

Proximal Point Method for Online Saddle Point Problem

Qing-xin Meng^[0000–0003–4014–7405] and Jian-wei Liu^{*}

Department of Automation, College of Artificial Intelligence
China University of Petroleum, Beijing, China
qingxin6174@gmail.com, liujw@cup.edu.cn

Abstract. This paper focuses on the online saddle point problem, which involves a sequence of two-player time-varying convex-concave games. Considering the nonstationarity of the environment, we adopt the duality gap and the dynamic Nash equilibrium regret as performance metrics for algorithm design. We present three variants of the proximal point method: the Online Proximal Point Method (OPPM), the Optimistic OPPM (OptOPPM), and the OptOPPM with multiple predictors. Each algorithm guarantees upper bounds for both the duality gap and dynamic Nash equilibrium regret, achieving near-optimality when measured against the duality gap. Specifically, in certain benign environments, such as sequences of stationary payoff functions, these algorithms maintain a nearly constant metric bound. Experimental results further validate the effectiveness of these algorithms. Lastly, this paper discusses potential reliability concerns associated with using dynamic Nash equilibrium regret as a performance metric. The technical appendix and code can be found at <https://github.com/qingxin6174/PPM-for-OSP>.

Keywords: Online Saddle Point Problem · Proximal Point Method · Duality Gap · Multiple Predictors.

1 Introduction

The Online Saddle Point (OSP) problem, initially introduced by [8], involves a sequence of two-player time-varying convex-concave games. In round t , Players 1 and 2 *jointly* select a strategy pair $(x_t, y_t) \in X \times Y$. Here, Player 1 minimizes his payoff, while Player 2 maximizes his payoff. Both players make decisions without prior knowledge of the current or future payoff functions. Upon finalizing the strategy pair, the environment reveals a continuous payoff function $f_t: X \times Y \rightarrow \mathbb{R}$, which satisfies the following conditions: $\forall y \in Y$, $f_t(\cdot, y)$ is convex on X , and $\forall x \in X$, $f_t(x, \cdot)$ is concave on Y . No additional assumptions are imposed on the environment, thereby allowing potential regularity or even adversarial behavior. The goal is to provide players with decision-making algorithms that approximate Nash equilibria, ensuring that the players' decisions in most rounds are close to saddle points.

^{*} Corresponding author

Table 1: Summary of our results. In this table, \tilde{O} hides poly-logarithmic factors, V_T represents the cumulative difference of the convex-concave payoff function series, V'_T denotes the cumulative error of one single predictor, and V_T^k indicates the cumulative error of the k -th predictor.

Algorithm	Upper Bound of D-Gap _T and NE-Reg _T
OPPM	$\tilde{O}(\min\{V_T, \sqrt{(1+C_T)T}\})$
OptOPPM	$\tilde{O}(\min\{V'_T, \sqrt{(1+C_T)T}\})$
OptOPPM with multiple predictors	$\tilde{O}(\min\{V_T^1, \dots, V_T^d, \sqrt{(1+C_T)T}\})$

Observe that the OSP problem extends the application of Online Convex Optimization (OCO, [33]) to include interactions among two players and the environment. Consequently, it becomes straightforward to identify online scenarios that are well-suited for OSP. Such scenarios include dynamic routing [4,15] and online advertising auctions [21], which fall under the broader category of budget-reward trade-offs.

Given the nonstationarity of the environment, there are two optional metrics available for evaluating the performance of an OSP algorithm:

- 1) The duality gap (D-Gap, [32]), which is given by

$$\text{D-Gap}_T := \sum_{t=1}^T \max_{y \in Y} f_t(x_t, y) - \sum_{t=1}^T \min_{x \in X} f_t(x, y_t). \quad (1)$$

- 2) The dynamic Nash equilibrium regret (NE-Reg, [32]), which is defined as

$$\text{NE-Reg}_T := \left| \sum_{t=1}^T f_t(x_t, y_t) - \sum_{t=1}^T \max_{y \in Y} \min_{x \in X} f_t(x, y) \right|. \quad (2)$$

[32] first presented an algorithm that, under bilinear payoff functions, simultaneously guarantees upper bounds for three performance metrics: the duality gap, dynamic Nash equilibrium regret, and static individual regret. In their commendable work, the authors have directed their focus towards the algorithm's adaptability to a spectrum of metrics, attenuating the pursuit of algorithmic optimality. Notably, their method yields a duality gap upper bound of $\tilde{O}(\sqrt{T})$ for sequences of stationary payoff functions, whereas we advocate for a sharper bound of $O(1)$ that we believe is optimal in such scenarios. Our assertion is inspired by [5], who demonstrated that proximal methods can incur $O(1)$ -level regret in online convex optimization when dealing with stationary loss function sequences.

Contributions In this paper, we propose three variants of the proximal point method for solving the OSP problem: the Online Proximal Point Method (OPPM),

the Optimistic OPPM (OptOPPM), which allows for an arbitrary predictor, and OptOPPM with multiple predictors, enhancing the algorithm to handle multiple predictors. Results are shown in Table 1. All algorithms are near-optimal, as they achieve a worst-case duality gap upper bound of $\tilde{O}(\sqrt{(1 + C_T)T})$. In particular, under favorable scenarios such as stationary payoff function sequences, these algorithms attain a sharp bound of $\tilde{O}(1)$. Notably, the OptOPPM with multiple predictors can autonomously select the most effective predictor from a set of d options. Even when all predictors underperform, it preserves the worst-case bound, significantly enhancing the algorithm's practical utility.

It is imperative to highlight that in time-varying games, relying on NE-Reg for performance evaluation of algorithms may lead to concerns over its reliability. A sublinear D-Gap suggests an approximation to a coarse correlated equilibrium, while a sublinear NE-Reg does not necessarily imply such a correlation. Consider a scenario where, in half of the iterations, the actual payoffs significantly exceed those at the Nash equilibrium, while in the other half, they are substantially lower. This results in a small NE-Reg, but it does not guarantee that the strategies in most rounds are close to the Nash equilibrium. See Example 1 for a more precise discussion.

Related Work The OSP problem is a time-varying version of the minimax problem. The first minimax theorem was proven by [30]. Subsequent to the seminal work of [14], which unveiled the connections between the minimax problem and online learning, there has been a burgeoning interest in the development of no-regret algorithms tailored for static environments [2,11,12,17,24,29]. In recent years, the research focus has expanded to encompass the OSP problem and its various variants [3,8,7,13,26,32].

[8] were the pioneers in explicitly addressing the OSP problem, introducing the saddle-point regret as $|\sum_{t=1}^T f_t(x_t, y_t) - \min_{x \in X} \max_{y \in Y} \sum_{t=1}^T f_t(x, y)|$. In a subsequent work, [7] redefined saddle-point regret as Nash equilibrium regret and developed an FTRL-like algorithm capable of achieving a Nash equilibrium regret bound of $\tilde{O}(\sqrt{T})$. Moreover, they proved that it is impossible for any algorithm to simultaneously attain sublinear Nash equilibrium regret and sublinear individual regret for both players. Building on these findings, [32] refined the notion of dynamic Nash equilibrium regret by reassessing the Nash equilibrium regret, moving the minimax operation inside the summation, and delineating it as Equation (2). They also proposed an algorithm predicated on the meta-expert framework, which ensures upper bounds for three performance metrics: the duality gap, dynamic Nash equilibrium regret, and static individual regret, effectively covering beliefs from stationary to highly nonstationary.

In contrast to the approach by [32], which targets a broad spectrum of non-stationarity levels, this paper takes beliefs that the environment exhibits non-stationarity and accordingly designs algorithms to achieve near-optimal performance. Moreover, this study underscores the potential issues with the reliability of dynamic Nash equilibrium regret as a metric for evaluating algorithmic performance in time-varying games.

2 Preliminaries

Throughout this paper, we define $(\cdot)_+ := \max(\cdot, 0)$, and use the abbreviated notation $1:n$ to represent $1, 2, \dots, n$. For asymptotic upper bounds, we employ big O notation, while \tilde{O} is utilized to hide poly-logarithmic factors. Denote by $\langle \cdot, \cdot \rangle: \mathcal{X}^* \times \mathcal{X} \rightarrow \mathbb{R}$ the canonical dual pairing, where \mathcal{X}^* represents the dual space of \mathcal{X} .

The Fenchel coupling [22,23] induced by a proper function φ is defined as $B_\varphi(x, z) := \varphi(x) + \varphi^*(z) - \langle z, x \rangle$, $\forall (x, z) \in \mathcal{X} \times \mathcal{X}^*$, where φ^* represents the convex conjugate of φ given by $\varphi^*(z) := \sup_{x \in \mathcal{X}} \langle z, x \rangle - \varphi(x)$. Fenchel coupling is the general form of Bregman divergence. By Fenchel-Young inequality, $B_\varphi(x, z) \geq 0$, and equality holds iff z is the subgradient of φ at x . We use the superscripted notation x^φ to abbreviate one subgradient of φ at x . Directly applying the definition of Fenchel coupling yields $B_\varphi(x, y^\varphi) + B_\varphi(y, z) - B_\varphi(x, z) = \langle z - y^\varphi, x - y \rangle$. A similar version can be found in [10].

φ is μ -strongly convex if $\forall x \in \mathcal{X}, \forall y^\varphi \in \partial\varphi(y): B_\varphi(x, y^\varphi) \geq \frac{\mu}{2} \|x - y\|^2$. B_φ is L -Lipschitz with respect to its first variable, if $\exists L$, such that $\forall y \in \mathcal{X}, \forall y^\varphi \in \partial\varphi(y)$, and $\forall x_1, x_2 \in \mathcal{X}: |B_\varphi(x_1, y^\varphi) - B_\varphi(x_2, y^\varphi)| \leq L \|x_1 - x_2\|$.

3 Main Results

This section begins by outlining the assumptions related to the OSP problem. It then evaluates performance metrics and establishes the lower bound for the duality gap, a crucial step in analyzing algorithmic optimality. Finally, this section introduces three variants of the proximal point method: the Online Proximal Point Method (OPPM), the Optimistic OPPM (OptOPPM) that incorporates an arbitrary predictor, and a variant of OptOPPM designed to accommodate multiple predictors, thereby enhancing the algorithm's practical utility. Refer to <https://github.com/qingxin6174/PPM-for-OSP> for all theorem proofs.

