

# Gradient Descent Ascent for Minimax Problems on Riemannian Manifolds

Feihu Huang, Shangqian Gao

**Abstract**—In the paper, we study a class of useful minimax problems on Riemannian manifolds and propose a class of effective Riemannian gradient-based methods to solve these minimax problems. Specifically, we propose an effective Riemannian gradient descent ascent (RGDA) algorithm for the deterministic minimax optimization. Moreover, we prove that our RGDA has a sample complexity of  $O(\kappa^2 \epsilon^{-2})$  for finding an  $\epsilon$ -stationary solution of the Geodesically-Nonconvex Strongly-Concave (GNSC) minimax problems, where  $\kappa$  denotes the condition number. At the same time, we present an effective Riemannian stochastic gradient descent ascent (RSGDA) algorithm for the stochastic minimax optimization, which has a sample complexity of  $O(\kappa^4 \epsilon^{-4})$  for finding an  $\epsilon$ -stationary solution. To further reduce the sample complexity, we propose an accelerated Riemannian stochastic gradient descent ascent (Acc-RSGDA) algorithm based on the momentum-based variance-reduced technique. We prove that our Acc-RSGDA algorithm achieves a lower sample complexity of  $\tilde{O}(\kappa^4 \epsilon^{-3})$  in searching for an  $\epsilon$ -stationary solution of the GNSC minimax problems. Extensive experimental results on the robust distributional optimization and robust Deep Neural Networks (DNNs) training over Stiefel manifold demonstrate efficiency of our algorithms.

**Index Terms**—Riemannian Manifolds, Minimax Optimization, Stiefel Manifold, Deep Neural Networks, Robust Optimization.



## 1 INTRODUCTION

IN this paper, we study a class of useful minimax optimization problems on the Riemannian manifold  $\mathcal{M}$ , defined as:

$$\min_{x \in \mathcal{M}} \max_{y \in \mathcal{Y}} f(x, y), \quad (1)$$

where function  $f(x, y) : \mathcal{M} \times \mathcal{Y} \rightarrow \mathbb{R}$  is  $\mu$ -strongly concave in  $y \in \mathcal{Y} \subseteq \mathbb{R}^d$  but possibly (geodesically) nonconvex in  $x \in \mathcal{M}$ . Here  $\mathcal{M}$  is a Riemannian manifold, and  $\mathcal{Y}$  is a convex and closed set in Euclidean space.  $f(\cdot, y) : \mathcal{M} \rightarrow \mathbb{R}$  for any  $y \in \mathcal{Y}$  is a smooth but possibly (geodesically) nonconvex real-valued function on manifold  $\mathcal{M}$ , and  $f(x, \cdot) : \mathcal{Y} \rightarrow \mathbb{R}$  for any  $x \in \mathcal{M}$  is a smooth and strongly-concave real-valued function. Note that a geodesically nonconvex function on Riemannian manifold also is nonconvex on Euclidean space, and a geodesically convex function on Riemannian manifold may be nonconvex on Euclidean space. In this paper, we also focus on the stochastic form of minimax problem (1), defined as

$$\min_{x \in \mathcal{M}} \max_{y \in \mathcal{Y}} \mathbb{E}_{\xi \sim \mathcal{D}} [f(x, y; \xi)], \quad (2)$$

where  $\xi$  is a random variable that follows an unknown distribution  $\mathcal{D}$ . In fact, Problems (1) and (2) are associated to many existing machine learning applications:

**1). Robust DNNs Training over Riemannian manifold.** Deep Neural Networks (DNNs) recently have been demonstrating exceptional performance on many machine learning

applications such as image classification. However, they are vulnerable to the adversarial example attacks, which show that a small perturbation in the data input can significantly change the output of DNNs. Thus, the security properties of DNNs have been widely studied. One of secured DNN research topics is to enhance the robustness of DNNs under the adversarial example attacks. Given the training sample  $\mathcal{D} := \{\xi_i = (a_i, b_i)\}_{i=1}^n$ , where  $a_i \in \mathbb{R}^d$  and  $b_i \in \mathbb{R}$  represent the features and label of sample  $\xi_i$  respectively. Then we train a robust DNN against a universal adversarial attack [1], [2], which can be formulated the following minimax problem:

$$\min_{x \in \mathbb{R}^q} \max_{y \in \mathcal{Y}} \frac{1}{n} \sum_{i=1}^n \ell(h(a_i + y; x), b_i), \quad (3)$$

where  $x \in \mathbb{R}^q$  denotes weight of the DNN, and  $h(\cdot; x)$  denotes the DNN parameterized by  $x$ , and  $\ell(\cdot)$  is the loss function. Here  $y$  denotes a small universal perturbation in the features  $\{a_i\}_{i=1}^n$ , and the constraint  $\mathcal{Y} = \{y : \|y\|_\infty \leq \epsilon\}$  indicates that the poisoned samples should not be too different from the original ones.

Recently, the orthonormality on weights of DNNs has gained much interest and has been found to be useful across different tasks such as person re-identification [3] and image classification [4]. In fact, the orthonormality constraints improve the performances of DNNs [5], [6], and reduce overfitting to improve generalization [7]. At the same time, the orthonormality can stabilize the distribution of activation over layers within DNNs [8]. Thus, we further consider the following robust DNN training over the Stiefel manifold  $\mathcal{M}$ :

$$\min_{x \in \mathcal{M}} \max_{y \in \mathcal{Y}} \frac{1}{n} \sum_{i=1}^n \ell(h(a_i + y; x), b_i), \quad (4)$$

• Feihu Huang is with College of Computer Science and Technology, Nanjing University of Aeronautics and Astronautics, Nanjing, China; and also with MIIT Key Laboratory of Pattern Analysis and Machine Intelligence, Nanjing, China. E-mail: huangfeihu2018@gmail.com; huangfeihu@nuaa.edu.cn  
Shangqian Gao is with Department of Electrical and Computer Engineering, University of Pittsburgh, Pittsburgh, USA. E-mail: shg84@pitt.edu

When data are continuously coming, we can rewrite the stochastic form of Problem (4) as follows:

$$\min_{x \in \mathcal{M}} \max_{y \in \mathcal{Y}} \mathbb{E}_{\xi} [f(x, y; \xi)], \quad (5)$$

where  $f(x, y; \xi) = \ell(h(a + y; x), b)$  with  $\xi = (a, b)$ .

**2). Distributionally Robust Optimization over Riemannian manifold.** Distributionally Robust Optimization (DRO) [9], [10] is an effective method to deal with the noisy data, adversarial data, and imbalanced data. In the paper, we consider the DRO over the Riemannian manifold that can be applied in many machine learning problems such as robust principal component analysis (PCA) and distributionally robust DNN training. To be more specific, given a set of data samples  $\{\xi_i\}_{i=1}^n$ , the DRO over Riemannian manifold  $\mathcal{M}$  can be written as the following minimax problem:

$$\min_{x \in \mathcal{M}} \max_{\mathbf{p} \in \mathcal{S}} \left\{ \sum_{i=1}^n p_i \ell(x; \xi_i) - \left\| \mathbf{p} - \frac{\mathbf{1}}{n} \right\|^2 \right\}, \quad (6)$$

where  $\mathbf{p} = (p_1, \dots, p_n)$ ,  $\mathcal{S} = \{\mathbf{p} \in \mathbb{R}^n : \sum_{i=1}^n p_i = 1, p_i \geq 0\}$ . Here  $\ell(x; \xi_i)$  denotes the loss function over Riemannian manifold  $\mathcal{M}$ , which applies to many machine learning problems such as PCA [11], dictionary learning [12], DNNs [8], structured low-rank matrix learning [13], [14], [15], among others. For example, the task of PCA can be cast on a Grassmann manifold.

Recently some algorithms [16], [17], [18] have been studied for variational inequalities on Riemannian manifolds, which are the implicit minimax problems on Riemannian manifolds. Meanwhile, some methods [19], [20] for computing the projection robust Wasserstein distance, which can be represented as a minimax optimization over the Stiefel manifold [21]. To the best of our knowledge, the existing explicitly minimax optimization methods such as gradient descent ascent method only focus on the minimax problems in Euclidean space.

To fill this gap, in the paper, we study the explicit minimax optimization problems over the general Riemannian manifold, and propose a class of efficient Riemannian gradient-based algorithms to solve the Geodesically-Nonconvex Strongly-Concave (GNSC) minimax problem (1) via using general retraction and vector transport. When Problem (1) is deterministic, we propose a new deterministic Riemannian gradient descent ascent algorithm. When Problem (1) is stochastic (i.e, Problem (2)), we propose two efficient stochastic Riemannian gradient descent ascent algorithms. Our main **contributions** can be summarized as follows:

- 1) We propose an effective Riemannian gradient descent ascent (RGDA) algorithm for the deterministic minimax Problem (1). Moreover, we prove that the RGDA has a sample complexity of  $O(\kappa^2 \epsilon^{-2})$  in finding an  $\epsilon$ -stationary solution of Problem (1).
- 2) Meanwhile, we present an effective Riemannian stochastic gradient descent ascent (RSGDA) algorithm for the stochastic minimax Problem (2), which has a sample complexity of  $O(\kappa^4 \epsilon^{-4})$  in searching for an  $\epsilon$ -stationary solution of Problem (2).

- 3) We further propose an accelerated Riemannian stochastic gradient descent ascent (Acc-RSGDA) algorithm based on the variance-reduced technique of STORM [22]. We prove our Acc-RSGDA achieves a lower sample complexity of  $\tilde{O}(\kappa^4 \epsilon^{-3})$ .
- 4) Extensive experimental results on the robust DNNs training and distributionally robust optimization over Stiefel manifold demonstrate the efficiency of our proposed algorithms.

## 2 RELATED WORKS

In this section, we briefly review the minimax optimization and Riemannian manifold optimization, respectively.

### 2.1 Minimax Optimization

Minimax optimization [23] recently has been widely applied in many machine learning problems such as adversarial training [24], reinforcement learning [25], and robust federated learning [26]. Meanwhile, many efficient minimax methods [27], [28], [29], [30], [31], [32], [33], [34], [35], [36], [37], [38], [39] have been proposed for solving these minimax optimization problems. For example, [29] proposed a class of efficient dual implicit accelerated gradient algorithms to solve smooth minimax optimization. [27] studied the convergence properties of the gradient descent ascent (GDA) methods for nonconvex minimax optimization. Subsequently, the accelerated GDA algorithms [30] have been proposed for minimax optimization. Meanwhile, [33] presented a catalyst accelerated framework for minimax optimization. Moreover, [36], [39] proposed some faster stochastic variance-reduced GDA algorithms to solve the stochastic nonconvex-strongly-concave minimax problems. [32] studied the convergence properties of GDA methods for solving a class of nonconvex-nonconcave minimax problems. More recently, a class of efficient mirror descent ascent algorithms [38] have been proposed for nonconvex nonsmooth minimax optimization.

### 2.2 Riemannian Manifold Optimization

Riemannian manifold optimization methods have been widely applied in machine learning problems including dictionary learning [12], low-rank matrix completion [14], [15], DNNs [8] and natural language processing [40]. Many Riemannian optimization methods have been recently proposed. *E.g.* [41], [42] proposed some efficient first-order gradient methods for geodesically convex functions. Subsequently, [43] presented fast stochastic variance-reduced methods to Riemannian manifold optimization. More recently, [44] proposed fast first-order gradient algorithms for Riemannian manifold optimization by using general retraction and vector transport. Subsequently, based on these retraction and vector transport, some fast Riemannian gradient-based methods [11], [45], [46], [47], [48] have been proposed for non-convex optimization. Riemannian Adam-type algorithms [49] have been introduced for matrix manifold optimization. Subsequently, [40] proposed an efficient Riemannian adaptive optimization algorithm to natural language processing. Meanwhile, some algorithms

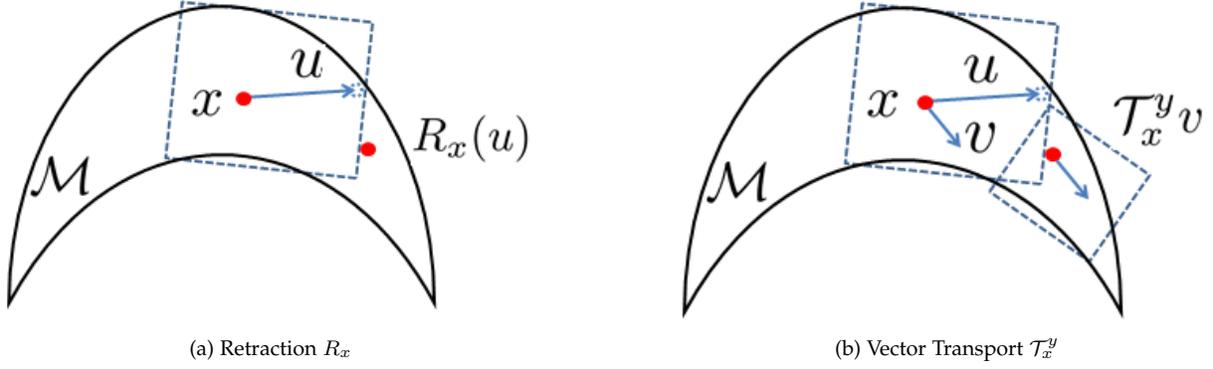


Fig. 1: Illustration of manifold operations.(a) A vector  $u$  in  $T_x\mathcal{M}$  is mapped to  $R_x(u)$  in  $\mathcal{M}$ ; (b) A vector  $v$  in  $T_x\mathcal{M}$  is transported to  $T_y\mathcal{M}$  by  $\mathcal{T}_x^y v$  (or  $\mathcal{T}_u v$ ), where  $y = R_x(u)$  and  $u \in T_x\mathcal{M}$ .

[16], [17], [18] have been studied for variational inequalities on Riemannian manifolds, which are the implicit minimax problems on Riemannian manifolds. More recently, [50] studied the stochastic composition optimization on Riemannian manifolds.

**Notations:**  $I_d$  denotes the identity matrix with  $d$  dimension.  $\text{diag}(a) \in \mathbb{R}^{d \times d}$  denotes a diagonal matrix, whose diagonal elements come from vector  $a \in \mathbb{R}^d$ .  $\text{sign}(\cdot)$  denotes the sign function, i.e., if  $x > 0$ ,  $\text{sign}(x) = 1$ ; if  $x = 0$ ,  $\text{sign}(x) = 0$ ; otherwise  $\text{sign}(x) = -1$ .  $\|\cdot\|$  denotes the  $\ell_2$  norm for vectors and Frobenius norm for matrices.  $\langle x, y \rangle$  denotes the inner product of two vectors  $x$  and  $y$ . For function  $f(x, y)$ ,  $f(x, \cdot)$  denotes function w.r.t. the second variable with fixing  $x$ , and  $f(\cdot, y)$  denotes function w.r.t. the first variable with fixing  $y$ . Given a convex closed set  $\mathcal{Y}$ , we define a projection operation on the set  $\mathcal{Y}$  as  $\mathcal{P}_{\mathcal{Y}}(y_0) = \arg \min_{y \in \mathcal{Y}} \frac{1}{2} \|y - y_0\|^2$ . We denote  $a = O(b)$  if  $a \leq Cb$  for some constant  $C > 0$ , and the notation  $\tilde{O}(\cdot)$  hides logarithmic terms. The operation  $\oplus$  denotes the Whitney sum that takes two vector bundles over a fixed space and produces a new vector bundle over the same space. Given function  $f(x)$ , let  $\text{grad}f(x)$  denote its Riemannian gradients at Riemannian manifold and  $\nabla f(x)$  denote its gradients at Euclidean space. Given  $\mathcal{B}_t = \{\xi_t^i\}_{i=1}^B$  for any  $t \geq 1$ , let  $\nabla f_{\mathcal{B}_t}(x) = \frac{1}{B} \sum_{i=1}^B \nabla f(x; \xi_t^i)$  and  $\text{grad}f_{\mathcal{B}_t}(x) = \frac{1}{B} \sum_{i=1}^B \text{grad}f(x; \xi_t^i)$ .

### 3 PRELIMINARIES

In this section, we first re-visit some basic information on the Riemannian manifold  $\mathcal{M}$ . In general, the manifold  $\mathcal{M}$  is endowed with a smooth inner product  $\langle \cdot, \cdot \rangle_x : T_x\mathcal{M} \cdot T_x\mathcal{M} \rightarrow \mathbb{R}$  on tangent space  $T_x\mathcal{M}$  for every  $x \in \mathcal{M}$ . The induced norm  $\|\cdot\|_x$  of a tangent vector in  $T_x\mathcal{M}$  is associated with the Riemannian metric. We first define a retraction  $R_x : T_x\mathcal{M} \rightarrow \mathcal{M}$  mapping tangent space  $T_x\mathcal{M}$  onto  $\mathcal{M}$  with a local rigidity condition that preserves the gradients at  $x \in \mathcal{M}$  (please see Fig.1 (a)). The retraction  $R_x$  satisfies all of the following: 1)  $R_x(0) = x$ , where  $0 \in T_x\mathcal{M}$ ; 2)  $DR_x(0) = id_{T_x\mathcal{M}}$ , where  $DR_x$  denotes the derivative of  $R_x$ , and  $id_{T_x\mathcal{M}}$  denotes an identity mapping on  $T_x\mathcal{M}$ . In fact, exponential mapping  $\text{Exp}_x$  is a special case of retraction  $R_x$  that locally approximates the exponential mapping  $\text{Exp}_x$  to the first order on the manifold.

Next, we define a vector transport  $\mathcal{T} : T\mathcal{M} \oplus T\mathcal{M} \rightarrow T\mathcal{M}$  (please see Fig.1 (b)) that satisfies all of the following 1)  $\mathcal{T}$  has an associated retraction  $R$ , i.e., for  $x \in \mathcal{M}$  and  $w, u \in T_x\mathcal{M}$ ,  $\mathcal{T}_u w$  is a tangent vector at  $R_x(w)$ ; 2)  $\mathcal{T}_0 v = v$ ; 3)  $\mathcal{T}_u(av + bw) = a\mathcal{T}_u v + b\mathcal{T}_u w$  for all  $a, b \in \mathbb{R}$  a  $u, v, w \in T\mathcal{M}$ . Vector transport  $\mathcal{T}_x^y v$  or equivalently  $\mathcal{T}_u v$  with  $y = R_x(u)$  transports  $v \in T_x\mathcal{M}$  along the retraction curve defined by direction  $u$ . Here we focus on the isometric vector transport  $\mathcal{T}_x^y$ , which satisfies  $\langle u, v \rangle_x = \langle \mathcal{T}_x^y u, \mathcal{T}_x^y v \rangle_y$  for all  $u, v \in T_x\mathcal{M}$ . Based on these definitions, we provide some standard assumptions about Problems (1) and (2).

**Assumption 1.**  $\mathcal{X} \subseteq \mathcal{M}$  is compact. Each component function  $f(x, y)$  is twice continuously differentiable in both  $x \in \mathcal{X}$  and  $y \in \mathcal{Y}$ , and there exist constants  $L_{11}$ ,  $L_{12}$ ,  $L_{21}$  and  $L_{22}$ , such that for every  $x, x_1, x_2 \in \mathcal{X}$  and  $y, y_1, y_2 \in \mathcal{Y}$ , we have

$$\begin{aligned} \|\text{grad}_x f(x_1, y; \xi) - \mathcal{T}_{x_2}^{x_1} \text{grad}_x f(x_2, y; \xi)\| &\leq L_{11} \|u\|, \\ \|\text{grad}_x f(x, y_1; \xi) - \text{grad}_x f(x, y_2; \xi)\| &\leq L_{12} \|y_1 - y_2\|, \\ \|\nabla_y f(x_1, y; \xi) - \nabla_y f(x_2, y; \xi)\| &\leq L_{21} \|u\|, \\ \|\nabla_y f(x, y_1; \xi) - \nabla_y f(x, y_2; \xi)\| &\leq L_{22} \|y_1 - y_2\|, \end{aligned}$$

where  $u \in T_{x_1}\mathcal{M}$  and  $x_2 = R_{x_1}(u)$ .

