{t,k} + \frac{\beta_{t,k}\gamma_{t,k}}{2} (\|x^k - x^*\|^2 - \|x^{k+1} - x^*\|^2) + \gamma_{t,k}D\|\nabla F(z^k) - v^k\| \end{aligned} \quad (\text{D.15})$$

where the second inequality follows from the definition of  $\alpha_{t,k}$ , and the last inequality from the fact that  $\beta_{t,k} \geq L\gamma_{t,k}$  together with Cauchy-Schwarz inequality. Then, we subtract  $F^*$  from both sides, take the expectation of both sides, and compute the telescopic sum over  $(t, k)$ . The first term vanishes due to the  $(1 - \gamma_{1,1}) = 0$  factor. Denoting by  $\mathcal{E}_{t,k} := \mathbb{E}\|x^k - x^*\|^2$ , and  $\mathcal{E}_{t,k+} := \mathbb{E}\|x^{k+1} - x^*\|^2$ , we get

$$\mathbb{E}[F(y^{k+1}) - F^*] \leq \sum_{(\tau,i)} \left[ \gamma_{\tau,i}\alpha_{\tau,i} + \frac{\beta_{\tau,i}\gamma_{\tau,i}}{2} (\mathcal{E}_{\tau,i} - \mathcal{E}_{\tau,i+}) + \gamma_{\tau,i}D\mathbb{E}\|\nabla F(z^k) - v^k\| \right] \prod_{(\tau',j)=(\tau,i)}^{(t,k)} (1 - \gamma_{\tau',j}). \quad (\text{D.16})$$

Remark that

$$\prod_{(\tau',j)=(\tau,i)}^{(t,k)} (1 - \gamma_{\tau',j}) = \prod_{r=i}^{K_\tau} \frac{s_{\tau,r} - 1}{s_{\tau,r} + 2} \prod_{\tau'=\tau+1}^{t-1} \prod_{j=1}^{K_{\tau'}} \frac{s_{\tau',j} - 1}{s_{\tau',j} + 2} \prod_{\ell=1}^k \frac{s_{t,\ell} - 1}{s_{t,\ell} + 2} = \frac{(s_{\tau,i} - 1)s_{\tau,i}(s_{\tau,i} + 1)}{s_{t,k}(s_{t,k} + 1)(s_{t,k} + 2)} \quad (\text{D.17})$$

Now we focus on the individual terms

$$\sum_{(\tau,i)} \gamma_{\tau,i} \alpha_{\tau,i} \frac{(s_{\tau,i} - 1)s_{\tau,i}(s_{\tau,i} + 1)}{s_{t,k}(s_{t,k} + 1)(s_{t,k} + 2)} \leq \sum_{(\tau,i)} \frac{3}{s_{\tau,i} + 2} \frac{2LD^2}{s_{\tau,i}^2} \frac{(s_{\tau,i} - 1)s_{\tau,i}(s_{\tau,i} + 1)}{s_{t,k}(s_{t,k} + 1)(s_{t,k} + 2)} \quad (\text{D.18})$$

Now, show that

$$\begin{aligned} \sum_{(\tau,i)} \frac{\beta_{\tau,i} \gamma_{\tau,i}}{2} (\mathcal{E}_{\tau,i} - \mathcal{E}_{\tau,i+}) \prod_{(\tau',j)=(\tau,i)}^{(t,k)} (1 - \gamma_{\tau',j}) &= \sum_{(\tau,i)} \frac{27L}{4} \frac{1}{(s_{\tau,i} + 2)^2} \frac{(s_{\tau,i} - 1)s_{\tau,i}(s_{\tau,i} + 1)}{s_{t,k}(s_{t,k} + 1)(s_{t,k} + 2)} (\mathcal{E}_{\tau,i} - \mathcal{E}_{\tau,i+}) \\ &= \frac{27L}{4s_{t,k}(s_{t,k} + 1)(s_{t,k} + 2)} \sum_{(\tau,i)} \frac{(s_{\tau,i} - 1)s_{\tau,i}(s_{\tau,i} + 1)}{(s_{\tau,i} + 2)^2} (\mathcal{E}_{\tau,i} - \mathcal{E}_{\tau,i+}) \end{aligned} \quad (\text{D.19})$$