3.1 Assumptions

Let $(\mathcal{X}, \|\cdot\|_{\mathcal{X}})$ and $(\mathcal{Y}, \|\cdot\|_{\mathcal{Y}})$ be normed vector spaces. We can omit norm subscripts without ambiguity when the norm is determined by the space to which the element belongs. Let $X \in \mathcal{X}$ and $Y \in \mathcal{Y}$. Now we introduce the following two assumptions:

- A1 The feasible sets X and Y are both compact and convex, with the diameter of X denoted as D_X , and the diameter of Y as D_Y .
- A2 All payoff functions satisfy the boundedness of subdifferentiation, i.e., $\exists G_X, G_Y < +\infty, \forall x \in X, \forall y \in Y, \forall t, \|\partial_x f_t(x, y)\| \leq G_X, \|\partial_y(-f_t)(x, y)\| \leq G_Y$.

In accordance with Theorem 3 by [20], A1 ensures the existence of a saddle point. Specifically, there exists a pair $(x_t^*, y_t^*) \in X \times Y$ such that for all $x \in X$ and $y \in Y: f_t(x_t^*, y) \leq f_t(x_t^*, y_t^*) \leq f_t(x, y_t^*)$. One may write $(x_t^*, y_t^*) = \arg \min_{x \in X} \max_{y \in Y} f_t(x, y)$.

3.2 Discussion on Performance Metrics

This subsection first confirms that the individual regret (Reg) can guarantee both the D-Gap and NE-Reg, and then proceeds to demonstrate the intrinsic issues associated with NE-Reg. Prior to this, let $x'_t = \arg \min_{x \in X} f_t(x, y_t)$ and $y'_t = \arg \max_{y \in Y} f_t(x_t, y)$. Individual regrets of Players 1 and 2 are defined as:

$$\text{Reg}_T^1 := \sum_{t=1}^T f_t(x_t, y_t) - \sum_{t=1}^T f_t(x'_t, y_t), \quad \text{Reg}_T^2 := \sum_{t=1}^T f_t(x_t, y'_t) - \sum_{t=1}^T f_t(x_t, y_t).$$

$\text{D-Gap}_T = \text{Reg}_T^1 + \text{Reg}_T^2$ implies the following lemma.

Lemma 1. *If an online algorithm guarantees individual regrets for both players, it also ensures a D-Gap guarantee.*

Considering that two individual regrets provide the same evaluative utility as the D-Gap, we choose to use only the D-Gap as the metric to assess an online algorithm's performance.

NE-Reg_T ≤ D-Gap_T (Proposition 11 of [32]) deduces the following lemma.

Lemma 2. *If an online algorithm offers individual regret guarantees for both players, it also ensures a guarantee of NE-Reg.*

It is noteworthy that using NE-Reg as a performance metric could raise questions about its reliability. The following extreme example illustrates that the environment may intentionally mislead players, resulting in an internal cancellation of the absolute value of NE-Reg.

Example 1. Consider the OSP problem defined on the feasible domain $[-1, 1]^2$. In round t , the environment feeds back the payoff function of $f_t(x, y) = (x - x_t^*)^2 - (y - y_t^*)^2$, where

$$x_t^* = x_t + \frac{(-1)^t + 1}{2}(2\chi_{x_t < 0} - 1), \quad y_t^* = y_t - \frac{(-1)^t - 1}{2}(2\chi_{y_t < 0} - 1),$$

χ_A is the 0/1 indicator function with $\chi_A = 1$ if and only if A is true. Regardless of the algorithm used, the players' strategy pair will never approach the saddle point, but $\text{NE-Reg}_T = \left| \sum_{t=1}^T (-1)^t \right| \leq 1$.

Although Remark 2 in [32] pointed out that performance metrics based on function values inherently possess certain limitations, the weakness of NE-Reg is easier to trigger.

3.3 Duality Gap Lower Bound

Given the reliability concerns of NE-Reg, this subsection only demonstrates the lower bound for the D-Gap.

Theorem 1 (Duality Gap Lower Bound). *Regardless of the strategy pairs adopted by the two players, there always exist a sequence of convex-concave payoff functions $f_{1:T}$ that adhere to Assumptions A1 and A2, ensuring a duality gap of at least $\Omega(\sqrt{(1 + C_T)T})$, where $C_T = \sum_{t=1}^T (\|x'_t - x'_{t-1}\| + \|y'_t - y'_{t-1}\|)$.*

An online algorithm that attains a D-Gap upper bound of $\tilde{O}(\sqrt{(1 + C_T)T})$ is considered near-optimal, signifying that the upper bound matches the lower bound up to poly-logarithmic terms. However, achieving this upper bound is trivial. Specifically, if Players 1 and 2 independently select OCO algorithms, such as Ader [31] or Online Mirror Descent with doubling trick, they can ensure that $\text{Reg}_T^1 \leq \tilde{O}(\sqrt{(1 + \sum_{t=1}^T \|x'_t - x'_{t-1}\|)T})$, $\text{Reg}_T^2 \leq \tilde{O}(\sqrt{(1 + \sum_{t=1}^T \|y'_t - y'_{t-1}\|)T})$, thereby securing $\text{D-Gap}_T \leq \tilde{O}(\sqrt{(1 + C_T)T})$.

In the following three subsections, we will design non-trivial algorithms that not only ensure an $\tilde{O}(\sqrt{(1 + C_T)T})$ D-Gap upper bound but also further reduce the D-Gap in certain benign environments.

3.4 Online Proximal Point Method

The proximal point method, initially introduced in the seminal work by [25], has established itself as a classic first-order method for solving minimax problems. To tailor it for the solution of the OSP problem, this subsection introduces and analyzes the Online Proximal Point Method (OPPM). OPPM can be formalized as

$$(x_{t+1}, y_{t+1}) = \arg \min_{x \in X} \max_{y \in Y} f_t(x, y) + \frac{1}{\eta_t} B_\phi(x, x_t^\phi) - \frac{1}{\gamma_t} B_\psi(y, y_t^\psi),$$

where $\eta_t > 0$ and $\gamma_t > 0$ represent learning rates of Players 1 and 2, respectively, $x_t^\phi \in \partial\phi(x_t)$, $y_t^\psi \in \partial\psi(y_t)$. To facilitate our analysis, we assume that the regularization terms satisfy the following property.

Property 1. The functions ϕ and ψ are both 1-strongly convex, and their Fenchel couplings satisfy Lipschitz continuity for the first variable. Specifically, $\exists L_\phi < +\infty$, $\exists L_\psi < +\infty$, $\forall \alpha, x, x' \in X$, $\forall \beta, y, y' \in Y$:

$$|B_\phi(x, \alpha^\phi) - B_\phi(x', \alpha^\phi)| \leq L_\phi \|x - x'\|, \quad |B_\psi(y, \beta^\psi) - B_\psi(y', \beta^\psi)| \leq L_\psi \|y - y'\|.$$

The following theorem provides the individual regret guarantee for OPPM.

Theorem 2 (Individual Regret for OPPM). *Under Assumptions A1 and A2, let regularizers satisfy Property 1, and let C be a preset upper bound of C_T . If the learning rates of two players follow from $\eta_t = L(2D + C) / (\epsilon + \sum_{\tau=1}^{t-2} \Delta_\tau) = \gamma_t$, where the constant $\epsilon > 0$ prevents initial learning rates from being infinite, $\Delta_t = (\Sigma_t - \max_{\tau \in 1:t-1} \Sigma_\tau)_+$, $\Sigma_t = \max\{\Sigma_t^1, \Sigma_t^2\}$, and*

$$\begin{aligned} \Sigma_T^1 &= \left(\sum_{t=1}^T (f_t(x_t, y_t) - f_t(x_{t+1}, y_{t+1}) + f_t(x'_{t+1}, y_{t+1}) - f_t(x'_t, y_t)) \right)_+, \\ \Sigma_T^2 &= \left(\sum_{t=1}^T (f_t(x_t, y'_t) - f_t(x_{t+1}, y'_{t+1}) + f_t(x_{t+1}, y_{t+1}) - f_t(x_t, y_t)) \right)_+. \end{aligned}$$

Algorithm 1 OPPM with Adaptive Learning Rates

Require: Feasible sets X and Y satisfy A1. The payoff function f_t satisfies A2. Regularizers satisfy Property 1. $\epsilon > 0$ prevents initial learning rates from being infinite

Initialize: C

- 1: **for** $t = 1, 2, \dots$ **do**
- 2: Output (x_t, y_t) , then observe a continuous convex-concave payoff function f_t
- 3: $x'_t = \arg \min_{x \in X} f_t(x, y_t)$, $y'_t = \max_{y \in Y} f_t(x_t, y)$
- 4: **if** $\sum_{\tau=1}^t (\|x'_\tau - x'_{\tau-1}\| + \|y'_\tau - y'_{\tau-1}\|) > C$ **then** $C \leftarrow 2C$ \triangleright Doubling trick
- 5: Update η_t and γ_t according to the setting of Theorem 2
- 6: $(x_{t+1}, y_{t+1}) = \arg \min_{x \in X} \max_{y \in Y} f_t(x, y) + B_\phi(x, \tilde{x}_t^\phi)/\eta_t - B_\psi(y, \tilde{y}_t^\psi)/\gamma_t$
- 7: **end for**

Then OPPM achieves $\text{Reg}_T^1, \text{Reg}_T^2 \leq O(\min\{\sum_{t=1}^T \rho(f_t, f_{t-1}), \sqrt{(1+C)T}\})$, where $\rho(f_t, f_{t-1}) = \max_{x \in X, y \in Y} |f_t(x, y) - f_{t-1}(x, y)|$ measures the distance between f_t and f_{t-1} .

Theorem 2 corresponds to Theorem 5.1 of [6], but involves more complex parameters and a complicated proof (see Appendix B) due to the mutual influence between players from the joint update of OPPM. This joint update is crucial. Independent execution of implicit online mirror descent, as per [6], only ensures the worst-case bound without addressing the temporal variability term.

In Theorem 2, learning rates depend on the preset value C , but this dependency can be avoided by using the doubling trick [27]. The pseudocode for OPPM is in Algorithm 1. Applying Lemmas 1 and 2, we deduce that Algorithm 1 offers assurances on both the D-Gap and NE-Reg.

Theorem 3 (Performance of Algorithm 1). *If two players output strategy pairs according to Algorithm 1, then*

$$\text{D-Gap}_T, \text{NE-Reg}_T \leq O\left(\min\left\{V_T, \sqrt{(1 + \log_2 C_T + C_T)T}\right\}\right),$$

where $V_T = \sum_{t=1}^T \rho(f_t, f_{t-1})$. Specially, in scenarios where the sequence of payoff functions is stationary, then $\text{D-Gap}_T, \text{NE-Reg}_T, \text{Reg}_T^1, \text{Reg}_T^2 = O(1)$.

The worst-case bound delineated in the first statement of Theorem 3 suggests that the OPPM approaches optimality. For sequences of stationary payoff functions, we observe that the D-Gap, NE-Reg, and two individual regrets all enjoy a sharp bound of $O(1)$, which surpasses the $\tilde{O}(\sqrt{T})$ D-Gap and slightly refines the $\tilde{O}(1)$ NE-Reg, both reported by [32], and matches the $O(1)$ static individual regret in [18].