Assumption 1 is commonly used in Riemannian optimization [11], [44], and minimax optimization [27], [36]. Here, the terms  $L_{11}$ ,  $L_{12}$  and  $L_{21}$  implicitly contain the curvature information as in [11], [44]. Specifically, Assumption 1 implies the partial Riemannian gradient  $\text{grad}_x f(\cdot, y; \xi)$  for all  $y \in \mathcal{Y}$  is  $L_{11}$ -Lipschitz continuous with respect to retraction as in [11] and the partial gradient  $\nabla_y f(x, \cdot; \xi)$  for all  $x \in \mathcal{X}$  is  $L_{22}$ -Lipschitz continuous as in [27].

To further verify the rationality of Assumption 1, we consider the Stiefel manifold  $\mathcal{M} = \{X \in \mathbb{R}^{d \times r} | X^T X = I_r\}$ . For notational simplicity, let matrix  $X$  instead of the variable  $x$  in Assumption 1. Let  $\nabla_X f(X, y)$  denote the gradient of  $f(X, y)$  on variable  $X$  in the Euclidean space, and  $\text{grad}_X f(X, y)$  denote the Riemannian gradient of  $f(X, y)$  on variable  $X$  in the Stiefel manifold. Following [5],  $\text{grad}_X f(X, y)$  can be seen as a projection onto the tangent space  $T_X\mathcal{M}$  of Riemannian  $\mathcal{M}$  at  $X$ , which can be computed

as follows:

$$\begin{aligned} \text{grad}_X f(X, y) &= P_{T_x}(\nabla_X f(X, y)) = WX, \\ W &= \hat{W} - \hat{W}^T, \\ \hat{W} &= \nabla_X f(X, y)X^T - \frac{1}{2}X(X^T \nabla_X f(X, y)X^T). \end{aligned} \quad (7)$$

Then we have for any  $X_1, X_2 \in \mathcal{M}$ ,

$$\begin{aligned} &\|\text{grad}_X f(X_1, y) - \mathcal{T}_{X_2}^{X_1} \text{grad}_X f(X_2, y)\| \\ &= \|P_{T_{X_1}}(\nabla_X f(X_1, y)) - \mathcal{T}_{X_2}^{X_1} P_{T_{X_2}}(\nabla_X f(X_2, y))\|, \\ &\leq \|\nabla_X f(X_1, y) - \nabla_X f(X_2, y)\| \leq L\|X_1 - X_2\|, \end{aligned} \quad (8)$$

where the last inequality holds by Lipschitz continuous for gradient in the Euclidean space. Let  $d(X_1, X_2)$  denote geodesic distance between  $X_1$  and  $X_2$  in  $\mathcal{M}$ , then we have  $d(X_1, X_2) = \zeta\|X_1 - X_2\|$ , where  $\zeta > 0$  denote curvature parameter of manifold  $\mathcal{M}$ . In our Assumption 1, due to  $X_2 = R_{X_1}(u)$ , we have  $\|u\| = d(X_1, X_2)$ . According to the above (8), we have

$$\begin{aligned} &\|\text{grad}_X f(X_1, y) - \mathcal{T}_{X_2}^{X_1} \text{grad}_X f(X_2, y)\| \\ &\leq L\|X_1 - X_2\| = \frac{L}{\zeta}d(X_1, X_2) = \frac{L}{\zeta}\|u\|, \end{aligned} \quad (9)$$

where  $X_2 = R_{X_1}(u)$ . This similarly holds for the other inequalities in our Assumption 1.

For the deterministic problem, let  $f(x, y)$  instead of  $f(x, y; \xi)$  in Assumption 1. In fact, these Lipschitz continuity assumptions are widely applicable to deep learning architectures [5]. Note that in the following experiments, given the DNNs using ReLU, the derivative of ReLU is Lipschitz continuous almost everywhere with an appropriate Lipschitz constant, except for a small neighbourhood around 0, whose measure tends to 0. Such cases do not affect either analysis in theory or training in practice.

Since  $f(x, y)$  is strongly concave in  $y \in \mathcal{Y}$ , there exists a unique solution to the problem  $\max_{y \in \mathcal{Y}} f(x, y)$  for any  $x$ . We define the function  $\Phi(x) = \max_{y \in \mathcal{Y}} f(x, y)$  and  $y^*(x) = \arg \max_{y \in \mathcal{Y}} f(x, y)$ .

**Assumption 2.** *The function  $\Phi(x) : \mathcal{M} \rightarrow \mathbb{R}$  is  $L$ -smooth. There exists a constant  $L > 0$ , for all  $x \in \mathcal{X}, z = R_x(u)$  with  $u \in T_x \mathcal{M}$ , such that*

$$\Phi(z) \leq \Phi(x) + \langle \text{grad} \Phi(x), u \rangle + \frac{L}{2}\|u\|^2.$$

**Assumption 3.** *The objective function  $f(x, y)$  is  $\mu$ -strongly concave w.r.t  $y$ , i.e., for any  $x \in \mathcal{M}, y_1, y_2 \in \mathcal{Y}$*

$$f(x, y_1) \leq f(x, y_2) + \langle \nabla_y f(x, y_2), y_1 - y_2 \rangle - \frac{\mu}{2}\|y_1 - y_2\|^2.$$

**Assumption 4.** *The function  $\Phi(x)$  is bounded from below in  $\mathcal{M}$ , i.e.,  $\Phi^* = \inf_{x \in \mathcal{M}} \Phi(x)$ .*

**Assumption 5.** *The variance of stochastic gradient is bounded, i.e., there exists a constant  $\sigma_1 > 0$  such that for all  $x$ , it follows  $\mathbb{E}_\xi \|\text{grad}_x f(x, y; \xi) - \text{grad}_x f(x, y)\|^2 \leq \sigma_1^2$ ; There exists a constant  $\sigma_2 > 0$  such that for all  $y$ , it follows  $\mathbb{E}_\xi \|\nabla_y f(x, y; \xi) - \nabla_y f(x, y)\|^2 \leq \sigma_2^2$ . We also define  $\sigma = \max\{\sigma_1, \sigma_2\}$ .*

Assumption 2 imposes the smooth of function  $\Phi(x)$  over Riemannian manifold  $\mathcal{M}$ , as in [11], [44], [48]. Assumption 3 imposes the strongly concave of  $f(x, y)$  on variable  $y$ , as in [27], [36]. Assumption 4 guarantees the feasibility of

---

### Algorithm 1 RGDA and RSGDA Algorithms

---

- 1: **Input:**  $T$ , parameters  $\{\gamma, \lambda, \eta_t\}_{t=1}^T$ , mini-batch size  $B$ , and initial input  $x_1 \in \mathcal{M}, y_1 \in \mathcal{Y}$ ;
- 2: **for**  $t = 1, 2, \dots, T$  **do**
- 3:     **(RGDA)** Compute deterministic gradients

$$v_t = \text{grad}_x f(x_t, y_t), \quad w_t = \nabla_y f(x_t, y_t);$$

- 4:     **(RSGDA)** Draw  $B$  i.i.d. samples  $\{\xi_t^i\}_{i=1}^B$ , then compute stochastic gradients

$$v_t = \frac{1}{B} \sum_{i=1}^B \text{grad}_x f(x_t, y_t; \xi_t^i),$$

$$w_t = \frac{1}{B} \sum_{i=1}^B \nabla_y f(x_t, y_t; \xi_t^i);$$

- 5:     Update:  $x_{t+1} = R_{x_t}(-\gamma \eta_t v_t)$ ;
  - 6:     Update:  $\tilde{y}_{t+1} = \mathcal{P}_Y(y_t + \lambda w_t)$  and  $y_{t+1} = y_t + \eta_t(\tilde{y}_{t+1} - y_t)$ ;
  - 7:     **end for**
  - 8: **Output:**  $x_\zeta$  and  $y_\zeta$  chosen uniformly random from  $\{x_t, y_t\}_{t=1}^T$ .
- 

the GNSC minimax problem (1), as the nonconvex-strongly-concave minimax optimization on Euclidean space used in [27], [36]. Assumption 5 imposes the bounded variance of stochastic (Riemannian) gradients, which is commonly used in the stochastic optimization [27], [36], [48].

## 4 RIEMANNIAN GRADIENT-BASED METHODS

In this section, we propose a class of Riemannian gradient-based methods to solve the deterministic and stochastic GNSC minimax problems (1) and (2), respectively.

### 4.1 RGDA and RSGDA Algorithms

In this subsection, we propose an efficient Riemannian gradient descent ascent (RGDA) algorithm to solve the deterministic minimax Problem (1). At the same time, we propose a standard Riemannian stochastic gradient descent ascent (RSGDA) algorithm to solve the stochastic minimax Problem (2). Algorithm 1 summarizes the algorithmic framework of our RGDA and RSGDA algorithms.

At the line 3 of Algorithm 1, we calculate the deterministic Riemannian gradient in variable  $x \in \mathcal{M}$ , and calculate the deterministic gradient in variable  $y \in \mathcal{Y}$ . At the line 4 of Algorithm 1, we calculate the stochastic Riemannian gradient for variable  $x \in \mathcal{M}$ , and calculate the stochastic gradient for variable  $y \in \mathcal{Y}$ .

At the line 5 of Algorithm 1, we use the Riemannian gradient descent to update variable  $x$  based on the retraction operator  $R_{x_t}(\cdot)$ , which guarantees the variable  $x_t$  for all  $t \geq 1$  in the manifold  $\mathcal{M}$ . Here  $R_{x_t}(\cdot)$  can be seen as a generalized projection operator, which can be competent to the general Riemannian manifolds. For example, we consider the popular Stiefel manifold  $\mathcal{M} = \text{St}(r, d) = \{X \in \mathbb{R}^{d \times r} : X^T X = I_r\}$  that is a nonconvex constraint set in the Euclidean space. Given  $g_t = -\gamma \eta_t v_t \in T_{x_t} \mathcal{M}$ , we can define a standard QR-based retraction:  $R_{x_t}(g_t) = QH$ ,

**Algorithm 2** Acc-RSGDA Algorithm

- 1: **Input:**  $T$ , parameters  $\{\gamma, \lambda, b, m, c_1, c_2\}$  and initial input  $x_1 \in \mathcal{M}$  and  $y_1 \in \mathcal{Y}$ ;
- 2: Draw  $B$  i.i.d. samples  $\mathcal{B}_1 = \{\xi_1^i\}_{i=1}^B$ , then compute  $v_1 = \text{grad}_x f_{\mathcal{B}_1}(x_1, y_1)$  and  $w_1 = \nabla_y f_{\mathcal{B}_1}(x_1, y_1)$ ;
- 3: **for**  $t = 1, 2, \dots, T$  **do**
- 4:   Update:  $x_{t+1} = R_{x_t}(-\gamma\eta_t v_t)$  with  $\eta_t = \frac{b}{(m+t)^{1/3}}$ ;
- 5:   Update:  $\tilde{y}_{t+1} = \mathcal{P}_{\mathcal{Y}}(y_t + \lambda w_t)$  and  $y_{t+1} = y_t + \eta_t(\tilde{y}_{t+1} - y_t)$ ;
- 6:   Draw  $B$  i.i.d. samples  $\mathcal{B}_{t+1} = \{\xi_{t+1}^i\}_{i=1}^B$ , then compute
 
$$v_{t+1} = \text{grad}_x f_{\mathcal{B}_{t+1}}(x_{t+1}, y_{t+1}) + (1 - \alpha_{t+1}) \cdot \mathcal{T}_{x_{t+1}}^{x_{t+1}} [v_t - \text{grad}_x f_{\mathcal{B}_{t+1}}(x_t, y_t)], \quad (10)$$

$$w_{t+1} = \nabla_y f_{\mathcal{B}_{t+1}}(x_{t+1}, y_{t+1}) + (1 - \beta_{t+1}) \cdot [w_t - \nabla_y f_{\mathcal{B}_{t+1}}(x_t, y_t)], \quad (11)$$
 where  $\alpha_{t+1} = c_1 \eta_t^2$  and  $\beta_{t+1} = c_2 \eta_t^2$ .
- 7: **end for**
- 8: **Output:**  $x_\zeta$  and  $y_\zeta$  chosen uniformly random from  $\{x_t, y_t\}_{t=1}^T$ .

where the matrices  $Q$  and  $H$  can be obtained from the QR decomposition of matrix  $x_t + g_t \in \mathbb{R}^{d \times r}$ , i.e.,  $x_t + g_t = QR$ , and  $H = \text{diag}(\{\text{sign}(R_{i,i})\}_{i=1}^r)$ . It is well known that the standard projected gradient methods with convergence guarantee require the convex constraint sets belonging to Euclidean space [51], while our Riemannian gradient-based methods with convergence guarantee do not need the convex constraint sets (Please see the following convergence analysis).

At the line 6 of Algorithm 1, we simultaneously use a projection iteration and a momentum iteration to update the variable  $y$ , where we use  $0 < \eta_t \leq 1$  to ensure the variable  $y_t$  for all  $t \geq 1$  in convex constraint  $\mathcal{Y}$ . Note that we use two learning rates  $\gamma$  and  $\eta_t$  at the line 5, where  $\gamma$  is a constant learning rate and  $\eta_t$  is a dynamic or constant learning rate with iteration  $t$ . Under this case, we can flexibly choose learning rates in practice, and can easily analyze the convergence properties of our algorithms, where simultaneously Riemannian gradient descent on the variable  $x \in \mathcal{M}$  and gradient ascent on the variable  $y \in \mathcal{Y}$ .

**4.2 Acc-RSGDA algorithm**

In this subsection, we propose an accelerated stochastic Riemannian gradient descent ascent (Acc-RSGDA) algorithm to solve the stochastic minimax Problem (2), which builds on the momentum-based variance reduction technique of STORM [22]. Algorithm 2 describes the algorithmic framework of Acc-RSGDA method.

At the line 4 of Algorithm 2, we use two learning rates  $\gamma$  and  $\eta_t$ , where  $\gamma$  is a constant learning rate and  $\eta_t = \frac{b}{(m+t)^{1/3}}$  is a decreasing learning rate with iteration  $t$ . Similarly, we can flexibly choose learning rates in practice, and can easily analyze the convergence properties of our algorithms, where simultaneously Riemannian gradient descent on the variable  $x \in \mathcal{M}$  and gradient ascent on the variable  $y \in \mathcal{Y}$ .

At the line 6 of Algorithm 2, we use the momentum-based variance-reduced technique of STORM to estimate stochastic Riemannian gradient  $v_t$  defined in (10).

where  $\alpha_{t+1} \in (0, 1]$ . When  $\alpha_{t+1} = 1$ ,  $v_{t+1} = \text{grad}_x f_{\mathcal{B}_{t+1}}(x_{t+1}, y_{t+1})$  will degenerate a vanilla stochastic Riemannian gradient estimator; When  $\alpha_{t+1} = 0$ ,  $v_{t+1} = \text{grad}_x f_{\mathcal{B}_{t+1}}(x_{t+1}, y_{t+1}) - \mathcal{T}_{x_t}^{x_{t+1}}(\text{grad}_x f_{\mathcal{B}_{t+1}}(x_t, y_t) - v_t)$  will degenerate a stochastic Riemannian gradient estimator based on variance-reduced technique of SPIDER [52]. Since our Acc-RSGDA algorithm uses variance-reduced technique of STORM to estimate the stochastic gradients, it does not rely on large mini-batch size to guarantee its convergence (Please see the following convergence analysis).

Riemannian gradient  $\text{grad}_x f_{\mathcal{B}_{t+1}}(x_{t+1}, y_{t+1})$  is over the tangent space  $T_{x_{t+1}}\mathcal{M}$ , while the Riemannian gradient estimator  $\text{grad}_x f_{\mathcal{B}_{t+1}}(x_t, y_t) - v_t$  is over the tangent space  $T_{x_t}\mathcal{M}$ . In order to feasibility of  $v_{t+1}$ , we use the vector transport  $\mathcal{T}_{x_t}^{x_{t+1}}$  to project the Riemannian gradient estimator  $\text{grad}_x f_{\mathcal{B}_{t+1}}(x_t, y_t) - v_t$  into the tangent space  $T_{x_{t+1}}\mathcal{M}$ . Thus, we can add the term  $\text{grad}_x f_{\mathcal{B}_{t+1}}(x_{t+1}, y_{t+1})$  and the term  $(1 - \alpha_{t+1})\mathcal{T}_{x_t}^{x_{t+1}}[v_t - \text{grad}_x f_{\mathcal{B}_{t+1}}(x_t, y_t)]$ .

**4.3 Novelities of Our Algorithms**

Compared with the existing Riemannian gradient algorithms [11], [45], [46] and minimax optimization algorithms [27], [36], our algorithms have the following main differences:

- 1) Compared with the existing Riemannian gradient algorithms, our algorithms simultaneously use a constant learning rate  $\gamma$  and a dynamic or constant learning rate  $\eta_t$  at each iteration. This dynamic/constant learning rate  $\eta_t$  is the same tuning parameter of the **momentum iteration** in updating variable  $y$  (i.e.,  $y_{t+1} = y_t + \eta_t(\tilde{y}_{t+1} - y_t)$ ). In other words, the learning rate in updating the variable  $x \in \mathcal{M}$  depends on the tuning parameter of the momentum iteration in updating dual variable  $y$ .
- 2) Compared with the existing minimax optimization algorithms, our algorithms simultaneously use a projection iteration and a momentum iteration to update the variable  $y$ . Meanwhile, our algorithms use the Riemannian gradients and retraction operator to update variable  $x \in \mathcal{M}$  instead of the standard gradients and projection operator used in the existing minimax algorithms.

**5 CONVERGENCE ANALYSIS**

In this section, we study the convergence properties of our RGDA, RSGDA, and Acc-RSGDA algorithms, respectively. The basic idea of our convergence analysis is given in Fig. 2. We first give some useful lemmas.

**Lemma 1.** *Under the above assumptions, the gradient of function  $\Phi(x) = \max_{y \in \mathcal{Y}} f(x, y)$  is  $G$ -Lipschitz with respect to retraction, and the mapping or function  $y^*(x) = \arg \max_{y \in \mathcal{Y}} f(x, y)$  is  $\kappa$ -Lipschitz with respect to retraction. Given any  $x_1, x_2 \in \mathcal{X} \subseteq \mathcal{M}$  and  $u \in T_{x_1}\mathcal{M}$ , we have:*

$$\|\text{grad}\Phi(x_1) - \mathcal{T}_{x_2}^{x_1} \text{grad}\Phi(x_2)\| \leq G\|u\|, \quad (12)$$

$$\|y^*(x_1) - y^*(x_2)\| \leq \kappa\|u\|, \quad (13)$$

where  $x_2 = R_{x_1}(u)$ , and  $G = \kappa L_{12} + L_{11}$ , and  $\kappa = L_{21}/\mu$  denotes the number condition of function  $f(x, y)$ .

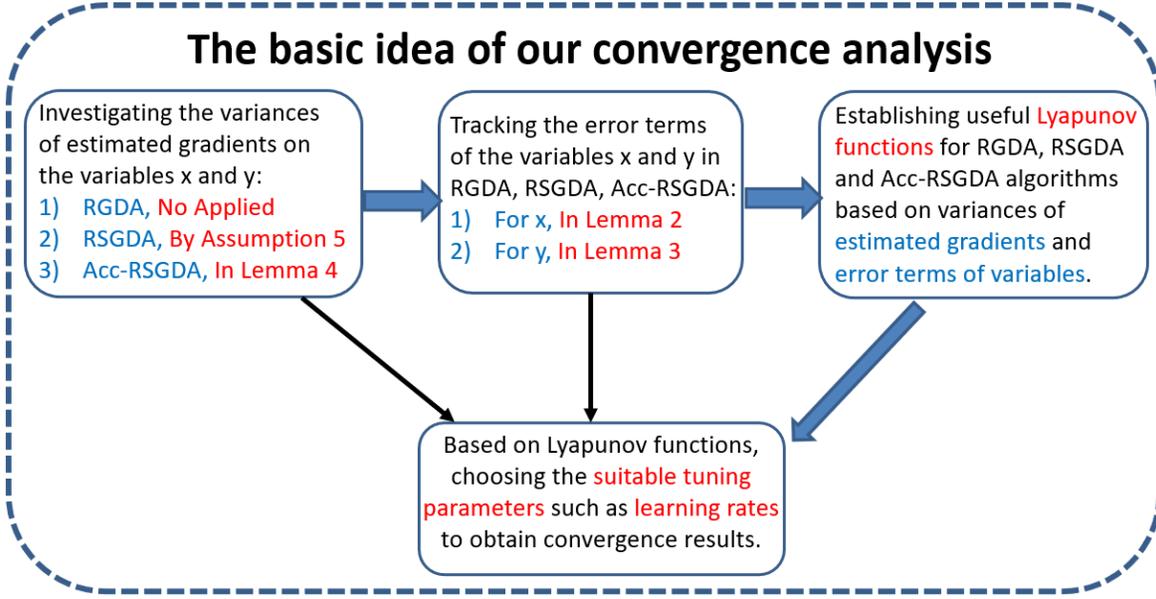


Fig. 2: The basic idea of our convergence analysis.

