

# On the Initialization for Convex-Concave Min-max Problems

**Mingrui Liu**

MINGRUIL@GMU.EDU

*Department of Computer Science, George Mason University, Fairfax, VA 22030, USA*

**Francesco Orabona**

FRANCESCO@ORABONA.COM

*Electrical & Computer Engineering, Boston University, Boston, Massachusetts 02215, USA*

**Editors:** Sanjoy Dasgupta and Nika Haghtalab

## Abstract

Convex-concave min-max problems are ubiquitous in machine learning, and people usually utilize first-order methods (e.g., gradient descent ascent) to find the optimal solution. One feature which separates convex-concave min-max problems from convex minimization problems is that the best known convergence rates for min-max problems have an explicit dependence on the size of the domain, rather than on the distance between initial point and the optimal solution. This means that the convergence speed does not have any improvement even if the algorithm starts from the optimal solution, and hence, is oblivious to the initialization. Here, we show that strict-convexity-strict-concavity is sufficient to get the convergence rate to depend on the initialization. We also show how different algorithms can asymptotically achieve initialization-dependent convergence rates on this class of functions. Furthermore, we show that the so-called “parameter-free” algorithms allow to achieve improved initialization-dependent asymptotic rates without any learning rate to tune. In addition, we utilize this particular parameter-free algorithm as a subroutine to design a new algorithm, which achieves a novel non-asymptotic fast rate for strictly-convex-strictly-concave min-max problems with a growth condition and Hölder continuous solution mapping. Experiments are conducted to verify our theoretical findings and demonstrate the effectiveness of the proposed algorithms.

**Keywords:** Convex-concave, Min-max, Initialization, Fast Rates, Parameter-Free

## 1. Introduction

In this paper, we are interested in the following problem

$$\min_{\mathbf{x} \in \mathcal{X}} \max_{\mathbf{y} \in \mathcal{Y}} F(\mathbf{x}, \mathbf{y}), \quad (1)$$

where  $\mathcal{X} \subset \mathbb{R}^m$ ,  $\mathcal{Y} \subset \mathbb{R}^n$  are convex and compact sets. This problem has broad applications in machine learning, e.g., stochastic AUC maximization (Ying et al., 2016), generative adversarial nets (Goodfellow et al., 2014), robust optimization (Ben-Tal et al., 2009), and adversarial training (Madry et al., 2017). Here, we will assume that  $F(\mathbf{x}, \mathbf{y})$  is convex in  $\mathbf{x}$  and concave in  $\mathbf{y}$ .

The canonical method for solving the convex-concave game is the primal-dual gradient method (a.k.a., gradient descent ascent) (Nemirovski and Yudin, 1983; Nemirovski et al., 2009). For example, in the Euclidean setup, primal-dual gradient method is performing gradient descent on the primal variable  $\mathbf{x}$  and gradient ascent on the dual variable  $\mathbf{y}$  simultaneously. It is proved (Nemirovski and Yudin, 1983; Nemirovski et al., 2009) that the averaged solution of primal-dual gradient method has good convergence guarantees in terms of the duality gap, which is

$$\max_{\mathbf{y} \in \mathcal{Y}} F(\bar{\mathbf{x}}_T, \mathbf{y}) - \min_{\mathbf{x} \in \mathcal{X}} F(\mathbf{x}, \bar{\mathbf{y}}_T) \leq \frac{D^2}{\eta T} + \eta G^2, \quad (2)$$

where  $(\bar{\mathbf{x}}_T, \bar{\mathbf{y}}_T) = (\frac{1}{T} \sum_{t=1}^T \mathbf{x}_t, \frac{1}{T} \sum_{t=1}^T \mathbf{y}_t)$  is the averaged solution,  $D$  is the diameter of the domain, and  $G$  is an upper bound on the norm of the gradient. To get the tightest bound for the RHS of (2), we can choose  $\eta = \frac{D}{G\sqrt{T}}$  and end up with  $\frac{DG}{\sqrt{T}}$  convergence rate for the duality gap.

Strangely enough, the optimal learning rate scheme in the primal-dual gradient method, as well as its corresponding convergence rate, have an explicit dependence on the size of the domain and are completely independent of the initialization. In other words, regardless of how close the initialization is to the optimal solution, this algorithm uses the same learning rate and ends up with the same complexity guarantees. This is counter-intuitive: we would expect to be able to obtain a faster convergence if the initialization is closer to the optimal solution. In other words, we aim to obtain an initialization-dependent convergence rate, i.e.,  $\tilde{\mathcal{O}}\left(\frac{f(\text{dist}(\mathbf{x}_0, \mathbf{x}_*) + f(\text{dist}(\mathbf{y}_0, \mathbf{y}_*)))}{\sqrt{T}}\right)$  rate,<sup>1</sup> where  $\text{dist}$  is some metric depending on the geometry of the problem, and  $f(\cdot) : \mathbb{R}^+ \mapsto \mathbb{R}^+$  is some non-decreasing function depending on the algorithm with  $f(0) = 0$ . Yet, the hardness results in the optimization literature show that  $\Omega(\frac{DG}{\sqrt{T}})$  complexity lower bound is unimprovable in general (Nemirovski and Yudin, 1983; Juditsky et al., 2011). This naturally motivates the following questions:

*What is the suitable class of the functions such that one can use first-order algorithms to get the initialization-dependent rates? In addition, how to utilize this fact to design better algorithms with faster rates for min-max problems?*

The main goal of this paper is to answer these questions. We show that strict-convexity-strict-concavity is the key to derive initialization-dependent rates as well as fast rates of first-order algorithms for solving min-max problems.

More in details, our contributions are summarized as follows.

- We identify a condition (i.e., the problem instance is a strictly-convex-strictly-concave min-max problem), and we show that this condition is sufficient to obtain asymptotic initialization-dependent rates for any first-order algorithm with a generic convergence guarantee.
- We show that invoking a parameter-free algorithm (Orabona and Pál, 2016; Cutkosky and Orabona, 2018) on both primal and dual variables can achieve asymptotic initialization-dependent rates without any learning rate to tune. Taking into account the known upper bounds of gradient descent ascent without additional assumptions on the curvature, it remains unclear how to obtain an initialization-dependent rate unless we have the knowledge of the distance between initialization and the optimal solution.
- When the function admits a growth condition and Hölder continuous solution mapping, we design a new algorithm by utilizing the parameter-free algorithm as a subroutine and periodically restarting the subroutine. Thanks to the initialization-dependent rate, we prove that our algorithm enjoys novel fast non-asymptotic rates. To the best of our knowledge, this is the first work leveraging the function growth condition and Hölder continuous solution mapping to obtain these improved non-asymptotic rates for min-max problems.
- We verify our theoretical results by conducting both synthetic experiments and distributionally robust optimization on benchmark datasets. We empirically show that our algorithms exhibit good performance in practice.

---

1. The  $\tilde{\mathcal{O}}$  notation hides poly-logarithmic terms. In the case that  $\mathbf{x}_0 = \mathbf{x}_*$ ,  $\mathbf{y}_0 = \mathbf{y}_*$ , we allow the big O notation to hide lower order terms such as  $\mathcal{O}(1/T)$ .

## 2. Related Work

**Convex-Concave Min-max Optimization** Convex-concave Min-max Optimization is widely studied in optimization literature, and it is closely related to the variational inequality. The work of [Korpelevich \(1976\)](#) proposed the extragradient method for solving variational inequalities, and this was later extended into non-euclidean space (e.g., mirror-prox ([Nemirovski, 2004](#)), dual extrapolation ([Nesterov, 2007](#))). The stochastic version of mirror-prox was proposed by [Juditsky et al. \(2011\)](#). [Hsieh et al. \(2020\)](#) proposed variable stepsize scaling for extragradient method to improve the algorithm’s performance. In nonsmooth case, [Nemirovski et al. \(2009\)](#) analyzed the primal-dual gradient method in non-euclidean space. [Nedić and Ozdaglar \(2009\)](#) considered subgradient methods for solving min-max problems and provided per-iteration convergence rate estimates on the solutions. [Monteiro and Svaiter \(2010\)](#) designed hybrid proximal extragradient methods with a different performance measure. [Bach and Levy \(2019\)](#) provided a universal algorithm for solving variational inequalities, which adapts to noise and smoothness.

There are several papers considering specific cases in convex-concave min-max optimization, including functions with a bilinear term ([Nesterov, 2005](#); [Chen et al., 2014, 2017](#); [He and Monteiro, 2016](#); [Liu et al., 2018](#); [Daskalakis et al., 2018](#); [Liang and Stokes, 2019](#); [Gidel et al., 2019](#); [Mokhtari et al., 2020](#); [Azizian et al., 2020](#); [Bailey et al., 2020](#)), smooth or strongly-convex-(strongly)-concave ([Nesterov and Scrimali, 2006](#); [Zhao, 2019](#); [Lin et al., 2020](#); [Yan et al., 2020](#)), the last-iterate convergence ([Abernethy et al., 2019](#); [Daskalakis and Panageas, 2018](#); [Golowich et al., 2020](#)), adapts to unknown smoothness parameter ([Diakonikolas, 2020](#)) . There are also some papers about establishing lower bounds in various cases ([Zhang et al., 2019](#); [Ibrahim et al., 2020](#); [Ouyang and Xu, 2021](#)).

However, none of these works provide an upper bound for duality gap which explicitly depends on the distance between the initialization and the optimal solution without assuming strong convexity/concavity ([Chambolle and Pock, 2011](#)).

**Parameter-Free Online Convex Optimization** In Online Convex Optimization (OCO) ([Gordon, 1999](#); [Zinkevich, 2003](#)), the aim of the learner is to minimize the regret w.r.t. any fixed predictor. Most of the OCO algorithm require some knowledge of the competitor, for example, its norm, in order to achieve the smallest regret (see, e.g., [Orabona, 2019](#)). Hence, it becomes impossible to compete *uniformly* with all competitors, unless the algorithm has some knowledge of the future. Morally speaking, the OCO setting is a strict generalization of the setting of stochastic optimization of convex functions and competing with any fixed predictor corresponds exactly to design convex optimization algorithms that have optimal dependency on the distance between the initial point and the optimal solution. Again, without some knowledge of the norm of the optimal solution, classic optimization algorithms fails to get the right dependency.

So-called *parameter-free* algorithms avoid setting step sizes completely and get the optimal dependency uniformly on all competitors ([Streeter and McMahan, 2012](#)). Note that “parameter-free” is a technical word that denotes only this ability of the algorithm and not a general lack of knowledge of *any* characteristic of the problem, see [Orabona \(2019, Section 9.7\)](#). Most of these algorithms are based on Follow-the-Regularized-Leader<sup>2</sup> (FTRL) ([Shalev-Shwartz, 2007](#); [Abernethy et al., 2008](#)) with a time-varying linearithmic (non-strongly convex) regularizer (e.g., [Streeter and McMahan, 2012](#); [Orabona, 2014](#); [Cutkosky and Boahen, 2016, 2017](#); [Kotłowski, 2020](#); [Kempka et al., 2019](#)). These algorithms can also be viewed as betting schemes through the duality between

---

2. Dual Averaging ([Nesterov, 2009](#)) is a specialization of FTRL to linear functions.

regret and reward (McMahan and Orabona, 2014; Orabona and Pál, 2016; Cutkosky and Orabona, 2018; Cutkosky and Sarlos, 2019).

As far as we know, we present the first application of parameter-free algorithms to convex-concave min-max problems.

### 3. Problem Setup

**Notation** Denote by  $\|\cdot\|$  the Euclidean norm. Define  $\langle \cdot, \cdot \rangle$  as the inner product in Euclidean space. Define  $\Pi_{\mathcal{S}}$  as the orthogonal projection operator onto the set  $\mathcal{S}$ . We denote vectors by bold letters, e.g.,  $\mathbf{x}, \mathbf{g}$ . We say that  $f(x) = o(g(x))$  iff  $\lim_{x \rightarrow \infty} \frac{f(x)}{g(x)} = 0$ .

**Setting and Assumptions** In (1), we assume that  $F$  is convex in the first argument and concave in the second one. Moreover, we assume to have access to a first-order black-box optimization oracle  $\mathbf{g} = (\mathbf{g}^x, \mathbf{g}^y)$  at any point  $(\mathbf{x}, \mathbf{y})$ , where  $\mathbf{g}^x \in \partial_{\mathbf{x}} F(\mathbf{x}, \mathbf{y})$ , the subgradient of  $F$  w.r.t. its first argument, and  $\mathbf{g}^y \in -\partial_{\mathbf{y}}(-F(\mathbf{x}, \mathbf{y}))$ , the negative subgradient of  $-F$  w.r.t. its second argument.

We will also assume that the subgradients  $\mathbf{g}^x$  and  $\mathbf{g}^y$  have bounded support. In particular, we assume w.l.o.g. that  $\|\mathbf{g}^x\| \leq 1$  and  $\|\mathbf{g}^y\| \leq 1$  for all  $\mathbf{x} \in \mathcal{X}$  and  $\mathbf{y} \in \mathcal{Y}$ . Note that it is reasonable because we assume bounded domains.

**An Sufficient Condition for Initialization-Dependent Rates: the Class of Strictly-Convex-Strictly-Concave Min-Max Problems** In this paper, we aim to show that we can achieve initialization-dependent convergence rates for the class of strictly-convex-strictly-concave min-max problems (which covers the class of strongly-convex-strongly-concave min-max problems as a subclass), which means that  $F(\mathbf{x}, \mathbf{y})$  is strictly convex in  $\mathbf{x}$  and strictly concave in  $\mathbf{y}$ . The definition of strictly convex function is the following.

