

---

# OPTIMISTIC REGRET MINIMIZATION FOR EXTENSIVE-FORM GAMES VIA DILATED DISTANCE-GENERATING FUNCTIONS\*

---

ARXIV PREPRINT

**Gabriele Farina**

Computer Science Department  
Carnegie Mellon University  
gfarina@cs.cmu.edu

**Christian Kroer**

IEOR Department  
Columbia University  
christian.kroer@columbia.edu

**Tuomas Sandholm**

Computer Science Department, CMU  
Strategic Machine, Inc.  
Strategy Robot, Inc.  
Optimized Markets, Inc.  
sandholm@cs.cmu.edu

October 27, 2019

## Abstract

We study the performance of optimistic regret-minimization algorithms for both minimizing regret in, and computing Nash equilibria of, zero-sum extensive-form games. In order to apply these algorithms to extensive-form games, a distance-generating function is needed. We study the use of the dilated entropy and dilated Euclidean distance functions. For the dilated Euclidean distance function we prove the first explicit bounds on the strong-convexity parameter for general treplexes. Furthermore, we show that the use of dilated distance-generating functions enable us to decompose the mirror descent algorithm, and its optimistic variant, into local mirror descent algorithms at each information set. This decomposition mirrors the structure of the counterfactual regret minimization framework, and enables important techniques in practice, such as distributed updates and pruning of cold parts of the game tree. Our algorithms provably converge at a rate of  $T^{-1}$ , which is superior to prior counterfactual regret minimization algorithms. We experimentally compare to the popular algorithm CFR+, which has a theoretical convergence rate of  $T^{-0.5}$  in theory, but is known to often converge at a rate of  $T^{-1}$ , or better, in practice. We give an example matrix game where CFR+ experimentally converges at a relatively slow rate of  $T^{-0.74}$ , whereas our optimistic methods converge faster than  $T^{-1}$ . We go on to show that our fast rate also holds in the Kuhn poker game, which is an extensive-form game. For games with deeper game trees however, we find that CFR+ is still faster. Finally we show that when the goal is minimizing regret, rather than computing a Nash equilibrium, our optimistic methods can outperform CFR+, even in deep game trees.

---

\*This paper was accepted for publication at NeurIPS 2019.

## 1. Introduction

Extensive-form games (EFGs) are a broad class of games that can model sequential interaction, imperfect information, and stochastic outcomes. To operationalize them they must be accompanied by techniques for computing game-theoretic equilibria such as Nash equilibrium. A notable success story of this is poker: Bowling et al. (2015) computed a near-optimal Nash equilibrium for heads-up limit Texas hold'em, while Brown & Sandholm (2017b) beat top human specialist professionals at the larger game of heads-up no-limit Texas hold'em. Solving extremely large EFGs relies on many methods for dealing with the scale of the problem: abstraction methods are sometimes used to create smaller games (Gilpin & Sandholm, 2007; Lanctot et al., 2012; Kroer & Sandholm, 2014; Ganzfried & Sandholm, 2014; Brown et al., 2015; Kroer & Sandholm, 2016), endgame solving is used to compute refined solutions to the end of the game in real time (Burch et al., 2014; Ganzfried & Sandholm, 2015; Moravčík et al., 2016), and recently depth-limited subgame solving has been very successfully used in real time (Moravčík et al., 2017; Brown et al., 2018; Brown & Sandholm, 2019b). At the core of all these methods is a reliance on a fast algorithm for computing approximate Nash equilibria of the abstraction, endgame, and/or depth-limited subgame (Moravčík et al., 2017; Brown et al., 2018; Brown & Sandholm, 2019b). In practice the most popular method has been the  $\text{CFR}^+$  algorithm (Zinkevich et al., 2007; Tammelin et al., 2015), which was used within all three two-player poker breakthroughs (Bowling et al., 2015; Moravčík et al., 2017; Brown & Sandholm, 2017b).  $\text{CFR}^+$  has been shown to converge to a Nash equilibrium at a rate of  $T^{-0.5}$ , but in practice it often performs much better, even outperforming faster methods that have a guaranteed rate of  $T^{-1}$  (Brown et al., 2017; Kroer et al., 2018b;a; Brown & Sandholm, 2019a).

Recently, another class of optimization algorithms has been shown to have appealing theoretical properties. *Online convex optimization* (OCO) algorithms are online variants of first-order methods: at each timestep  $t$  they receive some loss function  $\ell^t$  (often a linear loss which is a gradient of some underlying loss function), and must then recommend a point from some convex set based on the series of past points and losses. While these algorithms are generally known to have a  $T^{-0.5}$  rate of convergence when solving static problems, a recent series of papers showed that when two *optimistic* OCO algorithms are faced against each other, and they have some estimate of the next loss faced, a rate of  $T^{-1}$  can be achieved (Rakhlin & Sridharan, 2013a;b; Syrgkanis et al., 2015). In this paper we investigate the application of these algorithms to EFG solving, both in the regret-minimization setting, and for computing approximate Nash equilibria at the optimal rate of  $O(T^{-1})$ . The only prior attempt at using optimistic OCO algorithm in extensive-form games is due to Farina et al. (2019). In that paper, the authors show that by restricting to the weaker notion of *stable-predictive optimism*, one can mix and match local stable-predictive optimistic algorithm at every decision point in the game as desired and obtain an overall stable-predictive optimistic algorithm that enables  $O(T^{-0.75})$  convergence to Nash equilibrium. The approach we adopt in this paper is different from that of Farina et al. (2019) in that our construction does not allow one to pick different regret minimizers for different decision points; however, our algorithms converge to Nash equilibrium at the improved rate  $O(T^{-1})$ .

The main hurdle to overcome is that in all known OCO algorithms a *distance-generating function* (DGF) is needed to maintain feasibility via proximal operators and ensure that the stepsizes of the algorithms are appropriate for the convex set at hand. For the case of EFGs, the convex set is known as a *treplex*, and the so-called dilated DGFs are known to have appealing properties, including closed-form iterate updates and strong convexity properties (Hoda et al., 2010; Kroer et al., 2018b). In particular, the dilated entropy DGF, which applies the negative entropy at each information set, is known to lead to the state-of-the-art theoretical rate on convergence for iterative methods (Kroer et al., 2018b). Another potential DGF is the dilated Euclidean DGF, which applies the  $\ell_2$  norm as a DGF at each information set. We show the first explicit bounds on the strong-convexity parameter for the dilated Euclidean DGF when applied to the strategy space of an EFG. We go on to show that when a dilated DGF is paired with the *online mirror descent* (OMD) algorithm, or its optimistic variant, the resulting algorithm decomposes into a recursive application of local online mirror descent algorithms at each information set of the game. This decomposition is similar to the decomposition achieved in the counterfactual regret minimization framework, where a local regret minimizer is applied on the counterfactual regret at each information set. This localization of the updates

along the tree structure enables further techniques, such as distributing the updates (Brown & Sandholm, 2017b; Brown et al., 2015) or skipping updates on cold parts of the game tree (Brown & Sandholm, 2017a).

It is well-known that the entropy DGF is the theoretically superior DGF when applied to optimization over a simplex (Hoda et al., 2010). For the treplex case where the entropy DGF is used at each information set, Kroer et al. (2018b) showed that the strong theoretical properties of the simplex entropy DGF generalize to the dilated entropy DGF on a treplex (with earlier weaker results shown by Kroer et al. (2015)). Our results on the dilated Euclidean DGF confirm this finding, as the dilated Euclidean DGF has a similar strong convexity parameter, but with respect to the  $\ell_2$  norm, rather than the  $\ell_1$  norm for dilated entropy (having strong convexity with respect to the  $\ell_1$  norm leads to a tighter convergence-rate bound because it gives a smaller matrix norm, another important constant in the rate).

In contrast to these theoretical results, for the case of computing a Nash equilibrium in matrix games it has been found experimentally that the Euclidean DGF often performs much better than the entropy DGF. This was shown by Chambolle & Pock (2016) when using a particular accelerated primal-dual algorithm (Chambolle & Pock, 2011; 2016) and using the *last iterate* (as opposed to the uniformly-averaged iterate as the theory suggests). Kroer (2019) recently showed that this extends to the theoretically-sound case of using linear or quadratic averaging in the same primal-dual algorithm, or in mirror prox (Nemirovski, 2004) (the offline variant of optimistic OMD). In this paper we replicate these results when using OCO algorithms: first we show it on a particular matrix game, where we also exhibit a slow  $T^{-0.74}$  convergence rate of CFR<sup>+</sup> (the slowest CFR<sup>+</sup> rate seen to the best of our knowledge). We show that for the Kuhn poker game the last iterate of optimistic OCO algorithms with the dilated Euclidean DGF also converges extremely fast. In contrast to this, we show that for deeper EFGs CFR<sup>+</sup> is still faster. Finally we compare the performance of CFR<sup>+</sup> and optimistic OCO algorithms for minimizing regret, where we find that OCO algorithms perform better.

## 2. Regret Minimization Algorithms

In this section we present the regret-minimization algorithms that we will work with. We will operate within the framework of *online convex optimization* (Zinkevich, 2003). In this setting, a decision maker repeatedly plays against an unknown environment by making decision  $\mathbf{x}^1, \mathbf{x}^2, \dots \in \mathcal{X}$  for some convex compact set  $\mathcal{X}$ . After each decision  $\mathbf{x}^t$  at time  $t$ , the decision maker faces a *linear loss*  $\mathbf{x}^t \mapsto \langle \boldsymbol{\ell}^t, \mathbf{x}^t \rangle$ , where  $\boldsymbol{\ell}^t$  is a vector in  $\mathcal{X}$ . Summarizing, the decision maker makes a decision  $\mathbf{x}^{t+1}$  based on the sequence of losses  $\boldsymbol{\ell}^1, \dots, \boldsymbol{\ell}^t$  as well as the sequence of past iterates  $\mathbf{x}^1, \dots, \mathbf{x}^t$ .

The quality metric for a regret minimizer is its *cumulative regret*, which is the difference between the loss cumulated by the sequence of decisions  $\mathbf{x}^1, \dots, \mathbf{x}^T$  and the loss that would have been cumulated by playing the best-in-hindsight time-independent decision  $\hat{\mathbf{x}}$ . Formally, the cumulative regret up to time  $T$  is

$$R^T := \sum_{t=1}^T \langle \boldsymbol{\ell}^t, \mathbf{x}^t \rangle - \min_{\hat{\mathbf{x}} \in \mathcal{X}} \left\{ \sum_{t=1}^T \langle \boldsymbol{\ell}^t, \hat{\mathbf{x}} \rangle \right\}.$$

A “good” regret minimizer is such that the cumulative regret grows *sublinearly in  $T$* .

The algorithms we consider assume access to a *distance-generating function*  $d : \mathcal{X} \rightarrow \mathbb{R}$ , which is 1-strongly convex (with respect to some norm) and continuously differentiable on the interior of  $\mathcal{X}$ . Furthermore  $d$  should be such that the gradient of the convex conjugate  $\nabla d(\mathbf{g}) = \operatorname{argmax}_{\mathbf{x} \in \mathcal{X}} \langle \mathbf{g}, \mathbf{x} \rangle - d(\mathbf{x})$  is easy to compute. Following Hoda et al. (2010) we say that a DGF satisfying these properties is a *nice* DGF for  $\mathcal{X}$ . From  $d$  we also construct the *Bregman divergence*  $D(\mathbf{x} \parallel \mathbf{x}') := d(\mathbf{x}) - d(\mathbf{x}') - \langle \nabla d(\mathbf{x}'), \mathbf{x} - \mathbf{x}' \rangle$ .

First we present two classical regret minimization algorithms. The *online mirror descent* (OMD) algorithm produces iterates according to the rule

$$\mathbf{x}^{t+1} = \operatorname{argmin}_{\mathbf{x} \in \mathcal{X}} \left\{ \langle \boldsymbol{\ell}^t, \mathbf{x} \rangle + \frac{1}{\eta} D(\mathbf{x} \parallel \mathbf{x}^t) \right\}. \quad (1)$$

The *follow the regularized leader* (FTRL) algorithm produces iterates according to the rule (Shalev-Shwartz & Singer, 2007)

$$\mathbf{x}^{t+1} = \operatorname{argmin}_{\mathbf{x} \in \mathcal{X}} \left\{ \left\langle \sum_{\tau=1}^t \ell^\tau, \mathbf{x} \right\rangle + \frac{1}{\eta} d(\mathbf{x}) \right\}. \quad (2)$$

OMD and FTRL satisfy regret bounds of the form  $R^T \leq O\left(D(\mathbf{x}^* \parallel \mathbf{x}^1) L \sqrt{T}\right)$  (e.g. Hazan (2016)).

