

# Generalization Bounds of Nonconvex-(Strongly)-Concave Stochastic Minimax Optimization

Siqi Zhang <sup>\*†</sup>    Yifan Hu <sup>\*‡</sup>    Liang Zhang <sup>§</sup>    Niao He <sup>§</sup>

February 8, 2023

## Abstract

This paper takes an initial step to systematically investigate the generalization bounds of algorithms for solving nonconvex-(strongly)-concave (NC-SC / NC-C) stochastic minimax optimization measured by the stationarity of primal functions. We first establish *algorithm-agnostic generalization bounds* via *uniform convergence* between the empirical minimax problem and the population minimax problem. The sample complexities for achieving  $\epsilon$ -generalization are  $\tilde{O}(d\kappa^2\epsilon^{-2})$  and  $\tilde{O}(d\epsilon^{-4})$  for NC-SC and NC-C settings, respectively, where  $d$  is the dimension and  $\kappa$  is the condition number. We further study the *algorithm-dependent generalization bounds* via stability arguments of algorithms. In particular, we introduce a novel stability notion for minimax problems and build a connection between generalization bounds and the stability notion. As a result, we establish *algorithm-dependent generalization bounds* for *stochastic gradient descent ascent (SGDA)* algorithm and the more general *sampling-determined algorithms*.

## 1 Introduction

In this paper, we consider stochastic minimax problems:

$$\min_{x \in \mathcal{X}} \max_{y \in \mathcal{Y}} F(x, y) \triangleq \mathbb{E}_{\xi} [f(x, y; \xi)], \quad (1)$$

where  $\mathcal{X} \subseteq \mathbb{R}^d$  and  $\mathcal{Y} \subseteq \mathbb{R}^{d'}$  ( $d, d' \in \mathbb{N}_+$ ) are two nonempty closed convex sets,  $\xi \in \Xi$  is a random variable following an unknown distribution  $\mathcal{D}$ , and  $f : \mathcal{X} \times \mathcal{Y} \times \Xi \rightarrow \mathbb{R}$  is continuously differentiable and Lipschitz smooth jointly in  $x$  and  $y$  for any  $\xi$ . We denote the objective (1) as the *population minimax problem*. Throughout the paper, we focus on the case where  $F$  is nonconvex in  $x$  and (strongly)-concave in  $y$ , i.e., *nonconvex-(strongly)-concave (NC-SC / NC-C)*. Such minimax problems appear ubiquitously in practical applications, including adversarial training [Madry et al., 2018, Wang et al., 2019], generative adversarial networks (GANs) [Goodfellow et al., 2014, Sanjabi et al., 2018, Lei et al., 2020], reinforcement learning [Dai et al., 2017, 2018, Huang and Jiang, 2020] and robust training [Sinha et al., 2018].

Although the distribution  $\mathcal{D}$  often remains unknown, one generally has access to a dataset  $S = \{\xi_1, \dots, \xi_n\}$  consisting of  $n$  independently and identical distributed (i.i.d.) samples from  $\mathcal{D}$ . Correspondingly, researchers resort to solving an *empirical minimax problem*:

$$\min_{x \in \mathcal{X}} \max_{y \in \mathcal{Y}} F_S(x, y) \triangleq \frac{1}{n} \sum_{i=1}^n f(x, y; \xi_i). \quad (2)$$

A natural question arises:

---

\*Equal contribution.

†Department of Applied Mathematics and Statistics, Johns Hopkins University, USA. [szhan207@jhu.edu](mailto:szhan207@jhu.edu)

‡Risk Analytics and Optimization Chair, EPFL, Switzerland. [yifan.hu@epfl.ch](mailto:yifan.hu@epfl.ch)

§Department of Computer Science, ETH Zürich, Switzerland. [liang.zhang@inf.ethz.ch](mailto:liang.zhang@inf.ethz.ch), [niao.he@inf.ethz.ch](mailto:niao.he@inf.ethz.ch)

*How does the output of an algorithm  $\mathcal{A}$  for solving the empirical minimax problem generalize on the population minimax problem?*

We first specify the measurement. Since functions  $F$  and  $F_S$  are nonconvex in  $x$ , finding their global optimal solutions is generally intractable. Instead, one aims to design an algorithm  $\mathcal{A}$  that finds an  $\epsilon$ -stationary point Lin et al. [2020a], i.e.,

$$\|\nabla\Phi(\mathcal{A}_x(S))\| \leq \epsilon \quad \text{or} \quad \mathbf{dist}(0, \partial\Phi(\mathcal{A}_x(S))) \leq \epsilon,$$

where  $\Phi(x) \triangleq \max_{y \in \mathcal{Y}} F(x, y)$  and  $\Phi_S(x) \triangleq \max_{y \in \mathcal{Y}} F_S(x, y)$  are primal functions,  $\mathcal{A}_x(S)$  is the  $x$ -component of the output of any algorithm  $\mathcal{A}$  for solving (2),  $\mathbf{dist}(y, X) \triangleq \inf_{x \in X} \|y - x\|$  and  $\partial\Phi$  is the (Fréchet) subdifferential of  $\Phi$ . When  $\Phi$  is nonsmooth, we resort to the gradient norm of its Moreau envelope to measure the first-order stationarity as it provides an upper bound on  $\mathbf{dist}(0, \partial\Phi(\cdot))$  [Davis and Drusvyatskiy, 2019].

Taking the gradient norm as an example, the error for solving the population minimax problem (1) via solving its empirical counterpart (2) consists of two terms<sup>1</sup>:

$$\mathbb{E} \|\nabla\Phi(\mathcal{A}_x(S))\| \leq \underbrace{\mathbb{E} \|\nabla\Phi_S(\mathcal{A}_x(S))\|}_{\text{optimization error}} + \underbrace{\mathbb{E} \|\nabla\Phi(\mathcal{A}_x(S)) - \nabla\Phi_S(\mathcal{A}_x(S))\|}_{\text{generalization error}}. \quad (3)$$

Such decomposition on the gradient norm also appears in nonconvex minimization, e.g., Foster et al. [2018], Mei et al. [2018], Davis and Drusvyatskiy [2022], Lei [2022]. The optimization error corresponds to the error of solving the empirical minimax problem (2) and has been widely studied [Luo et al., 2020, Yang et al., 2020b]. On the other hand, the generalization error for minimax problems remains largely unexplored. A few recent works [Farnia and Ozdaglar, 2021, Zhang et al., 2021a, Lei et al., 2021, Ozdaglar et al., 2022] studied the generalization performances in minimax optimization measured by the function value-based gap. However, these do not fit well in the nonconvex setting since the optimization part remains intractable.

The goal of our paper is to characterize the generalization error

$$\mathbb{E} \|\nabla\Phi(\mathcal{A}_x(S)) - \nabla\Phi_S(\mathcal{A}_x(S))\|.$$

It is not easy as both  $\Phi_S(\cdot)$  and  $\mathcal{A}_x(S)$  depend on the dataset  $S$ , which induces correlation issues when taking expectation. To address such dependence issue, one may use *uniform convergence* or *stability arguments*.

Uniform convergence characterizes the difference between the empirical minimax optimization and the population minimax problem on worst  $x \in \mathcal{X}$ , i.e.,

$$\mathbb{E} \sup_{x \in \mathcal{X}} \|\nabla\Phi(x) - \nabla\Phi_S(x)\|.$$

Although uniform convergence has been extensively studied in stochastic optimization [Kleywegt et al., 2002, Mei et al., 2018, Davis and Drusvyatskiy, 2022], a key difference for stochastic minimax optimization is that the primal function is the average over  $n$  i.i.d. random functions. Thus existing techniques in uniform convergence for classical stochastic optimization are not directly applicable here. One needs to additionally characterize the differences between maximizers

$$\operatorname{argmax}_{y \in \mathcal{Y}} F_S(x, y) \quad \text{and} \quad \operatorname{argmax}_{y \in \mathcal{Y}} F(x, y).$$

Note that uniform convergence is invariant to the choice of algorithms and provides an upper bound on the generalization error for any  $\mathcal{A}_x(S) \in \mathcal{X}$ . Thus the derived generalization bound is *algorithm-agnostic* that applies to any algorithms. Since it is the worst case over all  $x \in \mathcal{X}$ , the derived bounds generally involve the dimension of  $x$ .

We further investigate generalization bounds via stability arguments. This approach analyzes the stability of specific algorithms and builds a connection between stability and generalization. It has been extensively

---

<sup>1</sup> For the simplicity of demonstration, here we assume there is no constraint, and primal functions are differentiable. We will formally introduce the detailed settings in Section 2.

| Setting <sup>1</sup> | Approach  | Stability Argument   |   |
|----------------------|---|--|---|
|                      | Uniform Convergence   | SGDA   | Sampling-determined Alg.  |
| NC-SC                | $\tilde{\mathcal{O}}\left(\kappa\sqrt{\frac{d}{n}}\right)$<br>Theorem 3.1       | $\mathcal{O}\left(\kappa^{1+\zeta_1}\left(\frac{T^{1-\zeta_1}}{n} + \frac{1}{\sqrt{n}}\right)\right)$<br>Corollary 4.1 | $\mathcal{O}\left(\kappa\left(\sqrt{\frac{T}{n}} + \frac{1}{\sqrt{n}}\right)\right)$<br>Corollary 4.3       |
| NC-C                 | $\tilde{\mathcal{O}}\left(\left(\frac{d}{n}\right)^{1/4}\right)$<br>Theorem 3.2 | $\mathcal{O}\left(\left(\frac{T^{1-\zeta_2}}{n}\right)^{1/6} + \left(\frac{1}{n}\right)^{1/8}\right)$<br>Corollary 4.2 | $\mathcal{O}\left(\left(\frac{T}{n}\right)^{1/12} + \left(\frac{1}{n}\right)^{1/8}\right)$<br>Corollary 4.4 |

<sup>1</sup>  $\tilde{\mathcal{O}}(\cdot)$  hides logarithmic factors,  $d$ : the dimension of  $\mathcal{X}$ ,  $n$ : sample size,  $\kappa$ : condition number  $\frac{L}{\mu}$   
 $L$ : Lipschitz smoothness parameter,  $\mu$ : strong concavity parameter,  $T$ : iteration number of algorithms  
 $\zeta_1, \zeta_2 \in (0, 1)$ : constants depending on stepsizes, refer to Corollary 4.1 and 4.2 for details. SGDA has specific requirements on stepsize, while sampling-determined algorithms do not have restrictions on stepsize.

Table 1: Summary of Generalization Bounds for Nonconvex Stochastic Minimax Optimization

studied in stochastic minimization [Bousquet and Elisseeff, 2002, Shalev-Shwartz et al., 2010, Hardt et al., 2016, Klochkov and Zhivotovskiy, 2021] and minimax problems recently [Farnia and Ozdaglar, 2021, Lei et al., 2021, Boob and Guzmán, 2021, Yang et al., 2022c]. Most of these work focuses on the measurement of the function-value gap. Under such measurement, generalization follows directly from stability. However, for the measurement of stationarity for nonconvex problems, building up a link between stability and generalization becomes significantly more challenging. Compared to uniform convergence, the stability-based generalization bound is generally independent of the dimension  $d$ . As it requires a case-by-case analysis of stability for different algorithms, it is *algorithm-dependent*. We particularly study the generalization of the widely used *stochastic gradient descent ascent (SGDA)* [Farnia and Ozdaglar, 2021] and a class of algorithms that extends SGDA called *sampling-determined algorithms* (see Definition 4.3) [Lei, 2022].

## 1.1 Contributions

In this paper, we initiate a systematic study on the generalization bounds (see Table 1) for nonconvex stochastic minimax problems from both *uniform convergence* and *stability argument* perspectives. Our contributions are two-fold:

- We establish the first uniform convergence results between the population and the empirical nonconvex minimax optimization in NC-SC and NC-C settings, measured by stationarity. Our results provide an *algorithm-agnostic* generalization bound for any algorithms that solve empirical nonconvex minimax problems. Specifically, the sample complexities to achieve an  $\epsilon$ -uniform convergence or an  $\epsilon$ -generalization error are  $\tilde{\mathcal{O}}(d\kappa^2\epsilon^{-2})$  and  $\tilde{\mathcal{O}}(d\epsilon^{-4})$  for the NC-SC and NC-C settings, respectively.
- From the stability argument perspective, we first introduce a novel stability measurement based on stationarity; then, we establish the connection between the stability and the generalization error of an algorithm in both NC-SC and NC-C settings. We further provide the *algorithm-dependent* generalization error measured by the stationarity for the classical SGDA algorithm and sampling-determined algorithms utilizing their stability.

## 1.2 Literature Review

**Nonconvex Minimax Optimization** Various algorithms have been proposed to solve NC-SC minimax optimization [Nouiehed et al., 2019, Lin et al., 2020a,b, Luo et al., 2020, Yang et al., 2020a, Boç and Böhm, 2020, Xu et al., 2020, Lu et al., 2020, Yan et al., 2020, Guo et al., 2021, Sharma et al., 2022, Zhang et al., 2022]. For stochastic NC-SC minimax problems, Zhang et al. [2022] achieves the best-known complexity of  $\mathcal{O}(\kappa\epsilon^{-4})$ ,

and  $\mathcal{O}(\kappa^2\epsilon^{-3})$  result with additional individual smoothness assumption. Also, Yang et al. [2022b] introduced a stochastic smoothed-AGDA algorithm which achieves the best-known  $\mathcal{O}(\kappa^2\epsilon^{-4})$  for single-loop algorithms. The lower bounds of NC-SC problems are extensively studied in several recent works [Zhang et al., 2021b, Han et al., 2021, Li et al., 2021].

The primal function for NC-SC problems is generally smooth, while the primal functions can be nonsmooth for NC-C problems [Thekumparampil et al., 2019, Lin et al., 2020a]. Recent years witnessed a surge of algorithms for NC-C problems in deterministic, finite-sum, and stochastic settings, e.g., [Zhang et al., 2020, Ostrovskii et al., 2021, Thekumparampil et al., 2019, Zhao, 2020, Nouiehed et al., 2019, Yang et al., 2020b, Lin et al., 2020a, Boş and Böhm, 2020, Rafique et al., 2021], to name a few. To the best of our knowledge, Thekumparampil et al. [2019], Yang et al. [2020b], Lin et al. [2020b] achieved the best known  $\tilde{\mathcal{O}}(\epsilon^{-3})$  complexity in the deterministic case, while Yang et al. [2020b] achieved the best known  $\tilde{\mathcal{O}}(n^{3/4}\epsilon^{-3})$  complexity in the finite-sum case, and Zhang et al. [2022] provided the best known  $\mathcal{O}(\epsilon^{-6})$  complexity in the stochastic case. These works differ from our paper in that we aim to characterize the generalization error of algorithms while they focus mainly on the optimization error of the algorithms.

**Uniform Convergence** A series of works from stochastic optimization and statistical learning theory studied uniform convergence on the worst-case differences between the population objective  $L(x)$  and its empirical objective  $L_S(x)$  constructed via sample average approximation (SAA, also known as empirical risk minimization). Interested readers may refer to prominent results in statistical learning [Fisher, 1922, Vapnik, 1999, Van der Vaart, 2000]. For finite-dimensional problem, Kleywegt et al. [2002] showed that the sample complexity is  $\mathcal{O}(d\epsilon^{-2})$  to achieve an  $\epsilon$ -uniform convergence in high probability, i.e.,  $\mathbb{P}(\sup_{x \in \mathcal{X}} |L(x) - L_S(x)| \geq \epsilon)$ . For nonconvex empirical objectives, Mei et al. [2018] and Davis and Drusvyatskiy [2022] established  $\tilde{\mathcal{O}}(d\epsilon^{-2})$  sample complexity of uniform convergence measured by the stationarity for nonconvex smooth and weakly convex functions, respectively. In addition, Wang et al. [2017] used uniform convergence to demonstrate the generalization and the gradient complexity of differentially private algorithms for stochastic optimization. Recently, Amir et al. [2022] demonstrated the generalization error of gradient descent on a generalized linear model using uniform convergence and showed that the stability argument is insufficient to achieve generalization. To the best of our knowledge, our paper is the first to study uniform convergence for nonconvex minimax optimization.

**Stability-Based Generalization Bounds** This line of research focuses on analyzing generalization bounds of stochastic optimization via the uniform stability property of specific algorithms, including SAA [Bousquet and Elisseeff, 2002, Shalev-Shwartz et al., 2009], stochastic gradient descent [Hardt et al., 2016, Bassily et al., 2020, Lei, 2022], and uniformly stable algorithms [Klochkov and Zhivotovskiy, 2021]. Recently, a series of works further studied the generalization performances measured by the function-value gap of various algorithms in minimax problems. Farnia and Ozdaglar [2021] gave the generalization bound for the outputs of gradient-descent-ascent (GDA) and proximal-point algorithm (PPA) in both (strongly)-convex-(strongly)-concave and nonconvex-nonconcave smooth minimax problems. Lei et al. [2021] focused on GDA and provided a comprehensive study for different settings of minimax problems with generalization measured by function-value gaps. Boob and Guzmán [2021] provided stability and generalization results of extragradient algorithm (EG) in the smooth convex-concave setting. On the other hand, Zhang et al. [2021a] studied stability and generalization of the empirical minimax problem under the (strongly)-convex-(strongly)-concave setting, assuming that one can find the optimal solution of the empirical minimax problem. Our work differs from those in that we propose a novel stability notion for minimax optimization measured by stationarity and build up the link between such stability and generalization error.

## 2 Problem Setting

**Notations** Throughout the paper, we use  $\|\cdot\|$  as the  $\ell_2$ -norm,  $\nabla f = (\nabla_x f, \nabla_y f)$  as the gradient of a function  $f$ , for nonnegative functions  $f$  and  $g$ , we say  $f = \mathcal{O}(g)$  if  $f(x) \leq cg(x)$  for some  $c > 0$ . We denote  $\mathbf{proj}_{\mathcal{X}}(x') \triangleq \operatorname{argmin}_{x \in \mathcal{X}} \|x - x'\|^2$  as the projection operator. Let  $\mathcal{A}(S) \triangleq (\mathcal{A}_x(S), \mathcal{A}_y(S))$  denote the output of an algorithm

$\mathcal{A}$  on the empirical minimax problem (2) with dataset  $S$ . Given  $\mu \geq 0$ , we say a function  $g : \mathcal{X} \rightarrow \mathbb{R}$  is  $\mu$ -strongly convex if  $g(x) - \frac{\mu}{2}\|x\|^2$  is convex, and it is  $\mu$ -strongly concave if  $-g$  is  $\mu$ -strongly convex. Function  $g$  is  $\mu$ -weakly convex if  $g(x) + \frac{\mu}{2}\|x\|^2$  is convex (see more notations and standard definitions in Appendix A).

**Definition 2.1 (Smooth Function)** *We say a function  $f : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$  is  $L$ -smooth jointly in  $(x, y)$  if the function is continuously differentiable, and there exists a constant  $L > 0$  such that for any  $(x_1, y_1), (x_2, y_2) \in \mathcal{X} \times \mathcal{Y}$ , we have*

$$\begin{aligned}\|\nabla_x f(x_1, y_1) - \nabla_x f(x_2, y_2)\| &\leq L(\|x_1 - x_2\| + \|y_1 - y_2\|), \\ \|\nabla_y f(x_1, y_1) - \nabla_y f(x_2, y_2)\| &\leq L(\|x_1 - x_2\| + \|y_1 - y_2\|).\end{aligned}$$

By definition, it is easy to find that an  $L$ -smooth function is also  $L$ -weakly convex. Next, we introduce the main assumptions used throughout the paper.

**Assumption 2.1 (Main Settings)** *We assume the following:*

- The function  $f(x, y; \xi)$  is  $L$ -smooth jointly in  $(x, y) \in \mathcal{X} \times \mathcal{Y}$  for any  $\xi$ .
- The function  $f(x, y; \xi)$  is  $\mu$ -strongly concave in  $y \in \mathcal{Y}$  for any  $x \in \mathcal{X}$  and any  $\xi$  where  $\mu \geq 0$ .
- The gradient norms of  $f(\cdot, \cdot; \xi)$  are bounded by  $G$  respectively for any  $\xi$ .
- The domains  $\mathcal{X}$  and  $\mathcal{Y}$  are compact convex sets, i.e., there exists constants  $D_{\mathcal{X}}, D_{\mathcal{Y}} > 0$  such that for any  $x \in \mathcal{X}$ ,  $\|x\|^2 \leq D_{\mathcal{X}}$  and for any  $y \in \mathcal{Y}$ ,  $\|y\|^2 \leq D_{\mathcal{Y}}$ , respectively.

Note that compact domain assumption appears widely in uniform convergence literature [Kleywegt et al., 2002, Davis and Drusvyatskiy, 2022]. Under Assumption 2.1, the objective function  $F$  is  $L$ -smooth in  $(x, y)$  and  $\mu$ -strongly concave for any  $\xi$ . When  $\mu > 0$ , we call the population minimax problem (1) a *nonconvex-strongly-concave* (NC-SC) minimax problem; when  $\mu = 0$ , we call it a *nonconvex-concave* (NC-C) minimax problem.

**Definition 2.2 (Moreau Envelope)** *For an  $L$ -weakly convex function  $\Phi$  and  $0 < \lambda < 1/L$ , we use  $\Phi^\lambda(x)$  and  $\text{prox}_{\lambda\Phi}(x)$  to denote the Moreau envelope of  $\Phi$  and the proximal point of  $\Phi$  for a given point  $x$ , defined as following:*

$$\Phi^\lambda(x) \triangleq \min_{z \in \mathcal{X}} \left\{ \Phi(z) + \frac{1}{2\lambda} \|z - x\|^2 \right\}, \quad \text{prox}_{\lambda\Phi}(x) \triangleq \underset{z \in \mathcal{X}}{\text{argmin}} \left\{ \Phi(z) + \frac{1}{2\lambda} \|z - x\|^2 \right\}. \quad (4)$$

Below we recall some important properties on the primal function  $\Phi$  and its Moreau envelope  $\Phi^\lambda(x)$  presented in the literature [Davis and Drusvyatskiy, 2019, Thekumparampil et al., 2019, Lin et al., 2020a,b].

