

Faster Rates for No-Regret Learning in General Games via Cautious Optimism

Ashkan Soleymani¹, Georgios Piliouras², and Gabriele Farina¹

¹MIT EECS, {ashkanso, gfarina}@mit.edu

²Google DeepMind, gpil@google.com

Abstract

We establish the first uncoupled learning algorithm that attains $O(n \log^2 d \log T)$ per-player regret in multi-player general-sum games, where n is the number of players, d is the number of actions available to each player, and T is the number of repetitions of the game. Our results exponentially improve the dependence on d compared to the $O(nd \log T)$ regret attainable by Log-Regularized Lifted Optimistic FTRL [Far+22c], and also reduce the dependence on the number of iterations T from $\log^4 T$ to $\log T$ compared to Optimistic Hedge, the previously well-studied algorithm with $O(n \log d \log^4 T)$ regret [DFG21]. Our algorithm is obtained by combining the classic Optimistic Multiplicative Weights Update (OMWU) with an adaptive, non-monotonic learning rate that paces the learning process of the players, making them more cautious when their regret becomes too negative.

Contents

1	Introduction	1
1.1	Overview of the Main Result and Techniques	2
1.2	Prior Work on Regularized Learning in General Games	4
1.3	Additional Related Work	5
2	Preliminaries	6
3	Dynamic Learning Rate Control	8
3.1	The Learning Rate Control Problem	10
3.2	Properties of the Learning Rate Control Objective	10
4	Extension to 0/1-Polyhedral Games via Kernels	12
5	Analysis of the Dynamic Learning Rate Control	18
5.1	Design of the DLRC-OMWU Algorithm and Equivalent Viewpoints	19
5.2	Proof Sketch	21

5.3	Detailed Analysis	24
5.3.1	Auxiliary Lemmas	24
5.3.2	Proofs for Learning Rate Control Step (Section 3.1)	24
5.3.3	Equivalent Viewpoints of DLRC-OMWU	30
5.3.4	Strong Spectral Properties of ψ	32
5.3.5	Positive Regret	36
5.3.6	Proofs for RVU Bounds (Section 5.2)	36
5.3.7	Proofs for Main Results	41
6	Conclusion	43
7	Acknowledgments	43

1 Introduction

The study of how multiple interacting agents learn to adapt their strategies is a well-established problem with significant foundations in game theory, online optimization, control theory, economics, and behavioral sciences [CL06; Nis+07; MS15; Gin00; Cam11]. This area has gained increased importance with the advent of machine learning, where multi-agent games are integral to many fundamental architectures and applications [Goo+14; Sil+17; Mor+17; BS19; Big+24]. Despite the intuitive appeal of simple uncoupled learning algorithms that globally converge to Nash equilibria in all games, strong impossibility results demonstrate that such algorithms do not exist—even in the classic and relatively constrained context of normal-form games [HM03; Mil+23]. From the perspective of centralized algorithms, Daskalakis, Goldberg, and Papadimitriou [DGP06] took a step further by linking the computational complexity of finding a Nash equilibrium in games to solving Brouwer’s fixed-point problem, showing that computing a Nash equilibrium in polynomial time is impossible unless $\text{PPAD} = \text{P}$, suggesting the problem is inherently computationally difficult.

These roadblocks have inspired the pursuit of alternative solution concepts, with the notion of *regret minimization* being arguably the most influential and well-studied among them [Sha12; Haz+16]. Originally defined within the context of single-agent optimization, regret measures the difference between the accumulated performance of an algorithm and that of the best fixed action in hindsight, under uncertain and possibly adversarial realizations of payoffs. Regret minimization in general games is sufficient to establish the time-average convergence of the empirical distribution of play to the set of Coarse Correlated Equilibria (CCE). This set is a relaxation of the Nash equilibrium concept—motivated by the intractability of Nash equilibrium—and possesses several desirable properties, including approximate optimality of the resulting social welfare, connections to Nash and correlated equilibria, broad applicability across a diverse range of games, and flexibility in strategy coordination [BEL06; NP10; Rou15; Cai+16; MP17; RST17]. Beyond convergence to CCE in game settings, regret minimization is recognized as a fundamental concept with wide-ranging applications, including learning and generalization [Lit88; Blu90; LN23], von Neumann’s minimax theorem and Blackwell approachability [FS96; ABH11], boosting [FS95; FS96], combinatorial optimization [PST95; AK07], complexity theory [KS99; BHK09], differential privacy [HR10; HRU13], prediction markets [CV10; ACV13], evolutionary dynamics [Cha+13], and more.

In the adversarial case, it is classically known that simple algorithms such as Hedge/Multiplicative Weights Update (MWU) suffice to establish the optimal regret rate of $O(T^{1/2})$ [FS97]. Nevertheless, an adversarial framework is often overly pessimistic for applications in game theory, as it may yield suboptimal results in more favorable and predictable environments. This is particularly evident in the context of learning in games, where player interactions are nonstationary but typically evolve slowly. Consequently, determining the tightest possible bounds for no-regret learning in general games remains a fundamental and unresolved problem. We enumerate major previous attempts in this line of work in Section 1.2.

By the *Optimism* framework of Syrgkanis, Agarwal, Luo, and Schapire [Syr+15], there exist algorithms for which the sum of players’ regrets in self-play remains constant over time. Unfortunately, controlling individual regrets—necessary for convergence to Coarse Correlated Equilibria (CCE)—is *much harder*, and despite nearly a decade of noteworthy progress, the optimal regret rate remains elusive. Intuitively, we seek to avoid agents with runaway negative regret. Our main high-level idea is to adapt optimism into a form of *Cautious Optimism*, where agents decrease their learning rates when their regret becomes too negative—that is, when they significantly outperform all fixed actions. Surprisingly, we show that it is possible to carefully apply this germ of an idea to

achieve new state-of-the-art regret bounds.

1.1 Overview of the Main Result and Techniques

We establish the first *uncoupled* no-regret learning algorithm that attains $O(n \log^2 d \log T)$ regret in multi-player general-sum games. Our algorithm is best understood as an Optimistic Multiplicative Weights Update (OMWU) paired with a *dynamic* learning rate that is adjusted based on the regret accumulated by the learner. For this reason, we coin our method *Dynamic Learning Rate Control OMWU (DLRC-OMWU)*. A key characteristic that sets our approach apart from standard adaptive learning-rate techniques is that our dynamic control does not produce monotonically decreasing learning rates. Instead, the goal of our dynamic learning rate is to properly pace the learner—*slowing it down when it is performing too well*—that is, when its maximum regret becomes too negative. We achieve this by solving an optimization problem at each iteration to determine the learning rate, based on the player’s current regret vector.

The idea of agents differentiating their behavior depending on whether they feel content or discontent is both simple and intuitive, and has inspired other game dynamics that provably concentrate around pure Nash equilibria [You09], as well as adjusted replicator dynamics in evolutionary game theory [Wei97]. In the context of regret minimization, this idea can be traced back to the work of Bowling and Veloso [BV02], who introduced the Win or Learn Fast (WoLF) principle. This principle increases the learning rate of agents when they are losing, thereby resulting in quicker adaptation to the environment and to the strategies of other agents. Bowling and Veloso [BV02] demonstrated the convergence of gradient ascent-descent with WoLF to a Nash equilibrium in the restricted case of two players with two actions. Bowling [Bow04] extended WoLF to multiplayer settings, showing it achieves no-regret dynamics with $O(d\sqrt{T})$ regret. Despite the success of WoLF and related heuristics in small-scale games [AL08; BP03; Blo+15; Kai+09; LP22], their theoretical guarantees in general games remain unclear.

On the negative side, it has recently been shown that the multiplicative weights update algorithm, even when equipped with a continuous learning rate (rather than the original fast and slow rates of WoLF [BV02]), exhibits chaotic behavior in nonatomic congestion games [VFP23]. While our motivation for adaptive learning rates and our methodology differ from the WoLF principle, to our knowledge, DLRC-OMWU is the *first algorithm to demonstrate theoretical benefits of such ideas in regret minimization for general games*.

The regret guarantees of DLRC-OMWU exponentially improve the dependence on d in the $O(nd \log T)$ regret attained by Log-Regularized Lifted Optimistic FTRL [Far+22c]. Compared to the regret analysis of OMWU in Daskalakis, Fishelson, and Golowich [DFG21], we improve the regret bound in several aspects. Not only do we reduce the dependence on the time horizon T from $\log^4 T$ to $\log T$, but our guarantees also hold across all regimes of n , d , and T , whereas theirs require $T > Cn \log^4 T$ for some constant C (see Section 5, and Lemmas 4.2 and C.4 in [DFG21]). Additionally, as a minor note, the regret analysis of OMWU [DFG21] hides extremely large constants in the asymptotic notation, especially in comparison to our analysis.

Technical Contributions. The technical novelties of this work, which lead to fast no-regret learning rates, are multifaceted and can be briefly summarized as follows.

First, as discussed previously, we conceptualize the idea of learning rate control for no-regret learning (see Sections 3 and 5.1) and formalize it to design optimization algorithms for learning rate control, resulting in DLRC-OMWU, a computationally efficient algorithm (see Section 3). The

Method	Games' Regret	Cost per Iteration	Adversarial Regret
OFTRL / OMD [Syr+15]	$O(\sqrt{n} \mathfrak{R} T^{1/4})$	Regularizer- dependent & oracle-	$\tilde{O}(\sqrt{T \log d})$
OMWU [CP20a]	$O(n \log^{5/6} d T^{1/6})$ †	$O(d)$	$\tilde{O}(\sqrt{T \log d})$
OMWU [DFG21]	$O(n \log d \log^4 T)$	$O(d)$	$\tilde{O}(\sqrt{T \log d})$
Clairvoyant MWU [PSS22]	$O(n \log d)$ for a subsequence only ‡	$O(d)$	No guarantees
LRL-OFTRL [Far+22c]	$O(n d \log T)$	$O(d \log \log T)$	$\tilde{O}(\sqrt{T \log d})$
DLRC-OMWU [This paper]	$O(n \log^2 d \log T)$	$O(d \log \log T)$	$\tilde{O}(\sqrt{T \log d})$

Table 1: Comparison of prior results on minimizing external regret in general games. For simplicity, we omit dependencies on the smoothness and range of the utility functions. We use n to denote the number of players, T the number of repetitions of the game, and d the number of actions. \mathfrak{R} denotes the maximum value attained by the regularizer. † Limited to two-player games only ($n = 2$). ‡ Unlike all other algorithms, the full sequence of iterates produced by Clairvoyant MWU (CMWU) is not known to achieve sublinear regret. Instead, after running CMWU for T iterations, only a smaller subsequence of length $\Theta(T/\log T)$ is known to attain the regret stated in the table.

concept of dynamic learning rates has the potential to be beneficial in other areas involving regret minimization, particularly in multi-agent settings.

Secondly, using analysis techniques based on the properties of self-concordant functions, we establish strong *multiplicative stability* results for the resulting learning rate, even when the actions themselves are not known to be multiplicatively stable (see Sections 3.1 and 5.3.2). This approach enables predictability of the dynamics beyond the constant learning rate setting. Previously, it was unclear how to design learning dynamics that evolve smoothly while simultaneously adapting to changes induced by the learning processes of other agents.

Thirdly, we demonstrate that combining the Optimistic Multiplicative Weights Update algorithm with our learning rate control can be viewed as an instantiation of optimistic Follow-the-Regularized-Leader (FTRL) with a novel regularizer ψ . We further explore equivalent viewpoints of DLRC-OMWU, revealing connections to the Lifted Optimistic FTRL algorithm proposed by Farina, Anagnostides, Luo, Lee, Kroer, and Sandholm [Far+22c], albeit with a different regularizer (see Sections 5.1 and 5.3.3), where the FTRL dynamics are executed over the lifted space.

Fourthly, we introduce a novel regularizer, ψ , whose spectral properties lie between those of the logarithmic and entropic regularizers. We establish *strong spectral properties* for this regularizer, including strong convexity and high curvature (see Sections 5.2 and 5.3.4). This approach allows us to depart from the traditional methodology of using intrinsic norms induced by the Hessian matrix, which is prevalent in the online learning literature. We note that the regularizer ψ may be of independent interest in future studies.

Lastly, it is important to mention that, in Section 4, inspired by Kernelized OMWU [Far+22b], we introduce a kernelized version of DLRC-OMWU, (denoted as KDLRC-OMWU), extended to convex 0/1-polyhedral games such as extensive-form games or flows on directed graphs. This way, we

demonstrate that DLRC-OMWU inherits the fundamental and intriguing properties of OMWU.

While we adopt the intriguing nonnegative RVU property idea from Log-Regularized Lifted OFTRL (LRL-OFTRL) [Far+22c] (see Sections 5 and 5.3.5), our construction and proof technique differ significantly. The analysis in LRL-OFTRL relies heavily on the intrinsic norm of the logarithmic regularizer to ensure multiplicative stability of the iterates. In contrast, our approach does not depend on multiplicative stability or intrinsic norms. Instead, it involves a more nuanced analysis, focusing on the dynamic learning rate control optimization problem and the strong spectral properties of our specifically chosen regularizer.

1.2 Prior Work on Regularized Learning in General Games

Starting with the seminal paper of Syrgkanis, Agarwal, Luo, and Schapire [Syr+15], which established the first $o(T^{1/2})$ regret bound for self-play in general games, a race was initiated to provide the tightest possible bounds in this setting. In particular, Syrgkanis, Agarwal, Luo, and Schapire [Syr+15] identified the *RVU property*, an adversarial regret bound applicable to a broad class of so-called *Optimistic* variants of standard no-regret learning algorithms, such as Follow-the-Regularized-Leader and Mirror Descent. Using this property, they demonstrated that the individual regret of each player grows as $O(T^{1/4})$.

In a more recent breakthrough, Daskalakis, Fishelson, and Golowich [DFG21] exponentially improved the regret guarantees for Optimistic MWU (OMWU), achieving the first polylogarithmic in T regret bound of $O(n \log d \log^4 T)$, where n is the number of players and d is the number of actions per player. Their analysis relies heavily on discrete-time Fourier transforms and the smoothness of higher-order discrete differentials of the learning dynamics in the frequency domain.

Clairvoyant MWU (CMWU) advanced the concept of prediction even further [PSS22], albeit at the cost of creating an algorithm that does not minimize regret in adversarial settings. The key intuition is that, given a perfect prediction of future payoff sequences, a bounded total payoff can be easily achieved by applying a Be-the-Leader type of algorithm. Although such approaches are not feasible in standard online learning, agents can implement such a sequence of play in game settings via uncoupled learning algorithms. The full sequence of iterates produced by CMWU does not achieve sublinear regret. Instead, after running CMWU for T iterations, only a smaller subsequence of length $\Theta(T/\log T)$ maintains a bounded total regret of $O(n \log d)$. Combining these results, the effective convergence rate toward CCE is $O\left(\frac{n \log d \log T}{T}\right)$. For a more detailed discussion of how CMWU differs from other algorithms in this line of work, see Section 1.3.

In the last major step prior to our work, Farina, Anagnostides, Luo, Lee, Kroer, and Sandholm [Far+22c] achieved logarithmic dependence on T over the entire history of play via the Log-Regularized Lifted Optimistic FTRL algorithm. Their learning dynamics are based on an instantiation of Optimistic Follow-the-Regularized-Leader over an appropriately *lifted* space using a logarithmic regularizer. Interestingly, their approach generalizes beyond normal-form games to encompass general convex games. However, this analysis comes at the cost of exponentially worse dependence on the number of actions d , with an overall regret guarantee of $O(nd \log T)$.

Table 1 provides a summary of prior work aimed at establishing optimal regret bounds for no-regret learners in games.

1.3 Additional Related Work

The evolution of regret minimization in games has seen significant breakthroughs, beginning with the pioneering work of Daskalakis, Deckelbaum, and Kim [DDK11] on zero-sum games. They developed *strongly uncoupled* learning dynamics that achieve a regret growth rate of $O(\log T)$. This milestone was further refined by Rakhlin and Sridharan [RS13b], who introduced *Optimistic Mirror Descent* (OMD), simplifying the implementation while maintaining robust performance. The practical benefits of recency bias in OMD have also been substantiated in behavioral economics [FP14].

Beyond their near-optimal performance, an additional advantage of optimistic learning algorithms is their incorporation of a straightforward and intuitive recency bias, which aligns well with real-world behavioral biases [NE12; EH13]. Each agent optimistically assumes that all other agents will play tomorrow exactly as they did today, and then applies their online learning algorithm, incorporating this extrapolation step into the decision-making process. This alignment with behavioral data has two key benefits. First, it suggests that the resulting theoretical analysis may offer insights into human behavior in everyday interactions. Second, it highlights the potential for behavioral science to inspire algorithmic modifications with even stronger performance guarantees.

Farina, Lee, Luo, and Kroer [Far+22b] extended these theoretical advancements beyond normal-form games, generalizing the $O(\text{polylog}(T))$ regret bounds to polyhedral games, which encompass extensive-form games. The introduction of *no-swap-regret* dynamics has also led to fast convergence to correlated equilibria in normal-form games [CP20a; Ana+22a; Ana+22d]. Wei and Luo [WL18] further advanced the field by leveraging optimism to achieve adaptive regret bounds in bandit settings.

The focus on *last-iterate* convergence has brought renewed interest, with significant contributions highlighting the performance of Optimistic Mirror Descent (OMD) algorithms [DP19; Mer+19; GPD20; Lei+21; HAM21; Wei+21; Azi+21]. These works examine the conditions under which OMD algorithms achieve favorable iterate convergence properties.

Parallel to these developments, extensive research has examined the dynamics of learning algorithms in games. This includes studies that do not incorporate optimism and often reveal complex behaviors such as divergence, recurrence, or chaos [Har+03; Das+10; KLP11; BCM12; PP18; BP18; MPP18; BP19; CP19; Vla+20; BGP20; CP20b]. These findings underscore the intricate and often unstable nature of learning dynamics in strategic environments, highlighting both the challenges and opportunities that remain in this vibrant area of research.

On Clairvoyant MWU. As a specific uncoupled learning algorithm, Clairvoyant MWU (CMWU) was introduced by Piliouras, Sim, and Skoulakis [PSS22] to compute the Coarse Correlated Equilibrium (CCE) of a normal-form game. It is known that, given a perfect prediction of future payoff sequences, a bounded total regret of $O(n \log d)$ can be easily achieved by applying a Be-the-Leader type algorithm. Although such approaches are not feasible in online learning, in self-play settings—under certain predetermined agreements among the players—it is possible to implement a sequence of play that is informed by hindsight of the game’s next action.

Farina, Kroer, Lee, and Luo [Far+22a] showed that CMWU is equivalent to Nemirovski’s Conceptual Proximal Method [Nem04]. Given the proximal operator of the game (in this setting, the online mirror descent optimization step of all players combined), if the players play the fixed point of this operator, they are guaranteed to achieve constant regret. In the setting of normal-form games, this operator is proven to be monotone for a sufficiently small choice of learning rate [PSS22; Far+22a]. Thus, its fixed point can be computed—either in a centralized or decentralized manner.

The key idea in Piliouras, Sim, and Skoulakis [PSS22] is to use self-play communication as a broadcast channel, allowing players to collaboratively find this fixed point by broadcasting their updates through in-game actions. In this setup, for $T - O(\log T)$ iterations, players are merely communicating indirectly. Only once every predetermined $O(\log T)$ rounds do they play the approximate fixed point they have collectively computed. Consequently, only $O(\log T)$ actual gameplay iterations occur, and the empirical distribution over these predetermined rounds converges to the CCE of the game, resulting in an effective convergence rate of $O\left(\frac{n \log d \log T}{T}\right)$.

While this approach yields faster rates, it comes with two caveats: (i) it does not guarantee sublinear regret in adversarial settings or even if a single player deviates slightly; (ii) players must have prior agreements on how to self-play and must strictly adhere to the protocol. For instance, if one player begins a round early or late, the communication process fails, as the true gameplay during the $O(\log T)$ periods becomes mixed with the fixed-point computation steps of other players. Due to these issues, this algorithm is not applicable to large-scale games in real-world applications.

2 Preliminaries

For any $d \in \mathbb{N}$, we denote the set $\{1, 2, \dots, d\}$ by $[d]$. We use bold letters to denote vectors, for example, $\mathbf{x} \in \mathbb{R}^d$. We let $\mathbf{x}[r]$ denote the r -th coordinate of the vector \mathbf{x} , for any $r \in [d]$. In our notation, players are generally indicated by subscripts, which we omit whenever the results involve a generic player and are clear from context, to avoid overloading the notation. The time index, represented by the variable t , is indicated using superscripts. For example, $\mathbf{x}_i^{(t)}$ denotes the action played by player i at time t . We denote the set of distributions supported on a finite set \mathcal{A} of size $d = |\mathcal{A}|$ by $\Delta^d := \Delta(\mathcal{A})$. For given vectors \mathbf{x} and \mathbf{u} , we denote their inner product by $\langle \mathbf{x}, \mathbf{u} \rangle$. For any vector $\mathbf{x} \in \mathbb{R}^d$, we write $\Lambda(\mathbf{x}) := \sum_{k=1}^d \mathbf{x}[k]$ for the sum of its elements. We also define the negative entropy by $H(\mathbf{x}) = \sum_{k=1}^d \mathbf{x}[k] \log \mathbf{x}[k]$ and the Kullback–Leibler (KL) divergence by $\text{KL}(\mathbf{x} \parallel \mathbf{x}') = \sum_{k=1}^d \mathbf{x}[k] \log \frac{\mathbf{x}[k]}{\mathbf{x}'[k]}$. Lastly, we use the notation $\mathbf{1}_d$ to denote the vector $[1, 1, \dots, 1] \in \mathbb{R}^d$.

General-sum Games. We consider n -player general-sum games and denote the set of players by $[n]$. In a normal-form game, each player $i \in [n]$ has a finite and nonempty set of deterministic strategies \mathcal{A}_i . The set of mixed strategies for player i is given by the probability simplex over \mathcal{A}_i , denoted $\mathcal{X}_i = \Delta(\mathcal{A}_i)$. The joint action space of the players is then $\times_{j=1}^n \mathcal{A}_j$. Each player i is associated with a utility function $\mathcal{U}_i : \times_{j=1}^n \mathcal{A}_j \rightarrow \mathbb{R}$ defined over joint deterministic strategies. Let the expected utility of player i under a joint mixed strategy profile $\mathbf{x} = (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n) \in \times_{j=1}^n \Delta(\mathcal{A}_j)$ be denoted by the multilinear extension:

$$\boldsymbol{\nu}_i(\mathbf{x}) := \mathbb{E}_{\mathbf{s} \sim \mathbf{x}}[\mathcal{U}_i(\mathbf{s})] = \langle \mathbf{x}_i, \nabla_{\mathbf{x}_i} \boldsymbol{\nu}_i(\mathbf{x}) \rangle.$$

Let $d := \max_{i \in [n]} |\mathcal{A}_i|$ be the maximum number of actions available to any player. Without loss of generality, we assume that $|\mathcal{A}_i| \geq 2$ for each player i , since players with only one action can be removed from the game without loss of generality. For simplicity of notation, we further assume that all players have exactly d actions. Finally, we analyze the game under the following standard assumption regarding the boundedness and smoothness of the utility functions.

Assumption 1. The utility function $\boldsymbol{\nu}_j$ of the game (for each player $j \in [n]$) satisfies, for any two strategy profiles $\mathbf{x}, \mathbf{x}' \in \times_{j=1}^n \Delta(\mathcal{A}_j)$,

- (Bounded Utilities) $\max_{s \in \times_{i=1}^n \mathcal{A}_i} |\mathcal{U}_j(s)| \leq 1$, for each player $j \in [n]$.
- (L -smoothness) there exists a positive number $L > 0$ such that for each player $j \in [n]$,

$$\|\nabla_{\mathbf{x}_j} \boldsymbol{\nu}_j(\mathbf{x}) - \nabla_{\mathbf{x}_j} \boldsymbol{\nu}_j(\mathbf{x}')\|_\infty \leq L \sum_{i \in [n]} \|\mathbf{x}_i - \mathbf{x}'_i\|_1.$$

This assumption is general and not restrictive, as any bounded utility function can be rescaled—without loss of generality—to satisfy the boundedness condition. Moreover, it is straightforward to show that under this assumption, the game is L -smooth with $L = 1$. However, depending on the structure of the game, the smoothness parameter L may be substantially smaller than 1. For example, in games with vanishing sensitivity ϵ_n , it can be observed that $L = \epsilon_n \ll 1$ [AS24]. Here, we distinguish L from the boundedness assumption in order to account for additional structural properties that may lead to faster convergence guarantees.

No-regret Learning. In the *online learning framework*, a learning agent is required to select a strategy $\mathbf{x}^{(t)} \in \mathcal{X} = \Delta^d$ at each time $t \in \mathbb{N}$. We consider the *full-information* model, in which the environment provides feedback in the form of a *linear* utility function $\mathbf{x} \mapsto \langle \mathbf{x}, \boldsymbol{\nu}^{(t)} \rangle$, where $\boldsymbol{\nu}^{(t)} \in \mathbb{R}^d$. The principal measure of the agent’s performance is *external regret* (also referred to as *regret*), defined over a time horizon $T \in \mathbb{N}$ as:

$$\text{Reg}^{(T)} := \max_{\mathbf{x}^* \in \Delta^d} \left\{ \sum_{t=1}^T \langle \boldsymbol{\nu}^{(t)}, \mathbf{x}^* \rangle \right\} - \sum_{t=1}^T \langle \boldsymbol{\nu}^{(t)}, \mathbf{x}^{(t)} \rangle. \quad (2.1)$$

This metric evaluates the agent’s performance against the best possible *fixed* strategy chosen in hindsight. The goal is to ensure that regret grows as slowly as possible with the time horizon T . It is well known that standard algorithms such as Online Mirror Descent (OMD) and Follow-the-Regularized-Leader (FTRL) achieve the optimal regret rate of $\text{Reg}^{(T)} = \Theta(\sqrt{T})$ in adversarial settings [Haz+16; Ora19]. However, in *game self-play* settings, faster rates are achievable, since the dynamics of each player’s online learning process follow the same learning algorithm, in contrast to the adversarial case. Specifically, each player $i \in [n]$ at time step t receives a utility vector $\boldsymbol{\nu}^{(t)} = \nabla_{\mathbf{x}_i} \boldsymbol{\nu}_i(\mathbf{x}^{(t)})$, where $\mathbf{x}^{(t)} = (\mathbf{x}_1^{(t)}, \mathbf{x}_2^{(t)}, \dots, \mathbf{x}_n^{(t)})$ denotes the joint play at time t . It is important to note that a player’s regret may be negative.