The *optimistic* variants of the classical regret minimization algorithms take as input an additional vector  $\mathbf{m}^{t+1}$ , which is an estimate of the loss faced at time  $t+1$  (Chiang et al., 2012; Rakhlin & Sridharan, 2013a). Optimistic OMD produces iterates according to the rule (Rakhlin & Sridharan, 2013a) (note that  $\mathbf{x}^{t+1}$  is produced before seeing  $\ell^{t+1}$ , while  $\mathbf{z}^{t+1}$  is produced after)

$$\mathbf{x}^{t+1} = \operatorname{argmin}_{\mathbf{x} \in \mathcal{X}} \left\{ \langle \mathbf{m}^{t+1}, \mathbf{x} \rangle + \frac{1}{\eta} D(\mathbf{x} \parallel \mathbf{z}^t) \right\}, \quad \mathbf{z}^{t+1} = \operatorname{argmin}_{\mathbf{z} \in \mathcal{X}} \left\{ \langle \ell^{t+1}, \mathbf{z} \rangle + \frac{1}{\eta} D(\mathbf{z} \parallel \mathbf{z}^t) \right\}. \quad (3)$$

Thus it is like OMD, except that  $\mathbf{x}^{t+1}$  is generated by an additional step taken using the loss estimate. This additional step is transient in the sense that  $\mathbf{x}^{t+1}$  is not used as a center for the next iterate. OFTRL produces iterates according to the rule (Rakhlin & Sridharan, 2013a; Syrgkanis et al., 2015)

$$\mathbf{x}^{t+1} = \operatorname{argmin}_{\mathbf{x} \in \mathcal{X}} \left\{ \left\langle \mathbf{m}^{t+1} + \sum_{\tau=1}^t \ell^\tau, \mathbf{x} \right\rangle + \frac{1}{\eta} d(\mathbf{x}) \right\}. \quad (4)$$

Again the loss estimate is used in a transient way: it is used as if we already saw the loss at time  $t+1$ , but then discarded and not used in future iterations.

## 2.1. Connection to Saddle Points

A *bilinear saddle-point problem* is a problem of the form  $\min_{\mathbf{x} \in \mathcal{X}} \max_{\mathbf{y} \in \mathcal{Y}} \{ \mathbf{x}^\top \mathbf{A} \mathbf{y} \}$ , where  $\mathcal{X}, \mathcal{Y}$  are closed convex sets. This general formulation allows us to capture, among other settings, several game-theoretical applications such as computing Nash equilibria in two-player zero-sum games. In that setting,  $\mathcal{X}$  and  $\mathcal{Y}$  are convex polytopes whose description is provided by the *sequence-form constraints*, and  $\mathbf{A}$  is a real payoff matrix (von Stengel, 1996).

The error metric that we use is the *saddle-point residual* (or *gap*)  $\xi$  of  $(\bar{\mathbf{x}}, \bar{\mathbf{y}})$ , defined as  $\xi(\bar{\mathbf{x}}, \bar{\mathbf{y}}) := \max_{\hat{\mathbf{y}} \in \mathcal{Y}} \langle \bar{\mathbf{x}}, \mathbf{A} \hat{\mathbf{y}} \rangle - \min_{\hat{\mathbf{x}} \in \mathcal{X}} \langle \hat{\mathbf{x}}, \mathbf{A} \bar{\mathbf{y}} \rangle$ . A well-known folk theorem shows that the average of a sequence of regret-minimizing strategies for the choice of losses  $\ell_{\mathcal{X}}^t : \mathcal{X} \ni \mathbf{x} \mapsto (-\mathbf{A} \mathbf{y}^t)^\top \mathbf{x}$ ,  $\ell_{\mathcal{Y}}^t : \mathcal{Y} \ni \mathbf{y} \mapsto (\mathbf{A}^\top \mathbf{x}^t)^\top \mathbf{y}$  leads to a bounded saddle-point residual, since one has

$$\xi(\bar{\mathbf{x}}, \bar{\mathbf{y}}) = \frac{1}{T} (R_{\mathcal{X}}^T + R_{\mathcal{Y}}^T). \quad (5)$$

When  $\mathcal{X}, \mathcal{Y}$  are the players' sequence-form strategy spaces, this implies that the average strategy profile produced by the regret minimizers is a  $1/T(R_{\mathcal{X}}^T + R_{\mathcal{Y}}^T)$ -Nash equilibrium. This also implies that by using online mirror descent or follow-the-regularizer-leader, one obtains an anytime algorithm for computing a Nash equilibrium. In particular, at each time  $T$ , the average strategy output by each of the two regret minimizers forms a  $\epsilon$ -Nash equilibrium, where  $\epsilon = O(T^{-0.5})$ .

## 2.2. RVU Property and Fast Convergence to Saddle Points

Both optimistic OMD and optimistic FTRL satisfy the *Regret bounded by Variation in Utilities* (RVU) property, as given by Syrgkanis et al.:

**Definition 1** (RVU property, (Syrgkanis et al., 2015)). *We say that a regret minimizer satisfies the RVU property if there exist constants  $\alpha > 0$  and  $0 < \beta \leq \gamma$ , as well as a pair of dual norms  $(\|\cdot\|, \|\cdot\|_*)$  such that, no matter what the loss functions  $\ell^1, \dots, \ell^T$  are,*

$$R^T \leq \alpha + \beta \sum_{t=1}^T \|\ell^t - \mathbf{m}^t\|_*^2 - \gamma \sum_{t=1}^T \|\mathbf{x}^t - \mathbf{x}^{t-1}\|^2. \quad (\text{RVU})$$

The definition given here is slightly more general than that of Syrgkanis et al. (2015): we allow a general estimate  $\mathbf{m}^t$  of  $\ell^t$ , whereas their definition requires using  $\mathbf{m}^t = \ell^{t-1}$ . While the choice  $\mathbf{m}^t = \ell^{t-1}$  is often reasonable, in some cases other definitions of the loss prediction are more natural (Farina et al., 2019). In practice, both optimistic OMD and optimistic FTRL satisfy a parametric notion of the RVU property, which depends on the value of the step-size parameter that was chosen to set up either algorithm.

**Theorem 1** (Syrgkanis et al. (2015)). *For all step-size parameters  $\eta > 0$ , Optimistic OMD satisfies the RVU conditions with respect to the primal-dual norm pair  $(\|\cdot\|_1, \|\cdot\|_\infty)$  with parameters  $\alpha = R/\eta$ ,  $\beta = \eta$ ,  $\gamma = 1/(8\eta)$ , where  $R$  is a constant that scales with the maximum allowed norm of any loss function  $\ell$ .*

**Theorem 2.** *For all step-size parameters  $\eta > 0$ , OFTRL satisfies the RVU conditions with respect to any primal-dual norm pair  $(\|\cdot\|, \|\cdot\|_*)$  with parameters  $\alpha = \Delta_d/\eta$ ,  $\beta = \eta$ ,  $\gamma = 1/(4\eta)$ , where  $\Delta_d := \max_{\mathbf{x}, \mathbf{y} \in \mathcal{X}} \{d(\mathbf{x}) - d(\mathbf{y})\}$ .*

Our proof, available in the appendix of the full paper, generalizes the work by Syrgkanis et al. (2015) by extending the proof beyond simplex domains and beyond the fixed choice  $\mathbf{m}^t = \ell^{t-1}$ .

It turns out that this is enough to accelerate the convergence to a saddle point in the construction of Section 2.1. In particular, by letting the predictions be defined as  $\mathbf{m}_x^t := \ell_x^{t-1}$ ,  $\mathbf{m}_y^t := \ell_y^{t-1}$ , we obtain that the residual  $\xi$  of the average decisions  $(\bar{\mathbf{x}}, \bar{\mathbf{y}})$  satisfies

$$\begin{aligned} T\xi(\bar{\mathbf{x}}, \bar{\mathbf{y}}) &\leq \frac{2\alpha'}{\eta} + \eta \sum_{t=1}^T \left( \|\mathbf{A}\mathbf{y}^t + \mathbf{A}\mathbf{y}^{t-1}\|_*^2 + \|\mathbf{A}^\top \mathbf{x}^t - \mathbf{A}^\top \mathbf{x}^{t-1}\|_*^2 \right) \\ &\quad - \frac{\gamma'}{\eta} \sum_{t=1}^T \left( \|\mathbf{x}^t - \mathbf{x}^{t-1}\|^2 + \|\mathbf{y}^t - \mathbf{y}^{t-1}\|^2 \right) \\ &\leq \frac{2\alpha'}{\eta} + \left( \eta \|\mathbf{A}\|_{\text{op}}^2 - \frac{\gamma'}{\eta} \right) \left( \sum_{t=1}^T \|\mathbf{x}^t - \mathbf{x}^{t-1}\|^2 + \sum_{t=1}^T \|\mathbf{y}^t - \mathbf{y}^{t-1}\|^2 \right), \end{aligned}$$

where the first inequality holds by plugging (RVU) into (5), and the second inequality by noting that the operator norm  $\|\cdot\|_{\text{op}}$  of a linear function is equal to the operator norm of its transpose. This implies that when the step-size parameter is chosen as  $\eta = \frac{\sqrt{\gamma'}}{\|\mathbf{A}\|_{\text{op}}}$ , the saddle-point gap  $\xi(\bar{\mathbf{x}}, \bar{\mathbf{y}})$  satisfies  $\xi(\bar{\mathbf{x}}, \bar{\mathbf{y}}) \leq \frac{2\alpha' \|\mathbf{A}\|_{\text{op}}}{T\sqrt{\gamma'}} = O(T^{-1})$ .

### 3. Treplexes and Sequence Form

We formalize a sequential decision process as follows. We assume that we have a set of decision points  $\mathcal{J}$ . Each decision point  $j \in \mathcal{J}$  has a set of actions  $A_j$  of size  $n_j$ . Given a specific action at  $j$ , the set of possible decision points that the agent may next face is denoted by  $\mathcal{C}_{j,a}$ . It can be an empty set if no more actions are taken after  $j$ ,  $a$ . We assume that the decision points form a tree, that is,  $\mathcal{C}_{j,a} \cap \mathcal{C}_{j',a'} = \emptyset$  for all other convex sets and action choices  $j', a'$ . This condition is equivalent to the perfect-recall assumption in extensive-form games, and to conditioning on the full sequence of actions and observations in a finite-horizon partially-observable decision process. In our definition, the decision space starts with a root decision point, whereas in practice multiple root decision points may be needed, for example in order to model different starting hands in card games. Multiple root decision points can be modeled by having a dummy root decision point with only a single action.

The set of possible next decision points after choosing action  $a \in A_j$  at decision point  $j \in \mathcal{J}$ , denoted  $\mathcal{C}_{j,a}$ , can be thought of as representing the different decision points that an agent may face after taking action  $a$  and then making an observation on which she can condition her next action choice. In addition to games, our model of sequential decision process captures, for example, partially-observable Markov decision processes and Markov decision processes where we condition on the entire history of observations and actions.

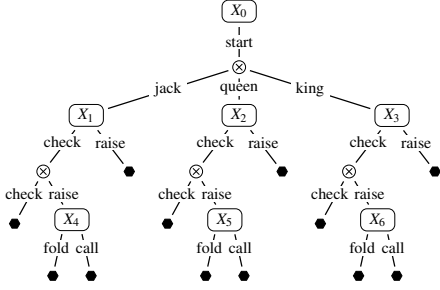


Figure 1. Sequential action space for the first player in the game of Kuhn poker.  $\otimes$  denotes an observation point;  $\bullet$  represents the end of the decision process.

As an illustration, consider the game of Kuhn poker (Kuhn, 1950). Kuhn poker consists of a three-card deck: king, queen, and jack. The action space for the first player is shown in Figure 1. For instance, we have:  $\mathcal{J} = \{0, 1, 2, 3, 4, 5, 6\}$ ;  $n_0 = 1$ ;  $n_j = 2$  for all  $j \in \mathcal{J} \setminus \{0\}$ ;  $A_0 = \{\text{start}\}$ ,  $A_1 = A_2 = A_3 = \{\text{check}, \text{raise}\}$ ,  $A_4 = A_5 = A_6 = \{\text{fold}, \text{call}\}$ ;  $\mathcal{C}_{0,\text{start}} = \{1, 2, 3\}$ ,  $\mathcal{C}_{1,\text{raise}} = \emptyset$ ,  $\mathcal{C}_{3,\text{check}} = \{6\}$ ; etc.