**Lemma 2.** Suppose the sequence  $\{x_t, y_t\}_{t=1}^T$  is generated from Algorithm 1 or 2. Given  $0 < \eta_t \leq \frac{1}{2\gamma L}$ , we have

$$\begin{aligned} \Phi(x_{t+1}) &\leq \Phi(x_t) + \gamma L_{12} \eta_t \|y^*(x_t) - y_t\|^2 - \frac{\gamma \eta_t}{4} \|v_t\|^2 \\ &\quad + \gamma \eta_t \|\text{grad}_x f(x_t, y_t) - v_t\|^2 - \frac{\gamma \eta_t}{2} \|\text{grad} \Phi(x_t)\|^2. \end{aligned} \quad (14)$$

**Lemma 3.** Suppose the sequence  $\{x_t, y_t\}_{t=1}^T$  is generated from Algorithm 1 or 2. Under the above assumptions, and set  $0 < \eta_t \leq 1$  and  $0 < \lambda \leq \frac{1}{6L}$ , we have

$$\begin{aligned} &\|y_{t+1} - y^*(x_{t+1})\|^2 \\ &\leq \left(1 - \frac{\eta_t \mu \lambda}{4}\right) \|y_t - y^*(x_t)\|^2 - \frac{3\eta_t}{4} \|\tilde{y}_{t+1} - y_t\|^2 \\ &\quad + \frac{25\eta_t \lambda}{6\mu} \|\nabla_y f(x_t, y_t) - w_t\|^2 + \frac{25\gamma^2 \kappa^2 \eta_t}{6\mu \lambda} \|v_t\|^2, \end{aligned} \quad (15)$$

where  $\kappa = L_{21}/\mu$  and  $\tilde{L} = \max(1, L_{11}, L_{12}, L_{21}, L_{22})$ .

Although Problems (1) and (2) are nonconvex, following [53], there exists a local solution or stationary point  $(x^*, y^*)$  satisfies the Nash Equilibrium, i.e.,  $f(x^*, y) \leq f(x^*, y^*) \leq f(x, y^*)$ , where  $x^* \in \mathcal{X} \subset \mathcal{M}$  and  $y^* \in \mathcal{Y}$ . Here  $\mathcal{X}$  is a neighbourhood around an optimal point  $x^*$ . Recall that the nonconvex minimax problem (1) is equivalent to minimizing the nonconvex function  $\Phi(x) = \max_{y \in \mathcal{Y}} f(x, y)$  for any  $x \in \mathcal{M}$ . It is NP hard to find the global minimum of  $\Phi(x)$  in general since  $\Phi(x)$  is nonconvex in  $x \in \mathcal{M}$ . Thus, we will find the stationary points of function  $\Phi(x)$ , which is equal to the stationary points of the minimax problem (1). Next we define an  $\epsilon$ -stationary point of  $\Phi(x)$  in  $x \in \mathcal{M}$ .

**Definition 1.** A point  $x \in \mathcal{M}$  is an  $\epsilon$ -stationary point ( $\epsilon \geq 0$ ) of a differentiable function  $\Phi(x)$  if  $\|\text{grad} \Phi(x)\| \leq \epsilon$ . If  $\epsilon = 0$ , then  $x$  is a stationary point.

## 5.1 Convergence Analysis of both RGDA and RSGDA Algorithms

In this subsection, we study the convergence properties of our RGDA and RSGDA algorithms, respectively.

Suppose the sequence  $\{x_t, y_t\}_{t=1}^T$  be generated from our RGDA Algorithm, we establish a useful *Lyapunov function* (i.e., potential function)  $\Lambda_t$  for convergence analysis of RGDA, defined as

$$\Lambda_t = \Phi(x_t) + \frac{6\gamma \tilde{L}^2}{\lambda \mu} \|y_t - y^*(x_t)\|^2, \quad \forall t \geq 1. \quad (16)$$

**Theorem 1.** Suppose the sequence  $\{x_t, y_t\}_{t=1}^T$  is generated from Algorithm 1 by using **deterministic** gradients. Given  $y_1 = y^*(x_1)$ ,  $\eta = \eta_t$  for all  $t \geq 1$ ,  $0 < \eta \leq \min(1, \frac{1}{2\gamma L})$ ,  $0 < \lambda \leq \frac{1}{6L}$  and  $0 < \gamma \leq \frac{\mu \lambda}{10L\kappa}$ , we have

$$\frac{1}{T} \sum_{t=1}^T \|\text{grad} \Phi(x_t)\| \leq \frac{2\sqrt{\Phi(x_1) - \Phi^*}}{\sqrt{\gamma \eta T}}. \quad (17)$$

**Remark 1.** Since  $0 < \eta \leq \min(1, \frac{1}{2\gamma L})$  and  $0 < \gamma \leq \frac{\mu \lambda}{10L\kappa}$ , we have  $0 < \eta \gamma \leq \min(\frac{\mu \lambda}{10L\kappa}, \frac{1}{2L})$ . Let  $\eta \gamma = \min(\frac{\mu \lambda}{10L\kappa}, \frac{1}{2L})$ , we have  $\eta \gamma = O(\frac{1}{\kappa^2})$ . The RGDA algorithm has convergence rate of  $O(\frac{\kappa}{T^{1/2}})$ . By  $\frac{\kappa}{T^{1/2}} \leq \epsilon$ , i.e.,  $\|\text{grad} \Phi(x_t)\| \leq \epsilon$ , we choose  $T \geq \kappa^2 \epsilon^{-2}$ . When our RGDA Algorithm solves the deterministic minimax Problem (1), we only need one sample to estimate the gradients  $v_t$  and  $w_t$  at each iteration, and need  $T$  iterations. Thus, our RGDA reaches a sample complexity of  $T = O(\kappa^2 \epsilon^{-2})$  for finding an  $\epsilon$ -stationary point of Problem (1). Note that since the function  $f(x, y)$  is  $\mu$ -strongly concave in  $y \in \mathcal{Y}$ , given any initial input  $x_1$ , we can easily obtain  $y_1 \approx y^*(x_1)$ . So we can assume  $y_1 = y^*(x_1)$ .

Suppose the sequence  $\{x_t, y_t\}_{t=1}^T$  be generated from our RSGDA Algorithm, we establish a useful *Lyapunov function*  $\Theta_t$  for convergence analysis of RSGDA, defined as

$$\Theta_t = \mathbb{E}[\Phi(x_t) + \frac{6\gamma \tilde{L}^2}{\lambda \mu} \|y_t - y^*(x_t)\|^2], \quad \forall t \geq 1. \quad (18)$$

**Theorem 2.** Suppose the sequence  $\{x_t, y_t\}_{t=1}^T$  is generated from Algorithm 1 by using **stochastic** gradients. Given  $y_1 = y^*(x_1)$ ,

$\eta = \eta_t$  for all  $t \geq 1$ ,  $0 < \eta \leq \min(1, \frac{1}{2\gamma L})$ ,  $0 < \lambda \leq \frac{1}{6L}$  and  $0 < \gamma \leq \frac{\mu\lambda}{10L\kappa}$ , we have

$$\frac{1}{T} \sum_{t=1}^T \mathbb{E} \|\text{grad } \Phi(x_t)\| \leq \frac{2\sqrt{\Phi(x_1) - \Phi^*}}{\sqrt{\gamma\eta T}} + (1 + \frac{5\tilde{L}}{\mu}) \frac{\sqrt{2}\sigma}{\sqrt{B}}. \quad (19)$$

**Remark 2.** Since  $0 < \eta \leq \min(1, \frac{1}{2\gamma L})$  and  $0 < \gamma \leq \frac{\mu\lambda}{10L\kappa}$ , we have  $0 < \eta\gamma \leq \min(\frac{\mu\lambda}{10L\kappa}, \frac{1}{2L})$ . Let  $\eta\gamma = \min(\frac{\mu\lambda}{10L\kappa}, \frac{1}{2L})$ , we have  $\eta\gamma = O(\frac{1}{\kappa^2})$ . Let  $B = T$ , the RSGDA algorithm has convergence rate of  $O(\frac{1}{T^{1/2}})$ . By  $\frac{\kappa}{T^{1/2}} \leq \epsilon$ , i.e.,  $\mathbb{E} \|\text{grad } \Phi(x_\zeta)\| \leq \epsilon$ , we choose  $T \geq \kappa^2 \epsilon^{-2}$ . When our RGDA Algorithm solves the stochastic minimax Problem (2), we need  $B$  samples to estimate the gradients  $v_t$  and  $w_t$  at each iteration, and need  $T$  iterations. Thus, the RSGDA reaches a sample complexity of  $BT = O(\kappa^4 \epsilon^{-4})$  for finding an  $\epsilon$ -stationary point of Problem (2).

## 5.2 Convergence Analysis of Acc-RSGDA Algorithm

In the subsection, we provide the convergence properties of our Acc-RSGDA algorithm.

**Lemma 4.** Suppose the stochastic gradients  $v_t$  and  $w_t$  is generated from Algorithm 2, given  $0 < \alpha_{t+1} \leq 1$  and  $0 < \beta_{t+1} \leq 1$ , we have

$$\begin{aligned} \mathbb{E} \|\text{grad}_x f(x_{t+1}, y_{t+1}) - v_{t+1}\|^2 &\leq 4(1 - \alpha_{t+1})^2 L_{11}^2 \gamma^2 \eta_t^2 \mathbb{E} \|v_t\|^2 \\ &+ (1 - \alpha_{t+1})^2 \mathbb{E} \|\text{grad}_x f(x_t, y_t) - v_t\|^2 \\ &+ 4(1 - \alpha_{t+1})^2 L_{12}^2 \eta_t^2 \mathbb{E} \|\tilde{y}_{t+1} - y_t\|^2 + \frac{2\alpha_{t+1}^2 \sigma^2}{B}, \end{aligned} \quad (20)$$

$$\begin{aligned} \mathbb{E} \|\nabla_y f(x_{t+1}, y_{t+1}) - w_{t+1}\|^2 &\leq 4(1 - \beta_{t+1})^2 L_{21}^2 \gamma^2 \eta_t^2 \mathbb{E} \|v_t\|^2 \\ &+ (1 - \beta_{t+1})^2 \mathbb{E} \|\nabla_y f(x_t, y_t) - w_t\|^2 \\ &+ 4(1 - \beta_{t+1})^2 L_{22}^2 \eta_t^2 \mathbb{E} \|\tilde{y}_{t+1} - y_t\|^2 + \frac{2\beta_{t+1}^2 \sigma^2}{B}. \end{aligned} \quad (21)$$

Assume the sequence  $\{x_t, y_t\}_{t=1}^T$  be generated from our Acc-RSGDA Algorithm, we establish a useful Lyapunov function  $\Omega_t$  for convergence analysis of Acc-RSGDA, defined as

$$\begin{aligned} \Omega_t &= \mathbb{E} [\Phi(x_t) + \frac{\gamma}{2\lambda\mu\eta_{t-1}} (\|\text{grad}_x f(x_t, y_t) - v_t\|^2 \\ &+ \|\nabla_y f(x_t, y_t) - w_t\|^2) + \frac{6\gamma\tilde{L}^2}{\lambda\mu} \|y_t - y^*(x_t)\|^2], \quad \forall t \geq 1. \end{aligned} \quad (22)$$

**Theorem 3.** Suppose the sequence  $\{x_t, y_t\}_{t=1}^T$  is generated from Algorithm 2. Given  $y_1 = y^*(x_1)$ ,  $c_1 \geq \frac{2}{3b^3} + 2\lambda\mu$ ,  $c_2 \geq \frac{2}{3b^3} + \frac{50\lambda\tilde{L}^2}{\mu}$ ,  $b > 0$ ,  $m \geq \max(2, (\tilde{c}b)^3)$ ,  $0 < \gamma \leq \frac{\mu\lambda}{2\kappa\tilde{L}\sqrt{25+4\mu\lambda}}$  and  $0 < \lambda \leq \frac{1}{6L}$ , we have

$$\frac{1}{T} \sum_{t=1}^T \mathbb{E} \|\text{grad } \Phi(x_t)\| \leq \frac{\sqrt{2M'}m^{1/6}}{T^{1/2}} + \frac{\sqrt{2M'}}{T^{1/3}}, \quad (23)$$

where  $\tilde{c} = \max(1, c_1, c_2, 2\gamma L)$  and  $M' = \frac{2(\Phi(x_1) - \Phi^*)}{\gamma^b} + \frac{2\sigma^2}{B\lambda\mu\eta_0 b} + \frac{2(c_1^2 + c_2^2)\sigma^2 b^2}{B\lambda\mu} \ln(m + T)$ .

**Remark 3.** Let  $c_1 = \frac{2}{3b^3} + 2\lambda\mu$ ,  $c_2 = \frac{2}{3b^3} + \frac{50\lambda\tilde{L}^2}{\mu}$ ,  $\lambda = \frac{1}{6L}$ ,  $\gamma = \frac{\mu\lambda}{2\kappa\tilde{L}\sqrt{25+4\mu\lambda}}$  and  $\eta_0 = \frac{b}{m^{1/3}}$ . It is easily verified that  $\gamma = O(\frac{1}{\kappa^2})$ ,  $\lambda = O(1)$ ,  $\lambda\mu = O(\frac{1}{\kappa})$ ,  $c_1 = O(1)$ ,  $c_2 = O(\kappa)$ ,

datasets	#samples	#dimension	#classes
MNIST	60,000	28 × 28	10
FashionMNIST	60,000	28 × 28	10
STL-10 (resized)	5,000	32 × 32 × 3	10
CIFAR-10	50,000	32 × 32 × 3	10

TABLE 1: Benchmark datasets used in our experiments

Inputs ( $d$ channels)
Conv $d \rightarrow 32$ , Batchnorm, ReLU
Conv $32 \rightarrow 64$ , Batchnorm, ReLU
Conv $64 \rightarrow 64$ , Batchnorm, ReLU
Max Pool
Linear $200 \rightarrow 200$ , ReLU
Linear $200 \rightarrow C$
Outputs

TABLE 2: The DNN used in our experiments.  $C$  is the number of classes, and  $d$  is the number of channels for inputs.

$m = O(\kappa^3)$  and  $\eta_0 = O(\frac{1}{\kappa})$ . Without loss of generality, let  $T \geq m = O(\kappa^3)$ , we have  $M' = O(\kappa^2 + \frac{\kappa^2}{B} + \frac{\kappa^3}{B} \ln(T))$ . When  $B = \kappa$ , we have  $M' = O(\kappa^2 \ln(T))$ . Thus, the Acc-RSGDA algorithm has a convergence rate of  $\tilde{O}(\frac{\kappa}{T^{1/3}})$ . By  $\frac{\kappa}{T^{1/3}} \leq \epsilon$ , i.e.,  $\mathbb{E} \|\text{grad } \Phi(x_\zeta)\| \leq \epsilon$ , we choose  $T \geq \kappa^3 \epsilon^{-3}$ . In Algorithm 2, we require  $B$  samples to estimate the stochastic gradients  $v_t$  and  $w_t$  at each iteration, and need  $T$  iterations. Thus, the Acc-RSGDA has a sample complexity of  $BT = \tilde{O}(\kappa^4 \epsilon^{-3})$  for finding an  $\epsilon$ -stationary point of Problem (2). **Since our Acc-RSGDA algorithm uses variance-reduced technique of STORM to estimate the stochastic gradients, it does not rely on large mini-batch size to guarantee its convergence.** When  $B = 1$ , our Acc-RSGDA algorithm has a convergence rate of  $\tilde{O}(\frac{\kappa^{3/2}}{T^{1/3}})$ , and has a sample complexity of  $BT = \tilde{O}(\kappa^{4.5} \epsilon^{-3})$  for finding an  $\epsilon$ -stationary point.

**Remark 4.** In the above theoretical analysis, we only assume the convexity of constraint set  $\mathcal{Y}$ , while [27] not only assume the convexity of set  $\mathcal{Y}$ , but also assume and use **its bounded** (i.e.,  $|\mathcal{Y}| \leq D$ , where  $D$  is a positive constant.) to guarantee convergence of the GDA and SGDA algorithms in [27] (Please see Assumption 4.2 in [27]). Clearly, our assumption is milder than [27]. When there does not exist a constraint set on parameter  $y$ , i.e.,  $\mathcal{Y} = \mathbb{R}^d$ , our RGDA and RSGDA algorithms and theoretical results still work, while [27] can not work.

## 6 EXPERIMENTS

In this section, we conduct experiments on two tasks: 1) robust DNNs training over Riemannian manifold and distributionally robust optimization over Riemannian manifold. In the experiment, we use the SGDA [27] and Acc-MDA [39] as the comparison baselines. Since the SGDA and Acc-MDA methods are not designed for optimization on Riemannian manifolds, we add the retraction operation (projection-like) at the end of parameter updates.

### 6.1 Robust DNNs Training

In this subsection, we focus on the robust DNNs training over Riemannian manifold defined in Problem (4), which

TABLE 3: Test accuracy against nature images and different attacks for **MNIST**. All comparison methods are test against PGD<sup>40</sup> [54]: PGD attack of 40 steps, and FGSM [55] attacks.

Methods	Nat. Img.	PGD <sup>40</sup> $L_\infty$				FGSM $L_\infty$			
		$\epsilon=0.1$	$\epsilon=0.2$	$\epsilon=0.3$	$\epsilon=0.4$	$\epsilon=0.1$	$\epsilon=0.2$	$\epsilon=0.3$	$\epsilon=0.4$
SGDA	98.94%	85.95%	82.10%	75.64%	61.95%	91.20%	89.06%	85.67%	78.01%
Acc-MDA	99.23%	86.12%	82.15%	75.22%	58.06%	92.25%	90.29%	87.11%	79.56%
RSGDA	99.22%	87.47%	84.17%	78.61%	64.92%	93.05%	91.26%	87.47%	80.51%
Acc-RSGDA	99.37%	<b>90.08%</b>	<b>87.29%</b>	<b>82.65%</b>	<b>73.66%</b>	<b>93.38%</b>	<b>91.67%</b>	<b>88.83%</b>	<b>82.81%</b>

TABLE 4: Test accuracy against nature images and different attacks for **FashionMNIST**. All comparison methods are test against PGD<sup>40</sup> and FGSM attacks.

Methods	Nat. Img.	PGD <sup>40</sup> $L_\infty$				FGSM $L_\infty$			
		$\epsilon=0.05$	$\epsilon=0.1$	$\epsilon=0.15$	$\epsilon=0.2$	$\epsilon=0.05$	$\epsilon=0.1$	$\epsilon=0.15$	$\epsilon=0.2$
SGDA	82.30%	69.12%	66.92%	64.83%	62.55%	73.83%	72.65%	71.65%	70.64%
Acc-MDA	83.89%	68.41%	65.86%	63.32%	60.76%	73.11%	71.77%	70.54%	69.05%
RSGDA	83.23%	70.23%	67.84%	65.57%	63.48%	75.85%	74.95%	74.33%	73.97%
Acc-RSGDA	84.15%	<b>71.03%</b>	<b>68.99%</b>	<b>66.07%</b>	<b>64.35%</b>	<b>76.01%</b>	<b>75.44%</b>	<b>75.08%</b>	<b>74.44%</b>

TABLE 5: Test accuracy against nature images and different attacks for **CIFAR10**. All comparison methods are tested against PGD<sup>40</sup> and FGSM attacks.