An important class of no-regret learning algorithms that achieve faster convergence guarantees in games are those that satisfy the *Regret bounded by Variation in Utilities (RVU)* property—a concept introduced by Syrgkanis, Agarwal, Luo, and Schapire [Syr+15], which we formalize below.

Definition 2 (RVU property [Syr+15]). A no-regret learning algorithm satisfies Regret bounded by Variation in Utilities (RVU) property with parameters a, b, c and a pair of dual norms $(\|\cdot\|, \|\cdot\|_*)$, if the regret on any sequence of utilities $\{\boldsymbol{\nu}^{(t)}\}_{t=1}^T$ is upper bounded by

$$\text{Reg}^{(T)} \leq a + b \sum_{t=1}^T \|\boldsymbol{\nu}^{(t)} - \boldsymbol{\nu}^{(t-1)}\|_*^2 - c \sum_{t=1}^T \|\mathbf{x}^{(t+1)} - \mathbf{x}^{(t)}\|^2,$$

where $a \geq 0$ and $0 < b \leq c$.

It has been shown that optimistic variants of OMD and FTRL algorithms—originally introduced by Rakhlin and Sridharan [RS13a]—satisfy the RVU property [Syr+15], and achieve faster rates in the self-play setting compared to the adversarial one. These algorithms also attain the optimal convergence rate for the *sum of regrets*; however, it remains an open question whether the RVU property alone suffices to guarantee optimal individual no-regret learning rates in games.

Coarse Correlated Equilibrium. A probability distribution $\sigma \in \Delta\left(\times_{j=1}^n \mathcal{A}_j\right)$ over the joint action space $\times_{j=1}^n \mathcal{A}_j$ is called an ϵ -*approximate coarse correlated equilibrium (CCE)* if, for every player $i \in [n]$ and every deterministic strategy $s'_i \in \mathcal{A}_i$, unilateral deviations do not increase the expected utility by more than ϵ , that is:

$$\mathbb{E}_{\mathbf{s} \sim \sigma} [\mathcal{U}_i(\mathbf{s})] \geq \mathbb{E}_{\mathbf{s} \sim \sigma} [\mathcal{U}_i(s'_i, \mathbf{s}_{-i})] - \epsilon,$$

where $\mathbf{s} = (s_1, s_2, \dots, s_n)$ denotes the joint action drawn from σ , and \mathbf{s}_{-i} represents the actions of all players other than i .

A folklore result in game theory establishes a connection between no-regret learning algorithms and coarse correlated equilibria (CCE) of a game [CL06]. When the players follow no-regret learning algorithms with regrets $\text{Reg}_i^{(T)}$ for all $i \in [n]$, the empirical play of the game, defined as $\sigma = \frac{1}{T} \sum_{t=1}^T \mathbf{x}^{(t)}$, constitutes an $\left(\frac{1}{T} \max_{i \in [n]} \text{Reg}_i^{(T)}\right)$ -approximate CCE of the underlying game. Hence, faster regret rates in self-play settings lead directly to faster convergence to the game’s coarse correlated equilibrium.

3 Dynamic Learning Rate Control

We restate that our algorithm is a variant of the Optimistic Multiplicative Weights Update algorithm (OMWU), in which the learning rate is selected adaptively based on the regret accumulated up to any point in time. For this reason, we refer to our method as *Dynamic Learning Rate Control OMWU (DLRC-OMWU)*.

In its standard version, OMWU picks the next distribution of play proportionally to the exponential of the optimistic regret accumulated on each action, that is, according to the formula

$$\mathbf{x}^{(t)}[k] := \frac{\exp\{\lambda^{(t)} \mathbf{r}^{(t)}[k]\}}{\sum_{k' \in \mathcal{A}} \exp\{\lambda^{(t)} \mathbf{r}^{(t)}[k']\}} \quad \forall k \in \mathcal{A}, \quad (3.1)$$

where $\lambda^{(t)} > 0$ is a learning rate and $\mathbf{r}^{(t)}$ is the vector of optimistically-corrected regrets accumulated by each action up to time t ,

$$\mathbf{r}^{(t)}[k] := (\boldsymbol{\nu}^{(t-1)}[k] - \langle \boldsymbol{\nu}^{(t-1)}, \mathbf{x}^{(t-1)} \rangle) + \sum_{\tau=1}^{t-1} [\boldsymbol{\nu}^{(\tau)}[k] - \langle \boldsymbol{\nu}^{(\tau)}, \mathbf{x}^{(\tau)} \rangle] \quad \forall k \in \mathcal{A}.$$

For simplicity, define the corrected reward signal $\mathbf{u}^{(t)} := \boldsymbol{\nu}^{(t)} - \langle \boldsymbol{\nu}^{(t)}, \mathbf{x}^{(t)} \rangle \mathbf{1}_d$ and the accumulated signal $\mathbf{U}^{(t)} := \sum_{\tau=1}^{t-1} [\boldsymbol{\nu}^{(\tau)} - \langle \boldsymbol{\nu}^{(\tau)}, \mathbf{x}^{(\tau)} \rangle \mathbf{1}_d]$.

It was recently discovered that when all players in an n -player game learn using OMWU with a suitable *constant* learning rate $\lambda^{(t)} = \eta$, the maximum regret accumulated by the players grows at most as $O(n \log d \log^4 T)$ as a function of the time horizon, for sufficiently large T [DFG21].

In this paper, we show that the previous result can be improved to a $\log T$ dependence by using a different approach based on a novel technique that we term *dynamic learning rate control*. Unlike most of the literature on adaptive learning rate schedules in online learning and optimization, our dynamic control does not produce monotonically decreasing learning rates and conceptually, it is *not* designed as a means of circumventing uncertainty about the problem’s conditioning (e.g., the Lipschitz constants). Rather, our dynamic learning rate control aims to pace the learner—*slowing it down when it is performing too well*—that is, *when its maximum regret is too negative*.

Such a goal might appear backwards—after all, if a learner is doing so well, why pace them down? The contradiction is resolved when considering the learning system as a whole. Intuitively, if one player is doing too well, the other players might be unable to catch up. Instead, when all players’ regrets remain nonnegative, one is able to show desirable overall properties of the learning process. These include not only small swap regret [Ana+22d], but also iterate convergence to equilibrium [Ana+22b], and discovery of strongly incentive-compatible equilibria [Ana+22c]. Alternatively, there is another way to look at this idea. As discussed in Section 2, no-regret dynamics lead to convergence to coarse correlated equilibria at a rate determined by the *worst-performing* player—the player with the maximum regret. Therefore, it seems that for faster convergence, enforcing some form of performance balance among the players of the game is natural.

We present our algorithm, *Dynamic Learning Rate Control Optimistic MWU (DLRC-OMWU)*, in Algorithm 1. We prove that when the players of the game follow DLRC-OMWU (i.e., in the self-play setting), each player experiences $O(n \log^2 d \log T)$ regret, while in the adversarial setting, each player suffers at most $\tilde{O}(\sqrt{T \log d})$ regret. Therefore, DLRC-OMWU is robust to adversarial behavior. The dynamics of DLRC-OMWU are simple and mirror those of OMWU, except that when the maximum regret becomes too negative (Line 6 in Algorithm 1), the learning rate in the Multiplicative Weights Update step is adjusted dynamically (Line 7 in Algorithm 1). Our theoretical guarantees in Theorem 3.1 hold across all regimes of n , d , and T , as a constant choice of η is sufficient for the result.

Theorem 3.1 (Informal; see Theorem 5.3 for the detailed version). *Suppose that n players self-play a general-sum multiplayer game with a finite set of d deterministic strategies per player over T rounds. Further, suppose that each player follows DLRC-OMWU to choose their action based on the history so far. Then, each player incurs $O(n \log^2 d \log T)$ regret. Moreover, when faced with adversarial utilities, each player playing DLRC-OMWU is guaranteed to experience $\tilde{O}(\sqrt{T \log d})$ regret.*

An immediate consequence of Theorem 3.1 is that the empirical distribution of play constitutes an $O\left(\frac{n \log^2 d \log T}{T}\right)$ -approximate coarse correlated equilibrium (CCE) of the game. This corollary follows from the fact that, in a general-sum multiplayer game, if each player’s regret is at most $\epsilon(T)$, then the empirical distribution of their joint strategies converges to a CCE at a rate of $O(\epsilon(T)/T)$.

Corollary 3.2. *If n players employ the uncoupled learning dynamics of DLRC-OMWU for T rounds in a general-sum multiplayer game with a finite set of d deterministic strategies per player, then the empirical distribution of play forms an $O\left(\frac{n \log^2 d \log T}{T}\right)$ -approximate coarse correlated equilibrium (CCE) of the game.*

Algorithm 1: Dynamic Learning Rate Control - Optimistic MWU (DLRC-OMWU)

Data: Learning rate η , parameters α and β

```
1 Set  $\mathbf{U}^{(1)}, \mathbf{u}^{(0)} \leftarrow \mathbf{0} \in \mathbb{R}^d$ 
2 for  $t = 1, 2, \dots, T$  do
3   Set  $\mathbf{r}^{(t)} \leftarrow \mathbf{U}^{(t)} + \mathbf{u}^{(t-1)}$  [▷ Optimism]
4   if  $\max_{k \in [d]} \{\mathbf{r}[k]\} \geq -\beta \log^2 d$  then
5     Set  $\lambda^{(t)} \leftarrow \eta$ 
6   else
7     /* Dynamic Learning Rate Control */
8     Set  $\lambda^{(t)} \leftarrow \arg \max_{\lambda \in (0, \eta]} \left\{ \log \left( \sum_{k=1}^d e^{\lambda \mathbf{r}^{(t)}[k]} \right) + (\alpha - 1) \log \lambda \right\}$ 
9     Play strategy  $\mathbf{x}^{(t)}[k] := \frac{\exp\{\lambda^{(t)} \mathbf{r}^{(t)}[k]\}}{\sum_{k'=1}^d \exp\{\lambda^{(t)} \mathbf{r}^{(t)}[k']\}}$  [▷ OMWU]
10    Observe  $\boldsymbol{\nu}^{(t)} \in \mathbb{R}^d$ 
11    Set  $\mathbf{u}^{(t)} \leftarrow \boldsymbol{\nu}^{(t)} - \langle \boldsymbol{\nu}^{(t)}, \mathbf{x}^{(t)} \rangle \mathbf{1}_d$ 
12    Set  $\mathbf{U}^{(t+1)} \leftarrow \mathbf{U}^{(t)} + \mathbf{u}^{(t)}$  [▷ Empirical cumulated regrets]
```

3.1 The Learning Rate Control Problem

At every time t , our learning algorithm outputs strategies according to (3.1), where the learning rate is carefully chosen as the solution to the following optimization problem:

$$\lambda^{(t)} := \arg \max_{\lambda \in (0, \eta]} \left\{ f(\lambda; \mathbf{r}^{(t)}) := (\alpha - 1) \log \lambda + \log \left(\sum_{k=1}^d e^{\lambda \mathbf{r}^{(t)}[k]} \right) \right\}, \quad (3.2)$$

where $\eta > 0$ is a constant that caps the maximum learning rate, and α is a key parameter of the algorithm. We set α to be on the order of $\Theta(\log^2 d)$ to achieve the guarantees mentioned in the introduction.

Under normal operating conditions, where the maximum regret $\max_k \{\mathbf{r}^{(t)}[k]\}$ accumulated on the actions is not too negative, it is immediate to observe that the optimal solution is $\lambda^{(t)} = \eta$ ¹ thus recovering the usual operating regime of a constant learning rate. However, when the maximum regret becomes sufficiently negative, the optimal value of $\lambda^{(t)}$ starts decreasing towards 0, causing the corresponding player to degrade performance by disregarding the history of the game. In the extreme case where $\lambda^{(t)} \rightarrow 0$, the player acts uniformly among its actions. Figure 1 illustrates the value of the learning rate $\lambda^{(t)}$ as a function of the optimistic cumulative regrets in a simple two-action case. It can be observed that $\lambda^{(t)}$ is a monotonically non-increasing function of $\{\mathbf{r}^{(t)}[1]\}$ and $\{\mathbf{r}^{(t)}[2]\}$.

3.2 Properties of the Learning Rate Control Objective

At first glance, it may not be immediately apparent that the maximization objective is tractable, or even concave in λ . However, the following result demonstrates that when α is chosen on the

¹Please refer to Lemma 5.16 for a concrete proof.

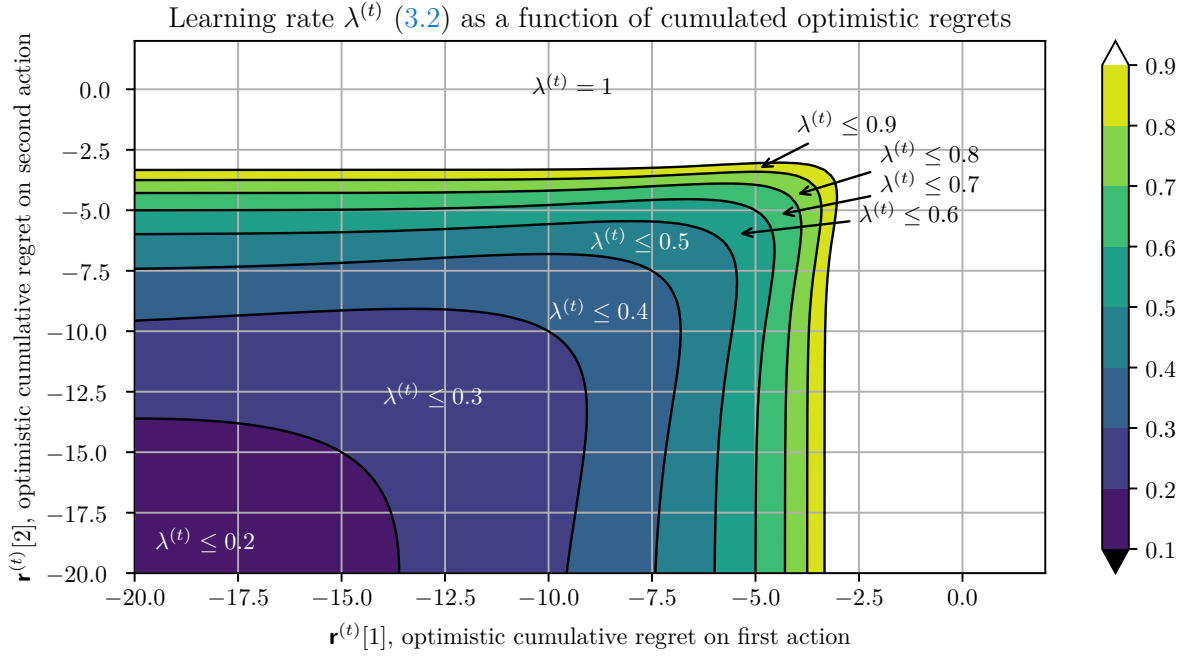


Figure 1: **Learning rate landscape:** Dependence of $\lambda^{(t)}$, as defined in (3.2), on the optimistic regrets cumulated in 2-action simplex. For the plot, the values $\eta = 1, \alpha = 4$ were chosen.

order of $\Omega(\log^2 d)$, $f(\lambda; \mathbf{r})$ becomes self-concordant (and thus strongly concave) in λ . A concrete proof is provided in Section 5.3.2.

Theorem 3.3. *For any $\mathbf{r} \in \mathbb{R}^d$, the rate control objective $f(\lambda; \mathbf{r})$ defined in (3.2) satisfies the following properties:*

- *Strong concavity:* $f''(\lambda; \mathbf{r}) \leq -(\alpha - \log^2 d - 1)/\lambda^2$ for all $\lambda \in (0, \infty)$.
- *Self-concordance:* $(f'''(\lambda; \mathbf{r}))^2 \leq -4f''(\lambda; \mathbf{r})^3$,

where all derivatives are with respect to λ .

An immediate byproduct of Theorem 3.3 is that the learning rate control optimization problem is well-defined and admits a unique solution for a sufficiently large choice of the hyperparameter α .

Remark 3.4. The learning rate control objective $f(\lambda; \mathbf{r})$ defined in (3.2) admits a unique solution $\lambda^{(t)} \in (0, \eta]$ for any given $\eta \geq 0$ and all $t \in \mathbb{N}$, whenever α is sufficiently large.

A key technical property used in the analysis of the regret accumulated by DLRC-OMWU is the sensitivity of the solution to (3.2) with respect to small perturbations in the regret vector $\mathbf{r}^{(t)}$. This property is crucial, as it ensures that the dynamic learning rate $\lambda^{(t)}$ for each player changes smoothly over time t , thereby allowing the self-play process to evolve in a stable and smooth manner.

Theorem 3.5. *(Sensitivity of learning rates on regrets) There exists a universal constant β ,² such that for $\alpha \geq 2 + 2 \log d + \beta \log^2 d$, the following property holds. Let $\mathbf{r}, \mathbf{r}' \in \mathbb{R}^d$ be such that*

²For concrete values, choosing any $\beta \geq 70$ suffices.

$\|\mathbf{r} - \mathbf{r}'\|_\infty \leq 2$, and let $\hat{\lambda}, \hat{\lambda}'$ be the corresponding learning rates, that is,

$$\hat{\lambda} = \arg \max_{t \in (0, \eta]} f(t; \mathbf{r}), \quad \hat{\lambda}' = \arg \max_{t \in (0, \eta]} f(t; \mathbf{r}').$$

Then, $\hat{\lambda}$ and $\hat{\lambda}'$ are multiplicatively stable; specifically,

$$\frac{7}{10} \leq \frac{\hat{\lambda}}{\hat{\lambda}'} \leq \frac{7}{5},$$

The proof of Theorem 3.5 is derived by combining an accurate analytical estimate, λ_0 , for the value of $\hat{\lambda}$ with techniques from the analysis of self-concordant functions. Specifically, by demonstrating that the intrinsic norm of the second-order ascent direction of f at λ_0 is small, we conclude that the solution $\hat{\lambda}$ lies within a small radius of λ_0 in the intrinsic norm. Moreover, using the bound on $f''(\lambda_0, \cdot)$ provided in Theorem 3.3, we establish the multiplicative proximity between λ_0 and $\hat{\lambda}$. In particular, we show that the specific choice $\lambda_0 = (\alpha - 1) / (-\max_{r \in [d]} \{\mathbf{r}[r]\})$ serves as a reasonable analytical estimate, as the *LogSumExp* function $\log\left(\sum_{k=1}^d e^{\lambda \mathbf{r}^{(t)}[k]}\right)$ in the dynamic learning rate control problem 3.2 behaves approximately like a softmax function over the regret vector $\mathbf{r}^{(t)}$. Applying a similar procedure to the regret vector \mathbf{r}' , we establish the multiplicative proximity between λ'_0 and $\hat{\lambda}'$.

Combining these results, we conclude the multiplicative closeness of $\hat{\lambda}$ and $\hat{\lambda}'$ in terms of \mathbf{r}, \mathbf{r}' , and β . This process is formalized in the proof of the Multiplicative Stability Lemma 5.17 in Section 5.3.2. Finally, by considering different cases for the values of $\max_{r \in [d]} \{\mathbf{r}[r]\}$ and $\max_{r \in [d]} \{\mathbf{r}'[r]\}$ relative to $-\beta \log^2 d$, we establish the concrete multiplicative stability of $\hat{\lambda}$ and $\hat{\lambda}'$, in light of Lemma 5.16. Full details are provided in Section 5.3.2.

Since the analytic guess λ_0 can be computed efficiently and guarantees a small norm of the Newton step, the standard analysis of Newton's method for self-concordant functions immediately implies the following.

Corollary 3.6. *Given any $\mathbf{r} \in \mathbb{R}^d$ and a desired relative accuracy $\epsilon > 0$, $O(\log \log 1/\epsilon)$ iterations of Newton's method, starting from the initialization point λ_0 , are sufficient to compute a point λ that approximates $\lambda^* := \arg \min_{\lambda \in (0, \eta]} f(\lambda; \mathbf{r})$ with relative error at most ϵ , meaning that $(1 - \epsilon)\lambda^* < \lambda < (1 + \epsilon)\lambda^*$.*

A byproduct of this analysis is that solving the learning rate control optimization problem 3.2 up to a sufficient accuracy requires $O(d \log \log T)$ computational cost per iteration, which is negligible given that reading a reward vector already has complexity $O(d)$.

4 Extension to 0/1-Polyhedral Games via Kernels

Before delving into the analysis of DLRC-OMWU, we remark that—similarly to the OMWU algorithm—our algorithm can be applied efficiently to certain classes of polyhedral convex games with 0/1-integral vertices.

Consider a convex 0/1-polyhedral set $\Omega \subseteq \mathbb{R}^d$, that is, a polytope whose set of vertices \mathcal{V}_Ω is a subset of $\{0, 1\}^d$. The OMWU algorithm can be directly applied on Ω by keeping track of a distribution $\chi^{(t)}$ over the vertices and updating the distribution multiplicatively depending on the utility $\mathcal{V}[\mathbf{v}] = \langle \boldsymbol{\nu}, \mathbf{v} \rangle$ scored by each vertex $\mathbf{v} \in \mathcal{V}_\Omega$. While this process takes time proportional

Algorithm 2: Kernelized Dynamic Learning Rate Control - Optimistic MWU (KDLRC-OMWU)

Data: Learning rate η , parameters α and β

```

1 Set  $\boldsymbol{\mu}^{(1)}, \boldsymbol{\nu}^{(0)}, \mathbf{x}^{(0)} \leftarrow \mathbf{0} \in \mathbb{R}^d, \sigma^{(1)}, \lambda^{(0)} \leftarrow 0 \in \mathbb{R}$ 
2 for  $t = 1, 2, \dots, T$  do
3   Set  $\boldsymbol{\mu}^{(t)} \leftarrow \boldsymbol{\mu}^{(t)} + \boldsymbol{\nu}^{(t-1)}$  [▷ Optimism for utility]
4   Set  $\sigma^{(t)} \leftarrow \sigma^{(t)} - \langle \boldsymbol{\nu}^{(t-1)}, \mathbf{x}^{(t-1)} \rangle$  [▷ Optimism for correction]
   /* Dynamic Learning Rate Control via Kernelized Newton */
5   Set  $\lambda \leftarrow \lambda^{(t-1)}$  [▷ Warm-start initialization for Newton]
6   repeat
7     for  $r = 1, 2, \dots, d$  do
8       Set  $\mathbf{b}[r] \leftarrow \exp\{\lambda \boldsymbol{\mu}^{(t)}[r]\}$  [▷ See (4.4)]
9       for  $i = 1, 2, \dots, d$  do
10        Set  $\mathbb{E}[\mathbf{v}]_i \leftarrow 1 - \frac{K_\Omega(\mathbf{b}, \bar{\mathbf{e}}_i)}{K_\Omega(\mathbf{b}, \mathbf{1}_d)}$  [▷ See (4.5)]
11        for  $i, j = 1, 2, \dots, d$  do
12          if  $i \neq j$  then
13            Set  $\mathbb{E}[\mathbf{v}\mathbf{v}^\top]_{ij} \leftarrow 1 + \frac{K_\Omega(\mathbf{b}, \bar{\mathbf{e}}_i)}{K_\Omega(\mathbf{b}, \mathbf{1}_d)} + \frac{K_\Omega(\mathbf{b}, \bar{\mathbf{e}}_j)}{K_\Omega(\mathbf{b}, \mathbf{1}_d)} - \frac{K_\Omega(\mathbf{b}, \bar{\mathbf{e}}_{ij})}{K_\Omega(\mathbf{b}, \mathbf{1}_d)}$  [▷ See (4.6)]
14          else
15            Set  $\mathbb{E}[\mathbf{v}\mathbf{v}^\top]_{ii} \leftarrow 1 + \frac{K_\Omega(\mathbf{b}, \bar{\mathbf{e}}_i)}{K_\Omega(\mathbf{b}, \mathbf{1}_d)}$  [▷ See (4.7)]
16          Set  $f'(\lambda; \mathcal{R}) \leftarrow (\boldsymbol{\mu}^{(t)})^\top \mathbb{E}[\mathbf{v}] + \sigma^t + \frac{\alpha - 1}{\lambda}$ 
17          Set  $f''(\lambda; \mathcal{R}) \leftarrow (\boldsymbol{\mu}^{(t)})^\top \mathbb{E}[\mathbf{v}\mathbf{v}^\top] \boldsymbol{\mu}^{(t)} - ((\boldsymbol{\mu}^{(t)})^\top \mathbb{E}[\mathbf{v}])^2 - \frac{\alpha - 1}{\lambda^2}$ 
          /* Newton Update */
18          Set  $\lambda \leftarrow \lambda + \frac{f'(\lambda; \mathcal{R})}{f''(\lambda; \mathcal{R})}$ 
19        until Convergence of  $\lambda$ 
20      Set  $\lambda^{(t)} \leftarrow \lambda$ 
   /* Kernelized OMWU */
21   for  $r = 1, 2, \dots, d$  do
22     Set  $\mathbf{b}^{(t)}[r] \leftarrow \exp\{\lambda^{(t)} \boldsymbol{\mu}^{(t)}[r]\}$  [▷ See (4.2)]
23   for  $r = 1, 2, \dots, d$  do
24     Set  $\mathbf{x}^{(t)}[r] \leftarrow 1 - \frac{K_\Omega(\mathbf{b}^{(t)}, \bar{\mathbf{e}}_r)}{K_\Omega(\mathbf{b}^{(t)}, \mathbf{1}_d)}$  [▷ See (4.3)]
25   Play strategy  $\mathbf{x}^{(t)}$  and observe  $\boldsymbol{\nu}^{(t)} \in \mathbb{R}^d$ 
26   Set  $\boldsymbol{\mu}^{(t)} \leftarrow \boldsymbol{\mu}^{(t)} + \boldsymbol{\nu}^{(t)} - \boldsymbol{\nu}^{(t-1)}$  [▷ Empirical cumulated corrections]
27   Set  $\sigma^{(t)} \leftarrow \sigma^{(t)} - \langle \boldsymbol{\nu}^{(t)}, \mathbf{x}^{(t)} \rangle + \langle \boldsymbol{\nu}^{(t-1)}, \mathbf{x}^{(t-1)} \rangle$  [▷ Empirical cumulated utilities]

```

to the number of vertices $|\mathcal{V}_\Omega|$ if implemented naively, Farina, Lee, Luo, and Kroer [Far+22b] prove that the process can sometimes be simulated in polynomial time per iteration even when $|\mathcal{V}_\Omega|$ is large. Specifically, they show that the update of $\boldsymbol{\chi}^{(t)}$ and computation of the expectation $\sum_{\boldsymbol{v} \in \mathcal{V}_\Omega} \boldsymbol{\chi}^{(t)}[\boldsymbol{v}] \boldsymbol{v}$ can be carried out using only $d + 1$ evaluations of a *0/1-polyhedral kernel*, which they demonstrate can be evaluated efficiently in extensive-form games and various other convex 0/1-polyhedral settings, including m -sets, unit cubes, and flows on directed acyclic graphs by building on a prior idea of Takimoto and Warmuth [TW03].