**Lemma 2.1 (Properties of  $\Phi$  and  $\Phi^\lambda$ )** *In the NC-SC setting ( $\mu > 0$ ), both  $\Phi(x)$  and  $\Phi_S(x)$  are  $\tilde{L} \triangleq L(1+\kappa)$ -smooth with the condition number  $\kappa \triangleq L/\mu$ . In the NC-C setting ( $\mu = 0$ ), the primal function  $\Phi$  is  $L$ -weakly convex, its Moreau envelope  $\Phi^\lambda(x)$  is Lipschitz smooth, also  $\nabla\Phi^\lambda(x) = \lambda^{-1}(x - \hat{x})$ ,  $\|\nabla\Phi^\lambda(x)\| \geq \text{dist}(0, \partial\Phi(\hat{x}))$ , where  $\hat{x} = \text{prox}_{\lambda\Phi}(x)$  and  $0 < \lambda < 1/L$ .*

**Performance Measurement** In the NC-SC setting, the primal functions  $\Phi$  and  $\Phi_S$  are both  $\tilde{L}$ -smooth. Regarding the constraint, we measure the difference between the population and empirical minimax problems using the *generalized gradient of the population and the empirical primal functions*, i.e.,  $\mathbb{E} \|\mathcal{G}_\Phi(\mathcal{A}_x(S)) - \mathcal{G}_{\Phi_S}(\mathcal{A}_x(S))\|$ , where  $\mathcal{G}_\Phi(x) \triangleq \tilde{L}(x - \text{proj}_{\mathcal{X}}(x - (1/\tilde{L})\nabla\Phi(x)))$ . The following inequality summarized the relationship of measurements in terms of generalized gradient and in terms of gradient used in Section 1.

$$\underbrace{\mathbb{E} \|\mathcal{G}_\Phi(\mathcal{A}_x(S)) - \mathcal{G}_{\Phi_S}(\mathcal{A}_x(S))\|}_{\text{generalization error of Algorithm } \mathcal{A}} \leq \mathbb{E} \|\nabla\Phi(\mathcal{A}_x(S)) - \nabla\Phi_S(\mathcal{A}_x(S))\| \leq \underbrace{\mathbb{E} \left[ \max_{x \in \mathcal{X}} \|\nabla\Phi(x) - \nabla\Phi_S(x)\| \right]}_{\text{algorithm-agnostic uniform convergence}}, \quad (5)$$

where the first inequality holds as projection is a non-expansive operator. The term in the left-hand side (LHS) above is the generalization error of an algorithm  $\mathcal{A}$  we desire in the NC-SC case.

For the NC-C case, the primal function  $\Phi(x)$  is  $L$ -weakly convex, we use the gradient of its Moreau Envelope to characterize the (near)-stationarity [Davis and Drusvyatskiy, 2019]. We measure the difference between the population and empirical problems using the difference between *the gradients of their respective Moreau envelopes*.

The generalization error and the uniform convergence in the NC-C case should be given as follows:

$$\underbrace{\mathbb{E} \left\| \nabla \Phi^{1/(2L)}(\mathcal{A}_x(S)) - \nabla \Phi_S^{1/(2L)}(\mathcal{A}_x(S)) \right\|}_{\text{generalization error of Algorithm } \mathcal{A}} \leq \underbrace{\mathbb{E} \left[ \max_{x \in \mathcal{X}} \left\| \nabla \Phi^{1/(2L)}(x) - \nabla \Phi_S^{1/(2L)}(x) \right\| \right]}_{\text{algorithm-agnostic uniform convergence}}. \quad (6)$$

The term in the LHS above is the generalization error of an algorithm  $\mathcal{A}$  we desire in the NC-C case.

### 3 Uniform Convergence and Generalization

In this section, we discuss the sample complexity for achieving  $\epsilon$ -uniform convergence and  $\epsilon$ -generalization error for NC-SC and NC-C stochastic minimax optimization.

#### 3.1 NC-SC Stochastic Minimax Optimization

Under the NC-SC setting, the next theorem demonstrates the uniform convergence between gradients of primal functions of the population and empirical minimax problem, which is a structural property of the empirical and population minimax problem. We defer the proof to Appendix B.

**Theorem 3.1 (Uniform Convergence, NC-SC)** *Under Assumption 2.1 with  $\mu > 0$ , we have*

$$\mathbb{E} \left[ \max_{x \in \mathcal{X}} \left\| \nabla \Phi(x) - \nabla \Phi_S(x) \right\| \right] = \tilde{\mathcal{O}} \left( d^{1/2} \kappa n^{-1/2} \right). \quad (7)$$

Furthermore, to achieve  $\epsilon$ -uniform convergence and  $\epsilon$ -generalization error for any algorithm  $\mathcal{A}$  such that  $\mathbb{E} \left\| \mathcal{G}_\Phi(\mathcal{A}_x(S)) - \mathcal{G}_{\Phi_S}(\mathcal{A}_x) \right\| \leq \epsilon$ , it suffices to have

$$n = n_{\text{NCSC}}^* \triangleq \tilde{\mathcal{O}} \left( d \kappa^2 \epsilon^{-2} \right). \quad (8)$$

To the best of our knowledge, it is the first uniform convergence and algorithm-agnostic generalization error bound result for NC-SC stochastic minimax problem. In comparison, existing works on generalization error analysis of minimax problems [Farnia and Ozdaglar, 2021, Lei et al., 2021] using stability arguments are algorithm-specific and can only handle function-value gap measurement. Zhang et al. [2021a] establish algorithm-agnostic stability and generalization in the strongly-convex-strongly-concave regime, yet their analysis does not extend to the nonconvex regime.

Error decomposition (3) and Theorem 3.1 imply that for any algorithm that achieves an  $\epsilon$ -stationarity point of the empirical minimax problem, its sample complexity for finding an  $\epsilon$ -stationary point of the population minimax problem is at most  $\tilde{\mathcal{O}} \left( d \kappa^2 \epsilon^{-2} \right)$ . Such an observation provides generalization guarantees for any algorithms that solve finite-sum (empirical) minimax problems. It is especially useful for some SOTA algorithms like Catalyst-SVRG [Zhang et al., 2021b] and finite-sum version SREDA [Luo et al., 2020] that are complicated and thus there lack stability analysis and generalization bounds for them.

**Proof Sketch** We briefly discuss the proof of Theorem 3.1.

**Step 1:** First, we use a  $v$ -net  $\{x_k\}_{k=1}^Q$  [Vapnik, 1999] to decompose the error and handle the dependence issue between  $\operatorname{argmax}_{x \in \mathcal{X}} \left\| \nabla \Phi_S(x) - \nabla \Phi(x) \right\|$  and  $\Phi_S(x)$ .

**Step 2:** For any  $x_k$  within the  $v$ -net, we have the error following decomposition

$$\left\| \nabla \Phi_S(x_k) - \nabla \Phi(x_k) \right\| \leq \left( \left\| \nabla \Phi_S(x_k) - \nabla \Phi(x_k) \right\| - \mathbb{E} \left\| \nabla \Phi_S(x) - \nabla \Phi(x_k) \right\| \right) + \mathbb{E} \left\| \nabla \Phi_S(x_k) - \nabla \Phi(x_k) \right\|.$$

When bounding  $\mathbb{E} \|\nabla\Phi_S(x_k) - \nabla\Phi(x_k)\|$  in the right-hand side (RHS), we need to characterize the difference between  $\operatorname{argmax}_{y \in \mathcal{Y}} F_S(x, y)$  and  $\operatorname{argmax}_{y \in \mathcal{Y}} F(x, y)$  using the stability argument of sample average approximation [Shalev-Shwartz et al., 2009]. This step appears uniquely for minimax optimization due to the special structure of the primal function  $\Phi_S(x) = \max_y \frac{1}{n} \sum_{i=1}^n f(x, y; \xi_i)$ , which is not the average over  $n$  random functions. Then we utilize the established stability argument to show that the first term in the RHS is sub-Gaussian and apply the concentration inequality, which leads to the result.  $\blacksquare$

### 3.2 NC-C Stochastic Minimax Optimization

In this subsection, we derive the uniform convergence and algorithm-agnostic generalization bounds for NC-C stochastic minimax problems. Recall that the primal function  $\Phi$  is  $L$ -weakly convex [Thekumparampil et al., 2019], and thus that  $\nabla\Phi$  is not well-defined. We use the gradient norm of the Moreau envelope of the primal function as the measurement [Davis and Drusvyatskiy, 2019].

**Theorem 3.2 (Uniform Convergence, NC-C)** *Under Assumption 2.1 with  $\mu = 0$ , we have*

$$\mathbb{E} \left[ \max_{x \in \mathcal{X}} \left\| \nabla\Phi_S^{1/(2L)}(x) - \nabla\Phi^{1/(2L)}(x) \right\| \right] = \tilde{\mathcal{O}}\left(d^{1/4}n^{-1/4}\right). \quad (9)$$

Furthermore, to achieve  $\epsilon$ -uniform convergence and  $\epsilon$ -generalization error for any algorithm  $\mathcal{A}$  such that  $\mathbb{E} \left[ \left\| \nabla\Phi^{1/(2L)}(\mathcal{A}_x(S)) - \nabla\Phi^{1/(2L)}(x) \right\| \right] \leq \epsilon$ , it suffices to have

$$n = n_{\text{NCC}}^* \triangleq \tilde{\mathcal{O}}(d\epsilon^{-4}). \quad (10)$$

We defer the detailed proof to Appendix C. To the best of our knowledge, this is the first algorithm-agnostic generalization error result in NC-C stochastic minimax optimization. Similar to the NC-SC setting, Theorem 3.2 with error decomposition (3) provides generalization guarantees for any algorithms that solve the NC-C empirical minimax problem, including the best-known Catalyst algorithm [Yang et al., 2020b]. More specifically, if an algorithm finds an  $\epsilon$ -stationarity point of the empirical minimax problem, with sample size  $n = \tilde{\mathcal{O}}(d\epsilon^{-4})$ , the point is also an  $\mathcal{O}(\epsilon)$ -stationarity point of the population minimax problem.

**Proof Sketch** The analysis of Theorem 3.2 builds up a link between NC-C and NC-SC settings and consists of three parts.

**Step 1:** By definition of the gradient of the Moreau envelope, it holds that

$$\|\nabla\Phi_S^\lambda(x) - \nabla\Phi^\lambda(x)\| \leq \frac{1}{\lambda} \|\mathbf{prox}_{\lambda\Phi}(x) - \mathbf{prox}_{\lambda\Phi_S}(x)\|.$$

We first use a  $v$ -net  $\{x_k\}_{k=1}^Q$  [Vapnik, 1999] to handle the dependence issue between  $\tilde{x}^* \in \operatorname{argmax}_{x \in \mathcal{X}} \|\mathbf{prox}_{\lambda\Phi}(x) - \mathbf{prox}_{\lambda\Phi_S}(x)\|$  and  $\Phi_S$ .

**Step 2:** We introduce the following  $\ell_2$ -regularized minimax problem:

$$\min_{x \in \mathcal{X}} \max_{y \in \mathcal{Y}} F(x, y) - \frac{\nu}{2} \|y\|^2.$$

Notice that this problem is NC-SC. We further build a connection between NC-C stochastic minimax optimization problems and the corresponding regularized NC-SC stochastic minimax optimization problems. Then we carefully choose the regularization parameter  $\nu$  to derive the uniform convergence.

The following lemma characterizes the distance between the proximal points of the primal function from the original NC-C problem and its regularized NC-SC problem. Note that the lemma may be of independent interest for the design and the analysis of gradient-based methods for NC-C problems.

**Lemma 3.1** For  $\nu > 0$ , denote  $\hat{\Phi}(x) = \max_{y \in \mathcal{Y}} F(x, y) - \frac{\nu}{2} \|y\|^2$  as the primal function of the regularized NC-C problem. It holds for  $\lambda \in (0, (L + \nu)^{-1})$  that

$$\|\mathbf{prox}_{\lambda\hat{\Phi}}(x) - \mathbf{prox}_{\lambda\hat{\Phi}_S}(x)\|^2 \leq \frac{\nu D_{\mathcal{Y}} \lambda}{1 - \lambda(L + \nu)}.$$

The above lemma implies that for small regularization parameter  $\nu$ , the difference between the proximal point of the primal function  $\Phi$  of the NC-C problem and the primal function  $\hat{\Phi}$  of the regularized NC-SC problem is small.

**Step 3:** It remains to characterize the distance between  $\mathbf{prox}_{\lambda\hat{\Phi}}(x)$  and  $\mathbf{prox}_{\lambda\hat{\Phi}_S}(x)$ , where  $\hat{\Phi}_S$  is the primal function of the regularized empirical minimax problem. By definition of  $\mathbf{prox}_{\lambda\hat{\Phi}}(x)$  and  $\mathbf{prox}_{\lambda\hat{\Phi}_S}(x)$ , the distance is equivalent to the difference between the optimal solutions on  $x$  of a strongly-convex strongly-concave (SC-SC) population minimax problem and the counterpart empirical minimax problem. We utilize the existing stability-based results for SC-SC minimax optimization [Zhang et al. \[2021a\]](#) to the upper bound such a distance. We further show that  $\|\mathbf{prox}_{\lambda\hat{\Phi}_S}(x) - \mathbf{prox}_{\lambda\hat{\Phi}}(x)\| - \mathbb{E} \|\mathbf{prox}_{\lambda\hat{\Phi}_S}(x) - \mathbf{prox}_{\lambda\hat{\Phi}}(x)\|$  is a sub-Gaussian random variable and apply concentration inequality.  $\blacksquare$

### 3.3 Comparing Uniform Convergence for Minimization, NC-SC, and NC-C Minimax Problems

For general stochastic nonconvex optimization  $\min_{x \in \mathcal{X}} \mathbb{E}[f(x; \xi)]$ , the sample complexity of achieving  $\epsilon$ -uniform convergence,

$$\mathbb{E} \max_{x \in \mathcal{X}} \left\| \frac{1}{n} \sum_{i=1}^n \nabla f(x; \xi_i) - \mathbb{E} \nabla f(x; \xi) \right\| \leq \epsilon,$$

between the gradient of the population problem and the empirical problem is  $\tilde{\mathcal{O}}(d\epsilon^{-2})$  [[Davis and Drusvyatskiy, 2022](#), [Mei et al., 2018](#)]. For nonconvex minimax optimization, if we care about the uniform convergence in terms of the gradient of  $F$ , i.e.,

$$\mathbb{E} \max_{x \in \mathcal{X}, y \in \mathcal{Y}} \left\| \frac{1}{n} \sum_{i=1}^n \nabla f(x, y; \xi_i) - \mathbb{E} \nabla f(x, y; \xi) \right\|,$$

where  $\nabla f$  denotes the full gradient with respect to  $x$  and  $y$ , existing analysis in [Mei et al. \[2018\]](#) directly gives a  $\tilde{\mathcal{O}}(d\epsilon^{-2})$  sample complexity. However, since we care about the gradient of the primal function, the analysis becomes more complicated, which we detail in the following.

1. In the NC-SC setting, to establish uniform convergence, we bound

$$\mathbb{E} \max_{x \in \mathcal{X}} \|\nabla \Phi_S(x) - \nabla \Phi(x)\| = \mathbb{E} \max_{x \in \mathcal{X}} \left\| \frac{1}{n} \sum_{i=1}^n \nabla_x f(x, y_S^*(x); \xi_i) - \mathbb{E} \nabla_x f(x, y^*(x); \xi_i) \right\|$$

where

$$y_S^*(x) \triangleq \operatorname{argmax}_{y \in \mathcal{Y}} F_S(x, y), \quad y^*(x) \triangleq \operatorname{argmax}_{y \in \mathcal{Y}} F(x, y) \quad (11)$$

The primal function  $\Phi_S$  is not in the form of averaging over  $n$  samples, and thus existing analysis for the minimization problem is not directly applicable. In addition, as the optimal point  $y_S^*(x)$  differs from  $y^*(x)$ , such difference brings in an additional error term. In the NC-SC case, the error is upper bounded by  $\mathcal{O}(n^{-1/2})$ , which is the same scale as the error from establishing uniform convergence on  $x$ . Thus, the final uniform convergence bound established in [Theorem 3.1](#) is of the same order as that for minimization problem [[Mei et al., 2018](#), [Davis and Drusvyatskiy, 2022](#)] except for an additional dependence on the condition number  $\kappa$ .

2. In the NC-C case, since there may exist multiple maximizers, we have

$$y^* \in \mathcal{Y}^* = \operatorname{argmax}_{y \in \mathcal{Y}} \mathbb{E}f(x, y; \xi), \quad y_S^* \in \mathcal{Y}_S^* = \operatorname{argmax}_{y \in \mathcal{Y}} \frac{1}{n} \sum_{i=1}^n f(x, y; \xi_i). \quad (12)$$

Thus, the distance between  $y^*$  and  $y_S^*$  may not be well-defined. Instead, we bound the distance between  $\hat{y}_S^*(x) \triangleq \operatorname{argmax}_{y \in \mathcal{Y}} F_S(x, y) - \frac{\nu}{2} \|y\|^2$  and  $\hat{y}^*(x) \triangleq \operatorname{argmax}_{y \in \mathcal{Y}} F(x, y) - \frac{\nu}{2} \|y\|^2$  for a small regularization parameter  $\nu = \mathcal{O}(n^{-1/2})$ . The distance can be controlled by  $\mathcal{O}(n^{-1/4})$ . Thus, the sample complexity for achieving  $\epsilon$ -uniform convergence for the NC-C case is large than that of the NC-SC case. We leave it for future investigation to see if one could achieve smaller sample complexity in the NC-C case via a better characterization of the extra error brought in by  $y$  in the NC-C setting.

## 4 Algorithmic Stability and Generalization Bounds

Notice that the uniform convergence in Theorems 3.1 and 3.2 has a dependence on the dimension  $d$ . It becomes less meaningful for high-dimensional problems [Lei, 2022, Feldman and Vondrak, 2019]. It remains interesting to build dimension-independent generalization results utilizing the special structure of the algorithms. In this section, we investigate the generalization performance of specific algorithms for nonconvex stochastic minimax optimization problems utilizing stability arguments.

### 4.1 Stability and Generalization

Existing literature on stability arguments in minimax optimization often rely on stability notions based on function values [Farnia and Ozdaglar, 2021, Lei et al., 2021, Zhang et al., 2021a]. In order to derive bounds on the generalization in terms of primal stationarity, we introduce the following novel notions of uniform stability on gradients of the primal function, called *uniform primal stability*.

**Definition 4.1 (Uniform Primal Stability)** *A randomized algorithm  $\mathcal{A}$  is  $\delta$ -uniformly primal stable if for every two neighboring dataset  $S, S'$  which differ in only one sample, for every  $\xi \in \Xi$  we have*

$$\sup_{\xi} \mathbb{E}_{\mathcal{A}} \|\nabla f(\mathcal{A}_x(S), y^*(\mathcal{A}_x(S)); \xi) - \nabla f(\mathcal{A}_x(S'), y^*(\mathcal{A}_x(S')); \xi)\|^2 \leq \delta^2, \quad (13)$$

where  $\nabla f = (\nabla_x f, \nabla_y f)^\top$  denotes the full gradient.

The following theorem connects stability and generalization in minimax optimization problems. We defer the proof to Appendix D.

**Theorem 4.1 (Stability and Generalization, NC-SC)** *Let  $\mathcal{A}$  be a  $\delta$ -uniformly primal stable algorithm. For any function  $f$  satisfying Assumption 2.1 with  $\mu > 0$ , we have*

$$\mathbb{E}_{\mathcal{A}, S} \|\nabla \Phi(\mathcal{A}_x(S)) - \nabla \Phi_S(\mathcal{A}_x(S))\| \leq (1 + \kappa) \left( 4\delta + \frac{G}{\sqrt{n}} \right). \quad (14)$$

To the best of our knowledge, this is the first result that connects uniformly stable algorithms and generalization errors in minimax optimization measured by primal stationarity. As a comparison, in the minimization case, Lei [2022, Theorem 2] proved that the gap between the empirical and population gradients is  $\mathcal{O}\left(\delta + \frac{1}{\sqrt{n}}\right)$ , while Theorem 4.1 has an additional dependence on the condition number  $\kappa$  that comes from the minimax structure.

In the NC-C case, the uniform primal stability is less meaningful as  $y^*(\cdot)$  is not well-defined. Instead, we use the following notion of *uniform primal argument stability*.

**Definition 4.2 (Uniform Primal Argument Stability)** *A randomized algorithm  $\mathcal{A}$  is  $\delta$ -uniformly primal*

argument stable if for every two dataset  $S, S'$  which differ in only one sample,

$$\mathbb{E}_{\mathcal{A}} \|\mathcal{A}_x(S) - \mathcal{A}_x(S')\|^2 \leq \delta^2.$$

The following theorem connects argument stability and generalization in NC-C case, measured by primal Moreau envelope stationarity.

**Theorem 4.2 (Stability and Generalization, NC-C)** *Let  $\mathcal{A}$  be a  $\delta$ -uniformly primal argument stable algorithm. For any function  $f$  satisfying Assumption 2.1 with  $\mu = 0$ , we have*

$$\mathbb{E}_{\mathcal{A}, S} \left\| \nabla \Phi^{1/(2L)}(\mathcal{A}_x(S)) - \nabla \Phi_S^{1/(2L)}(\mathcal{A}_x(S)) \right\| \leq \mathcal{O}\left(\delta^{1/6} + n^{-1/8}\right). \quad (15)$$

We defer the proof to Appendix E. Note that the analysis also leverages the idea of adding regularization to create a surrogate NC-SC problem, as we did in Section 3.2. This result yields the relationship between stability and generalization in NC-C problems measured by primal stationarity. Different from the minimization case, the perturbation on the dataset incurs errors on both the function gradients and the dual maximizers, which requires more careful analysis to derive the final generalization bound. With Theorems 4.1 and 4.2, to obtain the generalization bounds of algorithms designed for NC-SC and NC-C minimax optimization problems, it suffices to derive the stability of specific algorithms.