Methods	Nat. Img.	PGD <sup>40</sup> $L_\infty$				FGSM $L_\infty$			
		$\epsilon=0.005$	$\epsilon=0.01$	$\epsilon=0.015$	$\epsilon=0.02$	$\epsilon=0.005$	$\epsilon=0.01$	$\epsilon=0.015$	$\epsilon=0.02$
SGDA	64.73%	38.52%	33.89%	28.97%	22.50%	41.66%	37.92%	34.04%	28.75%
Acc-MDA	69.01%	42.78%	38.08%	32.84%	25.76%	46.18%	42.41%	38.42%	32.87%
RCG	64.64%	39.73%	35.47%	30.97%	24.67%	42.71%	39.28%	35.83%	31.11%
Acc-RSGDA	72.18%	<b>46.36%</b>	<b>41.70%</b>	<b>36.46%</b>	<b>29.16%</b>	<b>49.41%</b>	<b>45.62%</b>	<b>41.55%</b>	<b>35.94%</b>

TABLE 6: Test accuracy against nature images and different attacks for the first 10 classes of **CIFAR100**. All comparison methods are tested against PGD<sup>40</sup> and FGSM attacks.

Methods	Nat. Img.	PGD <sup>40</sup> $L_\infty$				FGSM $L_\infty$			
		$\epsilon=0.005$	$\epsilon=0.01$	$\epsilon=0.015$	$\epsilon=0.02$	$\epsilon=0.005$	$\epsilon=0.01$	$\epsilon=0.015$	$\epsilon=0.02$
SGDA	69.60%	45.62%	41.47%	37.65%	32.30%	47.48%	43.87%	40.40%	36.10%
Acc-MDA	70.70%	47.42%	43.20%	38.40%	31.90%	49.73%	46.10%	42.15%	37.10%
RCG	69.30%	47.27%	43.10%	38.60%	31.80%	49.98%	46.67%	43.25%	37.40%
Acc-RSGDA	70.10%	<b>48.12%</b>	<b>44.27%</b>	<b>39.60%</b>	<b>33.50%</b>	<b>50.52%</b>	<b>47.30%</b>	<b>43.60%</b>	<b>38.70%</b>

is a nonconvex and nonconcave minimax problem. Following [28], we cast the original robust training problem into the following nonconvex-(strongly)-concave problem:

$$\min_{x \in \mathcal{M}} \max_{u \in \mathcal{U}} \frac{1}{n} \sum_{i=1}^n \sum_{j=1}^C u_j \ell(h(a_{ij}^K; x), b_i) - r(u), \quad (24)$$

$$\text{s.t. } \mathcal{U} = \{u \in \mathbb{R}^C \mid u \geq 0, \|u\|_1 = 1\},$$

where  $a_{ij}^K$  is the permuted sample after  $K$  iterations of Projected Gradient Descent (PGD) [54] attack, and  $C$  is the number of classes for the dataset. Here  $r(u)$  is a (strongly) convex regularization term, e.g.,  $r(u) = \alpha \|u - 1/C\|^2$  or KL divergence  $r(u) = \alpha \sum_{i=1}^C u_i \log(u_i C)$ , where  $\alpha \geq 0$  is a tuning parameter. In the experiment, we use Stiefel manifold  $\mathcal{M} = \text{St}(r, d) = \{X \in \mathbb{R}^{d \times r} : X^T X = I_r\}$  on parameters  $x$  of DNNs (convolution layers and linear layers).

For robust training, we choose five datasets for this experiment: MNIST, FashionMNIST, CIFAR10, CIFAR100 and STL10. We use a 5 layer DNN as the target model, whose architecture is given in Tab. 2. For five datasets, we set  $\{\gamma, \lambda, \eta_t\} = \{0.1, 0.01, 0.1\}$  for RSGDA. For SGDA, we set the learning rates of both maximization and minimization as 0.01. For Acc-RSGDA, we set  $\{\gamma, \lambda, b, m, c_1, c_2\} =$

$\{1.0, 0.1, 0.5, 8, 512, 512\}$ , and we apply the same hyperparameters to Acc-MDA to ensure a fair comparison. We set  $K = 3$  for five datasets, and  $\epsilon$  for the robust training is set to 0.4, 0.2, 0.02, 0.02 and 0.02 for MNIST, FashionMNIST, CIFAR10, CIFAR100 and STL10 separately. We further set the mini-batch size as 512, and the model is trained for 200 epochs.

The training progress for robust training is shown in Fig. 3. From Fig. 3, we can see that our Acc-RSGDA method converges faster than the other comparison baselines, and it can achieve the best test accuracy with natural images for both datasets. RSGDA does not use momentum terms, but it reaches lower training loss compared to SGDA and Acc-MDA. This observation implies that our framework better utilizes the property of Riemannian manifold for robust DNN training. On the other hand, simply adding the retraction operation (Acc-MDA and SGDA) can not achieve the same effect.

The numeric results against different attacks (i.e., PGD attack [54] and Fast Gradient Sign Method (FGSM) attack [55]) are shown in Tab. 3, Tab. 4, Tab. 5, Tab. 6 and Tab. 7. Specifically, in the training progress, we report the numeric

TABLE 7: Test accuracy against nature images and different attacks for STL10. All comparison methods are tested against PGD<sup>40</sup> and FGSM attacks.

Methods	Nat. Img.	PGD <sup>40</sup> $L_\infty$				FGSM $L_\infty$			
		$\epsilon=0.005$	$\epsilon=0.01$	$\epsilon=0.015$	$\epsilon=0.02$	$\epsilon=0.005$	$\epsilon=0.01$	$\epsilon=0.015$	$\epsilon=0.02$
SGDA	51.24%	26.28%	22.53%	18.81%	14.12%	28.32%	25.06%	21.88%	17.81%
Acc-MDA	51.94%	28.22%	24.39%	20.66%	16.04%	30.28%	27.01%	23.93%	19.88%
RCG	51.86%	28.62%	24.83%	20.96%	16.54%	30.80%	27.56%	24.34%	20.31%
Acc-RSGDA	52.51%	<b>29.48%</b>	<b>25.64%</b>	<b>21.76%</b>	<b>16.76%</b>	<b>31.48%</b>	<b>28.20%</b>	<b>24.91%</b>	<b>20.69%</b>

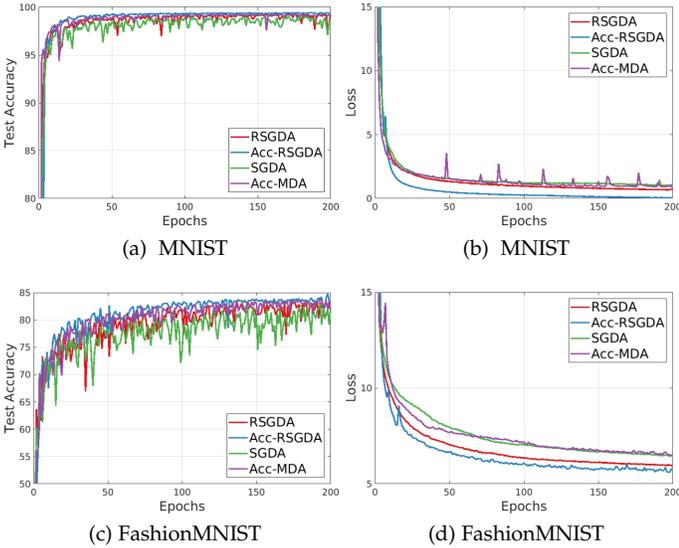


Fig. 3: Experimental results for the robust training task. (a, c) Test accuracy with natural images for MNIST and FashionMNIST datasets. (b, d) Training loss for MNIST and FashionMNIST datasets.

results against PGD attack of 40 steps and FGSM attack. For all settings, our Acc-RSGDA method achieves the best accuracy against PGD and FGSM attacks. Interestingly, the Acc-MDA method performs worse than SGDA under PGD and FGSM attacks, which suggests that the momentum may be not functional properly without considering the property of Riemannian manifold.

### 6.2 Distributionally Robust Optimization

In the subsection, we focus on distributionally robust optimization over Riemannian manifold defined in Problem (6). CIFAR-10 and STL-10 are selected as the datasets for this task. We use the same DNN architecture from the above robust DNN training for this task. We also apply Stiefel manifold  $\mathcal{M} = \text{St}(r, d) = \{X \in \mathbb{R}^{d \times r} : X^T X = I_r\}$  to the parameters of the DNN. We use the same hyper-parameter setting for RSGDA, Acc-RSGDA, SGDA and Acc-MDA from this task. The mini-batch size is also set 512, and the model is trained for 200 epochs. We report mean and variance across 3 runs for this experiment.

The results are reported in Fig. 4, and shaded areas represent variance. From the figure, we can see that our Acc-RSGDA achieves the best test accuracy and converges fastest. The difference between Acc-RSGDA and Acc-MDA is small, but due to using the property of Riemannian manifold, Acc-RSGDA is more stable compared to Acc-MDA.

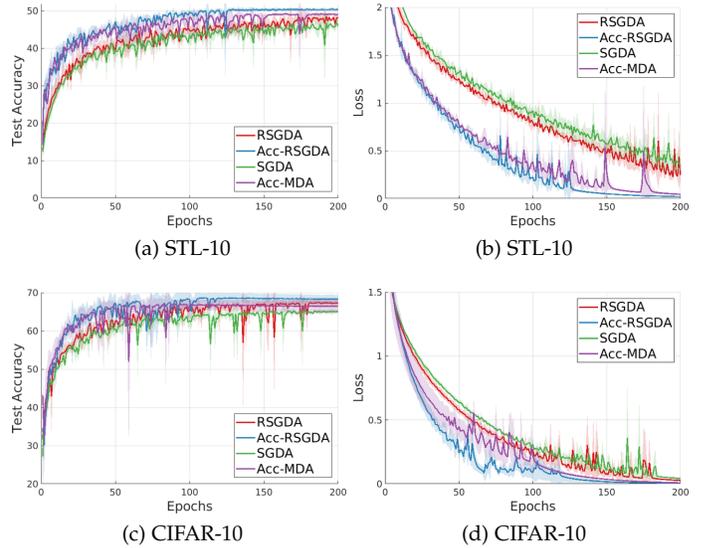


Fig. 4: Experimental results for the distributionally robust optimization task. (a, c) Test accuracy for STL-10 and CIFAR-10 datasets. (b, d) Training loss for STL-10 and CIFAR-10 datasets.

## 7 CONCLUSION

In the paper, we investigated a class of useful minimax optimization problems on Riemannian manifolds. Meanwhile, we proposed a class of effective and efficient Riemannian gradient descent ascent algorithms to solve these minimax problems. Moreover, we studied convergence properties of our proposed algorithms. To the best of our knowledge, our Riemannian gradient-based methods are the first to study the minimax optimization over the **general** Riemannian manifolds.

## ACKNOWLEDGMENTS

We thank the anonymous reviewers for their helpful comments. We also thank so much for the help of Prof. Heng Huang. This work was partially supported by NSFC under Grant No. 61806093. Feihu Huang is the corresponding author.

## REFERENCES

- [1] S.-M. Moosavi-Dezfooli, A. Fawzi, O. Fawzi, and P. Frossard, "Universal adversarial perturbations," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 1765–1773.
- [2] N. Akhtar, J. Liu, and A. Mian, "Defense against universal adversarial perturbations," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 3389–3398.

- [3] Y. Sun, L. Zheng, W. Deng, and S. Wang, "Svdnet for pedestrian retrieval," in *Proceedings of the IEEE International Conference on Computer Vision*, 2017, pp. 3800–3808.
- [4] D. Xie, J. Xiong, and S. Pu, "All you need is beyond a good init: Exploring better solution for training extremely deep convolutional neural networks with orthonormality and modulation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 6176–6185.
- [5] J. Li, L. Fuxin, and S. Todorovic, "Efficient riemannian optimization on the stiefel manifold via the cayley transform," in *International Conference on Learning Representations*, 2020.
- [6] N. Bansal, X. Chen, and Z. Wang, "Can we gain more from orthogonality regularizations in training deep networks?" in *Advances in Neural Information Processing Systems*, 2018, pp. 4261–4271.
- [7] M. Cogswell, F. Ahmed, R. Girshick, L. Zitnick, and D. Batra, "Reducing overfitting in deep networks by decorrelating representations," *arXiv preprint arXiv:1511.06068*, 2015.
- [8] L. Huang, X. Liu, B. Lang, A. W. Yu, Y. Wang, and B. Li, "Orthogonal weight normalization: Solution to optimization over multiple dependent stiefel manifolds in deep neural networks," in *32nd AAAI Conference on Artificial Intelligence, AAAI 2018*. AAAI Press, 2018, pp. 3271–3278.
- [9] R. S. Chen, B. Lucier, Y. Singer, and V. Syrgkanis, "Robust optimization for non-convex objectives," in *Advances in Neural Information Processing Systems*, 2017, pp. 4705–4714.
- [10] H. Rahimian and S. Mehrotra, "Distributionally robust optimization: A review," *arXiv preprint arXiv:1908.05659*, 2019.
- [11] A. Han and J. Gao, "Riemannian stochastic recursive momentum method for non-convex optimization," *arXiv preprint arXiv:2008.04555*, 2020.
- [12] J. Sun, Q. Qu, and J. Wright, "Complete dictionary recovery over the sphere ii: Recovery by riemannian trust-region method," *IEEE Transactions on Information Theory*, vol. 63, no. 2, pp. 885–914, 2016.
- [13] P. Jawanpuria and B. Mishra, "A unified framework for structured low-rank matrix learning," in *International Conference on Machine Learning*, 2018, pp. 2254–2263.
- [14] S. Mao, L. Xiong, L. Jiao, T. Feng, and S.-K. Yeung, "A novel riemannian metric based on riemannian structure and scaling information for fixed low-rank matrix completion," *IEEE Transactions on Cybernetics*, vol. 47, no. 5, pp. 1299–1312, 2016.
- [15] B. Vandereycken, "Low-rank matrix completion by riemannian optimization," *SIAM Journal on Optimization*, vol. 23, no. 2, pp. 1214–1236, 2013.
- [16] O. P. Ferreira, L. L. Pérez, and S. Z. Németh, "Singularities of monotone vector fields and an extragradient-type algorithm," *Journal of Global Optimization*, vol. 31, no. 1, pp. 133–151, 2005.
- [17] C. Li, G. López, and V. Martín-Márquez, "Monotone vector fields and the proximal point algorithm on hadamard manifolds," *Journal of the London Mathematical Society*, vol. 79, no. 3, pp. 663–683, 2009.
- [18] J. Wang, G. López, V. Martín-Márquez, and C. Li, "Monotone and accretive vector fields on riemannian manifolds," *Journal of optimization theory and applications*, vol. 146, no. 3, pp. 691–708, 2010.
- [19] T. Lin, C. Fan, N. Ho, M. Cuturi, and M. Jordan, "Projection robust wasserstein distance and riemannian optimization," *Advances in Neural Information Processing Systems*, vol. 33, pp. 9383–9397, 2020.
- [20] M. Huang, S. Ma, and L. Lai, "A riemannian block coordinate descent method for computing the projection robust wasserstein distance," in *International Conference on Machine Learning*. PMLR, 2021, pp. 4446–4455.
- [21] J. Weed and F. Bach, "Sharp asymptotic and finite-sample rates of convergence of empirical measures in wasserstein distance," *Bernoulli*, vol. 25, no. 4A, pp. 2620–2648, 2019.
- [22] A. Cutkosky and F. Orabona, "Momentum-based variance reduction in non-convex sgd," in *Advances in Neural Information Processing Systems*, 2019, pp. 15210–15219.
- [23] M. Razaviyayn, T. Huang, S. Lu, M. Nouiehed, M. Sanjabi, and M. Hong, "Nonconvex min-max optimization: Applications, challenges, and recent theoretical advances," *IEEE Signal Processing Magazine*, vol. 37, no. 5, pp. 55–66, 2020.
- [24] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative adversarial nets," in *Advances in neural information processing systems*, 2014, pp. 2672–2680.
- [25] K. Zhang, S. M. Kakade, T. Başar, and L. F. Yang, "Model-based multi-agent rl in zero-sum markov games with near-optimal sample complexity," *arXiv preprint arXiv:2007.07461*, 2020.
- [26] Y. Deng, M. M. Kamani, and M. Mahdavi, "Distributionally robust federated averaging," *Advances in Neural Information Processing Systems*, vol. 33, 2020.
- [27] T. Lin, C. Jin, and M. Jordan, "On gradient descent ascent for nonconvex-concave minimax problems," in *International Conference on Machine Learning*. PMLR, 2020, pp. 6083–6093.
- [28] M. Nouiehed, M. Sanjabi, T. Huang, J. D. Lee, and M. Razaviyayn, "Solving a class of non-convex min-max games using iterative first order methods," in *Advances in Neural Information Processing Systems*, 2019, pp. 14934–14942.
- [29] K. K. Thekumparampil, P. Jain, P. Netrapalli, and S. Oh, "Efficient algorithms for smooth minimax optimization," in *Advances in Neural Information Processing Systems*, 2019, pp. 12680–12691.
- [30] T. Lin, C. Jin, and M. I. Jordan, "Near-optimal algorithms for minimax optimization," in *Conference on Learning Theory*. PMLR, 2020, pp. 2738–2779.
- [31] S. Lu, I. Tsaknakis, M. Hong, and Y. Chen, "Hybrid block successive approximation for one-sided non-convex min-max problems: algorithms and applications," *IEEE Transactions on Signal Processing*, vol. 68, pp. 3676–3691, 2020.
- [32] J. Yang, N. Kiyavash, and N. He, "Global convergence and variance reduction for a class of nonconvex-nonconcave minimax problems," *Advances in Neural Information Processing Systems*, vol. 33, pp. 1153–1165, 2020.
- [33] J. Yang, S. Zhang, N. Kiyavash, and N. He, "A catalyst framework for minimax optimization," *Advances in Neural Information Processing Systems*, vol. 33, pp. 5667–5678, 2020.
- [34] S. Zhang, J. Yang, C. Guzmán, N. Kiyavash, and N. He, "The complexity of nonconvex-strongly-concave minimax optimization," in *Uncertainty in Artificial Intelligence*. PMLR, 2021, pp. 482–492.
- [35] Y. Yan, Y. Xu, Q. Lin, W. Liu, and T. Yang, "Sharp analysis of epoch stochastic gradient descent ascent methods for min-max optimization," *arXiv preprint arXiv:2022.05309*, 2020.
- [36] L. Luo, H. Ye, Z. Huang, and T. Zhang, "Stochastic recursive gradient descent ascent for stochastic nonconvex-strongly-concave minimax problems," *Advances in Neural Information Processing Systems*, vol. 33, 2020.
- [37] Z. Chen, Y. Zhou, T. Xu, and Y. Liang, "Proximal gradient descent-ascent: Variable convergence under kl geometry," in *Proc. International Conference on Learning Representations (ICLR)*, 2021.
- [38] F. Huang, X. Wu, and H. Huang, "Efficient mirror descent ascent methods for nonsmooth minimax problems," *Advances in Neural Information Processing Systems*, vol. 34, 2021.
- [39] F. Huang, S. Gao, J. Pei, and H. Huang, "Accelerated zeroth-order and first-order momentum methods from mini to minimax optimization," *Journal of Machine Learning Research*, vol. 23, no. 36, pp. 1–70, 2022.
- [40] H. Sakai and H. Iiduka, "Riemannian adaptive optimization algorithm and its application to natural language processing," *IEEE Transactions on Cybernetics*, 2021.
- [41] H. Zhang and S. Sra, "First-order methods for geodesically convex optimization," in *Conference on Learning Theory*, 2016, pp. 1617–1638.
- [42] Y. Liu, F. Shang, J. Cheng, H. Cheng, and L. Jiao, "Accelerated first-order methods for geodesically convex optimization on riemannian manifolds," in *Advances in Neural Information Processing Systems*, 2017, pp. 4868–4877.
- [43] H. Zhang, S. J. Reddi, and S. Sra, "Riemannian svrg: Fast stochastic optimization on riemannian manifolds," in *Advances in Neural Information Processing Systems*, 2016, pp. 4592–4600.
- [44] H. Sato, H. Kasai, and B. Mishra, "Riemannian stochastic variance reduced gradient algorithm with retraction and vector transport," *SIAM Journal on Optimization*, vol. 29, no. 2, pp. 1444–1472, 2019.
- [45] J. Zhang, H. Zhang, and S. Sra, "R-spider: A fast riemannian stochastic optimization algorithm with curvature independent rate," *arXiv preprint arXiv:1811.04194*, 2018.
- [46] H. Kasai, H. Sato, and B. Mishra, "Riemannian stochastic recursive gradient algorithm," in *International Conference on Machine Learning*, 2018, pp. 2516–2524.
- [47] P. Zhou, X.-T. Yuan, S. Yan, and J. Feng, "Faster first-order methods for stochastic non-convex optimization on riemannian manifolds," *IEEE transactions on pattern analysis and machine intelligence*, vol. 43, no. 2, pp. 459–472, 2021.