**Definition 1** *A function  $h$  is strictly convex if  $h(\lambda \mathbf{x} + (1 - \lambda)\mathbf{y}) < \lambda h(\mathbf{x}) + (1 - \lambda)h(\mathbf{y})$  for any  $0 < \lambda < 1$  and any  $\mathbf{x} \neq \mathbf{y}$ .*

Given the above arguments, in the following we will assume that the following assumption holds:

**Assumption 1**  *$F(\mathbf{x}, \mathbf{y})$  is strictly convex in  $\mathbf{x}$  and strictly concave in  $\mathbf{y}$  and we denote by  $(\mathbf{x}_*, \mathbf{y}_*)$  the optimal solution of the min-max problem (1).  $F(\mathbf{x}, \mathbf{y})$  is continuous in  $\mathbf{y}$  (resp.  $\mathbf{x}$ ) given  $\mathbf{x}$  (resp.  $\mathbf{y}$ ).*

## 4. Algorithms with Initialization-Dependent Convergence Rates

In this section, we first present a general convergence theory for first-order algorithms to obtain initialization-dependent convergence rates, under the assumption that the function is strictly-convex-strictly-concave (in Section 4.1). We also show conditions under which we get initialization-dependent convergence rates for gradient descent ascent (in Section 4.2). Then, we show a parameter-free algorithm that always enjoys an initialization-dependent convergence rate (in Section 4.3). Finally we discuss possible alternative approaches to obtain initialization-dependent rate (in Section 4.4).

### 4.1. General Convergence Theory for First-order Algorithms

In this section, we present our main theorem (Theorem 1), which characterizes a general convergence theory for first-order algorithms to achieve initialization-dependent asymptotic rates.

**Theorem 1** Let  $D_{\mathcal{X}}$  be the diameter of  $\mathcal{X}$  and  $D_{\mathcal{Y}}$  is the diameter of  $\mathcal{Y}$ , and  $D = \max(D_{\mathcal{X}}, D_{\mathcal{Y}})$ . Suppose Assumption 1 holds, and there is an algorithm  $\mathcal{A}$  which returns a solution  $(\tilde{\mathbf{x}}_T, \tilde{\mathbf{y}}_T)$  after  $T$  iterations with following guarantee:

$$\max_{\mathbf{y} \in \mathcal{Y}} F(\tilde{\mathbf{x}}_T, \mathbf{y}) - \min_{\mathbf{x} \in \mathcal{X}} F(\mathbf{x}, \tilde{\mathbf{y}}_T) \leq A(\|\mathbf{x}_0 - \mathbf{x}'_T\|)B(T) + A(\|\mathbf{y}_0 - \mathbf{y}'_T\|)B(T) + C(T), \quad (3)$$

where  $\mathbf{x}'_T = \arg \min_{\mathbf{x}} F(\mathbf{x}, \tilde{\mathbf{y}}_T)$ ,  $\mathbf{y}'_T = \arg \max_{\mathbf{y}} F(\tilde{\mathbf{x}}_T, \mathbf{y})$ ,  $A : \mathbb{R}_+ \rightarrow \mathbb{R}$  is non-decreasing,  $A(0) = 0$ ,  $B : \mathcal{N} \rightarrow \mathbb{R}_{++}$  and  $C : \mathcal{N} \rightarrow \mathbb{R}_{++}$  are non-increasing,  $\lim_{T \rightarrow \infty} B(T) = 0$ ,  $\lim_{T \rightarrow \infty} C(T) = 0$ ,  $A(x + y) \leq c(A(x) + A(y))$  with  $c$  being a positive constant independent of  $D$ . Then, we have

$$\max_{\mathbf{y} \in \mathcal{Y}} F(\tilde{\mathbf{x}}_T, \mathbf{y}) - \min_{\mathbf{x} \in \mathcal{X}} F(\mathbf{x}, \tilde{\mathbf{y}}_T) \leq cB(T)(A(\|\mathbf{x}_0 - \mathbf{x}_*\|) + A(\|\mathbf{y}_0 - \mathbf{y}_*\|)) + R(T) + C(T), \quad (4)$$

where the residual term  $R : \mathcal{N} \rightarrow \mathbb{R}_+$  has the following properties:

$$R(t) = o(B(t)) \quad \text{and} \quad R(T) \leq 2c \cdot A(D)B(T).$$

Theorem 1 indicates that if we have any algorithm with guarantees of form (3), then under Assumption 1, it automatically enjoys an initialization-dependent asymptotic convergence rate as illustrated by (4).

The Assumption 1 is crucial to get the desired bound due to the following reasons. First, it can guarantee that both  $\mathbf{x}'_T$  and  $\mathbf{y}'_T$  are unique. Second, we can utilize this assumption to show that  $\mathbf{x}'_T$  ( $\mathbf{y}'_T$ ) gets close to  $\mathbf{x}_*$  ( $\mathbf{y}_*$ ) when the algorithm runs a sufficiently large number of iterations, which ends up with the  $o(B(T))$  term in the bound (4).

**Proof** [Proof of Theorem 1] According to the fact that  $A(x)$  is non-decreasing and (3), and noting that  $\|\mathbf{x}_0 - \mathbf{x}'_T\| \leq \|\mathbf{x}_0 - \mathbf{x}_*\| + \|\mathbf{x}_* - \mathbf{x}'_T\|$ ,  $\|\mathbf{y}_0 - \mathbf{y}'_T\| \leq \|\mathbf{y}_0 - \mathbf{y}_*\| + \|\mathbf{y}_* - \mathbf{y}'_T\|$ , we have

$$\begin{aligned} & \max_{\mathbf{y} \in \mathcal{Y}} F(\tilde{\mathbf{x}}_T, \mathbf{y}) - \min_{\mathbf{x} \in \mathcal{X}} F(\mathbf{x}, \tilde{\mathbf{y}}_T) \\ & \leq A(\|\mathbf{x}_0 - \mathbf{x}_*\| + \|\mathbf{x}_* - \mathbf{x}'_T\|)B(T) + A(\|\mathbf{y}_0 - \mathbf{y}_*\| + \|\mathbf{y}_* - \mathbf{y}'_T\|)B(T) + C(T) \\ & \leq cB(T) [A(\|\mathbf{x}_0 - \mathbf{x}_*\|) + A(\|\mathbf{y}_0 - \mathbf{y}_*\|) + A(\|\mathbf{x}_* - \mathbf{x}'_T\|) + A(\|\mathbf{y}_* - \mathbf{y}'_T\|)] + C(T), \end{aligned} \quad (5)$$

where the last inequality holds due to the fact that  $A(x + y) \leq c(A(x) + A(y))$ . It remains to show that  $A(\|\mathbf{x}_* - \mathbf{x}'_T\|) + A(\|\mathbf{y}_* - \mathbf{y}'_T\|) \rightarrow 0$  when  $T \rightarrow \infty$ . Since  $A(0) = 0$ , it suffices to show that  $\|\mathbf{x}_* - \mathbf{x}'_T\| \rightarrow 0$  and  $\|\mathbf{y}_* - \mathbf{y}'_T\| \rightarrow 0$ .

Note that

$$\begin{aligned} F(\mathbf{x}_*, \mathbf{y}_*) - F(\mathbf{x}_*, \tilde{\mathbf{y}}_T) & \leq F(\mathbf{x}_*, \mathbf{y}_*) - \min_{\mathbf{x} \in \mathcal{X}} F(\mathbf{x}, \tilde{\mathbf{y}}_T) \leq \max_{\mathbf{y} \in \mathcal{Y}} F(\tilde{\mathbf{x}}_T, \mathbf{y}) - \min_{\mathbf{x} \in \mathcal{X}} F(\mathbf{x}, \tilde{\mathbf{y}}_T) + C(T) \\ & \leq A(\|\mathbf{x}_0 - \mathbf{x}'_T\|)B(T) + A(\|\mathbf{y}_0 - \mathbf{y}'_T\|)B(T) + C(T) \\ & \leq A(D_{\mathcal{X}})B(T) + A(D_{\mathcal{Y}})B(T) + C(T) \leq 2A(D)B(T) + C(T). \end{aligned} \quad (6)$$

(i). We now claim that  $\tilde{\mathbf{y}}_T \rightarrow \mathbf{y}_*$  when  $T \rightarrow \infty$ . Let's see why.

Note that  $F(\mathbf{x}, \mathbf{y})$  is strictly concave in terms of  $\mathbf{y}$ . Define  $\tilde{F}(\mathbf{y}) = F(\mathbf{x}_*, \mathbf{y})$ . Taking the limit for  $T \rightarrow \infty$  in (6), for the sequence  $\{\tilde{\mathbf{y}}_t\}_{t=1}^{\infty}$ , we have  $\tilde{F}(\tilde{\mathbf{y}}_t) \rightarrow \tilde{F}(\mathbf{y}_*)$ . Given that the domain is bounded, we can extract a convergent subsequence  $\{\tilde{\mathbf{y}}_{t_j}\} \subset \{\tilde{\mathbf{y}}_t\}$  and we assume that  $\tilde{\mathbf{y}}_{t_j} \rightarrow \tilde{\mathbf{y}}$ . By the continuity of  $\tilde{F}(\mathbf{y})$  in terms of  $\mathbf{y}$ , we know that  $\tilde{F}(\tilde{\mathbf{y}}_{t_j}) \rightarrow \tilde{F}(\tilde{\mathbf{y}})$ . Now,  $\tilde{F}(\tilde{\mathbf{y}}_{t_j})$  is also a

subsequence of the convergent sequence  $\tilde{F}(\tilde{\mathbf{y}}_t)$ , then  $\tilde{F}(\tilde{\mathbf{y}}_{t_j}) \rightarrow \tilde{F}(\mathbf{y}_*)$ . Since  $\mathbf{y}_*$  is uniquely defined, this implies that  $\tilde{\mathbf{y}} = \mathbf{y}_*$ . This means that any convergent subsequence of  $\{\tilde{\mathbf{y}}_t\}_{t=1}^\infty$  converges to  $\mathbf{y}_*$ , so  $\tilde{\mathbf{y}}_T \rightarrow \mathbf{y}_*$  when  $T \rightarrow \infty$ .

(ii). Our next claim is that the mapping  $\arg \min_{\mathbf{x} \in \mathcal{X}} F(\mathbf{x}, \mathbf{y})$  is a continuous function in  $\mathbf{y}$ .

First, define  $H(\mathbf{y}) = \arg \min_{\mathbf{x} \in \mathcal{X}} F(\mathbf{x}, \mathbf{y})$ . By the compactness of  $\mathcal{Y}$ , we have a sequence  $\mathbf{y}_k \rightarrow \mathbf{y}_*$ . Define  $\mathbf{x}_k = H(\mathbf{y}_k)$  and  $\mathbf{x}_* = H(\mathbf{y}_*)$ . By the compactness of  $\mathcal{X}$ , there exists a convergent subsequence  $\mathbf{x}_{k_i}$  of  $\mathbf{x}_k$ , and we denote its limit by  $\tilde{\mathbf{x}}$ . From the above, we have  $F(\mathbf{x}_{k_i}, \mathbf{y}_{k_i}) \leq F(\mathbf{x}, \mathbf{y}_{k_i})$  for all  $\mathbf{x} \in \mathcal{X}$ , that implies  $F(\tilde{\mathbf{x}}, \mathbf{y}_*) \leq F(\mathbf{x}, \mathbf{y}_*)$  for any  $\mathbf{x} \in \mathcal{X}$ . In particular, this implies  $F(\tilde{\mathbf{x}}, \mathbf{y}_*) \leq F(\mathbf{x}_*, \mathbf{y}_*)$ . By the uniqueness of the minimizer, we must have  $\tilde{\mathbf{x}} = \mathbf{x}_*$ . Given that any convergent subsequence converges to  $\mathbf{x}_*$ , this means that  $\mathbf{x}_k \rightarrow \mathbf{x}_*$  and hence  $\arg \min_{\mathbf{x} \in \mathcal{X}} F(\mathbf{x}, \mathbf{y})$  is a continuous function in terms of  $\mathbf{y}$ .