Typically, OPPM lacks a general solution and necessitates a case-by-case computation tailored to the specific nature of the payoff functions. In essence, each iteration is tasked with solving a strongly-convex-strongly-concave saddle point problem. In certain instances, OPPM can be expressed in a closed form, such as when the regularizers ϕ and ψ are both square norms, and f_t is bilinear or quadratic. When a closed form solution is unavailable, numerical techniques can be employed for efficient approximation [1, 9, 19].

3.5 Optimistic OPPM

The performance of the OPPM algorithm is partially contingent upon the stationarity of the sequence of payoff functions. However, taking a highly nonstationary periodic environment as an example, proactive prediction of the upcoming payoff function can significantly improve the performance of online algorithms. This approach, known as ‘optimism’, replaces the stationarity of the payoff function sequence with prediction accuracy. In this subsection, we explore the optimistic variant of OPPM and introduce the following update:

$$\begin{aligned}(x_t, y_t) &= \arg \min_{x \in X} \max_{y \in Y} h_t(x, y) + \frac{1}{\eta_t} B_\phi(x, \tilde{x}_t^\phi) - \frac{1}{\gamma_t} B_\psi(y, \tilde{y}_t^\psi), \\ \tilde{x}_{t+1} &= \arg \min_{x \in X} \eta_t f_t(x, y_t) + B_\phi(x, \tilde{x}_t^\phi), \\ \tilde{y}_{t+1} &= \arg \max_{y \in Y} \gamma_t f_t(x_t, y) - B_\psi(y, \tilde{y}_t^\psi),\end{aligned}$$

where h_t is an arbitrary convex-concave predictor. We refer to the above update as the Optimistic OPPM (OptOPPM).

The following theorem provides the individual regret guarantee for OptOPPM.

Theorem 4 (Individual Regret for OptOPPM). *Under Assumptions A1 and A2, let regularizers satisfy Property 1, let the predictor h_t satisfy Assumption A2, and let C^1 and C^2 be the preset upper bounds for C_T^1 and C_T^2 , respectively. If the learning rates of two players follow from*

$$\eta_t = L_\phi(D_X + C^1)/(\epsilon + \sum_{\tau=1}^{t-1} \delta_\tau^1), \quad \gamma_t = L_\psi(D_Y + C^2)/(\epsilon + \sum_{\tau=1}^{t-1} \delta_\tau^2),$$

where the constant $\epsilon > 0$ prevents initial learning rates from being infinite, and

$$\begin{aligned}0 \leq \delta_t^1 &= f_t(x_t, y_t) - h_t(x_t, y_t) + h_t(\tilde{x}_{t+1}, y_t) - f_t(\tilde{x}_{t+1}, y_t) - B_\phi(\tilde{x}_{t+1}, x_t^\phi)/\eta_t, \\ 0 \leq \delta_t^2 &= f_t(x_t, \tilde{y}_{t+1}) - h_t(x_t, \tilde{y}_{t+1}) + h_t(x_t, y_t) - f_t(x_t, y_t) - B_\psi(\tilde{y}_{t+1}, y_t^\psi)/\gamma_t.\end{aligned}$$

Then OptOPPM enjoys $\text{Reg}_T^i \leq O(\min \{\sum_{t=1}^T \rho(f_t, h_t), \sqrt{(1 + C^i)T}\})$, $i = 1, 2$.

Similar to Algorithm 1, the dependence on C^1 and C^2 for the learning rates in OptOPPM can be obviated by employing the doubling trick. Details on the adjustment of learning rates are delineated in Algorithm 2.

Utilizing Lemmas 1 and 2, it can be inferred that Algorithm 2 provides guarantees regarding both the D-Gap and the NE-Reg.

Theorem 5 (Performance of Algorithm 2). *If two players output strategy pairs according to Algorithm 2, then*

$$\text{D-Gap}_T, \text{NE-Reg}_T \leq O\left(\min \left\{V'_T, \sqrt{(1 + \log_2 C_T + C_T)T}\right\}\right),$$

where $V'_T = \sum_{t=1}^T \rho(f_t, h_t)$.

Algorithm 2 OptOPPM with Adaptive Learning Rates

Require: Feasible sets X and Y satisfy A1. Both the payoff function f_t and the predictor h_t satisfy A2. Regularizers satisfy Property 1. The constant $\epsilon > 0$ prevents initial learning rates from being infinite

Initialize: C^1, C^2

```

1: for  $t = 1, 2, \dots$  do
2:   Receive  $h_t$ , update  $\eta_t$  and  $\gamma_t$  according to the setting of Theorem 4
3:   Update  $(x_t, y_t) = \arg \min_{x \in X} \max_{y \in Y} h_t(x, y) + B_\phi(x, \tilde{x}_t^\phi)/\eta_t - B_\psi(y, \tilde{y}_t^\psi)/\gamma_t$ 
4:   Output  $(x_t, y_t)$ , then observe a continuous convex-concave payoff function  $f_t$ 
5:    $x'_t = \arg \min_{x \in X} f_t(x, y_t)$ ,  $y'_t = \arg \max_{y \in Y} f_t(x_t, y)$ 
6:   if  $\sum_{\tau=1}^t \|x'_\tau - x'_{\tau-1}\| > C^1$  then  $C^1 \leftarrow 2C^1$  ▷ Player 1 doubles his preset
7:   if  $\sum_{\tau=1}^t \|y'_\tau - y'_{\tau-1}\| > C^2$  then  $C^2 \leftarrow 2C^2$  ▷ Player 2 doubles his preset
8:    $\tilde{x}_{t+1} = \arg \min_{x \in X} f_t(x, y_t) + B_\phi(x, \tilde{x}_t^\phi)/\eta_t$ 
9:    $\tilde{y}_{t+1} = \arg \max_{y \in Y} f_t(x_t, y) - B_\psi(y, \tilde{y}_t^\psi)/\gamma_t$ 
10: end for

```

Algorithm 2 is designed to surpass Algorithm 1. By selecting $h_t = f_{t-n}$, where n is a positive integer, OptOPPM ensures a sharp metric bound of $O(1)$ in environments with a periodicity of k , provided that k is a factor of n .

In Algorithm 2, the update rules for \tilde{x}_{t+1} and \tilde{y}_{t+1} can be determined numerically [28], or transformed into a closed-form expression where feasible.

3.6 Optimistic OPPM with Multiple Predictors

The OptOPPM algorithm is nearly optimal and can reduce the D-Gap with an accurate predictor sequence. However, it currently supports only one predictor, h_t . Given the continuously changing and uncertain nature of the environment, multiple prediction models are often under consideration. An online algorithm that can track the best predictor among many and stay optimal even when all predictors perform poorly would be more versatile. This section extends OptOPPM to support multiple predictors.

Let's consider that there are d models, denoted by M^k , $k = 1:d$. In round t , each model M^k provides its predictor h_t^k . Given that OptOPPM is limited to a single predictor, a logical approach is to compute the weighted average of these d predictors, formulated as $h_t = \sum_{k=1}^d \omega_t^k h_t^k$, with the weights ω_t being derived from a specific implementation of "learning from expert advice".

Consider employing the "clipped" variant of the Hedge algorithm to generate the weight coefficients. The clipped Hedge can be formalized as follows:

$$\omega_{t+1} = \arg \min_{\omega \in \Delta_d^\alpha} \theta_t \langle L_t, \omega \rangle + \text{KL}(\omega, \omega_t),$$

where $\Delta_d^\alpha = \{w \in \mathbb{R}_+^d \mid \|w\|_1 = 1, w^i \geq \alpha/d, \forall i \in 1:d\}$, L_t is a linearized loss, $\theta_t > 0$ is the learning rate, and KL represents the Kullback-Leibler divergence.

Detailed parameter settings are delineated in Algorithm 3, supported by the performance guarantees outlined in Theorem 6.

Algorithm 3 Multi-Predictor Support Subroutine for OptOPPM**Require:** All predictors satisfy Assumption A2 and Property 2. $\epsilon > 0$.**Initialize:** T

```

1: for  $t = 1, 2, \dots$  do
2:   Receive a group of predictors  $h_t^1, h_t^2, \dots, h_t^d$ 
3:   Provide predictor  $h_t = \sum_{k=1}^d \omega_t^k h_t^k$  to OptOPPM
4:   Obtain  $f_t, x_t, \tilde{x}_{t+1}, y_t$ , and  $\tilde{y}_{t+1}$  from Algorithm 2
5:   if  $t > T$  then  $T \leftarrow 2T$  ▷ Doubling trick
6:    $L_t = \left[ \max_{k \in 1:d} \left\{ \begin{array}{l} |f_t(x_t, y_t) - h_t^k(x_t, y_t)|, \\ |f_t(\tilde{x}_{t+1}, y_t) - h_t^k(\tilde{x}_{t+1}, y_t)|, \\ |f_t(x_t, \tilde{y}_{t+1}) - h_t^k(x_t, \tilde{y}_{t+1})| \end{array} \right\} \right]_{k \in 1:d}, \theta_t = \frac{\ln T}{\epsilon + \sum_{\tau=1}^{t-1} \sigma_\tau}$ 
7:    $\omega_{t+1} = \arg \min_{\omega \in \Delta_d^{d/T}} \theta_t \langle L_t, \omega \rangle + \text{KL}(\omega, \omega_t), \sigma_t = \langle L_t, \omega_t - \omega_{t+1} \rangle - \text{KL}(\omega_{t+1}, \omega_t) / \theta_t$ 
8: end for

```

Property 2. All prediction errors are bounded, that is, $\exists L_\infty < +\infty, \forall t, \forall k \in 1:d$, we have that $|\rho(f_t, h_t^k)| \leq L_\infty$.

Theorem 6 (Performance of Algorithm 3). *Let there be d predictors satisfying Property 2. If two players output strategy pairs according to Algorithm 2 and integrate the d predictors following Algorithm 3, then*

$$\text{D-Gap}_T, \text{NE-Reg}_T \leq \tilde{O}\left(\min\left\{V_T^1, \dots, V_T^d, \sqrt{(1+C_T)T}\right\}\right),$$

where $V_T^k = \sum_{t=1}^T \rho(f_t, h_t^k)$ represents the cumulative error of the k -th predictor.

Algorithm 3 is designed to extend multi-predictor support for OptOPPM. We refer to the integration of Algorithms 2 and 3 as *OptOPPM with multiple predictors*.

For the clipped Hedge, an efficient solution involves a slight modification to the algorithm depicted in Figure 3 of [16]. See Algorithm 4 in Appendix G for more details.