Remark that

$$\begin{aligned} \sum_{s_{\tau,i}=1}^{s_{t,k}} \frac{(s_{\tau,i} - 1)s_{\tau,i}(s_{\tau,i} + 1)}{(s_{\tau,i} + 2)^2} \mathcal{E}_{\tau,i} - \sum_{s_{\tau,i}=1}^{s_{t,k}} \frac{(s_{\tau,i} - 1)s_{\tau,i}(s_{\tau,i} + 1)}{(s_{\tau,i} + 2)^2} \mathcal{E}_{\tau,i+} \\ &= \sum_{s_{\tau,i}=1}^{s_{t,k}-1} \frac{s_{\tau,i}(s_{\tau,i} + 1)(s_{\tau,i} + 2)}{(s_{\tau,i} + 3)^2} \mathcal{E}_{\tau,i+} - \sum_{s_{\tau,i}=1}^{s_{t,k}} \frac{(s_{\tau,i} - 1)s_{\tau,i}(s_{\tau,i} + 1)}{(s_{\tau,i} + 2)^2} \mathcal{E}_{\tau,i+} \\ &\leq \sum_{s_{\tau,i}=1}^{s_{t,k}} \frac{s_{\tau,i}(s_{\tau,i} + 1)(s_{\tau,i} + 2)}{(s_{\tau,i} + 3)^2} \mathcal{E}_{\tau,i+} - \sum_{s_{\tau,i}=1}^{s_{t,k}} \frac{(s_{\tau,i} - 1)s_{\tau,i}(s_{\tau,i} + 1)}{(s_{\tau,i} + 2)^2} \mathcal{E}_{\tau,i+} \\ &\leq D^2 \sum_{s_{\tau,i}=1}^{s_{t,k}} \left( \frac{s_{\tau,i}(s_{\tau,i} + 1)(s_{\tau,i} + 2)}{(s_{\tau,i} + 3)^2} - \frac{(s_{\tau,i} - 1)s_{\tau,i}(s_{\tau,i} + 1)}{(s_{\tau,i} + 2)^2} \right) \\ &\leq D^2 s_{t,k} \end{aligned} \quad (\text{D.20})$$

Hence, we get

$$\sum_{(\tau,i)} \frac{\beta_{\tau,i} \gamma_{\tau,i}}{2} (\mathcal{E}_{\tau,i} - \mathcal{E}_{\tau,i+}) \prod_{(\tau',j)=(\tau,i)}^{(t,k)} (1 - \gamma_{\tau',j}) \leq \frac{27LD^2}{4(s_{t,k} + 1)(s_{t,k} + 2)} \quad (\text{D.21})$$

Finally, we focus on the last term:

$$\begin{aligned} \sum_{(\tau,i)} \gamma_{\tau,i} D \mathbb{E} \|\nabla F(z^{\tau,i}) - v^{\tau,i}\| \frac{(s_{\tau,i} - 1)s_{\tau,i}(s_{\tau,i} + 1)}{s_{t,k}(s_{t,k} + 1)(s_{t,k} + 2)} &\leq LD^2 \sum_{s_{\tau,i}=1}^{s_{t,k}} \frac{6\sqrt{2}}{(s_{\tau,i} + 2)(s_{\tau,i} + 1)^2} \frac{(s_{\tau,i} - 1)s_{\tau,i}(s_{\tau,i} + 1)}{s_{t,k}(s_{t,k} + 1)(s_{t,k} + 2)} \\ &\leq \frac{6\sqrt{2}LD^2}{(s_{t,k} + 1)(s_{t,k} + 2)} \end{aligned} \quad (\text{D.22})$$

Combining these bounds, we obtain

$$\mathbb{E}[F(y^{k+1}) - F^*] = \mathcal{O} \left( \frac{LD^2}{(s_{t,k} + 1)(s_{t,k} + 2)} \right) \quad (\text{D.23})$$

Easy to verify by induction that  $K_t \leq s_{t,k} \leq 4K_t$ . Hence,  $s_{t,k} = \Theta(K_t) = \Theta(2^{t/2})$ .