Combining the above two claims (i) and (ii), and by the definition of  $\mathbf{x}_*$  and  $\mathbf{x}'_T$ , we obtain that  $\mathbf{x}'_T \rightarrow \mathbf{x}_*$  when  $T \rightarrow \infty$ . Hence  $\|\mathbf{x}_* - \mathbf{x}'_T\| \rightarrow 0$ . A parallel argument can guarantee that  $\|\mathbf{y}_* - \mathbf{y}'_T\| \rightarrow 0$ . This completes the proof.  $\blacksquare$

## 4.2. Initialization-dependent Rate by Gradient Descent Ascent

In this section, we show that the standard algorithm gradient descent ascent satisfies (3). The update rule of gradient descent ascent is

$$\begin{cases} \mathbf{x}_{t+1} = \Pi_{\mathcal{X}} [\mathbf{x}_t - \eta \mathbf{g}^{\mathbf{x}}(\mathbf{x}_t, \mathbf{y}_t)] \\ \mathbf{y}_{t+1} = \Pi_{\mathcal{Y}} [\mathbf{y}_t + \eta \mathbf{g}^{\mathbf{y}}(\mathbf{x}_t, \mathbf{y}_t)] \end{cases}$$

where  $\mathbf{g}^{\mathbf{x}}(\mathbf{x}_t, \mathbf{y}_t) \in \partial_{\mathbf{x}} F(\mathbf{x}_t, \mathbf{y}_t)$ ,  $\mathbf{g}^{\mathbf{y}}(\mathbf{x}_t, \mathbf{y}_t) \in -\partial_{\mathbf{y}} (-F(\mathbf{x}_t, \mathbf{y}_t))$ ,  $\Pi$  is the projection operator,  $\eta$  is the learning rate, and  $\partial_{\mathbf{x}}, \partial_{\mathbf{y}}$  denote the subdifferential in terms of  $\mathbf{x}$  and  $\mathbf{y}$  respectively. If we use  $\eta = \alpha/\sqrt{T}$  and run the gradient descent ascent for  $T$  iterations, the standard theory of gradient descent ascent (Nemirovski et al., 2009) shows that

$$\max_{\mathbf{y} \in \mathcal{Y}} F(\bar{\mathbf{x}}_T, \mathbf{y}) - \min_{\mathbf{x} \in \mathcal{X}} F(\mathbf{x}, \bar{\mathbf{y}}_T) \leq A(\|\mathbf{x}_0 - \mathbf{x}'_T\|)B(T) + A(\|\mathbf{y}_0 - \mathbf{y}'_T\|)B(T) + C(T),$$

where  $A(x) = \frac{x^2}{2\alpha}$ ,  $B(T) = \frac{1}{\sqrt{T}}$ ,  $C = \frac{\alpha}{\sqrt{T}}$ ,  $\mathbf{x}'_T = \arg \min_{\mathbf{x} \in \mathcal{X}} F(\mathbf{x}, \mathbf{y}_T)$ , and  $\mathbf{y}'_T = \arg \min_{\mathbf{y} \in \mathcal{Y}} F(\bar{\mathbf{x}}_T, \mathbf{y})$ . We can verify that it satisfies the premises of Theorem 1 (note that  $A(x+y) \leq 2A(x) + 2A(y)$ ). Hence, invoking Theorem 1 yields

$$\begin{aligned} \max_{\mathbf{y} \in \mathcal{Y}} F(\bar{\mathbf{x}}_T, \mathbf{y}) - \min_{\mathbf{x} \in \mathcal{X}} F(\mathbf{x}, \bar{\mathbf{y}}_T) &\leq 2A(\|\mathbf{x}_0 - \mathbf{x}_*\|)B(T) + 2A(\|\mathbf{y}_0 - \mathbf{y}_*\|)B(T) + R(T) + \frac{\alpha}{\sqrt{T}} \\ &= \frac{\|\mathbf{x}_0 - \mathbf{x}_*\|^2 + \|\mathbf{y}_0 - \mathbf{y}_*\|^2}{\alpha\sqrt{T}} + \frac{\alpha}{\sqrt{T}} + R(T), \end{aligned}$$

where  $R(t) = o(\frac{1}{\sqrt{t}})$  and  $R(T) \leq \frac{2D^2}{\alpha\sqrt{T}}$ .

From the above, we see that the hyperparameter  $\alpha$  in the learning rate plays a major role in the convergence rate. There are 3 possible choices for  $\alpha$ .

- The optimal setting of  $\alpha$  would be  $\sqrt{\|\mathbf{x}_0 - \mathbf{x}_*\|^2 + \|\mathbf{y}_0 - \mathbf{y}_*\|^2}$ , that could give an initialization-dependent convergence rate of  $\mathcal{O}(\frac{\sqrt{\|\mathbf{x}_0 - \mathbf{x}_*\|^2 + \|\mathbf{y}_0 - \mathbf{y}_*\|^2}}{\sqrt{T}})$ . However, it is easy to realize that

such setting is impossible in the black-box optimization setting: The algorithm has no way to estimate the initial condition, not even in the convex optimization case. Hence, for any fixed choice of  $\alpha$  there exists an optimization on which this choice is suboptimal. Indeed, there are little-known lower bounds in the OCO setting that hints to the fact that the rate  $\mathcal{O}\left(\frac{\sqrt{\|\mathbf{x}_0 - \mathbf{x}_*\|^2 + \|\mathbf{y}_0 - \mathbf{y}_*\|^2}}{\sqrt{T}}\right)$  might be actually *impossible* (Streeter and McMahan, 2012) and *an additional polylogarithmic price should be paid to estimate the initial condition*, for example trying a logarithmic number of possible settings of  $\alpha$  and selecting the best one.

- The worst-case choice and the standard choice for this algorithm (Nemirovski et al., 2009) is to set  $\alpha = D$  and obtain a finite-time rate of  $\frac{D}{\sqrt{T}}$ , without any asymptotic acceleration.
- The last choice we have is to set  $\alpha$  to any arbitrary number and have an asymptotic rate of  $\mathcal{O}\left(\frac{\|\mathbf{x}_0 - \mathbf{x}_*\|^2 + \|\mathbf{y}_0 - \mathbf{y}_*\|^2 + \alpha}{\sqrt{T}}\right)$ . This is the only viable choice that gives an initialization-dependent rate on all problems such that  $\alpha < \sqrt{\|\mathbf{x}_0 - \mathbf{x}_*\|^2 + \|\mathbf{y}_0 - \mathbf{y}_*\|^2}$ .

In the next section, we will show how parameter-free algorithms can get initialization-dependent asymptotic rate without any hyperparameter tuning.

### 4.3. Improved Initialization-Dependent Convergence by Parameter-free Algorithms

In this section, we establish an improved initialization-dependent convergence by parameter-free algorithms for solving min-max problems. The core idea of this approach is to decouple the primal problem and dual problem in (1), and by then utilize no-regret algorithms to perform the optimization. This is not a new idea by any means (see, for example, Abernethy et al., 2018). However, we need to be particularly careful of the dependency on the initial point. For this reason, we first state a Theorem to use no-regret algorithms in min-max problems that is a simple generalization of Theorem 9 in (Abernethy et al., 2018). In particular, here we care about considering the regret of the OCO algorithms w.r.t. specific points. This will be critical in obtaining initialization-dependent rates.

Suppose to use an OCO algorithm fed with losses  $\ell_t(\mathbf{x}) = F(\mathbf{x}, \mathbf{y}_t)$  that produces the iterates  $\mathbf{x}_t$  and another OCO algorithm fed with losses  $h_t(\mathbf{y}) = -F(\mathbf{x}_t, \mathbf{y})$  that produces the iterates  $\mathbf{y}_t$ . Then, we can state the following Theorem (for completeness the proof is in Appendix A).

**Theorem 2** *Let  $\bar{\mathbf{x}}_T = \frac{1}{T} \sum_{t=1}^T \mathbf{x}_t$ ,  $\bar{\mathbf{y}}_T = \frac{1}{T} \sum_{t=1}^T \mathbf{y}_t$ . Then, we have*

$$\max_{\mathbf{y} \in \mathcal{Y}} F(\bar{\mathbf{x}}_T, \mathbf{y}) - \min_{\mathbf{x} \in \mathcal{X}} F(\mathbf{x}, \bar{\mathbf{y}}_T) \leq \frac{R_T(\mathbf{x}'_T) + R_T(\mathbf{y}'_T)}{T},$$

where  $R_T(\mathbf{y}) = \sum_{t=1}^T h_t(\mathbf{y}_t) - \sum_{t=1}^T h_t(\mathbf{y})$ ,  $R_T(\mathbf{x}) = \sum_{t=1}^T \ell_t(\mathbf{x}_t) - \sum_{t=1}^T \ell_t(\mathbf{x})$ ,  $\mathbf{x}'_T \in \arg \min_{\mathbf{x} \in \mathcal{X}} F(\mathbf{x}, \bar{\mathbf{y}}_T)$ , and  $\mathbf{y}'_T \in \arg \max_{\mathbf{y} \in \mathcal{Y}} F(\bar{\mathbf{x}}_T, \mathbf{y})$ .

In words, the above theorem says that we can use two OCO algorithms to minimize the problem in (2). In particular, the convergence rate depends on the sum of the regret of the two OCO algorithms evaluated in specific points, divided by the sum of the weights.

Inspired by Orabona and Pál (2016); Cutkosky and Orabona (2018), we present a parameter-free algorithm for constrained optimization, in Algorithm 1. For it, we can prove the following regret guarantee in Theorem 3 (the proof is in Appendix A).

---

**Algorithm 1** Constrained Parameter-free OCO
 

---

**Input:** Convex and compact feasible set  $\mathcal{X}$ , initial point  $\mathbf{x}_0 \in \mathcal{X}$ , number of iterations  $T$ 

- 1:  $\tilde{\mathbf{x}}_1 = \mathbf{x}_0$
  - 2:  $\boldsymbol{\theta}_t = \mathbf{0}$
  - 3: **for**  $t = 1, \dots, T$  **do**
  - 4:    $\mathbf{x}_t = \Pi_{\mathcal{X}}(\tilde{\mathbf{x}}_t)$
  - 5:   Receive subgradients  $\hat{\mathbf{g}}_t \in \partial \ell_t(\mathbf{x}_t)$  such that  $\|\hat{\mathbf{g}}_t\| \leq 1$
  - 6:    $\mathbf{g}_t = \frac{1}{2} \left( \hat{\mathbf{g}}_t + \|\hat{\mathbf{g}}_t\| \cdot \frac{\tilde{\mathbf{x}}_t - \mathbf{x}_t}{\|\tilde{\mathbf{x}}_t - \mathbf{x}_t\|} \right)$  (Define  $\mathbf{0}/0 = \mathbf{0}$ )
  - 7:    $\tilde{\mathbf{x}}_{t+1} = \mathbf{x}_0 - \frac{\sum_{i=1}^{t-1} \mathbf{g}_i}{t+1} \left( 1 - \sum_{i=1}^{t-1} \langle \tilde{\mathbf{x}}_i, \mathbf{g}_i \rangle \right)$
  - 8: **end for**
- 

---

**Algorithm 2** CB-Min-Max( $\mathbf{x}_0, \mathbf{y}_0, T$ )
 

---

**Input:**  $\mathbf{x}_0 \in \mathcal{X}, \mathbf{y}_0 \in \mathcal{Y}$ 

- 1:  $\tilde{\mathbf{x}}_1 = \mathbf{x}_0, \tilde{\mathbf{y}}_1 = \mathbf{y}_0$
  - 2: **for**  $t = 1, \dots, T$  **do**
  - 3:    $\mathbf{x}_t = \Pi_{\mathcal{X}}(\tilde{\mathbf{x}}_t), \mathbf{y}_t = \Pi_{\mathcal{Y}}(\tilde{\mathbf{y}}_t)$
  - 4:   Receive subgradients  $\hat{\mathbf{g}}_t = (\hat{\mathbf{g}}_t^{\mathbf{x}}, \hat{\mathbf{g}}_t^{\mathbf{y}})$
  - 5:    $\mathbf{g}_t^{\mathbf{x}} = \frac{1}{2} \left( \hat{\mathbf{g}}_t^{\mathbf{x}} + \|\hat{\mathbf{g}}_t^{\mathbf{x}}\| \cdot \frac{\tilde{\mathbf{x}}_t - \mathbf{x}_t}{\|\tilde{\mathbf{x}}_t - \mathbf{x}_t\|} \right)$  (Define  $\mathbf{0}/0 = \mathbf{0}$ )
  - 6:    $\mathbf{g}_t^{\mathbf{y}} = \frac{1}{2} \left( \hat{\mathbf{g}}_t^{\mathbf{y}} + \|\hat{\mathbf{g}}_t^{\mathbf{y}}\| \cdot \frac{\tilde{\mathbf{y}}_t - \mathbf{y}_t}{\|\tilde{\mathbf{y}}_t - \mathbf{y}_t\|} \right)$  (Define  $\mathbf{0}/0 = \mathbf{0}$ )
  - 7:    $\tilde{\mathbf{x}}_{t+1} = \mathbf{x}_0 - \frac{\sum_{i=1}^{t-1} \mathbf{g}_i^{\mathbf{x}}}{t+1} \left( 1 - \sum_{i=1}^{t-1} \langle \tilde{\mathbf{x}}_i, \mathbf{g}_i^{\mathbf{x}} \rangle \right)$
  - 8:    $\tilde{\mathbf{y}}_{t+1} = \mathbf{y}_0 + \frac{\sum_{i=1}^{t-1} \mathbf{g}_i^{\mathbf{y}}}{t+1} \left( 1 + \sum_{i=1}^{t-1} \langle \tilde{\mathbf{y}}_i, \mathbf{g}_i^{\mathbf{y}} \rangle \right)$
  - 9: **end for**
  - 10: **Return**  $\bar{\mathbf{x}}_T = \frac{1}{T} \sum_{t=1}^T \mathbf{x}_t, \bar{\mathbf{y}}_T = \frac{1}{T} \sum_{t=1}^T \mathbf{y}_t$
- 

**Theorem 3** Assume  $\|\hat{\mathbf{g}}_t\| \leq 1$  for all  $t = 1, \dots, T$ . Then, for all  $\mathbf{u} \in \mathcal{X}$ , we have

$$\sum_{t=1}^T (\ell_t(\mathbf{x}_t) - \ell_t(\mathbf{u})) \leq 2 + 2\|\mathbf{x}_0 - \mathbf{u}\| \sqrt{T \ln(24T^2 \|\mathbf{x}_0 - \mathbf{u}\|^2 + 1)}.$$

As in other parameter-free algorithms, the regret depends on the optimal quantity  $\tilde{\mathcal{O}}(\|\mathbf{x}_0 - \mathbf{u}\|)$ , rather than on the worse  $\mathcal{O}(\|\mathbf{x}_0 - \mathbf{u}\|^2)$  or  $D$  of online gradient descent. Finally, the constrained set is dealt with the black-box reductions from unconstrained OCO to constrained OCO (Cutkosky and Orabona, 2018), that prescribes the use of the projections and surrogate gradients.