4 Experiments

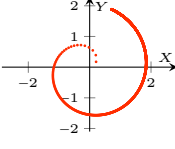
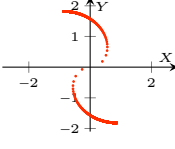
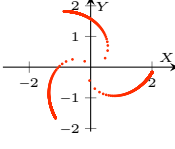
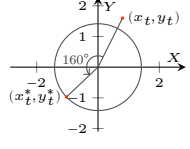
This section validates the effectiveness of our algorithms through experimentation and demonstrates the cancellation phenomenon that occurs within the absolute value of NE-Reg. Section 4.1 describes the experimental setup, and Section 4.2 presents the experimental results.

4.1 Setup

Consider the following synthesis problem: in round t , Players 1 and 2 jointly select a strategy pair $(x_t, y_t) \in X \times Y$, where $X = [-4, 4]$ and $Y = [-4, 4]$, and then the environment feeds back a convex-concave payoff function f_t :

$$f_t(x, y) = \frac{1}{2}(x - x_t^*)^2 - \frac{1}{2}(y - y_t^*)^2 + (x - x_t^*)(y - y_t^*),$$

Table 2: Four Settings for Synthesis Experiments. In this table, $z_1(t) = \ln(1+t)$ is a logarithmic growth function of t , and $z_2(t) = \ln \ln(e+t)$ is a log-logarithmic growth function of t . As time t progresses, the increments of z_1 and z_2 slow down. We represent the saddle point (x_t^*, y_t^*) in the complex form $x_t^* + iy_t^*$, where i is the imaginary unit and satisfying the equation $i^2 = -1$.

Cases	I	II	III	IV
$x_t^* + iy_t^*$	$z_2(t)e^{iz_1(t)}$	$z_2(t)e^{i\pi t + iz_2(t)}$	$z_2(t)e^{i\frac{2\pi}{3}t + iz_2(t)}$	$\sqrt{2}e^{i(\frac{8\pi}{9} + \arg(x_t + iy_t))}$
Trajectories				
Property	$\rho(f_t, f_{t-1}) \rightarrow 0$	$\rho(f_t, f_{t-2}) \rightarrow 0$	$\rho(f_t, f_{t-3}) \rightarrow 0$	Adversarial

where $(x_t^*, y_t^*) \in X \times Y$ is the saddle point of f_t . It is evident that Assumptions A1 and A2 are satisfied. To determine the payoff function f_t , it suffices to fix the saddle point (x_t^*, y_t^*) . We set up four cases, which are listed in Table 2. Case I is characterized by an asymptotic stability, where the saddle points of the payoff functions exhibit a decelerating trend over time. Cases II and III display traits of periodic oscillation; the saddle points under Case II vacillate between two branches, whereas those under Case III alternate among three branches. Case IV is indicative of an adversarial setting. As shown in its figure, upon the players' selection of the strategy pair (x_t, y_t) , the environment engenders the saddle point (x_t^*, y_t^*) by initially rotating the strategy pair by $8\pi/9$, followed by its projection onto a circle with a radius of $\sqrt{2}$. This setting results in the absence of any algorithm that approximates saddle points.

Next, we proceed to instantiate three algorithms: OPPM, OptOPPM, and OptOPPM with multiple predictors. Let $\phi(x) = x^2/2$ and $\psi(y) = y^2/2$. Consequently, both B_ϕ and B_ψ are bounded and exhibit Lipschitz continuity with respect to their first variables. All algorithms are initialized with random initial values, as previous analysis has shown that all metric bounds are invariant to initial conditions. To circumvent an infinitely large initial learning rate, we set $\epsilon = 0.1$. For the OptOPPM algorithm, we employ the predictor $h_t = f_{t-4}$, enabling it to attain a sharp metric bound of $O(1)$ in environments that are either stationary or exhibit periodicity with cycles of 2 or 4. For the OptOPPM with multiple predictors, we configure three predictors: $h_t^1 = f_{t-4}$, $h_t^2 = f_{t-5}$, and $h_t^3 = f_{t-6}$, allowing it to enjoy a sharp metric bound of $\tilde{O}(1)$ in environments that are stationary or have periodicity with cycles of 2, 3, 4, 5, or 6.

Finally, we adopt the algorithm delineated by [32] as the foundational benchmark for comparison.

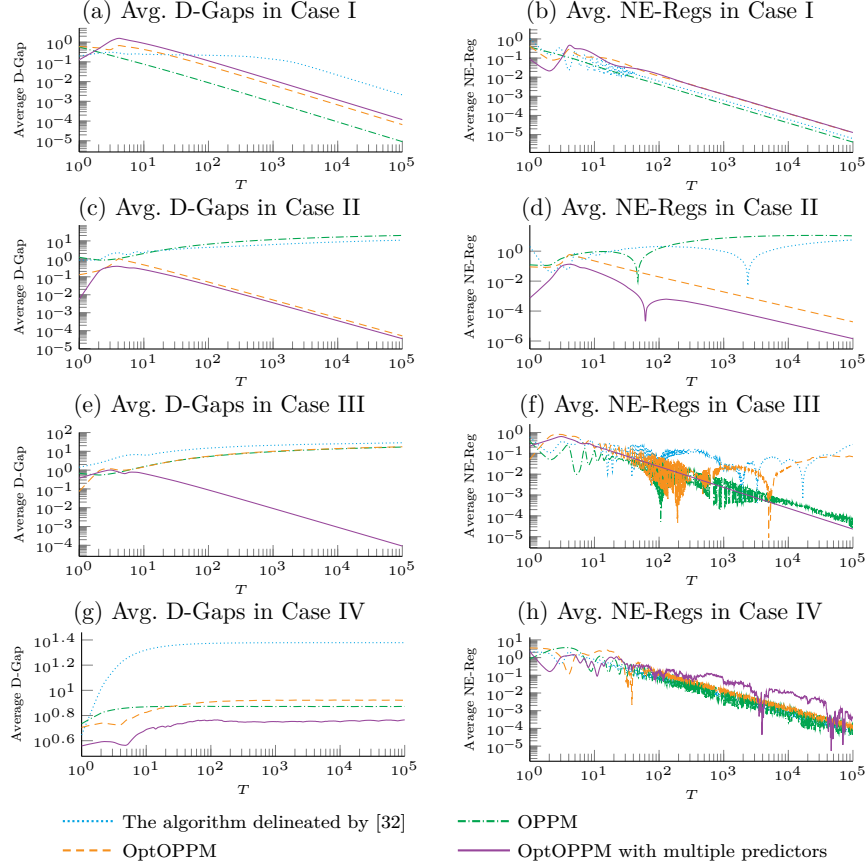


Fig. 1: Average D-Gaps and Average NE-Regs of Algorithms

4.2 Results

We run the repeated game for 10^5 rounds and record the time averaged duality gap $\frac{1}{T}\text{D-Gap}_T$ and the time averaged dynamic Nash equilibrium regret $\frac{1}{T}\text{NE-Reg}_T$. The trajectory of the average duality gap approaches zero, indicating that the players' decisions in most rounds are close to saddle points.

All experimental outcomes align with theoretical expectations. In the benign environment of Case I, the average D-Gap and NE-Reg trajectories of all algorithms converge (see Figure 1a and Figure 1b). In the periodically oscillatory environment of Case II, the average D-Gap and NE-Reg trajectories of OptOPPM and OptOPPM with multiple predictors gradually converge (see Figure 1c and Figure 1d). In the oscillatory scenario of Case III, the average D-Gap and NE-Reg trajectories of OptOPPM with multiple predictors converge (see Figure 1e and Figure 1f). Notably, the average NE-Reg trajectory of OPPM converges amidst pronounced fluctuations, indicating internal cancellation within the ab-

solute value of NE-Reg (see Figure 1f). In the adversarial setting of Case IV, no algorithm approximates the saddle points, resulting in non-converging average D-Gap trajectories (see Figure 1g). However, the average NE-Reg trajectories exhibit intense fluctuations, indicating internal cancellation within the absolute value of NE-Reg (see Figure 1h). The results show that OPPM and its variants outperform the algorithm by [32] in D-Gap performance. Notably, OptOPPM with multiple predictors consistently approximates saddle points in Cases I, II, and III.

Given the insights garnered from both theoretical analysis and empirical results, it seems prudent to reconsider the reliance on NE-Reg as a metric.

5 Conclusion

This study addresses the online saddle point problem by introducing three adaptations of the proximal point method: OPPM, OptOPPM, and OptOPPM with multiple predictors. These methods are crafted to secure upper bounds on D-Gap and NE-Reg, ensuring near-optimal performance in relation to D-Gap. In favorable conditions, such as stationary payoff functions, they maintain near-constant metric bounds. The study also questions the reliability of NE-Reg as a metric and validates the algorithms’ effectiveness through experiments.

Disclosure of Interests. The authors have no competing interests to declare that are relevant to the content of this article.

References

1. Abernethy, J., Lai, K.A., Levy, K.Y., Wang, J.K.: Faster rates for convex-concave games. In: Bubeck, S., Perchet, V., Rigollet, P. (eds.) *Proceedings of the 31st Conference On Learning Theory. Proceedings of Machine Learning Research*, vol. 75, pp. 1595–1625. PMLR (06–09 Jul 2018), <https://proceedings.mlr.press/v75/abernethy18a.html>
2. Anagnostides, I., Daskalakis, C., Farina, G., Fishelson, M., Golowich, N., Sandholm, T.: Near-optimal no-regret learning for correlated equilibria in multi-player general-sum games. In: *Proceedings of the 54th Annual ACM SIGACT Symposium on Theory of Computing*. pp. 736–749. STOC 2022, Association for Computing Machinery, New York, NY, USA (2022). <https://doi.org/10.1145/3519935.3520031>
3. Anagnostides, I., Panageas, I., Farina, G., Sandholm, T.: On the convergence of no-regret learning dynamics in time-varying games. In: *Thirty-seventh Conference on Neural Information Processing Systems* (2023)
4. Awerbuch, B., Kleinberg, R.: Online linear optimization and adaptive routing. *Journal of Computer and System Sciences* **74**(1), 97–114 (2008). <https://doi.org/10.1016/j.jcss.2007.04.016>, <https://www.sciencedirect.com/science/article/pii/S0022000007000621>, *learning Theory* 2004
5. Campolongo, N., Orabona, F.: Temporal variability in implicit online learning. In: Larochelle, H., Ranzato, M., Hadsell, R., Balcan, M., Lin, H. (eds.) *Advances in Neural Information Processing Systems*. vol. 33, pp. 12377–12387. Curran Associates, Inc. (2020), <https://proceedings.neurips.cc/paper/2020/hash/9239be5f9dc4058ec647f14fd04b1290-Abstract.html>