- [48] A. Han and J. Gao, "Improved variance reduction methods for riemannian non-convex optimization," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2021.
- [49] H. Kasai, P. Jawanpuria, and B. Mishra, "Riemannian adaptive stochastic gradient algorithms on matrix manifolds," in *International Conference on Machine Learning*. PMLR, 2019, pp. 3262–3271.
- [50] F. Huang and S. Gao, "Riemannian gradient methods for stochastic composition problems," *Neural Networks*, vol. 153, pp. 224–234, 2022.
- [51] W. Huang, "Optimization algorithms on riemannian manifolds with applications," Ph.D. dissertation, The Florida State University, 2013.
- [52] C. Fang, C. J. Li, Z. Lin, and T. Zhang, "Spider: Near-optimal non-convex optimization via stochastic path-integrated differential estimator," in *Advances in Neural Information Processing Systems*, 2018, pp. 689–699.
- [53] J. Von Neumann and O. Morgenstern, *Theory of games and economic behavior*. Princeton university press, 2007.
- [54] A. Kurakin, I. Goodfellow, and S. Bengio, "Adversarial machine learning at scale," *International Conference on Learning Representations*, 2017.
- [55] I. J. Goodfellow, J. Shlens, and C. Szegedy, "Explaining and harnessing adversarial examples," *arXiv preprint arXiv:1412.6572*, 2014.
- [56] Y. Nesterov, *Lectures on convex optimization*. Springer, 2018, vol. 137.

## APPENDIX A

### DETAILED PROOFS IN CONVERGENCE ANALYSIS

In this section, we provide the detailed convergence analysis of our algorithms. We first review some useful lemmas.

**Lemma 5.** [56] Assume that  $f(x)$  is a differentiable convex function and  $\mathcal{X}$  is a convex set.  $x^* \in \mathcal{X}$  is the solution of the constrained problem  $\min_{x \in \mathcal{X}} f(x)$ , if

$$\langle \nabla f(x^*), x - x^* \rangle \geq 0, \quad \forall x \in \mathcal{X}. \quad (25)$$

**Lemma 6.** [56] Assume the function  $f(x)$  is  $L$ -smooth, i.e.,  $\|\nabla f(x) - \nabla f(y)\| \leq L\|x - y\|$ , and then the following inequality holds

$$|f(y) - f(x) - \nabla f(x)^T(y - x)| \leq \frac{L}{2}\|x - y\|^2. \quad (26)$$

**Lemma 7.** (Restatement of Lemma 1) The gradient of function  $\Phi(x) = \max_{y \in \mathcal{Y}} f(x, y)$  is retraction  $G$ -Lipschitz, and the mapping or function  $y^*(x) = \arg \max_{y \in \mathcal{Y}} f(x, y)$  is retraction  $\kappa$ -Lipschitz. Given any  $x_1, x_2 = R_{x_1}(u) \in \mathcal{X} \subset \mathcal{M}$  and  $u \in T_{x_1}\mathcal{M}$ , we have

$$\begin{aligned} \|\text{grad}\Phi(x_1) - \mathcal{T}_{x_2}^{x_1} \text{grad}\Phi(x_2)\| &\leq G\|u\|, \\ \|y^*(x_1) - y^*(x_2)\| &\leq \kappa\|u\|, \end{aligned}$$

where  $G = \kappa L_{12} + L_{11}$  and  $\kappa = L_{21}/\mu$ , and vector transport  $\mathcal{T}_{x_2}^{x_1}$  transport the tangent space of  $x_1$  to that of  $x_2$ .

*Proof.* Given any  $x_1, x_2 = R_{x_1}(u) \in \mathcal{X}$  and  $u \in T_{x_1}\mathcal{M}$ , define  $y^*(x_1) = \arg \max_{y \in \mathcal{Y}} f(x_1, y)$  and  $y^*(x_2) = \arg \max_{y \in \mathcal{Y}} f(x_2, y)$ , by the above Lemma 5, we have

$$(y - y^*(x_1))^T \nabla_y f(x_1, y^*(x_1)) \leq 0, \quad \forall y \in \mathcal{Y} \quad (27)$$

$$(y - y^*(x_2))^T \nabla_y f(x_2, y^*(x_2)) \leq 0, \quad \forall y \in \mathcal{Y}. \quad (28)$$

Let  $y = y^*(x_2)$  in the inequality (27) and  $y = y^*(x_1)$  in the inequality (28), then summing these inequalities, we have

$$(y^*(x_2) - y^*(x_1))^T (\nabla_y f(x_1, y^*(x_1)) - \nabla_y f(x_2, y^*(x_2))) \leq 0. \quad (29)$$

Since the function  $f(x_1, \cdot)$  is  $\mu$ -strongly concave, we have

$$f(x_1, y^*(x_1)) \leq f(x_1, y^*(x_2)) + (\nabla_y f(x_1, y^*(x_2)))^T (y^*(x_1) - y^*(x_2)) - \frac{\mu}{2}\|y^*(x_1) - y^*(x_2)\|^2, \quad (30)$$

$$f(x_1, y^*(x_2)) \leq f(x_1, y^*(x_1)) + (\nabla_y f(x_1, y^*(x_1)))^T (y^*(x_2) - y^*(x_1)) - \frac{\mu}{2}\|y^*(x_1) - y^*(x_2)\|^2. \quad (31)$$

Combining the inequalities (30) with (31), we obtain

$$(y^*(x_2) - y^*(x_1))^T (\nabla_y f(x_1, y^*(x_2)) - \nabla_y f(x_1, y^*(x_1))) + \mu\|y^*(x_1) - y^*(x_2)\|^2 \leq 0. \quad (32)$$

By plugging the inequalities (29) into (32), we have

$$\begin{aligned} \mu\|y^*(x_1) - y^*(x_2)\|^2 &\leq (y^*(x_2) - y^*(x_1))^T (\nabla_y f(x_2, y^*(x_2)) - \nabla_y f(x_1, y^*(x_2))) \\ &\leq \|y^*(x_2) - y^*(x_1)\| \|\nabla_y f(x_2, y^*(x_2)) - \nabla_y f(x_1, y^*(x_2))\| \\ &\leq L_{21}\|u\| \|y^*(x_2) - y^*(x_1)\|, \end{aligned} \quad (33)$$

where the last inequality is due to Assumption 1. Thus, we have

$$\|y^*(x_1) - y^*(x_2)\| \leq \kappa\|u\|, \quad (34)$$

where  $\kappa = L_{21}/\mu$  and  $x_2 = R_{x_1}(u)$ ,  $u \in T_{x_1}\mathcal{M}$ .

Since  $\Phi(x) = f(x, y^*(x))$ , we have  $\text{grad}\Phi(x) = \text{grad}_x f(x, y^*(x))$ . Then we have

$$\begin{aligned} &\|\text{grad}\Phi(x_1) - \mathcal{T}_{x_2}^{x_1} \text{grad}\Phi(x_2)\| \\ &= \|\text{grad}_x f(x_1, y^*(x_1)) - \mathcal{T}_{x_2}^{x_1} \text{grad}_x f(x_2, y^*(x_2))\| \\ &\leq \|\text{grad}_x f(x_1, y^*(x_1)) - \text{grad}_x f(x_1, y^*(x_2))\| + \|\text{grad}_x f(x_1, y^*(x_2)) - \mathcal{T}_{x_2}^{x_1} \text{grad}_x f(x_2, y^*(x_2))\| \\ &\leq L_{12}\|y^*(x_1) - y^*(x_2)\| + L_{11}\|u\| \\ &\leq (\kappa L_{12} + L_{11})\|u\|, \end{aligned} \quad (35)$$

where  $u \in T_{x_1}\mathcal{M}$ . □

**Lemma 8.** (Restatement of Lemma 2) Suppose the sequence  $\{x_t, y_t\}_{t=1}^T$  is generated from Algorithm 1 or 2. Given  $0 < \eta_t \leq \frac{1}{2\gamma L}$ , we have

$$\begin{aligned} \Phi(x_{t+1}) &\leq \Phi(x_t) + \gamma L_{12} \eta_t \|y^*(x_t) - y_t\|^2 + \gamma \eta_t \|\text{grad}_x f(x_t, y_t) - v_t\|^2 - \frac{\gamma \eta_t}{2} \|\text{grad} \Phi(x_t)\|^2 \\ &\quad - \frac{\gamma \eta_t}{4} \|v_t\|^2. \end{aligned} \quad (36)$$

*Proof.* According to Assumption 2, i.e., the function  $\Phi(x)$  is retraction  $L$ -smooth, we have

$$\begin{aligned} \Phi(x_{t+1}) &\leq \Phi(x_t) - \gamma \eta_t \langle \text{grad} \Phi(x_t), v_t \rangle + \frac{\gamma^2 \eta_t^2 L}{2} \|v_t\|^2 \\ &= \Phi(x_t) + \frac{\gamma \eta_t}{2} \|\text{grad} \Phi(x_t) - v_t\|^2 - \frac{\gamma \eta_t}{2} \|\text{grad} \Phi(x_t)\|^2 + \left(\frac{\gamma^2 \eta_t^2 L}{2} - \frac{\gamma \eta_t}{2}\right) \|v_t\|^2 \\ &= \Phi(x_t) + \frac{\gamma \eta_t}{2} \|\text{grad} \Phi(x_t) - \text{grad}_x f(x_t, y_t) + \text{grad}_x f(x_t, y_t) - v_t\|^2 - \frac{\gamma \eta_t}{2} \|\text{grad} \Phi(x_t)\|^2 \\ &\quad + \left(\frac{\gamma^2 \eta_t^2 L}{2} - \frac{\gamma \eta_t}{2}\right) \|v_t\|^2 \\ &\leq \Phi(x_t) + \gamma \eta_t \|\text{grad} \Phi(x_t) - \text{grad}_x f(x_t, y_t)\|^2 + \gamma \eta_t \|\text{grad}_x f(x_t, y_t) - v_t\|^2 - \frac{\gamma \eta_t}{2} \|\text{grad} \Phi(x_t)\|^2 \\ &\quad + \left(\frac{L \gamma^2 \eta_t^2}{2} - \frac{\gamma \eta_t}{2}\right) \|v_t\|^2 \\ &\leq \Phi(x_t) + \gamma \eta_t \|\text{grad} \Phi(x_t) - \text{grad}_x f(x_t, y_t)\|^2 + \gamma \eta_t \|\text{grad}_x f(x_t, y_t) - v_t\|^2 - \frac{\gamma \eta_t}{2} \|\text{grad} \Phi(x_t)\|^2 \\ &\quad - \frac{\gamma \eta_t}{4} \|v_t\|^2, \end{aligned} \quad (37)$$

where the last inequality is due to  $0 < \eta_t \leq \frac{1}{2\gamma L}$ .

Considering an upper bound of  $\|\text{grad} \Phi(x_t) - \text{grad}_x f(x_t, y_t)\|^2$ , we have

$$\begin{aligned} \|\text{grad} \Phi(x_t) - \text{grad}_x f(x_t, y_t)\|^2 &= \|\text{grad}_x f(x_t, y^*(x_t)) - \text{grad}_x f(x_t, y_t)\|^2 \\ &\leq L_{12} \|y^*(x_t) - y_t\|^2. \end{aligned} \quad (38)$$

Then we have

$$\begin{aligned} \Phi(x_{t+1}) &\leq \Phi(x_t) + \gamma \eta_t L_{12} \|y^*(x_t) - y_t\|^2 + \gamma \eta_t \|\text{grad}_x f(x_t, y_t) - v_t\|^2 - \frac{\gamma \eta_t}{2} \|\text{grad} \Phi(x_t)\|^2 \\ &\quad - \frac{\gamma \eta_t}{4} \|v_t\|^2. \end{aligned} \quad (39)$$

□

**Lemma 9.** (Restatement of Lemma 3) Suppose the sequence  $\{x_t, y_t\}_{t=1}^T$  is generated from Algorithm 1 or 2. Under the above assumptions, and set  $0 < \eta_t \leq 1$  and  $0 < \lambda \leq \frac{1}{6\tilde{L}}$ , we have

$$\begin{aligned} \|y_{t+1} - y^*(x_{t+1})\|^2 &\leq \left(1 - \frac{\eta_t \mu \lambda}{4}\right) \|y_t - y^*(x_t)\|^2 - \frac{3\eta_t}{4} \|\tilde{y}_{t+1} - y_t\|^2 \\ &\quad + \frac{25\eta_t \lambda}{6\mu} \|\nabla_y f(x_t, y_t) - w_t\|^2 + \frac{25\gamma^2 \kappa^2 \eta_t}{6\mu \lambda} \|v_t\|^2, \end{aligned} \quad (40)$$

where  $\kappa = L_{21}/\mu$  and  $\tilde{L} = \max(1, L_{11}, L_{12}, L_{21}, L_{22})$ .

*Proof.* Since the constraint set  $\mathcal{Y}$  in Euclidean space, this proof can easily follow the proofs of Lemma 28 in [39]. According to Assumption 3, i.e., the function  $f(x, y)$  is  $\mu$ -strongly concave w.r.t  $y$ , we have

$$\begin{aligned} f(x_t, y) &\leq f(x_t, y_t) + \langle \nabla_y f(x_t, y_t), y - y_t \rangle - \frac{\mu}{2} \|y - y_t\|^2 \\ &= f(x_t, y_t) + \langle w_t, y - \tilde{y}_{t+1} \rangle + \langle \nabla_y f(x_t, y_t) - w_t, y - \tilde{y}_{t+1} \rangle \\ &\quad + \langle \nabla_y f(x_t, y_t), \tilde{y}_{t+1} - y_t \rangle - \frac{\mu}{2} \|y - y_t\|^2. \end{aligned} \quad (41)$$

According to the assumption 1, i.e., the function  $f(x, y)$  is  $L_{22}$ -smooth w.r.t  $y$ , and  $\tilde{L} \geq L_{22}$ , we have

$$\begin{aligned} f(x_t, \tilde{y}_{t+1}) - f(x_t, y_t) - \langle \nabla_y f(x_t, y_t), \tilde{y}_{t+1} - y_t \rangle &\geq -\frac{L_{22}}{2} \|\tilde{y}_{t+1} - y_t\|^2 \\ &\geq -\frac{\tilde{L}}{2} \|\tilde{y}_{t+1} - y_t\|^2. \end{aligned} \quad (42)$$

Combining the inequalities (41) with (42), we have

$$\begin{aligned} f(x_t, y) &\leq f(x_t, \tilde{y}_{t+1}) + \langle w_t, y - \tilde{y}_{t+1} \rangle + \langle \nabla_y f(x_t, y_t) - w_t, y - \tilde{y}_{t+1} \rangle \\ &\quad - \frac{\mu}{2} \|y - y_t\|^2 + \frac{\tilde{L}}{2} \|\tilde{y}_{t+1} - y_t\|^2. \end{aligned} \quad (43)$$

According to the line 6 of Algorithm 1 or 2, we have  $\tilde{y}_{t+1} = \mathcal{P}_{\mathcal{Y}}(y_t + \lambda w_t) = \arg \min_{y \in \mathcal{Y}} \frac{1}{2} \|y - y_t - \lambda w_t\|^2$ . Since  $\mathcal{Y}$  is a convex set and the function  $\frac{1}{2} \|y - y_t - \lambda w_t\|^2$  is convex, according to Lemma 5, we have

$$\langle \tilde{y}_{t+1} - y_t - \lambda w_t, y - \tilde{y}_{t+1} \rangle \geq 0, \quad y \in \mathcal{Y}. \quad (44)$$

Then we obtain

$$\begin{aligned} \langle w_t, y - \tilde{y}_{t+1} \rangle &\leq \frac{1}{\lambda} \langle \tilde{y}_{t+1} - y_t, y - \tilde{y}_{t+1} \rangle \\ &= \frac{1}{\lambda} \langle \tilde{y}_{t+1} - y_t, y_t - \tilde{y}_{t+1} \rangle + \frac{1}{\lambda} \langle \tilde{y}_{t+1} - y_t, y - y_t \rangle \\ &= -\frac{1}{\lambda} \|\tilde{y}_{t+1} - y_t\|^2 + \frac{1}{\lambda} \langle \tilde{y}_{t+1} - y_t, y - y_t \rangle. \end{aligned} \quad (45)$$

Combining the inequalities (43) with (45), we have

$$\begin{aligned} f(x_t, y) &\leq f(x_t, \tilde{y}_{t+1}) + \frac{1}{\lambda} \langle \tilde{y}_{t+1} - y_t, y - y_t \rangle + \langle \nabla_y f(x_t, y_t) - w_t, y - \tilde{y}_{t+1} \rangle \\ &\quad - \frac{1}{\lambda} \|\tilde{y}_{t+1} - y_t\|^2 - \frac{\mu}{2} \|y - y_t\|^2 + \frac{\tilde{L}}{2} \|\tilde{y}_{t+1} - y_t\|^2. \end{aligned} \quad (46)$$

Let  $y = y^*(x_t)$  and we obtain

$$\begin{aligned} f(x_t, y^*(x_t)) &\leq f(x_t, \tilde{y}_{t+1}) + \frac{1}{\lambda} \langle \tilde{y}_{t+1} - y_t, y^*(x_t) - y_t \rangle + \langle \nabla_y f(x_t, y_t) - w_t, y^*(x_t) - \tilde{y}_{t+1} \rangle \\ &\quad - \frac{1}{\lambda} \|\tilde{y}_{t+1} - y_t\|^2 - \frac{\mu}{2} \|y^*(x_t) - y_t\|^2 + \frac{\tilde{L}}{2} \|\tilde{y}_{t+1} - y_t\|^2. \end{aligned} \quad (47)$$

Due to the concavity of  $f(\cdot, y)$  and  $y^*(x_t) = \arg \max_{y \in \mathcal{Y}} f(x_t, y)$ , we have  $f(x_t, y^*(x_t)) \geq f(x_t, \tilde{y}_{t+1})$ . Thus, we obtain

$$\begin{aligned} 0 &\leq \frac{1}{\lambda} \langle \tilde{y}_{t+1} - y_t, y^*(x_t) - y_t \rangle + \langle \nabla_y f(x_t, y_t) - w_t, y^*(x_t) - \tilde{y}_{t+1} \rangle \\ &\quad - \left( \frac{1}{\lambda} - \frac{\tilde{L}}{2} \right) \|\tilde{y}_{t+1} - y_t\|^2 - \frac{\mu}{2} \|y^*(x_t) - y_t\|^2. \end{aligned} \quad (48)$$

By  $y_{t+1} = y_t + \eta_t(\tilde{y}_{t+1} - y_t)$ , we have

$$\begin{aligned} \|y_{t+1} - y^*(x_t)\|^2 &= \|y_t + \eta_t(\tilde{y}_{t+1} - y_t) - y^*(x_t)\|^2 \\ &= \|y_t - y^*(x_t)\|^2 + 2\eta_t \langle \tilde{y}_{t+1} - y_t, y_t - y^*(x_t) \rangle + \eta_t^2 \|\tilde{y}_{t+1} - y_t\|^2. \end{aligned} \quad (49)$$