Now, using Theorem 2, we can use two instantiations of Algorithm 1 on the primal  $\mathbf{x}$  and dual variable  $\mathbf{y}$  to solve min-max problems. Putting all together, we obtain Algorithm 2.

Using Theorem 3 and Theorem 2, we are able to show the following Corollary (proof is in Appendix A).



**Corollary 1** *Algorithm 2 guarantees that*

$$\begin{aligned} & \max_{\mathbf{y} \in \mathcal{Y}} F(\bar{\mathbf{x}}_T, \mathbf{y}) - \min_{\mathbf{x} \in \mathcal{X}} F(\mathbf{x}, \bar{\mathbf{y}}_T) \\ & \leq \frac{4}{T} + \frac{2\|\mathbf{x}_0 - \mathbf{x}'_T\| \sqrt{\ln(24T^2\|\mathbf{x}_0 - \mathbf{x}'_T\|^2 + 1)} + 2\|\mathbf{y}_0 - \mathbf{y}'_T\| \sqrt{\ln(24T^2\|\mathbf{y}_0 - \mathbf{y}'_T\|^2 + 1)}}{\sqrt{T}}, \end{aligned} \quad (7)$$

where  $\mathbf{x}'_T \in \arg \min_{\mathbf{x} \in \mathcal{X}} F(\mathbf{x}, \bar{\mathbf{y}}_T)$ , and  $\mathbf{y}'_T \in \arg \max_{\mathbf{y} \in \mathcal{Y}} F(\bar{\mathbf{x}}_T, \mathbf{y})$ .

Now, from Corollary 1 and Theorem 1, we show asymptotic initialization-dependent convergence rates for Algorithm 2. The proof of Theorem 4 is included in Appendix A.

**Theorem 4** *Suppose Assumption 1 holds. Then, Algorithm 2 for any  $T$  iterations satisfies*

$$\max_{\mathbf{y} \in \mathcal{Y}} F(\bar{\mathbf{x}}_T, \mathbf{y}) - \min_{\mathbf{x} \in \mathcal{X}} F(\mathbf{x}, \bar{\mathbf{y}}_T) \leq \frac{2(\|\mathbf{x}_0 - \mathbf{x}_*\| + \|\mathbf{y}_0 - \mathbf{y}_*\|) \sqrt{\ln(24T^2D^2 + 1)}}{\sqrt{T}} + o\left(\frac{1}{\sqrt{T}}\right).$$

**Remark:** Comparing Theorem 4 and the argument in Section 4.2, we can see that asymptotically CB-Min-Max has rate  $\frac{\sqrt{\log T}}{\sqrt{T}}$ , while gradient descent ascent has rate  $\frac{1}{\sqrt{T}}$ , so CB-Min-Max is worse than gradient descent ascent asymptotically. However, gradient descent ascent has initialization-dependent rate only if we have the knowledge of the initial condition, while CB-Min-Max always obtains initialization-dependent rate. We will show in Section 5 that this particular feature of CB-Min-Max is important to derive non-asymptotic fast rates on another class of min-max problems.

#### 4.4. Discussion on Possible Alternative Approaches

We have shown how to utilize the strict-convexity-strict-concavity to design algorithms with initialization-dependent rates. However, one may wonder about alternative approaches to achieve the same goal. For example, adding arbitrary small amount of  $\ell_2$  regularization to both primal and dual variables will induce strict-convexity-strict-concavity with negligible bias. However, note that any arbitrarily small bias would not be negligible when the number of iterations  $T$  gets large.

In alternative, one could think to trade off the bias with the rate. But, this approach seems to be able to recover only a  $DG/\sqrt{T}$  rate ( $D$  is the domain size,  $G$  is the gradient upper bound), so it would not be initialization-dependent. Let us consider the function  $F(\mathbf{x}, \mathbf{y})$  with  $\|\mathbf{x}\| \leq D$  and  $\|\mathbf{y}\| \leq D$ . In particular, denote the primal function by  $p(\mathbf{x}) = \max_{\mathbf{y} \in \mathcal{Y}} F(\mathbf{x}, \mathbf{y})$ , the dual function by  $d(\mathbf{y}) = \min_{\mathbf{x} \in \mathcal{X}} F(\mathbf{x}, \mathbf{y})$ . Define the regularized primal function by  $p_\mu(\mathbf{x}) = \max_{\mathbf{y} \in \mathcal{Y}} F(\mathbf{x}, \mathbf{y}) + \mu\|\mathbf{x}\|^2$ , and the regularized dual function by  $d_\mu(\mathbf{y}) = \min_{\mathbf{x} \in \mathcal{X}} F(\mathbf{x}, \mathbf{y}) - \mu\|\mathbf{y}\|^2$ , where  $\mu > 0$  is a very small constant. Since the domain is bounded (assume both  $\mathcal{X}$  and  $\mathcal{Y}$  have diameter  $D$ ), we have  $|d(\mathbf{y}) - d_\mu(\mathbf{y})| \leq \mu D^2$  and  $|p(\mathbf{x}) - p_\mu(\mathbf{x})| \leq \mu D^2$  for any  $\mathbf{x} \in \mathcal{X}$  and any  $\mathbf{y} \in \mathcal{Y}$ . By using the state-of-the-art solver for the strongly-convex-strongly-concave subproblem, we have  $p_\mu(\mathbf{x}_T) - d_\mu(\mathbf{y}_T) \leq \frac{G^2}{\mu T}$ , and hence for the original problem, the duality gap is  $p(\mathbf{x}_T) - d(\mathbf{y}_T) \leq \frac{G^2}{\mu T} + \mu D^2$ . Choosing the optimal  $\mu = \frac{G}{D\sqrt{T}}$ , the rate becomes  $DG/\sqrt{T}$ . Please note that the error term  $\mu D^2$  explicitly depends on the size of the domain, so it seems difficult to get around the domain diameter by the approach of adding a small  $\ell_2$  regularization. In contrast, directly assuming strict-convexity-strict-concavity we can get initialization-dependent asymptotic rates.

---

**Algorithm 3** Restart-CB-Min-Max
 

---

**Input:**  $\widehat{\mathbf{x}}_0, \widehat{\mathbf{y}}_0$   
 1: **for**  $s = 1, \dots, S$  **do**  
 2:    $(\widehat{\mathbf{x}}_s, \widehat{\mathbf{y}}_s) = \text{CB-Min-Max}(\widehat{\mathbf{x}}_{s-1}, \widehat{\mathbf{y}}_{s-1}, T_s)$   
 3: **end for**  
 4: **Return**  $\widehat{\mathbf{x}}_S, \widehat{\mathbf{y}}_S$

---

## 5. Fast Rates under Growth Condition and Hölder Continuous Solution Mapping

In this section, we consider an extension of Algorithm 2 when the function satisfies a growth condition and a Hölder continuous solution mapping, in which we establish improved rates.

We consider the following assumptions.

**Assumption 2 (Growth condition)** (i)  $\|\mathbf{x} - \mathbf{x}_*\| \leq c_1(F(\mathbf{x}, \mathbf{y}_*) - F(\mathbf{x}_*, \mathbf{y}_*))^{\theta_1}$  and  $\|\mathbf{y} - \mathbf{y}_*\| \leq c_2(F(\mathbf{x}_*, \mathbf{y}) - F(\mathbf{x}_*, \mathbf{y}_*))^{\theta_2}$ , where  $c_1 > 0, c_2 > 0, 0 < \theta_1 \leq 1, (\mathbf{x}_*, \mathbf{y}_*)$  is the optimal solution of the original problem (1).

**Assumption 3 (Hölder continuous solution mapping)** Define  $H_1(\mathbf{y}) = \arg \min_{\mathbf{x} \in \mathcal{X}} F(\mathbf{x}, \mathbf{y}), H_2(\mathbf{x}) = \arg \max_{\mathbf{y} \in \mathcal{Y}} F(\mathbf{x}, \mathbf{y})$ .  $\|H_1(\mathbf{y}_1) - H_1(\mathbf{y}_2)\| \leq L_y \|\mathbf{y}_1 - \mathbf{y}_2\|^{\theta_2}$  for any  $\mathbf{y}_1, \mathbf{y}_2 \in \mathcal{Y}$ , and  $\|H_2(\mathbf{x}_1) - H_2(\mathbf{x}_2)\| \leq L_x \|\mathbf{x}_1 - \mathbf{x}_2\|^{\theta_2}$  for any  $\mathbf{x}_1, \mathbf{x}_2 \in \mathcal{X}$ , where  $0 < \theta_2 \leq 1$ .

**Remark:** Assumption 2 is a generalization of growth condition in minimization problems (Li, 2013; Yang and Lin, 2018). Assumption 3 is closely related to the Aubin Property (Aubin, 1984), which is usually employed to characterize the Lipschitz behavior of solution set for convex optimization problems. Examples satisfying this assumption can be found in (Dontchev and Rockafellar, 2009) (e.g., Example 3B.6, Exercise 3C.5). We want to emphasize that it covers a wide range of min-max problems. For example, strongly-convex-strongly-concave min-max problems satisfy Assumption 2 and 3 with  $\theta_1 = 1/2$  and  $\theta_2 = 1$  ( $\theta_1 = 1/2$  is due to the definition of strong-convexity-strong-concavity,  $\theta_2 = 1$  is shown in Lemma 2.2 of Ghadimi and Wang (2018)). Other examples satisfying Assumption 1, 2 and 3 can be found in Appendix F.

Under Assumption 2 and 3, we can design a restart version of the algorithm, which is presented in Algorithm 3. We will prove that Algorithm 3 enjoys faster non-asymptotic convergence rates.

We first provide convergence guarantee for one-stage of Algorithm 3, which is presented in Theorem 5 (the proof is included in Appendix B).

**Theorem 5** Suppose Assumptions 1, 2 and 3 hold. Running Algorithm 2 for  $T$  iterations yields

$$\max_{\mathbf{y} \in \mathcal{Y}} F(\bar{\mathbf{x}}_T, \mathbf{y}) - \min_{\mathbf{x} \in \mathcal{X}} F(\mathbf{x}, \bar{\mathbf{y}}_T) \leq \mathcal{O} \left( \frac{1}{T} + \frac{1}{T^{\frac{\theta+1}{2}}} + \frac{\text{ObjGap}^{\theta_1}(\mathbf{x}_0, \mathbf{y}_0) \ln T}{\sqrt{T}} \right), \quad (8)$$

where  $\text{ObjGap}(\mathbf{x}_0, \mathbf{y}_0) \triangleq F(\mathbf{x}_0, \mathbf{y}_*) - F(\mathbf{x}_*, \mathbf{y}_*) + F(\mathbf{x}_*, \mathbf{y}_*) - F(\mathbf{x}_*, \mathbf{y}_0)$ , and  $\theta \triangleq \theta_1 \theta_2$ .

Based on Theorem 5, we then introduce the Theorem 6, which illustrates the improved rate achieved by Algorithm 3 when  $\theta_1 > 0$  and  $\theta_2 > 0$ . The proof of Theorem 6 is included in Appendix B.

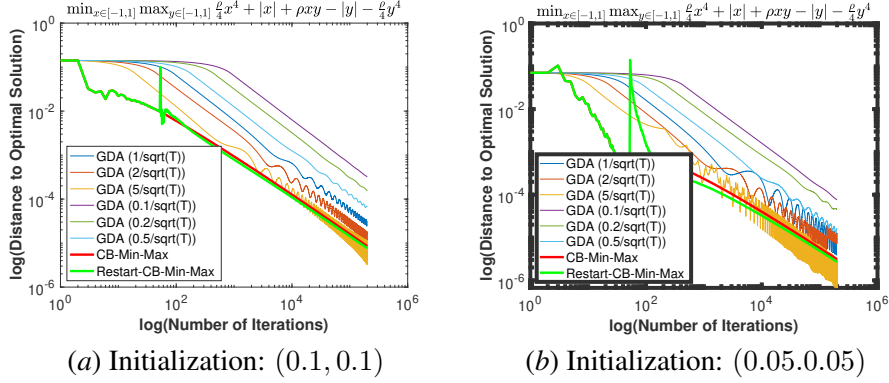


Figure 1: Comparison of different algorithms for the synthetic problem (9). GDA( $\cdot$ ) stands for gradient descent ascent (with learning rate), where  $T$  is total number of iterations. CB-Min-Max stands for Algorithm 2. Restart-CB-Min-Max stands for Algorithm 3.

