

Minimizing Regret on Reflexive Banach Spaces and Nash Equilibria in Continuous Zero-Sum Games

Authored by:

Walid Krichene
Alexandre Bayen
Maximilian Balandat
Claire Tomlin

Abstract

We study a general adversarial online learning problem, in which we are given a decision set X' in a reflexive Banach space X and a sequence of reward vectors in the dual space of X . At each iteration, we choose an action from X' , based on the observed sequence of previous rewards. Our goal is to minimize regret, defined as the gap between the realized reward and the reward of the best fixed action in hindsight. Using results from infinite dimensional convex analysis, we generalize the method of Dual Averaging (or Follow the Regularized Leader) to our setting and obtain upper bounds on the worst-case regret that generalize many previous results. Under the assumption of uniformly continuous rewards, we obtain explicit regret bounds in a setting where the decision set is the set of probability distributions on a compact metric space S . Importantly, we make no convexity assumptions on either the set S or the reward functions. We also prove a general lower bound on the worst-case regret for any online algorithm. We then apply these results to the problem of learning in repeated two-player zero-sum games on compact metric spaces. In doing so, we first prove that if both players play a Hannan-consistent strategy, then with probability 1 the empirical distributions of play weakly converge to the set of Nash equilibria of the game. We then show that, under mild assumptions, Dual Averaging on the (infinite-dimensional) space of probability distributions indeed achieves Hannan-consistency.

1 Paper Body

Regret analysis is a general technique for designing and analyzing algorithms for sequential decision problems in adversarial or stochastic settings (Shalev-Shwartz, 2012; Bubeck and Cesa-Bianchi, 2012). Online learning algorithms

have applications in machine learning (Xiao, 2010), portfolio optimization (Cover, 1991), online convex optimization (Hazan et al., 2007) and other areas. Regret analysis also plays an important role in the study of repeated play of finite games (Hart and MasColell, 2001). It is well known, for example, that in a two-player zero-sum finite game, if both players play according to a Hannan-consistent strategy (Hannan, 1957), their (marginal) empirical distributions of play almost surely converge to the set of Nash equilibria of the game (Cesa-Bianchi and Lugosi, 2006). Moreover, it can be shown that playing a strategy that achieves sublinear regret almost surely guarantees Hannan-consistency. A natural question then is whether a similar result holds for games with infinite action sets. In this article we provide a positive answer. In particular, we prove that in a continuous two-player zero sum game over compact (not necessarily convex) metric spaces, if both players follow a Hannan-consistent strategy, then with probability 1 their empirical distributions of play weakly converge to the set of Nash equilibria of the game. This in turn raises another important question: Do algorithms that ensure Hannan-consistency exist in such a setting? More generally, can one develop algorithms that guarantee sub-linear growth of the worst-case regret? We answer these questions affirmatively as well. To this end, we develop a general framework to study the Dual Averaging (or Follow the Regularized Leader) method on reflexive Banach spaces. This framework generalizes a wide range of existing 30th Conference on Neural Information Processing Systems (NIPS 2016), Barcelona, Spain.

results in the literature, including algorithms for online learning on finite sets (Arora et al., 2012) and finite-dimensional online convex optimization (Hazan et al., 2007). Given a convex subset X of a reflexive Banach space X , the generalized Dual Averaging (DA) method maximizes, at each iteration, the cumulative past rewards (which are elements of X^* , the dual space of X) minus a regularization term h . We show that under certain conditions, the maximizer in the DA update is the Fenchel gradient ∂h of the regularizer's conjugate function. In doing so, we develop a novel characterization of the duality between essential strong convexity of h and essential Fenchel differentiability of h^* in reflexive Banach spaces, which is of independent interest. We apply these general results to the problem of minimizing regret when the rewards are uniformly continuous functions over a compact metric space S . Importantly, we do not assume convexity of either S or the rewards, and show that it is possible to achieve sublinear regret under a mild geometric condition on S (namely, the existence of a locally Q -regular Borel measure). We provide explicit bounds for a class of regularizers, which guarantee sublinear worst-case regret. We also prove a general lower bound on the regret for any online algorithm and show that DA asymptotically achieves this bound up to a $\log t$ factor. Our results are related to work by Lehrer (2003) and Sridharan and Tewari (2010); Srebro et al. (2011). Lehrer (2003) gives necessary geometric conditions for Blackwell approachability in infinite-dimensional spaces, but no implementable algorithm guaranteeing Hannan-consistency. Sridharan and Tewari (2010) derive general regret bounds for Mirror Descent (MD) under the assumption that the strategy set is uniformly bounded in the norm of the Banach space. We do not make

such an assumption here. In fact, this assumption does not hold in general for our applications in Section 3. The paper is organized as follows: In Section 2 we introduce and provide a general analysis of Dual Averaging in reflexive Banach spaces. In Section 3 we apply these results to obtain explicit regret bounds on compact metric spaces with uniformly continuous reward functions. We use these results in Section 4 in the context of learning Nash equilibria in continuous two-player zero sum games, and provide a numerical example in Section 4. All proofs are given in the supplementary material.