Then we obtain

$$\langle \tilde{y}_{t+1} - y_t, y^*(x_t) - y_t \rangle \leq \frac{1}{2\eta_t} \|y_t - y^*(x_t)\|^2 + \frac{\eta_t}{2} \|\tilde{y}_{t+1} - y_t\|^2 - \frac{1}{2\eta_t} \|y_{t+1} - y^*(x_t)\|^2. \quad (50)$$

Consider the upper bound of the term  $\langle \nabla_y f(x_t, y_t) - w_t, y^*(x_t) - \tilde{y}_{t+1} \rangle$ , we have

$$\begin{aligned} &\langle \nabla_y f(x_t, y_t) - w_t, y^*(x_t) - \tilde{y}_{t+1} \rangle \\ &= \langle \nabla_y f(x_t, y_t) - w_t, y^*(x_t) - y_t \rangle + \langle \nabla_y f(x_t, y_t) - w_t, y_t - \tilde{y}_{t+1} \rangle \\ &\leq \frac{1}{\mu} \|\nabla_y f(x_t, y_t) - w_t\|^2 + \frac{\mu}{4} \|y^*(x_t) - y_t\|^2 + \frac{1}{\mu} \|\nabla_y f(x_t, y_t) - w_t\|^2 + \frac{\mu}{4} \|y_t - \tilde{y}_{t+1}\|^2 \\ &= \frac{2}{\mu} \|\nabla_y f(x_t, y_t) - w_t\|^2 + \frac{\mu}{4} \|y^*(x_t) - y_t\|^2 + \frac{\mu}{4} \|y_t - \tilde{y}_{t+1}\|^2. \end{aligned} \quad (51)$$

By plugging the inequalities (48), (50) to (51), we have

$$\begin{aligned}
 \frac{1}{2\eta_t\lambda}\|y_{t+1} - y^*(x_t)\|^2 &\leq \left(\frac{1}{2\eta_t\lambda} - \frac{\mu}{4}\right)\|y_t - y^*(x_t)\|^2 + \left(\frac{\eta_t}{2\lambda} + \frac{\mu}{4} + \frac{\tilde{L}}{2} - \frac{1}{\lambda}\right)\|\tilde{y}_{t+1} - y_t\|^2 \\
 &\quad + \frac{2}{\mu}\|\nabla_y f(x_t, y_t) - w_t\|^2 \\
 &\leq \left(\frac{1}{2\eta_t\lambda} - \frac{\mu}{4}\right)\|y_t - y^*(x_t)\|^2 + \left(\frac{3\tilde{L}}{4} - \frac{1}{2\lambda}\right)\|\tilde{y}_{t+1} - y_t\|^2 + \frac{2}{\mu}\|\nabla_y f(x_t, y_t) - w_t\|^2 \\
 &= \left(\frac{1}{2\eta_t\lambda} - \frac{\mu}{4}\right)\|y_t - y^*(x_t)\|^2 - \left(\frac{3}{8\lambda} + \frac{1}{8\lambda} - \frac{3\tilde{L}}{4}\right)\|\tilde{y}_{t+1} - y_t\|^2 \\
 &\quad + \frac{2}{\mu}\|\nabla_y f(x_t, y_t) - w_t\|^2 \\
 &\leq \left(\frac{1}{2\eta_t\lambda} - \frac{\mu}{4}\right)\|y_t - y^*(x_t)\|^2 - \frac{3}{8\lambda}\|\tilde{y}_{t+1} - y_t\|^2 + \frac{2}{\mu}\|\nabla_y f(x_t, y_t) - w_t\|^2, \tag{52}
 \end{aligned}$$

where the second inequality holds by  $\tilde{L} \geq L_{22} \geq \mu$  and  $0 < \eta_t \leq 1$ , and the last inequality is due to  $0 < \lambda \leq \frac{1}{6\tilde{L}}$ . It implies that

$$\|y_{t+1} - y^*(x_t)\|^2 \leq \left(1 - \frac{\eta_t\mu\lambda}{2}\right)\|y_t - y^*(x_t)\|^2 - \frac{3\eta_t}{4}\|\tilde{y}_{t+1} - y_t\|^2 + \frac{4\eta_t\lambda}{\mu}\|\nabla_y f(x_t, y_t) - w_t\|^2. \tag{53}$$

Next, we decompose the term  $\|y_{t+1} - y^*(x_{t+1})\|^2$  as follows:

$$\begin{aligned}
 \|y_{t+1} - y^*(x_{t+1})\|^2 &= \|y_{t+1} - y^*(x_t) + y^*(x_t) - y^*(x_{t+1})\|^2 \\
 &= \|y_{t+1} - y^*(x_t)\|^2 + 2\langle y_{t+1} - y^*(x_t), y^*(x_t) - y^*(x_{t+1}) \rangle + \|y^*(x_t) - y^*(x_{t+1})\|^2 \\
 &\leq \left(1 + \frac{\eta_t\mu\lambda}{4}\right)\|y_{t+1} - y^*(x_t)\|^2 + \left(1 + \frac{4}{\eta_t\mu\lambda}\right)\|y^*(x_t) - y^*(x_{t+1})\|^2 \\
 &\leq \left(1 + \frac{\eta_t\mu\lambda}{4}\right)\|y_{t+1} - y^*(x_t)\|^2 + \left(1 + \frac{4}{\eta_t\mu\lambda}\right)\eta_t^2\gamma^2\kappa^2\|v_t\|^2, \tag{54}
 \end{aligned}$$

where the first inequality holds by the Cauchy-Schwarz inequality and Young's inequality, and the last equality is due to Lemma 7.

By combining the above inequalities (53) and (54), we have

$$\begin{aligned}
 \|y_{t+1} - y^*(x_{t+1})\|^2 &\leq \left(1 + \frac{\eta_t\mu\lambda}{4}\right)\left(1 - \frac{\eta_t\mu\lambda}{2}\right)\|y_t - y^*(x_t)\|^2 - \left(1 + \frac{\eta_t\mu\lambda}{4}\right)\frac{3\eta_t}{4}\|\tilde{y}_{t+1} - y_t\|^2 \\
 &\quad + \left(1 + \frac{\eta_t\mu\lambda}{4}\right)\frac{4\eta_t\lambda}{\mu}\|\nabla_y f(x_t, y_t) - w_t\|^2 + \left(1 + \frac{4}{\eta_t\mu\lambda}\right)\eta_t^2\gamma^2\kappa^2\|v_t\|^2. \tag{55}
 \end{aligned}$$

Since  $0 < \eta_t \leq 1$ ,  $0 < \lambda \leq \frac{1}{6\tilde{L}}$  and  $\tilde{L} \geq L_{22} \geq \mu$ , we have  $\lambda \leq \frac{1}{6\tilde{L}} \leq \frac{1}{6\mu}$  and  $\eta_t \leq 1 \leq \frac{1}{6\mu\lambda}$ . Then we obtain

$$\begin{aligned}
 \left(1 + \frac{\eta_t\mu\lambda}{4}\right)\left(1 - \frac{\eta_t\mu\lambda}{2}\right) &= 1 - \frac{\eta_t\mu\lambda}{2} + \frac{\eta_t\mu\lambda}{4} - \frac{\eta_t^2\mu^2\lambda^2}{8} \leq 1 - \frac{\eta_t\mu\lambda}{4}, \\
 -\left(1 + \frac{\eta_t\mu\lambda}{4}\right)\frac{3\eta_t}{4} &\leq -\frac{3\eta_t}{4}, \\
 \left(1 + \frac{\eta_t\mu\lambda}{4}\right)\frac{4\eta_t\lambda}{\mu} &\leq \left(1 + \frac{1}{24}\right)\frac{4\eta_t\lambda}{\mu} = \frac{25\eta_t\lambda}{6\mu}, \\
 \left(1 + \frac{4}{\eta_t\mu\lambda}\right)\gamma^2\kappa^2\eta_t^2 &= \gamma^2\kappa^2\eta_t^2 + \frac{4\gamma^2\kappa^2\eta_t}{\mu\lambda} \leq \frac{\gamma^2\kappa^2\eta_t}{6\mu\lambda} + \frac{4\gamma^2\kappa^2\eta_t}{\mu\lambda} = \frac{25\gamma^2\kappa^2\eta_t}{6\mu\lambda}. \tag{56}
 \end{aligned}$$

Thus we have

$$\begin{aligned}
 \|y_{t+1} - y^*(x_{t+1})\|^2 &\leq \left(1 - \frac{\eta_t\mu\lambda}{4}\right)\|y_t - y^*(x_t)\|^2 - \frac{3\eta_t}{4}\|\tilde{y}_{t+1} - y_t\|^2 \\
 &\quad + \frac{25\eta_t\lambda}{6\mu}\|\nabla_y f(x_t, y_t) - w_t\|^2 + \frac{25\gamma^2\kappa^2\eta_t}{6\mu\lambda}\|v_t\|^2. \tag{57}
 \end{aligned}$$

□

### A.1 Convergence Analysis of RGDA and RSGDA Algorithms

In the subsection, we study the convergence properties of our RGDA and RSGDA algorithms, respectively.

**Theorem 4.** (Restatement of Theorem 1) Suppose the sequence  $\{x_t, y_t\}_{t=1}^T$  is generated from Algorithm 1 by using deterministic gradients. Given  $y_1 = y^*(x_1)$ ,  $\eta = \eta_t$  for all  $t \geq 1$ ,  $0 < \eta \leq \min(1, \frac{1}{2\gamma L})$ ,  $0 < \lambda \leq \frac{1}{6L}$  and  $0 < \gamma \leq \frac{\mu\lambda}{10L\kappa}$ , we have

$$\frac{1}{T} \sum_{t=1}^T \|\text{grad } \Phi(x_t)\| \leq \frac{2\sqrt{\Phi(x_1) - \Phi^*}}{\sqrt{\gamma\eta T}}, \quad (58)$$

where  $\tilde{L} = \max(1, L_{11}, L_{12}, L_{21}, L_{22})$ .

*Proof.* According to Lemma 9, we have

$$\begin{aligned} \|y_{t+1} - y^*(x_{t+1})\|^2 &\leq (1 - \frac{\eta_t\mu\lambda}{4})\|y_t - y^*(x_t)\|^2 - \frac{3\eta_t}{4}\|\tilde{y}_{t+1} - y_t\|^2 + \frac{25\eta_t\lambda}{6\mu}\|\nabla_y f(x_t, y_t) - w_t\|^2 \\ &\quad + \frac{25\gamma^2\kappa^2\eta_t}{6\mu\lambda}\|v_t\|^2. \end{aligned} \quad (59)$$

We first define a *Lyapunov* function  $\Lambda_t$ , for any  $t \geq 1$

$$\Lambda_t = \Phi(x_t) + \frac{6\gamma\tilde{L}^2}{\lambda\mu}\|y_t - y^*(x_t)\|^2. \quad (60)$$

According to Lemma 8, we have

$$\begin{aligned} \Lambda_{t+1} - \Lambda_t &= \Phi(x_{t+1}) - \Phi(x_t) + \frac{6\gamma\tilde{L}^2}{\lambda\mu}(\|y_{t+1} - y^*(x_{t+1})\|^2 - \|y_t - y^*(x_t)\|^2) \\ &\leq \gamma\eta_t L_{12}\|y_t - y^*(x_t)\|^2 + \gamma\eta_t\|\text{grad}_x f(x_t, y_t) - v_t\|^2 - \frac{\gamma\eta_t}{2}\|\text{grad}\Phi(x_t)\|^2 - \frac{\gamma\eta_t}{4}\|v_t\|^2 \\ &\quad + \frac{6\gamma\tilde{L}^2}{\lambda\mu}(-\frac{\mu\lambda\eta_t}{4}\|y_t - y^*(x_t)\|^2 - \frac{3\eta_t}{4}\|\tilde{y}_{t+1} - y_t\|^2 + \frac{25\lambda\eta_t}{6\mu}\|\nabla_y f(x_t, y_t) - w_t\|^2 \\ &\quad + \frac{25\gamma^2\kappa^2\eta_t}{6\mu\lambda}\|v_t\|^2) \\ &\leq -\frac{\tilde{L}^2\gamma\eta_t}{2}\|y_t - y^*(x_t)\|^2 - \frac{\gamma\eta_t}{2}\|\text{grad}\Phi(x_t)\|^2 - \frac{9\gamma\tilde{L}^2\eta_t}{2\lambda\mu}\|\tilde{y}_{t+1} - y_t\|^2 \\ &\quad - (\frac{1}{4} - \frac{25\kappa^2\tilde{L}^2\gamma^2}{\mu^2\lambda^2})\gamma\eta_t\|v_t\|^2 \\ &\leq -\frac{\tilde{L}^2\gamma\eta_t}{2}\|y_t - y^*(x_t)\|^2 - \frac{\gamma\eta_t}{2}\|\text{grad}\Phi(x_t)\|^2, \end{aligned} \quad (61)$$

where the first inequality holds by the inequality (59); the second last inequality is due to  $\tilde{L} = \max(1, L_{11}, L_{12}, L_{21}, L_{22})$  and  $v_t = \text{grad}_x f(x_t, y_t)$ ,  $w_t = \nabla_y f(x_t, y_t)$ , and the last inequality is due to  $0 < \gamma \leq \frac{\mu\lambda}{10L\kappa}$ . Thus, we obtain

$$\frac{\tilde{L}^2\gamma\eta_t}{2}\|y_t - y^*(x_t)\|^2 + \frac{\gamma\eta_t}{2}\|\text{grad}\Phi(x_t)\|^2 \leq \Lambda_t - \Lambda_{t+1}. \quad (62)$$

Since the initial solution satisfies  $y_1 = y^*(x_1) = \arg \max_{y \in \mathcal{Y}} f(x_1, y)$ , we have

$$\Lambda_1 = \Phi(x_1) + \frac{6\gamma\tilde{L}^2}{\lambda\mu}\|y_1 - y^*(x_1)\|^2 = \Phi(x_1). \quad (63)$$

Taking average over  $t = 1, 2, \dots, T$  on both sides of the inequality (62), we have

$$\frac{1}{T} \sum_{t=1}^T [\frac{\tilde{L}^2\eta_t}{2}\|y_t - y^*(x_t)\|^2 + \frac{\eta_t}{2}\|\text{grad}\Phi(x_t)\|^2] \leq \frac{\Lambda_1 - \Lambda_{T+1}}{\gamma T} \leq \frac{\Phi(x_1) - \Phi^*}{\gamma T}, \quad (64)$$

where the last equality is due to the above equality (63) and Assumption 4. Let  $\eta = \eta_1 = \dots = \eta_T$ , we have

$$\frac{1}{T} \sum_{t=1}^T [\tilde{L}^2\|y_t - y^*(x_t)\|^2 + \|\text{grad}\Phi(x_t)\|^2] \leq \frac{2(\Phi(x_1) - \Phi^*)}{\gamma\eta T}. \quad (65)$$

According to Jensen's inequality, we have

$$\begin{aligned} \frac{1}{T} \sum_{t=1}^T [\tilde{L}\|y_t - y^*(x_t)\| + \|\text{grad}\Phi(x_t)\|] &\leq \left( \frac{2}{T} \sum_{t=1}^T [\tilde{L}^2\|y_t - y^*(x_t)\|^2 + \|\text{grad}\Phi(x_t)\|^2] \right)^{1/2} \\ &\leq \left( \frac{4(\Phi(x_1) - \Phi^*)}{\gamma\eta T} \right)^{1/2} = \frac{2\sqrt{\Phi(x_1) - \Phi^*}}{\sqrt{\gamma\eta T}}. \end{aligned} \quad (66)$$

Since  $\tilde{L}\|y_t - y^*(x_t)\| + \|\text{grad}\Phi(x_t)\| \geq \|\text{grad}\Phi(x_t)\|$ , we can obtain

$$\frac{1}{T} \sum_{t=1}^T \|\text{grad}\Phi(x_t)\| \leq \frac{2\sqrt{\Phi(x_1) - \Phi^*}}{\sqrt{\gamma\eta T}}. \quad (67)$$

□

**Theorem 5.** (Restatement of Theorem 2) Suppose the sequence  $\{x_t, y_t\}_{t=1}^T$  is generated from Algorithm 1 by using stochastic gradients. Given  $y_1 = y^*(x_1)$ ,  $\eta = \eta_t$  for all  $t \geq 1$ ,  $0 < \eta \leq \min(1, \frac{1}{2\gamma\tilde{L}})$ ,  $0 < \lambda \leq \frac{1}{6\tilde{L}}$  and  $0 < \gamma \leq \frac{\mu\lambda}{10\tilde{L}\kappa}$ , we have

$$\frac{1}{T} \sum_{t=1}^T \mathbb{E}\|\text{grad}\Phi(x_t)\| \leq \frac{2\sqrt{\Phi(x_1) - \Phi^*}}{\sqrt{\gamma\eta T}} + \frac{\sqrt{2}\sigma}{\sqrt{B}} + \frac{5\sqrt{2}\tilde{L}\sigma}{\sqrt{B\mu}}. \quad (68)$$

*Proof.* According to Lemma 9, we have

$$\begin{aligned} \|y_{t+1} - y^*(x_{t+1})\|^2 &\leq \left(1 - \frac{\eta_t\mu\lambda}{4}\right)\|y_t - y^*(x_t)\|^2 - \frac{3\eta_t}{4}\|\tilde{y}_{t+1} - y_t\|^2 + \frac{25\eta_t\lambda}{6\mu}\|\nabla_y f(x_t, y_t) - w_t\|^2 \\ &\quad + \frac{25\gamma^2\kappa^2\eta_t}{6\mu\lambda}\|v_t\|^2. \end{aligned} \quad (69)$$

We first define a Lyapunov function  $\Theta_t$ , for any  $t \geq 1$

$$\Theta_t = \mathbb{E}[\Phi(x_t) + \frac{6\gamma\tilde{L}^2}{\lambda\mu}\|y_t - y^*(x_t)\|^2]. \quad (70)$$

By Assumption 5, we have

$$\mathbb{E}\|\text{grad}_x f(x_t, y_t) - v_t\|^2 = \mathbb{E}\|\text{grad}_x f(x_t, y_t) - \frac{1}{B} \sum_{i=1}^B \text{grad}_x f(x_t, y_t; \xi_t^i)\|^2 \leq \frac{\sigma^2}{B}, \quad (71)$$

$$\mathbb{E}\|\nabla_y f(x_t, y_t) - w_t\|^2 = \mathbb{E}\|\nabla_y f(x_t, y_t) - \frac{1}{B} \sum_{i=1}^B \nabla_y f(x_t, y_t; \xi_t^i)\|^2 \leq \frac{\sigma^2}{B}. \quad (72)$$