**Theorem 6** *Suppose Assumptions 1, 2 and 3 hold. Define  $\theta \triangleq \theta_1\theta_2$ . Assume  $\text{ObjGap}(\hat{\mathbf{x}}_0, \hat{\mathbf{y}}_0) \triangleq F(\hat{\mathbf{x}}_0, \mathbf{y}_*) - F(\mathbf{x}_*, \mathbf{y}_*) + F(\mathbf{x}_*, \hat{\mathbf{y}}_0) - F(\mathbf{x}_*, \mathbf{y}_*) \leq \epsilon_0$ . Define  $\epsilon_s = \epsilon_0/2^s$ . Run Algorithm 3 for  $S = \lceil \log(\epsilon_0/\epsilon) \rceil$  stages with  $T_s = \tilde{O}(\max(\frac{1}{2-2\theta_1}, \frac{1}{2/(1+\theta)}))$ . Then, we can guarantee that  $\max_{\mathbf{y} \in \mathcal{Y}} F(\hat{\mathbf{x}}_S, \mathbf{y}) - \min_{\mathbf{x} \in \mathcal{X}} F(\mathbf{x}, \hat{\mathbf{y}}_S) \leq \epsilon$  and the total iteration complexity is  $\tilde{O}(\max(\frac{1}{\epsilon^{2(1-\theta_1)}}, \frac{1}{\epsilon^{2/(1+\theta)}}))$ .*

**Remark** When the function is strongly-convex-strongly-concave, the state-of-the-art complexity is  $O(1/\epsilon)$  while the complexity of our Algorithm 3 is  $\tilde{O}(1/\epsilon^{4/3})$  ( $\theta_1 = 1/2$ ,  $\theta_2 = 1$ ,  $\theta = \theta_1\theta_2 = 1/2$ ), and hence our algorithm does not achieve state-of-the-art in this particular case. However, our analysis captures a more general class of min-max problems, and the complexity is strictly better than  $O(1/\epsilon^2)$  as long as  $\theta_1 > 0$  and  $\theta_2 > 0$ . Intuitively speaking, our algorithm can take advantage of the growth condition of the loss landscape and enjoy faster convergence than the standard algorithm (e.g., GDA) which is oblivious of such favorable structure. To the best of our knowledge, leveraging function growth condition for min-max problems is novel and does not appear in the previous literature.

## 6. Experiments

In this section, we conduct experiments to justify the effectiveness of our proposed algorithm. We first consider a synthetic problem. Another problem is distributionally robust optimization (DRO), which requires an algorithm which can deal with non-Euclidean space such as probability simplex. For lack of space, the algorithm, theory and experiments about DRO are included in Appendix C.

**Synthetic Problem** We consider the following synthetic min-max problem:

$$\min_{x \in \mathcal{X}} \max_{y \in \mathcal{Y}} F(x, y) := \frac{\rho}{4}x^4 + |x| + \rho xy - |y| - \frac{\rho}{4}y^4, \quad (9)$$

where  $\mathcal{X} = \{x \mid |x| \leq R_x\}$  and  $\mathcal{Y} = \{y \mid |y| \leq R_y\}$ ,  $\rho, R_x, R_y$  are all positive constants. One can show that this problem satisfies Assumptions 1, 2, 3 simultaneously with  $\theta_1 = 1$  and  $\theta_2 = 1/3$  (the proof of this claim is included in Appendix F). In addition, the optimal solution of (9) is  $(0, 0)$ .

In our experiment, we set  $R_x = R_y = 1$ ,  $\rho = 0.5$ . We consider two different initializations  $(x_0, y_0) = (0.1, 0.1)$  and  $(x_0, y_0) = (0.05, 0.05)$  to test our algorithms and other baselines. We

choose the same initial point for both primal-dual gradient method with and CB-Min-Max, and report the distance to the optimal solution versus the number of iterations. The learning rate of the gradient descent ascent method (GDA) is set to be  $\frac{c}{G\sqrt{T}}$ , where  $c$  is tuned from  $\{1, 2, 5, 0.1, 0.2, 0.5\}$ ,  $G$  is the gradient’s upper bound and  $T$  is the number of iterations ( $T = 2 \times 10^5$ ). For CB-Min-Max and Restart-CB-Min-Max, we set all gradients in Algorithm 2 to be scaled by its upper bound  $G$  to make sure that the scaled gradient has norm smaller than 1. For Restart-CB-Min-Max, the algorithm restarts at the  $ar^s$ -th iteration, where  $a = 50$ ,  $r = 10^4$ ,  $s = 0, 1, \dots$ . We plot the loglog curves for distance to the optimal solution versus the number of iterations for two different initializations, which are presented in Figure 1 (a) and (b).

From Figure 1, we can see that the algorithms CB-Min-Max and Restart-CB-Min-Max are better than the GDA with theoretically optimal learning rate ( $2/\sqrt{T}$ ), and are comparable to the GDA with best-tuned learning rate ( $5/\sqrt{T}$ ). Note that  $2/\sqrt{T}$  learning rate for GDA is the theoretically best since the domain’s size is  $D = 2$ , while the  $5/\sqrt{T}$  learning rate gives the best performance as shown in Figure 1. Another interesting observation is that Restart-CB-Min-Max is slightly better than CB-Min-Max at the very end, and it fluctuates a little bit in the middle. The reason is due to the restart. When the restart happens, the algorithm enters into a new stage and the update becomes a bit aggressive at the very beginning of this stage, and then it quickly converges to a good solution. Finally, when comparing (a) and (b) in Figure 1, we can see that our algorithms indeed have better performance when the initialization is better, and this is also consistent with our theory.

We want to emphasize that the problem instance (9) is 1 dimension and nonsmooth, so GDA with best-tuned learning rate is the strongest baseline in this case. Other popular algorithms are not applicable or not better than the well-tuned GDA. For example, extragradient method (Korpelevich, 1976) requires the function to be smooth and hence is not applicable to nonsmooth problems such as (9), and coordinate-wise adaptive gradient algorithm (Bach and Levy, 2019) is the same as the nonadaptive version since (9) is of dimension 1.

## 7. Conclusion

In this paper, we consider how to get initialization-dependent convergence rate of first-order algorithms for convex-concave min-max problems. We first identify a condition (i.e., strict-convexity-strict-concavity) and show that it is sufficient to get the initialization-dependent rate. We also take advantage of a parameter-free algorithm with this initialization-dependent rate to design a better algorithm with fast non-asymptotic convergence rate, for min-max problems with a growth condition and Hölder continuous solution mapping. Experimental results show the superior performance of the proposed algorithms. In the future, we plan to study the lower bound for parameter-free algorithms in min-max problems.

## Acknowledgments

This material is based upon work supported by the National Science Foundation under the grants no. 1908111 “AF: Small: Collaborative Research: New Representations for Learning Algorithms and Secure Computation” and no. 2046096 “CAREER: Parameter-free Optimization Algorithms for Machine Learning”.

## References

- J. Abernethy, K. A. Lai, K. Y. Levy, and J.-K. Wang. Faster rates for convex-concave games. In *Conference On Learning Theory*, pages 1595–1625. PMLR, 2018.
- J. D. Abernethy, E. Hazan, and A. Rakhlin. Competing in the dark: An efficient algorithm for bandit linear optimization. In Rocco A. Servedio and Tong Zhang, editors, *Proc. of Conference on Learning Theory (COLT)*, pages 263–274. Omnipress, 2008. URL [https://repository.upenn.edu/cgi/viewcontent.cgi?article=1492&context=statistics\\_papers](https://repository.upenn.edu/cgi/viewcontent.cgi?article=1492&context=statistics_papers).
- Jacob Abernethy, Kevin A Lai, and Andre Wibisono. Last-iterate convergence rates for min-max optimization. *arXiv preprint arXiv:1906.02027*, 2019.
- Jean-Pierre Aubin. Lipschitz behavior of solutions to convex minimization problems. *Mathematics of Operations Research*, 9(1):87–111, 1984.
- Waïss Azizian, Damien Scieur, Ioannis Mitliagkas, Simon Lacoste-Julien, and Gauthier Gidel. Accelerating smooth games by manipulating spectral shapes. In *International Conference on Artificial Intelligence and Statistics*, pages 1705–1715. PMLR, 2020.
- Francis Bach and Kfir Y Levy. A universal algorithm for variational inequalities adaptive to smoothness and noise. In *Conference on Learning Theory*, pages 164–194. PMLR, 2019.
- James P Bailey, Gauthier Gidel, and Georgios Piliouras. Finite regret and cycles with fixed step-size via alternating gradient descent-ascent. In *Conference on Learning Theory*, pages 391–407. PMLR, 2020.
- Aharon Ben-Tal, Laurent El Ghaoui, and Arkadi Nemirovski. *Robust optimization*. Princeton university press, 2009.
- Antonin Chambolle and Thomas Pock. A first-order primal-dual algorithm for convex problems with applications to imaging. *Journal of mathematical imaging and vision*, 40(1):120–145, 2011.
- Yunmei Chen, Guanghui Lan, and Yuyuan Ouyang. Optimal primal-dual methods for a class of saddle point problems. *SIAM Journal on Optimization*, 24(4):1779–1814, 2014.
- Yunmei Chen, Guanghui Lan, and Yuyuan Ouyang. Accelerated schemes for a class of variational inequalities. *Mathematical Programming*, 165(1):113–149, 2017.
- A. Cutkosky and K. Boahen. Online learning without prior information. In *Proc. of the 2017 Conference on Learning Theory*, volume 65 of *Proc. of Machine Learning Research*, pages 643–677, Amsterdam, Netherlands, 07–10 Jul 2017. PMLR. URL <http://proceedings.mlr.press/v65/cutkosky17a.html>.
- A. Cutkosky and K. A. Boahen. Online convex optimization with unconstrained domains and losses. In *Advances in Neural Information Processing Systems*, pages 748–756, 2016.
- A. Cutkosky and F. Orabona. Black-box reductions for parameter-free online learning in Banach spaces. In *Proc. of the Conference on Learning Theory (COLT)*, 2018. URL <https://arxiv.org/abs/1802.06293>.

- A. Cutkosky and T. Sarlos. Matrix-free preconditioning in online learning. In *International Conference on Machine Learning*, pages 1455–1464. PMLR, 2019.
- Constantinos Daskalakis and Ioannis Panageas. Last-iterate convergence: Zero-sum games and constrained min-max optimization. *arXiv preprint arXiv:1807.04252*, 2018.
- Constantinos Daskalakis, Andrew Ilyas, Vasilis Syrgkanis, and Haoyang Zeng. Training GANs with optimism. In *International Conference on Learning Representations*, 2018. URL <https://openreview.net/forum?id=SJJySbbAZ>.
- Jelena Diakonikolas. Halpern iteration for near-optimal and parameter-free monotone inclusion and strong solutions to variational inequalities. In *Conference on Learning Theory*, pages 1428–1451. PMLR, 2020.
- Asen L Dontchev and R Tyrrell Rockafellar. *Implicit functions and solution mappings*, volume 543. Springer, 2009.
- Saeed Ghadimi and Mengdi Wang. Approximation methods for bilevel programming. *arXiv preprint arXiv:1802.02246*, 2018.
- Gauthier Gidel, Reyhane Askari Hemmat, Mohammad Pezeshki, Rémi Le Priol, Gabriel Huang, Simon Lacoste-Julien, and Ioannis Mitliagkas. Negative momentum for improved game dynamics. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pages 1802–1811. PMLR, 2019.
- Noah Golowich, Sarath Pattathil, Constantinos Daskalakis, and Asuman Ozdaglar. Last iterate is slower than averaged iterate in smooth convex-concave saddle point problems. In *Conference on Learning Theory*, pages 1758–1784. PMLR, 2020.
- Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *Advances in neural information processing systems*, pages 2672–2680, 2014.
- G. J. Gordon. Regret bounds for prediction problems. In *Proc. of the twelfth annual conference on Computational learning theory (COLT)*, pages 29–40, 1999. URL <http://citeseerx.ist.psu.edu/viewdoc/download?doi=10.1.1.67.4767&rep=rep1&type=pdf>.
- Yunlong He and Renato DC Monteiro. An accelerated HPE-type algorithm for a class of composite convex-concave saddle-point problems. *SIAM Journal on Optimization*, 26(1):29–56, 2016.
- Y.-G. Hsieh, F. Iutzeler, J. Malick, and P. Mertikopoulos. Explore aggressively, update conservatively: Stochastic extragradient methods with variable stepsize scaling. In H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 16223–16234. Curran Associates, Inc., 2020. URL <https://proceedings.neurips.cc/paper/2020/file/ba9a56ce0a9bfa26e8ed9e10b2cc8f46-Paper.pdf>.
- Adam Ibrahim, Waiss Azizian, Gauthier Gidel, and Ioannis Mitliagkas. Linear lower bounds and conditioning of differentiable games. In *International Conference on Machine Learning*, pages 4583–4593. PMLR, 2020.