## 4.2 Generalization of SGDA

In this subsection, we study the generalization bounds of the classical *stochastic gradient descent ascent* (SGDA) for minimax optimization problems in both NC-SC and NC-C cases. Recall the procedures of SGDA: in each iteration  $t$ ,

$$\begin{cases} x_{t+1} = \mathbf{proj}_{\mathcal{X}}(x_t - \alpha_t^x \nabla_x f(x_t, y_t; \xi_t)), \\ y_{t+1} = \mathbf{proj}_{\mathcal{Y}}(y_t + \alpha_t^y \nabla_y f(x_t, y_t; \xi_t)), \end{cases} \quad (16)$$

where  $(\alpha_t^x, \alpha_t^y)$  are the stepsizes. Farnia and Ozdaglar [2021] investigated the  $\delta$ -stability of SGDA. Together with Theorems 4.1 and 4.2, we have the following generalization errors in NC-SC and NC-C cases, respectively.

**Corollary 4.1 (Generalization of SGDA, NC-SC)** *Assume the function  $f$  is NC-SC as defined in Assumption 2.1 with  $\mu > 0$ , then if we run SGDA for  $T$  iterations with stepsize  $(\alpha_t^x, \alpha_t^y) = \left(\frac{c}{t}, \frac{cr^2}{t}\right)$  for some constant  $c > 0$  and  $1 \leq r < \kappa$ , let  $\zeta_1 = (cL(r+1) + 1)^{-1}$ , we have*

$$\mathbb{E}_{S, \mathcal{A}} \|\nabla \Phi(\mathcal{A}_x(S)) - \nabla \Phi_S(\mathcal{A}_x(S))\| \leq \mathcal{O}\left(\kappa^{1+\zeta_1} \left(\frac{T^{1-\zeta_1}}{n} + \frac{1}{\sqrt{n}}\right)\right), \quad (17)$$

where  $(\mathcal{A}_x(S), \mathcal{A}_y(S)) = (x_T, y_T)$  is the output of SGDA.

**Corollary 4.2 (Generalization of SGDA, NC-C)** *Assume the function  $f$  is NC-C as defined in Assumption 2.1 with  $\mu = 0$ , then if we run SGDA for  $T$  iterations with stepsize  $\max\{\alpha_t^x, \alpha_t^y\} \leq \frac{c}{t}$  for some constant  $c > 0$ , let  $\zeta_2 = (cL + 1)^{-1}$  then we have*

$$\mathbb{E}_{S, \mathcal{A}} \left\| \nabla \Phi^{1/(2L)}(\mathcal{A}_x(S)) - \nabla \Phi_S^{1/(2L)}(\mathcal{A}_x(S)) \right\| \leq \mathcal{O}\left(\left(\frac{T^{1-\zeta_2}}{n}\right)^{1/6} + \left(\frac{1}{n}\right)^{1/8}\right), \quad (18)$$

where  $(\mathcal{A}_x(S), \mathcal{A}_y(S)) = (x_T, y_T)$  is the output of SGDA.

The proof relies on the stability results in Farnia and Ozdaglar [2021], which we defer to Appendix F. Compared to the generalization bounds in Theorems 3.1 and 3.2 that use uniform convergence, the generalization bounds of SGDA avoid the dependence on the dimension  $d$ . However, the dependence on  $n$  of generalization bounds of SGDA becomes worse compared to uniform convergence in the NC-C setting.

### 4.3 Generalization of Sampling-determined Algorithms

In this subsection, we consider an extension of *sampling-determined algorithm* (SDA) class proposed in [Lei \[2022\]](#). For completeness, we present its definition below.

**Definition 4.3 (Sampling-determined Algorithm)** *Let  $\mathcal{A}$  be an algorithm that randomly chooses an index sequence  $I(\mathcal{A}) = \{i_t\}$  from the dataset to build stochastic gradients. We say  $\mathcal{A}$  is sampling-determined if its output is independent of the sample  $\xi_i$  for any  $i \notin I(\mathcal{A})$ .*

SDA covers a wide range of algorithms, including classical SGDA, stochastic extragradient, and some adaptive variants of SGDA [[Yang et al., 2022a](#)] in minimax optimization literature. [Lei \[2022\]](#) derives  $\delta$ -stability of SDA leveraging its sampling-determined property. Following their techniques and combining with [Theorems 4.1 and 4.2](#), we obtain the following generalization bounds for SDA in both NC-SC and NC-C scenarios.

**Corollary 4.3 (Generalization of SDA, NC-SC)** *Assume the function  $f$  is NC-SC as defined in [Assumption 2.1](#) with  $\mu > 0$ . If we run a SDA algorithm  $\mathcal{A}$  for  $T$  iterations, we have*

$$\mathbb{E}_{S, \mathcal{A}} \|\nabla \Phi(\mathcal{A}_x(S)) - \nabla \Phi_S(\mathcal{A}_x(S))\| \leq \mathcal{O}\left(\kappa \left(\sqrt{\frac{T}{n}} + \frac{1}{\sqrt{n}}\right)\right) \quad (19)$$

Compared with [Corollary 4.1](#), the generalization bound of SDA does not require specific stepsizes and applies to a wider class of algorithms.

**Corollary 4.4 (Generalization of SDA, NC-C)** *Assume the function  $f$  is NC-C as defined in [Assumption 2.1](#) with  $\mu = 0$ . If we run a SDA algorithm  $\mathcal{A}$  for  $T$  iterations, we have*

$$\mathbb{E}_{S, \mathcal{A}} \left\| \nabla \Phi^{1/(2L)}(\mathcal{A}_x(S)) - \nabla \Phi_S^{1/(2L)}(\mathcal{A}_x(S)) \right\| \leq \mathcal{O}\left(\left(\frac{T}{n}\right)^{1/12} + \left(\frac{1}{n}\right)^{1/8}\right). \quad (20)$$

Compared with [Corollary 4.2](#), the generalization bound of SDA algorithm has a worse dependence on sample size  $n$ . Due to the  $T/n$  term in [Corollaries 4.3 and 4.4](#), to achieve good generalization, SDA should have less than one pass of the dataset. On the other hand, SGDA may use the stepsize to control  $\zeta$  and can do multiple pass over the dataset.

## 5 Conclusion and Future Directions

We take an initial step toward understanding the generalization performances of NC-SC and NC-C minimax problems measured by the first-order stationarity from both uniform convergence and stability argument perspectives. Several future directions are worth further investigation. It remains interesting to see whether we can improve the uniform convergence and stability results under the NC-C setting, particularly the dependence on sample size  $n$ . Another possible direction is to investigate the generalization performances for specific applications. Some studies in stochastic minimization show that specific machine learning models (e.g., generalized linear models) enjoy dimension-free uniform convergence bounds [[Amir et al., 2022](#), [Davis and Drusvyatskiy, 2022](#)]. It would be interesting to see if such dimension-free uniform convergence property also holds for some minimax applications.

## References

- Idan Amir, Roi Livni, and Nathan Srebro. Thinking outside the ball: Optimal learning with gradient descent for generalized linear stochastic convex optimization. *arXiv preprint arXiv:2202.13328*, 2022. (Cited on 1.2, 5)
- Raef Bassily, Vitaly Feldman, Cristóbal Guzmán, and Kunal Talwar. Stability of stochastic gradient descent on nonsmooth convex losses. *Advances in Neural Information Processing Systems*, 33, 2020. (Cited on 1.2)
- Digvijay Boob and Cristóbal Guzmán. Optimal algorithms for differentially private stochastic monotone variational inequalities and saddle-point problems. *arXiv preprint arXiv:2104.02988*, 2021. (Cited on 1, 1.2)
- Radu Ioan Boț and Axel Böhm. Alternating proximal-gradient steps for (stochastic) nonconvex-concave minimax problems. *arXiv preprint arXiv:2007.13605*, 2020. (Cited on 1.2)
- Olivier Bousquet and André Elisseeff. Stability and generalization. *The Journal of Machine Learning Research*, 2:499–526, 2002. (Cited on 1, 1.2)
- Bo Dai, Niao He, Yunpeng Pan, Byron Boots, and Le Song. Learning from conditional distributions via dual embeddings. In *Artificial Intelligence and Statistics*, pages 1458–1467, 2017. (Cited on 1)
- Bo Dai, Albert Shaw, Lihong Li, Lin Xiao, Niao He, Zhen Liu, Jianshu Chen, and Le Song. SBEED: Convergent reinforcement learning with nonlinear function approximation. In *International Conference on Machine Learning*, pages 1125–1134, 2018. (Cited on 1)
- Damek Davis and Dmitriy Drusvyatskiy. Stochastic model-based minimization of weakly convex functions. *SIAM Journal on Optimization*, 29(1):207–239, 2019. (Cited on 1, 2, 2, 3.2)
- Damek Davis and Dmitriy Drusvyatskiy. Graphical convergence of subgradients in nonconvex optimization and learning. *Mathematics of Operations Research*, 47(1):209–231, 2022. (Cited on 1, 1.2, 2, 3.3, 1, 5, C)
- Farzan Farnia and Asuman Ozdaglar. Train simultaneously, generalize better: Stability of gradient-based minimax learners. In *International Conference on Machine Learning*, pages 3174–3185. PMLR, 2021. (Cited on 1, 1.2, 3.1, 4.1, 4.2, 4.2, F, F)
- Vitaly Feldman and Jan Vondrak. High probability generalization bounds for uniformly stable algorithms with nearly optimal rate. In *Conference on Learning Theory*, pages 1270–1279. PMLR, 2019. (Cited on 4)
- Ronald A Fisher. On the mathematical foundations of theoretical statistics. *Philosophical transactions of the Royal Society of London. Series A, containing papers of a mathematical or physical character*, 222(594-604): 309–368, 1922. (Cited on 1.2)
- Dylan J Foster, Ayush Sekhari, and Karthik Sridharan. Uniform convergence of gradients for non-convex learning and optimization. *Advances in Neural Information Processing Systems*, 31, 2018. (Cited on 1)
- Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. *Advances in neural information processing systems*, 27, 2014. (Cited on 1)
- Zhishuai Guo, Yi Xu, Wotao Yin, Rong Jin, and Tianbao Yang. A novel convergence analysis for algorithms of the adam family and beyond. *arXiv preprint arXiv:2104.14840*, 2021. (Cited on 1.2)
- Yuze Han, Guangzeng Xie, and Zhihua Zhang. Lower complexity bounds of finite-sum optimization problems: The results and construction. *arXiv preprint arXiv:2103.08280*, 2021. (Cited on 1.2)
- Moritz Hardt, Ben Recht, and Yoram Singer. Train faster, generalize better: Stability of stochastic gradient descent. In *International Conference on Machine Learning*, pages 1225–1234. PMLR, 2016. (Cited on 1, 1.2, 4)

- Jiawei Huang and Nan Jiang. On the convergence rate of density-ratio based off-policy gradient methods. In *Neural Information Processing Systems Offline Reinforcement Learning Workshop*, 2020. (Cited on 1)
- Anton J Kleywegt, Alexander Shapiro, and Tito Homem-de Mello. The sample average approximation method for stochastic discrete optimization. *SIAM Journal on Optimization*, 12(2):479–502, 2002. (Cited on 1, 1.2, 2, B)
- Yegor Klochkov and Nikita Zhivotovskiy. Stability and deviation optimal risk bounds with convergence rate  $o(1/n)$ . *Advances in Neural Information Processing Systems*, 34, 2021. (Cited on 1, 1.2)
- Qi Lei, Jason Lee, Alex Dimakis, and Constantinos Daskalakis. Sgd learns one-layer networks in w-gans. In *International Conference on Machine Learning*, pages 5799–5808. PMLR, 2020. (Cited on 1)
- Yunwen Lei. Stability and generalization of stochastic optimization with nonconvex and nonsmooth problems. *arXiv preprint arXiv:2206.07082*, 2022. (Cited on 1, 1.2, 4, 4.1, 4.3, 4.3, D, E, E, E, 4, E, G)
- Yunwen Lei, Zhenhuan Yang, Tianbao Yang, and Yiming Ying. Stability and generalization of stochastic gradient methods for minimax problems. In *International Conference on Machine Learning*, pages 6175–6186. PMLR, 2021. (Cited on 1, 1.2, 3.1, 4.1)
- Haochuan Li, Yi Tian, Jingzhao Zhang, and Ali Jadbabaie. Complexity lower bounds for nonconvex-strongly-concave min-max optimization. *Advances in Neural Information Processing Systems*, 34, 2021. (Cited on 1.2)
- Tianyi Lin, Chi Jin, and Michael Jordan. On gradient descent ascent for nonconvex-concave minimax problems. In *International Conference on Machine Learning*, pages 6083–6093. PMLR, 2020a. (Cited on 1, 1.2, 2, 2, E, 4, E)
- Tianyi Lin, Chi Jin, and Michael I Jordan. Near-optimal algorithms for minimax optimization. In *Conference on Learning Theory*, pages 2738–2779. PMLR, 2020b. (Cited on 1.2, 2)
- Songtao Lu, Ioannis Tsaknakis, Mingyi Hong, and Yongxin Chen. Hybrid block successive approximation for one-sided non-convex min-max problems: algorithms and applications. *IEEE Transactions on Signal Processing*, 68:3676–3691, 2020. (Cited on 1.2)
- Luo Luo, Haishan Ye, Zhichao Huang, and Tong Zhang. Stochastic recursive gradient descent ascent for stochastic nonconvex-strongly-concave minimax problems. *Advances in Neural Information Processing Systems*, 33, 2020. (Cited on 1, 1.2, 3.1)
- Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards deep learning models resistant to adversarial attacks. In *International Conference on Learning Representations*, 2018. (Cited on 1)
- Song Mei, Yu Bai, and Andrea Montanari. The landscape of empirical risk for nonconvex losses. *The Annals of Statistics*, 46(6A):2747–2774, 2018. (Cited on 1, 1.2, 3.3, 1)
- Maher Nouiehed, Maziar Sanjabi, Tianjian Huang, Jason D Lee, and Meisam Razaviyayn. Solving a class of non-convex min-max games using iterative first order methods. *Advances in Neural Information Processing Systems*, 32:14934–14942, 2019. (Cited on 1.2)
- Dmitrii M Ostrovskii, Andrew Lowy, and Meisam Razaviyayn. Efficient search of first-order nash equilibria in nonconvex-concave smooth min-max problems. *SIAM Journal on Optimization*, 31(4):2508–2538, 2021. (Cited on 1.2)
- Asuman Ozdaglar, Sarath Pattathil, Jiawei Zhang, and Kaiqing Zhang. What is a good metric to study generalization of minimax learners? *arXiv preprint arXiv:2206.04502*, 2022. (Cited on 1)

- Hassan Rafique, Mingrui Liu, Qihang Lin, and Tianbao Yang. Weakly-convex-concave min-max optimization: provable algorithms and applications in machine learning. *Optimization Methods and Software*, pages 1–35, 2021. (Cited on [1.2](#))
- Maziar Sanjabi, Jimmy Ba, Meisam Razaviyayn, and Jason D Lee. On the convergence and robustness of training gans with regularized optimal transport. In *Proceedings of the 32nd International Conference on Neural Information Processing Systems*, pages 7091–7101, 2018. (Cited on [1](#))
- Shai Shalev-Shwartz, Ohad Shamir, Nathan Srebro, and Karthik Sridharan. Stochastic convex optimization. In *COLT*, volume 2, page 5, 2009. (Cited on [1.2](#), [3.1](#))
- Shai Shalev-Shwartz, Ohad Shamir, Nathan Srebro, and Karthik Sridharan. Learnability, stability and uniform convergence. *The Journal of Machine Learning Research*, 11:2635–2670, 2010. (Cited on [1](#))
- Pranay Sharma, Rohan Panda, Gauri Joshi, and Pramod K Varshney. Federated minimax optimization: Improved convergence analyses and algorithms. *arXiv preprint arXiv:2203.04850*, 2022. (Cited on [1.2](#))
- Aman Sinha, Hongseok Namkoong, and John Duchi. Certifying some distributional robustness with principled adversarial training. In *International Conference on Learning Representations*, 2018. (Cited on [1](#))
- Kiran K Thekumparampil, Prateek Jain, Praneeth Netrapalli, and Sewoong Oh. Efficient algorithms for smooth minimax optimization. *Advances in Neural Information Processing Systems*, 32:12680–12691, 2019. (Cited on [1.2](#), [2](#), [3.2](#), [C](#))
- Aad W Van der Vaart. *Asymptotic statistics*, volume 3. Cambridge university press, 2000. (Cited on [1.2](#))
- Vladimir N Vapnik. An overview of statistical learning theory. *IEEE transactions on neural networks*, 10(5): 988–999, 1999. (Cited on [1.2](#), [3.1](#), [3.2](#))
- Di Wang, Minwei Ye, and Jinhui Xu. Differentially private empirical risk minimization revisited: Faster and more general. *Advances in Neural Information Processing Systems*, 30, 2017. (Cited on [1.2](#))
- Yisen Wang, Xingjun Ma, James Bailey, Jinfeng Yi, Bowen Zhou, and Quanquan Gu. On the convergence and robustness of adversarial training. In *International Conference on Machine Learning*, pages 6586–6595. PMLR, 2019. (Cited on [1](#))
- Zi Xu, Huiling Zhang, Yang Xu, and Guanghui Lan. A unified single-loop alternating gradient projection algorithm for nonconvex-concave and convex-nonconcave minimax problems. *arXiv preprint arXiv:2006.02032*, 2020. (Cited on [1.2](#), [2](#))
- Yan Yan, Yi Xu, Qihang Lin, Wei Liu, and Tianbao Yang. Optimal epoch stochastic gradient descent ascent methods for min-max optimization. *Advances in Neural Information Processing Systems*, 33:5789–5800, 2020. (Cited on [1.2](#))
- Junchi Yang, Negar Kiyavash, and Niao He. Global convergence and variance reduction for a class of nonconvex-nonconcave minimax problems. *Advances in Neural Information Processing Systems*, 33, 2020a. (Cited on [1.2](#))
- Junchi Yang, Siqi Zhang, Negar Kiyavash, and Niao He. A catalyst framework for minimax optimization. In *Advances in Neural Information Processing Systems*, 2020b. (Cited on [1](#), [1.2](#), [3.2](#))
- Junchi Yang, Xiang Li, and Niao He. Nest your adaptive algorithm for parameter-agnostic nonconvex minimax optimization. *arXiv preprint arXiv:2206.00743*, 2022a. (Cited on [4.3](#))
- Junchi Yang, Antonio Orvieto, Aurelien Lucchi, and Niao He. Faster single-loop algorithms for minimax optimization without strong concavity. In *International Conference on Artificial Intelligence and Statistics*, pages 5485–5517. PMLR, 2022b. (Cited on [1.2](#), [2](#))

- Zhenhuan Yang, Shu Hu, Yunwen Lei, Kush R Varshney, Siwei Lyu, and Yiming Ying. Differentially private sgda for minimax problems. *arXiv preprint arXiv:2201.09046*, 2022c. (Cited on 1)
- Jiawei Zhang, Peijun Xiao, Ruoyu Sun, and Zhiquan Luo. A single-loop smoothed gradient descent-ascent algorithm for nonconvex-concave min-max problems. *Advances in Neural Information Processing Systems*, 33:7377–7389, 2020. (Cited on 1.2)
- Junyu Zhang, Mingyi Hong, Mengdi Wang, and Shuzhong Zhang. Generalization bounds for stochastic saddle point problems. In *International Conference on Artificial Intelligence and Statistics*, pages 568–576. PMLR, 2021a. (Cited on 1, 1.2, 3.1, 3.2, 4.1, C, C)
- Siqi Zhang, Junchi Yang, Cristóbal Guzmán, Negar Kiyavash, and Niao He. The complexity of nonconvex-strongly-concave minimax optimization. In *Uncertainty in Artificial Intelligence*, pages 482–492. PMLR, 2021b. (Cited on 1.2, 3.1)
- Xuan Zhang, Necdet Serhat Aybat, and Mert Gurbuzbalaban. Sapd+: An accelerated stochastic method for nonconvex-concave minimax problems. *arXiv preprint arXiv:2205.15084*, 2022. (Cited on 1.2)
- Renbo Zhao. A primal dual smoothing framework for max-structured nonconvex optimization. *arXiv preprint arXiv:2003.04375*, 2020. (Cited on 1.2)

## A Additional Definitions and Tools

For convenience, we summarize the notations commonly used throughout the paper.