6. Campolongo, N., Orabona, F.: A Closer Look at Temporal Variability in Dynamic Online Learning. arXiv e-prints arXiv:2102.07666 (Feb 2021). <https://doi.org/10.48550/arXiv.2102.07666>
7. Cardoso, A.R., Abernethy, J., Wang, H., Xu, H.: Competing against nash equilibria in adversarially changing zero-sum games. In: Chaudhuri, K., Salakhutdinov, R. (eds.) Proceedings of the 36th International Conference on Machine Learning. Proceedings of Machine Learning Research, vol. 97, pp. 921–930. PMLR (09–15 Jun 2019), <https://proceedings.mlr.press/v97/cardoso19a.html>
8. Cardoso, A.R., Wang, H., Xu, H.: The Online Saddle Point Problem and Online Convex Optimization with Knapsacks. arXiv e-prints (Jun 2018). <https://doi.org/10.48550/arXiv.1806.08301>
9. Carmon, Y., Jin, Y., Sidford, A., Tian, K.: Coordinate methods for matrix games. In: 2020 IEEE 61st Annual Symposium on Foundations of Computer Science (FOCS). pp. 283–293. IEEE Computer Society, Los Alamitos, CA, USA (nov 2020). <https://doi.org/10.1109/FOCS46700.2020.00035>
10. Chen, G., Teboulle, M.: Convergence analysis of a proximal-like minimization algorithm using bregman functions. *SIAM Journal on Optimization* **3**(3), 538–543 (1993). <https://doi.org/10.1137/0803026>
11. Daskalakis, C., Deckelbaum, A., Kim, A.: Near-optimal no-regret algorithms for zero-sum games. *Games and Economic Behavior* **92**, 327–348 (2015). <https://doi.org/10.1016/j.geb.2014.01.003>
12. Daskalakis, C.C., Fishelson, M., Golowich, N.: Near-optimal no-regret learning in general games. In: Beygelzimer, A., Dauphin, Y., Liang, P., Vaughan, J.W. (eds.) Advances in Neural Information Processing Systems (2021)
13. Fiez, T., Sim, R., SKOULAKIS, E.P., Piliouras, G., Ratliff, L.J.: Online learning in periodic zero-sum games. In: Beygelzimer, A., Dauphin, Y., Liang, P., Vaughan, J.W. (eds.) Advances in Neural Information Processing Systems (2021)
14. Freund, Y., Schapire, R.E.: Adaptive game playing using multiplicative weights. *Games and Economic Behavior* **29**(1), 79–103 (1999). <https://doi.org/10.1006/game.1999.0738>
15. Guo, Z., Zhang, Y., Lv, J., Liu, Y., Liu, Y.: An online learning collaborative method for traffic forecasting and routing optimization. *IEEE Transactions on Intelligent Transportation Systems* **22**(10), 6634–6645 (2021). <https://doi.org/10.1109/TITS.2020.2986158>
16. Herbster, M., Warmuth, M.K.: Tracking the best linear predictor. *Journal of Machine Learning Research* **1**, 281–309 (2001)
17. Ho-Nguyen, N., Kılınç-Karzan, F.: Exploiting problem structure in optimization under uncertainty via online convex optimization. *Mathematical Programming* **177**(1), 113–147 (Sep 2019). <https://doi.org/10.1007/s10107-018-1262-8>
18. Hsieh, Y.G., Antonakopoulos, K., Mertikopoulos, P.: Adaptive learning in continuous games: Optimal regret bounds and convergence to nash equilibrium. In: Belkin, M., Kpotufe, S. (eds.) Proceedings of Thirty Fourth Conference on Learning Theory. Proceedings of Machine Learning Research, vol. 134, pp. 2388–2422. PMLR (15–19 Aug 2021), <https://proceedings.mlr.press/v134/hsieh21a.html>
19. Jin, Y., Sidford, A., Tian, K.: Sharper rates for separable minimax and finite sum optimization via primal-dual extragradient methods. In: Loh, P.L., Raginsky, M. (eds.) Proceedings of Thirty Fifth Conference on Learning Theory. Proceedings of Machine Learning Research, vol. 178, pp. 4362–4415. PMLR (02–05 Jul 2022), <https://proceedings.mlr.press/v178/jin22b.html>
20. Kakutani, S.: A generalization of brouwer’s fixed point theorem. *Duke Mathematical Journal* **8**(3), 457–459 (1941). <https://doi.org/10.1215/S0012-7094-41-00838-4>

21. Lykouris, T., Syrgkanis, V., Tardos, E.: Learning and efficiency in games with dynamic population. In: Proceedings of the Twenty-Seventh Annual ACM-SIAM Symposium on Discrete Algorithms. pp. 120–129. SODA '16, Society for Industrial and Applied Mathematics, USA (2016). <https://doi.org/10.5555/2884435.2884444>
22. Mertikopoulos, P., Sandholm, W.H.: Learning in games via reinforcement and regularization. *Mathematics of Operations Research* **41**(4), 1297–1324 (2016). <https://doi.org/10.1287/moor.2016.0778>
23. Mertikopoulos, P., Zhou, Z.: Learning in games with continuous action sets and unknown payoff functions. arXiv e-prints arXiv:1608.07310 (Aug 2016). <https://doi.org/10.48550/arXiv.1608.07310>
24. Rakhlin, A., Sridharan, K.: Optimization, learning, and games with predictable sequences. In: Proceedings of the 26th International Conference on Neural Information Processing Systems - Volume 2. pp. 3066–3074. NIPS'13, Curran Associates Inc., Red Hook, NY, USA (2013). <https://doi.org/10.5555/2999792.2999954>
25. Rockafellar, R.T.: Monotone operators and the proximal point algorithm. *SIAM Journal on Control and Optimization* **14**(5), 877–898 (1976). <https://doi.org/10.1137/0314056>
26. Roy, A., Chen, Y., Balasubramanian, K., Mohapatra, P.: Online and Bandit Algorithms for Nonstationary Stochastic Saddle-Point Optimization. arXiv e-prints (Dec 2019). <https://doi.org/10.48550/arXiv.1912.01698>
27. Schapire, R., Cesa-Bianchi, N., Auer, P., Freund, Y.: Gambling in a rigged casino: The adversarial multi-armed bandit problem. In: 2013 IEEE 54th Annual Symposium on Foundations of Computer Science. p. 322. IEEE Computer Society, Los Alamitos, CA, USA (October 1995). <https://doi.org/10.1109/SFCS.1995.492488>
28. Song, C., Liu, J., Liu, H., Jiang, Y., Zhang, T.: Fully Implicit Online Learning. arXiv e-prints (Sep 2018). <https://doi.org/10.48550/arXiv.1809.09350>
29. Syrgkanis, V., Agarwal, A., Luo, H., Schapire, R.E.: Fast convergence of regularized learning in games. In: Proceedings of the 28th International Conference on Neural Information Processing Systems - Volume 2. pp. 2989–2997. NIPS'15, MIT Press, Cambridge, MA, USA (2015). <https://doi.org/10.5555/2969442.2969573>
30. von Neumann, J.: Zur theorie der gesellschaftsspiele. *Mathematische Annalen* **100**(1), 295–320 (Dec 1928). <https://doi.org/10.1007/BF01448847>
31. Zhang, L., Lu, S., Zhou, Z.H.: Adaptive online learning in dynamic environments. In: Bengio, S., Wallach, H., Larochelle, H., Grauman, K., Cesa-Bianchi, N., Garnett, R. (eds.) *Advances in Neural Information Processing Systems*. vol. 31, pp. 1323–1333. Curran Associates, Inc. (2018), <https://proceedings.neurips.cc/paper/2018/file/10a5ab2db37feedfdeaab192ead4ac0e-Paper.pdf>
32. Zhang, M., Zhao, P., Luo, H., Zhou, Z.H.: No-regret learning in time-varying zero-sum games. In: Chaudhuri, K., Jegelka, S., Song, L., Szepesvari, C., Niu, G., Sabato, S. (eds.) *Proceedings of the 39th International Conference on Machine Learning. Proceedings of Machine Learning Research*, vol. 162, pp. 26772–26808. PMLR (17–23 Jul 2022), <https://proceedings.mlr.press/v162/zhang22an.html>
33. Zinkevich, M.: Online convex programming and generalized infinitesimal gradient ascent. In: Fawcett, T., Mishra, N. (eds.) *Proceedings of the Twentieth International Conference on Machine Learning*. p. 928–935. ICML'03, AAAI Press (2003)

Technical Appendix

A Proof of Theorem 1

Proof (Proof of Theorem 1). Let \mathcal{F} be the set of all convex-concave functions satisfying Assumption A2, and let $\mathcal{L}_X^{G_X} = \{\ell \mid \ell(x) + 0y \in \mathcal{F}\}$, let $\mathcal{L}_Y^{G_Y} = \{\ell \mid 0x - \ell(y) \in \mathcal{F}\}$. Formally, we need to prove that

$$\exists f_{1:T} \in \mathcal{F}, \quad \text{such that} \quad \sum_{t=1}^T \max_{y \in Y} f_t(x_t, y) - \sum_{t=1}^T \min_{x \in X} f_t(x, y_t) \geq \Omega(\sqrt{(1 + C_T)T}).$$

Indeed, we may choose $f_t(x, y) = \alpha_t(x) - \beta_t(y) \in \mathcal{F}$, where α_t and β_t are both convex functions. $\forall C \in [0, T(D_X + D_Y)]$ and $\forall P \in [0, C]$, $\exists \alpha_{1:T}$ and $\exists \beta_{1:T}$, such that

$$\begin{aligned} & \sum_{t=1}^T \max_{y \in Y} f_t(x_t, y) - \sum_{t=1}^T \min_{x \in X} f_t(x, y_t) \\ &= \max_{\forall u_{1:T} \in X, \forall v_{1:T} \in Y} \sum_{t=1}^T (f_t(x_t, v_t) - f_t(u_t, y_t)) \\ &\geq \max_{\sum_{t=1}^T \|u_t - u_{t-1}\| \leq P, \sum_{t=1}^T \|v_t - v_{t-1}\| \leq C-P} \sum_{t=1}^T (f_t(x_t, v_t) - f_t(u_t, y_t)) \\ &= \max_{\sum_{t=1}^T \|u_t - u_{t-1}\| \leq P} \sum_{t=1}^T (\alpha_t(x_t) - \alpha_t(u_t)) + \max_{\sum_{t=1}^T \|v_t - v_{t-1}\| \leq C-P} \sum_{t=1}^T (\beta_t(y_t) - \beta_t(v_t)) \\ &\geq \Omega(\sqrt{(1+P)T}) + \Omega(\sqrt{(1+C-P)T}) \\ &= \Omega(\sqrt{(1+C)T}), \end{aligned} \tag{3}$$

where the last “ \leq ” follows from Lemma 3. Note that Equation (3) holds for arbitrary $C \in [0, T(D_X + D_Y)]$. Choosing $C = C_T$ yields the desired result.