The expected loss for a given strategy is non-linear in the vectors of probability masses for each decision point  $j$ . This non-linearity is due to the probability of reaching each  $j$ , which is computed as the product of the probabilities of all actions on the path to from the root to  $j$ . An alternative formulation which preserves linearity

is called the *sequence form*. In the sequence-form representation, the simplex strategy space at a generic decision point  $j \in \mathcal{J}$  is scaled by the decision variable associated with the last action in the path from the root of the process to  $j$ . In this formulation, the value of a particular action represents the probability of playing the whole *sequence* of actions from the root to that action. This allows each term in the expected loss to be weighted only by the sequence ending in the corresponding action. The sequence form has been used to instantiate linear programming (von Stengel, 1996) and first-order methods (Hoda et al., 2010; Kroer et al., 2015; 2018b) for computing Nash equilibria of zero-sum EFGs. Formally, the sequence-form representation  $\mathcal{X}$  of a sequential decision process can be obtained recursively, as follows: for every  $j \in \mathcal{J}$ ,  $a \in A_j$ , we let  $\mathcal{X}_{\downarrow j,a} := \prod_{j' \in \mathcal{C}_{j,a}} \mathcal{X}_{\downarrow j'}$ , where  $\Pi$  denotes Cartesian product; at every decision point  $j \in \mathcal{J}$ , we let

$$\mathcal{X}_{\downarrow j} := \{(\lambda_1, \dots, \lambda_{n_j}, \lambda_1 \mathbf{x}_{a_1}, \dots, \lambda_{n_j} \mathbf{x}_{a_{n_j}}) : (\lambda_1, \dots, \lambda_n) \in \Delta^{n_j}, \mathbf{x}_a \in \mathcal{X}_{\downarrow j,a} \forall a \in A_j\},$$

where we assumed  $A_j = \{a_1, \dots, a_{n_j}\}$ .

The sequence form strategy space for the whole sequential decision process is then  $\mathcal{X} := \{1\} \times \mathcal{X}_{\downarrow r}$ , where  $r$  is the root of the process. The first entry, identically equal to 1 for any point in  $\mathcal{X}$ , corresponds to what is called the *empty sequence*. Crucially,  $\mathcal{X}$  is a convex and compact set, and the expected loss of the process is a linear function over  $\mathcal{X}$ . With the sequence-form representation the problem of computing a Nash equilibrium in an EFG can be formulated as a *bilinear saddle-point problem* (see Section 2.1), where  $\mathcal{X}$  and  $\mathcal{Y}$  are the sequence-form strategy spaces of the sequential decision processes faced by the two players, and  $\mathbf{A}$  is a sparse matrix encoding the leaf payoffs of the game.

As we have already observed, vectors that pertain to the sequence form have one entry for each sequence of the decision process. We denote with  $v_\phi$  the entry in  $\mathbf{v}$  corresponding to the empty sequence, and  $v_{ja}$  the entry corresponding to any other sequence  $(j, a)$  where  $j \in \mathcal{J}$ ,  $a \in A_j$ . Sometimes, we will need to *slice* a vector  $\mathbf{v}$  and isolate only those entries that refer to all decision points  $j'$  and actions  $a' \in A_{j'}$  that are at or below some  $j \in \mathcal{J}$ ; we will denote such operation as  $\mathbf{v}_{\downarrow j}$ . Similarly, we introduce the syntax  $v_j$  to denote the subset of  $n_j = |A_j|$  entries of  $\mathbf{v}$  that pertain to all actions  $a \in A_j$  at decision point  $j \in \mathcal{J}$ . Finally, note that for any  $j \in \mathcal{J} - \{r\}$  there is a unique sequence  $(j', a')$ , denoted  $p_j$  and called the *parent sequence of decision point j*, such that  $j \in \mathcal{C}_{j'a'}$ . When  $j = r$  is the root decision point, we let  $p_r := \phi$ , the empty sequence.

#### 4. Dilated Distance Generating Functions

We will be interested in a particular type of DGF which is suitable for sequential decision-making problems: a *dilated DGF*. A dilated DGF is constructed by taking a sum over suitable local DGFs for each decision point, where each local DGF is dilated by the parent variable leading to the decision point:  $d(\mathbf{x}) = \sum_{j \in \mathcal{J}} x_{p_j} d_j \left( \frac{\mathbf{x}_j}{x_{p_j}} \right)$ . Each “local” DGF  $d_j$  is given the local variable  $\mathbf{x}_j$  divided by  $x_{p_j}$ , so that  $\frac{\mathbf{x}_j}{x_{p_j}} \in \Delta^{n_j}$ . The idea is that  $d_j$  can be any DGF suitable for  $\Delta^{n_j}$ ; by multiplying  $d_j$  by  $x_{p_j}$  and taking a sum over  $\mathcal{J}$

we construct a DGF for the whole treeplex from these local DGFs. Hoda et al. (2010) showed that dilated DGFs have many of the desired properties of a DGF for an optimization problem over a treeplex.

We now present two local DGFs for simplexes, that are by far the most common in practice. In the following we let  $\mathbf{b}$  be a vector in the  $n$ -dimensional simplex  $\Delta^n$ . First, the *Euclidean DGF*  $d(\mathbf{b}) = \|\mathbf{b}\|_2^2$ , which is 1-strongly convex with respect to the  $\ell_2$  norm; secondly, the *negative entropy DGF*  $d(\mathbf{b}) = \sum_{i=1}^n b_i \log(b_i)$  (we will henceforth drop the “negative” and simply refer to it as the entropy DGF), which is 1-strongly convex with respect to the  $\ell_1$  norm. The strong convexity properties of the dilated entropy DGF were shown by Kroer et al. (2018b) (with earlier weaker results shown by Kroer et al. (2015)). However, for the dilated Euclidean DGF a setup for achieving a strong-convexity parameter of 1 was unknown until now; Hoda et al. (2010) show that a strong-convexity parameter exists, but do not show what it is for the general case (they give specific results for a particular class of *uniform treeplexes*). We now show how to achieve this.

We are now ready to state our first result on dilated regularizers that are strongly convex with respect to the Euclidean norm:

**Theorem 3.** *Let  $d(\mathbf{x}) = \sum_{j \in \mathcal{J}} x_{p_j} d_j(\mathbf{x}_j/x_{p_j})$  where for all  $j$ ,  $d_j$  is  $\mu_j$ -strongly convex with respect to the Euclidean norm over  $\Delta^{n_j}$ . Furthermore, define  $\sigma_{j_a} := \frac{\mu_j}{2} - \sum_{j' \in \mathcal{C}_{j_a}} \mu_{j'}$ , and  $\bar{\sigma} := \min_{j_a} \sigma_{j_a}$ . Then,  $d$  is  $\bar{\sigma}$ -strongly convex with respect to the Euclidean norm over  $\mathcal{X}$ .*

We can immediately use Theorem 3 to prove the following corollary:

**Corollary 1.** *Let  $\bar{\sigma} > 0$  be arbitrary, and for all  $j$  let  $d_j$  be a  $\mu_j$ -strongly convex function over  $\Delta^{n_j}$  with respect to the Euclidean norm, where the  $\mu_j$ 's satisfy*

$$\mu_j = 2\bar{\sigma} + 2 \max_{a \in A_j} \sum_{j' \in \mathcal{C}_{j_a}} \mu_{j'}. \quad (6)$$

*Then,  $d(\mathbf{x}) = \sum_{j \in \mathcal{J}} x_{p_j} d_j(\mathbf{x}_j/x_{p_j})$  is  $\bar{\sigma}$ -strongly convex over  $\mathcal{X}$  with respect to the Euclidean norm.*

## 5. Local Regret Minimization

We now show that OMD and Optimistic OMD run on a treeplex  $\mathcal{X}$  with a dilated DGF can both be interpreted as locally minimizing a modified variant of loss at each information set, with correspondingly-modified loss predictions. The modified local loss at a given information set  $j$  takes into account the loss and DGF below  $j$  by adding the expectation with respect to the next iterate  $\mathbf{x}_{\downarrow j}^t$ . In practice this modified loss is easily handled by computing  $\mathbf{x}^t$  bottom-up, thereby visiting  $j$  after having visited the whole subtree below.

We first show that the problem of computing the *prox mapping*, the minimizer of a linear term plus the Bregman divergence, decomposes into local prox mappings at each simplex of a treeplex. This will then be used to show that OMD and Optimistic OMD can be viewed as a tree of local simplex-instantiations of the respective algorithms.

### 5.1. Decomposition into Local Prox Mappings with a Dilated DGF

We will be interested in solving the following prox mapping, which takes place in the sequence form:

$$\text{Prox}(\mathbf{g}, \hat{\mathbf{x}}) = \underset{\mathbf{x} \in \mathcal{X}}{\text{argmin}} \left\{ \langle \mathbf{g}, \mathbf{x} \rangle + D(\mathbf{x} \parallel \hat{\mathbf{x}}) \right\}. \quad (7)$$

The reason is that the update applied at each iteration of several OCO algorithms run on the sequence-form polytope of  $\mathcal{X}$  can be described as an instantiation of this prox mapping. We now show that this update can be interpreted as a local prox mapping at each decision point, but with a new loss  $\hat{g}_j$  that depends on the update applied in the subtree beneath  $j$ .

**Proposition 1** (Decomposition into local prox mappings). *A prox mapping (7) on a treeplex with a Bregman divergence constructed from a dilated DGF decomposes into local prox mappings at each decision point  $j$  where the solution is as follows:*

$$\mathbf{x}_j^* = x_{p_j} \cdot \underset{\mathbf{b}_j \in \Delta^{n_j}}{\text{argmin}} \left\{ \langle \hat{g}_j, \mathbf{b}_j \rangle + D_j \left( \mathbf{b}_j \parallel \frac{\hat{\mathbf{x}}_j}{\hat{x}_{p_j}} \right) \right\},$$

where

$$\hat{g}_{j,a} = g_{j,a} + \sum_{j' \in \mathcal{C}_{j,a}} \left[ d_{\downarrow j'}^* (-\mathbf{g}_{\downarrow j'} + \nabla d_{\downarrow j'}(\hat{\mathbf{x}}_{\downarrow j'})) - d_{j'} \left( \frac{\hat{\mathbf{x}}_j}{\hat{x}_{p_j}} \right) + \left\langle \nabla d_{j'} \left( \frac{\hat{\mathbf{x}}_{j'}}{\hat{x}_{p_{j'}}} \right), \frac{\hat{\mathbf{x}}_{j'}}{\hat{x}_{p_{j'}}} \right\rangle \right].$$

Hoda et al. (2010) and Kroer et al. (2018a) gave variations on a similar result: that the convex conjugate  $d_{\downarrow j}^*(-\mathbf{g})$  can be computed in bottom-up fashion similar to the recursion we show here. Proposition 1 is slightly different in that we additionally show that the Bregman divergence also survives the decomposition and can be viewed as a local Bregman divergence. This latter difference will be necessary for showing that OMD can be interpreted as a local RM.

## 5.2. Decomposition into Local Regret Minimizers

With Proposition 1 it follows almost directly that OMD and Optimistic OMD can be seen as a set of local regret minimizers, one for each simplex. Each produces iterates from their respective simplex, with the overall strategy produced by then applying the sequence-form transformation to these local iterates.

**Theorem 4.** *OMD with a dilated DGF for a treeplex  $\mathcal{X}$  corresponds to running OMD locally at each simplex  $j$ , with the local loss  $\ell^t$  constructed according to Proposition 1. Optimistic OMD corresponds to the optimistic variant of this local OMD with local loss predictions  $\hat{\ell}^t, \hat{\mathbf{m}}_j^{t+1}$  again constructed according to Proposition 1 using  $\mathbf{x}^t$  as Bregman divergence center and  $\mathbf{x}^{t+1}$  for aggregating losses below each simplex. Here the modified loss uses  $\mathbf{z}_{\downarrow j}^t$  and  $\mathbf{x}^{t+1}$  as Bregman divergence center and aggregating loss below, respectively. The prediction  $\hat{\mathbf{m}}_j^{t+1}$  uses  $\mathbf{z}_{\downarrow j}^t$  and  $\mathbf{z}^{t+1}$ .*

Unlike OMD and its optimistic variant, it is not the case that FTRL has a nice interpretation as a local regret minimizer. The reason is that the prox mapping in (2) or (4) minimizes the sum of losses, rather than the most recent loss. Because of this, the expected value  $\langle \sum_{\tau=1}^t \ell_{\downarrow j}^\tau, \mathbf{x}_{\downarrow j}^{t+1} \rangle$  at simplex  $j$ , which influences the modified loss at parent simplexes, is computed based on  $\mathbf{x}^{t+1}$  for all  $t$  losses. Thus there is no local modified loss that could be received at rounds 1 through  $t$  that accurately reflects the modified loss needed in Proposition 1.