- Anatoli Juditsky, Arkadi Nemirovski, Claire Tauvel, et al. Solving variational inequalities with stochastic mirror-prox algorithm. *Stochastic Systems*, 1(1):17–58, 2011.
- M. Kempka, W. Kotłowski, and M. K. Warmuth. Adaptive scale-invariant online algorithms for learning linear models. In K. Chaudhuri and R. Salakhutdinov, editors, *Proc. of the 36th International Conference on Machine Learning*, volume 97 of *Proc. of Machine Learning Research*, pages 3321–3330, Long Beach, California, USA, 09–15 Jun 2019. PMLR.
- GM Korpelevich. The extragradient method for finding saddle points and other problems. *Matecon*, 12:747–756, 1976.
- W. Kotłowski. Scale-invariant unconstrained online learning. *Theoretical Computer Science*, 808: 139–158, 2020.
- Guoyin Li. Global error bounds for piecewise convex polynomials. *Mathematical Programming*, pages 1–28, 2013.
- Tengyuan Liang and James Stokes. Interaction matters: A note on non-asymptotic local convergence of generative adversarial networks. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pages 907–915. PMLR, 2019.
- Tianyi Lin, Chi Jin, and Michael I Jordan. Near-optimal algorithms for minimax optimization. In *Conference on Learning Theory*, pages 2738–2779. PMLR, 2020.
- Mingrui Liu, Xiaoxuan Zhang, Zaiyi Chen, Xiaoyu Wang, and Tianbao Yang. Fast stochastic AUC maximization with  $O(1/n)$ -convergence rate. In *International Conference on Machine Learning*, pages 3195–3203, 2018.
- Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards deep learning models resistant to adversarial attacks. *arXiv preprint arXiv:1706.06083*, 2017.
- H. B. McMahan and F. Orabona. Unconstrained online linear learning in Hilbert spaces: Minimax algorithms and normal approximations. In *Proc of the Annual Conference on Learning Theory, COLT*, 2014. URL <https://arxiv.org/abs/1403.0628>.
- Aryan Mokhtari, Asuman Ozdaglar, and Sarath Pattathil. A unified analysis of extra-gradient and optimistic gradient methods for saddle point problems: Proximal point approach. In *International Conference on Artificial Intelligence and Statistics*, pages 1497–1507. PMLR, 2020.
- Renato DC Monteiro and Benar Fux Svaiter. On the complexity of the hybrid proximal extragradient method for the iterates and the ergodic mean. *SIAM Journal on Optimization*, 20(6):2755–2787, 2010.
- Angelia Nedić and Asuman Ozdaglar. Subgradient methods for saddle-point problems. *Journal of optimization theory and applications*, 142(1):205–228, 2009.
- Arkadi Nemirovski. Prox-method with rate of convergence  $O(1/t)$  for variational inequalities with lipschitz continuous monotone operators and smooth convex-concave saddle point problems. *SIAM Journal on Optimization*, 15(1):229–251, 2004.

- Arkadi Nemirovski and David Yudin. Problem complexity and method efficiency in optimization. 1983.
- Arkadi Nemirovski, Anatoli Juditsky, Guanghui Lan, and Alexander Shapiro. Robust stochastic approximation approach to stochastic programming. *SIAM Journal on optimization*, 19(4):1574–1609, 2009.
- Y. Nesterov. Primal-dual subgradient methods for convex problems. *Mathematical programming*, 120(1):221–259, 2009. URL <https://link.springer.com/content/pdf/10.1007/s10107-007-0149-x.pdf>.
- Yu Nesterov. Smooth minimization of non-smooth functions. *Mathematical programming*, 103(1):127–152, 2005.
- Yurii Nesterov. Dual extrapolation and its applications to solving variational inequalities and related problems. *Mathematical Programming*, 109(2):319–344, 2007.
- Yurii Nesterov and Laura Scramali. Solving strongly monotone variational and quasi-variational inequalities. Available at SSRN 970903, 2006.
- F. Orabona. Simultaneous model selection and optimization through parameter-free stochastic learning. In *Advances in Neural Information Processing Systems 27*, 2014.
- F. Orabona. A modern introduction to online learning. *arXiv preprint arXiv:1912.13213*, 2019.
- F. Orabona and D. Pál. Coin betting and parameter-free online learning. In D. D. Lee, M. Sugiyama, U. V. Luxburg, I. Guyon, and R. Garnett, editors, *Advances in Neural Information Processing Systems 29*, pages 577–585. Curran Associates, Inc., 2016. URL <https://arxiv.org/pdf/1602.04128.pdf>.
- Yuyuan Ouyang and Yangyang Xu. Lower complexity bounds of first-order methods for convex-concave bilinear saddle-point problems. *Mathematical Programming*, 185(1):1–35, 2021.
- S. Shalev-Shwartz. *Online Learning: Theory, Algorithms, and Applications*. PhD thesis, The Hebrew University, 2007. URL <https://www.cs.huji.ac.il/~shais/papers/ShalevThesis07.pdf>.
- M. Streeter and B. McMahan. No-regret algorithms for unconstrained online convex optimization. In *Advances in Neural Information Processing Systems 25*, pages 2402–2410. Curran Associates, Inc., 2012. URL <https://arxiv.org/pdf/1211.2260.pdf>.
- Yan Yan, Yi Xu, Qihang Lin, Wei Liu, and Tianbao Yang. Optimal epoch stochastic gradient descent ascent methods for min-max optimization. *Advances in Neural Information Processing Systems*, 33, 2020.
- Tianbao Yang and Qihang Lin. RSG: Beating subgradient method without smoothness and strong convexity. *The Journal of Machine Learning Research*, 19(1):236–268, 2018.
- Yiming Ying, Longyin Wen, and Siwei Lyu. Stochastic online AUC maximization. In *Advances in Neural Information Processing Systems*, pages 451–459, 2016.



Junyu Zhang, Mingyi Hong, and Shuzhong Zhang. On lower iteration complexity bounds for the saddle point problems. *arXiv preprint arXiv:1912.07481*, 2019.

Renbo Zhao. Optimal stochastic algorithms for convex-concave saddle-point problems. *arXiv preprint arXiv:1903.01687*, 2019.

M. Zinkevich. Online convex programming and generalized infinitesimal gradient ascent. In *Proc. of the International Conference on Machine Learning*, pages 928–936, 2003. URL <https://www.aaai.org/Papers/ICML/2003/ICML03-120.pdf>.

## Appendix

### Appendix A. Proofs in Section 4

#### A.1. Proof of Theorem 2

##### Proof

Note that  $\ell_t(\mathbf{x}) = F(\mathbf{x}, \mathbf{y}_t)$ , and  $h_t(\mathbf{y}) = -F(\mathbf{x}_t, \mathbf{y})$ . By Jensen's inequality, we have

$$\begin{aligned} F(\bar{\mathbf{x}}_T, \mathbf{y}) - F(\mathbf{x}, \bar{\mathbf{y}}_T) &\leq \frac{1}{T} \sum_{t=1}^T F(\mathbf{x}_t, \mathbf{y}) - \frac{1}{T} \sum_{t=1}^T F(\mathbf{x}, \mathbf{y}_t) \\ &= \frac{1}{T} \sum_{t=1}^T F(\mathbf{x}_t, \mathbf{y}) - \frac{1}{T} \sum_{t=1}^T F(\mathbf{x}_t, \mathbf{y}_t) + \frac{1}{T} \sum_{t=1}^T F(\mathbf{x}_t, \mathbf{y}_t) - \frac{1}{T} \sum_{t=1}^T F(\mathbf{x}, \mathbf{y}_t) \\ &= \frac{1}{T} \sum_{t=1}^T (h_t(\mathbf{y}_t) - h_t(\mathbf{y})) + \frac{1}{T} \sum_{t=1}^T (\ell_t(\mathbf{x}_t) - \ell_t(\mathbf{x})). \end{aligned}$$

Taking  $\mathbf{x} = \mathbf{x}'_T \in \arg \min_{\mathbf{x} \in \mathcal{X}} F(\mathbf{x}, \bar{\mathbf{y}}_T)$  and  $\mathbf{y} = \mathbf{y}'_T \in \arg \max_{\mathbf{y} \in \mathcal{Y}} F(\bar{\mathbf{x}}_T, \mathbf{y})$ , we get the stated result.  $\blacksquare$

#### A.2. Proof of Theorem 3

**Proof** Define  $\tilde{\ell}_t(\mathbf{x}) = \frac{1}{2} (\langle \hat{\mathbf{g}}_t, \mathbf{x} \rangle + \|\hat{\mathbf{g}}_t\| \|\mathbf{x} - \Pi_{\mathcal{X}}(\mathbf{x})\|)$ . Noting that  $\mathbf{g}_t \in \partial \tilde{\ell}_t(\tilde{\mathbf{x}}_t)$ , so Algorithm 1 is a coin-betting algorithm for minimizing regret defined by the function  $\tilde{\ell}$  over the sequence of points  $\{\tilde{\mathbf{x}}_t\}$ . By Corollary 5 in [Cutkosky and Orabona \(2018\)](#), we know that

$$\sum_{t=1}^T (\tilde{\ell}_t(\tilde{\mathbf{x}}_t) - \tilde{\ell}_t(\mathbf{u})) \leq 1 + \|\mathbf{x}_0 - \mathbf{u}\| \sqrt{T \ln(24T^2 \|\mathbf{x}_0 - \mathbf{u}\|^2 + 1)}.$$

By the black-box reduction argument (Theorem 3 in [Cutkosky and Orabona \(2018\)](#)), we have

$$\sum_{t=1}^T (\ell_t(\mathbf{x}_t) - \ell_t(\mathbf{u})) \leq 2 \left( \sum_{t=1}^T (\tilde{\ell}_t(\tilde{\mathbf{x}}_t) - \tilde{\ell}_t(\mathbf{u})) \right) \leq 2 + 2\|\mathbf{x}_0 - \mathbf{u}\| \sqrt{T \ln(24T^2 \|\mathbf{x}_0 - \mathbf{u}\|^2 + 1)},$$

that completes the proof.  $\blacksquare$

#### A.3. Proof of Corollary 1

**Proof** Using Theorem 3 twice (applying it with  $\ell_t = F(\cdot, \mathbf{y}_t)$  and  $\mathbf{u} = \mathbf{x}'_T \in \arg \min_{\mathbf{x} \in \mathcal{X}} F(\mathbf{x}, \bar{\mathbf{y}}_T)$ , and applying it with  $h_t = -F(\mathbf{x}_t, \cdot)$  with  $\mathbf{u} = \mathbf{y}'_T \in \arg \max_{\mathbf{y} \in \mathcal{Y}} F(\bar{\mathbf{x}}_T, \mathbf{y})$ ) and Theorem 2 concludes the proof.  $\blacksquare$

#### A.4. Proof of Theorem 4

**Proof** We can see that (7) satisfies the premises of Theorem 1 with  $A(x) = x$ ,  $B(T) = \frac{\sqrt{4\ln(24T^2D^2+1)}}{\sqrt{T}}$ ,  $C(T) = 4/\sqrt{T}$ . Note that  $A(x+y) \leq A(x) + A(y)$  for any  $0 \leq x \leq D, 0 \leq y \leq D$ , then invoking Theorem 1, we have

$$\begin{aligned} \max_{\mathbf{y} \in \mathcal{Y}} F(\bar{\mathbf{x}}_T, \mathbf{y}) - \min_{\mathbf{x} \in \mathcal{X}} F(\mathbf{x}, \bar{\mathbf{y}}_T) &\leq [A(\|\mathbf{x}_0 - \mathbf{x}_*\|)B(T) + A(\|\mathbf{y}_0 - \mathbf{y}_*\|)B(T)] + R(T) + \frac{4}{T} \\ &\leq \frac{2(\|\mathbf{x}_0 - \mathbf{x}_*\| + \|\mathbf{y}_0 - \mathbf{y}_*\|) \sqrt{\ln(24T^2D^2+1)}}{\sqrt{T}} + o\left(\frac{1}{\sqrt{T}}\right). \end{aligned}$$

■

## Appendix B. Proof in Section 5

### B.1. Proof of Theorem 5

**Proof** By Corollary 1, we have

$$\begin{aligned} \max_{\mathbf{y} \in \mathcal{Y}} F(\bar{\mathbf{x}}_T, \mathbf{y}) - \min_{\mathbf{x} \in \mathcal{X}} F(\mathbf{x}, \bar{\mathbf{y}}_T) \\ \leq \frac{4}{T} + \frac{\|\mathbf{x}_0 - \mathbf{x}'_T\| \sqrt{\ln(1 + 24\|\mathbf{x}_0 - \mathbf{x}'_T\|^2 T^2)}}{\sqrt{T}} + \frac{\|\mathbf{y}_0 - \mathbf{y}'_T\| \sqrt{\ln(1 + 24\|\mathbf{y}_0 - \mathbf{y}'_T\|^2 T^2)}}{\sqrt{T}}, \end{aligned} \quad (10)$$

where  $\mathbf{x}'_T = \arg \min_{\mathbf{x} \in \mathcal{X}} F(\mathbf{x}, \bar{\mathbf{y}}_T)$ ,  $\mathbf{y}'_T = \arg \min_{\mathbf{y} \in \mathcal{Y}} F(\bar{\mathbf{x}}_T, \mathbf{y})$ . By the growth condition in Assumption 2, we know that

$$\|\mathbf{x}_0 - \mathbf{x}_*\| \leq c_1 (F(\mathbf{x}_0, \mathbf{y}_*) - F(\mathbf{x}_*, \mathbf{y}_*))^{\theta_1}. \quad (11)$$

By the Hölder continuity of the solution mapping (Assumption 3), we have  $\|\mathbf{x}_* - \mathbf{x}'_T\| = \|\arg \min_{\mathbf{x}} F(\mathbf{x}, \mathbf{y}_*) - \arg \min_{\mathbf{x}} F(\mathbf{x}, \bar{\mathbf{y}}_T)\| \leq L_{\mathbf{y}} \|\bar{\mathbf{y}}_T - \mathbf{y}_*\|^{\theta_2}$ . Note that we have  $\|\bar{\mathbf{y}}_T - \mathbf{y}_*\| \leq c_2 (F(\mathbf{x}_*, \mathbf{y}_*) - F(\mathbf{x}_*, \bar{\mathbf{y}}_T))^{\theta_1}$ , so we have

$$\begin{aligned} \|\mathbf{x}_* - \mathbf{x}'_T\| &\leq c_2^{\theta_2} L_{\mathbf{y}} (F(\mathbf{x}_*, \mathbf{y}_*) - F(\mathbf{x}_*, \bar{\mathbf{y}}_T))^{\theta_1 \theta_2} \leq c_{\mathbf{y}} \left( F(\mathbf{x}_*, \mathbf{y}_*) - \min_{\mathbf{x}} F(\mathbf{x}, \bar{\mathbf{y}}_T) \right)^{\theta} \\ &\leq O\left( \frac{c_{\mathbf{y}} D_{\mathcal{Y}}^{\theta} \ln^{\frac{\theta}{2}}(D_{\mathcal{Y}} T)}{T^{\theta/2}} \right). \end{aligned} \quad (12)$$

where  $c_{\mathbf{y}} = c_2^{\theta_2} L_{\mathbf{y}}$ ,  $\theta = \theta_1 \theta_2$ , and the last inequality holds due to the convergence guarantee established in Corollary 1.