Although the dynamics of DLRC-OMWU are quite similar to those of OMWU, extending DLRC-OMWU to a kernelized version (KDLRC-OMWU) is not an immediate consequence of Farina, Lee, Luo, and Kroer [Far+22b]. This extension requires additional considerations, particularly due to the need to solve the dynamic learning rate control problem in Equation (3.2) at each time step t . We will show that this optimization problem can also be addressed using the kernel trick, with novel modifications. Let us begin with definitions of the 0/1-polyhedral feature mapping, the associated kernel, and the key results from Farina, Lee, Luo, and Kroer [Far+22b].

Definition 1 (0/1-polyhedral feature map and kernel [Far+22b]). Associated with a convex 0/1-polyhedral set $\Omega \subseteq \mathbb{R}^d$, define the *0/1-polyhedral feature map* $\phi_\Omega : \mathbb{R}^d \rightarrow \mathbb{R}^{\mathcal{V}_\Omega}$,

$$\phi_\Omega(\boldsymbol{x})[\boldsymbol{v}] := \prod_{k: \boldsymbol{v}[k]=1} \boldsymbol{x}[k] \quad \forall \boldsymbol{x} \in \mathbb{R}^d, \boldsymbol{v} \in \mathcal{V}_\Omega,$$

and the corresponding *0/1-polyhedral kernel* $K_\Omega : \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$,

$$K_\Omega(\boldsymbol{x}_1, \boldsymbol{x}_2) := \langle \phi_\Omega(\boldsymbol{x}_1), \phi_\Omega(\boldsymbol{x}_2) \rangle = \sum_{\boldsymbol{v} \in \mathcal{V}_\Omega} \prod_{k: \boldsymbol{v}[k]=1} \boldsymbol{x}_1[k] \boldsymbol{x}_2[k] \quad \forall \boldsymbol{x}_1, \boldsymbol{x}_2 \in \mathbb{R}^d.$$

We also define auxiliary indicator vectors $\bar{\boldsymbol{e}}_i, \bar{\boldsymbol{e}}_{ij} \in \mathbb{R}^d$ for all $i, j \in [d]$, such that

$$\bar{\boldsymbol{e}}_i[k] := \mathbb{I}_{k \neq i} = \begin{cases} 0 & \text{if } k = i, \\ 1 & \text{if } k \neq i, \end{cases} \quad \text{and} \quad \bar{\boldsymbol{e}}_{ij}[k] := \mathbb{I}_{k \neq i \wedge k \neq j} = \begin{cases} 0 & \text{if } k = i \text{ or } k = j, \\ 1 & \text{if } k \neq i \text{ and } k \neq j. \end{cases}$$

The embedding of these vectors is useful in kernel computations.

The following proposition ensures that Step (3.1),

$$\boldsymbol{\chi}^{(t)}[\boldsymbol{v}] := \frac{\exp\{\lambda^{(t)} \mathcal{R}^{(t)}[\boldsymbol{v}]\}}{\sum_{\boldsymbol{v}' \in \mathcal{V}_\Omega} \exp\{\lambda^{(t)} \mathcal{R}^{(t)}[\boldsymbol{v}']\}} \quad \forall \boldsymbol{v} \in \mathcal{V}_\Omega, \quad (4.1)$$

can be simulated using only $d + 1$ kernel evaluations, assuming that the dynamic learning rate $\lambda^{(t)}$ is available. The key idea is that by embedding a carefully constructed vector $\boldsymbol{b}^{(t)}$ into the feature mapping ϕ_Ω , the algorithm's updates—represented as the distribution over the vertices $\boldsymbol{\chi}^{(t)} \in \Delta(\mathcal{V}_\Omega)$ —and thus the actions $\boldsymbol{x}^{(t)} \in \Omega$, become computable via the kernel K_Ω .

Proposition 4.2 (Theorem 4.1 and 4.2 of [Far+22b]). *For all time steps $t \in T$, let $\boldsymbol{\mu}^{(t)} := \boldsymbol{\nu}^{(t)} + \sum_{\tau=1}^t \boldsymbol{\nu}^{(\tau)}$ be the optimistic sum of the utility vectors $\boldsymbol{\nu}^{(t)}$, and define the vector $\boldsymbol{b}^{(t)} \in \mathbb{R}^d$ as*

$$\boldsymbol{b}^{(t)}[k] := \exp\{\lambda^{(t)} \boldsymbol{\mu}^{(t)}[k]\} \quad \forall k \in [d]. \quad (4.2)$$

Then, the distributions $\boldsymbol{\chi}^{(t)}$ are proportional to $\phi_{\Omega}(\mathbf{b}^{(t)})$,

$$\boldsymbol{\chi}^{(t)} = \frac{\phi_{\Omega}(\mathbf{b}^{(t)})}{K_{\Omega}(\mathbf{b}^{(t)}, \mathbf{1}_d)},$$

and the iterates $\mathbf{x}^{(t)}$ produced by DLRC-OMWU are computed as

$$\mathbf{x}^{(t)} = \sum_{\mathbf{v} \in \mathcal{V}_{\Omega}} \boldsymbol{\chi}^{(t)}[\mathbf{v}] \mathbf{v} = \left[1 - \frac{K_{\Omega}(\mathbf{b}^{(t)}, \bar{\mathbf{e}}_1)}{K_{\Omega}(\mathbf{b}^{(t)}, \mathbf{1}_d)}, 1 - \frac{K_{\Omega}(\mathbf{b}^{(t)}, \bar{\mathbf{e}}_2)}{K_{\Omega}(\mathbf{b}^{(t)}, \mathbf{1}_d)}, \dots, 1 - \frac{K_{\Omega}(\mathbf{b}^{(t)}, \bar{\mathbf{e}}_d)}{K_{\Omega}(\mathbf{b}^{(t)}, \mathbf{1}_d)} \right]. \quad (4.3)$$

It remains to show that, for each time step t , the dynamic learning rate $\lambda^{(t)}$ can be determined via kernelization. By Corollary 3.6, we need to verify that the Newton steps of the optimization problem in Equation (3.2) can be simulated by the kernel K_{Ω} . The Newton algorithm requires the calculation of $f'(\lambda; \mathcal{R})$ and $f''(\lambda; \mathcal{R})$,³

$$f'(\lambda; \mathcal{R}) = \frac{\sum_{\mathbf{v} \in \mathcal{V}_{\Omega}} \mathcal{R}[\mathbf{v}] e^{\lambda \mathcal{R}[\mathbf{v}]}}{\sum_{\mathbf{v} \in \mathcal{V}_{\Omega}} e^{\lambda \mathcal{R}[\mathbf{v}]}} + \frac{\alpha - 1}{\lambda}$$

$$f''(\lambda; \mathcal{R}) = \frac{\sum_{\mathbf{v} \in \mathcal{V}_{\Omega}} \mathcal{R}[\mathbf{v}]^2 e^{\lambda \mathcal{R}[\mathbf{v}]}}{\sum_{\mathbf{v} \in \mathcal{V}_{\Omega}} e^{\lambda \mathcal{R}[\mathbf{v}]}} - \left(\frac{\sum_{\mathbf{v} \in \mathcal{V}_{\Omega}} \mathcal{R}[\mathbf{v}] e^{\lambda \mathcal{R}[\mathbf{v}]}}{\sum_{\mathbf{v} \in \mathcal{V}_{\Omega}} e^{\lambda \mathcal{R}[\mathbf{v}]}} \right)^2 - \frac{\alpha - 1}{\lambda^2},$$

where the vector \mathcal{R} is potentially of exponential size. We recall that by Equation (4.1), $\boldsymbol{\chi}$ can be seen as a discrete random variable. We revisit the representation of $f'(\lambda; \mathcal{R})$ and $f''(\lambda; \mathcal{R})$. Interestingly, they can be rewritten as

$$f'(\lambda; \mathcal{R}) = \mathbb{E}[\mathcal{R}[\mathbf{v}]] + \frac{\alpha - 1}{\lambda}$$

$$f''(\lambda; \mathcal{R}) = \mathbb{E}[\mathcal{R}[\mathbf{v}]^2] - \mathbb{E}[\mathcal{R}[\mathbf{v}]]^2 - \frac{\alpha - 1}{\lambda^2},$$

where the expectations are taken with respect to the distribution $\mathbf{v} \sim \boldsymbol{\chi}$. Consequently, it suffices to verify that the first and second moments of \mathcal{R} with respect to the distribution $\boldsymbol{\chi}$ can be computed via the kernelization approach. Given that $\mathcal{R}[\mathbf{v}] = \langle \boldsymbol{\mu}, \mathbf{v} \rangle + \sigma$, where

$$\boldsymbol{\mu} := \boldsymbol{\nu}^{(t)} + \sum_{\tau=1}^t \boldsymbol{\nu}^{(\tau)} \quad \text{and} \quad \sigma := -(\langle \boldsymbol{\nu}^{(t)}, \mathbf{x}^{(t)} \rangle + \sum_{\tau=1}^t \langle \boldsymbol{\nu}^{(\tau)}, \mathbf{x}^{(\tau)} \rangle),$$

we infer that

$$\mathbb{E}[\mathcal{R}[\mathbf{v}]] = \mathbb{E}[\langle \boldsymbol{\mu}, \mathbf{v} \rangle] + \sigma = \langle \boldsymbol{\mu}, \mathbb{E}[\mathbf{v}] \rangle + \sigma.$$

And,

$$\mathbb{E}[\mathcal{R}[\mathbf{v}]^2] = \mathbb{E}[(\langle \boldsymbol{\mu}, \mathbf{v} \rangle + \sigma)^2] = \mathbb{E}[\langle \boldsymbol{\mu}, \mathbf{v} \rangle^2] + 2\sigma \mathbb{E}[\langle \boldsymbol{\mu}, \mathbf{v} \rangle] + \sigma^2 = \boldsymbol{\mu}^{\top} \mathbb{E}[\mathbf{v} \mathbf{v}^{\top}] \boldsymbol{\mu} + 2\sigma \langle \boldsymbol{\mu}, \mathbb{E}[\mathbf{v}] \rangle + \sigma^2,$$

³To simplify notation, we omit the superscript t from this point onward in this section whenever it is clear from the context.

by multilinearity. We can simplify the $f''(\lambda; \mathcal{R})$ term further,

$$\begin{aligned} f''(\lambda; \mathcal{R}) &= \boldsymbol{\mu}^\top \mathbb{E}[\mathbf{v}\mathbf{v}^\top] \boldsymbol{\mu} + 2\sigma \boldsymbol{\mu}^\top \mathbb{E}[\mathbf{v}] + \sigma^2 - (\boldsymbol{\mu}^\top \mathbb{E}[\mathbf{v}] + \sigma)^2 - \frac{\alpha - 1}{\lambda^2} \\ &= \boldsymbol{\mu}^\top \mathbb{E}[\mathbf{v}\mathbf{v}^\top] \boldsymbol{\mu} - (\boldsymbol{\mu}^\top \mathbb{E}[\mathbf{v}])^2 - \frac{\alpha - 1}{\lambda^2} \end{aligned}$$

Hence, it is adequate to calculate $\mathbb{E}[\mathbf{v}]$ and $\mathbb{E}[\mathbf{v}\mathbf{v}^\top]$, which we formalize in the following Proposition.

Proposition 4.3. *Define the vector $\mathbf{b} \in \mathbb{R}^d$ as*

$$\mathbf{b}[k] := \exp\{\lambda \boldsymbol{\mu}[k]\} \quad \forall k \in [d]. \quad (4.4)$$

Then,

$$\mathbb{E}[\mathbf{v}] = \left[1 - \frac{K_\Omega(\mathbf{b}, \bar{\mathbf{e}}_1)}{K_\Omega(\mathbf{b}, \mathbf{1}_d)}, 1 - \frac{K_\Omega(\mathbf{b}, \bar{\mathbf{e}}_2)}{K_\Omega(\mathbf{b}, \mathbf{1}_d)}, \dots, 1 - \frac{K_\Omega(\mathbf{b}, \bar{\mathbf{e}}_d)}{K_\Omega(\mathbf{b}, \mathbf{1}_d)} \right], \quad (4.5)$$

And,

$$\mathbb{E}[\mathbf{v}\mathbf{v}^\top]_{ij} = 1 + \frac{K_\Omega(\mathbf{b}, \bar{\mathbf{e}}_i)}{K_\Omega(\mathbf{b}, \mathbf{1}_d)} + \frac{K_\Omega(\mathbf{b}, \bar{\mathbf{e}}_j)}{K_\Omega(\mathbf{b}, \mathbf{1}_d)} - \frac{K_\Omega(\mathbf{b}, \bar{\mathbf{e}}_{ij})}{K_\Omega(\mathbf{b}, \mathbf{1}_d)}, \quad (4.6)$$

for all $i, j \in [d]$, where $i \neq j$, and

$$\mathbb{E}[\mathbf{v}\mathbf{v}^\top]_{ii} = 1 + \frac{K_\Omega(\mathbf{b}, \bar{\mathbf{e}}_i)}{K_\Omega(\mathbf{b}, \mathbf{1}_d)}, \quad (4.7)$$

for all $i \in [d]$.

Proof. We start by $\mathbb{E}[\mathbf{v}]$.

$$\mathbb{E}[\mathbf{v}] = \sum_{\mathbf{v} \in \mathcal{V}_\Omega} \boldsymbol{\chi}[\mathbf{v}] \mathbf{v} = \mathbf{x},$$

the rest follows similar to Proposition 4.2. Next, we analyze $\mathbb{E}[\mathbf{v}\mathbf{v}^\top]$. First we show that, for every $i, j \in [d]$,

$$\phi_\Omega(\bar{\mathbf{e}}_i)[\mathbf{v}] = \prod_{k:\mathbf{v}[k]=1} \bar{\mathbf{e}}_i[k] = \prod_{k:\mathbf{v}[k]=1} \mathbb{I}_{k \neq i} = \mathbb{I}_{k \notin \mathbf{v}}$$

and by the fact that $\phi_\Omega(\mathbf{1}_d) = \mathbf{1}$,

$$\phi_\Omega(\mathbf{1}_d)[\mathbf{v}] - \phi_\Omega(\bar{\mathbf{e}}_i)[\mathbf{v}] = \mathbb{I}_{i \in \mathbf{v}}.$$

Similarly,

$$\phi_\Omega(\bar{\mathbf{e}}_{ij})[\mathbf{v}] = \prod_{k:\mathbf{v}[k]=1} \bar{\mathbf{e}}_{ij}[k] = \prod_{k:\mathbf{v}[k]=1} \mathbb{I}_{k \neq i \text{ and } k \neq j} = \mathbb{I}_{i, j \notin \mathbf{v}}$$

For every $i, j \in [d]$,

$$\begin{aligned}
\mathbb{E}[\mathbf{v}_i \mathbf{v}_j] &= \sum_{\mathbf{v} \in \mathcal{V}_\Omega} \chi[\mathbf{v}] \mathbf{v}_i \mathbf{v}_j \\
&= \sum_{\mathbf{v} \in \mathcal{V}_\Omega} \chi[\mathbf{v}] \cdot \mathbb{I}_{i, j \in \mathbf{v}} \\
&= \sum_{\mathbf{v} \in \mathcal{V}_\Omega} \chi[\mathbf{v}] \cdot (1 - \mathbb{I}_{i \notin \mathbf{y}} - \mathbb{I}_{j \notin \mathbf{y}} + \mathbb{I}_{i, j \notin \mathbf{y}}) \\
&= 1 - \sum_{\mathbf{v} \in \mathcal{V}_\Omega} \chi[\mathbf{v}] \cdot (-\phi_\Omega(\bar{\mathbf{e}}_i)[\mathbf{v}] - \phi_\Omega(\bar{\mathbf{e}}_j)[\mathbf{v}] + \phi_\Omega(\bar{\mathbf{e}}_{ij})[\mathbf{v}]) \\
&= 1 + \frac{K_\Omega(\mathbf{b}, \bar{\mathbf{e}}_i)}{K_\Omega(\mathbf{b}, \mathbf{1}_d)} + \frac{K_\Omega(\mathbf{b}, \bar{\mathbf{e}}_j)}{K_\Omega(\mathbf{b}, \mathbf{1}_d)} - \frac{K_\Omega(\mathbf{b}, \bar{\mathbf{e}}_{ij})}{K_\Omega(\mathbf{b}, \mathbf{1}_d)}
\end{aligned} \tag{4.8}$$

where line (4.8) follows from the inclusion–exclusion principle, and in the last line, we applied Proposition 4.2 multiple times. The case with $i = j$ is quite similar. \square

An immediate byproduct of Proposition 4.3 is that $d^2 + 1$ evaluations of Kernel K_Ω are sufficient for each iteration of the Newton optimization algorithm for Problem (3.2).

Corollary 4.4. *KDLRC-OMWU requires $(d + 1) + (d^2 + 1) O(\log \log T)$ kernel evaluations at each time step t , where the first $d + 1$ evaluations are used for Step (4.1), and the remaining $(d^2 + 1) O(\log \log T)$ evaluations are for computing the dynamic learning rate $\lambda^{(t)}$ via the Newton algorithm. KDLRC-OMWU achieves $O(n \log^2 |\mathcal{V}_\Omega| \log T)$ regret in the self-play setting, and $\tilde{O}(\sqrt{T \log |\mathcal{V}_\Omega|})$ regret in adversarial settings.*

We conclude this section with the final piece of the puzzle for Algorithm 2 (KDLRC-OMWU). Since computing the initial guess $\lambda_0 = (\alpha - 1) / (-\max_{\mathbf{v} \in \mathcal{V}_\Omega} \{\mathcal{R}^{(t)}[\mathbf{v}]\})$ —the initialization point for the Newton algorithm—is not necessarily efficient, we warm-start the Newton algorithm at each iteration t by using the previous solution $\lambda_0 = \lambda^{(t-1)}$ from time $t - 1$. It is straightforward to show that, due to the multiplicative stability property of the learning rates (Theorem 3.5), setting $\lambda_0 = \lambda^{(t-1)}$ provides a sufficiently accurate analytical estimate for initializing the Newton algorithm. We formalize this observation below.

Observation 4.5. Given the dynamic learning rate control optimization problem Equation (3.2) for the regret vector $\mathcal{R}^{(t)}$ of the game over $\Delta(\mathcal{V}_\Omega)$,

$$\lambda^{(t)} = \arg \max_{\lambda \in (0, \eta]} \left\{ f(\lambda; \mathcal{R}^{(t)}) := (\alpha - 1) \log \lambda + \log \left(\sum_{\mathbf{v} \in \mathcal{V}_\Omega} e^{\lambda \mathcal{R}^{(t)}[\mathbf{v}]} \right) \right\},$$

Newton’s method, warm-started from $\lambda_0 = \lambda^{(t-1)}$, converges in $O(\log \log 1/\epsilon)$ iterations to the optimal solution $\lambda^{(t)}$ with a relative error of at most ϵ .

Proof. By the analysis of Newton’s method for self-concordant functions [Ren01], it is sufficient to show that the size of the Newton step at initialization $\lambda_0 = \lambda^{(t-1)}$ is small in the local norm, i.e.,

$$\|n(\lambda_0)\|_{f''(\lambda_0; \mathcal{R}^{(t)})} = \frac{f'(\lambda_0; \mathcal{R}^{(t)})^2}{|f''(\lambda_0; \mathcal{R}^{(t)})|} \leq 1,$$

where $n(\lambda_0)$ denotes the Newton step. By the multiplicative stability of $\lambda^{(t)}$ and our choice of initialization $\lambda_0 = \lambda^{(t-1)}$ —as implied by the sensitivity of the learning rate to the regret vector in Theorem 3.5, and following similar reasoning to the proof of Theorem 3.5—this bound on the Newton step size can be readily established. We omit the details here in the interest of space. \square

5 Analysis of the Dynamic Learning Rate Control

In this section, we present our results and analysis of the regret for DLRC-OMWU. A cornerstone of this analysis is the following regret bound (Theorem 5.1), which follows the style of the *nonnegative RVU property*, originally introduced by Farina, Anagnostides, Luo, Lee, Kroer, and Sandholm [Far+22c]. This property differs from the original RVU property discussed in Section 2, and is stronger in the sense that the nonnegative RVU property directly implies the RVU property.

Theorem 5.1 (Nonnegative RVU bound of DLRC-OMWU). *Consider the cumulative regret $\tilde{\text{Reg}}^{(T)}$ accrued by DLRC-OMWU algorithm up to time T . Assuming that $\|\boldsymbol{\nu}^{(t)}\|_\infty \leq 1$ is satisfied for all $t \in [T]$, it follows that for any time $T \in \mathbb{N}$ and any learning rate $\eta \leq \frac{1}{50}$ and β high enough ($\beta \geq 70$),*

$$\max\{0, \text{Reg}^{(T)}\} \leq 3 + \frac{\alpha \log T + \log d}{\eta} + 6\eta \sum_{t=1}^{T-1} \|\boldsymbol{\nu}^{(t)} - \boldsymbol{\nu}^{(t-1)}\|_\infty^2 - \frac{1}{24\eta} \sum_{t=1}^{T-1} \|\mathbf{x}^{(t+1)} - \mathbf{x}^{(t)}\|_1^2.$$

The proof sketch of our nonnegative RVU bound is provided in Section 5.2. The use of positive regret, $\max\{0, \text{Reg}^{(T)}\}$, is pivotal in our analysis, as any upper bound on the sum $\sum_{i=1}^n \max\{0, \text{Reg}_i^{(T)}\}$ directly implies the same upper bound on the maximum regret, $\max_{i \in [n]} \text{Reg}_i^{(T)}$, among the n players due to nonnegativity. This, in turn, implies that the empirical distribution of joint strategies converges to an approximate CCE of the game.

With Theorem 5.1 at hand, the path forward becomes straightforward. The general plan is to use the nonnegativity of the sum $\sum_{i=1}^n \max\{0, \text{Reg}_i^{(T)}\}$ and to set a sufficiently small learning rate η in order to infer that the total path length of the play, i.e., $\sum_{i=1}^n \sum_{t=1}^{T-1} \|\mathbf{x}_i^{(t+1)} - \mathbf{x}_i^{(t)}\|_1^2$, is bounded by $O(n \log^2 d \log T)$. This procedure is formalized in Theorem 5.2.

Theorem 5.2 (Bound on total path length). *Under Assumption 1, if all the players follow DLRC-OMWU algorithm with learning rate $\eta \leq \min\{\frac{1}{50}, \frac{1}{12\sqrt{2}Ln}\}$, then*

$$\sum_{i=1}^n \sum_{t=1}^{T-1} \|\mathbf{x}_i^{(t+1)} - \mathbf{x}_i^{(t)}\|_1^2 \leq 144n\eta + 48n(\alpha \log T + \log d).$$

Given the RVU bound of Theorem 5.1 and the bound on the total path length from Theorem 5.2, we can adhere to the standard machinery of RVU bounds [Syr+15] and prove that the regret $\text{Reg}_i^{(T)}$ of each player $i \in [n]$ is bounded by $O(n \log^2 d \log T)$, as stated in Theorem 5.3. The idea is that, under Assumption 1, the variation in the utilities observed over time, $\sum_{t=1}^{T-1} \|\boldsymbol{\nu}^{(t)} - \boldsymbol{\nu}^{(t-1)}\|_\infty^2$, can be bounded by the total path length, thereby yielding the result.

Theorem 5.3 (Regret bound of DLRC-OMWU (Formal version of Theorem 3.1)). *Under Assumption 1, if all the players $i \in [n]$ follow DLRC-OMWU with a learning rate $\eta = \min\{\frac{1}{50}, \frac{1}{12\sqrt{2}Ln}\}$, then the regret $\text{Reg}_i^{(T)}$ of each player $i \in [n]$ is bounded by*

$$\text{Reg}_i^{(T)} \leq 6 + \max\{50 + 12\sqrt{2}Ln, 24\sqrt{2}Ln\} (\alpha \log T + \log d) = O(n \log^2 d \log T).$$

Additionally, the DLRC-OMWU algorithm for each player $i \in [n]$ is adaptive to adversarial utilities, meaning that the regret incurred in the face of adversarial utilities is $\text{Reg}_i^{(T)} = \tilde{O}(\sqrt{T \log d})$.

For the detailed proofs, please refer to Section 5.3.7. We conclude this section by noting that our regret guarantees in Theorem 5.3 hold *uniformly for all* $T \in \mathbb{N}^+$ *simultaneously*, since the algorithm DLRC-OMWU and the choice of learning rate η do not depend on the time horizon T . As a result, even when the horizon of self-play T is unknown, or miscoordination occurs among the participants, DLRC-OMWU still enjoys low regret guarantees in both self-play and adversarial settings—without the need for additional techniques such as the doubling trick.

This is in contrast to the analysis of Daskalakis, Fishelson, and Golowich [DFG21]—and indeed our *anytime convergence* provides a strictly stronger guarantee—where the choice of learning rate η depends on the time horizon T . Consequently, when T is unknown, additional techniques such as the doubling trick appear to be unavoidable. Furthermore, since the learning rate η in Daskalakis, Fishelson, and Golowich [DFG21] is restricted to the interval $1/T \leq \eta \leq 1/(Cn \log^4 T)$ for some constant C (see Lemmas 4.2 and C.4 in [DFG21]), their guarantees fail to apply in the regime where $T \leq Cn \log^4 T$, which frequently arises in games with many players or in short-horizon scenarios.

5.1 Design of the DLRC-OMWU Algorithm and Equivalent Viewpoints

Before delving into the proof sketch in Section 5.2, we first discuss how to mathematically formalize the concept of dynamic learning rate control, leading to Algorithm 1 (DLRC-OMWU). In this context, we also present alternative perspectives on DLRC-OMWU and highlight how they contribute to both its regret analysis and computational properties.