Lemma 3 (Theorem 2 of [31]). *In the context of Online Convex Optimization, let the feasible set X be compact and convex, and let \mathcal{L}_X^G be the set of all convex loss functions defined on X with subgradients bounded by G . Regardless of the strategy adopted by the player, there always exists a comparator sequence $u_{1:T}$ satisfying $\sum_{t=1}^T \|u_t - u_{t-1}\| \leq P$, and a sequence of loss functions $\ell_{1:T} \in \mathcal{L}_X^G$, ensuring that the regret is not less than $\Omega(\sqrt{(1+P)T})$.*

B Proof of Theorem 2

Proof (Proof of Theorem 2). The first-order optimality condition of OPPM implies that

$$\begin{aligned} \exists \nabla_x f_t(x_{t+1}, y_{t+1}), \quad \forall x' \in X, \quad \langle \eta_t \nabla_x f_t(x_{t+1}, y_{t+1}) + x_{t+1}^\phi - x_t^\phi, x_{t+1} - x' \rangle &\leq 0, \\ \exists \nabla_y(-f_t)(x_{t+1}, y_{t+1}), \quad \forall y' \in Y, \quad \langle \gamma_t \nabla_y(-f_t)(x_{t+1}, y_{t+1}) + y_{t+1}^\psi - y_t^\psi, y_{t+1} - y' \rangle &\leq 0. \end{aligned}$$

Let's take Player 1 as an example. The identity transformation on the instantaneous individual regret is as follows:

$$\begin{aligned} f_t(x_t, y_t) - f_t(x'_t, y_t) &= f_t(x_t, y_t) - f_t(x_{t+1}, y_{t+1}) + f_t(x'_{t+1}, y_{t+1}) - f_t(x'_t, y_t) \\ &\quad + \underbrace{f_t(x_{t+1}, y_{t+1}) - f_t(x'_{t+1}, y_{t+1})}_{(4a)}. \end{aligned} \tag{4}$$

Applying convexity and first-order optimality condition, we get

$$\begin{aligned} \text{Equation (4a)} &\leq \langle \nabla_x f_t(x_{t+1}, y_{t+1}), x_{t+1} - x'_{t+1} \rangle \leq \frac{1}{\eta_t} \langle x_t^\phi - x_{t+1}^\phi, x_{t+1} - x'_{t+1} \rangle \\ &= \frac{1}{\eta_t} [B_\phi(x'_{t+1}, x_t^\phi) - B_\phi(x'_{t+1}, x_{t+1}^\phi) - B_\phi(x_{t+1}, x_t^\phi)]. \end{aligned}$$

Summing Equation (4) over time yields

$$\begin{aligned} \text{Reg}_T^1 &\leq \underbrace{\sum_{t=1}^T \frac{1}{\eta_t} (B_\phi(x'_{t+1}, x_t^\phi) - B_\phi(x'_{t+1}, x_{t+1}^\phi))}_{(5a)} \\ &\quad + \underbrace{\sum_{t=1}^T (f_t(x_t, y_t) - f_t(x_{t+1}, y_{t+1}) + f_t(x'_{t+1}, y_{t+1}) - f_t(x'_t, y_t))}_{(5b)}. \end{aligned} \tag{5}$$

Since the learning rate η_t does not increase, B_ϕ is upper bounded by $L_\phi D_X$ and is L_ϕ -Lipschitz for the first variable, we have that

$$\begin{aligned} \text{Equation (5a)} &\leq \sum_{t=1}^T \frac{1}{\eta_t} (B_\phi(x'_{t+1}, x_t^\phi) - B_\phi(x'_t, x_t^\phi)) + \frac{B_\phi(x'_1, x_1^\phi)}{\eta_0} + \sum_{t=1}^T \left(\frac{1}{\eta_t} - \frac{1}{\eta_{t-1}} \right) B_\phi(x'_t, x_t^\phi) \\ &\leq \frac{L_\phi D_X}{\eta_T} + \sum_{t=1}^T \frac{L_\phi}{\eta_t} \|x'_{t+1} - x'_t\| \leq \frac{2L_\phi D_X}{\eta_T} + \sum_{t=1}^T \frac{L_\phi}{\eta_t} \|x'_t - x'_{t-1}\|. \end{aligned}$$

Let $\Sigma_T^1 = (\text{Equation (5b)})_+$, then Equation (5) can be rearranged as follows:

$$\text{Reg}_T^1 \leq L_\phi \left(2 \frac{D_X}{\eta_T} + \sum_{t=1}^T \frac{1}{\eta_t} \|x'_t - x'_{t-1}\| \right) + \Sigma_T^1.$$

Likewise,

$$\text{Reg}_T^2 \leq L_\psi \left(2 \frac{D_Y}{\gamma_T} + \sum_{t=1}^T \frac{1}{\gamma_t} \|y'_t - y'_{t-1}\| \right) + \Sigma_T^2,$$

where

$$\Sigma_T^2 = \left(\sum_{t=1}^T (f_t(x_t, y'_t) - f_t(x_{t+1}, y'_{t+1}) + f_t(x_{t+1}, y_{t+1}) - f_t(x_t, y_t)) \right)_+.$$

Let $L = \max\{L_\phi, L_\psi\}$, $D = \max\{D_X, D_Y\}$, $C_T = \sum_{t=1}^T (\|x'_t - x'_{t-1}\| + \|y'_t - y'_{t-1}\|)$, and let learning rates satisfy $\eta_t = \gamma_t$. Consequently, the individual regrets are subject to the following public upper bound:

$$\text{Reg}_T^1, \text{Reg}_T^2 \leq L \frac{2D + C_T}{\eta_T} + \max\{\Sigma_T^1, \Sigma_T^2\}.$$

Let $\Sigma_T = \max\{\Sigma_T^1, \Sigma_T^2\}$, $\Delta_t = (\Sigma_t - \max_{\tau \in 1:t-1} \Sigma_\tau)_+$. We claim that $\sum_{\tau=1}^t \Delta_\tau = \max_{\tau \in 1:t} \Sigma_\tau$. This claim can be proved by induction. It is obvious that $\Delta_1 = \Sigma_1$. Now we

assume the claim holds for $t-1$ and prove it for t :

$$\begin{aligned}\sum_{\tau=1}^t \Delta_\tau &= \Delta_t + \sum_{\tau=1}^{t-1} \Delta_\tau = \left(\Sigma_t - \max_{\tau \in 1:t-1} \Sigma_\tau \right)_+ + \max_{\tau \in 1:t-1} \Sigma_\tau \\ &= \begin{cases} \Sigma_t & \Sigma_t \geq \max_{\tau \in 1:t-1} \Sigma_\tau \\ \max_{\tau \in 1:t-1} \Sigma_\tau & \Sigma_t < \max_{\tau \in 1:t-1} \Sigma_\tau \end{cases} = \max_{\tau \in 1:t} \Sigma_\tau.\end{aligned}$$

Let C be a preset upper bound of C_T , then applying the prescribed learning rate yields

$$\text{Reg}_T^1, \text{Reg}_T^2 \leq \epsilon + 2 \sum_{t=1}^T \Delta_t = \epsilon + 2 \max_{t \in 1:T} \Sigma_t. \quad (6)$$

Next, we estimate the upper bound of the r.h.s. of the above inequality in two ways. The first one:

$$\max_{t \in 1:T} \Sigma_t \leq O(1) + 2 \sum_{t=1}^T \rho(f_t, f_{t-1}). \quad (7)$$

Indeed,

$$\begin{aligned}\Sigma_t^1 &= \left(\sum_{\tau=1}^t \left(f_\tau(x_\tau, y_\tau) - f_\tau(x_{\tau+1}, y_{\tau+1}) + f_\tau(x'_{\tau+1}, y_{\tau+1}) - f_\tau(x'_\tau, y_\tau) \right) \right)_+ \\ &\leq \left| f_0(x_1, y_1) - f_0(x'_1, y_1) - f_t(x_{t+1}, y_{t+1}) + f_t(x'_{t+1}, y_{t+1}) \right| \\ &\quad + \sum_{\tau=1}^t \left| f_\tau(x_\tau, y_\tau) - f_{\tau-1}(x_\tau, y_\tau) \right| + \sum_{\tau=1}^t \left| f_\tau(x'_\tau, y_\tau) - f_{\tau-1}(x'_\tau, y_\tau) \right| \\ &\leq 2(D_X G_X + D_Y G_Y) + 2 \sum_{\tau=1}^t \rho(f_\tau, f_{\tau-1}),\end{aligned}$$

where the last “ \leq ” follows from Lemma 4. Likewise,

$$\Sigma_t^2 \leq 2(D_X G_X + D_Y G_Y) + 2 \sum_{\tau=1}^t \rho(f_\tau, f_{\tau-1}).$$

So Equation (7) holds. The other one:

$$\sum_{t=1}^T \Delta_t \leq O\left(\sqrt{(1+C)T}\right). \quad (8)$$

The derivation is as follows. According to Lemma 5, $\|x_t - x_{t+1}\| \leq \min\{D_X, \eta_t G_X\}$, and $\|y_t - y_{t+1}\| \leq \min\{D_Y, \eta_t G_Y\}$. Let $G = \max\{G_X, G_Y\}$. Thus we have that

$$\begin{aligned}(\Sigma_t^1 - \Sigma_{t-1}^1)_+ &\leq (f_t(x_t, y_t) - f_t(x_{t+1}, y_{t+1}) + f_t(x'_{t+1}, y_{t+1}) - f_t(x'_t, y_t))_+ \\ &\leq |f_t(x_t, y_t) - f_t(x_{t+1}, y_{t+1})| + |f_t(x_{t+1}, y_{t+1}) - f_t(x'_{t+1}, y_{t+1})| \\ &\quad + |f_t(x'_{t+1}, y_{t+1}) - f_t(x'_t, y_t)| + |f_t(x'_t, y_t) - f_t(x_t, y_t)| \\ &\leq G_Y \|y_t - y_{t+1}\| + G_X \|x_t - x_{t+1}\| + G_Y \|y_{t+1} - y_t\| + G_X \|x'_t - x'_{t+1}\| \\ &\leq \min\{D_X G_X + 2D_Y G_Y, \eta_t (G_X^2 + 2G_Y^2)\} + G_X \|x'_t - x'_{t+1}\| \\ &\leq \min\{3DG, 3\eta_t G^2\} + G(\|x'_t - x'_{t+1}\| + \|y'_t - y'_{t+1}\|).\end{aligned}$$