2

Regret Minimization on Reflexive Banach Spaces

Consider a sequential decision problem in which we are to choose a sequence (x_1, x_2, \dots) of actions from some feasible subset X of a reflexive Banach space X , and seek to maximize a sequence $(u_1(x_1), u_2(x_2), \dots)$ of rewards, where the $u_i : X \rightarrow \mathbb{R}$ are elements of a given subset $U \subset X^*$, with X^* the dual space of X . We assume that x_t , the action chosen at time t , may only depend on the sequence of previously observed reward vectors (u_1, \dots, u_{t-1}) . We call any such algorithm an online algorithm. We consider the adversarial setting, i.e., we do not make any distributional assumptions on the rewards. In particular, they could be picked maliciously by some adversary. The notion of regret is a standard measure of performance for such a sequential decision problem. For a sequence (u_1, \dots, u_t) of reward vectors, and a sequence of decisions (x_1, \dots, x_t) produced by an algorithm, the regret of the algorithm w.r.t. a (fixed) $x \in X$ is the gap between the realized t -th decision x_t and the reward under x , i.e., $R_t(x) := \sum_{i=1}^t u_i(x_t) - \sum_{i=1}^t u_i(x)$. The regret is defined as $R_t := \sup_{x \in X} R_t(x)$. An algorithm is said to have sublinear regret if for any sequence $(u_t)_{t=1}^\infty$ in the set of admissible reward functions U , the regret grows sublinearly, i.e. $\limsup_{t \rightarrow \infty} R_t/t = 0$.

Example 1. Consider a finite action set $S = \{1, \dots, n\}$, let $X = X^* = \mathbb{R}^n$, and let $X^* = \mathbb{R}^n$, the probability simplex in \mathbb{R}^n . A reward function can be identified with a vector $u \in \mathbb{R}^n$, such that the i -th element u_i is the reward of action i . A choice $x \in X$ corresponds to a randomization over the n actions in S . This is the classic setting of many regret-minimizing algorithms in the literature. **Example 2.** Suppose S is a compact metric space with μ a finite measure on S . Consider $X = X^* = L^2(S, \mu)$ and let $X = \{x \in X : x \geq 0 \text{ a.e., } \|x\|_1 = 1\}$. A reward function is an L^2 integrable function on S , and each choice $x \in X$ corresponds to a probability distribution (absolutely continuous w.r.t. μ) over S . We will explore a more general variant of this problem in Section 3. In this Section, we prove a general bound on the worst-case regret for DA. DA was introduced by Nesterov (2009) for (finite dimensional) convex optimization, and has also been applied to online learning, e.g. by Xiao (2010). In the finite dimensional case, the

t -th method solves, at each iteration, the optimization problem $x_{t+1} = \arg \max_{x \in X} \sum_{i=1}^t u_i(x) - \eta \sum_{i=1}^t \langle u_i, x \rangle$, where h is a strongly convex 2

regularizer defined on $X \subset \mathbb{R}^n$ and $(\eta_t)_{t=0}^\infty$ is a sequence of learning rates. The regret analysis of the method relies on the duality between strong convexity and smoothness (Nesterov, 2009, Lemma 1). In order to generalize DA to our

Banach space setting, we develop an analogous duality result in Theorem 1. In particular, we show that the correct notion of strong convexity is (uniform) essential strong convexity. Equipped with this duality result, we analyze the regret of the Dual Averaging method and derive a general bound in Theorem 2.