We begin with the standard dynamics of Optimistic Follow-the-Regularized-Leader (OFTRL) algorithms with a potentially dynamic learning rate $\lambda^{(t)}$,

$$\mathbf{x}^{(t)} \leftarrow \arg \max_{\mathbf{x} \in \mathcal{X}} \left\{ \lambda^{(t)} \langle \mathbf{r}^{(t)}, \mathbf{x} \rangle + \phi(\mathbf{x}) \right\}, \quad (5.1)$$

where $\lambda^{(t)} \in (0, \eta]$ is chosen according to a separate dynamic that we will design later, and ϕ is a regularizer over the space \mathcal{X} . In the case of the negative entropy regularizer, $\phi(\mathbf{x}) = \sum_{k=1}^d \mathbf{x}[k] \log \mathbf{x}[k]$, Formulation (5.1) recovers the celebrated Optimistic MWU algorithm with learning rate $\lambda^{(t)}$.

The fundamental idea here is to also integrate the dynamic change of $\lambda^{(t)}$ into Equation (5.1). Naturally, in online learning problems, we aim to incentivize selecting actions with the best rewards. On the other hand, for the purpose of self-play, as discussed comprehensively in Section 3, we seek to additionally *pace down* the learner when it is *performing too well*. Thus, a seamless extension of Optimization Step 5.1, equipped with an automatic dynamic adjustment of $\lambda^{(t)}$, takes the form:

$$\begin{pmatrix} \lambda^{(t)} \\ \mathbf{x}^{(t)} \end{pmatrix} \leftarrow \arg \max_{\lambda \in (0, \eta], \mathbf{x} \in \mathcal{X}} \left\{ \lambda \langle \mathbf{r}^{(t)}, \mathbf{x} \rangle + \phi(\mathbf{x}) \right\}. \quad (5.2)$$

However, in this formulation, $\lambda^{(t)}$ exhibits behavior akin to that of a step function. When $\langle \mathbf{r}^{(t)}, \mathbf{x}^{(t)} \rangle > 0$, this formulation naively reduces to $\lambda^{(t)} = \eta$, which corresponds to the original OFTRL with a constant learning rate; otherwise, it trivially sets $\lambda^{(t)} = 0$ and $\mathbf{x}^{(t)} = \arg \max_{\mathbf{x} \in \mathcal{X}} \phi(\mathbf{x})$. Instead, to improve the predictability of the behavior of the dynamics $\mathbf{x}^{(t)}$ in the self-play setting, we need a smoother transition for the dynamic learning rate $\lambda^{(t)}$. To achieve this, we can incorporate an additional regularizer $\rho(\lambda)$ for $\lambda \in (0, \eta]$ that defines a smooth thresholding procedure for what

it means to “*perform too well*” relative to the regret vector $\mathbf{r}^{(t)}$ into the dynamics of Equation (5.2), which leads to

$$\begin{pmatrix} \lambda^{(t)} \\ \mathbf{x}^{(t)} \end{pmatrix} \leftarrow \arg \max_{\lambda \in (0, \eta], \mathbf{x} \in \mathcal{X}} \left\{ \lambda \langle \mathbf{r}^{(t)}, \mathbf{x} \rangle + \rho(\lambda) + \phi(\mathbf{x}) \right\}. \quad (5.3)$$

In other words, the term $\lambda \langle \mathbf{r}^{(t)}, \mathbf{x} \rangle$ is motivating *regret minimization* for the player, while the regularizer $\rho(\lambda)$ is prohibiting the player from having a *very low regret*, thereby creating a *balance among the performance of the players* of the game.

We note that the dynamics of Formulation (5.3) are not jointly concave in (λ, \mathbf{x}) for a general choice of ρ and ϕ ; therefore, its *computational* aspects and *regret analysis* remain unclear. Notably, for the special choice of $\rho(\lambda) := (\alpha - 1) \log \lambda$, $\phi(\mathbf{x}) = \sum_{k=1}^d \mathbf{x}[k] \log \mathbf{x}[k]$, and $\mathcal{X} = \Delta^d$, we will demonstrate that Formulation (5.3) is equivalent to DLRC-OMWU. Hence, it can be computed at each iteration by solving the 1-dimensional optimization problem of dynamic learning rate control (3.2) and the play of OMWU, as outlined in Algorithm 1 and Section 3.1.

To analyze the regret of DLRC-OMWU, we show that Algorithm 1 can be reformulated as a specific instance of OFTRL, beyond its original form in Formulation (5.3). A key technical step in this analysis is to express the iterates $\mathbf{x}^{(t)}$ produced by the algorithm through the OFTRL perspective. By defining $\mathbf{y}^{(t)} := \lambda^{(t)} \mathbf{x}^{(t)} \in (0, 1] \Delta^d = \left\{ \mathbf{y} \in \mathbb{R}_{\geq 0}^d \mid \sum_{k=1}^d \mathbf{y}[k] \leq 1 \right\}$, we will demonstrate that the iterates $\mathbf{y}^{(t)}$ satisfy the OFTRL rule,

$$\mathbf{y}^{(t)} := \arg \max_{\mathbf{y} \in (0, 1] \Delta^d} \left\{ \langle \mathbf{r}^{(t)}, \mathbf{y} \rangle + \psi(\mathbf{y}) \right\}, \quad (5.4)$$

where ψ is a special regularizer with strong spectral properties (to be discussed in Section 5.2),

$$\psi(\mathbf{y}) = -\frac{1}{\eta} \alpha \log \left(\sum_{k=1}^d \mathbf{y}[k] \right) + \frac{1}{\eta} \frac{1}{\sum_{k=1}^d \mathbf{y}[k]} \sum_{k=1}^d \mathbf{y}[k] \log \mathbf{y}[k]. \quad (5.5)$$

Interestingly, this formulation looks similar to the lifting idea of Farina, Anagnostides, Luo, Lee, Kroer, and Sandholm [Far+22c], where instead of applying OFTRL on the original space $\mathcal{X} = \Delta^n$, OFTRL is designed over the lifted space $\tilde{\mathcal{X}} = \{(\gamma, \mathbf{y}) : \gamma \in [0, 1], \mathbf{y} \in \gamma \mathcal{X}\} = [0, 1] \Delta^n$, but without the need to incorporate the lifting parameter γ into the dynamics of the algorithm (neither in the utility term nor in the regularizer). We note that despite this similar appearance, as will be seen in the next sections, the analysis of DLRC-OMWU is substantially different from that of Log-Regularized Lifted OFTRL (LRL-OFTRL) [Far+22c], where the log regularizer is chosen over the lifted space $\tilde{\mathcal{X}}$, i.e., $\psi(\gamma, \mathbf{y}) = -\frac{1}{\eta} \log \gamma - \frac{1}{\eta} \sum_{k=1}^d \log \mathbf{y}[k]$, to ensure elementwise multiplicative stability of actions due to the high curvature of the log regularizers and the structure of the induced intrinsic norms.

The proof sketch of the regret analysis using Formulation (5.4) is discussed in Section 5.2. For the moment, we formalize these three alternative formulations of DLRC-OMWU in the following theorem and defer the equivalence proofs to Section 5.3.3.

Theorem 5.4. *Algorithm 1 (DLRC-OMWU) can alternatively be viewed as DLRC-OFTRL (Formulation 5.3) with the choice of regularizers $\rho(\lambda) = (\alpha - 1) \log \lambda$ and $\phi(\mathbf{x}) = \sum_{k=1}^d \mathbf{x}[k] \log \mathbf{x}[k]$, i.e.,*

$$\begin{pmatrix} \lambda^{(t)} \\ \mathbf{x}^{(t)} \end{pmatrix} \leftarrow \arg \max_{\lambda \in (0, \eta], \mathbf{x} \in \Delta^d} \left\{ \lambda \langle \mathbf{r}^{(t)}, \mathbf{x} \rangle + (\alpha - 1) \log \lambda - \sum_{k=1}^d \mathbf{x}[k] \log \mathbf{x}[k] \right\}. \quad (5.6)$$

or additionally as OFTRL with the regularizer ψ defined in Equation (5.5), i.e.,

$$\mathbf{y}^{(t)} \leftarrow \arg \max_{\mathbf{y} \in (0,1]^{\Delta^d}} \left\{ \eta \langle \mathbf{r}^{(t)}, \mathbf{y} \rangle + \alpha \log \left(\sum_{k=1}^d \mathbf{y}[k] \right) - \frac{1}{\sum_{k=1}^d \mathbf{y}[k]} \sum_{k=1}^d \mathbf{y}[k] \log \mathbf{y}[k] \right\},$$

In other words, the three perspectives (Algorithm 1, Formulations 5.6 and 5.4) of DLRC-OMWU are equivalent and result in the same learning dynamics.

5.2 Proof Sketch

As discussed in the previous section, an essential step in the regret analysis of DLRC-OMWU is that, by Theorem 5.4, the play $\mathbf{x}^{(t)}$ generated by DLRC-OMWU is equivalent to $\mathbf{y}^{(t)}$ with a transformation by the dynamic learning rate $\lambda^{(t)}$, where the iterates of $\mathbf{y}^{(t)}$ follow

$$\mathbf{y}^{(t)} := \arg \max_{\mathbf{y} \in (0,1]^{\Delta^d}} \left\{ \langle \mathbf{r}^{(t)}, \mathbf{y} \rangle + \psi(\mathbf{y}) \right\},$$

with ψ being our regularizer $\psi(\mathbf{y}) : \Omega \rightarrow \mathbb{R}$ on the domain

$$\Omega := [0, 1]^{\Delta^d} = \left\{ \mathbf{y} \in \mathbb{R}_+^d \mid \sum_{k=1}^d \mathbf{y}[k] \leq 1 \right\}$$

defined as

$$\psi(\mathbf{y}) = -\frac{1}{\eta} \alpha \log \left(\sum_{k=1}^d \mathbf{y}[k] \right) + \frac{1}{\eta} \frac{1}{\sum_{k=1}^d \mathbf{y}[k]} \sum_{k=1}^d \mathbf{y}[k] \log \mathbf{y}[k],$$

where $\alpha := 2 + 2 \log d + \beta \log^2 d$ with $\beta \geq 70$ is a hyperparameter. We restate that the dynamic learning rate parameter $\lambda^{(t)}$ is the unique solution to the Dynamic Learning Rate Control Problem 3.2 for each time $t \in [T]$.

Notation. For the analysis of the algorithm, we introduce additional necessary notation for convenience. We use $\|\mathbf{z}\|_{\mathbf{B}}^2 := \mathbf{z}^\top \mathbf{B} \mathbf{z}$ to denote the quadratic norm induced by the matrix $\mathbf{B} \in \mathbb{R}^{d \times d}$. Let $\Lambda(\mathbf{y}) = \sum_{k=1}^d \mathbf{y}[k]$ be the sum of the elements of the vector $\mathbf{y} \in \mathbb{R}^d$. For two vectors $\mathbf{y}, \mathbf{z} \in (0, 1]^{\Delta^d}$ in $(0, 1]^{\Delta^d}$ space, define $\mathbf{x}, \boldsymbol{\theta} \in \Delta^d$ as the induced actions in the simplex, with $\mathbf{x}[r] = \frac{\mathbf{y}[r]}{\Lambda(\mathbf{y})}$ and $\boldsymbol{\theta}[r] = \frac{\mathbf{z}[r]}{\Lambda(\mathbf{z})}$ for every coordinate $r \in [d]$.

Spectral Properties of ψ . Prior to discussing the analysis of regret and deriving the nonnegative RVU bounds in Theorem 5.1, we examine the strong spectral properties of the regularizer ψ , which may be of independent interest for future studies. Proofs of these properties are provided in Section 5.3.4. In Theorem 5.5, we prove that the regularizer ψ is not only strongly convex, but its curvature is lower bounded, interestingly, by the diagonal matrix $\frac{1}{2\eta \sum_{k=1}^d \mathbf{y}[k]} \text{diag} \left(\frac{1}{\mathbf{y}[1]}, \frac{1}{\mathbf{y}[2]}, \dots, \frac{1}{\mathbf{y}[d]} \right)$. It is notable that these strong spectral properties of the regularizer ψ hold only when $\alpha \geq 2 + 2 \log d + 2 \log^2 d$, and it is not possible to eliminate the dependence of the parameter α on $\log^2 d$, which makes this regime particularly interesting.

Theorem 5.5. For $\alpha \geq 2 + 2 \log d + 2 \log^2 d$, the function $\psi(\mathbf{y})$ is strongly convex, and its positive definite Hessian at any point satisfies the bound

$$\nabla^2 \psi(\mathbf{y}) \succeq \frac{1}{2\eta} \begin{pmatrix} \frac{1}{\mathbf{y}^{[1]} \cdot \sum_{k=1}^d \mathbf{y}^{[k]}} & 0 & \cdots & 0 \\ 0 & \frac{1}{\mathbf{y}^{[2]} \cdot \sum_{k=1}^d \mathbf{y}^{[k]}} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \frac{1}{\mathbf{y}^{[d]} \cdot \sum_{k=1}^d \mathbf{y}^{[k]}} \end{pmatrix}.$$

Consequently, by strong convexity of the regularizer ψ (see Theorem 5.5, we can define its induced Bregman divergence D_ψ ,

$$D_\psi(\mathbf{y} \parallel \mathbf{z}) = \psi(\mathbf{y}) - \psi(\mathbf{z}) - \langle \nabla \psi(\mathbf{z}), \mathbf{y} - \mathbf{z} \rangle.$$

In turn, we establish the following spectral properties of the induced Bregman divergence D_ψ , which are pivotal for the regret analysis to be discussed next.

Theorem 5.6. The Bregman divergence $D_\psi(\cdot \parallel \cdot)$ induced by the regularizer $\psi(\cdot)$ satisfies the following properties:

- **Curvature in the lifted space** $(0, 1]^d$: $D_\psi(\cdot \parallel \cdot)$ is lower bounded by a term proportional to the ℓ_1 norm on the lifted simplex $(0, 1]^d$:

$$D_\psi(\mathbf{y} \parallel \mathbf{z}) \geq \frac{1}{2\eta} \|\mathbf{y} - \mathbf{z}\|_1^2.$$

- **Curvature in the action simplex** Δ^d : $D_\psi(\cdot \parallel \cdot)$ is lower bounded by a term proportional to the ℓ_1 norm on the action simplex Δ^d :

$$D_\psi(\mathbf{z} \parallel \mathbf{y}) \geq \frac{1}{4\eta} (1 - \epsilon) \|\boldsymbol{\theta} - \mathbf{x}\|_1^2,$$

under the multiplicative stability assumption over the sum of actions, $\omega := \frac{\Lambda(\mathbf{z})}{\Lambda(\mathbf{y})} \in [1 - \epsilon, 1 + \epsilon]$ for a constant $\epsilon \in (0, \frac{2}{5})$.

This theorem is proved in Propositions 5.25 and 5.26 in Section 5.3.4. The proof of Proposition 5.25 follows from observing that Theorem 5.5 implies that the regularizer ψ is strongly convex with respect to the ℓ_1 norm. The proof of Proposition 5.26 is more involved and relies on the special representation of $D_\psi(\cdot \parallel \cdot)$ shown in Proposition 5.24, accompanied by the multiplicative stability property, finite differences of entropies, and Pinsker's inequality.

Nonnegative Regret. Recall that the dynamics of DLRC-OMWU are equivalent to those of Equation (5.4). Thus, we begin by analyzing the regret of the OFTRL in Equation (5.4). We denote this regret by $\tilde{\text{Reg}}^{(T)}$,

$$\tilde{\text{Reg}}^{(T)} := \max_{\mathbf{y}^* \in [0, 1]^d} \sum_{t=1}^T \langle \mathbf{u}^{(t)}, \mathbf{y}^* - \mathbf{y}^{(t)} \rangle.$$

Recall that our nonnegative RVU bounds in Theorem 5.1 are on the positive regret $\max\{0, \text{Reg}^{(T)}\}$. By the following proposition, we connect the $\tilde{\text{Reg}}^{(T)}$ to $\text{Reg}^{(T)}$ defined in Equation (2.1).

Proposition 5.7. *For any time horizon $T \in \mathbb{N}$, we have that $\tilde{\text{Reg}}^{(T)} = \max\{0, \text{Reg}^{(T)}\}$. As a result, $\tilde{\text{Reg}}^{(T)} \geq 0$ and $\tilde{\text{Reg}}^{(T)} \geq \text{Reg}^{(T)}$.*

This proposition is proved in Section 5.3.5. An immediate consequence of this proposition is that any RVU bounds for the OFTRL in Equation (5.4) immediately imply nonnegative RVU bounds for DLRC-OMWU. Hence, we proceed by analyzing $\tilde{\text{Reg}}^{(T)}$.

Nonnegative RVU. With the equivalent formulation of DLRC-OMWU as the OFTRL algorithm in Equation (5.4), the strong spectral properties of ψ , and the nonnegative Regret, the analysis follows the standard machinery of Optimistic Follow-the-Regularized-Leader algorithms.

Consequently, we define $\mathbf{y}^{(t)}$ as the outputs generated by the OFTRL in Equation (5.4):

$$\mathbf{y}^{(t)} = \arg \max_{\mathbf{y} \in \Omega} -F_t(\mathbf{y}) = \arg \min_{\mathbf{y} \in \Omega} F_t(\mathbf{y}), \text{ where } F_t(\mathbf{y}) = -\langle \mathbf{U}^{(t)} + \mathbf{u}^{(t-1)}, \mathbf{y} \rangle + \psi(\mathbf{y}). \quad (5.7)$$

Moreover, we define the auxiliary sequence $\mathbf{z}^{(t)}$ as the solutions to the corresponding standard FTRL for each time t :

$$\mathbf{z}^{(t)} = \arg \max_{\mathbf{z} \in \Omega} -G_t(\mathbf{z}) = \arg \min_{\mathbf{z} \in \Omega} G_t(\mathbf{z}), \text{ where } G_t(\mathbf{z}) = -\langle \mathbf{U}^{(t)}, \mathbf{z} \rangle + \psi(\mathbf{z}). \quad (5.8)$$

Functions G_t and F_t are strongly convex, as shown by Theorem 5.5. We continue the analysis with the following standard lemma from OFTRL analysis.

Lemma 5.8. *For any $\mathbf{y} \in \Omega$, and the sequence $\{\mathbf{y}^{(t)}\}_{t=1}^T$ generated by Equation (5.7), it holds that*

$$\begin{aligned} \sum_{t=1}^T \langle \mathbf{y} - \mathbf{y}^{(t)}, \mathbf{u}^{(t)} \rangle &\leq \psi(\mathbf{y}) - \psi(\mathbf{y}^{(1)}) + \underbrace{\sum_{t=1}^T \langle \mathbf{z}^{(t+1)} - \mathbf{y}^{(t)}, \mathbf{u}^{(t)} - \mathbf{u}^{(t-1)} \rangle}_{\text{(I)}} \\ &\quad - \underbrace{\sum_{t=1}^T \left(D_\psi(\mathbf{y}^{(t)} \parallel \mathbf{z}^{(t)}) + D_\psi(\mathbf{z}^{(t+1)} \parallel \mathbf{y}^{(t)}) \right)}_{\text{(II)}}. \end{aligned}$$

The general plan for our RVU bounds is to carefully upper bound term (I) to obtain the beta terms $O(\eta) \sum_{t=1}^{T-1} \|\boldsymbol{\nu}^{(t)} - \boldsymbol{\nu}^{(t-1)}\|_\infty^2$, and to upper bound term (II) to obtain the gamma terms $-\Omega \left(\frac{1}{\eta} \right) \sum_{t=1}^{T-1} \|\mathbf{x}^{(t+1)} - \mathbf{x}^{(t)}\|_1^2$. This procedure is formalized in the proof of Theorem 5.1 in Section 5.3.6.

A key observation in this regard is that, by Theorem 5.6 and the multiplicative stability of the Dynamic Learning Rate Control Step in Theorem 3.5, the Bregman divergence $D_\psi(\mathbf{y} \parallel \mathbf{z})$ is not only lower bounded by the squared norm of the difference of actions in the lifted space, $\Omega \left(\frac{1}{\eta} \right) \|\mathbf{y} - \mathbf{z}\|_1^2$ (Lemma 5.30), but also, due to our novel choice of regularizer, it is lower bounded by the squared norm of the difference of actions in the original space (action simplex), $\Omega \left(\frac{1}{\eta} \right) \|\boldsymbol{\theta} - \mathbf{x}\|_1^2$ (Lemma 5.31). Combining these two bounds, we obtain that

$$D_\psi(\mathbf{y} \parallel \mathbf{z}) \geq \Omega \left(\frac{1}{\eta} \right) \left(\|\mathbf{y} - \mathbf{z}\|_1^2 + \|\boldsymbol{\theta} - \mathbf{x}\|_1^2 \right). \quad (5.9)$$

By choosing a small learning rate η , the $\Omega\left(\frac{1}{\eta}\right) \|\mathbf{y} - \mathbf{z}\|_1^2$ terms in Equation (5.9) compensate for the $O(\eta) \|\mathbf{y} - \mathbf{z}\|_1^2$ that results from converting term (I) into the beta terms $O(\eta) \sum_{t=1}^{T-1} \|\boldsymbol{\nu}^{(t)} - \boldsymbol{\nu}^{(t-1)}\|_\infty^2$. This conversion process follows from the relationship between the \mathbf{u} s and $\boldsymbol{\nu}$ s, basic calculations involving the Cauchy–Schwarz and Hölder’s inequalities, and a sufficiently small choice of η . Details are provided in the proof of Theorem 5.1 in Section 5.3.6.

5.3 Detailed Analysis

In this section, we first state the auxiliary lemmas used in the analysis, and then provide the detailed analysis and formal proofs for the sketches presented in the previous sections.

5.3.1 Auxiliary Lemmas

Theorem 5.9 (Theorem 2.2.5 of [Ren01]). *Given a self-concordant function $f : \mathbb{R}^d \rightarrow \mathbb{R}$, and the induced local norm $\|\cdot\|_{\mathbf{x}} := \|\cdot\|_{\nabla^2 f(\mathbf{x})}$, if for some \mathbf{x} in the domain of f we have $\|n(\mathbf{x})\|_{\mathbf{x}} \leq \frac{1}{4}$, then f has a minimizer \mathbf{z} that*

$$\|z - x\|_x \leq 3\|n(\mathbf{x})\|_x,$$

where $n(\mathbf{x})$ is the Newton update, $n(\mathbf{x}) = -[\nabla^2 f(\mathbf{x})]^{-1} \nabla f(\mathbf{x})$. By

Lemma 5.10. *For any set of numbers s_1, s_2, \dots, s_d , the log-sum-exp function $g(t) := h(ts_1, ts_2, \dots, ts_d) = \log\left(\sum_{k=1}^d \exp\{ts_k\}\right)$ satisfies,*

$$\max\{s_1, s_2, \dots, s_d\} \leq \frac{g(t)}{t} \leq \max\{s_1, s_2, \dots, s_d\} + \frac{\log d}{t}.$$

Proof. The LHS follows by lowerbounding $\sum_{k=1}^d \exp\{to_k\}$ by $\max\{\exp\{to_k\}\}$. The RHS follows by upperbounding $\sum_{k=1}^d \exp\{to_k\}$ by $d \cdot \max\{\exp\{to_k\}\}$. \square

Lemma 5.11 (Pinsker’s inequality). *Given two discrete random variables \mathbf{p} and \mathbf{q} ,*

$$\|\mathbf{p} - \mathbf{q}\|_1^2 \leq 2\text{KL}(\mathbf{p}\|\mathbf{q}).$$

Lemma 5.12 (Entropy difference). *Given two discrete random variables \mathbf{p} and \mathbf{q} with support size d ,*

$$|\text{H}(\mathbf{p}) - \text{H}(\mathbf{q})| \leq (\log d) \sqrt{2\text{KL}(\mathbf{p}\|\mathbf{q})}.$$

5.3.2 Proofs for Learning Rate Control Step (Section 3.1)

We start proving Theorem 3.3, by first proving strong concavity of optimization problem (3.2).

Lemma 5.13. *(Strong concavity of learning rate control Step) The objective of the optimization problem,*

$$\widehat{\lambda} := \arg \max_{\lambda \in (0, \eta]} \left\{ f(\lambda; \mathbf{r}) := (\alpha - 1) \log \lambda + \log \left(\sum_{k=1}^d e^{\lambda \mathbf{r}^{[k]}} \right) \right\},$$

is $\frac{\alpha - \log^2 d - 1}{\lambda^2}$ -strongly concave.

Proof. Taking derivatives from $f(\lambda; \mathbf{r})$ w.r.t. λ ,

$$\begin{aligned}\frac{\partial f(\lambda; \mathbf{r})}{\partial \lambda} &= \frac{\sum_{k=1}^d \mathbf{r}[k] e^{\lambda \mathbf{r}[k]}}{\sum_{k=1}^d e^{\lambda \mathbf{r}[k]}} + \frac{\alpha - 1}{\lambda}. \\ \frac{\partial^2 f(\lambda; \mathbf{r})}{\partial \lambda^2} &= \frac{\left(\sum_{k=1}^d \mathbf{r}[k]^2 e^{\lambda \mathbf{r}[k]} \right) \left(\sum_{k=1}^d e^{\lambda \mathbf{r}[k]} \right) - \left(\sum_{k=1}^d \mathbf{r}[k] e^{\lambda \mathbf{r}[k]} \right)^2}{\left(\sum_{k=1}^d e^{\lambda \mathbf{r}[k]} \right)^2} - \frac{\alpha - 1}{\lambda^2} \\ &= \frac{\sum_{k=1}^d \mathbf{r}[k]^2 e^{\lambda \mathbf{r}[k]}}{\sum_{k=1}^d e^{\lambda \mathbf{r}[k]}} - \left(\frac{\sum_{k=1}^d \mathbf{r}[k] e^{\lambda \mathbf{r}[k]}}{\sum_{k=1}^d e^{\lambda \mathbf{r}[k]}} \right)^2 - \frac{\alpha - 1}{\lambda^2}.\end{aligned}$$

Substituting $\mathbf{r}[r] = \frac{\log \mathbf{x}[r] + \log \Gamma}{\lambda}$ for every coordinate $r \in [d]$, where $\Gamma = \sum_{k=1}^d \exp\{\lambda \mathbf{r}[k]\}$ and $\mathbf{x} \in \Delta^d$, we get

$$\begin{aligned}\frac{\partial^2 f(\lambda; \mathbf{r})}{\partial \lambda^2} &= \frac{1}{\lambda^2} \sum_{k=1}^d (\mathbf{x}[k] (\log \mathbf{x}[k] + \log \Gamma)^2) - \frac{1}{\lambda^2} \left(\sum_{k=1}^d \mathbf{x}[k] (\log \mathbf{x}[k] + \log \Gamma) \right)^2 - \frac{\alpha - 1}{\lambda^2} \\ &= \frac{1}{\lambda^2} \sum_{k=1}^d (\mathbf{x}[k] (\log \mathbf{x}[k] + \log \Gamma)^2) - \frac{1}{\lambda^2} \left(\sum_{k=1}^d \mathbf{x}[k] \log \mathbf{x}[k] + \log \Gamma \right)^2 - \frac{\alpha - 1}{\lambda^2} \\ &= \frac{1}{\lambda^2} \left(\sum_{k=1}^d \mathbf{x}[k] \log^2 \mathbf{x}[k] - \left(\sum_{k=1}^d \mathbf{x}[k] \log \mathbf{x}[k] \right)^2 - (\alpha - 1) \right) \\ &\leq \frac{1}{\lambda^2} (\log^2 d - (\alpha - 1)).\end{aligned}$$

Therefore, the function $f(\lambda; \mathbf{r})$ is $\frac{\alpha - \log^2 d - 1}{\lambda^2}$ -strongly concave. \square

Secondly, we show that the optimization problem (3.2) is self-concordant.