- Population minimax problem and its primal function<sup>2</sup>

$$F(x, y) \triangleq \mathbb{E}_\xi f(x, y; \xi), \quad \Phi(x) \triangleq \max_{y \in \mathcal{Y}} F(x, y), \quad y^*(x) \triangleq \operatorname{argmax}_{y \in \mathcal{Y}} F(x, y).$$

- Empirical minimax problem and its primal function

$$F_S(x, y) \triangleq \frac{1}{n} \sum_{i=1}^n f(x, y; \xi_i), \quad \Phi_S(x) \triangleq \max_{y \in \mathcal{Y}} F_S(x, y), \quad y_S^*(x) \triangleq \operatorname{argmax}_{y \in \mathcal{Y}} F_S(x, y).$$

- Moreau envelope and corresponding proximal point:

$$\begin{aligned} \Phi^\lambda(x) &\triangleq \min_{z \in \mathcal{X}} \left\{ \Phi(z) + \frac{1}{2\lambda} \|z - x\|^2 \right\}, & \mathbf{prox}_{\lambda\Phi}(x) &\triangleq \operatorname{argmin}_{z \in \mathcal{X}} \left\{ \Phi(z) + \frac{1}{2\lambda} \|z - x\|^2 \right\}, \\ \Phi_S^\lambda(x) &\triangleq \min_{z \in \mathcal{X}} \left\{ \Phi_S(z) + \frac{1}{2\lambda} \|z - x\|^2 \right\}, & \mathbf{prox}_{\lambda\Phi_S}(x) &\triangleq \operatorname{argmin}_{z \in \mathcal{X}} \left\{ \Phi_S(z) + \frac{1}{2\lambda} \|z - x\|^2 \right\}. \end{aligned}$$

- $\mathcal{G}_\Phi(x)$ : gradient mapping (generalized gradient) of a function  $\Phi$ .
- $\|\cdot\|$ :  $\ell_2$ -norm.
- $\nabla f = (\nabla_x f, \nabla_y f)$ : the gradient of a function  $f$ .
- $\mathbf{proj}_{\mathcal{X}}(x')$ : the projection operator.
- $\mathcal{A}(S) \triangleq (\mathcal{A}_x(S), \mathcal{A}_y(S))$ : the output of an algorithm  $\mathcal{A}$  on the empirical minimax problem (2) with dataset  $S$ .
- NC / WC: nonconvex, weakly convex.
- NC-SC / NC-C: nonconvex-(strongly)-concave.
- SOTA: state-of-the-art.
- $d$ : dimension number of  $\mathcal{X}$ .
- $\kappa$ : condition number  $\frac{L}{\mu}$ ,  $L$ : Lipschitz smoothness parameter,  $\mu$ : strong concavity parameter.
- $\tilde{O}(\cdot)$  hides poly-logarithmic factors.
- $f = \Omega(g)$  if  $f(x) \geq cg(x)$  for some  $c > 0$  and nonnegative functions  $f$  and  $g$ .
- We say a function  $g : \mathcal{X} \rightarrow \mathbb{R}$  is convex if  $\forall x_1, x_2 \in \mathcal{X}$  and  $p \in [0, 1]$ , we have  $g(px_1 + (1-p)x_2) \geq pg(x_1) + (1-p)g(x_2)$ .
- A function  $h : \mathcal{X} \rightarrow \mathbb{R}$  is  $L$ -smooth<sup>3</sup> if  $h$  is continuously differentiable in  $\mathcal{X}$  and there exists a constant  $L > 0$  such that  $\|\nabla h(x_1) - \nabla h(x_2)\| \leq L\|x_1 - x_2\|$  holds for any  $x_1, x_2$ .

<sup>2</sup> Another commonly used convergence criterion in minimax optimization is the *first-order stationarity of  $F$* , i.e.,  $\|\nabla_x F\| \leq \epsilon$  and  $\|\nabla_y F\| \leq \epsilon$  (or its corresponding gradient mapping) [Lin et al., 2020a, Xu et al., 2020]. We refer readers to Lin et al. [2020a], Yang et al. [2022b] for a thorough comparison of these two measurements. In this paper, we always stick to the convergence measured by the stationarity of the primal function.

<sup>3</sup> Here the smoothness definition for single-variable functions is subtly different from that of two-variable functions in Definition 2.1, so we list it here for completeness.

For completeness, we introduce the definition of a sub-Gaussian random variable and related lemma, which are important tools in the analysis.

**Definition A.1 (Sub-Gaussian Random Variable)** *A random variable  $\eta$  is a zero-mean sub-Gaussian random variable with variance proxy  $\sigma_\eta^2$  if  $\mathbb{E}\eta = 0$  and either of the following two conditions hold:*

$$(a) \mathbb{E}[\exp(s\eta)] \leq \exp\left(\frac{\sigma_\eta^2 s^2}{2}\right) \text{ for any } s \in \mathbb{R}; \quad (b) \mathbb{P}(|\eta| \geq t) \leq 2 \exp\left(-\frac{t^2}{2\sigma_\eta^2}\right) \text{ for any } t > 0.$$

We use the following McDiarmid's inequality to show that a random variable is sub-Gaussian.

**Lemma A.1 (McDiarmid's inequality)** *Let  $\eta_1, \dots, \eta_n \in \mathbb{R}$  be independent random variables. Let  $h: \mathbb{R}^n \rightarrow \mathbb{R}$  be any function with the  $(c_1, \dots, c_n)$ -bounded differences property: for every  $i = 1, \dots, n$  and every  $(\eta_1, \dots, \eta_n)$ , and  $(\eta'_1, \dots, \eta'_n)$  that differ only in the  $i$ -th coordinate ( $\eta_j = \eta'_j$  for all  $j \neq i$ ), we have*

$$|h(\eta_1, \dots, \eta_n) - h(\eta'_1, \dots, \eta'_n)| \leq c_i.$$

For any  $t > 0$ , it holds that

$$\mathbb{P}(|h(\eta_1, \dots, \eta_n) - \mathbb{E}h(\eta_1, \dots, \eta_n)| \geq t) \leq 2 \exp\left(-\frac{2t^2}{\sum_{i=1}^n c_i^2}\right).$$

**Lemma A.2 (Properties of  $\Phi$  and  $\Phi^\lambda$ , Full Version)** *In the NC-SC setting ( $\mu > 0$ ), both  $\Phi(x)$  and  $\Phi_S(x)$  are  $\tilde{L} \triangleq L(1 + \kappa)$ -smooth with the condition number  $\kappa \triangleq L/\mu$ , both  $y^*(x)$  and  $y_S^*(x)$  are  $\kappa$ -Lipschitz continuous and  $\nabla\Phi(x) = \nabla_x F(x, y^*(x))$ ,  $\nabla\Phi_S(x) = \nabla_x F_S(x, y_S^*(x))$ . In the NC-C setting ( $\mu = 0$ ), the primal function  $\Phi$  is  $L$ -weakly convex, and its Moreau envelope  $\Phi^\lambda(x)$  is differentiable, Lipschitz smooth, also*

$$\nabla\Phi^\lambda(x) = \lambda^{-1}(x - \hat{x}), \quad \|\nabla\Phi^\lambda(x)\| \geq \text{dist}(0, \partial\Phi(\hat{x})), \quad (21)$$

where  $\hat{x} = \text{prox}_{\lambda\Phi}(x)$  and  $0 < \lambda < 1/L$ .

For completeness, we formally define the stationary point here. Note that the generalized gradient is defined on  $\mathcal{X}$  while the Moreau envelope is defined on the whole domain  $\mathbb{R}^d$ .

**Definition A.2 (Stationary Point)** *Let  $\epsilon > 0$ , for an  $\tilde{L}$ -smooth function  $\Phi: \mathcal{X} \rightarrow \mathbb{R}$ , we call a point  $x$  an  $\epsilon$ -stationary point of  $\Phi$  if  $\|\mathcal{G}_\Phi(x)\| \leq \epsilon$ , where  $\mathcal{G}_\Phi$  is the gradient mapping (or generalized gradient) defined as  $\mathcal{G}_\Phi(x) \triangleq \tilde{L}(x - \text{proj}_{\mathcal{X}}(x - (1/\tilde{L})\nabla\Phi(x)))$ ; for an  $L$ -weakly convex function  $\Phi$ , we say a point  $x$  an  $\epsilon$ -(nearly)-stationary point of  $\Phi$  if  $\|\nabla\Phi^{1/(2L)}(x)\| \leq \epsilon$ .*

## B Proof of Theorem 3.1

**Proof** To derive the desired generalization bounds, we take an  $v$ -net  $\{x_k\}_{k=1}^Q$  on  $\mathcal{X}$  so that there exists a  $k \in \{1, \dots, Q\}$  for any  $x \in \mathcal{X}$  such that  $\|x - x_k\| \leq v$ . Note that such  $v$ -net exists with  $Q = \mathcal{O}(v^{-d})$  for compact  $\mathcal{X}$  [Kleywegt et al., 2002]. Utilizing the definition of the  $v$ -net, we have

$$\begin{aligned} & \mathbb{E} \max_{x \in \mathcal{X}} \|\nabla\Phi_S(x) - \nabla\Phi(x)\| \\ & \leq \mathbb{E} \max_{x \in \mathcal{X}} [\|\nabla\Phi_S(x) - \nabla\Phi_S(x_k)\| + \|\nabla\Phi_S(x_k) - \nabla\Phi(x_k)\| + \|\nabla\Phi(x_k) - \nabla\Phi(x)\|] \\ & \leq \mathbb{E} \max_{k \in [Q]} \|\nabla\Phi_S(x_k) - \nabla\Phi(x_k)\| + 2L(1 + \kappa)v, \end{aligned} \quad (22)$$

where the last inequality holds as  $\Phi$  and  $\Phi_S$  are  $L(1 + \kappa)$ -smooth following Lemma 2.1. For any  $s > 0$ , we have

$$\begin{aligned}
& \exp\left(s \mathbb{E} \max_{x \in \mathcal{X}} \|\nabla \Phi_S(x) - \nabla \Phi(x)\|\right) \\
& \leq \exp\left(s \left[ \mathbb{E} \max_{k \in [Q]} \|\nabla \Phi_S(x_k) - \nabla \Phi(x_k)\| + 2L(1 + \kappa)v \right]\right) \\
& \leq \mathbb{E} \max_{k \in [Q]} \exp\left(s \left[ \|\nabla \Phi_S(x_k) - \nabla \Phi(x_k)\| + 2L(1 + \kappa)v \right]\right) \\
& \leq \mathbb{E} \sum_{k \in [Q]} \exp\left(s \left[ \|\nabla \Phi_S(x_k) - \nabla \Phi(x_k)\| + 2L(1 + \kappa)v \right]\right) \\
& = \sum_{k \in [Q]} \mathbb{E} \exp\left(s \left[ \|\nabla \Phi_S(x_k) - \nabla \Phi(x_k)\| + 2L(1 + \kappa)v \right]\right),
\end{aligned} \tag{23}$$

where the second inequality uses Jensen's inequality and monotonicity of exponential function, and the third inequality uses summation over  $k \in [Q]$  to handle the dependence issue, i.e., the  $x_k$  in the last line is independent of  $S$ . We use the exponential function as an intermediate step so that the final sample complexity depends on  $\log(Q)$  rather than  $Q$ , which is of order  $\mathcal{O}(v^{-d})$ . Without loss of generality, selecting  $v$  such that  $2L(1 + \kappa)v = \frac{\epsilon}{2}$ , we have

$$\begin{aligned}
& \mathbb{E} \max_{x \in \mathcal{X}} \|\nabla \Phi_S(x) - \nabla \Phi(x)\| \\
& \leq \frac{1}{s} \log \left( \sum_{k \in [Q]} \mathbb{E} \exp(s \|\nabla \Phi(x_k) - \nabla \Phi_S(x_k)\| - \mathbb{E} \|\nabla \Phi(x_k) - \nabla \Phi_S(x_k)\|) \right. \\
& \quad \left. \cdot \exp(s \mathbb{E} \|\nabla \Phi(x_k) - \nabla \Phi_S(x_k)\|) \exp\left(\frac{s\epsilon}{2}\right) \right).
\end{aligned} \tag{24}$$

To upper bound  $\mathbb{E} \|\nabla \Phi(x_k) - \nabla \Phi_S(x_k)\|$ , we use the following observation. Define  $y_{S^{(i)}}^*(x) \triangleq \operatorname{argmax}_{y \in \mathcal{Y}} F_{S^{(i)}}(x, y)$  where  $S = \{\xi_i\}_{i=1}^n$ ,  $S^{(i)} = \{\xi_1, \dots, \xi_{i-1}, \xi'_i, \xi_{i+1}, \dots, \xi_n\}$  and  $\xi'_i$  is i.i.d. from  $\xi_i$ . Since  $x$  is independent of  $S$  or  $S^{(i)}$  for any  $i$ , by Danskin's theorem, we have

$$\begin{aligned}
& \mathbb{E} \|\nabla \Phi(x) - \nabla \Phi_S(x)\| = \mathbb{E} \left\| \mathbb{E}_\xi \nabla_x f(x, y^*(x); \xi) - \frac{1}{n} \sum_{i=1}^n \nabla_x f(x, y_S^*(x); \xi_i) \right\| \\
& = \mathbb{E} \left\| \mathbb{E}_\xi \nabla_x f(x, y^*(x); \xi) - \frac{1}{n} \sum_{i=1}^n \nabla_x f(x, y^*(x); \xi_i) \right. \\
& \quad \left. + \frac{1}{n} \sum_{i=1}^n \nabla_x f(x, y^*(x); \xi_i) - \frac{1}{n} \sum_{i=1}^n \nabla_x f(x, y_S^*(x); \xi_i) \right\| \\
& \leq \mathbb{E} \left\| \mathbb{E}_\xi \nabla_x f(x, y^*(x); \xi) - \frac{1}{n} \sum_{i=1}^n \nabla_x f(x, y^*(x); \xi_i) \right\| \\
& \quad + \mathbb{E} \left\| \frac{1}{n} \sum_{i=1}^n \nabla_x f(x, y^*(x); \xi_i) - \frac{1}{n} \sum_{i=1}^n \nabla_x f(x, y_S^*(x); \xi_i) \right\| \\
& \leq \mathbb{E} \left\| \mathbb{E}_\xi \nabla_x f(x, y^*(x); \xi) - \frac{1}{n} \sum_{i=1}^n \nabla_x f(x, y^*(x); \xi_i) \right\| + L \|y^*(x) - y_S^*(x)\| \\
& \leq \sqrt{\frac{\operatorname{Var}(\nabla_x f)}{n}} + L \|y^*(x) - y_S^*(x)\|,
\end{aligned} \tag{25}$$

where  $\operatorname{Var}(\nabla_x f)$  is the variance of  $\nabla_x f(\cdot, \cdot; \xi)$  and the second inequality holds by smoothness of  $f$ . Since the variance is upper bounded by the second moment:

$$\operatorname{Var}(\nabla_x f) \leq \mathbb{E} \|\nabla_x f(x, y^*(x); \xi)\|^2 \leq G^2, \tag{26}$$

it further holds that

$$\mathbb{E}\|\nabla\Phi(x) - \nabla\Phi_S(x)\| \leq \frac{G}{\sqrt{n}} + L\|y^*(x) - y_S^*(x)\|. \quad (27)$$

To derive an upper bound on  $\|y^*(x) - y_S^*(x)\|$ , we first bound  $\|y_{S^{(i)}}^*(x) - y_S^*(x)\|$  and utilize the stability argument. Since  $f(x, y; \xi)$  is  $\mu$ -strongly concave in  $y$  for any  $x$  and  $\xi$  and  $y_S^*(x)$  is the maximizer of  $F_S(x, \cdot)$ , we have

$$(-F_S(x, y_{S^{(i)}}^*(x))) - (-F_S(x, y_S^*(x))) \geq \frac{\mu}{2}\|y_{S^{(i)}}^*(x) - y_S^*(x)\|^2, \quad (28)$$

On the other hand, we have

$$\begin{aligned} & F_S(x, y_S^*(x)) - F_S(x, y_{S^{(i)}}^*(x)) \\ &= F_{S^{(i)}}(x, y_S^*(x)) - F_{S^{(i)}}(x, y_{S^{(i)}}^*(x)) \\ & \quad + \frac{1}{n} \left[ f(x, y_S^*(x); \xi_i) - f(x, y_{S^{(i)}}^*(x); \xi_i) + f(x, y_{S^{(i)}}^*(x); \xi'_i) - f(x, y_S^*(x); \xi'_i) \right] \\ &\leq F_{S^{(i)}}(x, y_S^*(x)) - F_{S^{(i)}}(x, y_{S^{(i)}}^*(x)) \\ & \quad + \frac{1}{n} \left| f(x, y_{S^{(i)}}^*(x); \xi_i) - f(x, y_S^*(x); \xi_i) \right| + \frac{1}{n} \left| f(x, y_{S^{(i)}}^*(x); \xi'_i) - f(x, y_S^*(x); \xi'_i) \right| \\ &\leq \frac{2G}{n} \|y_{S^{(i)}}^*(x) - y_S^*(x)\|, \end{aligned}$$

where the last inequality holds by Lipschitz continuity and the optimality of  $y_{S^{(i)}}^*(x)$ . Combined with (28), it holds that

$$\|y_{S^{(i)}}^*(x) - y_S^*(x)\| \leq \frac{4G}{\mu n}.$$

In addition, we have

$$\begin{aligned} & \mathbb{E}[F(x, y^*(x)) - F(x, y_S^*(x))] \\ &= \mathbb{E}[F(x, y^*(x)) - F_S(x, y^*(x))] + \mathbb{E}[F_S(x, y^*(x)) - F_S(x, y_S^*(x))] \\ & \quad + \mathbb{E}[F_S(x, y_S^*(x)) - F(x, y_S^*(x))] \\ &\leq \mathbb{E}[F_S(x, y_S^*(x)) - F(x, y_S^*(x))] \\ &= \mathbb{E} \left[ \frac{1}{n} \sum_{i=1}^n f(x, y_S^*(x); \xi_i) - \frac{1}{n} \sum_{i=1}^n \mathbb{E}_\xi f(x, y_S^*(x); \xi) \right] \\ &= \mathbb{E} \left[ \frac{1}{n} \sum_{i=1}^n f(x, y_S^*(x); \xi_i) - \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{\xi_i} f(x, y_{S^{(i)}}^*(x); \xi_i) \right] \\ &= \mathbb{E} \left[ \frac{1}{n} \sum_{i=1}^n f(x, y_S^*(x); \xi_i) - \frac{1}{n} \sum_{i=1}^n f(x, y_{S^{(i)}}^*(x); \xi_i) \right] \\ &\leq G \mathbb{E} \|y_S^*(x) - y_{S^{(i)}}^*(x)\| \\ &\leq \frac{4G^2}{\mu n} \end{aligned} \quad (29)$$

where the first inequality holds as  $y_S^*(x) = \operatorname{argmax}_{y \in \mathcal{Y}} F_S(x, y_S^*(x))$  and  $\mathbb{E}[F(x, y^*(x)) - F_S(x, y^*(x))] = 0$ , the third equality holds as  $y_S^*(x)$  and  $y_{S^{(i)}}^*(x)$  are identical distributed and  $y_{S^{(i)}}^*(x)$  is independent of  $\xi$  by definition, the second inequality holds by Lipschitz continuity of  $f$  on  $y$ , and the last inequality holds by plugging the upper bound on  $\|y_S^*(x) - y_{S^{(i)}}^*(x)\|$ . On the other hand, since  $F(x, y)$  is strongly concave in  $y$  and  $y^*(x) = \operatorname{argmax}_{y \in \mathcal{Y}} F(x, y)$ , it holds that

$$F(x, y^*(x)) - F(x, y_S^*(x)) \geq \frac{\mu}{2} \|y^*(x) - y_S^*(x)\|^2.$$

Therefore, we have

$$\mathbb{E} \|y^*(x) - y_S^*(x)\| \leq \sqrt{\frac{8G^2}{\mu^2 n}}.$$

Plugging into (27), it holds that

$$\mathbb{E} \|\nabla\Phi(x_k) - \nabla\Phi_S(x_k)\| \leq L\sqrt{\frac{8G^2}{\mu^2 n}} + \frac{G}{\sqrt{n}}. \quad (30)$$

Next we show that  $\|\nabla\Phi(x) - \nabla\Phi_S(x)\| - \mathbb{E} \|\nabla\Phi(x) - \nabla\Phi_S(x)\|$  is zero-mean sub-Gaussian. Notice that for any  $\xi'_i$ , we have

$$\begin{aligned} & \|\nabla\Phi(x) - \nabla\Phi_S(x)\| - \|\nabla\Phi(x) - \nabla\Phi_{S^{(i)}}(x)\| \\ & \leq \|\nabla\Phi_S(x) - \nabla\Phi_{S^{(i)}}(x)\| \\ & = \left\| \frac{1}{n} \sum_{j=1}^n \nabla_x f(x, y_S^*(x), \xi_j) - \frac{1}{n} \sum_{j \neq i}^n \nabla_x f(x, y_{S^{(i)}}^*(x), \xi_j) - \frac{1}{n} \nabla_x f(x, y_{S^{(i)}}^*(x), \xi'_i) \right\| \\ & \leq L \|y_{S^{(i)}}^*(x) - y_S^*(x)\| + \frac{1}{n} \|\nabla_x f(x, y_{S^{(i)}}^*(x); \xi'_i) - \nabla_x f(x, y_{S^{(i)}}^*(x); \xi_i)\| \\ & \leq \frac{4LG/\mu + 2G}{n}, \end{aligned} \quad (31)$$

where the first inequality uses triangle inequality, the first equality uses the definition of  $\Phi_S$  and  $\Phi_{S^{(i)}}$ , the third inequality uses the assumption that  $G$  is the uniform upper bound of  $\nabla f(x, y; \xi)$  on  $\mathcal{X} \times \mathcal{Y}$  for any  $\xi$ . By McDiarmid's inequality (Lemma A.1) and the definition of sub-Gaussian random variables, it holds that  $\|\nabla\Phi(x_k) - \nabla\Phi_S(x_k)\| - \mathbb{E} \|\nabla\Phi(x_k) - \nabla\Phi_S(x_k)\|$  is a zero-mean sub-Gaussian random variable with variance proxy  $\sigma^2 \triangleq (2LG/\mu + G)^2/n$ . By the definition of zero-mean sub-Gaussian random variables, it holds that

$$\mathbb{E} \exp(s[\|\nabla\Phi(x_k) - \nabla\Phi_S(x_k)\| - \mathbb{E} \|\nabla\Phi(x_k) - \nabla\Phi_S(x_k)\|]) \leq \exp\left(\frac{s^2 \sigma^2}{2}\right). \quad (32)$$

Plugging (30) and (32) into (24), we have

$$\mathbb{E} \|\nabla\Phi_S(x) - \nabla\Phi(x)\| \leq \frac{\log(Q)}{s} + \frac{s\sigma^2}{2} + L\sqrt{\frac{8G^2}{\mu^2 n}} + \frac{G}{\sqrt{n}} + \frac{\epsilon}{2} \quad (33)$$

Minimizing the right-hand side over  $s$ , we have

$$\begin{aligned} \mathbb{E} \|\nabla\Phi_S(x) - \nabla\Phi(x)\| & \leq 2\sqrt{\frac{\log(Q)\sigma^2}{2}} + L\sqrt{\frac{8G^2}{\mu^2 n}} + \frac{G}{\sqrt{n}} + \frac{\epsilon}{2} \\ & = \sqrt{\frac{2\log(Q)(2LG/\mu + G)^2}{n}} + L\sqrt{\frac{8G^2}{\mu^2 n}} + \frac{G}{\sqrt{n}} + \frac{\epsilon}{2}. \end{aligned} \quad (34)$$

Recall that  $Q = \mathcal{O}(v^{-d})$  with  $v = \epsilon/(4L(1 + \kappa))$ , thus  $\log(Q) = \mathcal{O}(d \log(4L(1 + \kappa)\epsilon^{-1}))$ , which verifies the first statement in the theorem. For the sample complexity, following the discussion on the performance measurement in Section 2, it is easy to derive that it requires

$$n = \mathcal{O}\left(2d\epsilon^{-2}(2LG/\mu + G)^2 \log(4L(1 + \kappa)\epsilon^{-1})\right) = \tilde{\mathcal{O}}(d\kappa^2\epsilon^{-2}) \quad (35)$$

to guarantee that  $\mathbb{E} \|\nabla\Phi_S(x) - \nabla\Phi(x)\| \leq \epsilon$  for any  $x \in \mathcal{X}$ , which concludes the proof.  $\blacksquare$

## C Proof of Theorem 3.2

We first provide the proof of Lemma 3.1.