## 6. Experimental Evaluation

We experimentally evaluate the performance of optimistic regret minimization methods instantiated with dilated distance-generating functions. We experiment on three games:

- *Smallmatrix*, a small  $2 \times 2$  matrix game. Given a mixed strategy  $\mathbf{x} = (x_1, x_2) \in \Delta^2$  for Player 1 and a mixed strategy  $\mathbf{y} = (y_1, y_2) \in \Delta^2$  for Player 2, the payoff function for player 1 is  $u(\mathbf{x}, \mathbf{y}) = 5x_1y_1 - x_1y_2 + x_2y_2$ .
- *Kuhn poker*, already introduced in Section 3. In Kuhn poker, each player first has to put a payment of 1 into the pot. Each player is then dealt one of the three cards, and the third is put aside unseen. A single round of betting then occurs: first, Player 1 can check or bet 1. Then,
  - If Player 1 checks Player 2 can check or raise 1.
    - \* If Player 2 checks a showdown occurs; if Player 2 raises Player 1 can fold or call.
      - If Player 1 folds Player 2 takes the pot; if Player 1 calls a showdown occurs.
  - If Player 1 raises Player 2 can fold or call.
    - \* If Player 2 folds Player 1 takes the pot; if Player 2 calls a showdown occurs.

If no player has folded, a showdown occurs where the player with the higher card wins.

- *Leduc poker*, a standard benchmark in imperfect-information game solving (Southey et al., 2005). The game is played with a deck consisting of 5 unique cards with 2 copies of each, and consists of two rounds. In the first round, each player places an ante of 1 in the pot and receives a single private card. A round of betting then takes place with a two-bet maximum, with Player 1 going first. A public shared card is then dealt face up and another round of betting takes place. Again, Player 1 goes first, and there is a two-bet maximum. If one of the players has a pair with the public card, that player wins. Otherwise, the player



with the higher card wins. All bets in the first round are 1, while all bets in the second round are 2. This game has 390 decision points and 911 sequences per player.

**Fast Last-Iterate Convergence.** In the first set of experiments (Figure 2, top row), we compare the saddle-point gap of the strategy profiles produced by optimistic OMD and optimistic FTRL to that produced by CFR and CFR<sup>+</sup>. Optimistic OMD and optimistic FTRL were set up with the step-size parameter  $\eta = 0.1$  in Smallmatrix and  $\eta = 2$  in Kuhn Poker, and the plots show the last-iterate convergence for the optimistic algorithms, which has recently received attention in the works by Chambolle & Pock (2016) and Kroer (2019). Finally, we instantiated optimistic OMD and optimistic FTRL with the Euclidean distance generating function as constructed in Corollary 1. The plots show that—at least in these shallow games—optimistic methods are able to produce even up to 12 orders of magnitude better-approximate saddle-points than CFR and CFR<sup>+</sup>.

Interestingly, Smallmatrix appears to be a hard instance for CFR<sup>+</sup>: linear regression on the first 20 000 iterations of CFR<sup>+</sup> shows, with a coefficient of determination of roughly 0.96, that  $\log \xi(\mathbf{x}_*^T, \mathbf{y}_*^T) \approx -0.7375 \cdot \log(T) - 2.1349$ , where  $(\mathbf{x}_*^T, \mathbf{y}_*^T)$  is the average strategy profile (computed using linear averaging, as per CFR<sup>+</sup>'s construction) up to time  $T$ . In other words, we have evidence of at least one game in which the approximate saddle-point computed by CFR<sup>+</sup> experimentally has residual bounded below by  $\Omega(T^{-0.74})$ . This observation suggests that the analysis of CFR<sup>+</sup> might actually be quite tight, and that CFR<sup>+</sup> is *not* an accelerated method.

Figure 2 (bottom left) shows the performance of OFTRL in Leduc Poker, compared to CFR and CFR<sup>+</sup> (we do not show optimistic OMD, which we found to have worse performance than OFTRL). Here OFTRL performs worse than CFR<sup>+</sup>. This shows that in deeper games, more work has to be done to fully exploit the accelerated bounds of optimistic regret minimization methods.

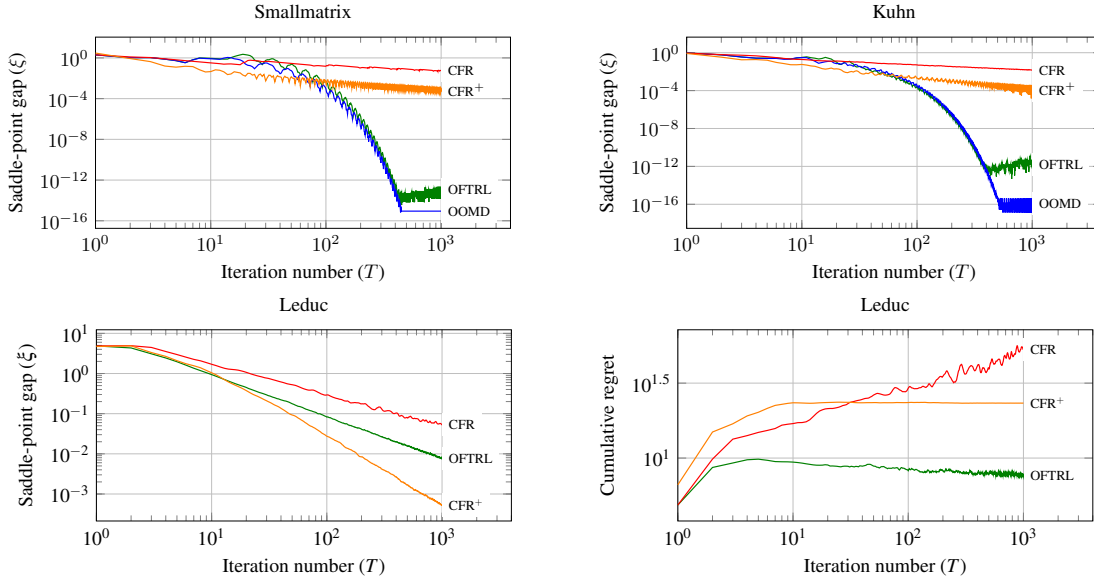


Figure 2. (Left and upper right) Saddle-point gap as a function of the number of iterations. The plots show the last-iterate convergence for OMD and OFTRL. (Lower right) Sum of cumulative regret for both players in Leduc. Optimistic OMD (OMD) and OFTRL use step-size parameter  $\eta = 0.1$  in Smallmatrix and  $\eta = 2$  in Kuhn. OFTRL uses step-size parameter  $\eta = 200$  in Leduc.

**Comparing the Cumulative Regret.** We also compared the algorithms based on the sum of cumulative regrets (again we omit optimistic OMD, which performed worse than OFTRL). In all three games, OFTRL leads to lower sum of cumulative regrets. Figure 2 (bottom right) shows the performance of OFTRL

in Leduc Poker. Here, we used the usual average of iterates  $\bar{x} := 1/T \sum_{t=1}^T x^t$  (note that the choice of averaging strategy has no effect on the bottom right plot.)

OFTRL’s performance matches the theory from Theorem 2 and Section 2.2. In particular, we observe that while OFTRL does not beat the state-of-the-art CFR<sup>+</sup> in terms of saddle-point gap, it beats it according to the regret sum metric. The fact that CFR<sup>+</sup> performs worse with respect to the regret sum metric is somewhat surprising: the entire derivation of CFR and CFR<sup>+</sup> is based on showing bounds on the regret sum. However, the connection between regret and saddle-point gap (or exploitability) is one-way: if the two regret minimizers (one per player) have regret  $R_1$  and  $R_2$ , then the saddle point gap can be easily shown to be less than or equal to  $(R_1 + R_2)/T$ . However, nothing prevents it from being much smaller than  $(R_1 + R_2)/T$ . What we empirically find is that for CFR<sup>+</sup> this bound is very loose. We are not sure why this is the case, and it potentially warrants further investigation in the future.

## 7. Conclusions

We studied how optimistic regret minimization can be applied in the context of extensive-form games, and introduced the first instantiations of regret-based techniques that achieve  $T^{-1}$  convergence to Nash equilibrium in extensive-form games. These methods rely crucially on having a tractable regularizer to maintain feasibility and control the stepsizes on the domain at hand—in our case, the sequence-form polytope. We provided the first explicit bound on the strong convexity properties of dilated distance-generating functions with respect to the Euclidean norm. We also showed that when optimistic regret minimization methods are instantiated with dilated distance-generating functions, the regret updates are local to each information set in the game, mirroring the structure of the counterfactual regret minimization framework. This localization of the updates along the tree structure enables further techniques, such as distributing the updates or skipping updates on cold parts of the game tree. Finally, when used in self play, these optimistic regret minimization methods guarantee an optimal  $T^{-1}$  convergence rate to Nash equilibrium.

We demonstrate that in shallow games, methods based on optimistic regret minimization can significantly outperform CFR and CFR<sup>+</sup>—even up to 12 orders of magnitude. In deeper games, more work has to be done to fully exploit the accelerated bounds of optimistic regret minimization methods. However, while the strong CFR<sup>+</sup> performance in large games remains a mystery, we elucidate some points about its performance—including showing that its theoretically slow convergence bound is somewhat tight. Finally, we showed that when the goal is minimizing regret, rather than computing a Nash equilibrium, optimistic methods can outperform CFR<sup>+</sup> even in deep game trees.

## Acknowledgments

This material is based on work supported by the National Science Foundation under grants IIS-1718457, IIS-1617590, and CCF-1733556, and the ARO under award W911NF-17-1-0082. Gabriele Farina is supported by a Facebook fellowship.

## References

- Bowling, M., Burch, N., Johanson, M., and Tammelin, O. Heads-up limit hold’em poker is solved. *Science*, 347(6218), January 2015.
- Brown, N. and Sandholm, T. Reduced space and faster convergence in imperfect-information games via pruning. In *International Conference on Machine Learning (ICML)*, 2017a.
- Brown, N. and Sandholm, T. Superhuman AI for heads-up no-limit poker: Libratus beats top professionals. *Science*, pp. eaao1733, Dec. 2017b.
- Brown, N. and Sandholm, T. Solving imperfect-information games via discounted regret minimization. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pp. 1829–1836, 2019a.