A parallel argument in terms of  $\mathbf{y}$  gives

$$\|\mathbf{y}_0 - \mathbf{y}_*\| \leq c_2 (F(\mathbf{x}_*, \mathbf{y}_*) - F(\mathbf{x}_*, \mathbf{y}_0))^{\theta_1}, \quad (13)$$

and

$$\begin{aligned} \|\mathbf{y}_* - \mathbf{y}'_T\| &\leq c_1^{\theta_1} L_{\mathbf{x}} (F(\bar{\mathbf{x}}_T, \mathbf{y}_*) - F(\mathbf{x}_*, \mathbf{y}_*))^{\theta_1 \theta_2} \leq c_{\mathbf{x}} (F(\bar{\mathbf{x}}_T, \mathbf{y}_*) - F(\mathbf{x}_*, \mathbf{y}_*))^{\theta} \\ &\leq O\left( \frac{c_{\mathbf{x}} D_{\mathcal{X}}^{\theta} \ln^{\frac{\theta}{2}}(D_{\mathcal{X}} T)}{T^{\theta/2}} \right), \end{aligned} \quad (14)$$

where  $c_{\mathbf{x}} = c_1^{\theta_1} L_{\mathbf{x}}$ .

Due to triangle inequality and (11) and (12), we have

$$\|\mathbf{x}_0 - \mathbf{x}'_T\| \leq \|\mathbf{x}_0 - \mathbf{x}_*\| + \|\mathbf{x}_* - \mathbf{x}'_T\| \leq c_1(F(\mathbf{x}_0, \mathbf{y}_*) - F(\mathbf{x}_*, \mathbf{y}_*))^{\theta_1} + O(\ln T/T^{\theta/2}). \quad (15)$$

Similarly, we have

$$\|\mathbf{y}_0 - \mathbf{y}'_T\| \leq \|\mathbf{y}_0 - \mathbf{y}_*\| + \|\mathbf{y}_* - \mathbf{y}'_T\| \leq c_2(F(\mathbf{x}_*, \mathbf{y}_*) - F(\mathbf{x}_*, \mathbf{y}_0))^{\theta_1} + O(\ln T/T^{\theta/2}). \quad (16)$$

Hence, we have

$$\|\mathbf{x}_0 - \mathbf{x}'_T\| + \|\mathbf{y}_0 - \mathbf{y}'_T\| \leq O\left(\text{ObjGap}^{\theta_1}(\mathbf{x}_0, \mathbf{y}_0)\right) + O(\ln T/T^{\theta/2}). \quad (17)$$

Combining (10) with (17), we have

$$\max_{\mathbf{y} \in \mathcal{Y}} F(\bar{\mathbf{x}}_T, \mathbf{y}) - \min_{\mathbf{x} \in \mathcal{X}} F(\mathbf{x}, \bar{\mathbf{y}}_T) \leq O\left(\frac{1}{T} + \frac{\ln T}{T^{\frac{\theta+1}{2}}} + \frac{\text{ObjGap}^{\theta_1}(\mathbf{x}_0, \mathbf{y}_0) \ln T}{\sqrt{T}}\right), \quad (18)$$

where  $\text{ObjGap}(\mathbf{x}_0, \mathbf{y}_0) = F(\mathbf{x}_0, \mathbf{y}_*) - F(\mathbf{x}_*, \mathbf{y}_*) + F(\mathbf{x}_*, \mathbf{y}_*) - F(\mathbf{x}_*, \mathbf{y}_0)$ .  $\blacksquare$

## B.2. Proof of Theorem 6

**Proof** Define  $\text{DualityGap}(\hat{\mathbf{x}}_0, \hat{\mathbf{y}}_0) = \max_{\mathbf{y} \in \mathcal{Y}} F(\hat{\mathbf{x}}_0, \mathbf{y}) - \min_{\mathbf{x} \in \mathcal{X}} F(\mathbf{x}, \hat{\mathbf{y}}_0)$ . We know that  $\text{ObjGap}(\hat{\mathbf{x}}_0, \hat{\mathbf{y}}_0) \leq \text{DualityGap}(\hat{\mathbf{x}}_0, \hat{\mathbf{y}}_0)$ . By invoking the subroutine Algorithm 2 to run  $T_0 = \tilde{O}\left(\max\left(\frac{1}{\epsilon_0^{2-2\theta_1}}, \frac{1}{\epsilon_0^{2/(1+\theta)}}\right)\right)$  iterations, we know that the duality gap at the new point will be decreased to  $\epsilon_1 = \epsilon_0/2$ . Then the Algorithm 3 restarts by setting the new point as the initial point, and then invokes the subroutine Algorithm 2 to run  $T_1 = \tilde{O}\left(\max\left(\frac{1}{\epsilon_1^{2-2\theta_1}}, \frac{1}{\epsilon_1^{2/(1+\theta)}}\right)\right)$  number of iterations, and then it restarts again. Algorithm 3 repeats this process under it reaches  $\epsilon$ -duality gap. We know that we have  $S = \lfloor \log(\epsilon_0/\epsilon) \rfloor$  stages and hence the total complexity is  $\sum_{s=0}^S \tilde{O}\left(\max\left(\frac{1}{\epsilon_s^{2-2\theta_1}}, \frac{1}{\epsilon_s^{2/(1+\theta)}}\right)\right) = \tilde{O}\left(\max\left(\frac{1}{\epsilon^{2(1-\theta_1)}}, \frac{1}{\epsilon^{2/(1+\theta)}}\right)\right)$ .  $\blacksquare$

## Appendix C. Non-Euclidean Space: Simplex Setup

Min-max optimization in Non-Euclidean spaces is an important topic, which has broad applications in machine learning (e.g., distributionally robust optimization). In this section, we focus on the simplex setup by considering the following problem

$$\min_{\mathbf{x} \in \mathcal{X}} \max_{\mathbf{p} \in \Delta_n} F(\mathbf{x}, \mathbf{p}),$$

where

$$\Delta_n = \left\{ (p_1, p_2, \dots, p_n) \in \mathbb{R}^n \mid 0 \leq p_i \leq 1, \sum_{i=1}^n p_i = 1 \right\}.$$

Our algorithm design shares the similar spirit as in the Euclidean case (Algorithm 2), in which we split the original problem into the primal and dual problems, and employ corresponding coin-betting

---

**Algorithm 4** Simplex-Constrained Coin-betting OCO
 

---

**Input:** Simplex  $\Delta_n \in \mathbb{R}^n$ , prior distribution  $\mathbf{p}_0 \in \Delta_n$

- 1: **for**  $t = 1, \dots, T$  **do**
  - 2:    $\mathbf{w}_{t,i} = \frac{\sum_{j=1}^{t-1} \tilde{\mathbf{g}}_{j,i}}{t} (1 + \sum_{j=1}^{t-1} \tilde{\mathbf{g}}_{j,i} \mathbf{w}_{j,i}), i = 1, \dots, n$
  - 3:    $\hat{\mathbf{p}}_{t,i} = \mathbf{p}_{0,i} \max(\mathbf{w}_{t,i}, 0), i = 1, \dots, n$
  - 4:    $\mathbf{p}_t = \frac{\hat{\mathbf{p}}_t}{\|\hat{\mathbf{p}}_t\|_1}$  if  $\|\hat{\mathbf{p}}_t\|_1 \neq 0$  else  $\mathbf{p}_t = \mathbf{p}_0$
  - 5:   Receive subgradient  $\hat{\mathbf{g}}_t^{\mathbf{P}}$
  - 6:    $\tilde{\mathbf{g}}_{t,i} = \hat{\mathbf{g}}_{t,i}^{\mathbf{P}} - \langle \hat{\mathbf{g}}_t^{\mathbf{P}}, \mathbf{p}_t \rangle$  if  $\mathbf{w}_{t,i} > 0$  else  $\tilde{\mathbf{g}}_{t,i} = \max(\hat{\mathbf{g}}_{t,i}^{\mathbf{P}} - \langle \hat{\mathbf{g}}_t^{\mathbf{P}}, \mathbf{p}_t \rangle, 0), i = 1, \dots, n$
  - 7: **end for**
  - 8: **Return**  $\bar{\mathbf{p}}_T = \frac{1}{T} \sum_{t=1}^T \mathbf{p}_t$
- 

---

**Algorithm 5** CB-Min-Max-Simplex
 

---

**Input:**  $\mathbf{x}_0 \in \mathcal{X}$ , prior distribution  $\mathbf{p}_0 \in \Delta_n$

- 1:  $\tilde{\mathbf{x}}_0 = \mathbf{x}_0$
  - 2: **for**  $t = 0, \dots, T$  **do**
  - 3:    $\mathbf{x}_t = \Pi_{\mathcal{X}}(\tilde{\mathbf{x}}_t)$
  - 4:   Receive Subgradient  $\hat{\mathbf{g}}_t^{\mathbf{x}}$
  - 5:    $\mathbf{g}_t^{\mathbf{x}} = \frac{1}{2} \left( \hat{\mathbf{g}}_t^{\mathbf{x}} + \|\hat{\mathbf{g}}_t^{\mathbf{x}}\| \cdot \frac{\tilde{\mathbf{x}}_t - \mathbf{x}_t}{\|\tilde{\mathbf{x}}_t - \mathbf{x}_t\|} \right)$  (Define  $\mathbf{0}/0 = \mathbf{0}$ )
  - 6:    $\mathbf{w}_{t,i} = \frac{\sum_{j=1}^{t-1} \tilde{\mathbf{g}}_{j,i}}{t} (1 + \sum_{j=1}^{t-1} \tilde{\mathbf{g}}_{j,i} \mathbf{w}_{j,i}), i = 1, \dots, n$
  - 7:    $\hat{\mathbf{p}}_{t,i} = \mathbf{p}_{0,i} \max(\mathbf{w}_{t,i}, 0), i = 1, \dots, n$
  - 8:    $\mathbf{p}_t = \frac{\hat{\mathbf{p}}_t}{\|\hat{\mathbf{p}}_t\|_1}$  if  $\|\hat{\mathbf{p}}_t\|_1 \neq 0$  else  $\mathbf{p}_t = \mathbf{p}_0$
  - 9:   Receive subgradient  $\hat{\mathbf{g}}_t^{\mathbf{P}}$
  - 10:    $\tilde{\mathbf{g}}_{t,i} = \hat{\mathbf{g}}_{t,i}^{\mathbf{P}} - \langle \hat{\mathbf{g}}_t^{\mathbf{P}}, \mathbf{p}_t \rangle$  if  $\mathbf{w}_{t,i} > 0$  else  $\tilde{\mathbf{g}}_{t,i} = \max(\hat{\mathbf{g}}_{t,i}^{\mathbf{P}} - \langle \hat{\mathbf{g}}_t^{\mathbf{P}}, \mathbf{p}_t \rangle, 0), i = 1, \dots, n$
  - 11:    $\tilde{\mathbf{x}}_{t+1} = \mathbf{x}_0 - \frac{\sum_{j=1}^t \mathbf{g}_j^{\mathbf{x}}}{t+1} (1 - \sum_{j=1}^t \langle \mathbf{g}_j^{\mathbf{x}}, \tilde{\mathbf{x}}_j \rangle)$
  - 12: **end for**
  - 13: **Return**  $\bar{\mathbf{x}}_T = \frac{1}{T} \sum_{t=1}^T \mathbf{x}_t, \bar{\mathbf{p}}_T = \frac{1}{T} \sum_{t=1}^T \mathbf{p}_t$
- 

algorithms for solving them simultaneously. The algorithm is presented in Algorithm 5, which is an synthesis of Algorithm 1 and Algorithm 4. It is worth mentioning that Algorithm 4 is adapted from Algorithm 2 in Orabona and Pál (2016), which is a parameter-free coin-betting algorithm in the probability simplex.

The main difficulty of analyzing parameter-free algorithms for min-max problems in a non-Euclidean space is that we have to use a different distance-generating function, and generalize the proofs done in Euclidean space in Section 4 to non-Euclidean spaces.

We first present two key Lemmas (Lemma 1 and Lemma 2)) which are useful for our analysis. The proofs are included in Appendix D. Lemma 1 plays the role of the triangle inequality for the KL-divergence, and Lemma 2 provides the non-Euclidean variant of Theorem 3.

**Lemma 1** *Assume  $\pi$  is a discrete probability distribution over  $\mathbb{R}^n$  with  $\pi_i > 0$ . Let  $\mathbf{p}, \mathbf{q}$  arbitrary vectors in the probability simplex. Then,*

$$KL(\mathbf{p}, \pi) \leq KL(\mathbf{p}, \mathbf{q}) + \max_i \max \left( \ln \frac{q_i}{\pi_i}, 0 \right) \|\mathbf{p} - \mathbf{q}\|_1 + KL(\mathbf{q}, \pi). \quad (19)$$

**Lemma 2** Define  $h_t(\mathbf{p}) = -F(\mathbf{x}_t, \mathbf{p})$ . Then, Algorithm 5 guarantees that for any  $\mathbf{p} \in \Delta_n$ , we have

$$\frac{1}{T} \sum_{t=1}^T h_t(\mathbf{p}_t) - \frac{1}{T} \sum_{t=1}^T h_t(\mathbf{p}) \leq \frac{1}{T} + \frac{\sqrt{\ln T + 3KL(\mathbf{p}, \mathbf{p}_0)}}{\sqrt{T}}. \quad (20)$$

Here, we present our main theorem in this section.