**Proof** Since  $F(x, y)$  is  $L$ -smooth, it is obvious that  $F(x, y) - \frac{\nu}{2}\|y\|^2$  is  $(L + \nu)$ -smooth. By [Thekumparampil et al. \[2019, Lemma 3\]](#),  $\hat{\Phi}(x)$  is  $(L + \nu)$ -weakly convex in  $x$ . Therefore,  $\hat{\Phi}(x) + \frac{1}{2\lambda}\|x - x'\|^2$  is  $(\frac{1}{\lambda} - (L + \nu))$ -strongly convex in  $x$  for any fixed  $x'$ . Denote  $\hat{y}(x) \triangleq \operatorname{argmax}_{y \in \mathcal{Y}} F(x, y) - \frac{\nu}{2}\|y\|^2$ ,  $y^*(x) \triangleq \operatorname{argmax}_{y \in \mathcal{Y}} F(x, y)$ . It holds that

$$\begin{aligned}
& \frac{1}{2}(1/\lambda - (L + \nu))\|\mathbf{prox}_{\lambda\Phi}(x) - \mathbf{prox}_{\lambda\hat{\Phi}}(x)\|^2 \\
& \leq \hat{\Phi}(\mathbf{prox}_{\lambda\Phi}(x)) + \frac{1}{2\lambda}\|\mathbf{prox}_{\lambda\Phi}(x) - x\|^2 - \hat{\Phi}(\mathbf{prox}_{\lambda\hat{\Phi}}(x)) - \frac{1}{2\lambda}\|\mathbf{prox}_{\lambda\hat{\Phi}}(x) - x\|^2 \\
& = F(\mathbf{prox}_{\lambda\Phi}(x), \hat{y}(\mathbf{prox}_{\lambda\Phi}(x))) - \frac{\nu}{2}\|\hat{y}(\mathbf{prox}_{\lambda\Phi}(x))\|^2 + \frac{1}{2\lambda}\|\mathbf{prox}_{\lambda\Phi}(x) - x\|^2 \\
& \quad - F(\mathbf{prox}_{\lambda\hat{\Phi}}(x), \hat{y}(\mathbf{prox}_{\lambda\hat{\Phi}}(x))) + \frac{\nu}{2}\|\hat{y}(\mathbf{prox}_{\lambda\hat{\Phi}}(x))\|^2 - \frac{1}{2\lambda}\|\mathbf{prox}_{\lambda\hat{\Phi}}(x) - x\|^2 \\
& \leq F(\mathbf{prox}_{\lambda\Phi}(x), y^*(\mathbf{prox}_{\lambda\Phi}(x))) + \frac{1}{2\lambda}\|\mathbf{prox}_{\lambda\Phi}(x) - x\|^2 - \frac{\nu}{2}\|\hat{y}(\mathbf{prox}_{\lambda\Phi}(x))\|^2 \\
& \quad - F(\mathbf{prox}_{\lambda\hat{\Phi}}(x), \hat{y}(\mathbf{prox}_{\lambda\hat{\Phi}}(x))) - \frac{1}{2\lambda}\|\mathbf{prox}_{\lambda\hat{\Phi}}(x) - x\|^2 + \frac{\nu}{2}\|\hat{y}(\mathbf{prox}_{\lambda\hat{\Phi}}(x))\|^2 \\
& \leq F(\mathbf{prox}_{\lambda\Phi}(x), y^*(\mathbf{prox}_{\lambda\Phi}(x))) + \frac{1}{2\lambda}\|\mathbf{prox}_{\lambda\Phi}(x) - x\|^2 - \frac{\nu}{2}\|\hat{y}(\mathbf{prox}_{\lambda\Phi}(x))\|^2 \\
& \quad - F(\mathbf{prox}_{\lambda\hat{\Phi}}(x), y^*(\mathbf{prox}_{\lambda\hat{\Phi}}(x))) - \frac{1}{2\lambda}\|\mathbf{prox}_{\lambda\hat{\Phi}}(x) - x\|^2 + \frac{\nu}{2}\|y^*(\mathbf{prox}_{\lambda\hat{\Phi}}(x))\|^2 \\
& = \Phi(\mathbf{prox}_{\lambda\Phi}(x)) + \frac{1}{2\lambda}\|\mathbf{prox}_{\lambda\Phi}(x) - x\|^2 - \Phi(\mathbf{prox}_{\lambda\hat{\Phi}}(x)) - \frac{1}{2\lambda}\|\mathbf{prox}_{\lambda\hat{\Phi}}(x) - x\|^2 \\
& \quad + \frac{\nu}{2}\|y^*(\mathbf{prox}_{\lambda\hat{\Phi}}(x))\|^2 - \frac{\nu}{2}\|\hat{y}(\mathbf{prox}_{\lambda\Phi}(x))\|^2 \\
& \leq \frac{\nu}{2}\|y^*(\mathbf{prox}_{\lambda\hat{\Phi}}(x))\|^2 - \frac{\nu}{2}\|\hat{y}(\mathbf{prox}_{\lambda\Phi}(x))\|^2 \\
& \leq \frac{\nu D_{\mathcal{Y}}}{2},
\end{aligned} \tag{36}$$

where the first inequality holds by strong convexity of  $\hat{\Phi}(z) + \frac{1}{2\lambda}\|z - x\|^2$  and optimality of  $\mathbf{prox}_{\lambda\hat{\Phi}}(x)$  for  $\min_{z \in \mathcal{X}} \hat{\Phi}(z) + \frac{1}{2\lambda}\|z - x\|^2$ , the first equality holds by definition of  $\hat{\Phi}$ , the second inequality holds by optimality of  $y^*(\mathbf{prox}_{\lambda\Phi}(x)) = \operatorname{argmax}_{y \in \mathcal{Y}} F(\mathbf{prox}_{\lambda\Phi}(x), y)$ , the third inequality holds by optimality of  $\hat{y}(\mathbf{prox}_{\lambda\Phi}(x)) = \operatorname{argmax}_{y \in \mathcal{Y}} F(\mathbf{prox}_{\lambda\Phi}(x), y) - \frac{\nu}{2}\|y\|^2$ , the second equality holds by definition of  $\Phi$ , the fourth inequality holds by optimality of  $\mathbf{prox}_{\lambda\Phi}(x) = \operatorname{argmin}_{x \in \mathcal{X}} \{\Phi(z) + \frac{1}{2\lambda}\|z - x\|^2\}$ , the last inequality holds by the compactness of domain  $\mathcal{Y}$ .  $\blacksquare$

Next, we demonstrate the proof of Theorem 3.2.

**Proof** By Lemma 3.1, we have

$$\begin{aligned}
\|\mathbf{prox}_{\lambda\Phi}(x) - \mathbf{prox}_{\lambda\hat{\Phi}}(x)\| & \leq \sqrt{\frac{\lambda\nu D_{\mathcal{Y}}}{1 - \lambda(L + \nu)}}; \\
\|\mathbf{prox}_{\lambda\Phi_S}(x) - \mathbf{prox}_{\lambda\hat{\Phi}_S}(x)\| & \leq \sqrt{\frac{\lambda\nu D_{\mathcal{Y}}}{1 - \lambda(L + \nu)}}.
\end{aligned}$$

To derive the desired uniform convergence, similar to the proof of Theorem 3.1, we take an  $\nu$ -net  $\{x_k\}_{k=1}^Q$  on  $\mathcal{X}$  so that there exists a  $k \in \{1, \dots, Q\}$  for any  $x \in \mathcal{X}$  such that  $\|x - x_k\| \leq \nu$ . Note that such  $\nu$ -net exists with  $Q = \mathcal{O}(\nu^{-d})$  for compact  $\mathcal{X}$ . We first decompose the error as the approximation error from NC-SC minimax problems to NC-C minimax problems. Then we utilize the  $\nu$ -net to address the dependence between  $S$  and  $\operatorname{argmax}_{x \in \mathcal{X}} \|\nabla\Phi_S^\lambda(x) - \nabla\Phi^\lambda(x)\|$ . First, note that

$$\begin{aligned}
& \mathbb{E} \max_{x \in \mathcal{X}} \|\nabla \Phi_S^\lambda(x) - \nabla \Phi^\lambda(x)\| \\
&= \frac{1}{\lambda} \mathbb{E} \max_{x \in \mathcal{X}} \|\mathbf{prox}_{\lambda\Phi_S}(x) - \mathbf{prox}_{\lambda\Phi}(x)\| \\
&\leq \frac{1}{\lambda} \mathbb{E} \max_{x \in \mathcal{X}} \|\mathbf{prox}_{\lambda\Phi_S}(x) - \mathbf{prox}_{\lambda\hat{\Phi}_S}(x)\| + \|\mathbf{prox}_{\lambda\hat{\Phi}_S}(x) - \mathbf{prox}_{\lambda\hat{\Phi}}(x)\| \\
&\quad + \|\mathbf{prox}_{\lambda\hat{\Phi}}(x) - \mathbf{prox}_{\lambda\Phi}(x)\| \\
&\leq \frac{2}{\lambda} \sqrt{\frac{\lambda\nu D y}{1 - \lambda(L + \nu)}} + \frac{1}{\lambda} \mathbb{E} \max_{x \in \mathcal{X}} \|\mathbf{prox}_{\lambda\hat{\Phi}_S}(x) - \mathbf{prox}_{\lambda\hat{\Phi}}(x)\| \\
&\leq \frac{2}{\lambda} \sqrt{\frac{\lambda\nu D y}{1 - \lambda(L + \nu)}} + \frac{1}{\lambda} \mathbb{E} \max_{x \in \mathcal{X}} [\|\mathbf{prox}_{\lambda\hat{\Phi}_S}(x) - \mathbf{prox}_{\lambda\hat{\Phi}_S}(x_k)\| \\
&\quad + \|\mathbf{prox}_{\lambda\hat{\Phi}_S}(x_k) - \mathbf{prox}_{\lambda\hat{\Phi}}(x_k)\| + \|\mathbf{prox}_{\lambda\hat{\Phi}}(x_k) - \mathbf{prox}_{\lambda\hat{\Phi}}(x)\|] \\
&\leq 2 \sqrt{\frac{\nu D y}{\lambda(1 - \lambda(L + \nu))}} + \frac{1}{\lambda} \mathbb{E} \max_{k \in [Q]} \|\mathbf{prox}_{\lambda\hat{\Phi}_S}(x_k) - \mathbf{prox}_{\lambda\hat{\Phi}}(x_k)\| + \frac{2\nu}{\lambda(1 - \lambda(L + \nu))} \\
&\leq 2 \sqrt{\frac{\nu D y}{\lambda(1 - \lambda(L + \nu))}} + \frac{1}{\lambda s} \log \left( \sum_{k \in [Q]} \mathbb{E} \exp \left( s \|\mathbf{prox}_{\lambda\hat{\Phi}_S}(x_k) - \mathbf{prox}_{\lambda\hat{\Phi}}(x_k)\| \right) \right) \\
&\quad + \frac{2\nu}{\lambda(1 - \lambda(L + \nu))},
\end{aligned} \tag{37}$$

where the first and the third inequality use the triangle inequality, the second inequality uses Lemma 3.1 for  $\Phi$  and  $\Phi_S$ ,  $x_k$  is the closest point to  $x$  in the  $\nu$ -net, the fourth inequality holds by  $(1 - \lambda(L + \nu))^{-1}$ -Lipschitz continuity of proximal operator [Davis and Drusvyatskiy, 2022, Lemma 4.3] since  $F(x, y) - \frac{\nu}{2}\|y\|^2$  is a  $(L + \nu)$ -smooth function, and the last inequality follows a similar argument in (23). All that remains is to bounding  $\mathbb{E} \exp \left( s \|\mathbf{prox}_{\lambda\hat{\Phi}_S}(x) - \mathbf{prox}_{\lambda\hat{\Phi}}(x)\| \right)$  for  $x \in \mathcal{X}$  that is independent of  $S$ . Notice that

$$\begin{aligned}
& \mathbb{E} \exp \left( s \|\mathbf{prox}_{\lambda\hat{\Phi}_S}(x_k) - \mathbf{prox}_{\lambda\hat{\Phi}}(x_k)\| \right) \\
&= \mathbb{E} \exp \left( s \left[ \|\mathbf{prox}_{\lambda\hat{\Phi}_S}(x_k) - \mathbf{prox}_{\lambda\hat{\Phi}}(x_k)\| - \mathbb{E} \|\mathbf{prox}_{\lambda\hat{\Phi}_S}(x_k) - \mathbf{prox}_{\lambda\hat{\Phi}}(x_k)\| \right] \right) \\
&\quad \cdot \exp \left( s \mathbb{E} \|\mathbf{prox}_{\lambda\hat{\Phi}_S}(x_k) - \mathbf{prox}_{\lambda\hat{\Phi}}(x_k)\| \right)
\end{aligned}$$

Next, we show that  $\|\mathbf{prox}_{\lambda\hat{\Phi}_S}(x_k) - \mathbf{prox}_{\lambda\hat{\Phi}}(x_k)\| - \mathbb{E} \|\mathbf{prox}_{\lambda\hat{\Phi}_S}(x_k) - \mathbf{prox}_{\lambda\hat{\Phi}}(x_k)\|$  is a zero-mean sub-Gaussian random variable and  $\mathbb{E} \|\mathbf{prox}_{\lambda\hat{\Phi}_S}(x_k) - \mathbf{prox}_{\lambda\hat{\Phi}}(x_k)\|$  is bounded. Since  $x_k$  is independent of  $S$ , it is sufficient to show an upper bound of the following term where  $x \in \mathcal{X}$  is independent of  $S$ .

$$\mathbb{E} \|\mathbf{prox}_{\lambda\hat{\Phi}_S}(x) - \mathbf{prox}_{\lambda\hat{\Phi}}(x)\|.$$

Recall the definition that

$$\mathbf{prox}_{\lambda\hat{\Phi}}(x) = \operatorname{argmin}_{z \in \mathcal{X}} \left\{ \max_{y \in \mathcal{Y}} \mathbb{E}_\xi f(z, y; \xi) - \frac{\nu}{2} \|y\|^2 + \frac{1}{2\lambda} \|z - x\|^2 \right\}, \tag{38}$$

$$\mathbf{prox}_{\lambda\hat{\Phi}_S}(x) = \operatorname{argmin}_{z \in \mathcal{X}} \left\{ \max_{y \in \mathcal{Y}} \frac{1}{n} \sum_{i=1}^n \left[ f(z, y; \xi_i) - \frac{\nu}{2} \|y\|^2 + \frac{1}{2\lambda} \|z - x\|^2 \right] \right\}. \tag{39}$$

Denote the solution of (38) as  $(z^*(x), y^*(x))$  and the solution of (39) as  $(z_S(x), y_S(x))$ . We need to bound the distance between  $z^*(x)$  and  $z_S(x)$ , note that this  $(z^*(x), y^*(x))$  comes from a strongly-convex-strongly-concave stochastic minimax problem, where the modulus is  $\frac{1-\lambda L}{\lambda}$  and  $\nu$ , respectively; while the other comes from the

sample average approximation counterpart. By Zhang et al. [2021a, Theorem 1 and Appendix A.1], we have the following results:

$$\frac{1-\lambda L}{2\lambda} \mathbb{E} \|z_S(x) - z^*(x)\|^2 + \frac{\nu}{2} \mathbb{E} \|y_S(x) - y^*(x)\|^2 \leq \frac{2\sqrt{2}}{n} \left( \frac{\hat{L}_x^2 \lambda}{1-\lambda L} + \frac{\hat{L}_y^2}{\nu} \right),$$

where  $\hat{L}_x$  is the Lipschitz continuity parameter of  $f(z, y; \xi) + \frac{1}{2\lambda} \|z - x\|^2$  in  $z \in \mathcal{X}$  for any given  $y \in \mathcal{Y}$  and  $\xi$ , and  $\hat{L}_y$  is the Lipschitz continuity parameter of  $f(z, y; \xi) - \frac{\nu}{2} \|y\|^2$  in  $y \in \mathcal{Y}$  for any given  $z \in \mathcal{X}$  and  $\xi$ . More specifically, since  $f(\cdot, \cdot; \xi)$  is  $G$ -Lipschitz continuous for any  $\xi$ , we have

$$\hat{L}_x \leq G + \frac{2\sqrt{D_{\mathcal{X}}}}{\lambda}, \quad \hat{L}_y \leq G + \nu\sqrt{D_{\mathcal{Y}}}.$$

Therefore, we have

$$\begin{aligned} \mathbb{E} \|\mathbf{prox}_{\lambda\hat{\Phi}_S}(x) - \mathbf{prox}_{\lambda\hat{\Phi}}(x)\| &= \mathbb{E} \|z_S(x) - z^*(x)\| \\ &\leq \sqrt{\mathbb{E} \|z_S(x) - z^*(x)\|^2} \leq \sqrt{\frac{2\lambda}{1-\lambda L} \frac{2\sqrt{2}}{n} \left( \frac{\hat{L}_x^2 \lambda}{1-\lambda L} + \frac{\hat{L}_y^2}{\nu} \right)}. \end{aligned} \quad (40)$$

Next, we show that  $\|z_S(x) - z^*(x)\| - \mathbb{E} \|z_S(x) - z^*(x)\|$  is a zero-mean sub-Gaussian random variable. Replacing one sample  $\xi_i$  in  $S$  with an i.i.d. sample  $\xi'_i$  and denote the new dataset as  $S^{(i)}$ , by Zhang et al. [2021a, Lemma 2], it holds that

$$\|z_S(x) - z^*(x)\| - \|z_{S^{(i)}}(x) - z^*(x)\| \leq \|z_S(x) - z_{S^{(i)}}(x)\| \leq \frac{2}{n} \sqrt{\frac{\hat{L}_x^2 \lambda^2}{(1-\lambda L)^2} + \frac{\hat{L}_y^2 \lambda}{\nu(1-\lambda L)}},$$

where  $z_{S^{(i)}}$  follows a similar definition of  $z_S$  but with a different dataset  $S^{(i)}$ . By McDiarmid's inequality (Lemma A.1) and the definition of sub-Gaussian random variables, it holds that  $\|z_S(x) - z^*(x)\| - \mathbb{E} \|z_S(x) - z^*(x)\|$  is a zero-mean sub-Gaussian random variable with variance proxy  $\frac{1}{n} \left( \frac{\hat{L}_x^2 \lambda^2}{(1-\lambda L)^2} + \frac{\hat{L}_y^2 \lambda}{\nu(1-\lambda L)} \right)$ . By the definition of sub-Gaussian random variable and (40), it holds that

$$\begin{aligned} &\mathbb{E} \exp \left( s \left\| \mathbf{prox}_{\lambda\hat{\Phi}_S}(x_k) - \mathbf{prox}_{\lambda\hat{\Phi}}(x_k) \right\| \right) \\ &= \mathbb{E} \exp \left( s \left[ \left\| \mathbf{prox}_{\lambda\hat{\Phi}_S}(x_k) - \mathbf{prox}_{\lambda\hat{\Phi}}(x_k) \right\| - \mathbb{E} \left\| \mathbf{prox}_{\lambda\hat{\Phi}_S}(x_k) - \mathbf{prox}_{\lambda\hat{\Phi}}(x_k) \right\| \right] \right) \\ &\quad \cdot \exp \left( s \mathbb{E} \left\| \mathbf{prox}_{\lambda\hat{\Phi}_S}(x_k) - \mathbf{prox}_{\lambda\hat{\Phi}}(x_k) \right\| \right) \\ &\leq \mathbb{E} \exp \left( s \left[ \left\| \mathbf{prox}_{\lambda\hat{\Phi}_S}(x_k) - \mathbf{prox}_{\lambda\hat{\Phi}}(x_k) \right\| - \mathbb{E} \left\| \mathbf{prox}_{\lambda\hat{\Phi}_S}(x_k) - \mathbf{prox}_{\lambda\hat{\Phi}}(x_k) \right\| \right] \right) \\ &\quad \cdot \exp \left( s \sqrt{\frac{2\lambda}{1-\lambda L} \frac{2\sqrt{2}}{n} \left( \frac{\hat{L}_x^2 \lambda}{1-\lambda L} + \frac{\hat{L}_y^2}{\nu} \right)} \right) \\ &\leq \exp \left( \frac{s^2}{2n} \left( \frac{\hat{L}_x^2 \lambda^2}{(1-\lambda L)^2} + \frac{\hat{L}_y^2 \lambda}{\nu(1-\lambda L)} \right) \right) \exp \left( s \sqrt{\frac{2\lambda}{1-\lambda L} \frac{2\sqrt{2}}{n} \left( \frac{\hat{L}_x^2 \lambda}{1-\lambda L} + \frac{\hat{L}_y^2}{\nu} \right)} \right), \end{aligned} \quad (41)$$

where the second inequality uses definition of zero-mean sub-Gaussian random variable. Combining (41) with (37), for

$$\lambda = \frac{1}{2L}, \quad v = \frac{\epsilon\lambda(1-\lambda L)}{8} = \frac{\epsilon}{32L}, \quad s = \sqrt{2n \log(Q) \left( \frac{\hat{L}_x^2 \lambda^2}{(1-\lambda L)^2} + \frac{\hat{L}_y^2 \lambda}{\nu(1-\lambda L)} \right)^{-1}}, \quad (42)$$