Lemma 5.14 (Self-concordance of learning rate control step). *The objective of the optimization problem,*

$$\hat{\lambda} := \arg \max_{\lambda \in (0, \eta]} \left\{ f(\lambda; \mathbf{r}) := (\alpha - 1) \log \lambda + \log \left(\sum_{k=1}^d e^{\lambda \mathbf{r}[k]} \right) \right\},$$

is self-concordant., i.e.,

$$\left(\frac{\partial^3 f(\lambda; \mathbf{r})}{\partial \lambda^3} \right)^2 \leq -4 \left(\frac{\partial^2 f(\lambda; \mathbf{r})}{\partial \lambda^2} \right)^3.$$

Proof. Recall that from proof of Lemma 5.13,

$$\frac{\partial^2 f(\lambda; \mathbf{r})}{\partial \lambda^2} = \frac{\sum_{k=1}^d \mathbf{r}[k]^2 e^{\lambda \mathbf{r}[k]}}{\sum_{k=1}^d e^{\lambda \mathbf{r}[k]}} - \left(\frac{\sum_{k=1}^d \mathbf{r}[k] e^{\lambda \mathbf{r}[k]}}{\sum_{k=1}^d e^{\lambda \mathbf{r}[k]}} \right)^2 - \frac{\alpha - 1}{\lambda^2},$$

Furthermore,

$$-\frac{\partial^2 f(\lambda; \mathbf{r})}{\partial \lambda^2} \geq \frac{\alpha - \log^2 d - 1}{\lambda^2}.$$

Taking derivatives,

$$\begin{aligned} \frac{\partial^3 f(\lambda; \mathbf{r})}{\partial \lambda^3} &= \frac{\sum_{k=1}^d \mathbf{r}[k]^3 e^{\lambda \mathbf{r}[k]}}{\sum_{k=1}^d e^{\lambda \mathbf{r}[k]}} - \frac{\left(\sum_{k=1}^d \mathbf{r}[k]^2 e^{\lambda \mathbf{r}[k]}\right) \left(\sum_{k=1}^d \mathbf{r}[k] e^{\lambda \mathbf{r}[k]}\right)}{\left(\sum_{k=1}^d e^{\lambda \mathbf{r}[k]}\right)^2} \\ &\quad - 2 \left(\frac{\sum_{k=1}^d \mathbf{r}[k] e^{\lambda \mathbf{r}[k]}}{\sum_{k=1}^d e^{\lambda \mathbf{r}[k]}} \right) \left(\frac{\sum_{k=1}^d \mathbf{r}[k]^2 e^{\lambda \mathbf{r}[k]}}{\sum_{k=1}^d e^{\lambda \mathbf{r}[k]}} - \left(\frac{\sum_{k=1}^d \mathbf{r}[k] e^{\lambda \mathbf{r}[k]}}{\sum_{k=1}^d e^{\lambda \mathbf{r}[k]}} \right)^2 \right) + \frac{2(\alpha - 1)}{\lambda^3}. \end{aligned}$$

Substituting $\mathbf{r}[r] = \frac{\log \mathbf{x}[r] + \log \Gamma}{\lambda}$ for every coordinate $r \in [d]$, where $\Gamma = \sum_{k=1}^d \exp\{\lambda \mathbf{r}[k]\}$ and $\mathbf{x} \in \Delta^d$, we get

$$\begin{aligned} &\frac{\partial^3 f(\lambda; \mathbf{r})}{\partial \lambda^3} \\ &= \frac{1}{\lambda^3} \sum_{k=1}^d \mathbf{x}[k] (\log \mathbf{x}[k] + \log \Gamma)^3 - 3 \frac{1}{\lambda^3} \left(\sum_{k=1}^d \mathbf{x}[k] (\log \mathbf{x}[k] + \log \Gamma)^2 \right) \left(\sum_{k=1}^d \mathbf{x}[k] (\log \mathbf{x}[k] + \log \Gamma) \right) \\ &\quad + 2 \frac{1}{\lambda^3} \left(\sum_{k=1}^d \mathbf{x}[k] (\log \mathbf{x}[k] + \log \Gamma) \right)^3 + \frac{2(\alpha - 1)}{\lambda^3} \\ &= \frac{1}{\lambda^3} \left(\sum_{k=1}^d \mathbf{x}[k] \log^3 \mathbf{x}[k] - 3 \left(\sum_{k=1}^d \mathbf{x}[k] \log^2 \mathbf{x}[k] \right) \left(\sum_{k=1}^d \mathbf{x}[k] \log \mathbf{x}[k] \right) + 2 \left(\sum_{k=1}^d \mathbf{x}[k] \log \mathbf{x}[k] \right)^3 + 2(\alpha - 1) \right) \\ &\leq \frac{1}{\lambda^3} (3 \log^3 d + 2(\alpha - 1)) \end{aligned}$$

Thus,

$$\left(\frac{\partial^3 f(\lambda; \mathbf{r})}{\partial \lambda^3} \right)^2 \leq \left(\frac{3 \log^3 d + 2(\alpha - 1)}{\lambda^3} \right)^2 \leq 4 \left(\frac{\alpha - \log^2 d - 1}{\lambda^2} \right)^3 \leq -4 \left(\frac{\partial^2 f(\lambda; \mathbf{r})}{\partial \lambda^3} \right)^3,$$

for all $d \geq 1$,⁴ where $\alpha = \beta \log^2 d + 2 \log d + 2$ and $\beta \geq 2$, hence the proof is concluded. \square

Theorem 3.3 (Properties of learning rate control step). *For any $\mathbf{r} \in \mathbb{R}^d$, the rate control objective $f(\lambda; \mathbf{r})$, defined in*

$$\hat{\lambda} := \arg \max_{\lambda \in (0, \eta]} \left\{ f(\lambda; \mathbf{r}) := (\alpha - 1) \log \lambda + \log \left(\sum_{k=1}^d e^{\lambda \mathbf{r}[k]} \right) \right\},$$

satisfies the following properties:

- *Strong concavity:* $f''(\lambda; \mathbf{r}) \leq -(\alpha - \log^2 d - 1)/\lambda^2$ for all $\lambda \in (0, \infty)$.
- *Self-concordance:* $(f'''(\lambda; \mathbf{r}))^2 \leq -4f''(\lambda; \mathbf{r})^3$,

where all derivatives are with respect to λ .

Proof. We show strong concavity in Lemma 5.13 and self-concordance in Lemma 5.14. \square

⁴the equality happens at $d = 1$.

Next, we need to prove Theorem 3.5. To do so, we first prove and state the following lemma, which demonstrate that when the the maximum regret $\max_k \{\mathbf{r}^{(t)}[k]\}$ accumulated on the actions is not too negative, the optimal solution to optimization problem (3.2) is $\lambda^{(t)} = \eta$.

Lemma 5.16. *Given an arbitrary vector $\mathbf{r} \in \mathbb{R}^d$, and $\hat{\lambda}$ as the solution to*

$$\hat{\lambda} := \arg \max_{\lambda \in (0, \eta]} \left\{ f(\lambda; \mathbf{r}) := (\alpha - 1) \log \lambda + \log \left(\sum_{k=1}^d e^{\lambda \mathbf{r}[k]} \right) \right\}, \quad (5.10)$$

then as long as $\max_{r \in [d]} \{\mathbf{r}[r]\} \geq -\beta \log^2 d$, we have $\hat{\lambda} = \eta$.

Proof. By KKT conditions, concavity and uniqueness of the solution, it is obvious that whenever $f'(\eta; \mathbf{r}) \geq 0$, then $\hat{\lambda} = \eta$. Thus, with the same arguments as proof of Lemma 5.17,

$$\begin{aligned} f'(\eta; \mathbf{r}) &= \frac{\alpha - 1}{\eta} + \frac{1}{\eta} \mathbf{H}(\mathbf{x}) + \frac{1}{\eta} \log \left(\sum_{k=1}^d e^{\lambda \mathbf{r}[k]} \right) \\ &\geq \frac{\alpha - 1}{\eta} - \frac{\log d}{\eta} + \frac{1}{\eta} \max_{r \in [d]} \{\mathbf{r}[r]\} \\ &\geq \frac{1}{\eta} \left(\beta \log^2 d + \max_{r \in [d]} \{\mathbf{r}[r]\} \right). \end{aligned}$$

Thus, whenever $\max_{r \in [d]} \{\mathbf{r}[r]\} \geq -\beta \log^2 d$, then $\hat{\lambda} = \eta$. \square

In sequel, to prove Theorem 3.5 (sensitivity of learning rates on regrets), we state and prove the following lemma on stability of learning rates in multiplicative sense. Combining a good analytic guess λ_0 for the value of λ with techniques for self-concordant function analysis, we establish that the intrinsic norm of the second-order ascent direction of f is small at λ_0 . This allows us to conclude that the solution λ must be within a small radius from λ_0 in intrinsic norm. Furthermore, using the bound on $f''(\lambda_0; \cdot)$ given in Theorem 3.3, we can finally conclude proximity in the multiplicative sense.

Lemma 5.17. *(Multiplicative Stability) Given vectors $\mathbf{r}, \mathbf{r}' \in \mathbb{R}^d$, then $\hat{\lambda}, \hat{\lambda}'$, the solutions to*

$$\hat{\lambda} = \arg \max_{\lambda \in (0, \eta]} \left\{ f(\lambda; \mathbf{r}) := (\alpha - 1) \log \lambda + \log \left(\sum_{k=1}^d e^{\lambda \mathbf{r}[k]} \right) \right\}, \quad (5.11)$$

and

$$\hat{\lambda}' = \arg \max_{\lambda \in (0, \eta]} \left\{ f(\lambda; \mathbf{r}') := (\alpha - 1) \log \lambda + \log \left(\sum_{k=1}^d e^{\lambda \mathbf{r}'[k]} \right) \right\} \quad (5.12)$$

respectively, are multiplicatively stable, i.e.,

$$\frac{\beta - 5}{\beta + 3} \left(\frac{\min_{r \in d} (-\mathbf{r}'[r])}{\min_{r \in d} (-\mathbf{r}[r])} \right) \leq \frac{\hat{\lambda}}{\hat{\lambda}'} \leq \frac{\beta + 3}{\beta - 5} \left(\frac{\min_{r \in d} (-\mathbf{r}'[r])}{\min_{r \in d} (-\mathbf{r}[r])} \right).$$

Proof. We use suboptimality bound entailed by Newton update for self-concordance functions to prove stability. Recall that from proof of Lemma 5.13,

$$\begin{aligned}\frac{\partial f(\lambda; \mathbf{r})}{\partial \lambda} &= \frac{\sum_{k=1}^d \mathbf{r}[k] e^{\lambda \mathbf{r}[k]}}{\sum_{k=1}^d e^{\lambda \mathbf{r}[k]}} + \frac{\alpha - 1}{\lambda}. \\ \frac{\partial^2 f(\lambda; \mathbf{r})}{\partial \lambda^2} &= \frac{\sum_{k=1}^d \mathbf{r}[k]^2 e^{\lambda \mathbf{r}[k]}}{\sum_{k=1}^d e^{\lambda \mathbf{r}[k]}} - \left(\frac{\sum_{k=1}^d \mathbf{r}[k] e^{\lambda \mathbf{r}[k]}}{\sum_{k=1}^d e^{\lambda \mathbf{r}[k]}} \right)^2 - \frac{\alpha - 1}{\lambda^2}.\end{aligned}$$

Let us choose the primary guess $\lambda_0 = \frac{\alpha - 1}{\min_{r \in [d]} (-\mathbf{r}[r])}$. By change of variables, we know that

$$\frac{\partial f}{\partial \lambda}(\lambda_0; \mathbf{r}) = \frac{\alpha - 1}{\lambda_0} + \frac{1}{\lambda_0} \mathbf{H}(\mathbf{x}) + \frac{1}{\lambda_0} \log \left(\sum_{k=1}^d e^{\lambda_0 \mathbf{r}[k]} \right),$$

where $\mathbf{H}(\mathbf{x})$ is the negative entropy of the vector $\mathbf{x}[r] = \frac{\exp\{\lambda_0 \mathbf{r}[r]\}}{\sum_{k=1}^d \exp\{\lambda_0 \mathbf{r}[k]\}}$. Therefore, $0 \geq H(\mathbf{x}) \geq -\log d$. On the other hand, by Lemma 5.10,

$$\max_{r \in [d]} \{\mathbf{r}[r]\} \leq \frac{1}{\lambda_0} \log \left(\sum_{k=1}^d e^{\lambda_0 \mathbf{r}[k]} \right) \leq \max_{r \in [d]} \{\mathbf{r}[r]\} + \frac{\log d}{\lambda_0}.$$

Hence,

$$-\frac{\log d}{\lambda_0} = \min_{r \in [d]} \{-\mathbf{r}[r]\} + \max_{r \in [d]} \{\mathbf{r}[r]\} - \frac{\log d}{\lambda_0} \leq f'(\lambda_0, \mathbf{r}) \leq \min_{r \in [d]} \{-\mathbf{r}[r]\} + \max_{r \in [d]} \{\mathbf{r}[r]\} + \frac{\log d}{\lambda_0} = \frac{\log d}{\lambda_0}.$$

On the other hand, from Lemma 5.13 we know that

$$|f''(\lambda_0; \mathbf{r})| \geq \frac{1}{\lambda_0^2} (\alpha - \log^2 d - 1).$$

Therefore, the local norm of Newton step $n(\lambda_0)$ is

$$\begin{aligned}\|n(\lambda_0)\|_{f''(\lambda_0; \mathbf{r})}^2 &= \frac{f'(\lambda_0, \mathbf{r})^2}{|f''(\lambda_0, \mathbf{r})|} \\ &\leq \frac{\log^2 d}{(\beta - 1) \log^2 d + 2 \log d + 1} \\ &\leq \frac{1}{\beta - 1},\end{aligned}$$

controlled by the hyperparameter $\beta \geq 17$.⁵ In turn, by Theorem 5.9,

$$\begin{aligned}\|\hat{\lambda} - \lambda_0\|_{f''(\lambda_0; \mathbf{r})} &\leq \frac{3}{\sqrt{\beta - 1}} \\ \Rightarrow |\hat{\lambda} - \lambda_0| &\leq \frac{3}{\sqrt{(\beta - 1)|f''(\lambda_0; \mathbf{r})|}} \leq \frac{4\lambda_0}{\beta - 1} \\ \Rightarrow \left| \frac{\hat{\lambda}}{\lambda_0} - 1 \right| &\leq \frac{4}{\beta - 1},\end{aligned}\tag{5.13}$$

⁵recall that $\alpha = \beta \log^2 d + 2 \log d + 2$.

for all $d \geq 2$. In a similar manner, we choose the primary guess $\lambda'_0 = \frac{\alpha - 1}{\min_{r \in d} \{-\mathbf{r}'[r]\}}$ and we infer,

$$\left| \frac{\widehat{\lambda}'}{\lambda'_0} - 1 \right| \leq \frac{4}{\beta - 1}. \quad (5.14)$$

Combining Equations (5.13) and (5.14),

$$\begin{aligned} (1 - \frac{4}{\beta - 1})\lambda_0 &\leq \widehat{\lambda} \leq (1 + \frac{4}{\beta - 1})\lambda_0 \\ (1 - \frac{4}{\beta - 1})\lambda'_0 &\leq \widehat{\lambda}' \leq (1 + \frac{4}{\beta - 1})\lambda'_0 \\ \Rightarrow \frac{\beta - 5}{\beta + 3} \left(\frac{\min_{r \in d} (-\mathbf{r}'[r])}{\min_{r \in d} (-\mathbf{r}[r])} \right) &= \frac{\beta - 5}{\beta + 3} \cdot \frac{\lambda_0}{\lambda'_0} \leq \frac{\widehat{\lambda}}{\widehat{\lambda}'} \leq \frac{\beta + 3}{\beta - 5} \cdot \frac{\lambda_0}{\lambda'_0} = \frac{\beta + 3}{\beta - 5} \left(\frac{\min_{r \in d} (-\mathbf{r}'[r])}{\min_{r \in d} (-\mathbf{r}[r])} \right). \end{aligned} \quad (5.15)$$

□

Now, we are ready to state and prove Theorem 3.5. To this end, we consider three cases based on whether $\max_{r \in [d]} \{\mathbf{r}[r]\}$ and $\max_{r \in [d]} \{\mathbf{r}'[r]\}$ are higher or lower than $-\beta \log^2 d$. In the regimes, where Lemma 5.16 is satisfied the proof is immediate. For other cases, we use careful analysis with the help of Lemma 5.17.

Theorem 3.5. (*Sensitivity of learning rates on regrets*) *There exists a universal constant β ,⁶ such that for $\alpha \geq 2 + 2 \log d + \beta \log^2 d$, the following property holds. Let $\mathbf{r}, \mathbf{r}' \in \mathbb{R}^d$ be such that $\|\mathbf{r} - \mathbf{r}'\|_\infty \leq 2$, and let $\widehat{\lambda}, \widehat{\lambda}'$ the corresponding learning rates, that is,*

$$\widehat{\lambda} = \arg \max_{t \in (0, \eta]} f(t; \mathbf{r}), \quad \widehat{\lambda}' = \arg \max_{t \in (0, \eta]} f(t; \mathbf{r}').$$

Then, $\widehat{\lambda}$ and $\widehat{\lambda}'$ are multiplicatively stable; specifically,

$$\frac{7}{10} \leq \frac{\widehat{\lambda}}{\widehat{\lambda}'} \leq \frac{7}{5},$$

Proof. Here we show the lemma for $\|\mathbf{r} - \mathbf{r}'\|_\infty \leq 2$. The extension to general $\|\mathbf{r} - \mathbf{r}'\|_\infty \leq o(1)$ is easy to infer by choosing β large enough. We have three cases depending on the size of $\max_{r \in [d]} \{\mathbf{r}[r]\}$.

1. If $\max_{r \in [d]} \{\mathbf{r}[r]\} \geq -\beta \log^2 d$ and $\max_{r \in [d]} \{\mathbf{r}'[r]\} \geq -\beta \log^2 d$, then by Lemma 5.16, we conclude that $\widehat{\lambda} = \widehat{\lambda}' = \eta$.
2. If $\max_{r \in [d]} \{\mathbf{r}[r]\} < -\beta \log^2 d$ and $\max_{r \in [d]} \{\mathbf{r}'[r]\} < -\beta \log^2 d$, then by Lemma 5.17 and $\|\mathbf{r} - \mathbf{r}'\|_\infty \leq 2$,

$$\frac{4}{5} \frac{\beta - 5}{\beta + 3} \leq \frac{\beta - 5}{\beta + 3} \left(\frac{\beta \log^2 d}{2 + \beta \log^2 d} \right) \leq \frac{\widehat{\lambda}}{\widehat{\lambda}'} \leq \frac{\beta + 3}{\beta - 5} \left(\frac{2 + \beta \log^2 d}{\beta \log^2 d} \right) \leq \frac{6}{5} \frac{\beta + 3}{\beta - 5},$$

for $\beta \geq 20$ and $d \geq 2$.

⁶For concrete values, choosing any $\beta \geq 70$ suffices.

3. And otherwise, if without loss of generality assume that $\max_{r \in [d]} \{\mathbf{r}[r]\} \geq -\beta \log^2 d$ and $\max_{r \in [d]} \{\mathbf{r}'[r]\} < -\beta \log^2 d$. Then, $\widehat{\lambda} = 1$. On the other hand, based on Equation (5.15) in the proof of Lemma 5.17,

$$\left(1 - \frac{4}{\beta - 1}\right) \lambda'_0 \leq \widehat{\lambda}' \leq \left(1 + \frac{4}{\beta - 1}\right) \lambda'_0.$$

Hence,

$$\frac{\alpha - 1}{2 + \beta \log^2 d} \frac{\beta - 5}{\beta - 1} \leq \frac{\beta - 5}{\beta - 1} \frac{\alpha - 1}{\min_{r \in d} \{-\alpha'[r]\}} \leq \frac{\widehat{\lambda}'}{\widehat{\lambda}} \leq \frac{\beta + 3}{\beta - 1} \frac{\alpha - 1}{\min_{r \in d} \{-\alpha'[r]\}} \leq \frac{\alpha - 1}{\beta \log^2 d} \frac{\beta + 3}{\beta - 1}$$

since $\|\mathbf{r} - \mathbf{r}'\|_\infty \leq 2$ and thus $\min_{r \in d} \{-\mathbf{r}'[r]\} \leq 2 + \beta \log^2 d$. Consequently,

$$\frac{4\beta - 5}{5\beta - 1} \leq \frac{\alpha - 1}{2 + \beta \log^2 d} \frac{\beta - 5}{\beta - 1} \leq \frac{\widehat{\lambda}'}{\widehat{\lambda}} \leq \frac{\alpha - 1}{\beta \log^2 d} \frac{\beta + 3}{\beta - 1} \leq \frac{6\beta + 3}{5\beta - 1}$$

for $\beta \geq 50$ and $d \geq 2$,

Putting all cases together, for $\beta \geq 50$,

$$\frac{4\beta - 5}{5\beta + 3} \leq \frac{\widehat{\lambda}'}{\widehat{\lambda}} \leq \frac{6\beta + 3}{5\beta - 5}.$$

Now, by choosing any $\beta \geq 70$,

$$\frac{7}{10} = \frac{4}{5} \times \frac{7}{8} \leq \frac{\widehat{\lambda}'}{\widehat{\lambda}} \leq \frac{6}{5} \times \frac{8}{7} < \frac{7}{5} \quad \Rightarrow \quad \frac{\widehat{\lambda}'}{\widehat{\lambda}} \in \left[\frac{21}{20} - \frac{7}{20}, \frac{21}{20} + \frac{7}{20} \right],$$

and the proof is completed. \square

5.3.3 Equivalent Viewpoints of DLRC-OMWU

In this section, we discuss equivalent formulations of DLRC-OMWU. First, to conceptualize the idea of dynamic learning rate control, we show that DLRC-OMWU (as presented in Algorithm 1) is a special case of Formulation 5.3 with the choice of regularizers $\rho(\lambda) = (\alpha - 1) \log \lambda$ and $\phi(\mathbf{x}) = \sum_{k=1}^d \mathbf{x}[k] \log \mathbf{x}[k]$, making it a special, easily computable instance. Secondly, we show that, for the purpose of regret analysis, DLRC-OMWU is equivalent to OFTRL with the regularizer ψ on the space $(0, 1]^d$ defined in Equation (5.5). This equivalence allows us to leverage standard techniques to analyze the regret, given the strong spectral properties of the regularizer ψ , which we will establish later.

Theorem 5.4. *Algorithm 1 (DLRC-OMWU) can alternatively be viewed as DLRC-OFTRL (Formulation 5.3) with the choice of regularizers $\rho(\lambda) = (\alpha - 1) \log \lambda$ and $\phi(\mathbf{x}) = \sum_{k=1}^d \mathbf{x}[k] \log \mathbf{x}[k]$, i.e.,*

$$\begin{pmatrix} \lambda^{(t)} \\ \mathbf{x}^{(t)} \end{pmatrix} \leftarrow \arg \max_{\lambda \in (0, \eta], \mathbf{x} \in \Delta^d} \left\{ \lambda \langle \mathbf{r}^{(t)}, \mathbf{x} \rangle + (\alpha - 1) \log \lambda - \sum_{k=1}^d \mathbf{x}[k] \log \mathbf{x}[k] \right\}.$$

or additionally as *OFTRL* with the regularizer ψ defined in Equation (5.5), i.e.,

$$\mathbf{y}^{(t)} \leftarrow \arg \max_{\mathbf{y} \in (0,1]^{\Delta^d}} \left\{ \eta \langle \mathbf{r}^{(t)}, \mathbf{y} \rangle + \alpha \log \left(\sum_{k=1}^d \mathbf{y}[k] \right) - \frac{1}{\sum_{k=1}^d \mathbf{y}[k]} \sum_{k=1}^d \mathbf{y}[k] \log \mathbf{y}[k] \right\},$$

In other words, the three perspectives (Algorithm 1, Formulations 5.6 and 5.4) of *DLRC-OMWU* are equivalent and result in the same learning dynamics.

Proof. The proof follows by combining Proposition 5.20, Lemma 5.21, and Corollary 5.22. \square

Proposition 5.20. *The optimization step of the following OFTRL,*

$$\mathbf{y}^{(t)} \leftarrow \arg \max_{\mathbf{y} \in (0,1]^{\Delta^d}} \left\{ \eta \langle \mathbf{U}^{(t)} + \mathbf{u}^{(t-1)}, \mathbf{y} \rangle + \alpha \log \left(\sum_{k=1}^d \mathbf{y}[k] \right) - \frac{1}{\sum_{k=1}^d \mathbf{y}[k]} \sum_{k=1}^d \mathbf{y}[k] \log \mathbf{y}[k] \right\} \quad (5.16)$$

is equivalent to optimization step of the following special instance of Formulation 5.3,

$$\begin{pmatrix} \lambda^{(t)} \\ \mathbf{x}^{(t)} \end{pmatrix} \leftarrow \arg \max_{\lambda \in (0,1], \mathbf{x} \in \Delta^d} \left\{ \eta \lambda \langle \mathbf{U}^{(t)} + \mathbf{u}^{(t-1)}, \mathbf{x} \rangle + (\alpha - 1) \log \lambda - \sum_{k=1}^d \mathbf{x}[k] \log \mathbf{x}[k] \right\} \quad (5.17)$$

under the change of variables $\lambda^{(t)} = \sum_{k=1}^d \mathbf{y}[k]$ and $\mathbf{x}^{(t)}[r] = \frac{\mathbf{y}[r]}{\sum_{k=1}^d \mathbf{y}[k]}$.