**Theorem 7** Suppose Assumption 1 holds and  $\mathbf{p}_{0,i} > 0$  for every  $i = 1, \dots, n$ . Then, Algorithm 5 guarantees that

$$\begin{aligned} & \max_{\mathbf{p} \in \Delta_n} F(\bar{\mathbf{x}}_T, \mathbf{p}) - \min_{\mathbf{x} \in \mathcal{X}} F(\mathbf{x}, \bar{\mathbf{p}}_T) & (21) \\ & \leq \frac{4}{T} + \frac{\|\mathbf{x}_0 - \mathbf{x}_*\| \sqrt{\ln(1 + 24\|\mathbf{x}_0 - \mathbf{x}_*\|T^2)}}{\sqrt{T}} + \frac{\sqrt{\ln T + 3KL(\mathbf{p}_*, \mathbf{p}_0)}}{\sqrt{T}} + o\left(\frac{1}{\sqrt{T}}\right). & (22) \end{aligned}$$

The high-level proof idea of Theorem 7 is to extend the proof of Theorem 4 from the Euclidean setting to the probability simplex. The main technical ingredient is utilizing the expansion of the KL divergence in Lemma 1 and making use of the objective gap bound in Lemma 2. The detailed proof can be found in Appendix D.

## Appendix D. Proofs in Appendix C

### D.1. Proof of Lemma 1

**Proof** We have

$$\begin{aligned} KL(\mathbf{p}, \boldsymbol{\pi}) &= \sum_{i=1}^d p_i \ln \frac{p_i}{\pi_i} = \sum_{i=1}^d p_i \ln \frac{p_i}{q_i} + \sum_{i=1}^d p_i \ln \frac{q_i}{\pi_i} \\ &= KL(\mathbf{q}, \boldsymbol{\pi}) + \sum_{i=1}^d q_i \ln \frac{q_i}{\pi_i} + \sum_{i=1}^d (p_i - q_i) \ln \frac{q_i}{\pi_i} & (23) \\ &\leq KL(\mathbf{q}, \boldsymbol{\pi}) + \sum_{i=1}^d q_i \ln \frac{q_i}{\pi_i} + \max_i \max(\ln \frac{q_i}{\pi_i}, 0) \sum_{i=1}^d |p_i - q_i|. \end{aligned}$$

■

### D.2. Proof of Lemma 2

**Proof** The proof follows Corollary 6 in Orabona and Pál (2016) and Jensen's inequality in terms of  $\mathbf{p}$ . ■

### D.3. Proof of Theorem 7

**Proof** By combining Lemma 2, Theorem 3 and Lemma 2, we can get the following bound

$$\begin{aligned} & \max_{\mathbf{p} \in \Delta_n} F(\bar{\mathbf{x}}_T, \mathbf{p}) - \min_{\mathbf{x} \in \mathcal{X}} F(\mathbf{x}, \bar{\mathbf{p}}_T) \\ & \leq \frac{4}{T} + \frac{\|\mathbf{x}_0 - \mathbf{x}'_T\| \sqrt{\ln(1 + 24\|\mathbf{x}_0 - \mathbf{x}_*\|^2 T^2)}}{\sqrt{T}} + \frac{\sqrt{\ln T + 3KL(\mathbf{p}'_T, \mathbf{p}_0)}}{\sqrt{T}}, \end{aligned} \quad (24)$$

where  $\mathbf{x}'_T = \arg \min_{\mathbf{x} \in \mathcal{X}} F(\mathbf{x}, \bar{\mathbf{p}}_T)$  and  $\mathbf{p}'_T = \arg \max_{\mathbf{p} \in \Delta_n} F(\bar{\mathbf{x}}_T, \mathbf{p})$ . Note that (24) is the counterpart of Corollary 1 in the probability simplex setup. By Lemma 2 (taking  $\boldsymbol{\pi} = \mathbf{p}_0$ ,  $\mathbf{p} = \mathbf{p}'_T$  and  $\mathbf{q} = \mathbf{p}_*$ ), we know that

$$KL(\mathbf{p}'_T, \mathbf{p}_0) \leq KL(\mathbf{p}'_T, \mathbf{p}_*) + \max_i \max \left( \ln \frac{\mathbf{p}_{*,i}}{\mathbf{p}_{0,i}}, 0 \right) \|\mathbf{p}'_T - \mathbf{p}_*\|_1 + KL(\mathbf{p}_*, \mathbf{p}_0). \quad (25)$$

Then, we are able to follow the proof of Theorem 1 to show that  $KL(\mathbf{p}'_T, \mathbf{p}_*) \rightarrow 0$  (this implies that  $\max_i \max \left( \ln \frac{\mathbf{p}_{*,i}}{\mathbf{p}_{0,i}}, 0 \right) \|\mathbf{p}'_T - \mathbf{p}_*\|_1 \rightarrow 0$  by Pinsker's inequality). In addition, it is proved in Theorem 1 that  $\mathbf{x}'_T \rightarrow \mathbf{x}_*$ . Combining these facts with (25), the theorem is proved.  $\blacksquare$

## Appendix E. More Experimental Results

### Detailed Experimental Settings

- SensIT Vehicle (combined): all the classes with label 2 and 3 are regarded as class  $-1$ , and the class 1 is regarded as class 1.
- dna: all the classes with label 2 and 3 are regarded as class  $-1$ , and the class 1 is regarded as class 1.
- gisette: we use the original dataset without any preprocessing.
- protein: all the classes with label 0 and 2 are regarded as class  $-1$ , and the class 1 is regarded as class 1.
- letter: all the classes with label 1 to label 25 are regarded as class  $-1$ , and the class 26 is regarded as class 1.
- mnist: all the classes with label 0 to label 8 are regarded as class  $-1$ , and the class 9 is regarded as class 1.
- madelon: we use the original dataset without any preprocessing.
- pendigits: all the classes with label 0 to label 8 are regarded as class  $-1$ , and the class 9 is regarded as class 1.

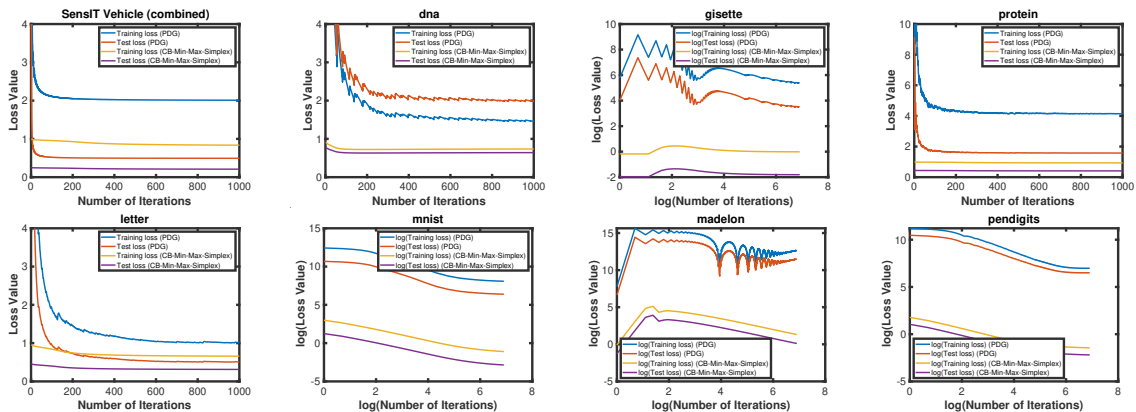


Figure 2: Comparison of different algorithms for the distributionally robust optimization problem (26) on SensIT Vehicle (combined), dna, gisette, protein, letter, mnist, madelon, pendigits datasets. PDG stands for primal-dual gradient method, CB-Min-Max-Simplex stands for Algorithm 5.

**Distributionally Robust Optimization** We consider the following distributionally robust optimization problem:

$$\min_{\|\mathbf{w}\| \leq R} \max_{\mathbf{p} \in \Delta_n} \sum_{i=1}^n p_i \ell_i(\mathbf{w}) + \frac{\lambda}{2} \left\| \mathbf{p} - \frac{\mathbf{1}_n}{n} \right\|^2 + \frac{\rho}{2} \|\mathbf{w}\|^2, \quad (26)$$

where  $\ell_i(\mathbf{w}) = \max(0, 1 - y_i \mathbf{w}^\top \mathbf{x}_i)$  with  $(\mathbf{x}_i, y_i) \in \mathbb{R}^d \times \{-1, +1\}$  being the feature-label pair and  $\mathbf{w} \in \mathbb{R}^m$  being the model parameter,  $\mathbf{p}$  is a probability vector,  $n$  is the number of training examples,  $\mathbf{1}_n = [1; 1; \dots; 1]$  with length  $n$ . In our experiment, we set  $R = 10^5$ ,  $\lambda = \rho = 10^{-4}$ . Because of the added regularizer, the function becomes strictly-convex-strictly-concave and hence satisfies our Assumption 1. We compare two algorithms: the primal-dual gradient (Nemirovski et al., 2009) and our Algorithm 5 (CB-Min-Max-Simplex). In this setup, primal-dual gradient method updates  $\mathbf{w}$  by gradient descent and updates  $\mathbf{p}$  by exponential gradient ascent simultaneously. The learning rates are set to be  $\frac{2R}{G_{\mathbf{w}}\sqrt{T}}$  and  $\frac{\log(n)}{G_{\mathbf{p}}\sqrt{T}}$  respectively for primal variable and dual variable, where  $G_{\mathbf{w}}$  is the 2-norm of gradient in terms of  $\mathbf{w}$ ,  $G_{\mathbf{p}}$  is the infinity-norm of the gradient in terms of  $\mathbf{p}$ , and  $T$  is the number of iterations. Both algorithms start from  $\mathbf{w}_0 = 0$ ,  $\mathbf{p}_0 = [1/n; 1/n, \dots, 1/n]$  and run  $T = 1000$  iterations (each iteration amounts to one pass of the training set). We test our algorithms on four benchmark datasets from libsvm website<sup>3</sup> (SensIT Vehicle (combined), dna, gisette, protein, letter, mnist, madelon, pendigits). We report the training loss and test loss, as shown in the Figure 2. It can be observed that our proposed parameter-free algorithm (i.e., Algorithm 5) significantly outperforms the standard primal-dual gradient method.

## Appendix F. Examples Satisfying Strict-Convexity-Strict-Concavity, Growth Condition and Hölder Continuous Solution Mapping

**Example 1:** We show that problem (9)  $F(x, y) = \frac{\rho}{4}x^4 + |x| + \rho xy - |y| - \frac{\rho}{4}y^4$ , where  $|x| \leq R_x$  and  $|y| \leq R_y$ , satisfies the Assumptions 1, 2 and 3.

**Proof**

3. <https://www.csie.ntu.edu.tw/~cjlin/libsvmtools/datasets/>



We first show that the function is strictly-convex-strictly-concave. The function  $F(x, y)$  is strictly-convex in  $x$ , since  $\frac{\rho}{4}x^4$  is strictly convex in terms of  $x$  and  $|x| + \rho xy$  is convex in terms of  $x$ , and the summation of a strictly convex function and a convex function is strictly convex. The strict-concavity in terms of  $y$  can be proved using a parallel argument.

We then show that Assumption 2 holds with  $\theta_1 = 1$ . We know that the optimal solution of the problem (9) is  $(x_*, y_*) = (0, 0)$ , so we have that

$$\begin{cases} |x - x_*| = |x| \leq \frac{\rho}{4}x^4 + |x| = F(x, y_*) - F(x_*, y_*), \\ |y - y_*| = |y| \leq \frac{\rho}{4}y^4 + |y| = F(x_*, y_*) - F(x_*, y). \end{cases}$$

and hence we have shown that  $\theta_1 = 1$ .

Finally, we show that Assumption 3 holds with  $\theta_2 = 1/3$ . One can show that

$$\arg \min_x F(x, y) = \begin{cases} 0, & \text{if } |y| \leq \frac{1}{\rho} \\ \max \left( (-y + \frac{1}{\rho})^{1/3}, -R_x \right), & \text{if } y > \frac{1}{\rho} \text{ and } y \leq R_y \\ \min \left( (-y - \frac{1}{\rho})^{1/3}, R_x \right), & \text{if } y < -\frac{1}{\rho} \text{ and } y \geq -R_y. \end{cases}$$

We have  $|\arg \min_x F(x, y_1) - \arg \min_x F(x, y_2)| \leq |y_1 - y_2|^{1/3}$ . The parallel argument of the Hölder continuity can also be made in terms of  $\arg \min_y F(x, y)$ . ■

**Example 2** The problem  $F(x, y) = \frac{\rho}{2n}x^{2n} + |x| + \rho xy - |y| - \frac{\rho}{2n}y^{2n}$ , where  $n = 1, 2, \dots$ , satisfies Assumptions 1, 2 and 3.

**Proof** We can use a similar argument as in Example 1 to show that the problem is strictly-convex-strictly-concave, Assumption 2 holds with  $\theta_1 = 1$  and Assumption 3 holds with  $\theta_2 = \frac{1}{2n-1}$ . ■