- Brown, N. and Sandholm, T. Superhuman AI for multiplayer poker. *Science*, 365(6456):885–890, 2019b. ISSN 0036-8075. doi: 10.1126/science.aay2400. URL <https://science.sciencemag.org/content/365/6456/885>.
- Brown, N., Ganzfried, S., and Sandholm, T. Hierarchical abstraction, distributed equilibrium computation, and post-processing, with application to a champion no-limit Texas Hold’em agent. In *International Conference on Autonomous Agents and Multi-Agent Systems (AAMAS)*, 2015.
- Brown, N., Kroer, C., and Sandholm, T. Dynamic thresholding and pruning for regret minimization. In *AAAI Conference on Artificial Intelligence (AAAI)*, 2017.
- Brown, N., Sandholm, T., and Amos, B. Depth-limited solving for imperfect-information games. *arXiv preprint arXiv:1805.08195*, 2018.
- Burch, N., Johanson, M., and Bowling, M. Solving imperfect information games using decomposition. In *AAAI Conference on Artificial Intelligence (AAAI)*, 2014.
- Chambolle, A. and Pock, T. A first-order primal-dual algorithm for convex problems with applications to imaging. *Journal of Mathematical Imaging and Vision*, 2011.
- Chambolle, A. and Pock, T. On the ergodic convergence rates of a first-order primal–dual algorithm. *Mathematical Programming*, 159(1-2):253–287, 2016.
- Chiang, C.-K., Yang, T., Lee, C.-J., Mahdavi, M., Lu, C.-J., Jin, R., and Zhu, S. Online optimization with gradual variations. In *Conference on Learning Theory*, pp. 6–1, 2012.
- Farina, G., Kroer, C., Brown, N., and Sandholm, T. Stable-predictive optimistic counterfactual regret minimization. In *International Conference on Machine Learning (ICML)*, 2019.
- Ganzfried, S. and Sandholm, T. Potential-aware imperfect-recall abstraction with earth mover’s distance in imperfect-information games. In *AAAI Conference on Artificial Intelligence (AAAI)*, 2014.
- Ganzfried, S. and Sandholm, T. Endgame solving in large imperfect-information games. In *International Conference on Autonomous Agents and Multi-Agent Systems (AAMAS)*, 2015.
- Gilpin, A. and Sandholm, T. Lossless abstraction of imperfect information games. *Journal of the ACM*, 54(5), 2007.
- Hazan, E. Introduction to online convex optimization. *Foundations and Trends in Optimization*, 2(3-4): 157–325, 2016.
- Hoda, S., Gilpin, A., Peña, J., and Sandholm, T. Smoothing techniques for computing Nash equilibria of sequential games. *Mathematics of Operations Research*, 35(2), 2010.
- Kroer, C. First-order methods with increasing iterate averaging for solving saddle-point problems. *arXiv preprint arXiv:1903.10646*, 2019.
- Kroer, C. and Sandholm, T. Extensive-form game abstraction with bounds. In *Proceedings of the ACM Conference on Economics and Computation (EC)*, 2014.
- Kroer, C. and Sandholm, T. Imperfect-recall abstractions with bounds in games. In *Proceedings of the ACM Conference on Economics and Computation (EC)*, 2016.
- Kroer, C., Waugh, K., Kılınç-Karzan, F., and Sandholm, T. Faster first-order methods for extensive-form game solving. In *Proceedings of the ACM Conference on Economics and Computation (EC)*, 2015.
- Kroer, C., Farina, G., and Sandholm, T. Solving large sequential games with the excessive gap technique. In *Proceedings of the Annual Conference on Neural Information Processing Systems (NIPS)*, 2018a.
- Kroer, C., Waugh, K., Kılınç-Karzan, F., and Sandholm, T. Faster algorithms for extensive-form game solving via improved smoothing functions. *Mathematical Programming*, pp. 1–33, 2018b.

- Kuhn, H. W. A simplified two-person poker. In Kuhn, H. W. and Tucker, A. W. (eds.), *Contributions to the Theory of Games*, volume 1 of *Annals of Mathematics Studies*, 24, pp. 97–103. Princeton University Press, Princeton, New Jersey, 1950.
- Lanctot, M., Gibson, R., Burch, N., Zinkevich, M., and Bowling, M. No-regret learning in extensive-form games with imperfect recall. In *International Conference on Machine Learning (ICML)*, 2012.
- Moravcik, M., Schmid, M., Ha, K., Hladik, M., and Gaukrodger, S. Refining subgames in large imperfect information games. In *AAAI Conference on Artificial Intelligence (AAAI)*, 2016.
- Moravčík, M., Schmid, M., Burch, N., Lisý, V., Morrill, D., Bard, N., Davis, T., Waugh, K., Johanson, M., and Bowling, M. Deepstack: Expert-level artificial intelligence in heads-up no-limit poker. *Science*, 356 (6337), May 2017.
- Nemirovski, A. Prox-method with rate of convergence  $O(1/t)$  for variational inequalities with Lipschitz continuous monotone operators and smooth convex-concave saddle point problems. *SIAM Journal on Optimization*, 15(1), 2004.
- Rakhlin, A. and Sridharan, K. Online learning with predictable sequences. In *Conference on Learning Theory*, pp. 993–1019, 2013a.
- Rakhlin, S. and Sridharan, K. Optimization, learning, and games with predictable sequences. In *Advances in Neural Information Processing Systems*, pp. 3066–3074, 2013b.
- Shalev-Shwartz, S. and Singer, Y. A primal-dual perspective of online learning algorithms. *Machine Learning*, 69(2-3):115–142, 2007.
- Southey, F., Bowling, M., Larson, B., Piccione, C., Burch, N., Billings, D., and Rayner, C. Bayes’ bluff: Opponent modelling in poker. In *Proceedings of the 21st Annual Conference on Uncertainty in Artificial Intelligence (UAI)*, July 2005.
- Syrkanis, V., Agarwal, A., Luo, H., and Schapire, R. E. Fast convergence of regularized learning in games. In *Advances in Neural Information Processing Systems*, pp. 2989–2997, 2015.
- Tammelin, O., Burch, N., Johanson, M., and Bowling, M. Solving heads-up limit Texas hold’em. In *Proceedings of the 24th International Joint Conference on Artificial Intelligence (IJCAI)*, 2015.
- von Stengel, B. Efficient computation of behavior strategies. *Games and Economic Behavior*, 14(2): 220–246, 1996.
- Zinkevich, M. Online convex programming and generalized infinitesimal gradient ascent. In *International Conference on Machine Learning (ICML)*, pp. 928–936, Washington, DC, USA, 2003.
- Zinkevich, M., Bowling, M., Johanson, M., and Piccione, C. Regret minimization in games with incomplete information. In *Proceedings of the Annual Conference on Neural Information Processing Systems (NIPS)*, 2007.

## A. Proofs: Optimistic Follow-the-Regularized-Leader

### A.1. Continuity of the Argmin-Function

Intuitively, the role of the regularizer  $d$  is to *smooth out* the linear objective function  $\langle \cdot, \mathbf{L} \rangle$ . So, it seems only reasonable to expect that, the higher the constant that multiplies  $d$ , the less the argmin  $\tilde{x}(\mathbf{L}) = \operatorname{argmin}_{\mathbf{x} \in \mathcal{X}} \langle \mathbf{L}, \mathbf{x} \rangle + \frac{1}{\eta} d(\mathbf{x})$  is affected by small changes in  $\mathbf{L}$ . In fact, the following holds:

**Lemma 1.** *Let  $d$  be 1-strongly convex with respect to a norm  $\| \cdot \|$ . The argmin-function  $\tilde{x}$  is  $\eta$ -Lipschitz continuous with respect to the dual norm  $\| \cdot \|_*$ , that is*

$$\|\tilde{x}(\mathbf{L}) - \tilde{x}(\mathbf{L}')\| \leq \eta \|\mathbf{L} - \mathbf{L}'\|_*.$$

*Proof.* The variational inequality for the optimality of  $\tilde{x}(\mathbf{L})$  implies

$$\left\langle \mathbf{L} + \frac{1}{\eta} \nabla d(\tilde{x}(\mathbf{L})), \tilde{x}(\mathbf{L}') - \tilde{x}(\mathbf{L}) \right\rangle \geq 0. \quad (8)$$

Symmetrically for  $\tilde{x}(\mathbf{L}')$ , we find that

$$\left\langle \mathbf{L}' + \frac{1}{\eta} d(\tilde{x}(\mathbf{L}')), \tilde{x}(\mathbf{L}) - \tilde{x}(\mathbf{L}') \right\rangle \geq 0. \quad (9)$$

Summing inequalities (8) and (9), we obtain

$$\frac{1}{\eta} \langle \nabla d(\tilde{x}(\mathbf{L})) - \nabla d(\tilde{x}(\mathbf{L}')), \tilde{x}(\mathbf{L}) - \tilde{x}(\mathbf{L}') \rangle \leq \langle \mathbf{L}' - \mathbf{L}, \tilde{x}(\mathbf{L}) - \tilde{x}(\mathbf{L}') \rangle.$$

Using the 1-strong convexity of  $d(\cdot)$  on the left-hand side and the generalized Cauchy-Schwarz inequality on the right-hand side, we obtain

$$\frac{1}{\eta} \|\tilde{x}(\mathbf{L}) - \tilde{x}(\mathbf{L}')\|^2 \leq \|\tilde{x}(\mathbf{L}) - \tilde{x}(\mathbf{L}')\| \|\mathbf{L} - \mathbf{L}'\|_*,$$

and dividing by  $\|\tilde{x}(\mathbf{L}) - \tilde{x}(\mathbf{L}')\|$  we obtain the Lipschitz continuity of the argmin-function  $\tilde{x}$ .  $\square$

Another observation that will be crucial in the analysis is that the objective function  $\langle \mathbf{L}, \mathbf{x} \rangle + (1/\eta)d(\mathbf{x})$  is  $(1/\eta)$ -strongly convex. Hence, for all  $\mathbf{L}$  and  $\hat{\mathbf{x}} \in \mathcal{X}$ ,

$$\left( \langle \mathbf{L}, \hat{\mathbf{x}} \rangle + \frac{1}{\eta} d(\hat{\mathbf{x}}) \right) \geq \left( \langle \mathbf{L}, \tilde{x}(\mathbf{L}) \rangle + \frac{1}{\eta} d(\tilde{x}(\mathbf{L})) \right) + \frac{1}{2\eta} \|\hat{\mathbf{x}} - \tilde{x}(\mathbf{L})\|^2. \quad (10)$$

#### A.1.1. THE OMNISCIENT CASE

The RVU property (Definition 1) implies that if  $\mathbf{m}^t = \boldsymbol{\ell}^t$  for all  $t = 1, \dots, T$ , then the regret cumulated is bounded above by a constant, independent on the time horizon  $T$ , and can only go down over time. We now show that this indeed holds for OFTRL. In what follows, we will use the notation  $\mathbf{L}^t$  to denote  $\mathbf{L}^t := \sum_{\tau=1}^t \boldsymbol{\ell}^\tau$  for all  $t \geq 1$ , and  $\mathbf{L}^0 := \mathbf{0}$ .

**Lemma 2.** *Let  $T > 0$ , and assume that OFTRL is set up so that  $\boldsymbol{\ell}^t = \mathbf{m}^t$  (i.e., the prediction is omniscient) for all  $t = 1, \dots, T$ . Furthermore, denote  $\mathbf{x}_o^t := \tilde{x}(\mathbf{L}^t)$  the decisions produces by OFTRL at all times  $t \geq 0$ . Then, the regret against any strategy  $\hat{\mathbf{x}} \in \mathcal{X}$  is bounded above as*

$$R^T(\hat{\mathbf{x}}) := \sum_{t=1}^T \langle \boldsymbol{\ell}^t, \mathbf{x}_o^t - \hat{\mathbf{x}} \rangle \leq \frac{1}{\eta} (d(\hat{\mathbf{x}}) - d(\mathbf{x}_o^0)) - \frac{1}{2\eta} \sum_{t=1}^T \|\mathbf{x}_o^t - \mathbf{x}_o^{t-1}\|^2.$$

As a direct consequence, if OFTRL is fed with exact predictions, the cumulated regret is bounded as

$$R^T = \max_{\hat{\mathbf{x}} \in \mathcal{X}} R^T(\hat{\mathbf{x}}) \leq \frac{\Delta_d}{\eta} - \frac{1}{2\eta} \sum_{t=1}^T \|\mathbf{x}_o^t - \mathbf{x}_o^{t-1}\|^2,$$

where  $\Delta_d := \max_{\mathbf{x}, \mathbf{y} \in \mathcal{X}} \{d(\mathbf{x}) - d(\mathbf{y})\}$ .