it holds that

$$\begin{aligned}
& \mathbb{E} \max_{x \in \mathcal{X}} \|\nabla \Phi_S^\lambda(x) - \nabla \Phi^\lambda(x)\| \\
& \leq 2 \sqrt{\frac{\nu D_y}{\lambda(1-\lambda(L+\nu))}} + \frac{2\nu}{\lambda(1-\lambda(L+\nu))} \\
& \quad + \frac{1}{\lambda s} \log \left( Q \exp \left( \frac{s^2}{2n} \left( \frac{\hat{L}_x^2 \lambda^2}{(1-\lambda L)^2} + \frac{\hat{L}_y^2 \lambda}{\nu(1-\lambda L)} \right) \right) \right) \\
& \quad + \frac{1}{\lambda s} \log \left( \exp \left( s \sqrt{\frac{2\lambda}{1-\lambda L} \frac{2\sqrt{2}}{n} \left( \frac{\hat{L}_x^2 \lambda}{1-\lambda L} + \frac{\hat{L}_y^2}{\nu} \right)} \right) \right) \\
& \leq 2 \sqrt{\frac{\nu D_y}{\lambda(1-\lambda L)}} + \frac{1}{\lambda s} \log(Q) + \frac{1}{\lambda s} \frac{s^2}{2n} \left( \frac{\hat{L}_x^2 \lambda^2}{(1-\lambda L)^2} + \frac{\hat{L}_y^2 \lambda}{\nu(1-\lambda L)} \right) \\
& \quad + \frac{1}{\lambda s} s \sqrt{\frac{2\lambda}{1-\lambda L} \frac{2\sqrt{2}}{n} \left( \frac{\hat{L}_x^2 \lambda}{1-\lambda L} + \frac{\hat{L}_y^2}{\nu} \right)} + \frac{2\nu}{\lambda(1-\lambda L)} \\
& = 2 \sqrt{\frac{\nu D_y}{\lambda(1-\lambda L)}} + \frac{\log(Q)}{\lambda s} + \frac{1}{\lambda} \frac{s}{2n} \left( \frac{\hat{L}_x^2 \lambda^2}{(1-\lambda L)^2} + \frac{\hat{L}_y^2 \lambda}{\nu(1-\lambda L)} \right) \\
& \quad + \frac{1}{\lambda} \sqrt{\frac{2\lambda}{1-\lambda L} \frac{2\sqrt{2}}{n} \left( \frac{\hat{L}_x^2 \lambda}{1-\lambda L} + \frac{\hat{L}_y^2}{\nu} \right)} + \frac{\epsilon}{4} \\
& = 2 \sqrt{4L\nu D_y} + 4L \sqrt{\frac{\log(Q)}{2n} \left( \frac{\hat{L}_x^2}{L^2} + \frac{\hat{L}_y^2}{\nu L} \right)} + 2L \sqrt{\frac{4\sqrt{2}}{Ln} \left( \frac{\hat{L}_x^2}{L} + \frac{\hat{L}_y^2}{\nu} \right)} + \frac{\epsilon}{4} \\
& = 2 \sqrt{4L\nu D_y} + 4L \sqrt{\frac{\log(Q)}{2n} \left( \frac{\hat{L}_x^2}{L^2} + \frac{\hat{L}_y^2}{\nu L} \right)} \\
& \quad + 2L \sqrt{\frac{4\sqrt{2}}{Ln} \left( \frac{(G+4L\sqrt{D_x})^2}{L} + \frac{(G+\nu\sqrt{D_y})^2}{\nu} \right)} + \frac{\epsilon}{4}.
\end{aligned} \tag{43}$$

Here the first equality holds by the selection of  $\nu$ , the second equality holds by the selection of  $\lambda$  and  $s$ , and the last equality holds by plugging in  $\hat{L}_x$  and  $\hat{L}_y$ . Note that  $\nu$ ,  $s$ , and  $\nu$  are only used for analysis purposes, and  $\lambda$  is only used in the definition of gradient mapping. Thus one has free choices on these parameters. Since  $Q = \mathcal{O}\left(\left(\frac{D_x}{\nu}\right)^d\right)$ , then we choose  $\nu = \tilde{\mathcal{O}}\left(\sqrt{\frac{d}{n}}\right)$  in the right-hand side above, which verifies the first statement. For the sample complexity result, to make sure that the right-hand side of (43) of order  $\mathcal{O}(\epsilon)$ , it suffices to have

$$\nu = \mathcal{O}(\epsilon^2), \quad n = \mathcal{O}\left(\frac{\log(Q)}{\nu} \epsilon^{-2}\right) = \mathcal{O}(d\epsilon^{-4} \log(\epsilon^{-1})), \tag{44}$$

which concludes the proof.  $\blacksquare$

## D Proof of Theorem 4.1

For simplicity we define the following notations:

$$\begin{aligned}
F(x, y) & \triangleq \mathbb{E}_\xi [f(x, y; \xi)], \quad \Phi(x) \triangleq \max_y F(x, y), \quad F_S(x, y) \triangleq \frac{1}{n} \sum_{i=1}^n [f(x, y; \xi_i)], \quad \Phi_S(x) \triangleq \max_y F_S(x, y), \\
y^*(x) & \triangleq \operatorname{argmax}_y F(x, y), \quad y_S^*(x) \triangleq \operatorname{argmax}_y F_S(x, y), \quad \Phi(x; \xi) \triangleq \max_y f(x, y; \xi),
\end{aligned} \tag{45}$$

and the Moreau envelope of a function  $\Phi$ :

$$\Phi^\lambda(x) \triangleq \min_{z \in \mathcal{X}} \left\{ \Phi(x) + \frac{1}{2\lambda} \|z - x\|_2^2 \right\}, \quad \text{prox}_{\lambda\Phi}(x) \triangleq \underset{z \in \mathcal{X}}{\text{argmin}} \left\{ \Phi(x) + \frac{1}{2\lambda} \|z - x\|_2^2 \right\}, \quad (46)$$

similar notations can be defined for  $\Phi_S$ , which we do not repeat here.

**Definition D.1 (Uniform Stability)** *We say a randomized algorithm  $\mathcal{A}$  is  $\delta$ -uniformly stable in  $x$ -gradients if for every two dataset  $S, S'$  which differ in only one sample, for every  $\xi \in \Xi$  we have*

$$\sup_{\xi} \mathbb{E}_{\mathcal{A}} \|\nabla_x f(\mathcal{A}_x(S), \mathcal{A}_y(S); \xi) - \nabla_x f(\mathcal{A}_x(S'), \mathcal{A}_y(S'); \xi)\|^2 \leq \delta^2. \quad (47)$$

**Lemma D.1 (Concentration of Optimizers)** *For  $y^*$  and  $y_S^*$  defined above, with Assumption 2.1, we have for any  $x \in \mathcal{X}$ ,*

$$\|y^*(x) - y_S^*(x)\| \leq \frac{1}{\mu} \|\nabla_y F_S(x, y^*(x)) - \nabla_y F(x, y^*(x))\|. \quad (48)$$

**Proof** By the optimality of  $y^*(x)$  and  $y_S^*(x)$ , we have for any  $y \in \mathcal{Y}$

$$\begin{aligned} \langle y - y^*(x), \nabla_y F(x, y^*(x)) \rangle &\leq 0 \\ \langle y - y_S^*(x), \nabla_y F_S(x, y_S^*(x)) \rangle &\leq 0. \end{aligned} \quad (49)$$

Setting  $y = y_S^*(x)$  and  $y = y^*(x)$  in the above inequalities respectively, we have

$$\langle y_S^*(x) - y^*(x), \nabla_y F(x, y^*(x)) - \nabla_y F_S(x, y_S^*(x)) \rangle \leq 0. \quad (50)$$

In addition, by strong concavity of  $F_S(x, \cdot)$ , we have

$$\langle y_S^*(x) - y^*(x), \nabla_y F_S(x, y_S^*(x)) - \nabla_y F_S(x, y^*(x)) \rangle + \mu \|y_S^*(x) - y^*(x)\|^2 \leq 0. \quad (51)$$

Combining (50) and (51), we have

$$\langle y_S^*(x) - y^*(x), \nabla_y F(x, y^*(x)) - \nabla_y F_S(x, y^*(x)) \rangle + \mu \|y_S^*(x) - y^*(x)\|^2 \leq 0. \quad (52)$$

Rearranging terms, it holds that

$$\begin{aligned} \mu \|y_S^*(x) - y^*(x)\|^2 &\leq \langle y_S^*(x) - y^*(x), \nabla_y F_S(x, y^*(x)) - \nabla_y F(x, y^*(x)) \rangle \\ &\leq \|y_S^*(x) - y^*(x)\| \cdot \|\nabla_y F_S(x, y^*(x)) - \nabla_y F(x, y^*(x))\|, \end{aligned} \quad (53)$$

which implies

$$\|y_S^*(x) - y^*(x)\| \leq \frac{1}{\mu} \|\nabla_y F_S(x, y^*(x)) - \nabla_y F(x, y^*(x))\|. \quad (54)$$

It concludes the proof. ■

**Lemma D.2 (Stability of Optimizers)** *For  $y_S^*$  and  $y_{S'}^*$  defined above where  $S$  and  $S'$  are two dataset differing in only one sample ( $\xi_i$  and  $\xi'_i$ ), with Assumption 2.1 while  $\mu > 0$ , we have for any  $x \in \mathcal{X}$ ,*

$$\|y_S^*(x) - y_{S'}^*(x)\| \leq \frac{1}{\mu} \|\nabla_y F_S(x, y_{S'}^*(x)) - \nabla_y F_{S'}(x, y_{S'}^*(x))\| \leq \frac{2G}{n\mu}. \quad (55)$$

**Proof** The proof is similar to that of Lemma D.1. By the optimality of  $y_S^*(x)$  and  $y_{S'}^*(x)$ , we have for any

$y \in \mathcal{Y}$

$$\begin{aligned} \langle y - y_S^*(x), \nabla_y F_S(x, y_S^*(x)) \rangle &\leq 0 \\ \langle y - y_{S'}^*(x), \nabla_y F_{S'}(x, y_{S'}^*(x)) \rangle &\leq 0. \end{aligned} \quad (56)$$

Setting  $y = y_{S'}^*(x)$  and  $y = y_S^*(x)$  in the above inequalities respectively, we have

$$\langle y_S^*(x) - y_{S'}^*(x), \nabla_y F_{S'}(x, y_{S'}^*(x)) - \nabla_y F_S(x, y_S^*(x)) \rangle \leq 0. \quad (57)$$

In addition, by strong concavity of  $F_S(x, \cdot)$ , we have

$$\langle y_S^*(x) - y_{S'}^*(x), \nabla_y F_S(x, y_S^*(x)) - \nabla_y F_S(x, y_{S'}^*(x)) \rangle + \mu \|y_S^*(x) - y_{S'}^*(x)\|^2 \leq 0. \quad (58)$$

Combining (57) and (58), we have

$$\langle y_S^*(x) - y_{S'}^*(x), \nabla_y F_S(x, y_S^*(x)) - \nabla_y F_S(x, y_{S'}^*(x)) \rangle + \mu \|y_S^*(x) - y_{S'}^*(x)\|^2 \leq 0. \quad (59)$$

Rearranging terms, it holds that

$$\begin{aligned} \mu \|y_S^*(x) - y_{S'}^*(x)\|^2 &\leq \langle y_S^*(x) - y_{S'}^*(x), \nabla_y F_S(x, y_S^*(x)) - \nabla_y F_{S'}(x, y_{S'}^*(x)) \rangle \\ &\leq \|y_S^*(x) - y_{S'}^*(x)\| \cdot \|\nabla_y F_S(x, y_S^*(x)) - \nabla_y F_{S'}(x, y_{S'}^*(x))\|, \end{aligned} \quad (60)$$

which implies

$$\begin{aligned} \|y_S^*(x) - y_{S'}^*(x)\| &\leq \frac{1}{\mu} \|\nabla_y F_S(x, y_S^*(x)) - \nabla_y F_{S'}(x, y_{S'}^*(x))\| \\ &= \frac{1}{\mu} \left\| \frac{1}{n} (\nabla_y f(x, y_S^*(x); \xi_i) - \nabla_y f(x, y_{S'}^*(x); \xi'_i)) \right\| \leq \frac{2G}{n\mu}, \end{aligned} \quad (61)$$

which concludes the proof. Here the equality above is due to the variables being the same  $(x, y_S^*(x))$ , while  $S$  and  $S'$  differ in only one sample.  $\blacksquare$

**Theorem D.1 (Stability and Generalization, NC-SC)** *Let  $\mathcal{A}$  be an  $\delta$ -uniformly primal stable algorithm, for any function  $f$  satisfying Assumption 2.1 with  $\mu > 0$ , we have*

$$\mathbb{E}_{\mathcal{A}, S} \|\nabla \Phi(\mathcal{A}_x(S)) - \nabla \Phi_S(\mathcal{A}_x(S))\| \leq (1 + \kappa) \left( 4\delta + \frac{G}{\sqrt{n}} \right). \quad (62)$$

**Proof** Following the definition, we have

$$\begin{aligned} &\nabla \Phi(\mathcal{A}_x(S)) - \nabla \Phi_S(\mathcal{A}_x(S)) \\ &= \nabla_x F(\mathcal{A}_x(S), y^*(\mathcal{A}_x(S))) - \nabla_x F_S(\mathcal{A}_x(S), y_S^*(\mathcal{A}_x(S))) \\ &= \nabla_x F(\mathcal{A}_x(S), y^*(\mathcal{A}_x(S))) - \nabla_x F_S(\mathcal{A}_x(S), y^*(\mathcal{A}_x(S))) + \nabla_x F_S(\mathcal{A}_x(S), y^*(\mathcal{A}_x(S))) - \nabla_x F_S(\mathcal{A}_x(S), y_S^*(\mathcal{A}_x(S))), \end{aligned} \quad (63)$$

so we know that

$$\begin{aligned} &\|\nabla \Phi(\mathcal{A}_x(S)) - \nabla \Phi_S(\mathcal{A}_x(S))\| \\ &\leq \|\nabla_x F(\mathcal{A}_x(S), y^*(\mathcal{A}_x(S))) - \nabla_x F_S(\mathcal{A}_x(S), y^*(\mathcal{A}_x(S)))\| \\ &\quad + \|\nabla_x F_S(\mathcal{A}_x(S), y^*(\mathcal{A}_x(S))) - \nabla_x F_S(\mathcal{A}_x(S), y_S^*(\mathcal{A}_x(S)))\|, \end{aligned} \quad (64)$$

for the first term above, by [Lei \[2022, Theorem 2\]](#) (i.e., regarding  $(\mathcal{A}_x(S), y^*(\mathcal{A}_x(S)))$  as one single variable to

recover their conclusion), we have

$$\mathbb{E}_{\mathcal{A},S} \|\nabla_x F(\mathcal{A}_x(S), y^*(\mathcal{A}_x(S))) - \nabla_x F_S(\mathcal{A}_x(S), y^*(\mathcal{A}_x(S)))\| \leq 4\delta + \sqrt{\frac{\text{Var}(\nabla_x f)}{n}} \leq 4\delta + \frac{G}{\sqrt{n}}, \quad (65)$$

for the second term above, by Lemma D.1, we have

$$\begin{aligned} & \mathbb{E}_{\mathcal{A},S} \|\nabla_x F_S(\mathcal{A}_x(S), y^*(\mathcal{A}_x(S))) - \nabla_x F_S(\mathcal{A}_x(S), y_S^*(\mathcal{A}_x(S)))\| \\ & \leq L \mathbb{E}_{\mathcal{A},S} \|y^*(\mathcal{A}_x(S)) - y_S^*(\mathcal{A}_x(S))\| \\ & \leq \kappa \mathbb{E}_{\mathcal{A},S} \|\nabla_y F_S(\mathcal{A}_x(S), y^*(\mathcal{A}_x(S))) - \nabla_y F(\mathcal{A}_x(S), y^*(\mathcal{A}_x(S)))\| \\ & \leq \kappa \left( 4\delta + \sqrt{\frac{\text{Var}(\nabla_y f)}{n}} \right) \\ & \leq \kappa \left( 4\delta + \frac{G}{\sqrt{n}} \right), \end{aligned} \quad (66)$$

where the third inequality applies the same argument as that in (65). We conclude the proof by combining the two bounds above together.  $\blacksquare$

## E Proof of Theorem 4.2

The proof uses the idea from Lei [2022, Theorem 3] and our proof of Theorem 3.2. Unlike Lei [2022] which considers the minimization case, with  $\Phi(x) \neq \mathbb{E}[\Phi(x; \xi)]$ , we need some modification in the proof. To address the non-uniqueness of  $y^*(x)$  in the NC-C case, similar to the uniform convergence analysis in the NC-C case (Theorem 3.2), we resort to the regularized objective in the proof to characterize corresponding distances.

For convenience, we recall the definition of regularized objective functions here.

$$\begin{aligned} \widehat{\Phi}(x) &= \max_{y \in \mathcal{Y}} F(x, y) - \frac{\nu}{2} \|y\|^2, & \widehat{\Phi}_S(x) &= \max_{y \in \mathcal{Y}} F_S(x, y) - \frac{\nu}{2} \|y\|^2, \\ \widehat{y}^*(x) &= \operatorname{argmax}_{y \in \mathcal{Y}} F(x, y) - \frac{\nu}{2} \|y\|^2, & \widehat{y}_S^*(x) &= \operatorname{argmax}_{y \in \mathcal{Y}} F_S(x, y) - \frac{\nu}{2} \|y\|^2. \end{aligned} \quad (67)$$

In addition, following the notation in Lei [2022], we define

$$\begin{aligned} \widetilde{w}_S &= \mathbf{prox}_{\frac{\widehat{\Phi}}{2L}}(\mathcal{A}_x(S)) = \operatorname{argmin}_{x \in \mathcal{X}} \left\{ \widehat{\Phi}(x) + L \|x - \mathcal{A}_x(S)\|^2 \right\}, \\ w_S &= \mathbf{prox}_{\frac{\widehat{\Phi}_S}{2L}}(\mathcal{A}_x(S)) = \operatorname{argmin}_{x \in \mathcal{X}} \left\{ \widehat{\Phi}_S(x) + L \|x - \mathcal{A}_x(S)\|^2 \right\}. \end{aligned} \quad (68)$$

As discussed in Appendix C, the function  $F(x, y) - \frac{\nu}{2} \|y\|^2$  is  $(L + \nu)$ -smooth, and the function  $\widehat{\Phi}(x)$  is  $(L + \nu)$ -weakly-convex (the same hold for  $F_S(x, y) - \frac{\nu}{2} \|y\|^2$  and  $\widehat{\Phi}_S(x)$ ).

First, we build up a connection between algorithm stability and proximal operators to facilitate the analysis.

**Lemma E.1 (Algorithm Stability and Proximal Operators)** *Let  $\mathcal{A}$  be an algorithm. For any function  $f$  satisfying Assumption 2.1 with  $\mu = 0$ , we have the following inequalities for any two neighboring dataset  $S$  and  $S'$ , we have*

$$\begin{aligned} \|\widetilde{w}_S - \widetilde{w}_{S'}\| &\leq \frac{2L}{L - \nu} \|\mathcal{A}(S) - \mathcal{A}(S')\| \\ \|w_S - w_{S'}\| &\leq \frac{2L}{L - \nu} \|\mathcal{A}(S) - \mathcal{A}(S')\| + \frac{2G}{n(L - \nu)} + \frac{2L(G + \nu\sqrt{D_Y})}{n\nu(L - \nu)}, \end{aligned} \quad (69)$$

where  $\widehat{\Phi}$  and  $\widehat{\Phi}_S$  follows the definitions in (67) and (68).

The proof basically follows the proof of Lei [2022, Lemma 15 and 16] with some differences in detailed parameters.

**Proof** For the first result, note that  $\widehat{\Phi}(x)$  is  $(L + \nu)$ -weakly-convex and differentiable, so we have

$$\left\langle \widetilde{w}_S - \widetilde{w}_{S'}, \nabla \widehat{\Phi}(\widetilde{w}_S) - \nabla \widehat{\Phi}(\widetilde{w}_{S'}) \right\rangle \geq -(L + \nu) \|\widetilde{w}_S - \widetilde{w}_{S'}\|^2. \quad (70)$$

On the other hand, by the optimality of  $\widetilde{w}_S$ , we have

$$-2L(\widetilde{w}_S - \mathcal{A}_x(S)) - \nabla \widehat{\Phi}(\widetilde{w}_S) \in \partial \mathcal{I}_{\mathcal{X}}(\widetilde{w}_S), \quad -2L(\widetilde{w}_{S'} - \mathcal{A}_x(S')) - \nabla \widehat{\Phi}(\widetilde{w}_{S'}) \in \partial \mathcal{I}_{\mathcal{X}}(\widetilde{w}_{S'}), \quad (71)$$

where  $\mathcal{I}_{\mathcal{X}}(x)$  is the indicator function of the set  $\mathcal{X}$ , i.e.,  $\mathcal{I}_{\mathcal{X}}(x) = 0$  if  $x \in \mathcal{X}$  and  $\mathcal{I}_{\mathcal{X}}(x) = \infty$  otherwise. Since  $\mathcal{X}$  is convex, the subgradient  $\partial \mathcal{I}_{\mathcal{X}}$  is monotone, and thus

$$\begin{aligned} \left\langle \widetilde{w}_S - \widetilde{w}_{S'}, 2L(\widetilde{w}_{S'} - \mathcal{A}_x(S')) - 2L(\widetilde{w}_S - \mathcal{A}_x(S)) + \nabla \widehat{\Phi}(\widetilde{w}_{S'}) - \nabla \widehat{\Phi}(\widetilde{w}_S) \right\rangle &= \langle \widetilde{w}_S - \widetilde{w}_{S'}, \partial \mathcal{I}_{\mathcal{X}}(\widetilde{w}_S) - \partial \mathcal{I}_{\mathcal{X}}(\widetilde{w}_{S'}) \rangle \\ &\geq 0. \end{aligned} \quad (72)$$

Combining (70) and (72), it follows that

$$\langle \widetilde{w}_S - \widetilde{w}_{S'}, 2L(\widetilde{w}_{S'} - \mathcal{A}_x(S')) - 2L(\widetilde{w}_S - \mathcal{A}_x(S)) \rangle \geq -(L + \nu) \|\widetilde{w}_S - \widetilde{w}_{S'}\|^2. \quad (73)$$

Rearranging the terms, we have

$$(L - \nu) \|\widetilde{w}_S - \widetilde{w}_{S'}\|^2 \leq 2L \langle \widetilde{w}_S - \widetilde{w}_{S'}, \mathcal{A}_x(S) - \mathcal{A}_x(S') \rangle \leq 2L \|\widetilde{w}_S - \widetilde{w}_{S'}\| \|\mathcal{A}_x(S) - \mathcal{A}_x(S')\|. \quad (74)$$

We obtain the first result by dividing both sides by  $(L - \nu) \|\widetilde{w}_S - \widetilde{w}_{S'}\|$ .