Likewise,

$$(\Sigma_t^2 - \Sigma_{t-1}^2)_+ \leq \min \{3DG, 3\eta_t G^2\} + G(\|x'_t - x'_{t+1}\| + \|y'_t - y'_{t+1}\|).$$

Note that $\Sigma_t - \Sigma_{t-1} = \max \{\Sigma_t^1 - \Sigma_{t-1}^1, \Sigma_t^2 - \Sigma_{t-1}^2\} \leq \max \{\Sigma_t^1 - \Sigma_{t-1}^1, \Sigma_t^2 - \Sigma_{t-1}^2\}$, so we get

$$\begin{aligned} \Delta_t &= \left(\Sigma_t - \max_{\tau \in 1:t-1} \Sigma_\tau \right)_+ \leq (\Sigma_t - \Sigma_{t-1})_+ \\ &\leq \min \{3DG, 3\eta_t G^2\} + G(\|x'_t - x'_{t+1}\| + \|y'_t - y'_{t+1}\|). \end{aligned}$$

Let $\xi_t = (\Delta_t - G(\|x'_t - x'_{t+1}\| + \|y'_t - y'_{t+1}\|))_+$, then $\xi_t \leq \min \{3DG, 3\eta_t G^2\}$, and thus,

$$\begin{aligned} \left(\sum_{t=1}^{T-1} \xi_t \right)^2 &= \sum_{t=1}^{T-1} \xi_t^2 + 2 \sum_{t=1}^{T-1} \xi_t \sum_{\tau=1}^{t-2} \xi_\tau \leq \sum_{t=1}^{T-1} \xi_t^2 + 2 \sum_{t=1}^{T-1} \xi_t \sum_{\tau=1}^{t-2} \Delta_\tau \\ &= \sum_{t=1}^{T-1} \xi_t^2 + 2 \sum_{t=1}^{T-1} \xi_t \left(\frac{L(2D+C)}{\eta_t} - \epsilon \right) \leq \sum_{t=1}^{T-1} 9D^2 G^2 + \sum_{t=1}^{T-1} 6G^2 L(2D+C), \end{aligned}$$

where the last “ \leq ” uses the first and second terms in the minimum of the bound for ξ_t in turn. Now we get

$$\sum_{t=1}^T \Delta_t \leq O(C_T) + \sum_{t=1}^T \xi_t \leq O(C_T) + O(\sqrt{(1+C)T}) \leq O(\sqrt{(1+C)T}).$$

So Equation (8) holds. Substituting Equations (7) and (8) into Equation (6) gives the following conclusion:

$$\text{Reg}_T^1, \text{Reg}_T^2 \leq O \left(\min \left\{ \sum_{t=1}^T \rho(f_t, f_{t-1}), \sqrt{(1+C)T} \right\} \right).$$

Lemma 4. *Assumptions A1 and A2 guarantee the finiteness of the instantaneous duality gap.*

Proof (Proof of Lemma 4). $\forall x, x' \in X, \forall y, y' \in Y$,

$$\begin{aligned} \langle \partial_x f_t(x', y), x - x' \rangle &\leq f_t(x, y) - f_t(x', y) \leq \langle \partial_x f_t(x, y), x - x' \rangle, \\ \langle \partial_y(-f_t)(x, y'), y - y' \rangle &\leq f_t(x, y') - f_t(x, y) \leq \langle \partial_y(-f_t)(x, y), y - y' \rangle, \end{aligned}$$

which implies that $|f_t(x, y) - f_t(x', y)| \leq D_X G_X$, $|f_t(x, y') - f_t(x, y)| \leq D_Y G_Y$, and $|f_t(x, y') - f_t(x', y)| \leq |f_t(x, y') - f_t(x, y)| + |f_t(x, y) - f_t(x', y)| \leq D_X G_X + D_Y G_Y$.

Lemma 5. *Under Assumption A2 and Property 1, OPPM guarantees that $\|x_t - x_{t+1}\| \leq \eta_t G_X$ and $\|y_t - y_{t+1}\| \leq \gamma_t G_Y$.*

Proof (Proof of Lemma 5). Let $F_t(x, y) = f_t(x, y) + B_\phi(x, x_t^\phi)/\eta_t - B_\psi(y, y_t^\psi)/\gamma_t$. Let $\Phi = F_t(\cdot, y_{t+1})$. Note that Φ is η_t^{-1} -strongly convex, so we have that

$$B_\Phi(x_{t+1}, x_t^\Phi) \geq \frac{1}{2\eta_t} \|x_t - x_{t+1}\|^2, \quad B_\Phi(x_t, x_{t+1}^\Phi) \geq \frac{1}{2\eta_t} \|x_t - x_{t+1}\|^2.$$

Choosing $x_{t+1}^\Phi = 0$, and adding the above two inequalities yields

$$\frac{1}{\eta_t} \|x_t - x_{t+1}\|^2 \leq \langle x_t^\Phi - x_{t+1}^\Phi, x_t - x_{t+1} \rangle \leq \|x_t^\Phi\| \|x_t - x_{t+1}\|.$$

Note that $\|x_t^\Phi\| = \|\nabla_x f_t(x_t, y_{t+1})\| \leq G_X$, Thus we have that $\|x_t - x_{t+1}\| \leq \eta_t G_X$. Likewise, $\|y_t - y_{t+1}\| \leq \gamma_t G_Y$.

C Proof of Theorem 3

Proof (Proof of Theorem 3). Consider the Algorithm 1 initially at stage 0. Each time the value C is doubled, the algorithm advances to the next stage. We use C_n to represent the value of C at stage n , and let S_n be the set of indices for all rounds in the n -th stage. Assume that the game pauses after the T -th round is completed, at which point the OPPM is at stage s . Applying Theorem 2 yields

$$\begin{aligned} \text{Reg}_T^i &= \sum_{n=1}^s \text{Reg}_{S_n}^i \leq \sum_{n=1}^s O \left(\min \left\{ \sum_{t \in S_n} \rho(f_t, f_{t-1}), \sqrt{(1 + 2^n C_0) |S_n|} \right\} \right) \\ &\leq O \left(\min \left\{ \sum_{t=1}^T \rho(f_t, f_{t-1}), \sqrt{s + \sum_{n=1}^s 2^n C_0 \sqrt{\sum_{n=1}^s |S_n|}} \right\} \right), \quad \forall i = 1, 2. \end{aligned}$$

Note that $2^{s-1} C_0 = C_{s-1} < C_T \leq C_s = 2^s C_0$, and $T = \sum_{n=1}^s |S_n|$. We have that

$$\text{Reg}_T^i \leq O \left(\min \left\{ \sum_{t=1}^T \rho(f_t, f_{t-1}), \sqrt{(1 + \log_2 C_T + 4C_T) T} \right\} \right), \quad \forall i = 1, 2.$$

Applying Lemmas 1 and 2 yields the conclusion to be proved.

D Proof of Theorem 4

Proof (Proof of Theorem 4). The first-order optimality condition of OptOPPM implies that

$$\begin{aligned} \exists \nabla_x h_t(x_t, y_t), \quad \forall x' \in X, \quad & \langle \eta_t \nabla_x h_t(x_t, y_t) + x_t^\phi - \tilde{x}_t^\phi, x_t - x' \rangle \leq 0, \\ \exists \nabla_y(-h_t)(x_t, y_t), \quad \forall y' \in Y, \quad & \langle \gamma_t \nabla_y(-h_t)(x_t, y_t) + y_t^\psi - \tilde{y}_t^\psi, y_t - y' \rangle \leq 0, \\ \exists \nabla_x f_t(\tilde{x}_{t+1}, y_t), \quad \forall x' \in X, \quad & \langle \eta_t \nabla_x f_t(\tilde{x}_{t+1}, y_t) + \tilde{x}_{t+1}^\phi - \tilde{x}_t^\phi, \tilde{x}_{t+1} - x' \rangle \leq 0, \\ \exists \nabla_y(-f_t)(x_t, \tilde{y}_{t+1}), \quad \forall y' \in Y, \quad & \langle \gamma_t \nabla_y(-f_t)(x_t, \tilde{y}_{t+1}) + \tilde{y}_{t+1}^\psi - \tilde{y}_t^\psi, \tilde{y}_{t+1} - y' \rangle \leq 0. \end{aligned}$$

Let's take Player 1 as an example. We first perform identity transformation on the instantaneous individual regret:

$$\begin{aligned} f_t(x_t, y_t) - f_t(x'_t, y_t) &= \underbrace{f_t(x_t, y_t) - h_t(x_t, y_t) + h_t(\tilde{x}_{t+1}, y_t) - f_t(\tilde{x}_{t+1}, y_t)}_{(9a)} \\ &\quad + \underbrace{h_t(x_t, y_t) - h_t(\tilde{x}_{t+1}, y_t) + f_t(\tilde{x}_{t+1}, y_t) - f_t(x'_t, y_t)}_{(9b)}. \end{aligned} \tag{9}$$

By using convexity and first-order optimality conditions, we get

$$\begin{aligned} \text{Equation (9b)} &\leq \langle \nabla_x h_t(x_t, y_t), x_t - \tilde{x}_{t+1} \rangle + \langle \nabla_x f_t(\tilde{x}_{t+1}, y_t), \tilde{x}_{t+1} - x'_t \rangle \\ &\leq \langle \tilde{x}_t^\phi - x_t^\phi, x_t - \tilde{x}_{t+1} \rangle / \eta_t + \langle \tilde{x}_t^\phi - \tilde{x}_{t+1}^\phi, \tilde{x}_{t+1} - x'_t \rangle / \eta_t \\ &= [B_\phi(\tilde{x}_{t+1}, \tilde{x}_t^\phi) - B_\phi(\tilde{x}_{t+1}, x_t^\phi) - B_\phi(x_t, \tilde{x}_t^\phi)] / \eta_t \\ &\quad + \underbrace{[B_\phi(x'_t, \tilde{x}_t^\phi) - B_\phi(x'_t, \tilde{x}_{t+1}^\phi)] / \eta_t}_{=:\Phi_t} - B_\phi(\tilde{x}_{t+1}, \tilde{x}_t^\phi) / \eta_t. \end{aligned} \tag{10}$$