*Proof.* Equation (10) immediately implies that for all  $\hat{\mathbf{x}} \in \mathcal{X}$  and time  $t$ ,

$$\langle \mathbf{L}^t, \hat{\mathbf{x}} - \mathbf{x}_o^t \rangle \geq -\frac{1}{\eta} (d(\hat{\mathbf{x}}) - d(\mathbf{x}_o^t)) + \frac{1}{2\eta} \|\hat{\mathbf{x}} - \mathbf{x}_o^t\|^2. \quad (11)$$

Consequently, for all  $\hat{\mathbf{x}} \in \mathcal{X}$  we have

$$\begin{aligned} \langle \mathbf{L}^T, \mathbf{x}_o^T - \hat{\mathbf{x}} \rangle &= \langle \ell^T, \mathbf{x}_o^T - \hat{\mathbf{x}} \rangle + \langle \mathbf{L}^{T-1}, \mathbf{x}_o^{T-1} - \hat{\mathbf{x}} \rangle + \langle \mathbf{L}^{T-1}, \mathbf{x}_o^T - \mathbf{x}_o^{T-1} \rangle \\ &\geq \langle \mathbf{L}^{T-1}, \mathbf{x}_o^{T-1} - \hat{\mathbf{x}} \rangle + \langle \ell^T, \mathbf{x}_o^T - \hat{\mathbf{x}} \rangle - \frac{1}{\eta} (d(\mathbf{x}_o^T) - d(\mathbf{x}_o^{T-1})) + \frac{1}{2\eta} \|\mathbf{x}_o^T - \mathbf{x}_o^{T-1}\|^2, \end{aligned} \quad (12)$$

where the first inequality is by (11) applied to  $\hat{\mathbf{x}} = \mathbf{x}^T$  and  $t = T - 1$ . Using (12) inductively (recursively expanding terms of the form  $\langle \mathbf{L}^t, \mathbf{x}_o^t - \hat{\mathbf{x}} \rangle$  for all  $t = T - 1, T - 2, \dots, 1$ ), we obtain

$$\langle \mathbf{L}^T, \mathbf{x}_o^T - \hat{\mathbf{x}} \rangle \geq \left( \sum_{t=1}^T \langle \ell^t, \mathbf{x}_o^t - \hat{\mathbf{x}} \rangle \right) - \frac{1}{\eta} (d(\mathbf{x}_o^T) - d(\mathbf{x}_o^0)) + \frac{1}{2\eta} \sum_{t=1}^T \|\mathbf{x}_o^t - \mathbf{x}_o^{t-1}\|^2,$$

Finally, inverting the sides of the inequality, we find

$$\begin{aligned} R^T(\hat{\mathbf{x}}) &= \sum_{t=1}^T \langle \ell^t, \mathbf{x}_o^t - \hat{\mathbf{x}} \rangle \\ &\leq \langle \mathbf{L}^T, \mathbf{x}_o^T - \hat{\mathbf{x}} \rangle + \frac{1}{\eta} (d(\mathbf{x}_o^T) - d(\mathbf{x}_o^0)) - \frac{1}{2\eta} \sum_{t=1}^T \|\mathbf{x}_o^t - \mathbf{x}_o^{t-1}\|^2 \\ &\leq \frac{1}{\eta} (d(\hat{\mathbf{x}}) - d(\mathbf{x}_o^T)) + \frac{1}{\eta} (d(\mathbf{x}_o^T) - d(\mathbf{x}_o^0)) - \frac{1}{2\eta} \sum_{t=1}^T \|\mathbf{x}_o^t - \mathbf{x}_o^{t-1}\|^2 \\ &= \frac{1}{\eta} (d(\hat{\mathbf{x}}) - d(\mathbf{x}_o^0)) - \frac{1}{2\eta} \sum_{t=1}^T \|\mathbf{x}_o^t - \mathbf{x}_o^{t-1}\|^2 \end{aligned}$$

where the second inequality is by (11) applied to  $t = T$ .  $\square$

## A.2. Bridging the Gap: The Non-Omniscient Case

We are now ready to complete the proof of Theorem 2. The idea of the proof is to use the Lipschitz continuity of the argmin-function (Lemma 1) to show that the regret generated by the non-omniscience of the OFTRL (that is, the fact that  $\mathbf{m}^t \neq \ell^t$ ) is proportional to the distance  $\sum_t \|\mathbf{m}^t - \ell^t\|_*^2$ . A naive attempt would could be as follows:

$$\begin{aligned} \langle \ell^t, \mathbf{x}^t \rangle &= \langle \ell^t, \tilde{\mathbf{x}}(\mathbf{L}^t) \rangle + \langle \ell^t, \tilde{\mathbf{x}}(\mathbf{L}^{t-1} + \mathbf{m}^t) - \tilde{\mathbf{x}}(\mathbf{L}^t) \rangle \\ &\leq \langle \ell^t, \tilde{\mathbf{x}}(\mathbf{L}^t) \rangle + \eta \|\ell^t\|_* \|\ell^t - \mathbf{m}^t\|_*, \end{aligned}$$

where the inequality follows from from the generalized Cauchy-Schwarz inequality and from Lemma 1. Unfortunately, the approach above is not powerful enough, as it would imply that the increase in regret is proportional to  $\sum_t \|\mathbf{m}^t - \ell^t\|_*$  (notice the different exponent). In order to obtain the better bound, we write instead

$$\begin{aligned} \langle \ell^t, \mathbf{x}^t \rangle &= \langle \ell^t, \tilde{\mathbf{x}}(\mathbf{L}^t) \rangle + \langle \ell^t - \mathbf{m}^t, \tilde{\mathbf{x}}(\mathbf{L}^{t-1} + \mathbf{m}^t) - \tilde{\mathbf{x}}(\mathbf{L}^t) \rangle + \langle \mathbf{m}^t, \tilde{\mathbf{x}}(\mathbf{L}^{t-1} + \mathbf{m}^t) - \tilde{\mathbf{x}}(\mathbf{L}^t) \rangle \\ &\leq \langle \ell^t, \tilde{\mathbf{x}}(\mathbf{L}^t) \rangle + \eta \|\ell^t - \mathbf{m}^t\|_*^2 + \langle \mathbf{m}^t, \tilde{\mathbf{x}}(\mathbf{L}^{t-1} + \mathbf{m}^t) - \tilde{\mathbf{x}}(\mathbf{L}^t) \rangle. \end{aligned} \quad (13)$$

The first term corresponds to the case of omniscient predictions, which was already analyzed in the previous section. Hence, we are left with the task of bounding the loss from the last term. To this end, we will use the following lemma:

**Lemma 3.** Let  $\mathbf{x}^t := \tilde{\mathbf{x}}(\mathbf{L}^{t-1} + \mathbf{m}^t)$  and  $\mathbf{x}_o^t := \tilde{\mathbf{x}}(\mathbf{L}^t)$ . For all  $T$ ,

$$\sum_{t=1}^T \langle \mathbf{m}^t, \mathbf{x}^t - \mathbf{x}_o^t \rangle \leq -R_o^T + \frac{\Delta d}{\eta} - \frac{1}{4\eta} \sum_{t=1}^{T-1} \|\mathbf{x}^{t+1} - \mathbf{x}^t\|^2$$

where

$$R_o^T := \max_{\hat{\mathbf{x}} \in \mathcal{X}} \sum_{t=1}^T \langle \ell^t, \mathbf{x}_o^t - \hat{\mathbf{x}} \rangle$$

is the regret of the omniscient minimizer.

*Proof.* We will make use of the following two inequalities, which are direct applications of (10) where  $L$  is set to  $\mathbf{L}^{t-1} + \mathbf{m}^t$  and  $\mathbf{L}^{t-1}$ , respectively, and  $\hat{\mathbf{x}}$  is set to  $\mathbf{x}_o^t$  and  $\mathbf{x}^t$ , respectively:

$$\langle \mathbf{L}^{t-1} + \mathbf{m}^t, \mathbf{x}^t \rangle \leq \langle \mathbf{L}^{t-1} + \mathbf{m}^t, \mathbf{x}_o^t \rangle - \frac{1}{\eta} (d(\mathbf{x}^t) - d(\mathbf{x}_o^t)) - \frac{1}{2\eta} \|\mathbf{x}^t - \mathbf{x}_o^t\|^2, \quad (14)$$

$$\langle \mathbf{L}^{t-1}, \mathbf{x}_o^{t-1} \rangle \leq \langle \mathbf{L}^{t-1}, \mathbf{x}^t \rangle - \frac{1}{\eta} (d(\mathbf{x}_o^{t-1}) - d(\mathbf{x}^t)) - \frac{1}{2\eta} \|\mathbf{x}^t - \mathbf{x}_o^{t-1}\|^2. \quad (15)$$

With the above, we can write, for all timesteps  $t$ ,

$$\begin{aligned} \langle \mathbf{m}^t, \mathbf{x}^t - \mathbf{x}_o^t \rangle &= -\langle \mathbf{m}^t, \mathbf{x}_o^t \rangle + \langle \mathbf{L}^{t-1} + \mathbf{m}^t, \mathbf{x}^t \rangle - \langle \mathbf{L}^{t-1}, \mathbf{x}^t \rangle \\ &\leq -\langle \mathbf{m}^t, \mathbf{x}_o^t \rangle + \langle \mathbf{L}^{t-1} + \mathbf{m}^t, \mathbf{x}^t \rangle - \langle \mathbf{L}^{t-1}, \mathbf{x}_o^{t-1} \rangle - \frac{1}{\eta} (d(\mathbf{x}_o^{t-1}) - d(\mathbf{x}^t)) - \frac{1}{2\eta} \|\mathbf{x}^t - \mathbf{x}_o^{t-1}\|^2 \\ &\leq -\langle \mathbf{m}^t, \mathbf{x}_o^t \rangle + \langle \mathbf{L}^{t-1} + \mathbf{m}^t, \mathbf{x}_o^t \rangle - \langle \mathbf{L}^{t-1}, \mathbf{x}_o^{t-1} \rangle - \frac{1}{\eta} (d(\mathbf{x}^t) - d(\mathbf{x}_o^t)) - \frac{1}{\eta} (d(\mathbf{x}_o^{t-1}) - d(\mathbf{x}^t)) \\ &\quad - \frac{1}{2\eta} \|\mathbf{x}^t - \mathbf{x}_o^t\|^2 - \frac{1}{2\eta} \|\mathbf{x}^t - \mathbf{x}_o^{t-1}\|^2 \\ &= \langle \mathbf{L}^{t-1}, \mathbf{x}_o^t - \mathbf{x}_o^{t-1} \rangle + \frac{1}{\eta} (d(\mathbf{x}_o^t) - d(\mathbf{x}_o^{t-1})) - \frac{1}{2\eta} \|\mathbf{x}^t - \mathbf{x}_o^t\|^2 - \frac{1}{2\eta} \|\mathbf{x}^t - \mathbf{x}_o^{t-1}\|^2, \end{aligned}$$

where the first inequality follows from (15) and the second inequality follows from (14). Note that

$$\langle \mathbf{L}^{t-1}, \mathbf{x}_o^t - \mathbf{x}_o^{t-1} \rangle = \langle \mathbf{L}^t, \mathbf{x}_o^t \rangle - \langle \mathbf{L}^{t-1}, \mathbf{x}_o^{t-1} \rangle - \langle \ell^t, \mathbf{x}_o^t \rangle.$$

Hence,

$$\langle \mathbf{m}^t, \mathbf{x}^t - \mathbf{x}_o^t \rangle \leq \langle \mathbf{L}^t, \mathbf{x}_o^t \rangle - \langle \mathbf{L}^{t-1}, \mathbf{x}_o^{t-1} \rangle - \langle \ell^t, \mathbf{x}_o^t \rangle + \frac{1}{\eta} (d(\mathbf{x}_o^t) - d(\mathbf{x}_o^{t-1})) - \frac{1}{2\eta} \|\mathbf{x}^t - \mathbf{x}_o^t\|^2 - \frac{1}{2\eta} \|\mathbf{x}^t - \mathbf{x}_o^{t-1}\|^2.$$