As a direct consequence,  $\mathbb{E}[F(y^{k+1}) - F^*] = \mathcal{O}(LD^2 2^{-t})$ . Therefore, to get  $\epsilon$ -solution, we set  $T_\epsilon$  as

$$T_\epsilon = \Theta \left( \log_2 \left( \frac{LD^2}{\epsilon} \right) \right) \quad (\text{D.24})$$

Then, to achieve this accuracy, we can calculate *ifo* complexity as

$$\begin{aligned} \#(ifo) &= \sum_{t=1}^{T_\epsilon} \left( Q_t + \sum_{k=2}^{K_t} S_{t,k} \right) = \mathcal{O} \left( \sum_{t=1}^{T_\epsilon} \left( n + \sum_{k=2}^{K_t} 2^{3t/2} \right) \right) = \mathcal{O} \left( nT_\epsilon + \sum_{t=1}^{T_\epsilon} 2^{2t} \right) = \mathcal{O} (nT_\epsilon + 2^{2T_\epsilon}) \\ &= \mathcal{O} \left( n \ln \left( \frac{LD^2}{\epsilon} \right) + \frac{L^2 D^4}{\epsilon^2} \right) \end{aligned} \quad (\text{D.25})$$

Finally, we can find *lmo* complexity by using Lemma 3:

$$\#(lmo) = \mathcal{O} \left( \sum_{t=1}^{T_\epsilon} \sum_{k=1}^{K_t} \frac{4\beta_{t,k} D^2}{\alpha_{t,k}} \right) = \mathcal{O} \left( \sum_{t=1}^{T_\epsilon} K_t^2 \right) = \mathcal{O} (2^{T_\epsilon}) = \mathcal{O} \left( \frac{LD^2}{\epsilon} \right) \quad (\text{D.26})$$

### Proof for Theorem 6 and Corollary 6

Similar to the proof of finite-sum setting, but we use (D.2) instead of (D.1), hence the constants at (D.22) change.

To get  $\epsilon$ -solution, we can calculate *sfo* complexity as

$$\begin{aligned} \#(sfo) &= \sum_{t=1}^{T_\epsilon} \left( Q_t + \sum_{k=2}^{K_t} S_{t,k} \right) = \mathcal{O} \left( \sum_{t=1}^{T_\epsilon} \left( \frac{\sigma^2 K_t^4}{L^2 D^2} + \sum_{k=2}^{K_t} 2^{3t/2} \right) \right) = \mathcal{O} \left( \left( \frac{\sigma^2}{L^2 D^2} + 1 \right) \sum_{t=1}^{T_\epsilon} 2^{2t} \right) \\ &= \mathcal{O} \left( \left( \frac{\sigma^2}{L^2 D^2} + 1 \right) 2^{T_\epsilon} \right) = \mathcal{O} \left( \frac{\sigma^2 D^2 + L^2 D^4}{\epsilon^2} \right) \end{aligned} \quad (\text{D.27})$$

The *lmo* complexity is same as the convex finite-sum case.

### Proof for Theorem 7 and Corollary 7

We start by (B.3), and rearrange it to obtain

$$F(y^{k+1}) \leq F(y^k) - \gamma \mathcal{G}(y^k) + \gamma \langle \nabla F(y^k) - v^k, x^{k+1} - w_\star^k \rangle + \gamma \langle v^k, x^{k+1} - w_\star^k \rangle + \gamma^2 \frac{L}{2} D^2 \quad (\text{D.28})$$

where  $w_\star^k = \arg \max_{x \in \Omega} \langle x, -\nabla F(y^k) \rangle$ .