Proof. By choice of $\mathbf{y} = \lambda \mathbf{x}$, for any feasible point in the first optimization problem,

$$\begin{aligned} & \eta \langle \mathbf{U}^{(t)} + \mathbf{u}^{(t-1)}, \lambda \mathbf{x} \rangle + \alpha \log \left(\sum_{k=1}^d \lambda \mathbf{x}[k] \right) - \frac{1}{\sum_{k=1}^d \lambda \mathbf{x}[k]} \sum_{k=1}^d \lambda \mathbf{x}[k] \log (\lambda \mathbf{x}[k]) \\ &= \lambda \langle \mathbf{U}^{(t)} + \mathbf{u}^{(t-1)}, \mathbf{x} \rangle + \alpha \log(\lambda) + \alpha \log \left(\sum_{k=1}^d \mathbf{x}[k] \right) - \frac{1}{\sum_{k=1}^d \mathbf{x}[k]} \sum_{k=1}^d \mathbf{x}[k] (\log \mathbf{x}[k] + \log \lambda) \\ &= \eta \lambda \langle \mathbf{U}^{(t)} + \mathbf{u}^{(t-1)}, \mathbf{x} \rangle + (\alpha - 1) \log \lambda - \sum_{k=1}^d \mathbf{x}[k] \log \mathbf{x}[k], \end{aligned}$$

where the last line follows since $\sum_{k=1}^d \mathbf{x}[k] = 1$. The other direction holds similarly. \square

Lemma 5.21. *The optimal solution $(\lambda^{(t)}, \mathbf{x}^{(t)})$ to the following special instance of *DLRC-OFTRL*,*

$$\begin{pmatrix} \lambda^{(t)} \\ \mathbf{x}^{(t)} \end{pmatrix} \leftarrow \arg \max_{\lambda \in (0,\eta], \mathbf{x} \in \Delta^d} \left\{ \lambda \langle \mathbf{U}^{(t)} + \mathbf{u}^{(t-1)}, \mathbf{x} \rangle + (\alpha - 1) \log \lambda - \sum_{k=1}^d \mathbf{x}[k] \log \mathbf{x}[k] \right\} \quad (5.18)$$

satisfies,⁷

$$\mathbf{x}^{(t)}[r] = \text{softmax}(\lambda^{(t)} \mathbf{r}^{(t)})[r] = \frac{\exp\{\lambda^{(t)} \mathbf{r}^{(t)}[r]\}}{\sum_{k=1}^d \exp\{\lambda^{(t)} \mathbf{r}^{(t)}[k]\}}, \quad (5.19)$$

for every coordinate $r \in [d]$, where $\mathbf{r}^{(t)}[r] = \mathbf{U}^{(t)}[r] + \mathbf{u}^{(t-1)}[r]$.

⁷Note that this formulation is the same as Equation (5.17). It is entailed by simply scaling λ by η .

Proof. By KKT conditions,

$$\lambda^{(t)} \mathbf{r}^{(t)} - [\log x^{(t)}[1], \log x^{(t)}[2], \dots, \log x^{(t)}[d]]^\top = \mu \in \mathbb{R}.$$

Therefore, $\mathbf{x}^{(t)}[r] \propto \exp\{\mathbf{r}^{(t)}[r]\}$. Now, after renormalization since $\mathbf{x} \in \Delta^d$,

$$\mathbf{x}^{(t)}[r] = \frac{\exp\{\lambda^{(t)} \mathbf{r}^{(t)}[r]\}}{\sum_{k=1}^d \exp\{\lambda^{(t)} \mathbf{r}^{(t)}[k]\}}.$$

□

Corollary 5.22. *We know that the optimal solution $\lambda^{(t)}$ can be written as,*

$$\lambda^{(t)} = \arg \max_{\lambda \in (0, \eta]} \left\{ \log \left(\sum_{k=1}^d e^{\lambda \mathbf{r}^{(t)}[k]} \right) + (\alpha - 1) \log \lambda \right\}, \quad (5.20)$$

where $\mathbf{r}^{(t)}[r] = \mathbf{U}^{(t)}[r] + \mathbf{u}^{(t-1)}[r]$.

Proof. Given Lemma 5.21, plugging Equation (5.19) into Equation (5.18) entails,

$$\begin{aligned} & \arg \max_{\lambda \in [0, \eta]} \\ & \left\{ \lambda \sum_{r=1}^d \mathbf{r}^{(t)}[r] \frac{e^{\lambda \mathbf{r}^{(t)}[r]}}{\sum_{k=1}^d e^{\lambda \mathbf{r}^{(t)}[k]}} + (\alpha - 1) \log \lambda - \sum_{r=1}^d \frac{e^{\lambda \mathbf{r}^{(t)}[r]}}{\sum_{r=1}^d e^{\lambda \mathbf{r}^{(t)}[k]}} \log \left(\frac{e^{\lambda \mathbf{r}^{(t)}[r]}}{\sum_{k=1}^d e^{\lambda \mathbf{r}^{(t)}[k]}} \right) \right\} \\ & = \arg \max_{\lambda \in [0, \eta]} \\ & \left\{ \lambda \sum_{r=1}^d \mathbf{r}^{(t)}[r] \frac{e^{\lambda \mathbf{r}^{(t)}[r]}}{\sum_{k=1}^d e^{\lambda \mathbf{r}^{(t)}[k]}} + (\alpha - 1) \log \lambda - \lambda \sum_{r=1}^d \mathbf{r}^{(t)}[r] \frac{e^{\lambda \mathbf{r}^{(t)}[r]}}{\sum_{k=1}^d e^{\lambda \mathbf{r}^{(t)}[k]}} + \frac{\sum_{r=1}^d e^{\lambda \mathbf{r}^{(t)}[r]}}{\sum_{k=1}^d e^{\lambda \mathbf{r}^{(t)}[k]}} \log \left(\sum_{k=1}^d e^{\lambda \mathbf{r}^{(t)}[k]} \right) \right\} \\ & = \arg \max_{\lambda \in [0, \eta]} \left\{ \log \left(\sum_{k=1}^d e^{\lambda \mathbf{r}^{(t)}[k]} \right) + (\alpha - 1) \log \lambda \right\}. \end{aligned}$$

□

5.3.4 Strong Spectral Properties of ψ

In this section, we prove strong convexity of ψ and hence we demonstrate that Bregman divergence of ψ is well-defined.

Theorem 5.5. *The regularizer $\psi(\mathbf{y})$ is strongly convex and furthermore,*

$$\nabla^2 \psi(\mathbf{y}) \succeq \frac{1}{2\eta} \begin{pmatrix} \frac{1}{\mathbf{y}[1] \cdot \sum_{k=1}^d \mathbf{y}[k]} & 0 & \cdots & 0 \\ 0 & \frac{1}{\mathbf{y}[2] \cdot \sum_{k=1}^d \mathbf{y}[k]} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \frac{1}{\mathbf{y}[d] \cdot \sum_{k=1}^d \mathbf{y}[k]} \end{pmatrix}.$$

Proof. The partial derivatives of the regularizer are

$$\eta \frac{\partial \psi}{\partial \mathbf{y}[i]} = -\frac{\alpha}{\sum_{k=1}^d \mathbf{y}[k]} - \frac{1}{(\sum_{k=1}^d \mathbf{y}[k])^2} \left(\sum_{k=1}^d \mathbf{y}[k] \log \mathbf{y}[k] \right) + \frac{1 + \log \mathbf{y}[i]}{\sum_{k=1}^d \mathbf{y}[k]}$$

For convenience, let $\Lambda(\mathbf{y}) := \sum_{k=1}^d \mathbf{y}[k]$ and $\mathbf{x}[i] := \frac{\mathbf{y}[i]}{\Lambda(\mathbf{y})}$, so that $\mathbf{x} \in \Delta^d$. Let $\mathbf{1}_{i=j}$ be the indicator function that is evaluated to one when $i = j$ and zero elsewhere. The second order derivatives are as follows,

$$\begin{aligned} \eta \frac{\partial^2 \psi}{\partial \mathbf{y}[i] \partial \mathbf{y}[j]} &= \frac{\alpha}{\Lambda(\mathbf{y})^2} + \frac{2}{\Lambda(\mathbf{y})^3} \left(\sum_{k=1}^d \mathbf{y}[k] \log \mathbf{y}[k] \right) - \frac{1 + \log \mathbf{y}[j]}{\Lambda(\mathbf{y})^2} - \frac{1 + \log \mathbf{y}[i]}{\Lambda(\mathbf{y})^2} + \frac{\mathbf{1}_{i=j}}{\mathbf{y}[i] \Lambda(\mathbf{y})} \\ &= \frac{\alpha}{\Lambda(\mathbf{y})^2} + \frac{2}{\Lambda(\mathbf{y})^2} \left(\sum_{k=1}^d \mathbf{x}[k] \log(\Lambda(\mathbf{y}) \mathbf{x}[k]) \right) - \frac{2}{\Lambda(\mathbf{y})^2} \\ &\quad - \frac{\log(\Lambda(\mathbf{y}) \mathbf{x}[i]) + \log(\Lambda(\mathbf{y}) \mathbf{x}[j])}{\Lambda(\mathbf{y})^2} + \frac{\mathbf{1}_{i=j}}{\mathbf{y}[i] \Lambda(\mathbf{y})} \\ &= \frac{\alpha - 2 + 2 \sum_{k=1}^d \mathbf{x}[k] \log \mathbf{x}[k]}{\Lambda(\mathbf{y})^2} + \frac{2 \log \Lambda(\mathbf{y})}{\Lambda(\mathbf{y})^2} - \frac{2 \log \Lambda(\mathbf{y})}{\Lambda(\mathbf{y})^2} \\ &\quad - \frac{\log \mathbf{x}[i] + \log \mathbf{x}[j]}{\Lambda(\mathbf{y})^2} + \frac{\mathbf{1}_{i=j}}{\mathbf{y}[i] \Lambda(\mathbf{y})} \\ &= \frac{\alpha - 2 + 2 \sum_{k=1}^d \mathbf{x}[k] \log \mathbf{x}[k]}{\Lambda(\mathbf{y})^2} - \frac{\log \mathbf{x}[i] + \log \mathbf{x}[j]}{\Lambda(\mathbf{y})^2} + \frac{\mathbf{1}_{i=j}}{\mathbf{y}[i] \Lambda(\mathbf{y})}. \end{aligned}$$

Let now $\alpha := 2 + 2 \log d + \alpha'$; we can therefore guarantee that for any vector $\mathbf{v} \in \mathbb{R}^d$,

$$\begin{aligned} &\eta \Lambda(\mathbf{y})^2 [\mathbf{v}^\top \nabla^2 \psi(\mathbf{y}) \mathbf{v}] \\ &\geq \alpha' \left(\sum_{k=1}^k \mathbf{v}[k] \right)^2 + \left(\sum_{k=1}^d \frac{\mathbf{v}[k]^2}{\mathbf{x}[i]} \right) + \left(\sum_{k=1}^d -2\mathbf{v}[k] \log \mathbf{x}[i] \right) \left(\sum_{k=1}^d \mathbf{v}[k] \right) \\ &= \left(\sum_{k=1}^d \frac{\mathbf{v}[k]^2}{2\mathbf{x}[i]} \right) + \alpha' \left(\sum_{k=1}^k \mathbf{v}[k] \right)^2 + \left(\sum_{k=1}^d \frac{\mathbf{v}[k]^2}{2\mathbf{x}[i]} \right) + \left(\sum_{k=1}^d -2\mathbf{v}[k] \log \mathbf{x}[i] \right) \left(\sum_{k=1}^d \mathbf{v}[k] \right) \\ &= \left(\sum_{k=1}^d \frac{\mathbf{v}[k]^2}{2\mathbf{x}[i]} \right) + \alpha' \left(\sum_{k=1}^k \mathbf{v}[k] \right)^2 \\ &\quad + \sum_{k=1}^d \left[\left(\sqrt{\frac{1}{2\mathbf{x}[k]}} \mathbf{v}[k] - \left(\sum_{j=1}^d \mathbf{v}[j] \right) \sqrt{2\mathbf{x}[k]} \log \mathbf{x}[k] \right)^2 - 2 \left(\sum_{j=1}^d \mathbf{v}[k] \right)^2 \mathbf{x}[k] \log^2 \mathbf{x}[k] \right] \\ &\geq \left(\sum_{k=1}^d \frac{\mathbf{v}[k]^2}{2\mathbf{x}[i]} \right) + \alpha' \left(\sum_{k=1}^k \mathbf{v}[k] \right)^2 - 2 \left(\sum_{k=1}^d \mathbf{v}[k] \right)^2 \left(\sum_{k=1}^d \mathbf{x}[k] \log^2 \mathbf{x}[k] \right) \\ &\geq \left(\sum_{k=1}^d \frac{\mathbf{v}[k]^2}{2\mathbf{x}[i]} \right) + \alpha' \left(\sum_{k=1}^k \mathbf{v}[k] \right)^2 - 2 \log^2 d \left(\sum_{k=1}^k \mathbf{v}[k] \right)^2. \end{aligned}$$

This shows that by setting any $\alpha' \geq 2 \log^2 d$, we can obtain

$$\eta \mathbf{v}^\top \nabla^2 \psi(\mathbf{y}) \mathbf{v} \geq \frac{1}{2} \mathbf{v}^\top \begin{pmatrix} \frac{1}{\mathbf{y}[1] \cdot \sum_{k=1}^d \mathbf{y}[k]} & 0 & \cdots & 0 \\ 0 & \frac{1}{\mathbf{y}[2] \cdot \sum_{k=1}^d \mathbf{y}[k]} & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \frac{1}{\mathbf{y}[d] \cdot \sum_{k=1}^d \mathbf{y}[k]} \end{pmatrix} \mathbf{v},$$

which concludes the proof. \square

Proposition 5.24. *The Bregman divergence $D_\psi(\cdot \| \cdot)$ induced by the regularizer $\psi(\cdot)$ has the following representation:*

$$\eta D_\psi(\mathbf{z} \| \mathbf{y}) = (\alpha - 1) D_{\log}(\Lambda(\mathbf{z}) \| \Lambda(\mathbf{y})) + \frac{\Lambda(\mathbf{z})}{\Lambda(\mathbf{y})} \text{KL}(\boldsymbol{\theta} \| \mathbf{x}) + \left(1 - \frac{\Lambda(\mathbf{z})}{\Lambda(\mathbf{y})}\right) \text{H}(\boldsymbol{\theta}) - \left(1 - \frac{\Lambda(\mathbf{z})}{\Lambda(\mathbf{y})}\right) \text{H}(\mathbf{x}),$$

where D_{\log} is the Bregman divergence induced by the log regularizer.

Proof. We know that

$$\begin{aligned} \eta \frac{\partial \psi}{\partial \mathbf{y}[i]} &= -\frac{\alpha}{\sum_{k=1}^d \mathbf{y}[k]} - \frac{1}{(\sum_{k=1}^d \mathbf{y}[k])^2} \left(\sum_{k=1}^d \mathbf{y}[k] \log \mathbf{y}[k] \right) + \frac{1 + \log \mathbf{y}[i]}{\sum_{k=1}^d \mathbf{y}[k]} \\ &= -\frac{\alpha}{\Lambda(\mathbf{y})} - \frac{1}{\Lambda(\mathbf{y})} \left(\sum_{k=1}^d \mathbf{x}[k] \log \mathbf{x}[k] \right) - \frac{\log \Lambda(\mathbf{y})}{\Lambda(\mathbf{y})} + \frac{1}{\Lambda(\mathbf{y})} + \frac{\log \mathbf{x}[i]}{\Lambda(\mathbf{y})} + \frac{\log \Lambda(\mathbf{y})}{\Lambda(\mathbf{y})} \\ &= -\frac{\alpha - 1}{\Lambda(\mathbf{y})} - \frac{1}{\Lambda(\mathbf{y})} \left(\sum_{k=1}^d \mathbf{x}[k] \log \mathbf{x}[k] \right) + \frac{\log \mathbf{x}[i]}{\Lambda(\mathbf{y})}. \end{aligned}$$

For the Bregman divergence, by definition, we get that

$$D_\psi(\mathbf{z} \| \mathbf{y}) = \psi(\mathbf{z}) - \psi(\mathbf{y}) - \sum_{k=1}^d \frac{\partial \psi(\mathbf{y})}{\partial \mathbf{y}[k]} (\mathbf{z}[k] - \mathbf{y}[k]).$$

Hence,

$$\begin{aligned} \eta D_\psi(\mathbf{z} \| \mathbf{y}) &= \left(-(\alpha - 1) \log \Lambda(\mathbf{z}) + \sum_{k=1}^d \boldsymbol{\theta}[k] \log \boldsymbol{\theta}[k] \right) - \left(-(\alpha - 1) \log \Lambda(\mathbf{y}) + \sum_{k=1}^d \mathbf{x}[k] \log \mathbf{x}[k] \right) \\ &\quad - \left(-\frac{\alpha - 1}{\Lambda(\mathbf{y})} - \frac{1}{\Lambda(\mathbf{y})} \left(\sum_{k=1}^d \mathbf{x}[k] \log \mathbf{x}[k] \right) \right) (\Lambda(\mathbf{z}) - \Lambda(\mathbf{y})) - \left(\frac{1}{\Lambda(\mathbf{y})} \sum_{k=1}^d \log \mathbf{x}[k] (\Lambda(\mathbf{z}) \boldsymbol{\theta}[k] - \Lambda(\mathbf{y}) \mathbf{x}[k]) \right) \\ &= -(\alpha - 1) \log \frac{\Lambda(\mathbf{z})}{\Lambda(\mathbf{y})} + (\alpha - 1) \left(\frac{\Lambda(\mathbf{z})}{\Lambda(\mathbf{y})} - 1 \right) + \left(\frac{\Lambda(\mathbf{z})}{\Lambda(\mathbf{y})} - 1 \right) \left(\sum_{k=1}^d \mathbf{x}[k] \log \mathbf{x}[k] \right) + \sum_{k=1}^d \boldsymbol{\theta}[k] \log \boldsymbol{\theta}[k] \\ &\quad - \frac{\Lambda(\mathbf{z})}{\Lambda(\mathbf{y})} \left(\sum_{k=1}^d \boldsymbol{\theta}[k] \log \mathbf{x}[k] \right) \\ &= (\alpha - 1) D_{\log}(\Lambda(\mathbf{z}) \| \Lambda(\mathbf{y})) + \frac{\Lambda(\mathbf{z})}{\Lambda(\mathbf{y})} \text{KL}(\boldsymbol{\theta} \| \mathbf{x}) + \left(1 - \frac{\Lambda(\mathbf{z})}{\Lambda(\mathbf{y})}\right) \text{H}(\boldsymbol{\theta}) - \left(1 - \frac{\Lambda(\mathbf{z})}{\Lambda(\mathbf{y})}\right) \text{H}(\mathbf{x}), \end{aligned}$$

where D_{\log} is the Bregman divergence induced by log regularizer. Thus, the proof is concluded. \square

Proposition 5.25. *The Bregman divergence $D_\psi(\cdot \| \cdot)$ induced by the regularizer $\psi(\cdot)$ is lower bounded by a term proportional to the ℓ_1 norm on the lifted simplex $(0, 1]^{\Delta^d}$,*

$$D_\psi(\mathbf{y} \| \mathbf{z}) \geq \frac{1}{2\eta} \|\mathbf{y} - \mathbf{z}\|_1^2.$$

Proof. We first show that $\psi(\mathbf{y})$ is strongly convex w.r.t. ℓ_1 norm. By Theorem 5.5, for any vector $\boldsymbol{\nu} \in \mathbb{R}^d$,

$$\begin{aligned} \boldsymbol{\nu}^\top \nabla^2 \psi(\mathbf{y}) \boldsymbol{\nu} &\geq \frac{1}{2\eta} \sum_{i=1}^d \frac{\boldsymbol{\nu}[i]^2}{\mathbf{y}[i] \cdot \sum_{k=1}^d \mathbf{y}[k]} \\ &\geq \frac{1}{2\eta} \sum_{i=1}^d \frac{\boldsymbol{\nu}[i]^2}{\mathbf{y}[i]} \end{aligned} \quad (5.21)$$

$$\begin{aligned} &\geq \frac{1}{2\eta} \left(\sum_{k=1}^d \mathbf{y}[k] \right) \cdot \sum_{i=1}^d \frac{\boldsymbol{\nu}[i]^2}{\mathbf{y}[i]} \\ &\geq \frac{1}{2\eta} \left(\sum_{i=1}^d \boldsymbol{\nu}[i] \right)^2 = \frac{1}{2\eta} \|\boldsymbol{\nu}\|_1^2, \end{aligned} \quad (5.22)$$

where Equations (5.21) and (5.22) follow since $\sum_{k=1}^d \mathbf{y}[k] \leq 1$ and the last line is derived by Cauchy–Schwarz. Next by definition of the Bregman divergence and strong convexity of ψ ,

$$\begin{aligned} D_\psi(\mathbf{y} \| \mathbf{z}) &= \psi(\mathbf{y}) - \psi(\mathbf{z}) - \langle \nabla \psi(\mathbf{z}), \mathbf{y} - \mathbf{z} \rangle \\ &\geq \frac{1}{2\eta} \|\mathbf{y} - \mathbf{z}\|_1^2. \end{aligned}$$

and the proof is completed. \square

Proposition 5.26. *Under the multiplicative stability assumption of $\omega := \frac{\Lambda(\mathbf{z})}{\Lambda(\mathbf{y})} \in [1 - \epsilon, 1 + \epsilon]$ for a constant $\epsilon \in (0, \frac{2}{5})$, the Bregman divergence $D_\psi(\cdot \| \cdot)$ induced by the regularizer $\psi(\cdot)$ is lower bounded by a term proportional to the ℓ_1 norm on the action simplex Δ^d :*

$$\eta D_\psi(\mathbf{z} \| \mathbf{y}) \geq \frac{1}{4} (1 - \epsilon) \|\boldsymbol{\theta} - \mathbf{x}\|_1^2.$$

Proof. By Proposition 5.24 and the multiplicative stability assumption $\omega := \frac{\Lambda(\mathbf{z})}{\Lambda(\mathbf{y})} \in [1 - \epsilon, 1 + \epsilon]$, we infer that

$$\begin{aligned} \eta D_\psi(\mathbf{z} \| \mathbf{y}) &\geq \beta \log^2 d \left(\log \left(\frac{1}{\omega} \right) + \omega - 1 \right) + (1 - \omega) (\mathbb{H}(\boldsymbol{\theta}) - \mathbb{H}(\mathbf{x})) + \omega \text{KL}(\boldsymbol{\theta} \| \mathbf{x}) \\ &\geq \frac{1}{4} \beta \log^2 d \left(1 - \frac{1}{\omega} \right)^2 + (\omega - 1) \log d \sqrt{2 \text{KL}(\boldsymbol{\theta} \| \mathbf{x})} + \frac{\omega^2}{\beta} \text{KL}(\boldsymbol{\theta} \| \mathbf{x}) + \left(\omega - \frac{\omega^2}{\beta} \right) \text{KL}(\boldsymbol{\theta} \| \mathbf{x}) \\ &\geq \left(\sqrt{\beta} \log d \left(1 - \frac{1}{\omega} \right) + \frac{\omega}{\sqrt{\beta}} \sqrt{\text{KL}(\boldsymbol{\theta} \| \mathbf{x})} \right)^2 + \frac{\omega}{2} \text{KL}(\boldsymbol{\theta} \| \mathbf{x}) \\ &\geq \frac{1}{4} (1 - \epsilon) \|\boldsymbol{\theta} - \mathbf{x}\|_1^2, \end{aligned} \quad (5.23)$$

where Equation (5.23) follows from the inequality $|\mathbb{H}(\mathbf{p}) - \mathbb{H}(\mathbf{q})| \leq (\log d)\sqrt{2\text{KL}(\mathbf{p}\|\mathbf{q})}$ for discrete probability distributions \mathbf{p} and \mathbf{q} with support size d (Lemma 5.12), and the fact that $\omega - \frac{\omega^2}{\beta} \geq \frac{\epsilon}{2}$ for all choices of $\omega \in [1 - \epsilon, 1 + \epsilon]$ and $\beta \geq 20$. The last line follows from Pinsker's inequality, i.e., $\|\mathbf{p} - \mathbf{q}\|_1^2 \leq 2\text{KL}(\mathbf{p}\|\mathbf{q})$ for discrete random variables \mathbf{p} and \mathbf{q} (Lemma 5.11). \square

5.3.5 Positive Regret

To analyze DLRC-OMWU equivalently as shown in Section 5.3.3, we start the analysis by a closer look at Equation (5.16). To analyze the Reg^T , we first study the nonnegative regret,

$$\tilde{\text{Reg}}^{(T)} := \max_{\mathbf{y}^* \in [0,1]^{\Delta^d}} \sum_{t=1}^T \langle \mathbf{u}^{(t)}, \mathbf{y}^* - \mathbf{y}^{(t)} \rangle.$$

Proposition 5.7. *For any time horizon $T \in \mathbb{N}$, we have that $\tilde{\text{Reg}}^{(T)} = \max\{0, \text{Reg}^{(T)}\}$. As a result, $\tilde{\text{Reg}}^{(T)} \geq 0$ and $\tilde{\text{Reg}}^{(T)} \geq \text{Reg}^{(T)}$.*

Proof. By definition of the reward signal $\mathbf{u}^{(t)} = \boldsymbol{\nu}^{(t)} - \langle \boldsymbol{\nu}^{(t)}, \mathbf{x}^{(t)} \rangle \mathbf{1}_d$ and induced action $\mathbf{x}^{(t)} = \frac{\mathbf{y}^{(t)}}{\langle \mathbf{y}^{(t)}, \mathbf{1} \rangle}$, for the regret we infer

$$\begin{aligned} \tilde{\text{Reg}}^{(T)} &= \max_{\mathbf{y}^* \in [0,1]^{\Delta^d}} \sum_{t=1}^T \langle \mathbf{u}^{(t)}, \mathbf{y}^* - \mathbf{y}^{(t)} \rangle \\ &= \max_{\mathbf{y}^* \in [0,1]^{\Delta^d}} \sum_{t=1}^T \langle \boldsymbol{\nu}^{(t)} - \langle \boldsymbol{\nu}^{(t)}, \frac{\mathbf{y}^{(t)}}{\langle \mathbf{y}^{(t)}, \mathbf{1} \rangle} \rangle \mathbf{1}_d, \mathbf{y}^* - \mathbf{y}^{(t)} \rangle \\ &= \max_{\mathbf{y}^* \in [0,1]^{\Delta^d}} \sum_{t=1}^T \langle \boldsymbol{\nu}^{(t)}, \mathbf{y}^* \rangle - \left\langle \langle \boldsymbol{\nu}^{(t)}, \frac{\mathbf{y}^{(t)}}{\langle \mathbf{y}^{(t)}, \mathbf{1} \rangle} \rangle \mathbf{1}_d, \mathbf{y}^* \right\rangle \\ &\geq \max_{\mathbf{y}^* \in \Delta^d} \sum_{t=1}^T \langle \boldsymbol{\nu}^{(t)}, \mathbf{y}^* \rangle - \langle \boldsymbol{\nu}^{(t)}, \mathbf{x}^{(t)} \rangle \cdot \langle \mathbf{1}_d, \mathbf{y}^* \rangle \\ &\geq \max_{\mathbf{y}^* \in \Delta^d} \sum_{t=1}^T \langle \boldsymbol{\nu}^{(t)}, \mathbf{y}^* \rangle - \langle \boldsymbol{\nu}^{(t)}, \mathbf{x}^{(t)} \rangle \\ &= \text{Reg}^{(T)}, \end{aligned} \tag{5.24}$$

where Equation (5.24) follows because of orthogonality $\mathbf{u}^{(t)} \perp \mathbf{y}^{(t)}$. On the other hand, it is clear that $\tilde{\text{Reg}}^{(T)} \geq 0$ by choosing 0 as the comparator. \square

This proposition is important as it implies that any RVU bounds on $\tilde{\text{Reg}}^{(T)}$ directly translate into nonnegative RVU bounds on Reg^T .