Let $\delta_t^1 = \text{Equation (9a)} - B_\phi(\tilde{x}_{t+1}, x_t^\phi)/\eta_t$, so we have that $f_t(x_t, y_t) - f_t(x'_t, y_t) \leq \Phi_t + \delta_t^1$. Note that η_t is non-increasing over time, B_ϕ is L_ϕ -Lipschitz w.r.t. the first variable, and $L_\phi D_X$ is the supremum of B_ϕ . Thus,

$$\begin{aligned} \sum_{t=1}^T \Phi_t &\leq \frac{B_\phi(x'_0, \tilde{x}_1^\phi)}{\eta_0} + \sum_{t=1}^T \frac{1}{\eta_t} (B_\phi(x'_t, \tilde{x}_t^\phi) - B_\phi(x'_{t-1}, \tilde{x}_t^\phi)) + \sum_{t=1}^T \left(\frac{1}{\eta_t} - \frac{1}{\eta_{t-1}} \right) B_\phi(x'_{t-1}, \tilde{x}_t^\phi) \\ &\leq \frac{L_\phi D_X}{\eta_T} + \sum_{t=1}^T \frac{L_\phi}{\eta_t} \|x'_t - x'_{t-1}\|, \end{aligned}$$

Now we get

$$\text{Reg}_T^1 \leq \frac{L_\phi D_X}{\eta_T} + \sum_{t=1}^T \frac{L_\phi}{\eta_t} \|x'_t - x'_{t-1}\| + \sum_{t=1}^T \delta_t^1.$$

Likewise,

$$\text{Reg}_T^2 \leq \frac{L_\psi D_Y}{\gamma_T} + \sum_{t=1}^T \frac{L_\psi}{\gamma_t} \|y'_t - y'_{t-1}\| + \sum_{t=1}^T \delta_t^2,$$

where $\delta_t^2 = f_t(x_t, \tilde{y}_{t+1}) - h_t(x_t, \tilde{y}_{t+1}) + h_t(x_t, y_t) - f_t(x_t, y_t) - B_\psi(\tilde{y}_{t+1}, y_t^\psi)/\gamma_t$. Next, we verify $\delta_t^1, \delta_t^2 \geq 0$. To verify $\delta_t^1 \geq 0$, it suffices to combine the following two inequalities:

$$\begin{aligned} f_t(x_t, y_t) + B_\phi(x_t, \tilde{x}_t^\phi)/\eta_t &\geq f_t(\tilde{x}_{t+1}, y_t) + B_\phi(\tilde{x}_{t+1}, \tilde{x}_t^\phi)/\eta_t, \\ -h_t(x_t, y_t) + h_t(\tilde{x}_{t+1}, y_t) &\geq -[B_\phi(\tilde{x}_{t+1}, \tilde{x}_t^\phi) - B_\phi(\tilde{x}_{t+1}, x_t^\phi) - B_\phi(x_t, \tilde{x}_t^\phi)]/\eta_t. \end{aligned}$$

The first inequality takes advantage of the optimality condition, and the second inequality is part of Equation (10). Likewise, $\delta_t^2 \geq 0$. So all learning rates are non-increasing. Let's go back to the focus on Player 1. The prescribed learning rate guarantees that

$$\text{Reg}_T^1 \leq \frac{L_\phi}{\eta_T} (D_X + C_T^1) + \sum_{t=1}^T \delta_t^1 \leq \epsilon + 2 \sum_{t=1}^T \delta_t^1.$$

On the one hand, $\delta_t^1 \leq 2\rho(f_t, h_t)$ causes

$$\text{Reg}_T^1 \leq \epsilon + 4 \sum_{t=1}^T \rho(f_t, h_t). \quad (11)$$

On the other hand, notice that

$$\begin{aligned} \delta_t^1 &\leq \langle \nabla_x f_t(x_t, y_t) - \nabla_x h_t(\tilde{x}_{t+1}, y_t), x_t - \tilde{x}_{t+1} \rangle - B_\phi(\tilde{x}_{t+1}, x_t^\phi)/\eta_t \\ &\leq 2G_X \|x_t - \tilde{x}_{t+1}\| - B_\phi(\tilde{x}_{t+1}, x_t^\phi)/\eta_t \leq \min\{2D_X G_X, 2\eta_t G_X^2\}, \end{aligned} \quad (12)$$

which implies that

$$\begin{aligned} \left(\sum_{t=1}^T \delta_t^1 \right)^2 &= \sum_{t=1}^T (\delta_t^1)^2 + 2 \sum_{t=1}^T \delta_t^1 \sum_{\tau=1}^{t-1} \delta_\tau^1 = \sum_{t=1}^T (\delta_t^1)^2 + 2 \sum_{t=1}^T \delta_t^1 \left(\frac{L_\phi (D_X + C^1)}{\eta_t} - \epsilon \right) \\ &\leq \sum_{t=1}^T 4G_X^2 D_X^2 + \sum_{t=1}^T 4G_X^2 L_\phi (D_X + C^1). \end{aligned}$$

This results in the following regret bound:

$$\text{Reg}_T^1 \leq \epsilon + 4G_X \sqrt{(D_X^2 + L_\phi D_X + L_\phi C^1) T}. \quad (13)$$

Combining Equations (11) and (13) yields

$$\text{Reg}_T^1 \leq \epsilon + 4 \min \left\{ \sum_{t=1}^T \rho(f_t, h_t), G_X \sqrt{(D_X^2 + L_\phi D_X + L_\phi C^1) T} \right\}.$$

Likewise, the individual regret of Player 2 satisfies

$$\text{Reg}_T^2 \leq \epsilon + 4 \min \left\{ \sum_{t=1}^T \rho(f_t, h_t), G_Y \sqrt{(D_Y^2 + L_\psi D_Y + L_\psi C^2) T} \right\}.$$

E Proof of Theorem 5

Refer to the proof of Theorem 3.

F Proof of Theorem 6

Proof (Proof of Theorem 6). Note that in the proof of Theorem 4, the relaxed form inequalities $\delta_t^1, \delta_t^2 \leq 2\rho(f_t, h_t)$ are employed. In fact, L_t induces a tighter upper bound.

$$\begin{aligned} \delta_t^1 &\leq f_t(x_t, y_t) - h_t(x_t, y_t) + h_t(\tilde{x}_{t+1}, y_t) - f_t(\tilde{x}_{t+1}, y_t) \\ &\leq \sum_{k=1}^{\kappa} \omega_t^k \left(|f_t(x_t, y_t) - h_t^k(x_t, y_t)| + |f_t(\tilde{x}_{t+1}, y_t) - h_t^k(\tilde{x}_{t+1}, y_t)| \right) \leq 2 \langle L_t, \omega_t \rangle. \end{aligned}$$

Likewise, $\delta_t^2 \leq 2 \langle L_t, \omega_t \rangle$. Now the benign bound V_T' in Theorem 5 can be substituted with $\sum_{t=1}^T \langle L_t, \omega_t \rangle$. Drawing from Lemma 6, we have that

$$\begin{aligned} \sum_{t=1}^T \langle L_t, \omega_t \rangle &\leq \sum_{t=1}^T \langle L_t, 1_k \rangle + 2 \sqrt{(1 + \ln T) L_\infty \sum_{t=1}^T \langle L_t, 1_k \rangle} + O(\ln T) \\ &= \sum_{t=1}^T \langle L_t, 1_k \rangle + O(\sqrt{\ln T}) \sqrt{\sum_{t=1}^T \langle L_t, 1_k \rangle} + O(\ln T) \\ &\leq 2 \sum_{t=1}^T \langle L_t, 1_k \rangle + O(\ln T) \leq 2 \sum_{t=1}^T \rho(f_t, h_t^k) + O(\ln T), \quad \forall k \in 1:d, \end{aligned}$$

where 1_k denotes the d -dimensional one-hot vector with the k -th element being 1. Given the arbitrariness of k , we obtain

$$\sum_{t=1}^T \langle L_t, \omega_t \rangle \leq 2 \min_{k \in 1:d} \sum_{t=1}^T \rho(f_t, h_t^k) + O(\ln T),$$

In conclusion, the benign bound V_T' in Theorem 5 can be substituted with

$$\min_{k \in 1:d} \sum_{t=1}^T \rho(f_t, h_t^k) + O(\ln T).$$

Ignoring the dependence on poly-logarithmic factors yields the conclusion to be proved.

Algorithm 4 Program to Solve $w^* = \arg \min_{w \in \Delta_d^\alpha} \langle \ln(w/W), w \rangle$

Input: W, α

```

1:  $d \leftarrow |W|, \quad I \leftarrow 1:d, \quad C_\# \leftarrow 0, \quad C_\% \leftarrow 0$ 
2: while  $I \neq \emptyset$  do
3:    $w \leftarrow$  the median of  $W_I$ 
4:    $L \leftarrow \{i \mid i \in I, W_i < w\}, \quad M \leftarrow \{i \mid i \in I, W_i = w\}, \quad H \leftarrow \{i \mid i \in I, W_i > w\}$ 
5:   if  $w \frac{d - (C_\# + |L|)\alpha}{\|W\|_1 - (C_\% + \|W_L\|_1)} < \alpha$  then
6:      $C_\# \leftarrow C_\# + |L| + |M|, \quad C_\% \leftarrow C_\% + \|W_L\|_1 + \|W_M\|_1$ 
7:     if  $H = \emptyset$  then
8:        $w \leftarrow \min\{W_i \mid i \in I, W_i > w\}$ 
9:     end if
10:     $I \leftarrow H$ 
11:  else
12:     $I \leftarrow L$ 
13:  end if
14: end while

```

Output: $\forall i \in 1:d, \quad w_i^* \leftarrow \begin{cases} \frac{\alpha}{d} & W_i < w \\ \frac{W_i}{d} \frac{d - C_\#\alpha}{\|W\|_1 - C_\%} & W_i \geq w \end{cases}$

Lemma 6 (Static Regret for Clipped Hedge, Static Version of Corollary B.0.1 of [6]). Assume that $L_t \geq 0$, $\max_{t \in 1:T} \|L_t\|_\infty = L_\infty$ and $\alpha = d/T < 1$. If the learning rate follow from $\theta_t = (\ln T)/(\epsilon + \sum_{\tau=1}^{t-1} \sigma_\tau)$, where the constant $\epsilon > 0$ prevent θ_1 from being infinite, and $\sigma_t = \langle L_t, \omega_t - \omega_{t+1} \rangle - \text{KL}(\omega_{t+1}, \omega_t)/\theta_t$, Then the clipped Hedge enjoys the following static regret:

$$\text{Reg}_T u \leq 2 \sqrt{(1 + \ln T) L_\infty \sum_{t=1}^T \langle L_t, u \rangle} + O(\ln T), \quad \forall u \in \Delta_d^0.$$

G Solver for Clipped Hedge

The clipped Hedge equivalent to the following update:

$$\omega_{t+1} = \arg \min_{\omega \in \Delta_d^\alpha} \left\langle \ln \frac{\omega}{\omega_t \cdot \exp(-\theta_t L_t)}, \omega \right\rangle,$$

Thus, an efficient solution is attainable by minor adjustments to the algorithm depicted in Figure 3 of [16]. The modified algorithm is elaborated in Algorithm 4, wherein the primary alteration is the removal of the constraint $\|W\|_1 = 1$.