Finally, summing over all  $t$ , we obtain

$$\begin{aligned} \sum_{t=1}^T \langle \mathbf{m}^t, \mathbf{x}^t - \mathbf{x}_o^t \rangle &\leq \langle \mathbf{L}^T, \mathbf{x}_o^T \rangle - \sum_{t=1}^T \langle \ell^t, \mathbf{x}_o^t \rangle + \frac{1}{\eta} (d(\mathbf{x}_o^T) - d(\mathbf{x}_o^0)) - \frac{1}{2\eta} \left( \sum_{t=1}^T \|\mathbf{x}^t - \mathbf{x}_o^t\|^2 + \sum_{t=1}^T \|\mathbf{x}_o^t - \mathbf{x}_o^{t-1}\|^2 \right) \\ &\leq \langle \mathbf{L}^T, \mathbf{x}_o^T \rangle - \sum_{t=1}^T \langle \ell^t, \mathbf{x}_o^t \rangle + \frac{1}{\eta} (d(\mathbf{x}_o^T) - d(\mathbf{x}_o^0)) - \frac{1}{2\eta} \left( \sum_{t=1}^{T-1} \|\mathbf{x}^t - \mathbf{x}_o^t\|^2 + \sum_{t=1}^{T-1} \|\mathbf{x}^{t+1} - \mathbf{x}_o^t\|^2 \right) \\ &\leq \langle \mathbf{L}^T, \mathbf{x}_o^T \rangle - \sum_{t=1}^T \langle \ell^t, \mathbf{x}_o^t \rangle + \frac{1}{\eta} (d(\mathbf{x}_o^T) - d(\mathbf{x}_o^0)) - \frac{1}{4\eta} \sum_{t=1}^{T-1} \|\mathbf{x}^{t+1} - \mathbf{x}^t\|^2, \quad (16) \end{aligned}$$

where the first inequality comes from simplifying telescopic sums, the second by removing some term from the last parenthesis, and the third from the parallelogram inequality

$$\|\mathbf{a}\|^2 + \|\mathbf{b}\|^2 \geq \frac{1}{2} \|\mathbf{a} - \mathbf{b}\|^2,$$

which holds for all choice of vectors  $\mathbf{a}, \mathbf{b}, \mathbf{c}$  and norm  $\|\cdot\|$ . The last step of the proof is to notice that by (11)

$$\begin{aligned} \langle \mathbf{L}^T, \mathbf{x}_o^T \rangle - \sum_{t=1}^T \langle \boldsymbol{\ell}^t, \mathbf{x}_o^t \rangle &\leq \langle \mathbf{L}^T, \hat{\mathbf{x}} \rangle - \sum_{t=1}^T \langle \boldsymbol{\ell}^t, \mathbf{x}_o^t \rangle + \frac{1}{\eta} (d(\hat{\mathbf{x}}) - d(\mathbf{x}_o^T)) \\ &= -R_o^T(\hat{\mathbf{x}}) + \frac{1}{\eta} (d(\hat{\mathbf{x}}) - d(\mathbf{x}_o^T)) \end{aligned}$$

for all  $\hat{\mathbf{x}} \in \mathcal{X}$ . Substituting into (16) yields

$$\begin{aligned} \sum_{t=1}^T \langle \mathbf{m}^t, \mathbf{x}^t - \mathbf{x}_o^t \rangle &\leq -R_o^T(\hat{\mathbf{x}}) + \frac{1}{\eta} (d(\hat{\mathbf{x}}) - d(\mathbf{x}_o^0)) - \frac{1}{4\eta} \sum_{t=1}^{T-1} \|\mathbf{x}^{t+1} - \mathbf{x}^t\|^2 \\ &= -R_o^T + \frac{\Delta_d}{\eta} - \frac{1}{4\eta} \sum_{t=1}^{T-1} \|\mathbf{x}^{t+1} - \mathbf{x}^t\|^2, \end{aligned}$$

as we wanted to show.  $\square$

Finally, using Lemma 3 together with (13) and Lemma 2, we obtain the proof of Theorem 2.

**Theorem 2.** *For all step-size parameters  $\eta > 0$ , OFTRL satisfies the RVU conditions with respect to any primal-dual norm pair  $(\|\cdot\|, \|\cdot\|_*)$  with parameters  $\alpha = \Delta_d/\eta, \beta = \eta, \gamma = 1/(4\eta)$ , where  $\Delta_d := \max_{\mathbf{x}, \mathbf{y} \in \mathcal{X}} \{d(\mathbf{x}) - d(\mathbf{y})\}$ .*

*Proof of Theorem 2.* As we already mentioned, the fundamental idea of the proof is to bound how much regret is generated from by the fact that the predictions  $\mathbf{m}^t$  are *not* perfect. Hence, the first step is to relate the notion of regret to the fact that the decisions produced are  $\mathbf{x}^t$  instead of  $\mathbf{x}_o^t$ :

$$\begin{aligned} R^T &= \max_{\hat{\mathbf{x}} \in \mathcal{X}} \left\{ \sum_{t=1}^T \langle \boldsymbol{\ell}^t, \mathbf{x}^t - \hat{\mathbf{x}} \rangle \right\} \\ &= \max_{\hat{\mathbf{x}} \in \mathcal{X}} \left\{ \sum_{t=1}^T \langle \boldsymbol{\ell}^t, \mathbf{x}_o^t - \hat{\mathbf{x}} \rangle \right\} + \sum_{t=1}^T \langle \boldsymbol{\ell}^t - \mathbf{m}^t, \mathbf{x}^t - \mathbf{x}_o^t \rangle + \sum_{t=1}^T \langle \mathbf{m}^t, \mathbf{x}^t - \mathbf{x}_o^t \rangle. \end{aligned}$$

The sum inside of the max is bounded according to the analysis of the omniscient case (Lemma 2). Using the (generalized) Cauchy-Schwarz inequality and the Lipschitz continuity of  $\tilde{x}$ , the sum in the middle can be easily bounded as

$$\sum_{t=1}^T \langle \boldsymbol{\ell}^t - \mathbf{m}^t, \mathbf{x}^t - \mathbf{x}_o^t \rangle \leq \eta \sum_{t=1}^T \|\boldsymbol{\ell}^t - \mathbf{m}^t\|_*^2.$$

Finally, the third summation was studied in Lemma 3. Putting everything together, we conclude that

$$\begin{aligned} R^T &\leq R_o^T + \eta \sum_{t=1}^T \|\boldsymbol{\ell}^t - \mathbf{m}^t\|_*^2 - R_o^T + \frac{\Delta_d}{\eta} - \frac{1}{4\eta} \sum_{t=1}^T \|\mathbf{x}^t - \mathbf{x}^{t-1}\|^2 \\ &\leq \frac{\Delta_d}{\eta} + \eta \sum_{t=1}^T \|\boldsymbol{\ell}^t - \mathbf{m}^t\|_*^2 - \frac{1}{4\eta} \sum_{t=1}^T \|\mathbf{x}^t - \mathbf{x}^{t-1}\|^2. \end{aligned}$$

$\square$

## B. Proofs: Dilated Entropy Distance Generating Function

In the proof of Theorem 3, we will make use of the following, useful fact.

**Lemma 4.** *Let  $d(\mathbf{x}) = \sum_{j \in \mathcal{J}} x_{p_j} d_j(x_j/x_{p_j})$  be a dilated DGF. Then, for all  $\mathbf{x}$  we have  $\langle \nabla d(\mathbf{x}), \mathbf{x} \rangle = d(\mathbf{x})$ .*



*Proof.* With some simple algebra, it is easy to show that the partial derivatives with respect to  $x_{ja}$  is

$$\frac{\partial d}{\partial x_{ja}}(\mathbf{x}) = \frac{\partial r_j}{\partial x_{ja}}\left(\frac{\mathbf{x}_j}{x_{p_j}}\right) + \sum_{j' \in C_{ja}} d_{j'}\left(\frac{\mathbf{x}_{j'}}{x_{ja}}\right) - \sum_{j' \in C_{ja}} \left\langle \nabla d_{j'}\left(\frac{\mathbf{x}_{j'}}{x_{ja}}\right), \frac{\mathbf{x}_{j'}}{x_{ja}} \right\rangle,$$

and that the partial derivative of  $d$  with respect to  $\mathbf{x}_\phi$  is

$$\frac{\partial d}{\partial \mathbf{x}_\phi}(\mathbf{x}) = d_r\left(\frac{\mathbf{x}_r}{x_\phi}\right) - \left\langle \nabla d_r\left(\frac{\mathbf{x}_r}{x_\phi}\right), \frac{\mathbf{x}_r}{x_\phi} \right\rangle.$$

As a consequence, for all  $j \in \mathcal{J}$ ,

$$\langle \nabla_{\mathbf{x}_j} d(\mathbf{x}), \mathbf{x}_j \rangle = \left\langle \nabla d_j\left(\frac{\mathbf{x}_j}{x_{p_j}}\right), \mathbf{x}_j \right\rangle - \sum_{a \in A_j} \sum_{j' \in C_{ja}} \left\langle \nabla d_{j'}\left(\frac{\mathbf{x}_{j'}}{x_{ja}}\right), \mathbf{x}_{j'} \right\rangle + \sum_{a \in A_j} \sum_{j' \in C_{ja}} x_{ja} d_{j'}\left(\frac{\mathbf{x}_{j'}}{x_{ja}}\right)$$

Hence,

$$\begin{aligned} \langle \nabla d(\mathbf{x}), \mathbf{x} \rangle &= \frac{\partial d}{\partial x_\phi}(\mathbf{x}) + \sum_j \langle \nabla_{\mathbf{x}_j} d(\mathbf{x}), \mathbf{x}_j \rangle \\ &= d_r\left(\frac{\mathbf{x}_r}{x_\phi}\right) + \sum_{j \in \mathcal{J}} \sum_{a \in A_j} \sum_{j' \in C_{ja}} x_{ja} d_{j'}\left(\frac{\mathbf{x}_{j'}}{x_{ja}}\right) \\ &= \sum_{j \in \mathcal{J}} x_{p_j} d_j\left(\frac{\mathbf{x}_j}{x_{p_j}}\right) = d(\mathbf{x}), \end{aligned}$$

where the second inequality comes from the observation that the inner products with the gradients cancel out.  $\square$

**Theorem 3.** Let  $d(\mathbf{x}) = \sum_{j \in \mathcal{J}} x_{p_j} d_j(\mathbf{x}_j/x_{p_j})$  where for all  $j$ ,  $d_j$  is  $\mu_j$ -strongly convex with respect to the Euclidean norm over  $\Delta^{n_j}$ . Furthermore, define  $\sigma_{ja} := \frac{\mu_j}{2} - \sum_{j' \in C_{ja}} \mu_{j'}$ , and  $\bar{\sigma} := \min_{ja} \sigma_{ja}$ . Then,  $d$  is  $\bar{\sigma}$ -strongly convex with respect to the Euclidean norm over  $\mathcal{X}$ .

*Proof.* Let  $x, y$  be arbitrary vectors in  $\in \mathcal{X}$ . Using the strong convexity of  $d_j$  and the fact that  $\mathbf{x} \in \mathcal{X}$ , we obtain for all  $j \in \mathcal{J}$

$$\begin{aligned} \sum_{a \in A_j} \sum_{j' \in C_{ja}} y_{ja} d_{j'}\left(\frac{\mathbf{y}_{j'}}{y_{ja}}\right) &\geq \sum_{a \in A_j} \sum_{j' \in C_{ja}} y_{ja} \cdot \left( d_{j'}\left(\frac{\mathbf{x}_{j'}}{x_{ja}}\right) + \left\langle \nabla d_{j'}\left(\frac{\mathbf{x}_{j'}}{x_{ja}}\right), \frac{\mathbf{y}_{j'}}{y_{ja}} - \frac{\mathbf{x}_{j'}}{x_{ja}} \right\rangle + \frac{\mu_{j'}}{2} \left\| \frac{\mathbf{y}_{j'}}{y_{ja}} - \frac{\mathbf{x}_{j'}}{x_{ja}} \right\|_2^2 \right) \\ &= \left( \sum_{a \in A_j} \sum_{j' \in C_{ja}} y_{ja} d_{j'}\left(\frac{\mathbf{x}_{j'}}{x_{ja}}\right) - \sum_{a \in A_j} \sum_{j' \in C_{ja}} y_{ja} \left\langle \nabla d_{j'}\left(\frac{\mathbf{x}_{j'}}{x_{ja}}\right), \frac{\mathbf{x}_{j'}}{x_{ja}} \right\rangle \right) \\ &\quad - \left\langle \nabla d_j\left(\frac{\mathbf{x}_j}{x_{p_j}}\right), \mathbf{y}_j \right\rangle + \sum_{a \in A_j} \sum_{j' \in C_{ja}} \left\langle \nabla d_{j'}\left(\frac{\mathbf{x}_{j'}}{x_{ja}}\right), \mathbf{y}_{j'} \right\rangle \\ &\quad + \sum_{a \in A_j} \sum_{j' \in C_{ja}} \frac{\mu_{j'}}{2} y_{ja} \left\| \frac{\mathbf{y}_{j'}}{y_{ja}} - \frac{\mathbf{x}_{j'}}{x_{ja}} \right\|_2^2 \\ &= \langle \nabla_{\mathbf{x}_j} d(\mathbf{x}), \mathbf{y}_j \rangle - \left\langle \nabla d_j\left(\frac{\mathbf{x}_j}{x_{p_j}}\right), \mathbf{y}_j \right\rangle + \sum_{a \in A_j} \sum_{j' \in C_{ja}} \left\langle \nabla d_{j'}\left(\frac{\mathbf{x}_{j'}}{x_{ja}}\right), \mathbf{y}_{j'} \right\rangle \\ &\quad + \sum_{a \in A_j} \sum_{j' \in C_{ja}} \frac{\mu_{j'}}{2} y_{ja} \left\| \frac{\mathbf{y}_{j'}}{y_{ja}} - \frac{\mathbf{x}_{j'}}{x_{ja}} \right\|_2^2. \end{aligned}$$