2.1 Preliminaries

Let $(X, \|\cdot\|)$ be a reflexive Banach space, and denote by $\langle \cdot, \cdot \rangle : X \times X \rightarrow \mathbb{R}$ the canonical pairing between X and its dual space X^* , so that $\langle x, y \rangle := \langle y, x \rangle$ for all $x \in X, y \in X^*$. By the effective domain of an extended real-valued function $f : X \rightarrow [-\infty, +\infty]$ we mean the set $\text{dom } f = \{x \in X : f(x) < +\infty\}$. A function f is proper if $f \not\equiv +\infty$ and $\text{dom } f$ is non-empty. The conjugate or Legendre-Fenchel transform of f is the function $f^* : X^* \rightarrow [-\infty, +\infty]$ given by $f^*(y) = \sup_{x \in X} \langle y, x \rangle - f(x)$.

for all $y \in X^*$. If f is proper, lower semicontinuous and convex, its subdifferential ∂f is the set-valued mapping $\partial f(x) = \{y \in X^* : f(y) \leq f(x) + \langle y, x \rangle\}$ for all $x \in X$. We define $\text{dom } \partial f := \{x \in X : \partial f(x) \neq \emptyset\}$. Let \mathcal{C} denote the set of all convex, lower semicontinuous functions $f : [0, \infty) \rightarrow [0, \infty]$ such that $f(0) = 0$, and let

$$\mathcal{U} := \{f \in \mathcal{C} : f(r) \leq 0, f(r) \leq 0, L := \{f \in \mathcal{C} : f(r)/r \leq 0, \text{ as } r \rightarrow 0\}$$

(2) We now introduce some definitions. Additional results are reviewed in the supplementary material. Definition 1 (Strömberg, 2011). A proper convex lower semicontinuous function $f : X \rightarrow [-\infty, \infty]$ is essentially strongly convex if (i) f is strictly convex on every convex subset of $\text{dom } f$ (ii) $(\partial f)^{-1}$ is locally bounded on its domain (iii) for every $x_0 \in \text{dom } \partial f$ there exists $\eta_0 \in X^*$ and $\gamma \in \mathcal{U}$ such that $f(x) \geq f(x_0) + \langle \eta_0, x - x_0 \rangle - \gamma(\|x - x_0\|)$, $\forall x \in X$.

(3)

If (3) holds with γ independent of x_0 , f is uniformly essentially strongly convex with modulus γ . Definition 2 (Strömberg, 2011). A proper convex lower semicontinuous function $f : X \rightarrow [-\infty, \infty]$ is essentially Fréchet differentiable if $\text{int dom } f \neq \emptyset$, f is Fréchet differentiable on $\text{int dom } f$ with Fréchet derivative Df , and $\|Df(x_j) - Df(x)\| \leq \gamma \|x_j - x\|$ for any sequence $(x_j)_j$ in $\text{int dom } f$ converging to some boundary point of $\text{dom } f$. Definition 3. A proper Fréchet differentiable function $f : X \rightarrow [-\infty, \infty]$ is essentially strongly smooth if $\exists x_0 \in \text{dom } \partial f, \eta_0 \in X^*, \gamma \in \mathcal{U}$ such that $f(x) \leq f(x_0) + \langle \eta_0, x - x_0 \rangle + \gamma(\|x - x_0\|)$, $\forall x \in X$. (4) If (4) holds with γ independent of x_0 , f is uniformly essentially strongly smooth with modulus γ . With this we are now ready to give our main duality result: Theorem 1. Let $f : X \rightarrow [-\infty, +\infty]$ be proper, lower semicontinuous and uniformly essentially strongly convex with modulus $\gamma \in \mathcal{U}$. Then (i) f^* is proper and essentially Fréchet differentiable with Fréchet derivative $Df^*(y) = \arg \max_{x \in X} \langle y, x \rangle - f(x)$.

(5)

$x \in X$

If, in addition, $\gamma(r) := \gamma(r)/r$ is strictly increasing, then

$$\|Df^*(y_1) - Df^*(y_2)\| \leq \gamma(r_1) \|y_1 - y_2\|, \quad r_1 = \|y_1\|, r_2 = \|y_2\|.$$

In other words, Df is uniformly continuous with modulus of continuity $\gamma(r) = \gamma(r)/r$.

(6) $\frac{1}{r}$
 $\frac{1}{r}$.

(ii) f is uniformly essentially smooth with modulus $\frac{1}{r}$.

Corollary 1. If $\frac{1}{r} \in C^{1+\alpha}$, $\alpha > 0$ then $\|Df\| \leq \frac{1}{r}$ and $\|Df\| \leq \frac{1}{r}$. In particular, with $\frac{1}{r} = K^2 r$, Definition 1 becomes the classic definition of K -strong convexity, and (6) yields the result familiar from the finite-dimensional case that the gradient Df is $1/K$ Lipschitz with respect to the dual norm (Nesterov, 2009, Lemma 1).