5.3.6 Proofs for RVU Bounds (Section 5.2)

We state and prove the following standard lemma from Optimistic FTRL analysis.

Lemma 5.8. For any $\mathbf{y} \in \Omega$, the sequences $\{\mathbf{y}^{(t)}\}_{t=1}^T$ generated by Equation (5.7), it holds that

$$\begin{aligned} \sum_{t=1}^T \langle \mathbf{y} - \mathbf{y}^{(t)}, \mathbf{u}^{(t)} \rangle &\leq \psi(\mathbf{y}) - \psi(\mathbf{y}^{(1)}) + \sum_{t=1}^T \langle \mathbf{z}^{(t+1)} - \mathbf{y}^{(t)}, \mathbf{u}^{(t)} - \mathbf{u}^{(t-1)} \rangle \\ &\quad - \sum_{t=1}^T \left(D_\psi(\mathbf{y}^{(t)} \parallel \mathbf{z}^{(t)}) + D_\psi(\mathbf{z}^{(t+1)} \parallel \mathbf{y}^{(t)}) \right). \end{aligned}$$

Proof. By Lemma 5.29, and optimality of $\mathbf{z}^{(t)}$,

$$\begin{aligned} G_t(\mathbf{z}^{(t)}) &\leq G_t(\mathbf{y}^{(t)}) - D_\psi(\mathbf{y}^{(t)} \parallel \mathbf{z}^{(t)}) \\ &\leq F_t(\mathbf{y}^{(t)}) + \langle \mathbf{y}^{(t)}, \mathbf{u}^{(t-1)} \rangle - D_\psi(\mathbf{y}^{(t)} \parallel \mathbf{z}^{(t)}). \end{aligned}$$

Similarly by optimality of $\mathbf{y}^{(t)}$,

$$\begin{aligned} F_t(\mathbf{y}^{(t)}) &\leq F_t(\mathbf{z}^{(t+1)}) - D_\psi(\mathbf{z}^{(t+1)} \parallel \mathbf{y}^{(t)}) \\ &\leq G_{t+1}(\mathbf{z}^{(t+1)}) + \langle \mathbf{z}^{(t+1)}, \mathbf{u}^{(t)} - \mathbf{u}^{(t-1)} \rangle - D_\psi(\mathbf{z}^{(t+1)} \parallel \mathbf{y}^{(t)}). \end{aligned}$$

By merging the inequalities and aggregating over all t , we derive

$$\begin{aligned} G_1(\mathbf{z}^{(1)}) &\leq G_{T+1}(\mathbf{z}^{(T+1)}) + \sum_{t=1}^T (\langle \mathbf{y}^{(t)}, \mathbf{u}^{(t)} \rangle + \langle \mathbf{z}^{(t+1)} - \mathbf{y}^{(t)}, \mathbf{u}^{(t)} - \mathbf{u}^{(t-1)} \rangle) \\ &\quad - \sum_{t=1}^T (D_\psi(\mathbf{y}^{(t)} \parallel \mathbf{z}^{(t)}) + D_\psi(\mathbf{z}^{(t+1)} \parallel \mathbf{y}^{(t)}). \end{aligned}$$

Plugging in $G_{T+1}(\mathbf{z}^{(T+1)}) \leq -\langle \mathbf{y}, \mathbf{U}^{(T+1)} \rangle + \psi(\mathbf{y})$ and $G_1(\mathbf{z}^{(1)}) = \psi(\mathbf{y}^{(1)})$, entails the proof,

$$\begin{aligned} \sum_{t=1}^T \langle \mathbf{y} - \mathbf{y}^{(t)}, \mathbf{u}^{(t)} \rangle &\leq \psi(\mathbf{y}) - \psi(\mathbf{y}^{(1)}) + \sum_{t=1}^T \langle \mathbf{z}^{(t+1)} - \mathbf{y}^{(t)}, \mathbf{u}^{(t)} - \mathbf{u}^{(t-1)} \rangle \\ &\quad - \sum_{t=1}^T \left(D_\psi(\mathbf{y}^{(t)} \parallel \mathbf{z}^{(t)}) + D_\psi(\mathbf{z}^{(t+1)} \parallel \mathbf{y}^{(t)}) \right). \end{aligned}$$

□

Lemma 5.29. Given any convex function $F : \Omega \rightarrow \mathbb{R}$ defined on the compact set Ω , the minimizer $\mathbf{z}^* = \arg \min_{\mathbf{z} \in \Omega} F(\mathbf{z})$ satisfies

$$F(\mathbf{z}^*) \leq F(\mathbf{z}) - D_F(\mathbf{z} \parallel \mathbf{z}^*) \quad \forall \mathbf{z} \in \Omega,$$

where D_F is the Bregman divergence induced by function F .

Proof. By definition,

$$F(\mathbf{z}^*) = F(\mathbf{z}) - \langle \nabla F(\mathbf{z}^*), \mathbf{z} - \mathbf{z}^* \rangle - D_F(\mathbf{z} \parallel \mathbf{z}^*) \leq F(\mathbf{z}) - D_F(\mathbf{z} \parallel \mathbf{z}^*),$$

which follows by the first order optimality condition of \mathbf{z}^* .

□

Lemma 5.30. *If $\eta \leq \frac{1}{50}$ and β is large enough ($\beta \geq 70$), then*

$$\sum_{t=1}^T \left(D_\psi(\mathbf{y}^{(t)} \parallel \mathbf{z}^{(t)}) + D_\psi(\mathbf{z}^{(t+1)} \parallel \mathbf{y}^{(t)}) \right) \geq \sum_{t=1}^T \frac{1}{2\eta} (\|\mathbf{y}^{(t)} - \mathbf{z}^{(t)}\|_1^2 + \|\mathbf{z}^{(t+1)} - \mathbf{y}^{(t)}\|_1^2).$$

Proof. By multiple usage of Proposition 5.25,

$$D_\psi(\mathbf{y}^{(t)} \parallel \mathbf{z}^{(t)}) + D_\psi(\mathbf{z}^{(t+1)} \parallel \mathbf{y}^{(t)}) \geq \frac{1}{2\eta} \|\mathbf{y}^{(t)} - \mathbf{z}^{(t)}\|_1^2 + \frac{1}{2\eta} \|\mathbf{z}^{(t+1)} - \mathbf{y}^{(t)}\|_1^2,$$

and summing over $t \in [T]$ concludes the proof. \square

Lemma 5.31. *If $\eta \leq \frac{1}{50}$ and β is large enough ($\beta \geq 70$), then*

$$\sum_{t=1}^T \left(D_\psi(\mathbf{y}^{(t)} \parallel \mathbf{z}^{(t)}) + D_\psi(\mathbf{z}^{(t+1)} \parallel \mathbf{y}^{(t)}) \right) \geq \sum_{t=1}^{T-1} \frac{1}{10\eta} (\|\mathbf{x}^{(t+1)} - \boldsymbol{\theta}^{(t+1)}\|_1^2 + \|\boldsymbol{\theta}^{(t+1)} - \mathbf{x}^{(t)}\|_1^2).$$

Proof. By Theorem 3.5 in Section 3.1, we have multiplicative stability in the learning rate as the solution to the Dynamic Learning Rate Control Problem, i.e., $\omega := \frac{\Lambda(\mathbf{z})}{\Lambda(\mathbf{y})} \in [1 - \epsilon, 1 + \epsilon]$ for $\epsilon = \frac{2}{5}$. Consequently, by Proposition 5.26, we infer that

$$\eta D_\psi(\mathbf{z}, \mathbf{y}) \geq \frac{1}{4} (1 - \epsilon) \|\boldsymbol{\theta} - \mathbf{x}\|_1^2,$$

Next, by setting $\mathbf{z} := \mathbf{z}^{(t+1)}$ and $\mathbf{y} := \mathbf{y}^{(t)}$, we obtain

$$D_\psi(\mathbf{z}^{(t+1)}, \mathbf{y}^{(t)}) \geq \frac{3}{20\eta} \|\boldsymbol{\theta}^{(t+1)} - \mathbf{x}^{(t)}\|_1^2 > \frac{1}{10\eta} \|\boldsymbol{\theta}^{(t+1)} - \mathbf{x}^{(t)}\|_1^2.$$

Similarly,

$$D_\psi(\mathbf{y}^{t+1}, \mathbf{z}^{t+1}) \geq \frac{3}{20\eta} \|\mathbf{x}^{(t+1)} - \boldsymbol{\theta}^{(t+1)}\|_1^2 > \frac{1}{10\eta} \|\mathbf{x}^{(t+1)} - \boldsymbol{\theta}^{(t+1)}\|_1^2.$$

To conclude,

$$\begin{aligned} \sum_{t=1}^T \left(D_\psi(\mathbf{y}^{(t)} \parallel \mathbf{z}^{(t)}) + D_\psi(\mathbf{z}^{(t+1)} \parallel \mathbf{y}^{(t)}) \right) &\geq \sum_{t=1}^{T-1} \left(D_\psi(\mathbf{y}^{(t+1)} \parallel \mathbf{z}^{(t+1)}) + D_\psi(\mathbf{z}^{(t+1)} \parallel \mathbf{y}^{(t)}) \right) \\ &\geq \sum_{t=1}^{T-1} \frac{1}{10\eta} (\|\mathbf{x}^{(t+1)} - \boldsymbol{\theta}^{(t+1)}\|_1^2 + \|\boldsymbol{\theta}^{(t+1)} - \mathbf{x}^{(t)}\|_1^2). \end{aligned}$$

\square

Lemma 5.32. *Assuming that $\|\boldsymbol{\nu}^{(t)}\|_\infty \leq 1$ is satisfied for all $t \in [T]$, we have*

$$\|\mathbf{u}^{(t)} - \mathbf{u}^{(t-1)}\|_\infty^2 \leq 6\|\boldsymbol{\nu}^{(t)} - \boldsymbol{\nu}^{(t-1)}\|_\infty^2 + 4\|\mathbf{x}^{(t)} - \mathbf{x}^{(t-1)}\|_1^2.$$

Proof.

$$\begin{aligned} \|\mathbf{u}^{(t)} - \mathbf{u}^{(t-1)}\|_\infty^2 &= \|(\boldsymbol{\nu}^{(t)} - \langle \boldsymbol{\nu}^{(t)}, \mathbf{x}^{(t)} \rangle \mathbf{1}) - (\boldsymbol{\nu}^{(t-1)} - \langle \boldsymbol{\nu}^{(t-1)}, \mathbf{x}^{(t-1)} \rangle \mathbf{1})\|_\infty^2 \\ &\leq \left(\|\boldsymbol{\nu}^{(t)} - \boldsymbol{\nu}^{(t-1)}\|_\infty + \|\langle \boldsymbol{\nu}^{(t)}, \mathbf{x}^{(t)} \rangle \mathbf{1} - \langle \boldsymbol{\nu}^{(t-1)}, \mathbf{x}^{(t-1)} \rangle \mathbf{1}\|_\infty \right)^2 \end{aligned} \quad (5.25)$$

$$\begin{aligned} &= \left(\|\boldsymbol{\nu}^{(t)} - \boldsymbol{\nu}^{(t-1)}\|_\infty + |\langle \boldsymbol{\nu}^{(t)}, \mathbf{x}^{(t)} \rangle - \langle \boldsymbol{\nu}^{(t-1)}, \mathbf{x}^{(t-1)} \rangle| \right)^2 \\ &\leq 2\|\boldsymbol{\nu}^{(t)} - \boldsymbol{\nu}^{(t-1)}\|_\infty^2 + 2|\langle \boldsymbol{\nu}^{(t)}, \mathbf{x}^{(t)} \rangle - \langle \boldsymbol{\nu}^{(t-1)}, \mathbf{x}^{(t-1)} \rangle|^2 \end{aligned} \quad (5.26)$$

$$\begin{aligned} &\leq 2\|\boldsymbol{\nu}^{(t)} - \boldsymbol{\nu}^{(t-1)}\|_\infty^2 + 2|(\langle \boldsymbol{\nu}^{(t)}, \mathbf{x}^{(t)} \rangle - \langle \boldsymbol{\nu}^{(t)}, \mathbf{x}^{(t-1)} \rangle) + (\langle \boldsymbol{\nu}^{(t)}, \mathbf{x}^{(t-1)} \rangle - \langle \boldsymbol{\nu}^{(t-1)}, \mathbf{x}^{(t-1)} \rangle)|^2 \\ &\leq 2\|\boldsymbol{\nu}^{(t)} - \boldsymbol{\nu}^{(t-1)}\|_\infty^2 + 4|\langle \boldsymbol{\nu}^{(t)}, \mathbf{x}^{(t)} - \mathbf{x}^{(t-1)} \rangle|^2 + 4|\langle \boldsymbol{\nu}^{(t)} - \boldsymbol{\nu}^{(t-1)}, \mathbf{x}^{(t-1)} \rangle|^2 \end{aligned} \quad (5.27)$$

$$\begin{aligned} &\leq 2\|\boldsymbol{\nu}^{(t)} - \boldsymbol{\nu}^{(t-1)}\|_\infty^2 + 4\|\mathbf{x}^{(t)} - \mathbf{x}^{(t-1)}\|_1^2 + 4\|\boldsymbol{\nu}^{(t)} - \boldsymbol{\nu}^{(t-1)}\|_\infty^2 \\ &= 6\|\boldsymbol{\nu}^{(t)} - \boldsymbol{\nu}^{(t-1)}\|_\infty^2 + 4\|\mathbf{x}^{(t)} - \mathbf{x}^{(t-1)}\|_1^2, \end{aligned} \quad (5.28)$$

where Equation (5.25) uses the triangle inequality, Equations (5.26) and (5.27) apply Young's inequality, and Equation (5.28) utilizes Hölder's inequality. \square

Theorem 5.1 (RVU bound of DLRC-OMWU). *Consider the cumulative regret $\tilde{\text{Reg}}^{(T)}$ accrued by the internal OFTRL algorithm up to time T . Assuming that $\|\boldsymbol{\nu}^{(t)}\|_\infty \leq 1$ is satisfied for all $t \in [T]$, it follows that for any time $T \in \mathbb{N}$ and any learning rate $\eta \leq \frac{1}{50}$ and β high enough ($\beta \geq 70$),*

$$\tilde{\text{Reg}}^{(T)} \leq 3 + \frac{\alpha \log T + \log d}{\eta} + 6\eta \sum_{t=1}^{T-1} \|\boldsymbol{\nu}^{(t)} - \boldsymbol{\nu}^{(t-1)}\|_\infty^2 - \frac{1}{24\eta} \sum_{t=1}^{T-1} \|\mathbf{x}^{(t+1)} - \mathbf{x}^{(t)}\|_1^2.$$

Proof. For any choice of comparator $\mathbf{y} \in \Omega$, let $\mathbf{y}' = \frac{T-1}{T}\mathbf{y} + \frac{1}{T}\mathbf{y}^{(1)} \in \Omega$. Recall that $\mathbf{y}^{(1)} = \arg \min_{\mathbf{y} \in \Omega} F_1(\mathbf{y}) = \arg \min_{\mathbf{y} \in \Omega} \psi(\mathbf{y})$. By straightforward calculations,

$$\begin{aligned} \sum_{t=1}^T \langle \mathbf{y} - \mathbf{y}^{(t)}, \mathbf{u}^{(t)} \rangle &= \sum_{t=1}^T \langle \mathbf{y} - \mathbf{y}', \mathbf{u}^{(t)} \rangle + \sum_{t=1}^T \langle \mathbf{y}' - \mathbf{y}^{(t)}, \mathbf{u}^{(t)} \rangle \\ &= \frac{1}{T} \sum_{t=1}^T \langle \mathbf{y} - \mathbf{y}^{(1)}, \mathbf{u}^{(t)} \rangle + \sum_{t=1}^T \langle \mathbf{y}' - \mathbf{y}^{(t)}, \mathbf{u}^{(t)} \rangle \\ &\leq 2 + \sum_{t=1}^T \langle \mathbf{y}' - \mathbf{y}^{(t)}, \mathbf{u}^{(t)} \rangle, \end{aligned}$$

where the last line follows because of Hölder's inequality and $\|\mathbf{u}^{(t)}\|_\infty \leq 1$. In turn, we need to upperbound the $\sum_{t=1}^T \langle \mathbf{y}' - \mathbf{y}^{(t)}, \mathbf{u}^{(t)} \rangle$ term. By Lemma 5.8,

$$\begin{aligned} \sum_{t=1}^T \langle \mathbf{y}' - \mathbf{y}^{(t)}, \mathbf{u}^{(t)} \rangle &\leq \underbrace{\psi(\mathbf{y}') - \psi(\mathbf{y}^{(1)})}_{\text{(I)}} + \underbrace{\sum_{t=1}^T \langle \mathbf{z}^{(t+1)} - \mathbf{y}^{(t)}, \mathbf{u}^{(t)} - \mathbf{u}^{(t-1)} \rangle}_{\text{(II)}} \\ &\quad - \underbrace{\sum_{t=1}^T (D_\psi(\mathbf{y}^{(t)} \parallel \mathbf{z}^{(t)}) + D_\psi(\mathbf{z}^{(t+1)} \parallel \mathbf{y}^{(t)}))}_{\text{(III)}}. \end{aligned}$$

For the term (I), after some calculations,

$$\begin{aligned}
(\text{I}) &= -\frac{1}{\eta}\alpha \log \left(\sum_{k=1}^d \mathbf{y}'[k] \right) + \frac{1}{\eta} \frac{1}{\sum_{k=1}^d \mathbf{y}'[k]} \sum_{k=1}^d \mathbf{y}'[k] \log \mathbf{y}'[k] \\
&\quad + \frac{1}{\eta}\alpha \log \left(\sum_{k=1}^d \mathbf{y}^{(1)}[k] \right) - \frac{1}{\eta} \frac{1}{\sum_{k=1}^d \mathbf{y}^{(1)}[k]} \sum_{k=1}^d \mathbf{y}^{(1)}[k] \log \mathbf{y}^{(1)}[k] \\
&\leq -\frac{1}{\eta}\alpha \log \left(\frac{\sum_{k=1}^d \mathbf{y}'[k]}{\sum_{k=1}^d \mathbf{y}^{(1)}[k]} \right) + \frac{1}{\eta} \frac{1}{\sum_{k=1}^d \mathbf{y}'[k]} \sum_{k=1}^d \mathbf{y}'[k] \log \mathbf{y}'[k] \\
&\leq \frac{1}{\eta}(\alpha \log T + \log d)
\end{aligned}$$

By Hölder's and Young's inequalities, we entail that term (II) is upper bounded by

$$\begin{aligned}
(\text{II}) &\leq \sum_{t=1}^T \langle \mathbf{z}^{(t+1)} - \mathbf{y}^{(t)}, \mathbf{u}^{(t)} - \mathbf{u}^{(t-1)} \rangle \\
&\leq \sum_{t=1}^T \|\mathbf{z}^{(t+1)} - \mathbf{y}^{(t)}\|_1 \cdot \|\mathbf{u}^{(t)} - \mathbf{u}^{(t-1)}\|_\infty \\
&\leq \sum_{t=1}^T \left(\frac{1}{4\eta} \|\mathbf{z}^{(t+1)} - \mathbf{y}^{(t)}\|_1^2 + \eta \|\mathbf{u}^{(t)} - \mathbf{u}^{(t-1)}\|_\infty^2 \right) \\
&\leq \sum_{t=1}^T \left(\frac{1}{4\eta} \|\mathbf{z}^{(t+1)} - \mathbf{y}^{(t)}\|_1^2 + 6\eta \|\boldsymbol{\nu}^{(t)} - \boldsymbol{\nu}^{(t-1)}\|_\infty^2 + 4\eta \|\mathbf{x}^{(t)} - \mathbf{x}^{(t-1)}\|_1^2 \right),
\end{aligned}$$

where we used Lemma 5.32. In turn, for term (III),

$$\begin{aligned}
(\text{III}) &= -\frac{1}{2} \sum_{t=1}^T \left(D_\psi(\mathbf{y}^{(t)} \parallel \mathbf{z}^{(t)}) + D_\psi(\mathbf{z}^{(t+1)} \parallel \mathbf{y}^{(t)}) \right) - \frac{1}{2} \sum_{t=1}^T \left(D_\psi(\mathbf{y}^{(t)} \parallel \mathbf{z}^{(t)}) + D_\psi(\mathbf{z}^{(t+1)} \parallel \mathbf{y}^{(t)}) \right) \\
&\leq -\sum_{t=1}^T \frac{1}{4\eta} (\|\mathbf{y}^{(t)} - \mathbf{z}^{(t)}\|_1^2 + \|\mathbf{z}^{(t+1)} - \mathbf{y}^{(t)}\|_1^2) - \sum_{t=1}^{T-1} \frac{1}{20\eta} (\|\mathbf{x}^{(t+1)} - \boldsymbol{\theta}^{(t+1)}\|_1^2 + \|\boldsymbol{\theta}^{(t+1)} - \mathbf{x}^{(t)}\|_1^2) \\
&\leq -\sum_{t=1}^T \frac{1}{4\eta} (\|\mathbf{y}^{(t)} - \mathbf{z}^{(t)}\|_1^2 + \|\mathbf{z}^{(t+1)} - \mathbf{y}^{(t)}\|_1^2) - \sum_{t=1}^{T-1} \frac{1}{20\eta} \|\mathbf{x}^{(t+1)} - \mathbf{x}^{(t)}\|_1^2,
\end{aligned} \tag{5.29}$$

where Equation (5.29) is obtained by applying Lemmas 5.30 and 5.31 and the last line is yielded by triangle inequality. Assembling the complete picture,

$$\begin{aligned}
(\text{II}) + (\text{III}) &\leq \sum_{t=1}^T 6\eta \|\boldsymbol{\nu}^{(t)} - \boldsymbol{\nu}^{(t-1)}\|_\infty^2 + \sum_{t=1}^T 4\eta \|\mathbf{x}^{(t)} - \mathbf{x}^{(t-1)}\|_1^2 - \sum_{t=1}^{T-1} \frac{1}{20\eta} \|\mathbf{x}^{(t+1)} - \mathbf{x}^{(t)}\|_1^2 \\
&\leq 24\eta + \sum_{t=1}^{T-1} 6\eta \|\boldsymbol{\nu}^{(t)} - \boldsymbol{\nu}^{(t-1)}\|_\infty^2 + 16\eta - \sum_{t=1}^{T-1} \left(\frac{1}{20\eta} - 4\eta \right) \|\mathbf{x}^{(t+1)} - \mathbf{x}^{(t)}\|_1^2 \\
&\leq 1 + \sum_{t=1}^{T-1} 6\eta \|\boldsymbol{\nu}^{(t)} - \boldsymbol{\nu}^{(t-1)}\|_\infty^2 - \sum_{t=1}^{T-1} \frac{1}{24\eta} \|\mathbf{x}^{(t+1)} - \mathbf{x}^{(t)}\|_1^2.
\end{aligned}$$

and this completes the proof. \square

5.3.7 Proofs for Main Results

Theorem 5.2 (Bound on total path length). *Under Assumption 1, if all the players follow DLRC-OMWU algorithm with learning rate $\eta \leq \min\{\frac{1}{50}, \frac{1}{12\sqrt{2}Ln}\}$, then*

$$\sum_{i=1}^n \sum_{t=1}^{T-1} \|\mathbf{x}_i^{(t+1)} - \mathbf{x}_i^{(t)}\|_1^2 \leq 144n\eta + 48n(\alpha \log T + \log d).$$

Proof. Assumption 1 implies that,

$$\begin{aligned} \|\boldsymbol{\nu}_i^{(t+1)} - \boldsymbol{\nu}_i^{(t)}\|_\infty^2 &\leq L^2 \left(\sum_{i=1}^n \|\mathbf{x}_i^{(t+1)} - \mathbf{x}_i^{(t)}\|_1 \right)^2 \\ &\leq L^2 n \sum_{i=1}^n \|\mathbf{x}_i^{(t+1)} - \mathbf{x}_i^{(t)}\|_1^2, \end{aligned}$$

where the last line is obtained by Jensen's inequality. Next, we combine this result with the RUV bound on $\tilde{\text{Reg}}^{(T)}$ for the i th player (Theorem 5.1),

$$\begin{aligned} \tilde{\text{Reg}}_i^{(T)} &\leq 3 + \frac{\alpha \log T + \log d}{\eta} + 6\eta \sum_{t=1}^{T-1} \|\boldsymbol{\nu}_i^{(t+1)} - \boldsymbol{\nu}_i^{(t)}\|_\infty^2 - \frac{1}{24\eta} \sum_{t=1}^{T-1} \|\mathbf{x}_i^{(t+1)} - \mathbf{x}_i^{(t)}\|_1^2 \\ &\leq 3 + \frac{\alpha \log T + \log d}{\eta} + (6L^2 n) \eta \sum_{j=1}^n \sum_{t=1}^{T-1} \|\mathbf{x}_j^{(t+1)} - \mathbf{x}_j^{(t)}\|_1^2 - \frac{1}{24\eta} \sum_{t=1}^{T-1} \|\mathbf{x}_i^{(t+1)} - \mathbf{x}_i^{(t)}\|_1^2. \end{aligned}$$

Summing over all the players $i \in [n]$,

$$\begin{aligned} \sum_{i=1}^n \tilde{\text{Reg}}_i^{(T)} &\leq 3n + n \frac{\alpha \log T + \log d}{\eta} + \left(6L^2 n^2 \eta - \frac{1}{24\eta} \right) \sum_{j=1}^n \sum_{t=1}^{T-1} \|\mathbf{x}_j^{(t+1)} - \mathbf{x}_j^{(t)}\|_1^2 \\ &\leq 3n + n \frac{\alpha \log T + \log d}{\eta} - \frac{1}{48\eta} \sum_{j=1}^n \sum_{t=1}^{T-1} \|\mathbf{x}_j^{(t+1)} - \mathbf{x}_j^{(t)}\|_1^2, \end{aligned}$$

since $\eta^2 \leq \frac{1}{288L^2 n^2}$. Now, by recalling that $\tilde{\text{Reg}}_i^{(T)} \geq 0$, we get,

$$0 \leq 3n + n \frac{\alpha \log T + \log d}{\eta} - \frac{1}{48\eta} \sum_{j=1}^n \sum_{t=1}^{T-1} \|\mathbf{x}_j^{(t+1)} - \mathbf{x}_j^{(t)}\|_1^2,$$

implying

$$\sum_{j=1}^n \sum_{t=1}^{T-1} \|\mathbf{x}_j^{(t+1)} - \mathbf{x}_j^{(t)}\|_1^2 \leq 144n\eta + 48n(\alpha \log T + \log d).$$

□

Theorem 5.3 (Regret bound of DLRC-OMWU). *Under Assumption 1, if all the players $i \in [n]$ follows DLRC-OMWU with learning rate $\eta = \min\{\frac{1}{50}, \frac{1}{12\sqrt{2}Ln}\}$, then the regret of each player $i \in [n]$ is bounded as,*

$$\text{Reg}_i^{(T)} \leq 6 + \max\{50 + 12\sqrt{2}Ln, 24\sqrt{2}Ln\}(\alpha \log T + \log d),$$

and the algorithm for each player $i \in [n]$ is adaptive to adversarial utilities, i.e., the regret that each player incurs is $\text{Reg}_i^{(T)} = \tilde{O}(\sqrt{T \log d})$.