Focus on the last inner-product term

$$\begin{aligned} \gamma \langle v^k, x^{k+1} - w_\star^k \rangle &= \gamma \langle v^k + \beta(x^{k+1} - x^k), x^{k+1} - w_\star^k \rangle - \gamma \beta \langle x^{k+1} - x^k, x^{k+1} - w_\star^k \rangle \\ &\leq \gamma \alpha - \gamma \beta \langle x^{k+1} - x^k, x^{k+1} - w_\star^k \rangle \\ &\leq \gamma \alpha + \gamma \beta D^2 \end{aligned} \quad (\text{D.29})$$

The first inequality follows from the role of  $\alpha$  in CndG, and the second from Cauchy-Schwarz.

Now we use

$$\begin{aligned} \gamma \langle \nabla F(y^k) - v^k, x^{k+1} - w_\star^k \rangle &\leq \gamma D \|\nabla F(y^k) - v^k\| \\ &\leq \gamma D \|\nabla F(y^k) - \nabla F(z^k)\| + \gamma D \|\nabla F(z^k) - v^k\| \\ &\leq \gamma LD \|y^k - z^k\| + \gamma D \|\nabla F(z^k) - v^k\| \\ &= \gamma^2 LD \|x^k - y^k\| + \gamma D \|\nabla F(z^k) - v^k\| \\ &\leq \gamma^2 LD^2 + \gamma D \|\nabla F(z^k) - v^k\| \end{aligned} \quad (\text{D.30})$$

Combining these bounds, we obtain

$$\begin{aligned} F(y^{k+1}) &\leq F(y^k) - \gamma \mathcal{G}(y^k) + \gamma D \|\nabla F(z^k) - v^k\| + \gamma \alpha + \gamma \beta D^2 + \frac{3}{2} \gamma^2 LD^2 \\ &= F(y^k) - \frac{1}{\sqrt{TK}} \mathcal{G}(y^k) + \frac{D}{\sqrt{TK}} \|\nabla F(z^k) - v^k\| + \frac{4LD^2}{TK} \end{aligned} \quad (\text{D.31})$$

Now we take expectation of both sides and use (D.3), and we obtain

$$\mathbb{E}[F(y^{k+1})] \leq \mathbb{E}[F(y^k)] - \frac{1}{\sqrt{TK}} \mathbb{E}[\mathcal{G}(y^k)] + \frac{6LD^2}{TK} \quad (\text{D.32})$$

Rearranging, and summing over all  $(t, k)$  pairs up to  $(T, K)$ , we get

$$\frac{1}{\sqrt{TK}} \sum_{(\tau, i)} \mathbb{E}[\mathcal{G}(y^{\tau, i})] \leq F(\bar{x}^1) - F(x^*) + 6LD^2 \quad (\text{D.33})$$

Hence, by definition of  $y^{\text{out}}$ , we get

$$\mathbb{E}[\mathcal{G}(y^{\text{out}})] \leq \frac{\mathcal{E} + 6LD^2}{\sqrt{TK}} \quad (\text{D.34})$$

This completes the convergence rate proof.

Proof for the *ifo* complexity follows similarly to the one for SPIDER-FW.

To show *lmo* complexity, we use

$$\#(lmo) = \mathcal{O} \left( \sum_{t=1}^{T_\epsilon} \sum_{k=1}^K \frac{4\beta D^2}{\alpha} \right) = \mathcal{O} \left( \sum_{t=1}^{T_\epsilon} \sum_{k=1}^K 6 \right) = \mathcal{O}(6KT_\epsilon) = \mathcal{O} \left( (\mathcal{E}^2 + L^2 D^4) \frac{1}{\epsilon^2} \right) \quad (\text{D.35})$$

### Proof for Theorem 8 and Corollary 8

Follows similarly as in the non-convex finite-sum case, but we use (D.4) instead of (D.3), hence we have additional  $\epsilon/2$  term on the right-hand-side of (D.34). *sfo* complexity follows similarly as in SPIDER-FW. *lmo* complexity is the same as the finite-sum case.