For the second statement, applying the fact that  $\widehat{\Phi}_S$  is weakly-convex and differentiable,

$$\left\langle w_S - w_{S'}, \nabla \widehat{\Phi}_S(w_S) - \nabla \widehat{\Phi}_S(w_{S'}) \right\rangle \geq -(L + \nu) \|w_S - w_{S'}\|^2. \quad (75)$$

Similar as (72), by the optimality condition of  $w_S$  and  $w_{S'}$ ,

$$\left\langle w_S - w_{S'}, 2L(w_{S'} - \mathcal{A}_x(S')) - 2L(w_S - \mathcal{A}_x(S)) + \nabla \widehat{\Phi}_{S'}(w_{S'}) - \nabla \widehat{\Phi}_S(w_S) \right\rangle \geq 0. \quad (76)$$

Therefore, by the above two equations, we obtain that

$$-(L + \nu) \|w_S - w_{S'}\|^2 \leq \left\langle w_S - w_{S'}, 2L(w_{S'} - \mathcal{A}_x(S')) - 2L(w_S - \mathcal{A}_x(S)) + \nabla \widehat{\Phi}_{S'}(w_{S'}) - \nabla \widehat{\Phi}_S(w_S) \right\rangle. \quad (77)$$

By the definition of  $\widehat{\Phi}_S$  and  $w_S$ , we rewrite the additional term  $\nabla \widehat{\Phi}_{S'}(w_{S'}) - \nabla \widehat{\Phi}_S(w_S)$  as

$$\begin{aligned} &\nabla \widehat{\Phi}_{S'}(w_{S'}) - \nabla \widehat{\Phi}_S(w_S) \\ &= \nabla_x \left( F_{S'}(w_{S'}; \widehat{y}_{S'}^*(w_{S'})) - \frac{\nu}{2} \|\widehat{y}_{S'}^*(w_{S'})\|^2 \right) - \nabla \widehat{\Phi}_S(w_S) \\ &= \nabla_x F_{S'}(w_{S'}; \widehat{y}_{S'}^*(w_{S'})) - \nabla \widehat{\Phi}_S(w_S) \\ &= \nabla_x F_{S'}(w_{S'}; \widehat{y}_{S'}^*(w_{S'})) - \nabla_x F_S(w_{S'}; \widehat{y}_{S'}^*(w_{S'})) + \nabla_x F_S(w_{S'}; \widehat{y}_{S'}^*(w_{S'})) - \nabla_x F_S(w_S; \widehat{y}_S^*(w_S)) \\ &= \underbrace{\frac{1}{n} \nabla_x f(w_{S'}; \widehat{y}_{S'}^*(w_{S'}); \xi_i') - \frac{1}{n} \nabla_x f(w_S; \widehat{y}_S^*(w_S); \xi_i)}_{E_1} + \underbrace{\nabla_x F_S(w_{S'}; \widehat{y}_{S'}^*(w_{S'})) - \nabla_x F_S(w_S; \widehat{y}_S^*(w_S))}_{E_2}, \end{aligned} \quad (78)$$

where the third equation holds since  $\nabla \widehat{\Phi}_S(w_{S'}) = \nabla_x F_S(w_{S'}; \widehat{y}_S^*(w_{S'}))$ . Thus it holds that

$$-(L + \nu)\|w_S - w_{S'}\|^2 \leq \langle w_S - w_{S'}, -2L(w_S - \mathcal{A}_x(S)) + 2L(w_{S'} - \mathcal{A}_x(S')) + E_1 + E_2 \rangle. \quad (79)$$

Rearranging terms, we have

$$\begin{aligned} & (L - \nu)\|w_S - w_{S'}\|^2 \\ & \leq \langle w_S - w_{S'}, 2L(\mathcal{A}_x(S) - \mathcal{A}_x(S')) + E_1 + E_2 \rangle \\ & \leq \|w_S - w_{S'}\| \|2L(\mathcal{A}_x(S) - \mathcal{A}_x(S')) + E_1 + E_2\| \\ & \leq \|w_S - w_{S'}\| (2L\|\mathcal{A}_x(S) - \mathcal{A}_x(S')\| + \|E_1\| + \|E_2\|) \\ & \leq \|w_S - w_{S'}\| \left( 2L\|\mathcal{A}_x(S) - \mathcal{A}_x(S')\| + \frac{2G}{n} + \frac{2L(G + \nu\sqrt{Dy})}{n\nu} \right), \end{aligned} \quad (80)$$

where the last inequality uses the fact that that  $\|E_1\| \leq 2G/n$  via Lipschitz continuity, and

$$\|E_2\| \leq L\|\widehat{y}_{S'}^*(w_{S'}) - \widehat{y}_S^*(w_{S'})\| \stackrel{\text{Lemma D.2}}{\leq} \frac{2L(G + \nu\sqrt{Dy})}{n\nu}.$$

It concludes the proof by diving  $(L - \nu)\|w_S - w_{S'}\|$  on both sides of (80).  $\blacksquare$

**Lemma E.2** *Let  $\mathcal{A}$  be an  $\delta$ -uniformly primal argument stable algorithm. For any function  $f$  satisfying Assumption 2.1 with  $\mu = 0$ , we have*

$$\mathbb{E} \left[ \widehat{\Phi}_S(\tilde{w}_S) - \widehat{\Phi}(\tilde{w}_S) \right] \leq \frac{2GL(L + 2\nu)}{\nu(L - \nu)}\delta + \frac{G}{\nu} \left( 4\sqrt{\frac{8L^4(L + 2\nu)^2}{\nu^2(L - \nu)^2}}\delta + \frac{G}{\sqrt{n}} \right) + \frac{\nu}{2}Dy. \quad (81)$$

**Proof** Note that

$$\begin{aligned} & \mathbb{E} \left[ \widehat{\Phi}_S(\tilde{w}_S) - \widehat{\Phi}(\tilde{w}_S) \right] \\ & = \mathbb{E} \left[ F_S(\tilde{w}_S, \widehat{y}_S^*(\tilde{w}_S)) - F(\tilde{w}_S, \widehat{y}^*(\tilde{w}_S)) - \frac{\nu}{2}\|\widehat{y}_S^*(\tilde{w}_S)\|^2 + \frac{\nu}{2}\|\widehat{y}^*(\tilde{w}_S)\|^2 \right] \\ & \leq \mathbb{E} \left[ \underbrace{F_S(\tilde{w}_S, \widehat{y}_S^*(\tilde{w}_S)) - F_S(\tilde{w}_S, \widehat{y}^*(\tilde{w}_S))}_{H_1} + \underbrace{F_S(\tilde{w}_S, \widehat{y}^*(\tilde{w}_S)) - F(\tilde{w}_S, \widehat{y}^*(\tilde{w}_S))}_{H_2} \right] + \frac{\nu}{2}Dy. \end{aligned} \quad (82)$$

We bound  $H_2$  via the stability argument of  $f(\tilde{w}_S, \widehat{y}^*(\tilde{w}_S); \xi)$ , i.e., regarding  $(\tilde{w}_S, \widehat{y}_S^*(\tilde{w}_S))$  as one single variable.

$$\begin{aligned} & \mathbb{E} [f(\tilde{w}_S, \widehat{y}^*(\tilde{w}_S); \xi)] - \mathbb{E} [f(\tilde{w}_{S'}, \widehat{y}^*(\tilde{w}_{S'}); \xi)] \\ & \leq G \mathbb{E} [\|\tilde{w}_S - \tilde{w}_{S'}\| + \|\widehat{y}^*(\tilde{w}_S) - \widehat{y}^*(\tilde{w}_{S'})\|] \\ & \leq G \mathbb{E} \left[ \|\tilde{w}_S - \tilde{w}_{S'}\| + \frac{L + \nu}{\nu}\|\tilde{w}_S - \tilde{w}_{S'}\| \right] \\ & \leq G \mathbb{E} \left[ \left( 1 + \frac{L + \nu}{\nu} \right) \cdot \frac{2L}{L - \nu}\|\mathcal{A}_x(S) - \mathcal{A}_x(S')\| \right] \\ & \leq \frac{L + 2\nu}{\nu} \cdot \frac{2GL}{L - \nu} \mathbb{E} [\|\mathcal{A}_x(S) - \mathcal{A}_x(S')\|] \\ & \leq \frac{2GL(L + 2\nu)}{\nu(L - \nu)}\delta, \end{aligned} \quad (83)$$

where the second inequality uses Lin et al. [2020a, Lemma 4.3], and the fact that  $\widehat{y}^*$  is the optimal solution of a  $(L + \nu)$ -smooth and  $\nu$ -strongly concave maximization problem defined in (67); the third inequality is due to Lemma E.1, and the last inequality follows the definition of  $\delta$ -uniform primal argument stability. So we have

the ‘‘composed algorithm’’  $\tilde{w}_S$  is stable<sup>4</sup> in function values, which implies [Hardt et al., 2016]

$$\mathbb{E} [F_S(\tilde{w}_S, \hat{y}^*(\tilde{w}_S)) - F(\tilde{w}_S, \hat{y}^*(\tilde{w}_S))] \leq \frac{2GL(L+2\nu)}{\nu(L-\nu)} \delta. \quad (84)$$

For the term  $H_1$  above, we have

$$\begin{aligned} & \mathbb{E} [F_S(\tilde{w}_S, \hat{y}_S^*(\tilde{w}_S)) - F_S(\tilde{w}_S, \hat{y}^*(\tilde{w}_S))] \\ & \leq G \mathbb{E} \|\hat{y}_S^*(\tilde{w}_S) - \hat{y}^*(\tilde{w}_S)\| \\ & \leq \frac{G}{\nu} \mathbb{E} \|\nabla_y F_S(\tilde{w}_S, \hat{y}^*(\tilde{w}_S)) - \nu \hat{y}^*(\tilde{w}_S) - \nabla_y F(\tilde{w}_S, \hat{y}^*(\tilde{w}_S)) + \nu \hat{y}^*(\tilde{w}_S)\| \\ & = \frac{G}{\nu} \mathbb{E} \|\nabla_y F_S(\tilde{w}_S, \hat{y}^*(\tilde{w}_S)) - \nabla_y F(\tilde{w}_S, \hat{y}^*(\tilde{w}_S))\|, \end{aligned} \quad (85)$$

where the second inequality applies Lemma D.1. We further upper bound the RHS above using the stability argument. For  $\nabla_y f(\tilde{w}_S, \hat{y}^*(\tilde{w}_S); \xi)$ , similar to the same argument as in (83), we have

$$\begin{aligned} & \mathbb{E} \|\nabla_y f(\tilde{w}_S, \hat{y}^*(\tilde{w}_S); \xi) - \nabla_y f(\tilde{w}_{S'}, \hat{y}^*(\tilde{w}_{S'}); \xi)\|^2 \\ & \leq 2L^2 \mathbb{E} \left[ \|\tilde{w}_S - \tilde{w}_{S'}\|^2 + \|\hat{y}^*(\tilde{w}_S) - \hat{y}^*(\tilde{w}_{S'})\|^2 \right] \\ & \leq 2L^2 \mathbb{E} \left[ \left( 1 + \left( \frac{L+\nu}{\nu} \right)^2 \right) \|\tilde{w}_S - \tilde{w}_{S'}\|^2 \right] \\ & \leq 2L^2 \left( 1 + \left( \frac{L+\nu}{\nu} \right)^2 \right) \cdot \left( \frac{2L}{L-\nu} \right)^2 \mathbb{E} \left[ \|\mathcal{A}_x(S) - \mathcal{A}_x(S')\|^2 \right] \\ & \leq \frac{8L^4(L+2\nu)^2}{\nu^2(L-\nu)^2} \delta^2, \end{aligned} \quad (86)$$

where the second inequality comes from Lin et al. [2020a, Lemma 4.3]. It concludes that algorithm  $\mathcal{A}$  is  $\delta$ -uniformly primal stable. Applying Lei [2022, Theorem 2] to (85), we have

$$\begin{aligned} \mathbb{E} [F_S(\tilde{w}_S, \hat{y}_S^*(\tilde{w}_S)) - F_S(\tilde{w}_S, \hat{y}^*(\tilde{w}_S))] & \leq \frac{G}{\nu} \mathbb{E} \|\nabla_y F_S(\tilde{w}_S, \hat{y}^*(\tilde{w}_S)) - \nabla_y F(\tilde{w}_S, \hat{y}^*(\tilde{w}_S))\| \\ & \leq \frac{G}{\nu} \left( 4\sqrt{\frac{8L^4(L+2\nu)^2}{\nu^2(L-\nu)^2}} \delta + \sqrt{\frac{\text{Var}(\nabla_y f)}{n}} \right) \\ & \leq \frac{G}{\nu} \left( 4\sqrt{\frac{8L^4(L+2\nu)^2}{\nu^2(L-\nu)^2}} \delta + \frac{G}{\sqrt{n}} \right), \end{aligned} \quad (87)$$

which concludes the proof.  $\blacksquare$

**Lemma E.3** *Let  $\mathcal{A}$  be an  $\delta$ -uniformly primal argument stable algorithm. For any function  $f$  satisfying Assumption 2.1 with  $\mu = 0$ , we have*

$$\begin{aligned} \mathbb{E} \left[ \hat{\Phi}(w_S) - \hat{\Phi}_S(w_S) \right] & \leq \frac{G}{\nu} \left( 4\sqrt{\frac{8L^2(L+2\nu)^2}{\nu^2} \left( \frac{4L^2}{(L-\nu)^2} \delta^2 + \frac{4G^2}{n^2(L-\nu)^2} + \frac{2L^2(G+\nu\sqrt{Dy})^2}{n^2\nu^2(L-\nu)^2} \right)} + \frac{G}{\sqrt{n}} \right) \\ & \quad + \frac{G(L+2\nu)}{\nu} \left( \frac{2L}{L-\nu} \delta + \frac{2G}{n(L-\nu)} + \frac{2L(G+\nu\sqrt{Dy})}{n\nu(L-\nu)} \right) + \frac{2G(G+\nu\sqrt{Dy})}{n\nu} + \frac{\nu}{2} Dy. \end{aligned} \quad (88)$$

<sup>4</sup> Here we call the iteration  $\tilde{w}_S = \mathbf{prox}_{\frac{\hat{\Phi}}{2L}}(\mathcal{A}_x(S)) = \arg\min_{x \in \mathcal{X}} \left\{ \hat{\Phi}(x) + L\|x - \mathcal{A}_x(S)\|^2 \right\}$  as an algorithm regarding that it is a composition of the algorithm  $\mathcal{A}$  and the proximal operator.

**Proof** Note that

$$\begin{aligned}
& \mathbb{E} \left[ \widehat{\Phi}(w_S) - \widehat{\Phi}_S(w_S) \right] \\
&= \mathbb{E} \left[ F(w_S, \widehat{y}^*(w_S)) - F_S(w_S, \widehat{y}_S^*(w_S)) - \frac{\nu}{2} \|\widehat{y}^*(w_S)\|^2 + \frac{\nu}{2} \|\widehat{y}_S^*(w_S)\|^2 \right] \\
&\leq \mathbb{E} \left[ \underbrace{F(w_S, \widehat{y}^*(w_S)) - F(w_S, \widehat{y}_S^*(w_S))}_{J_1} + \underbrace{F(w_S, \widehat{y}_S^*(w_S)) - F_S(w_S, \widehat{y}_S^*(w_S))}_{J_2} \right] + \frac{\nu}{2} D_{\mathbf{y}}.
\end{aligned} \tag{89}$$

For  $J_2$ , by Lemma E.1, similar to the analysis of  $H_2$  in the proof of Lemma E.2, we have

$$\begin{aligned}
& \mathbb{E} [f(w_S, \widehat{y}_S^*(w_S); \xi) - \mathbb{E} [f(w_{S'}, \widehat{y}_{S'}^*(w_{S'}); \xi)]] \\
&\leq G \mathbb{E} [\|w_S - w_{S'}\| + \|\widehat{y}_S^*(w_S) - \widehat{y}_{S'}^*(w_{S'})\| + \|\widehat{y}_{S'}^*(w_S) - \widehat{y}_{S'}^*(w_{S'})\|] \\
&\leq G \mathbb{E} \left[ \|w_S - w_{S'}\| + \frac{L + \nu}{\nu} \|w_S - w_{S'}\| \right] + \frac{2G(G + \nu\sqrt{D_{\mathbf{y}}})}{n\nu} \\
&\leq G \mathbb{E} \left[ \left(1 + \frac{L + \nu}{\nu}\right) \cdot \left( \frac{2L}{L - \nu} \|\mathcal{A}_x(S) - \mathcal{A}_x(S')\| + \frac{2G}{n(L - \nu)} + \frac{2L(G + \nu D_{\mathbf{y}})}{n\nu(L - \nu)} \right) \right] + \frac{2G(G + \nu\sqrt{D_{\mathbf{y}}})}{n\nu} \\
&\leq \frac{G(L + 2\nu)}{\nu} \cdot \left( \frac{2L}{L - \nu} \mathbb{E} [\|\mathcal{A}_x(S) - \mathcal{A}_x(S')\|] + \frac{2G}{n(L - \nu)} + \frac{2L(G + \nu D_{\mathbf{y}})}{n\nu(L - \nu)} \right) + \frac{2G(G + \nu\sqrt{D_{\mathbf{y}}})}{n\nu} \\
&\leq \frac{G(L + 2\nu)}{\nu} \left( \frac{2L}{L - \nu} \delta + \frac{2G}{n(L - \nu)} + \frac{2L(G + \nu D_{\mathbf{y}})}{n\nu(L - \nu)} \right) + \frac{2G(G + \nu\sqrt{D_{\mathbf{y}}})}{n\nu}.
\end{aligned} \tag{90}$$

It further holds that

$$\mathbb{E} [F(w_S, \widehat{y}_S^*(w_S)) - F_S(w_S, \widehat{y}_S^*(w_S))] \leq \frac{G(L + 2\nu)}{\nu} \left( \frac{2L}{L - \nu} \delta + \frac{2G}{n(L - \nu)} + \frac{2L(G + \nu\sqrt{D_{\mathbf{y}}})}{n\nu(L - \nu)} \right) + \frac{2G(G + \nu\sqrt{D_{\mathbf{y}}})}{n\nu}. \tag{91}$$

For  $J_1$ , similar to the analysis of  $H_1$  in the proof of Lemma E.2, we have

$$\begin{aligned}
& \mathbb{E} \|\nabla_{\mathbf{y}} f(w_S, \widehat{y}^*(w_S); \xi) - \nabla_{\mathbf{y}} f(w_{S'}, \widehat{y}^*(w_{S'}); \xi)\|^2 \\
&\leq 2L^2 \mathbb{E} \left[ \|w_S - w_{S'}\|^2 + \|\widehat{y}^*(w_S) - \widehat{y}^*(w_{S'})\|^2 \right] \\
&\leq 2L^2 \mathbb{E} \left[ \left(1 + \left(\frac{L + \nu}{\nu}\right)^2\right) \|w_S - w_{S'}\|^2 \right] \\
&\leq 2L^2 \left(1 + \left(\frac{L + \nu}{\nu}\right)^2\right) \cdot \mathbb{E} \left[ 4 \left(\frac{2L}{L - \nu}\right)^2 \|\mathcal{A}_x(S) - \mathcal{A}_x(S')\|^2 + 4 \frac{4G^2}{n^2(L - \nu)^2} + 2 \frac{4L^2(G + \nu\sqrt{D_{\mathbf{y}}})^2}{n^2\nu^2(L - \nu)^2} \right] \\
&\leq \frac{8L^2(L + 2\nu)^2}{\nu^2} \left( \frac{4L^2}{(L - \nu)^2} \delta^2 + \frac{4G^2}{n^2(L - \nu)^2} + \frac{2L^2(G + \nu\sqrt{D_{\mathbf{y}}})^2}{n^2\nu^2(L - \nu)^2} \right).
\end{aligned} \tag{92}$$

Combined with Lei [2022, Theorem 2], we have

$$\begin{aligned}
& \mathbb{E} [F(w_S, \widehat{y}_S^*(w_S)) - F(w_S, \widehat{y}^*(w_S))] \\
&\leq \frac{G}{\nu} \mathbb{E} \|\nabla_{\mathbf{y}} F_S(w_S, \widehat{y}^*(w_S)) - \nabla_{\mathbf{y}} F(w_S, \widehat{y}^*(w_S))\| \\
&\leq \frac{G}{\nu} \left( 4 \sqrt{\frac{8L^2(L + 2\nu)^2}{\nu^2} \left( \frac{4L^2}{(L - \nu)^2} \delta^2 + \frac{4G^2}{n^2(L - \nu)^2} + \frac{2L^2(G + \nu\sqrt{D_{\mathbf{y}}})^2}{n^2\nu^2(L - \nu)^2} \right)} + \sqrt{\frac{\text{Var}(\nabla_{\mathbf{y}} f)}{n}} \right) \\
&\leq \frac{G}{\nu} \left( 4 \sqrt{\frac{8L^2(L + 2\nu)^2}{\nu^2} \left( \frac{4L^2}{(L - \nu)^2} \delta^2 + \frac{4G^2}{n^2(L - \nu)^2} + \frac{2L^2(G + \nu\sqrt{D_{\mathbf{y}}})^2}{n^2\nu^2(L - \nu)^2} \right)} + \frac{G}{\sqrt{n}} \right),
\end{aligned} \tag{93}$$

which concludes the proof.  $\blacksquare$

Next, we formally demonstrate the proof for the generalization bounds in the NC-C setting.