Proof. Similar to the proof of Theorem 5.2,

$$\begin{aligned} \|\boldsymbol{\nu}_i^{(t+1)} - \boldsymbol{\nu}_i^{(t)}\|_\infty^2 &\leq L^2 \left(\sum_{i=1}^n \|\mathbf{x}_i^{(t+1)} - \mathbf{x}_i^{(t)}\|_1 \right)^2 \\ &\leq L^2 n \sum_{i=1}^n \|\mathbf{x}_i^{(t+1)} - \mathbf{x}_i^{(t)}\|_1^2. \end{aligned}$$

Summing over t from 1 to $T-1$,

$$\begin{aligned} \sum_{t=1}^{T-1} \|\boldsymbol{\nu}_i^{(t+1)} - \boldsymbol{\nu}_i^{(t)}\|_\infty^2 &\leq L^2 n \sum_{t=1}^{T-1} \sum_{i=1}^n \|\mathbf{x}_i^{(t+1)} - \mathbf{x}_i^{(t)}\|_1^2, \\ &\leq L^2 n (144n\eta + 48n(\alpha \log T + \log d)) \\ &= 144L^2 n^2 \eta + 48L^2 n^2 (\alpha \log T + \log d), \end{aligned} \tag{5.30}$$

where we leveraged Theorem 5.2. By Proposition 5.7 and Theorem 5.1 we infer that,

$$\begin{aligned} \text{Reg}_i^{(T)} &\leq \tilde{\text{Reg}}_i^{(T)} \\ &\leq 3 + \frac{\alpha \log T + \log d}{\eta} + 6\eta \sum_{t=1}^{T-1} \|\boldsymbol{\nu}^{(t)} - \boldsymbol{\nu}^{(t-1)}\|_\infty^2 - \frac{1}{24\eta} \sum_{t=1}^{T-1} \|\mathbf{x}^{(t+1)} - \mathbf{x}^{(t)}\|_1^2 \\ &\leq 3 + \frac{\alpha \log T + \log d}{\eta} + 6\eta \sum_{t=1}^{T-1} \|\boldsymbol{\nu}^{(t)} - \boldsymbol{\nu}^{(t-1)}\|_\infty^2 \\ &\leq 3 + \frac{\alpha \log T + \log d}{\eta} + 864L^2 n^2 \eta^2 + 288L^2 n^2 \eta (\alpha \log T + \log d) \\ &\leq 6 + \frac{\alpha \log T + \log d}{\eta} + 12\sqrt{2}Ln(\alpha \log T + \log d) \\ &\leq 6 + \max\{50, 12\sqrt{2}Ln\}(\alpha \log T + \log d) + 12\sqrt{2}Ln(\alpha \log T + \log d) \\ &\leq 6 + \max\{50 + 12\sqrt{2}Ln, 24\sqrt{2}Ln\}(\alpha \log T + \log d), \end{aligned} \tag{5.31}$$

where Equation (5.31) is due to Theorem 5.2, and the last lines are because $\eta = \min\{\frac{1}{50}, \frac{1}{\sqrt{288}Ln}\}$.

To prove the adversarial bound for each player $i \in [n]$, player i simply check if there exists a time $t \in [T]$, such that the

$$\sum_{\tau=1}^{t-1} \|\boldsymbol{\nu}_i^{(\tau+1)} - \boldsymbol{\nu}_i^{(\tau)}\|_\infty^2 > 144L^2 n^2 \eta + 48L^2 n^2 (\alpha \log t + \log d),$$

and if noticed that, start to switch to any no-regret learning algorithm, e.g., MWU [CL06] and get $O(\sqrt{T \log d})$ regret. The argument is based on the fact that if all the players follow the DLRC-OMWU dynamics, then Equation (5.30) should be satisfied. \square

6 Conclusion

We introduced an uncoupled online learning algorithm that achieves near-constant regret of $O(n \log^2 d \log T)$ in multi-player general-sum games. This significantly improves upon the $O(d \log T)$ regret achieved by Log-Regularized Lifted Optimistic FTRL, exponentially reducing the dependence on the number of actions d [Far+22c]. Furthermore, our algorithm reduces the dependence on the number of iterations T from $O(\log^4 T)$ in the Optimistic Hedge algorithm to $O(\log T)$, improving upon the regret bound of $O(n \log d \log^4 T)$ [DFG21]. At the heart of these improvements lies a dynamic, nonmonotonic pacing of the learning rate. Specifically, players slow down their learning when their regret becomes too negative—that is, when they are significantly outperforming all fixed actions.

While our algorithm achieves near-constant regret guarantees, it remains an interesting open question whether constant regret is achievable for regularized learning in general games. Another natural direction for future research is to explore how our ideas can be applied more broadly across regularized learning algorithms. Our adaptive learning rate framework may be fruitfully combined with other FTRL-based or OMD-based methods—beyond Optimistic Multiplicative Weight Updates—to push the boundaries of performance and potentially provide a unified perspective on accelerated no-regret learning in games.

Additionally, dynamic learning rate ideas could prove valuable in minimizing stronger notions of regret, such as swap regret, within game-theoretic settings. Beyond regret minimization, a particularly compelling challenge lies in understanding the *day-to-day* dynamics of learning with adaptive pacing in structured games—offering a finer-grained view of convergence behavior and opening the door to new theoretical insights and practical strategies.

More broadly, this work contributes to a shift in multi-agent learning: rather than relying on pre-specified schedules for learning rates (e.g., fixed or monotonically decreasing steps such as $1/\sqrt{t}$), we advocate for dynamically adaptive learning rates that respond to real-time performance. Although step-size tuning is widely recognized as a critical component of single-agent learning—particularly in neural network training (see, e.g., Bengio [Ben12])—such considerations have received far less attention in multi-agent and game-theoretic settings.

We hope that our work stimulates further discussion and research at this intersection. In particular, we believe that developing game-aware adaptive schemes opens up a rich and open-ended research direction—one that bridges online learning, control theory, dynamical systems, and game theory, and may ultimately lead to fundamentally new learning dynamics tailored to strategic multi-agent environments.

7 Acknowledgments

The authors appreciate Patrick Jaillet for his insightful comments and valuable suggestions. A.S. was partially supported by the National Research Foundation Singapore and DSO National Laboratories under the AI Singapore Programme AISG Award No: AISG2-RP-2020-018, and by the Office of Naval Research (ONR) grant N00014-24-1-2470. G.F. acknowledges the support of NSF Award CCF-244306.

References

- [ABH11] Jacob Abernethy, Peter L Bartlett, and Elad Hazan. “Blackwell approachability and no-regret learning are equivalent”. In: *Conference on Learning Theory (COLT)*. 2011, pp. 27–46.
- [ACV13] Jacob Abernethy, Yiling Chen, and Jennifer Wortman Vaughan. “Efficient market making via convex optimization, and a connection to online learning”. In: *ACM Transactions on Economics and Computation (TEAC)* 1.2 (2013), pp. 1–39.
- [AK07] Sanjeev Arora and Satyen Kale. “A combinatorial, primal-dual approach to semidefinite programs”. In: *Proceedings of the thirty-ninth annual ACM symposium on Theory of computing*. 2007.
- [AL08] Sherief Abdallah and Victor Lesser. “A multiagent reinforcement learning algorithm with non-linear dynamics”. In: *Journal of Artificial Intelligence Research* 33 (2008), pp. 521–549.
- [Ana+22a] Ioannis Anagnostides et al. “Near-optimal no-regret learning for correlated equilibria in multi-player general-sum games”. In: *Proceedings of the Annual Symposium on Theory of Computing (STOC)*. 2022.
- [Ana+22b] Ioannis Anagnostides et al. “On Last-Iterate Convergence Beyond Zero-Sum Games”. In: *International Conference on Machine Learning (ICML)*. 2022.
- [Ana+22c] Ioannis Anagnostides et al. “Optimistic Mirror Descent Either Converges to Nash or to Strong Coarse Correlated Equilibria in Bimatrix Games”. In: *Neural Information Processing Systems (NeurIPS)*. 2022.
- [Ana+22d] Ioannis Anagnostides et al. “Uncoupled Learning Dynamics with $O(\log T)$ Swap Regret in Multiplayer Games”. In: *Neural Information Processing Systems (NeurIPS)*. 2022.
- [AS24] Ioannis Anagnostides and Tuomas Sandholm. “On the interplay between social welfare and tractability of equilibria”. In: *Advances in Neural Information Processing Systems* 36 (2024).
- [Azi+21] Waïss Azizian et al. “The Last-Iterate Convergence Rate of Optimistic Mirror Descent in Stochastic Variational Inequalities”. In: *Conference on Learning Theory (COLT)*. 2021.
- [BCM12] Maria-Florina Balcan, Florin Constantin, and Ruta Mehta. “The Weighted Majority Algorithm does not Converge in Nearly Zero-sum Games”. en. In: *ICML Workshop on Markets, Mechanisms, and Multi-Agent Models*. 2012.
- [BEL06] Avrim Blum, Eyal Even-Dar, and Katrina Ligett. “Routing without Regret: On Convergence to Nash Equilibria of Regret-Minimizing Algorithms in Routing Games”. In: *Proceedings of the ACM Symposium on Principles of Distributed Computing*. 2006.
- [Ben12] Yoshua Bengio. “Practical recommendations for gradient-based training of deep architectures”. In: *Neural networks: Tricks of the trade: Second edition*. Springer, 2012, pp. 437–478.

- [BGP20] 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 (COLT)*. 2020.
- [BHK09] Boaz Barak, Moritz Hardt, and Satyen Kale. “The uniform hardcore lemma via approximate bregman projections”. In: *Proceedings of the twentieth annual ACM-SIAM symposium on Discrete algorithms*. SIAM. 2009, pp. 1193–1200.
- [Big+24] Ariyan Bighashdel et al. “Policy Space Response Oracles: A Survey”. In: *arXiv preprint arXiv:2403.02227* (2024).
- [Blo+15] Daan Bloembergen et al. “Evolutionary dynamics of multi-agent learning: A survey”. In: *Journal of Artificial Intelligence Research* 53 (2015), pp. 659–697.
- [Blu90] Avrim Blum. “Learning boolean functions in an infinite attribute space”. In: *Proceedings of the twenty-second annual ACM symposium on Theory of computing*. 1990, pp. 64–72.
- [Bow04] Michael Bowling. “Convergence and no-regret in multiagent learning”. In: *Neural Information Processing Systems (NIPS)* 17 (2004).
- [BP03] Bikramjit Banerjee and Jing Peng. “Adaptive policy gradient in multiagent learning”. In: *Proceedings of the second international joint conference on Autonomous agents and multiagent systems*. 2003, pp. 686–692.
- [BP18] James P. Bailey and Georgios Piliouras. “Multiplicative Weights Update in Zero-Sum Games”. In: *ACM Conference on Economics and Computation (EC)*. 2018.
- [BP19] James P. Bailey and Georgios Piliouras. “Fast and Furious Learning in Zero-Sum Games: Vanishing Regret with Non-Vanishing Step Sizes”. In: *Neural Information Processing Systems (NeurIPS)*. 2019.
- [BS19] Noam Brown and Tuomas Sandholm. “Superhuman AI for multiplayer poker”. In: *Science* 365.6456 (2019), pp. 885–890.
- [BV02] Michael Bowling and Manuela Veloso. “Multiagent learning using a variable learning rate”. In: *Artificial intelligence* 136.2 (2002), pp. 215–250.
- [Cai+16] Yang Cai et al. “Zero-Sum Polymatrix Games: A Generalization of Minmax”. In: *Mathematics of Operations Research* 41.2 (2016), pp. 648–655.
- [Cam11] Colin F. Camerer. *Behavioral game theory: Experiments in strategic interaction*. Princeton University Press, 2011.
- [Cha+13] Erick Chastain et al. “Multiplicative updates in coordination games and the theory of evolution”. In: *Proceedings of the 4th conference on Innovations in Theoretical Computer Science*. 2013.
- [CL06] Nicolo Cesa-Bianchi and Gabor Lugosi. *Prediction, learning, and games*. Cambridge University Press, 2006.
- [CP19] Yun Kuen Cheung and Georgios Piliouras. “Vortices Instead of Equilibria in MinMax Optimization: Chaos and Butterfly Effects of Online Learning in Zero-Sum Games”. In: *Conference on Learning Theory (COLT)*. 2019.
- [CP20a] Xi Chen and Binghui Peng. “Hedging in games: Faster convergence of external and swap regrets”. In: *Neural Information Processing Systems (NeurIPS)*. 2020.

- [CP20b] Yun Kuen Cheung and Georgios Piliouras. “Chaos, extremism and optimism: Volume analysis of learning in games”. In: *Neural Information Processing Systems (NeurIPS)*. 2020.
- [CV10] Yiling Chen and Jennifer Wortman Vaughan. “A new understanding of prediction markets via no-regret learning”. In: *Proceedings of the 11th ACM conference on Electronic commerce*. 2010, pp. 189–198.
- [Das+10] Constantinos Daskalakis et al. “On Learning Algorithms for Nash Equilibria”. In: *International Symposium on Algorithmic Game Theory (SAGT)*. 2010.
- [DDK11] Constantinos Daskalakis, Alan Deckelbaum, and Anthony Kim. “Near-optimal no-regret algorithms for zero-sum games”. In: *Annual ACM-SIAM Symposium on Discrete Algorithms (SODA)*. 2011.
- [DFG21] Constantinos Daskalakis, Maxwell Fishelson, and Noah Golowich. “Near-optimal no-regret learning in general games”. In: *Neural Information Processing Systems (NeurIPS)*. 2021.
- [DGP06] Constantinos Daskalakis, Paul W. Goldberg, and Christos H. Papadimitriou. “The Complexity of Computing a Nash Equilibrium”. In: *Proceedings of the Annual Symposium on Theory of Computing (STOC)*. 2006.
- [DP19] Constantinos Daskalakis and Ioannis Panageas. “Last-Iterate Convergence: Zero-Sum Games and Constrained Min-Max Optimization”. In: *Innovations in Theoretical Computer Science (ITCS)*. 2019.
- [EH13] Ido Erev and Ernan Haruvy. “Learning and the economics of small decisions”. In: *The handbook of experimental economics 2* (2013), pp. 638–700.
- [Far+22a] Gabriele Farina et al. “Clairvoyant regret minimization: Equivalence with nemirovski’s conceptual prox method and extension to general convex games”. In: *arXiv preprint arXiv:2208.14891* (2022).
- [Far+22b] Gabriele Farina et al. “Kernelized multiplicative weights for 0/1-polyhedral games: Bridging the gap between learning in extensive-form and normal-form games”. In: *International Conference on Machine Learning*. PMLR. 2022, pp. 6337–6357.
- [Far+22c] Gabriele Farina et al. “Near-optimal no-regret learning dynamics for general convex games”. In: *Neural Information Processing Systems (NeurIPS)*. 2022.
- [FP14] Drew Fudenberg and Alexander Peysakhovich. “Recency, Records and Recaps: Learning and Non-Equilibrium Behavior in a Simple Decision Problem”. In: *ACM Conference on Economics and Computation (EC)*. 2014.
- [FS95] Yoav Freund and Robert E Schapire. “A decision-theoretic generalization of on-line learning and an application to boosting”. In: *European conference on computational learning theory*. Springer. 1995, pp. 23–37.
- [FS96] Yoav Freund and Robert E Schapire. “Game theory, on-line prediction and boosting”. In: *Proceedings of the ninth annual conference on Computational learning theory*. 1996, pp. 325–332.
- [FS97] Yoav Freund and Robert Schapire. “A decision-theoretic generalization of on-line learning and an application to boosting”. In: *Journal of Computer and System Sciences* 55.1 (1997), pp. 119–139.

- [Gin00] Herbert Gintis. *Game theory evolving: A problem-centered introduction to modeling strategic behavior*. Princeton University Press, 2000.
- [Goo+14] Ian J. Goodfellow et al. “Generative Adversarial Nets”. In: *Neural Information Processing Systems (NIPS)*. 2014.
- [GPD20] Noah Golowich, Sarath Pattathil, and Constantinos Daskalakis. “Tight last-iterate convergence rates for no-regret learning in multi-player games”. In: *Neural Information Processing Systems (NeurIPS)*. 2020.
- [HAM21] Yu-Guan Hsieh, Kimon Antonakopoulos, and Panayotis Mertikopoulos. “Adaptive Learning in Continuous Games: Optimal Regret Bounds and Convergence to Nash Equilibrium”. In: *Conference on Learning Theory (COLT)*. 2021.
- [Har+03] Hart et al. “Uncoupled dynamics do not lead to Nash equilibrium”. In: *American Economic Review* (2003), pp. 1830–1836.
- [Haz+16] Elad Hazan et al. “Introduction to online convex optimization”. In: *Foundations and Trends in Optimization* 2.3-4 (2016), pp. 157–325.
- [HM03] Sergiu Hart and Andreu Mas-Colell. “Uncoupled dynamics do not lead to Nash equilibrium”. In: *American Economic Review* 93 (2003), pp. 1830–1836.
- [HR10] Moritz Hardt and Guy N Rothblum. “A multiplicative weights mechanism for privacy-preserving data analysis”. In: *2010 IEEE 51st annual symposium on foundations of computer science*. IEEE. 2010, pp. 61–70.
- [HRU13] Justin Hsu, Aaron Roth, and Jonathan Ullman. “Differential privacy for the analyst via private equilibrium computation”. In: *Proceedings of the forty-fifth annual ACM symposium on Theory of computing*. 2013, pp. 341–350.
- [Kai+09] Michael Kaisers et al. “An evolutionary model of multi-agent learning with a varying exploration rate”. In: *Proceedings of The 8th International Conference on Autonomous Agents and Multiagent Systems-Volume 2*. 2009, pp. 1255–1256.
- [KLP11] Robert Kleinberg, Katrina Ligett, and Georgios Piliouras. “Beyond the Nash equilibrium barrier”. In: *Innovations in Theoretical Computer Science (ITCS)*. 2011.
- [KS99] Adam R Klivans and Rocco A Servedio. “Boosting and hard-core sets”. In: *40th Annual Symposium on Foundations of Computer Science (Cat. No. 99CB37039)*. IEEE. 1999.
- [Lei+21] Qi Lei et al. “Last iterate convergence in no-regret learning: constrained min-max optimization for convex-concave landscapes”. In: *International Conference on Artificial Intelligence and Statistics (AISTATS)*. 2021.
- [Lit88] Nick Littlestone. “Learning quickly when irrelevant attributes abound: A new linear-threshold algorithm”. In: *Machine learning* 2.4 (1988), pp. 285–318.
- [LN23] Gábor Lugosi and Gergely Neu. “Online-to-PAC conversions: Generalization bounds via regret analysis”. In: *arXiv preprint arXiv:2305.19674* (2023).
- [LP22] Stefanos Leonardos and Georgios Piliouras. “Exploration-exploitation in multi-agent learning: Catastrophe theory meets game theory”. In: *Artificial Intelligence* 304 (2022), p. 103653.

- [Mer+19] Panayotis Mertikopoulos et al. “Optimistic mirror descent in saddle-point problems: Going the extra(-gradient) mile”. In: *International Conference on Learning Representations (ICLR)*. 2019.
- [Mil+23] Jason Milionis et al. “An impossibility theorem in game dynamics”. In: *Proceedings of the National Academy of Sciences* 120.41 (2023).
- [Mor+17] Matej Moravčík et al. “Deepstack: Expert-level artificial intelligence in heads-up no-limit poker”. In: *Science* 356.6337 (2017), pp. 508–513.
- [MP17] Barnabé Monnot and Georgios Piliouras. “Limits and limitations of no-regret learning in games”. In: *The Knowledge Engineering Review* 32 (2017), e21.
- [MPP18] Panayotis Mertikopoulos, Christos Papadimitriou, and Georgios Piliouras. “Cycles in Adversarial Regularized Learning”. In: *Annual ACM-SIAM Symposium on Discrete Algorithms (SODA)*. 2018.
- [MS15] Jason R Marden and Jeff S Shamma. “Game theory and distributed control”. In: *Handbook of game theory with economic applications*. Vol. 4. Elsevier, 2015, pp. 861–899.
- [NE12] Iris Nevo and Ido Erev. “On surprise, change, and the effect of recent outcomes”. In: *Frontiers in psychology* 3 (2012), p. 24.
- [Nem04] 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”. In: *SIAM Journal on Optimization* 15.1 (2004), pp. 229–251.
- [Nis+07] Noam Nisan et al. *Algorithmic Game Theory*. New York, NY, USA: Cambridge University Press, 2007. ISBN: 0521872820.
- [NP10] Uri Nadav and Georgios Piliouras. “No regret learning in oligopolies: Cournot vs. Bertrand”. In: *International Symposium on Algorithmic Game Theory (SAGT)*. 2010.
- [Ora19] Francesco Orabona. “A modern introduction to online learning”. In: *arXiv preprint arXiv:1912.13213* (2019).
- [PP18] Christos Papadimitriou and Georgios Piliouras. “From Nash Equilibria to Chain Recurrent Sets: An Algorithmic Solution Concept for Game Theory”. In: *Entropy* 20.10 (2018).
- [PSS22] Georgios Piliouras, Ryann Sim, and Stratis Skoulakis. “Beyond Time-Average Convergence: Near-Optimal Uncoupled Online Learning via Clairvoyant Multiplicative Weights Update”. In: *Neural Information Processing Systems (NeurIPS)*. 2022.
- [PST95] Serge A Plotkin, David B Shmoys, and Éva Tardos. “Fast approximation algorithms for fractional packing and covering problems”. In: *Mathematics of Operations Research* 20.2 (1995), pp. 257–301.
- [Ren01] James Renegar. *A mathematical view of interior-point methods in convex optimization*. SIAM, 2001.
- [Rou15] Tim Roughgarden. “Intrinsic Robustness of the Price of Anarchy”. In: *J. ACM* 62.5 (2015), 32:1–32:42.
- [RS13a] Alexander Rakhlin and Karthik Sridharan. “Online learning with predictable sequences”. In: *Conference on Learning Theory*. PMLR. 2013, pp. 993–1019.

- [RS13b] Alexander Rakhlin and Karthik Sridharan. “Optimization, learning, and games with predictable sequences”. In: *Neural Information Processing Systems (NIPS)*. 2013.
- [RST17] Tim Roughgarden, Vasilis Syrgkanis, and Eva Tardos. “The price of anarchy in auctions”. In: *Journal of Artificial Intelligence Research* 59 (2017), pp. 59–101.
- [Sha12] Shai Shalev-Shwartz. “Online Learning and Online Convex Optimization”. In: *Foundations and Trends in Machine Learning* 4.2 (2012). ISSN: 1935-8237.
- [Sil+17] David Silver et al. “Mastering the game of go without human knowledge”. In: *nature* 550.7676 (2017), pp. 354–359.
- [Syr+15] Vasilis Syrgkanis et al. “Fast convergence of regularized learning in games”. In: *Neural Information Processing Systems (NIPS)*. 2015.
- [TW03] Eiji Takimoto and Manfred K Warmuth. “Path kernels and multiplicative updates”. In: *The Journal of Machine Learning Research* 4 (2003), pp. 773–818.
- [VFP23] Emmanouil-Vasileios Vlatakis-Gkaragkounis, Lampros Flokas, and Georgios Piliouras. “Chaos persists in large-scale multi-agent learning despite adaptive learning rates”. In: *arXiv preprint arXiv:2306.01032* (2023).
- [Vla+20] Emmanouil-Vasileios Vlatakis-Gkaragkounis et al. “No-Regret Learning and Mixed Nash Equilibria: They Do Not Mix”. In: *Neural Information Processing Systems (NeurIPS)*. 2020.
- [Wei+21] Chen-Yu Wei et al. “Linear Last-iterate Convergence in Constrained Saddle-point Optimization”. In: *International Conference on Learning Representations (ICLR)*. 2021.
- [Wei97] Jörgen W Weibull. *Evolutionary game theory*. MIT press, 1997.
- [WL18] Chen-Yu Wei and Haipeng Luo. “More Adaptive Algorithms for Adversarial Bandits”. In: *Conference on Learning Theory (COLT)*. 2018.
- [You09] H Peyton Young. “Learning by trial and error”. In: *Games and economic behavior* 65.2 (2009), pp. 626–643.