According to Lemma 8, we have

$$\begin{aligned} \Theta_{t+1} - \Theta_t &= \mathbb{E}[\Phi(x_{t+1})] - \mathbb{E}[\Phi(x_t)] + \frac{6\gamma\tilde{L}^2}{\lambda\mu}(\mathbb{E}\|y_{t+1} - y^*(x_{t+1})\|^2 - \mathbb{E}\|y_t - y^*(x_t)\|^2) \\ &\leq \gamma\eta_t L_{12} \mathbb{E}\|y_t - y^*(x_t)\|^2 + \gamma\eta_t \mathbb{E}\|\text{grad}_x f(x_t, y_t) - v_t\|^2 - \frac{\gamma\eta_t}{2} \mathbb{E}\|\text{grad}\Phi(x_t)\|^2 - \frac{\gamma\eta_t}{4} \mathbb{E}\|v_t\|^2 \\ &\quad + \frac{6\gamma\tilde{L}^2}{\lambda\mu} \left( -\frac{\mu\lambda\eta_t}{4} \mathbb{E}\|y_t - y^*(x_t)\|^2 - \frac{3\eta_t}{4} \mathbb{E}\|\tilde{y}_{t+1} - y_t\|^2 + \frac{25\lambda\eta_t}{6\mu} \mathbb{E}\|\nabla_y f(x_t, y_t) - w_t\|^2 \right) \\ &\quad + \frac{25\gamma^2\kappa^2\eta_t}{6\mu\lambda} \mathbb{E}\|v_t\|^2 \\ &\leq -\frac{\tilde{L}^2\gamma\eta_t}{2} \mathbb{E}\|y_t - y^*(x_t)\|^2 - \frac{\gamma\eta_t}{2} \mathbb{E}\|\text{grad}\Phi(x_t)\|^2 - \frac{9\gamma\tilde{L}^2\eta_t}{2\lambda\mu} \mathbb{E}\|\tilde{y}_{t+1} - y_t\|^2 \\ &\quad - \left( \frac{1}{4} - \frac{25\kappa^2\tilde{L}^2\gamma^2}{\mu^2\lambda^2} \right) \gamma\eta_t \mathbb{E}\|v_t\|^2 + \gamma\eta_t \mathbb{E}\|\text{grad}_x f(x_t, y_t) - v_t\|^2 + \frac{25\tilde{L}^2\gamma\eta_t}{\mu^2} \mathbb{E}\|\nabla_y f(x_t, y_t) - w_t\|^2 \\ &\leq -\frac{\tilde{L}^2\gamma\eta_t}{2} \mathbb{E}\|y_t - y^*(x_t)\|^2 - \frac{\gamma\eta_t}{2} \mathbb{E}\|\text{grad}\Phi(x_t)\|^2 + \frac{\gamma\eta_t\sigma^2}{B} + \frac{25\tilde{L}^2\gamma\eta_t\sigma^2}{B\mu^2}, \end{aligned} \quad (73)$$

where the first inequality holds by the inequality (69); the second last inequality is due to  $\tilde{L} = \max(1, L_{11}, L_{12}, L_{21}, L_{22})$ , and the last inequality is due to  $0 < \gamma \leq \frac{\mu\lambda}{10\tilde{L}\kappa}$  and Assumption 5. Thus, we obtain

$$\frac{\tilde{L}^2\gamma\eta_t}{2} \mathbb{E}\|y_t - y^*(x_t)\|^2 + \frac{\gamma\eta_t}{2} \mathbb{E}\|\text{grad}\Phi(x_t)\|^2 \leq \Theta_t - \Theta_{t+1} + \frac{\gamma\eta_t\sigma^2}{B} + \frac{25\tilde{L}^2\gamma\eta_t\sigma^2}{B\mu^2}. \quad (74)$$

Since the initial solution satisfies  $y_1 = y^*(x_1) = \arg \max_{y \in \mathcal{Y}} f(x_1, y)$ , we have

$$\Theta_1 = \Phi(x_1) + \frac{6\gamma\tilde{L}^2}{\lambda\mu} \|y_1 - y^*(x_1)\|^2 = \Phi(x_1). \quad (75)$$

Taking average over  $t = 1, 2, \dots, T$  on both sides of the inequality (74), we have

$$\begin{aligned} \frac{1}{T} \sum_{t=1}^T \mathbb{E} \left[ \frac{\tilde{L}^2 \eta_t}{2} \|y_t - y^*(x_t)\|^2 + \frac{\eta_t}{2} \|\text{grad}\Phi(x_t)\|^2 \right] &\leq \frac{\Theta_t - \Theta_{t+1}}{\gamma T} + \frac{1}{T} \sum_{t=1}^T \frac{\eta_t \sigma^2}{B} + \frac{1}{T} \sum_{t=1}^T \frac{25\tilde{L}^2 \eta_t \sigma^2}{B\mu^2} \\ &= \frac{\Phi(x_1) - \Phi^*}{\gamma T} + \frac{1}{T} \sum_{t=1}^T \frac{\eta_t \sigma^2}{B} + \frac{1}{T} \sum_{t=1}^T \frac{25\tilde{L}^2 \eta_t \sigma^2}{B\mu^2}, \end{aligned} \quad (76)$$

where the last equality is due to the above equality (75). Let  $\eta = \eta_1 = \dots = \eta_T$ , we have

$$\frac{1}{T} \sum_{t=1}^T \mathbb{E} [\tilde{L}^2 \|y_t - y^*(x_t)\|^2 + \|\text{grad}\Phi(x_t)\|^2] \leq \frac{2(\Phi(x_1) - \Phi^*)}{\gamma \eta T} + \frac{\sigma^2}{B} + \frac{25\tilde{L}^2 \sigma^2}{B\mu^2}. \quad (77)$$

According to Jensen's inequality, we have

$$\begin{aligned} \frac{1}{T} \sum_{t=1}^T \mathbb{E} [\tilde{L} \|y_t - y^*(x_t)\| + \|\text{grad}\Phi(x_t)\|] &\leq \left( \frac{2}{T} \sum_{t=1}^T \mathbb{E} [\tilde{L}^2 \|y_t - y^*(x_t)\|^2 + \|\text{grad}\Phi(x_t)\|^2] \right)^{1/2} \\ &\leq \left( \frac{4(\Phi(x_1) - \Phi^*)}{\gamma \eta T} + \frac{2\sigma^2}{B} + \frac{50\tilde{L}^2 \sigma^2}{B\mu^2} \right)^{1/2} \\ &\leq \frac{2\sqrt{\Phi(x_1) - \Phi^*}}{\sqrt{\gamma \eta T}} + \frac{\sqrt{2}\sigma}{\sqrt{B}} + \frac{5\sqrt{2}\tilde{L}\sigma}{\sqrt{B}\mu}, \end{aligned} \quad (78)$$

where the last inequality is due to  $(a_1 + a_2 + a_3)^{1/2} \leq a_1^{1/2} + a_2^{1/2} + a_3^{1/2}$  for all  $a_1, a_2, a_3 > 0$ . Thus, we have

$$\frac{1}{T} \sum_{t=1}^T \mathbb{E} \|\text{grad}\Phi(x_t)\| \leq \frac{2\sqrt{\Phi(x_1) - \Phi^*}}{\sqrt{\gamma \eta T}} + \frac{\sqrt{2}\sigma}{\sqrt{B}} + \frac{5\sqrt{2}\tilde{L}\sigma}{\sqrt{B}\mu}. \quad (79)$$

□

## A.2 Convergence Analysis of the Acc-RSGDA Algorithm

In the subsection, we study the convergence properties of the Acc-RSGDA algorithm.

**Lemma 10.** (Restatement of Lemma 4) Suppose the stochastic gradients  $v_t$  and  $w_t$  is generated from Algorithm 2, given  $0 < \alpha_{t+1} \leq 1$  and  $0 < \beta_{t+1} \leq 1$ , we have

$$\begin{aligned} \mathbb{E} \|\text{grad}_x f(x_{t+1}, y_{t+1}) - v_{t+1}\|^2 &\leq (1 - \alpha_{t+1})^2 \mathbb{E} \|\text{grad}_x f(x_t, y_t) - v_t\|^2 + 4(1 - \alpha_{t+1})^2 L_{11}^2 \gamma^2 \eta_t^2 \mathbb{E} \|v_t\|^2 \\ &\quad + 4(1 - \alpha_{t+1})^2 L_{12}^2 \eta_t^2 \mathbb{E} \|\tilde{y}_{t+1} - y_t\|^2 + \frac{2\alpha_{t+1}^2 \sigma^2}{B}. \end{aligned} \quad (80)$$

$$\begin{aligned} \mathbb{E} \|\nabla_y f(x_{t+1}, y_{t+1}) - w_{t+1}\|^2 &\leq (1 - \beta_{t+1})^2 \mathbb{E} \|\nabla_y f(x_t, y_t) - w_t\|^2 + 4(1 - \beta_{t+1})^2 L_{21}^2 \gamma^2 \eta_t^2 \mathbb{E} \|v_t\|^2 \\ &\quad + 4(1 - \beta_{t+1})^2 L_{22}^2 \eta_t^2 \mathbb{E} \|\tilde{y}_{t+1} - y_t\|^2 + \frac{2\beta_{t+1}^2 \sigma^2}{B}. \end{aligned} \quad (81)$$

*Proof.* We first prove the inequality (80). According to the definition of  $v_t$  in Algorithm 2, we have

$$\begin{aligned} v_{t+1} - \mathcal{T}_{x_t}^{x_{t+1}} v_t &= -\alpha_{t+1} \mathcal{T}_{x_t}^{x_{t+1}} v_t + (1 - \alpha_{t+1}) (\text{grad}_x f_{\mathcal{B}_{t+1}}(x_{t+1}, y_{t+1}) - \mathcal{T}_{x_t}^{x_{t+1}} \text{grad}_x f_{\mathcal{B}_{t+1}}(x_t, y_t)) \\ &\quad + \alpha_{t+1} \text{grad}_x f_{\mathcal{B}_{t+1}}(x_{t+1}, y_{t+1}). \end{aligned} \quad (82)$$

Then we have

$$\begin{aligned}
 & \mathbb{E}\|\text{grad}_x f(x_{t+1}, y_{t+1}) - v_{t+1}\|^2 \\
 &= \mathbb{E}\|\text{grad}_x f(x_{t+1}, y_{t+1}) - \mathcal{T}_{x_t}^{x_{t+1}} v_t - (v_{t+1} - \mathcal{T}_{x_t}^{x_{t+1}} v_t)\|^2 \\
 &= \mathbb{E}\|\text{grad}_x f(x_{t+1}, y_{t+1}) - \mathcal{T}_{x_t}^{x_{t+1}} v_t + \alpha_{t+1} \mathcal{T}_{x_t}^{x_{t+1}} v_t - \alpha_{t+1} \text{grad}_x f_{\mathcal{B}_{t+1}}(x_{t+1}, y_{t+1}) \\
 &\quad - (1 - \alpha_{t+1})(\text{grad}_x f_{\mathcal{B}_{t+1}}(x_{t+1}, y_{t+1}) - \mathcal{T}_{x_t}^{x_{t+1}} \text{grad}_x f_{\mathcal{B}_{t+1}}(x_t, y_t))\|^2 \\
 &= \mathbb{E}\|(1 - \alpha_{t+1})\mathcal{T}_{x_t}^{x_{t+1}}(\text{grad}_x f(x_t, y_t) - v_t) + (1 - \alpha_{t+1})(\text{grad}_x f(x_{t+1}, y_{t+1}) - \mathcal{T}_{x_t}^{x_{t+1}} \text{grad}_x f(x_t, y_t) \\
 &\quad - \text{grad}_x f_{\mathcal{B}_{t+1}}(x_{t+1}, y_{t+1}) + \mathcal{T}_{x_t}^{x_{t+1}} \text{grad}_x f_{\mathcal{B}_{t+1}}(x_t, y_t)) \\
 &\quad + \alpha_{t+1}(\text{grad}_x f(x_{t+1}, y_{t+1}) - \text{grad}_x f_{\mathcal{B}_{t+1}}(x_{t+1}, y_{t+1}))\|^2 \\
 &= (1 - \alpha_{t+1})^2 \mathbb{E}\|\text{grad}_x f(x_t, y_t) - v_t\|^2 + \alpha_{t+1}^2 \mathbb{E}\|\text{grad}_x f(x_{t+1}, y_{t+1}) - \text{grad}_x f_{\mathcal{B}_{t+1}}(x_{t+1}, y_{t+1})\|^2 \\
 &\quad + (1 - \alpha_{t+1})^2 \mathbb{E}\|\text{grad}_x f(x_{t+1}, y_{t+1}) - \mathcal{T}_{x_t}^{x_{t+1}} \text{grad}_x f(x_t, y_t) - \text{grad}_x f_{\mathcal{B}_{t+1}}(x_{t+1}, y_{t+1}) \\
 &\quad + \mathcal{T}_{x_t}^{x_{t+1}} \text{grad}_x f_{\mathcal{B}_{t+1}}(x_t, y_t)\|^2 + 2\alpha_{t+1}(1 - \alpha_{t+1})\langle \text{grad}_x f(x_{t+1}, y_{t+1}) - \mathcal{T}_{x_t}^{x_{t+1}} \text{grad}_x f(x_t, y_t) \\
 &\quad - \text{grad}_x f_{\mathcal{B}_{t+1}}(x_{t+1}, y_{t+1}) + \mathcal{T}_{x_t}^{x_{t+1}} \text{grad}_x f_{\mathcal{B}_{t+1}}(x_t, y_t), \text{grad}_x f(x_{t+1}, y_{t+1}) - \text{grad}_x f_{\mathcal{B}_{t+1}}(x_{t+1}, y_{t+1}) \rangle \\
 &\leq (1 - \alpha_{t+1})^2 \mathbb{E}\|\text{grad}_x f(x_t, y_t) - v_t\|^2 + 2\alpha_{t+1}^2 \mathbb{E}\|\text{grad}_x f(x_{t+1}, y_{t+1}) - \text{grad}_x f_{\mathcal{B}_{t+1}}(x_{t+1}, y_{t+1})\|^2 \\
 &\quad + 2(1 - \alpha_{t+1})^2 \mathbb{E}\|\text{grad}_x f(x_{t+1}, y_{t+1}) - \mathcal{T}_{x_t}^{x_{t+1}} \text{grad}_x f(x_t, y_t) - \text{grad}_x f_{\mathcal{B}_{t+1}}(x_{t+1}, y_{t+1}) \\
 &\quad + \mathcal{T}_{x_t}^{x_{t+1}} \text{grad}_x f_{\mathcal{B}_{t+1}}(x_t, y_t)\|^2 \\
 &\leq (1 - \alpha_{t+1})^2 \mathbb{E}\|\text{grad}_x f(x_t, y_t) - v_t\|^2 + \frac{2\alpha_{t+1}^2 \sigma^2}{B} \\
 &\quad + 2(1 - \alpha_{t+1})^2 \underbrace{\mathbb{E}\|\text{grad}_x f_{\mathcal{B}_{t+1}}(x_{t+1}, y_{t+1}) - \mathcal{T}_{x_t}^{x_{t+1}} \text{grad}_x f_{\mathcal{B}_{t+1}}(x_t, y_t)\|^2}_{=T_1},
 \end{aligned} \tag{83}$$

where the fourth equality follows by  $\mathbb{E}[\text{grad}_x f_{\mathcal{B}_{t+1}}(x_{t+1}, y_{t+1})] = \text{grad}_x f(x_{t+1}, y_{t+1})$  and  $\mathbb{E}[\text{grad}_x f_{\mathcal{B}_{t+1}}(x_{t+1}, y_{t+1}) - \text{grad}_x f_{\mathcal{B}_{t+1}}(x_t, y_t)] = \text{grad}_x f(x_{t+1}, y_{t+1}) - \text{grad}_x f(x_t, y_t)$ ; the first inequality holds by Young's inequality; the last inequality is due to the equality  $\mathbb{E}\|\zeta - \mathbb{E}[\zeta]\|^2 = \mathbb{E}\|\zeta\|^2 - \|\mathbb{E}[\zeta]\|^2$  and Assumption 5.

Next, we consider an upper bound of the above term  $T_1$  as follows:

$$\begin{aligned}
 T_1 &= \mathbb{E}\|\text{grad}_x f_{\mathcal{B}_{t+1}}(x_{t+1}, y_{t+1}) - \mathcal{T}_{x_t}^{x_{t+1}} \text{grad}_x f_{\mathcal{B}_{t+1}}(x_t, y_t)\|^2 \\
 &= \mathbb{E}\|\text{grad}_x f_{\mathcal{B}_{t+1}}(x_{t+1}, y_{t+1}) - \mathcal{T}_{x_t}^{x_{t+1}} \text{grad}_x f(x_t, y_{t+1}; \xi_{t+1}) + \mathcal{T}_{x_t}^{x_{t+1}} \text{grad}_x f(x_t, y_{t+1}; \xi_{t+1}) \\
 &\quad - \mathcal{T}_{x_t}^{x_{t+1}} \text{grad}_x f_{\mathcal{B}_{t+1}}(x_t, y_t)\|^2 \\
 &\leq 2\mathbb{E}\|\text{grad}_x f_{\mathcal{B}_{t+1}}(x_{t+1}, y_{t+1}) - \mathcal{T}_{x_t}^{x_{t+1}} \text{grad}_x f(x_t, y_{t+1}; \xi_{t+1})\|^2 \\
 &\quad + 2\mathbb{E}\|\text{grad}_x f(x_t, y_{t+1}; \xi_{t+1}) - \text{grad}_x f_{\mathcal{B}_{t+1}}(x_t, y_t)\|^2 \\
 &\leq 2L_{11}^2 \gamma^2 \eta_t^2 \mathbb{E}\|v_t\|^2 + 2L_{12}^2 \mathbb{E}\|y_{t+1} - y_t\|^2 \\
 &= 2L_{11}^2 \gamma^2 \eta_t^2 \mathbb{E}\|v_t\|^2 + 2L_{12}^2 \eta_t^2 \mathbb{E}\|\tilde{y}_{t+1} - y_t\|^2,
 \end{aligned} \tag{84}$$

where the last inequality is due to Assumption 1. Thus, we have

$$\begin{aligned}
 \mathbb{E}\|\text{grad}_x f(x_{t+1}, y_{t+1}) - v_{t+1}\|^2 &\leq (1 - \alpha_{t+1})^2 \mathbb{E}\|\text{grad}_x f(x_t, y_t) - v_t\|^2 + 4(1 - \alpha_{t+1})^2 L_{11}^2 \gamma^2 \eta_t^2 \mathbb{E}\|v_t\|^2 \\
 &\quad + 4(1 - \alpha_{t+1})^2 L_{12}^2 \eta_t^2 \mathbb{E}\|\tilde{y}_{t+1} - y_t\|^2 + \frac{2\alpha_{t+1}^2 \sigma^2}{B}.
 \end{aligned} \tag{85}$$

We apply a similar analysis to prove the above inequality (81). We obtain

$$\begin{aligned}
 \mathbb{E}\|\nabla_y f(x_{t+1}, y_{t+1}) - w_{t+1}\|^2 &\leq (1 - \beta_{t+1})^2 \mathbb{E}\|\nabla_y f(x_t, y_t) - w_t\|^2 + 4(1 - \beta_{t+1})^2 L_{21}^2 \gamma^2 \eta_t^2 \mathbb{E}\|v_t\|^2 \\
 &\quad + 4(1 - \beta_{t+1})^2 L_{22}^2 \eta_t^2 \mathbb{E}\|\tilde{y}_{t+1} - y_t\|^2 + \frac{2\beta_{t+1}^2 \sigma^2}{B}.
 \end{aligned} \tag{86}$$

□

**Theorem 6.** (Restatement of Theorem 3) Suppose the sequence  $\{x_t, y_t\}_{t=1}^T$  is generated from Algorithm 2. Given  $y_1 = y^*(x_1)$ ,  $c_1 \geq \frac{2}{3b^3} + 2\lambda\mu$ ,  $c_2 \geq \frac{2}{3b^3} + \frac{50\lambda\tilde{L}^2}{\mu}$ ,  $b > 0$ ,  $m \geq \max(2, (\tilde{c}b)^3)$ ,  $0 < \gamma \leq \frac{\mu\lambda}{2\kappa\tilde{L}\sqrt{25+4\mu\lambda}}$  and  $0 < \lambda \leq \frac{1}{6\tilde{L}}$ , we have

$$\frac{1}{T} \sum_{t=1}^T \mathbb{E}\|\text{grad } \Phi(x_t)\| \leq \frac{\sqrt{2M'}m^{1/6}}{T^{1/2}} + \frac{\sqrt{2M'}}{T^{1/3}}, \tag{87}$$

where  $\tilde{c} = \max(2\gamma L, c_1, c_2, 1)$  and  $M' = \frac{2(\Phi(x_1) - \Phi^*)}{\gamma b} + \frac{2\sigma^2}{\lambda\mu\eta_0 b B} + \frac{2(c_1^2 + c_2^2)\sigma^2 b^2}{\lambda\mu B} \ln(m + T)$ .