Summing over all  $j$ , we obtain

$$\begin{aligned}
 d(\mathbf{y}) &\geq \langle \nabla d(\mathbf{x}), \mathbf{y} \rangle + \sum_j \frac{\mu_j}{2} y_{p_j} \left\| \frac{\mathbf{y}_j}{y_{p_j}} - \frac{\mathbf{x}_j}{x_{p_j}} \right\|_2^2 \\
 &= d(\mathbf{x}) + \langle \nabla d(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \sum_j \frac{\mu_j}{2} y_{p_j} \left\| \frac{\mathbf{y}_j}{y_{p_j}} - \frac{\mathbf{x}_j}{x_{p_j}} \right\|_2^2 \\
 &\geq d(\mathbf{x}) + \langle \nabla d(\mathbf{x}), \mathbf{y} - \mathbf{x} \rangle + \sum_j \frac{\mu_j}{2} y_{p_j}^2 \left\| \frac{\mathbf{y}_j}{y_{p_j}} - \frac{\mathbf{x}_j}{x_{p_j}} \right\|_2^2, \tag{17}
 \end{aligned}$$

where the equality comes from Lemma 4 and the second inequality from the fact that  $y_{p_j} \in [0, 1]$ . Note that for all  $j$  and for all  $\mathbf{x}_j, \mathbf{y}_j$ ,

$$\begin{aligned}
 \frac{1}{2} \|\mathbf{y}_j - \mathbf{x}_j\|_2^2 &= \frac{1}{2} \left\| y_{p_j} \left( \frac{\mathbf{y}_j}{y_{p_j}} - \frac{\mathbf{x}_j}{x_{p_j}} \right) + (y_{p_j} - x_{p_j}) \frac{\mathbf{x}_j}{x_{p_j}} \right\|_2^2 \\
 &\leq y_{p_j}^2 \left\| \frac{\mathbf{y}_j}{y_{p_j}} - \frac{\mathbf{x}_j}{x_{p_j}} \right\|_2^2 + (y_{p_j} - x_{p_j})^2 \left\| \frac{\mathbf{x}_j}{x_{p_j}} \right\|_2^2 \\
 &\leq y_{p_j}^2 \left\| \frac{\mathbf{y}_j}{y_{p_j}} - \frac{\mathbf{x}_j}{x_{p_j}} \right\|_2^2 + (y_{p_j} - x_{p_j})^2, \tag{18}
 \end{aligned}$$

where the first inequality is by parallelogram inequality, while the second inequality follows from the fact that  $\mathbf{x}_j/x_{p_j} \in \Delta^{n_j}$  and therefore it has norm upper-bounded by 1. By inverting (18), we obtain

$$\begin{aligned}
 \sum_j \mu_j y_{p_j}^2 \left\| \frac{\mathbf{y}_j}{y_{p_j}} - \frac{\mathbf{x}_j}{x_{p_j}} \right\|_2^2 &\geq \frac{1}{2} \sum_j \mu_j \|\mathbf{y}_j - \mathbf{x}_j\|_2^2 - \sum_j \mu_j (y_{p_j} - x_{p_j})^2 \\
 &= \sum_j \sum_{a \in A_j} \left( \frac{\mu_j}{2} - \sum_{j' \in C_{j,a}} \mu_{j'} \right) (y_{j,a} - x_{j,a})^2.
 \end{aligned}$$

Finally, plugging into (17)

$$\sum_j \mu_j y_{p_j} \left\| \frac{\mathbf{y}_j}{y_{p_j}} - \frac{\mathbf{x}_j}{x_{p_j}} \right\|_2^2 \geq \bar{\sigma} \sum_j \sum_{a \in A_j} (x_{j,a} - y_{j,a})^2 = \bar{\sigma} \|\mathbf{y} - \mathbf{x}\|_2^2,$$

and we conclude that  $d$  is  $\bar{\sigma}$ -strongly convex with respect to the Euclidean norm over  $\mathcal{X}$ , like we wanted to show.  $\square$

### C. Proofs: Local Regret Minimization

**Proposition 1** (Decomposition into local prox mappings). *A prox mapping (7) on a treeplex with a Bregman divergence constructed from a dilated DGF decomposes into local prox mappings at each decision point  $j$  where the solution is as follows:*

$$\mathbf{x}_j^* = x_{p_j} \cdot \operatorname{argmin}_{\mathbf{b}_j \in \Delta^{n_j}} \left\{ \langle \hat{\mathbf{g}}_j, \mathbf{b}_j \rangle + D_j \left( \mathbf{b}_j \left\| \frac{\hat{\mathbf{x}}_j}{\hat{\mathbf{x}}_{p_j}} \right\| \right) \right\},$$

where

$$\hat{g}_{j,a} = g_{j,a} + \sum_{j' \in C_{j,a}} \left[ d_{\downarrow j'}^* (-\mathbf{g}_{\downarrow j'} + \nabla d_{\downarrow j'}(\hat{\mathbf{x}}_{\downarrow j'})) - d_{j'} \left( \frac{\hat{\mathbf{x}}_j}{\hat{\mathbf{x}}_{p_j}} \right) + \left\langle \nabla d_{j'} \left( \frac{\hat{\mathbf{x}}_{j'}}{\hat{\mathbf{x}}_{p_{j'}}} \right), \frac{\hat{\mathbf{x}}_{j'}}{\hat{\mathbf{x}}_{p_{j'}}} \right\rangle \right].$$

*Proof.* We will show the decomposition assuming that  $j$  is the root decision point, the general case follows by induction. The prox mapping (7) can be written as

$$\begin{aligned} \text{Prox}(\mathbf{g}, \hat{\mathbf{x}}) &= \underset{\mathbf{x} \in \mathcal{X}}{\text{argmin}} \{ \langle \mathbf{g}, \mathbf{x} \rangle + D(\mathbf{x} \parallel \hat{\mathbf{x}}) \} \\ &= \underset{\mathbf{x} \in \mathcal{X}}{\text{argmin}} \{ \langle \mathbf{g}, \mathbf{x} \rangle + d(\mathbf{x}) - d(\hat{\mathbf{x}}) - \langle \nabla d(\hat{\mathbf{x}}), \mathbf{x} - \hat{\mathbf{x}} \rangle \} \\ &= \underset{\mathbf{x} \in \mathcal{X}}{\text{argmin}} \{ \langle \mathbf{g} - \nabla d(\hat{\mathbf{x}}), \mathbf{x} \rangle + d(\mathbf{x}) \}. \end{aligned}$$

Now we can use the fact that we are using a dilated DGF and the sequential structure of SDM problems to write the problem of finding just the minimizer for  $j$  as follows

$$\begin{aligned} &\underset{\mathbf{b}_j \in \Delta^{n_j}}{\text{argmin}} \left\{ \langle \mathbf{g}_j - \nabla_j d(\hat{\mathbf{x}}), \mathbf{b}_j \rangle + d_j(\mathbf{b}_j) + \sum_{a \in A_j} \sum_{j' \in \mathcal{C}_{j,a}} b_{j,a} \left[ \min_{\mathbf{x} \in \mathcal{X}_{\downarrow j'}} \langle \mathbf{g}_{\downarrow j'} - \nabla d_{\downarrow j'}(\hat{\mathbf{x}}_{\downarrow j'}), \mathbf{x} \rangle + d_{\downarrow j'}(\mathbf{x}) \right] \right\} \\ &= \underset{\mathbf{b}_j \in \Delta^{n_j}}{\text{argmin}} \left\{ \langle \mathbf{g}_j - \nabla_j d(\hat{\mathbf{x}}), \mathbf{b}_j \rangle + d_j(\mathbf{b}_j) - \sum_{a \in A_j} \sum_{j' \in \mathcal{C}_{j,a}} b_{j,a} d_{\downarrow j'}^* (-\mathbf{g}_{\downarrow j'} + \nabla d_{\downarrow j'}(\hat{\mathbf{x}}_{\downarrow j'})) \right\}. \end{aligned}$$

Now we note that the index of  $\nabla_j d(\hat{\mathbf{x}})$  corresponding to each  $a \in A_j$  can be expanded as

$$\nabla_{j,a} d(\hat{\mathbf{x}}) = \nabla_a d_j \left( \frac{\hat{\mathbf{x}}_j}{\hat{x}_{p_j}} \right) + \sum_{j' \in \mathcal{C}_{j,a}} \left( d_{j'} \left( \frac{\hat{\mathbf{x}}_{j'}}{\hat{x}_{p_{j'}}} \right) - \left\langle \nabla d_{j'} \left( \frac{\hat{\mathbf{x}}_{j'}}{\hat{x}_{p_{j'}}} \right), \frac{\hat{\mathbf{x}}_{j'}}{\hat{x}_{p_{j'}}} \right\rangle \right).$$

Plugging this into our expression for the minimizer for  $j$ , and noting that we can add and remove  $d_j \left( \frac{\hat{\mathbf{x}}_j}{\hat{x}_{p_j}} \right)$  without changing the argmin, gives

$$\begin{aligned} &\underset{\mathbf{b}_j \in \Delta^{n_j}}{\text{argmin}} \left\{ \langle \mathbf{g}_j, \mathbf{b}_j \rangle + D_j \left( \mathbf{b}_j \parallel \frac{\hat{\mathbf{x}}_j}{\hat{x}_{p_j}} \right) - \sum_{a \in A_j} \sum_{j' \in \mathcal{C}_{j,a}} b_{j,a} \left[ d_{\downarrow j'}^* (-\mathbf{g}_{\downarrow j'} + \nabla d_{\downarrow j'}(\hat{\mathbf{x}}_{\downarrow j'})) - d_{j'} \left( \frac{\hat{\mathbf{x}}_j}{\hat{x}_{p_j}} \right) \right. \right. \\ &\quad \left. \left. + \left\langle \nabla d_{j'} \left( \frac{\hat{\mathbf{x}}_{j'}}{\hat{x}_{p_{j'}}} \right), \frac{\hat{\mathbf{x}}_{j'}}{\hat{x}_{p_{j'}}} \right\rangle \right] \right\} \\ &= \underset{\mathbf{b}_j \in \Delta^{n_j}}{\text{argmin}} \langle \hat{\mathbf{g}}_j, \mathbf{b}_j \rangle + D_j \left( \mathbf{b}_j \parallel \frac{\hat{\mathbf{x}}_j}{\hat{x}_{p_j}} \right) \end{aligned}$$

as we wanted to show.  $\square$

**Theorem 4.** *OMD with a dilated DGF for a treplex  $\mathcal{X}$  corresponds to running OMD locally at each simplex  $j$ , with the local loss  $\hat{\ell}^t$  constructed according to Proposition 1. Optimistic OMD corresponds to the optimistic variant of this local OMD with local loss predictions  $\hat{\ell}^t, \hat{\mathbf{m}}_j^{t+1}$  again constructed according to Proposition 1 using  $\mathbf{x}^t$  as Bregman divergence center and  $\mathbf{x}^{t+1}$  for aggregating losses below each simplex. Here the modified loss uses  $\mathbf{z}_{\downarrow j}^t$  and  $\mathbf{x}^{t+1}$  as Bregman divergence center and aggregating loss below, respectively. The prediction  $\hat{\mathbf{m}}_j^{t+1}$  uses  $\mathbf{z}_{\downarrow j}^t$  and  $\mathbf{z}^{t+1}$ .*

*Proof.* Note that the OMD update in (1) is a prox mapping as given in Proposition 1. The update is:

$$\mathbf{x}^{t+1} = \underset{\mathbf{x} \in \mathcal{X}}{\text{argmin}} \left\{ \langle \hat{\ell}^t, \mathbf{x} \rangle + \frac{1}{\eta} D(\mathbf{x} \parallel \mathbf{x}^t) \right\} \Rightarrow \mathbf{x}_j^{t+1} = \mathbf{x}_{p_j}^{t+1} \underset{\mathbf{b}_j \in \Delta^{n_j}}{\text{argmin}} \left\{ \langle \hat{\ell}^t, \mathbf{x} \rangle + \frac{1}{\eta} D_j \left( \mathbf{b}_j \parallel \frac{\mathbf{x}^t}{x_{p_j}^t} \right) \right\}. \quad (19)$$

But this is exactly the same as the OMD update resulting from running OMD on simplex  $j$  with the modified loss, using DGF  $d_j$  and its associated Bregman divergence. The same logic shows that Optimistic OMD is a local variant of itself, since it is a repeated sequence of two prox mappings.  $\square$