**Theorem E.1 (Stability and Generalization, NC-C, repeat Theorem 4.2)** *Let  $\mathcal{A}$  be an  $\delta$ -uniformly primal argument stable algorithm, for any function  $f$  satisfying Assumption 2.1 with  $\mu = 0$ , we have*

$$\mathbb{E}_{\mathcal{A}, S} \left\| \nabla \Phi^{1/(2L)}(\mathcal{A}_x(S)) - \nabla \Phi_S^{1/(2L)}(\mathcal{A}_x(S)) \right\| \leq \mathcal{O} \left( \delta^{\frac{1}{6}} + \left( \frac{1}{n} \right)^{\frac{1}{12}} \right). \quad (94)$$

**Proof** Recall that

$$\nabla \Phi^{1/(2L)}(\mathcal{A}_x(S)) = 2L \left( \mathcal{A}_x(S) - \mathbf{prox}_{\frac{\Phi}{2L}}(\mathcal{A}(S)) \right), \quad \nabla \Phi_S^{1/(2L)}(\mathcal{A}_x(S)) = 2L \left( \mathcal{A}_x(S) - \mathbf{prox}_{\frac{\Phi_S}{2L}}(\mathcal{A}(S)) \right). \quad (95)$$

Since  $\Phi$  is  $L$ -weakly-convex and  $G$ -Lipschitz [Lin et al., 2020a, Lemma 4.7], it holds that

$$\left\| \nabla \Phi^{1/(2L)}(\mathcal{A}_x(S)) - \nabla \Phi_S^{1/(2L)}(\mathcal{A}_x(S)) \right\| = 2L \left\| \mathbf{prox}_{\frac{\Phi}{2L}}(\mathcal{A}(S)) - \mathbf{prox}_{\frac{\Phi_S}{2L}}(\mathcal{A}(S)) \right\|. \quad (96)$$

Utilizing the regularized objective function, we have

$$\begin{aligned} & \left\| \mathbf{prox}_{\frac{\Phi}{2L}}(\mathcal{A}(S)) - \mathbf{prox}_{\frac{\Phi_S}{2L}}(\mathcal{A}(S)) \right\| \\ \leq & \left\| \mathbf{prox}_{\frac{\Phi}{2L}}(\mathcal{A}(S)) - \mathbf{prox}_{\frac{\hat{\Phi}}{2L}}(\mathcal{A}(S)) \right\| + \left\| \mathbf{prox}_{\frac{\hat{\Phi}}{2L}}(\mathcal{A}(S)) - \mathbf{prox}_{\frac{\hat{\Phi}_S}{2L}}(\mathcal{A}(S)) \right\| + \left\| \mathbf{prox}_{\frac{\hat{\Phi}_S}{2L}}(\mathcal{A}(S)) - \mathbf{prox}_{\frac{\Phi_S}{2L}}(\mathcal{A}(S)) \right\| \\ \leq & 2\sqrt{\frac{\nu D_{\mathbf{y}}}{L - \nu}} + \|\tilde{w}_S - w_S\|, \end{aligned} \quad (97)$$

where the second inequality comes from Lemma 3.1 with  $\lambda = \frac{1}{2L}$ . So now the problem is transformed to characterizing the distance between  $\tilde{w}_S$  and  $w_S$  coming from the regularized surrogate objective which is NC-SC.

Since the function  $\hat{\Phi}(x) + L\|x - \mathcal{A}(S)\|^2$  is  $(L - \nu)$ -strongly convex, and by the definition of  $\tilde{w}_S$ , we have

$$\begin{aligned} & \frac{L - \nu}{2} \mathbb{E} \|w_S - \tilde{w}_S\|^2 \\ \leq & \mathbb{E} \hat{\Phi}(w_S) + L\|w_S - \mathcal{A}(S)\|^2 - \left( \hat{\Phi}(\tilde{w}_S) + L\|\tilde{w}_S - \mathcal{A}(S)\|^2 \right) \\ = & \mathbb{E} \hat{\Phi}_S(w_S) + L\|w_S - \mathcal{A}(S)\|^2 - \left( \hat{\Phi}_S(\tilde{w}_S) + L\|\tilde{w}_S - \mathcal{A}(S)\|^2 \right) + \left( \hat{\Phi}(w_S) - \hat{\Phi}_S(w_S) \right) + \left( \hat{\Phi}_S(\tilde{w}_S) - \hat{\Phi}(\tilde{w}_S) \right) \\ \leq & \mathbb{E} \left( \hat{\Phi}(w_S) - \hat{\Phi}_S(w_S) \right) + \left( \hat{\Phi}_S(\tilde{w}_S) - \hat{\Phi}(\tilde{w}_S) \right) \\ \leq & \frac{G}{\nu} \left( 4\sqrt{\frac{8L^2(L + 2\nu)^2}{\nu^2} \left( \frac{4L^2}{(L - \nu)^2} \delta^2 + \frac{4G^2}{n^2(L - \nu)^2} + \frac{2L^2(G + \nu\sqrt{D_{\mathbf{y}}})^2}{n^2\nu^2(L - \nu)^2} \right)} + \frac{G}{\sqrt{n}} \right) \\ & + \frac{G(L + 2\nu)}{\nu} \left( \frac{2L}{L - \nu} \delta + \frac{2G}{n(L - \nu)} + \frac{2L(G + \nu\sqrt{D_{\mathbf{y}}})}{n\nu(L - \nu)} \right) + \frac{2G(G + \nu\sqrt{D_{\mathbf{y}}})}{n\nu} \\ & + \frac{2GL(L + 2\nu)}{\nu(L - \nu)} \delta + \frac{G}{\nu} \left( 4\sqrt{\frac{8L^4(L + 2\nu)^2}{\nu^2(L - \nu)^2}} \delta + \frac{G}{\sqrt{n}} \right) + \nu D_{\mathbf{y}}, \end{aligned} \quad (98)$$

where the second inequality uses the optimality of  $w_S$  and  $\hat{w}_S$ , the last inequality is due to Lemma E.2 and E.3. Now we choose  $\nu$  to simplify the RHS above. For simplicity, first we set  $\nu \leq \frac{L}{2}$ , so  $L - \nu \geq \frac{L}{2}$ ,  $L + 2\nu \leq 2L$ .

The RHS above simplifies to

$$\begin{aligned}
& \frac{L-\nu}{2} \mathbb{E} \|w_S - \tilde{w}_S\|^2 \\
& \leq \frac{G}{\nu} \left( 4\sqrt{\frac{32L^4}{\nu^2} \left( 16\delta^2 + \frac{16G^2}{n^2L^2} + \frac{16(G+\nu\sqrt{Dy})^2}{n^2\nu^2} \right)} + \frac{G}{\sqrt{n}} \right) + \frac{2GL}{\nu} \left( 4\delta + \frac{4G}{nL} + \frac{4(G+\nu\sqrt{Dy})}{n\nu} \right) \\
& \quad + \frac{2G(G+\nu\sqrt{Dy})}{n\nu} + \frac{8GL}{\nu} \delta + \frac{G}{\nu} \left( 4\sqrt{\frac{128L^4}{\nu^2}} \delta + \frac{G}{\sqrt{n}} \right) + \nu Dy \\
& \leq \frac{G}{\nu} \left( \frac{128L^2}{\nu} \sqrt{\delta^2 + \frac{G^2}{n^2L^2} + \frac{(G+\nu\sqrt{Dy})^2}{n^2\nu^2}} + \frac{G}{\sqrt{n}} \right) + \frac{8GL}{\nu} \left( 2\delta + \frac{G}{nL} + \frac{G+\nu\sqrt{Dy}}{n\nu} \right) \\
& \quad + \frac{G}{\nu} \left( \frac{64L^2}{\nu} \delta + \frac{2(G+\nu\sqrt{Dy})}{n} + \frac{G}{\sqrt{n}} \right) + \nu Dy \tag{99} \\
& \leq \frac{G}{\nu} \left( \frac{128L^2}{\nu} \left( \delta + \frac{G}{nL} + \frac{G+\nu\sqrt{Dy}}{n\nu} \right) + \frac{G}{\sqrt{n}} \right) + \frac{8GL}{\nu} \left( 2\delta + \frac{G}{nL} + \frac{G+\nu\sqrt{Dy}}{n\nu} \right) \\
& \quad + \frac{G}{\nu} \left( \frac{64L^2}{\nu} \delta + \frac{2(G+\nu\sqrt{Dy})}{n} + \frac{G}{\sqrt{n}} \right) + \nu Dy \\
& = 64G \left( 3\delta + \frac{2G}{nL} + \frac{2G+2\nu\sqrt{Dy}}{n\nu} \right) \frac{L^2}{\nu^2} + 2G \left( 8\delta + \frac{G}{\sqrt{n}L} + \frac{4G}{nL} + \frac{G+\nu\sqrt{Dy}}{n} \left( \frac{4}{\nu} + \frac{1}{L} \right) \right) \frac{L}{\nu} + \nu Dy \\
& = \mathcal{O}\left(\frac{1}{n}\right) \cdot \mathcal{O}\left(\frac{1}{\nu^3} + \frac{1}{\nu^2} + \frac{1}{\nu} + 1\right) + \mathcal{O}\left(\frac{1}{\sqrt{n}}\right) \cdot \mathcal{O}\left(\frac{1}{\nu}\right) + \mathcal{O}(\delta) \cdot \mathcal{O}\left(\frac{1}{\nu^2} + \frac{1}{\nu}\right) + \mathcal{O}(1)\nu,
\end{aligned}$$

where the last step hides all other dependence on parameters except  $\delta$  and  $n$ . Let

$$\frac{1}{\nu} = \mathcal{O}\left(\min\left(\delta^{-\frac{1}{3}}, n^{\frac{1}{4}}\right)\right), \tag{100}$$

with  $\delta$  and  $1/n$  small enough such that  $\nu \leq L/2$  holds. The setting of  $\nu$  implies that

$$\begin{aligned}
\mathbb{E} \|w_S - \tilde{w}_S\|^2 & \leq \mathcal{O}\left(\frac{1}{n}\right) \cdot \mathcal{O}\left(\frac{1}{\nu^3}\right) + \mathcal{O}\left(\frac{1}{\sqrt{n}}\right) \cdot \mathcal{O}\left(\frac{1}{\nu}\right) + \mathcal{O}(\delta) \cdot \mathcal{O}\left(\frac{1}{\nu^2}\right) + \mathcal{O}(1)\nu \\
& \leq \mathcal{O}\left(\frac{1}{n}\right) \cdot \mathcal{O}\left(n^{\frac{3}{4}}\right) + \mathcal{O}\left(\frac{1}{\sqrt{n}}\right) \cdot \mathcal{O}\left(n^{\frac{1}{4}}\right) + \mathcal{O}(\delta) \cdot \mathcal{O}\left(\delta^{-\frac{2}{3}}\right) + \mathcal{O}\left(\delta^{\frac{1}{3}} + n^{-\frac{1}{4}}\right) \\
& = \mathcal{O}\left(\delta^{\frac{1}{3}} + n^{-\frac{1}{4}}\right).
\end{aligned} \tag{101}$$

As a result, we have

$$\mathbb{E} \|w_S - \tilde{w}_S\| \leq \mathcal{O}\left(\delta^{\frac{1}{6}} + n^{-\frac{1}{8}}\right). \tag{102}$$

Further incorporating (97), we have

$$\begin{aligned}
& \mathbb{E}_{\mathcal{A}, S} \left\| \nabla \Phi^{1/(2L)}(\mathcal{A}_x(S)) - \nabla \Phi_S^{1/(2L)}(\mathcal{A}_x(S)) \right\| \\
& \leq 2\sqrt{\frac{\nu Dy}{L-\nu}} + \mathbb{E} \|\tilde{w}_S - w_S\| \\
& \leq \mathcal{O}(\sqrt{\nu}) + \mathbb{E} \|\tilde{w}_S - w_S\| \\
& = \mathcal{O}\left(\delta^{\frac{1}{6}} + n^{-\frac{1}{8}}\right),
\end{aligned} \tag{103}$$

which concludes the proof. ■

## F Proof of Corollary 4.1 and 4.2

**Corollary F.1** Assume the function  $f$  is NC-SC as defined in Assumption 2.1, then if we run SGDA for  $T$  iterations with stepsize  $(\alpha_x, \alpha_y) = \left(\frac{c}{t}, \frac{cr^2}{t}\right)$  for some constant  $c > 0$  and  $1 \leq r < \kappa$ , we have

$$\mathbb{E}_{S, \mathcal{A}} \|\nabla\Phi(\mathcal{A}_x(S)) - \nabla\Phi_S(\mathcal{A}_x(S))\| \leq (1 + \kappa) \left( \frac{8G \left(1 + \frac{1}{cL(r+1)}\right)}{n} (24\kappa cL(r+1))^{\frac{1}{cL(r+1)+1}} T^{\frac{cL(r+1)}{cL(r+1)+1}} + \frac{G}{\sqrt{n}} \right). \quad (104)$$

**Proof** Denote  $\Delta_t \triangleq \sqrt{\|x_t - x'_t\|^2 + \|y_t - y'_t\|^2}$ , and the event  $E_{t_0} = \mathbf{1}(\Delta_{t_0} = 0)$ , we have for the full gradient  $\nabla f = (\nabla_x f, \nabla_y f)^\top$ ,

$$\begin{aligned} & \mathbb{E} \|\nabla f(x_t, y^*(x_t); \xi) - \nabla f(x'_t, y^*(x'_t); \xi)\| \\ &= \mathbb{P}(E_{t_0}) \mathbb{E} [\|\nabla f(x_t, y^*(x_t); \xi) - \nabla f(x'_t, y^*(x'_t); \xi)\| | E_{t_0}] \\ & \quad + \mathbb{P}(E_{t_0}^C) \mathbb{E} [\|\nabla f(x_t, y^*(x_t); \xi) - \nabla f(x'_t, y^*(x'_t); \xi)\| | E_{t_0}^C] \\ &\leq \mathbb{E} [\|\nabla f(x_t, y^*(x_t); \xi) - \nabla f(x'_t, y^*(x'_t); \xi)\| | E_{t_0}] + 2G\mathbb{P}(E_{t_0}^C) \\ &\leq \mathbb{E} [\|\nabla_x f(x_t, y^*(x_t); \xi) - \nabla_x f(x'_t, y^*(x'_t); \xi)\| + \|\nabla_y f(x_t, y^*(x_t); \xi) - \nabla_y f(x'_t, y^*(x'_t); \xi)\| | E_{t_0}] + 2G\mathbb{P}(E_{t_0}^C) \\ &\leq 2L\mathbb{E} [\|x_t - x'_t\| + \|y^*(x_t) - y^*(x'_t)\| | E_{t_0}] + 2G\mathbb{P}(E_{t_0}^C) \\ &\leq 2(1 + \kappa)L\mathbb{E} [\|x_t - x'_t\| | E_{t_0}] + 2G\frac{t_0}{n} \\ &\leq 4\kappa L\mathbb{E} [\Delta_t | \Delta_{t_0} = 0] + 2G\frac{t_0}{n}, \end{aligned} \quad (105)$$

the remaining steps aims to bound  $\mathbb{E}[\Delta_t | \Delta_{t_0} = 0]$ , which are the same as those in [Farnia and Ozdaglar, 2021, Appendix B.8], with that we will get

$$\mathbb{E} \|\nabla f(x_T, y^*(x_T); \xi) - \nabla f(x'_T, y^*(x'_T); \xi)\| \leq \frac{4\kappa L \cdot 12G}{nL} \left(\frac{T}{t_0}\right)^{cL(r+1)} + \frac{2G}{n} t_0, \quad (106)$$

to minimize the RHS above over  $t_0$ , we set

$$t_0 = \left( \frac{\frac{4\kappa L \cdot 12G}{nL} \cdot cL(r+1)}{\frac{2G}{n}} \right)^{\frac{1}{cL(r+1)+1}} \cdot T^{\frac{cL(r+1)}{cL(r+1)+1}} = (24\kappa cL(r+1))^{\frac{1}{cL(r+1)+1}} T^{\frac{cL(r+1)}{cL(r+1)+1}} \quad (107)$$

and we get

$$\mathbb{E} \|\nabla f(x_T, y^*(x_T); \xi) - \nabla f(x'_T, y^*(x'_T); \xi)\| \leq \frac{2G \left(1 + \frac{1}{cL(r+1)}\right)}{n} (24\kappa cL(r+1))^{\frac{1}{cL(r+1)+1}} T^{\frac{cL(r+1)}{cL(r+1)+1}}. \quad (108)$$

We conclude the proof by incorporating the above bound with Theorem 4.1.  $\blacksquare$

**Corollary F.2** Assume the function  $f$  is NC-C as defined in Assumption 2.1 with  $\mu = 0$ , then if we run SGDA for  $T$  iterations with stepsize  $\max\{\alpha_x, \alpha_y\} \leq \frac{c}{t}$  for some constant  $c > 0$ , we have

$$\mathbb{E}_{S, \mathcal{A}} \left\| \nabla\Phi^{1/(2L)}(\mathcal{A}_x(S)) - \nabla\Phi_S^{1/(2L)}(\mathcal{A}_x(S)) \right\| \leq \mathcal{O} \left( \left( \frac{T^{\frac{cL}{cL+1}}}{n} \right)^{1/6} + \left( \frac{1}{n} \right)^{1/8} \right). \quad (109)$$

**Proof** Denote  $\Delta_t \triangleq \sqrt{\|x_t - x'_t\|^2 + \|y_t - y'_t\|^2}$ , and the event  $E_{t_0} = \mathbf{1}(\Delta_{t_0} = 0)$ , we have that

$$\begin{aligned} \mathbb{E}\|x_t - x'_t\| &\leq \mathbb{E} \Delta_t \\ &= \mathbb{P}(E_{t_0})\mathbb{E}[\Delta_t|E_{t_0}] + \mathbb{P}(E_{t_0}^C)\mathbb{E}[\Delta_t|E_{t_0}^C] \\ &\leq \mathbb{E}[\Delta_t|E_{t_0}] + 2\sqrt{D_{\mathcal{X}} + D_{\mathcal{Y}}}\mathbb{P}(E_{t_0}^C) \\ &\leq \mathbb{E}[\Delta_t|\Delta_{t_0} = 0] + 2\sqrt{D_{\mathcal{X}} + D_{\mathcal{Y}}}\frac{t_0}{n}, \end{aligned} \quad (110)$$

the remaining steps aims to bound  $\mathbb{E}[\Delta_t|\Delta_{t_0} = 0]$ , with the results in [Farnia and Ozdaglar, 2021, Appendix B.9], we get

$$\mathbb{E}\|x_T - x'_T\| \leq \frac{2G}{nL} \left(\frac{T}{t_0}\right)^{cL} + 2\sqrt{D}\frac{t_0}{n}, \quad (111)$$

where we let  $D = D_{\mathcal{X}} + D_{\mathcal{Y}}$ . To minimize the RHS above over  $t_0$ , we set

$$t_0 = \left(\frac{cG}{\sqrt{D}}\right)^{1/(cL+1)} T^{\frac{cL}{cL+1}}, \quad (112)$$

and we get

$$\mathbb{E}\|x_T - x'_T\| \leq 2 \left( \frac{G}{L} \left(\frac{1}{cG}\right)^{\frac{cL}{cL+1}} + (cG)^{\frac{1}{cL+1}} \right) D^{\frac{cL}{2(cL+1)}} \frac{T^{\frac{cL}{cL+1}}}{n}. \quad (113)$$

The proof is complete by incorporating the above bound with Theorem 4.2. ■

## G Proof of Corollary 4.3 and 4.4

**Proof** For the NC-SC case, by Lei [2022, Corollary 6], we know the algorithm is  $\delta$  uniformly primal stable in gradients with  $\delta = 2G\sqrt{T/n}$ , the proof is complete by Theorem 4.1.

For the NC-C case, we want to derive the uniform primal argument stability, the flow here is almost the same as the proof of Lei [2022, Corollary 6], let  $\Omega_t = \|x_t - x'_t\|^2$ , define the event  $E_{\Omega}$  as that the only different data point  $\xi_i$  is selected by the algorithm  $\mathcal{A}$ , so we have

$$\mathbb{E}[\Omega_t] \leq \mathbb{E}[\Omega_t | E_{\Omega}]P(E_{\Omega}) + \mathbb{E}[\Omega_t | E_{\Omega}^C]\mathbb{P}(E_{\Omega}^C) \leq \mathbb{E}[\Omega_t | E_{\Omega}^C] \frac{T}{n} \leq \frac{4D_{\mathcal{X}}T}{n}, \quad (114)$$

so the algorithm is  $\sqrt{4D_{\mathcal{X}}T/n}$ -uniformly primal argument stable. Then we conclude the proof by substituting the above stability results into Theorem 4.2, i.e.,

$$\begin{aligned} &\mathbb{E}_{S, \mathcal{A}} \left\| \nabla \Phi^{1/(2L)}(\mathcal{A}_x(S)) - \nabla \Phi_S^{1/(2L)}(\mathcal{A}_x(S)) \right\| \\ &= \mathcal{O}\left(\delta^{\frac{1}{6}} + n^{-\frac{1}{8}}\right) = \mathcal{O}\left(\left(\frac{T}{n}\right)^{\frac{1}{12}} + \left(\frac{1}{n}\right)^{\frac{1}{8}}\right) = \mathcal{O}\left(\left(\frac{T}{n}\right)^{\frac{1}{12}} + \left(\frac{1}{n}\right)^{\frac{1}{8}}\right). \end{aligned} \quad (115)$$

■