*Proof.* Since  $\eta_t$  is decreasing and  $m \geq b^3$ , we have  $\eta_t \leq \eta_0 = \frac{b}{m^{1/3}} \leq 1$ . Similarly, due to  $m \geq (2\gamma Lb)^3$ , we have  $\eta_t \leq \eta_0 = \frac{b}{m^{1/3}} \leq \frac{1}{2\gamma L}$ . Due to  $0 < \eta_t \leq 1$  and  $m \geq \max((c_1 b)^3, (c_2 b)^3)$ , we have  $\alpha_{t+1} = c_1 \eta_t^2 \leq c_1 \eta_t \leq \frac{c_1 b}{m^{1/3}} \leq 1$  and  $\beta_{t+1} = c_2 \eta_t^2 \leq c_2 \eta_t \leq \frac{c_2 b}{m^{1/3}} \leq 1$ . According to Lemma 10, we have

$$\begin{aligned}
 & \frac{1}{\eta_t} \mathbb{E} \|\text{grad}_x f(x_{t+1}, y_{t+1}) - v_{t+1}\|^2 - \frac{1}{\eta_{t-1}} \mathbb{E} \|\text{grad}_x f(x_t, y_t) - v_t\|^2 \\
 & \leq \left( \frac{(1 - \alpha_{t+1})^2}{\eta_t} - \frac{1}{\eta_{t-1}} \right) \mathbb{E} \|\text{grad}_x f(x_t, y_t) - v_t\|^2 + 4(1 - \alpha_{t+1})^2 L_{11}^2 \gamma^2 \eta_t \mathbb{E} \|v_t\|^2 \\
 & \quad + 4(1 - \alpha_{t+1})^2 L_{12}^2 \eta_t \mathbb{E} \|\tilde{y}_{t+1} - y_t\|^2 + \frac{2\alpha_{t+1}^2 \sigma^2}{\eta_t B} \\
 & \leq \left( \frac{1 - \alpha_{t+1}}{\eta_t} - \frac{1}{\eta_{t-1}} \right) \mathbb{E} \|\text{grad}_x f(x_t, y_t) - v_t\|^2 + 4L_{11}^2 \gamma^2 \eta_t \mathbb{E} \|v_t\|^2 + 4L_{12}^2 \eta_t \mathbb{E} \|\tilde{y}_{t+1} - y_t\|^2 + \frac{2\alpha_{t+1}^2 \sigma^2}{\eta_t B} \\
 & = \left( \frac{1}{\eta_t} - \frac{1}{\eta_{t-1}} - c_1 \eta_t \right) \mathbb{E} \|\text{grad}_x f(x_t, y_t) - v_t\|^2 + 4L_{11}^2 \gamma^2 \eta_t \mathbb{E} \|v_t\|^2 + 4L_{12}^2 \eta_t \mathbb{E} \|\tilde{y}_{t+1} - y_t\|^2 + \frac{2\alpha_{t+1}^2 \sigma^2}{\eta_t B},
 \end{aligned} \tag{89}$$

where the second inequality is due to  $0 < \alpha_{t+1} \leq 1$ . By a similar way, we also obtain

$$\begin{aligned}
 & \frac{1}{\eta_t} \mathbb{E} \|\nabla_y f(x_{t+1}, y_{t+1}) - w_{t+1}\|^2 - \frac{1}{\eta_{t-1}} \mathbb{E} \|\nabla_y f(x_t, y_t) - w_t\|^2 \\
 & \leq \left( \frac{1}{\eta_t} - \frac{1}{\eta_{t-1}} - c_2 \eta_t \right) \mathbb{E} \|\nabla_y f(x_t, y_t) - w_t\|^2 + 4L_{21}^2 \gamma^2 \eta_t \mathbb{E} \|v_t\|^2 + 4L_{22}^2 \eta_t \mathbb{E} \|\tilde{y}_{t+1} - y_t\|^2 + \frac{2\beta_{t+1}^2 \sigma^2}{\eta_t B}.
 \end{aligned} \tag{90}$$

By  $\eta_t = \frac{b}{(m+t)^{1/3}}$ , we have

$$\begin{aligned}
 \frac{1}{\eta_t} - \frac{1}{\eta_{t-1}} &= \frac{1}{b} \left( (m+t)^{\frac{1}{3}} - (m+t-1)^{\frac{1}{3}} \right) \\
 &\leq \frac{1}{3b(m+t-1)^{2/3}} \leq \frac{1}{3b(m/2+t)^{2/3}} \\
 &\leq \frac{2^{2/3}}{3b(m+t)^{2/3}} = \frac{2^{2/3}}{3b^3} \frac{b^2}{(m/2+t)^{2/3}} = \frac{2^{2/3}}{3b^3} \eta_t^2 \leq \frac{2}{3b^3} \eta_t,
 \end{aligned} \tag{91}$$

where the first inequality holds by the concavity of function  $f(x) = x^{1/3}$ , i.e.,  $(x+y)^{1/3} \leq x^{1/3} + \frac{y}{3x^{2/3}}$ ; the second inequality is due to  $m \geq 2$ , and the last inequality is due to  $0 < \eta_t \leq 1$ . Let  $c_1 \geq \frac{2}{3b^3} + 2\lambda\mu$ , we have

$$\begin{aligned}
 & \frac{1}{\eta_t} \mathbb{E} \|\text{grad}_x f(x_{t+1}, y_{t+1}) - v_{t+1}\|^2 - \frac{1}{\eta_{t-1}} \mathbb{E} \|\text{grad}_x f(x_t, y_t) - v_t\|^2 \\
 & \leq -2\lambda\mu \eta_t \mathbb{E} \|\text{grad}_x f(x_t, y_t) - v_t\|^2 + 4L_{11}^2 \gamma^2 \eta_t \mathbb{E} \|v_t\|^2 + 4L_{12}^2 \eta_t \mathbb{E} \|\tilde{y}_{t+1} - y_t\|^2 + \frac{2\alpha_{t+1}^2 \sigma^2}{\eta_t B}.
 \end{aligned} \tag{92}$$

Let  $c_2 \geq \frac{2}{3b^3} + \frac{50\lambda\tilde{L}^2}{\mu}$ , we have

$$\begin{aligned}
 & \frac{1}{\eta_t} \mathbb{E} \|\nabla_y f(x_{t+1}, y_{t+1}) - w_{t+1}\|^2 - \frac{1}{\eta_{t-1}} \mathbb{E} \|\nabla_y f(x_t, y_t) - w_t\|^2 \\
 & \leq -\frac{50\lambda\tilde{L}^2}{\mu} \eta_t \mathbb{E} \|\nabla_y f(x_t, y_t) - w_t\|^2 + 4L_{21}^2 \gamma^2 \eta_t \mathbb{E} \|v_t\|^2 + 4L_{22}^2 \eta_t \mathbb{E} \|\tilde{y}_{t+1} - y_t\|^2 + \frac{2\beta_{t+1}^2 \sigma^2}{\eta_t B}.
 \end{aligned} \tag{93}$$

According to Lemma 9, we have

$$\begin{aligned}
 \|y_{t+1} - y^*(x_{t+1})\|^2 - \|y_t - y^*(x_t)\|^2 &\leq -\frac{\eta_t \mu \lambda}{4} \|y_t - y^*(x_t)\|^2 - \frac{3\eta_t}{4} \|\tilde{y}_{t+1} - y_t\|^2 \\
 &\quad + \frac{25\lambda\eta_t}{6\mu} \|\nabla_y f(x_t, y_t) - w_t\|^2 + \frac{25\gamma^2 \kappa^2 \eta_t}{6\mu\lambda} \|v_t\|^2.
 \end{aligned} \tag{94}$$

Next, we define a *Lyapunov* function  $\Omega_t$ , for any  $t \geq 1$

$$\Omega_t = \mathbb{E} [\Phi(x_t) + \frac{\gamma}{2\lambda\mu\eta_{t-1}} (\|\text{grad}_x f(x_t, y_t) - v_t\|^2 + \|\nabla_y f(x_t, y_t) - w_t\|^2) + \frac{6\gamma\tilde{L}^2}{\lambda\mu} \|y_t - y^*(x_t)\|^2]. \tag{95}$$

Then we have

$$\begin{aligned}
 \Omega_{t+1} - \Omega_t &= \mathbb{E}[\Phi(x_{t+1})] - \mathbb{E}[\Phi(x_t)] + \frac{6\gamma\tilde{L}^2}{\lambda\mu} (\mathbb{E}\|y_{t+1} - y^*(x_{t+1})\|^2 - \mathbb{E}\|y_t - y^*(x_t)\|^2) \\
 &\quad + \frac{\gamma}{2\lambda\mu} \left( \frac{1}{\eta_t} \mathbb{E}\|\text{grad}_x f(x_{t+1}, y_{t+1}) - v_{t+1}\|^2 - \frac{1}{\eta_{t-1}} \mathbb{E}\|\text{grad}_x f(x_t, y_t) - v_t\|^2 \right) \\
 &\quad + \frac{1}{\eta_t} \mathbb{E}\|\nabla_y f(x_{t+1}, y_{t+1}) - w_{t+1}\|^2 - \frac{1}{\eta_{t-1}} \mathbb{E}\|\nabla_y f(x_t, y_t) - w_t\|^2 \\
 &\leq L_{12}\gamma\eta_t \mathbb{E}\|y_t - y^*(x_t)\|^2 + \gamma\eta_t \mathbb{E}\|\text{grad}_x f(x_t, y_t) - v_t\|^2 - \frac{\gamma\eta_t}{2} \mathbb{E}\|\text{grad}\Phi(x_t)\|^2 - \frac{\gamma\eta_t}{4} \mathbb{E}\|v_t\|^2 \\
 &\quad + \frac{6\gamma\tilde{L}^2}{\lambda\mu} \left( -\frac{\mu\lambda\eta_t}{4} \mathbb{E}\|y_t - y^*(x_t)\|^2 - \frac{3\eta_t}{4} \mathbb{E}\|\tilde{y}_{t+1} - y_t\|^2 + \frac{25\lambda\eta_t}{6\mu} \mathbb{E}\|\nabla_y f(x_t, y_t) - w_t\|^2 + \frac{25\gamma^2\kappa^2\eta_t}{6\mu\lambda} \mathbb{E}\|v_t\|^2 \right) \\
 &\quad + \frac{\gamma}{2\lambda\mu} \left( -2\lambda\mu\eta_t \mathbb{E}\|\text{grad}_x f(x_t, y_t) - v_t\|^2 + 4L_{11}^2\gamma^2\eta_t \mathbb{E}\|v_t\|^2 + 4L_{12}^2\eta_t \mathbb{E}\|\tilde{y}_{t+1} - y_t\|^2 + \frac{2\alpha_{t+1}^2\sigma^2}{\eta_t B} \right. \\
 &\quad \left. - \frac{50\lambda\tilde{L}^2}{\mu} \eta_t \mathbb{E}\|\nabla_y f(x_t, y_t) - w_t\|^2 + 4L_{21}^2\gamma^2\eta_t \mathbb{E}\|v_t\|^2 + 4L_{22}^2\eta_t \mathbb{E}\|\tilde{y}_{t+1} - y_t\|^2 + \frac{2\beta_{t+1}^2\sigma^2}{\eta_t B} \right) \\
 &\leq -\frac{\gamma\tilde{L}^2\eta_t}{2} \mathbb{E}\|y_t - y^*(x_t)\|^2 - \frac{\gamma\eta_t}{2} \mathbb{E}\|\text{grad}\Phi(x_t)\|^2 - \frac{\gamma\tilde{L}^2\eta_t}{2\lambda\mu} \mathbb{E}\|\tilde{y}_{t+1} - y_t\|^2 - \left( \frac{\gamma}{4} - \frac{25\gamma^3\kappa^2\tilde{L}^2}{\mu^2\lambda^2} - \frac{4\gamma^3\tilde{L}^2}{\mu\lambda} \right) \eta_t \mathbb{E}\|v_t\|^2 \\
 &\quad + \frac{\gamma\alpha_{t+1}^2\sigma^2}{\lambda\mu\eta_t B} + \frac{\gamma\beta_{t+1}^2\sigma^2}{\lambda\mu\eta_t B} \\
 &\leq -\frac{\gamma\tilde{L}^2\eta_t}{2} \mathbb{E}\|y_t - y^*(x_t)\|^2 - \frac{\gamma\eta_t}{2} \mathbb{E}\|\text{grad}\Phi(x_t)\|^2 + \frac{\gamma\alpha_{t+1}^2\sigma^2}{\lambda\mu\eta_t B} + \frac{\gamma\beta_{t+1}^2\sigma^2}{\lambda\mu\eta_t B}, \tag{96}
 \end{aligned}$$

where the first inequality holds by Lemmas 8 and the above inequalities (92), (93) and (94); the second inequality is due to  $\tilde{L} = \max(1, L_{11}, L_{12}, L_{21}, L_{22})$ ; the last inequality is due to  $0 \leq \gamma \leq \frac{\mu\lambda}{2\kappa\tilde{L}\sqrt{25+4\mu\lambda}}$  and  $\kappa \geq 1$ .

According to the above inequality (96), we have

$$\frac{\gamma\eta_t}{2} (\mathbb{E}\|\text{grad}\Phi(x_t)\|^2 + \tilde{L}^2 \mathbb{E}\|y_t - y^*(x_t)\|^2) \leq \Omega_t - \Omega_{t+1} + \frac{\gamma\alpha_{t+1}^2\sigma^2}{\lambda\mu\eta_t B} + \frac{\gamma\beta_{t+1}^2\sigma^2}{\lambda\mu\eta_t B}. \tag{97}$$

Taking average over  $t = 1, 2, \dots, T$  on both sides of the inequality (97), we have

$$\frac{1}{T} \sum_{t=1}^T \eta_t \mathbb{E}(\|\text{grad}\Phi(x_t)\|^2 + \tilde{L}^2 \|y_t - y^*(x_t)\|^2) \leq \sum_{t=1}^T \frac{2(\Omega_t - \Omega_{t+1})}{\gamma T} + \frac{1}{T} \sum_{t=1}^T \left( \frac{2\alpha_{t+1}^2\sigma^2}{\lambda\mu\eta_t B} + \frac{2\beta_{t+1}^2\sigma^2}{\lambda\mu\eta_t B} \right).$$

Since the initial solution satisfies  $y_1 = y^*(x_1) = \arg \max_{y \in \mathcal{Y}} f(x_1, y)$ , we have

$$\begin{aligned}
 \Omega_1 &= \Phi(x_1) + \frac{6\gamma\tilde{L}^2}{\lambda\mu} \|y_1 - y^*(x_1)\|^2 + \frac{\gamma}{2\lambda\mu} \left( \frac{1}{\eta_0} \mathbb{E}\|\text{grad}_x f(x_1, y_1) - v_1\|^2 + \frac{1}{\eta_0} \mathbb{E}\|\nabla_y f(x_1, y_1) - w_1\|^2 \right) \\
 &= \Phi(x_1) + \frac{\gamma}{2\lambda\mu} \left( \frac{1}{\eta_0} \mathbb{E}\|\text{grad}_x f(x_1, y_1) - \text{grad}_x f_{\mathcal{B}_1}(x_1, y_1)\|^2 + \frac{1}{\eta_0} \mathbb{E}\|\nabla_y f(x_1, y_1) - \nabla_y f_{\mathcal{B}_1}(x_1, y_1)\|^2 \right) \\
 &\leq \Phi(x_1) + \frac{\gamma\sigma^2}{\lambda\mu\eta_0 B}, \tag{98}
 \end{aligned}$$

where the last inequality holds by Assumption 5.

Consider  $\eta_t$  is decreasing, i.e.,  $\eta_T^{-1} \geq \eta_t^{-1}$  for any  $0 \leq t \leq T$ , we have

$$\begin{aligned}
 & \frac{1}{T} \sum_{t=1}^T \mathbb{E}(\|\text{grad}\Phi(x_t)\|^2 + \tilde{L}^2 \|y_t - y^*(x_t)\|^2) \\
 & \leq \sum_{t=1}^T \frac{2(\Omega_t - \Omega_{t+1})}{T\gamma\eta_T} + \frac{1}{T\eta_T} \sum_{t=1}^T \left( \frac{2\alpha_{t+1}^2\sigma^2}{\lambda\mu\eta_t B} + \frac{2\beta_{t+1}^2\sigma^2}{\lambda\mu\eta_t B} \right) \\
 & \leq \frac{1}{T\eta_T} \left( \frac{2\Phi(x_1)}{\gamma} + \frac{2\sigma^2}{\lambda\mu\eta_0 B} - \frac{2\Phi^*}{\gamma} \right) + \frac{1}{T\eta_T} \sum_{t=1}^T \left( \frac{2\alpha_{t+1}^2\sigma^2}{\lambda\mu\eta_t B} + \frac{2\beta_{t+1}^2\sigma^2}{\lambda\mu\eta_t B} \right) \\
 & = \frac{2(\Phi(x_1) - \Phi^*)}{T\gamma\eta_T} + \frac{2\sigma^2}{T\lambda\mu\eta_0\eta_TB} + \frac{2(c_1^2 + c_2^2)\sigma^2}{T\eta_T\lambda\mu B} \sum_{t=1}^T \eta_t^3 \\
 & \leq \frac{2(\Phi(x_1) - \Phi^*)}{T\gamma\eta_T} + \frac{2\sigma^2}{T\lambda\mu\eta_0\eta_TB} + \frac{2(c_1^2 + c_2^2)\sigma^2}{T\eta_T\lambda\mu B} \int_1^T \frac{b^3}{m+t} dt \\
 & \leq \frac{2(\Phi(x_1) - \Phi^*)}{T\gamma\eta_T} + \frac{2\sigma^2}{T\lambda\mu\eta_0\eta_TB} + \frac{2(c_1^2 + c_2^2)\sigma^2 b^3}{T\eta_T\lambda\mu B} \ln(m+T) \\
 & = \frac{2(\Phi(x_1) - \Phi^*)}{T\gamma b} (m+T)^{1/3} + \frac{2\sigma^2}{T\lambda\mu\eta_0 b B} (m+T)^{1/3} + \frac{2(c_1^2 + c_2^2)\sigma^2 b^2}{T\lambda\mu B} \ln(m+T)(m+T)^{1/3},
 \end{aligned} \tag{99}$$

where the third inequality holds by  $\sum_{t=1}^T \eta_t^3 \leq \int_1^T \eta_t^3 dt$ . Let  $M' = \frac{2(\Phi(x_1) - \Phi^*)}{\gamma b} + \frac{2\sigma^2}{\lambda\mu\eta_0 b B} + \frac{2(c_1^2 + c_2^2)\sigma^2 b^2}{\lambda\mu B} \ln(m+T)$ , we rewrite the above inequality as follows:

$$\frac{1}{T} \sum_{t=1}^T \mathbb{E}(\|\text{grad}\Phi(x_t)\|^2 + \tilde{L}^2 \|y_t - y^*(x_t)\|^2) \leq \frac{M'}{T} (m+T)^{1/3}. \tag{100}$$

According to Jensen's inequality, we have

$$\begin{aligned}
 \frac{1}{T} \sum_{t=1}^T \mathbb{E}(\|\text{grad}\Phi(x_t)\| + \tilde{L} \|y_t - y^*(x_t)\|) & \leq \left( \frac{2}{T} \sum_{t=1}^T \mathbb{E}(\|\text{grad}\Phi(x_t)\|^2 + \tilde{L}^2 \|y_t - y^*(x_t)\|^2) \right)^{1/2} \\
 & \leq \frac{\sqrt{2M'}}{T^{1/2}} (m+T)^{1/6} \leq \frac{\sqrt{2M'}m^{1/6}}{T^{1/2}} + \frac{\sqrt{2M'}}{T^{1/3}},
 \end{aligned} \tag{101}$$

where the last inequality is due to  $(a_1 + a_2)^{1/6} \leq a_1^{1/6} + a_2^{1/6}$  for all  $a_1, a_2 > 0$ . Thus, we have

$$\frac{1}{T} \sum_{t=1}^T \mathbb{E}\|\text{grad}\Phi(x_t)\| \leq \frac{\sqrt{2M'}m^{1/6}}{T^{1/2}} + \frac{\sqrt{2M'}}{T^{1/3}}. \tag{102}$$

□