2.2

Dual Averaging in Reflexive Banach Spaces

We call a proper convex function $h : X \rightarrow (-\infty, +\infty]$ a regularizer function on a set $X \subseteq X$ if h is essentially strongly convex and $\text{dom } h = X$. We emphasize that we do not assume h to be Fréchet-differentiable. Definition 1 in conjunction with Lemma S.1 (supplemental material) implies that for any regularizer h , the supremum of any function of the form $h(x) + \langle x, u \rangle$ over X , where $u \in X^*$, will be attained at a unique element of X , namely $Dh(u)$, the Fréchet gradient of h at u .

DA with regularizer h and a sequence of learning rates $(\eta_t)_{t \geq 1}$ generates a sequence of decisions p_t using the simple update rule $x_{t+1} = Dh(\eta_t U_t)$, where $U_t = \sum_{i=1}^t u_i$ and $U_0 := 0$. Theorem 2. Let h be a uniformly essentially strongly convex regularizer on X with modulus $\frac{1}{r}$ and let $(\eta_t)_{t \geq 1}$ be a positive non-increasing sequence of learning rates. Then, for any sequence of payoff functions $(u_t)_{t \geq 1}$ in X^* for which there exists $M > 0$ such that $\sup_{t \geq 1} \|u_t\| \leq M$ for all t , the sequence of plays $(x_t)_{t \geq 0}$ given by

$$p_t x_{t+1} = Dh(\eta_t U_t) \quad (7)$$

$\sup_{t \geq 1} \|x_t\| \leq \frac{1}{r} \sum_{i=1}^t \|u_i\|$, $\|x_t\| \leq \frac{1}{r} \sum_{i=1}^t \|u_i\|$, $\|x_t\| \leq \frac{1}{r} \sum_{i=1}^t \|u_i\|$. (8) $\eta_t \sum_{i=1}^t \|u_i\| \leq 1$ where $h = \inf_{x \in X} h(x)$, $\eta(r) := \eta(r)/r$ and $\eta_0 := \eta_1$. It is possible to obtain a regret bound similar to (8) also in a continuous-time setting. In fact, following Kwon and Mertikopoulos (2014), we derive the bound (8) by first proving a bound on a suitably defined notion of continuous-time regret, and then bounding the difference between the continuous-time and discrete-time regrets. This analysis is detailed in the supplementary material. Note that the condition that $\sup_{t \geq 1} \|u_t\| \leq M$ in Theorem 2 is weaker than the one in Sridharan and Tewari (2010), as it does not imply a uniformly bounded strategy set (e.g., if $X = L^2(\mathbb{R})$ and X is the set of distributions on X , then X is unbounded in L^2 , but the condition may still hold).

Theorem 2 provides a regret bound for a particular choice $x \in X$. Recall that $R_t := \sup_{x \in X} R_t(x)$. In Example 1 the set X is compact, so any continuous regularizer h will be bounded, and hence taking the supremum over x in (8) poses no issue. However, this is not the case in our general setting, as the regularizer may R be unbounded on X . For instance, consider Example 2 with the entropy regularizer $h(x) = \int x(s) \log(x(s)) ds$, which is easily seen to be unbounded on X . As a consequence, obtaining a worst-case bound will in general require additional assumptions on the reward functions and the decision set X . This will be investigated in detail in Section 3. Corollary 2. Suppose that $\frac{1}{r} \in C$

if s is an isolated point of some $S \subset \mathbb{R}^n$ and μ is the Lebesgue measure, in which case $B(s, r) \cap S = \{s\}$. However, under an additional regularity assumption on the measure μ we can avoid such degenerate situations. Definition 4 (Heinonen et al., 2015). A Borel measure μ on a metric space (S, d) is Ahlfors Q -regular if there exist $0 < c_0 \leq C_0 < \infty$ such that for any open ball $B(s, r) \subset S$ we have $c_0 r^Q \leq \mu(B(s, r)) \leq C_0 r^Q$. (12) We say that μ is r_0 -locally Q -regular if (12) holds for all $0 < r \leq r_0$. Intuitively, under an r_0 -locally Q -regular measure, the mass in the neighborhood of any point of S is uniformly bounded from above and below. This will allow, at each iteration t , to assign sufficient probability mass around the maximizer(s) of the cumulative reward function. Example 3. The canonical example for a Q -regular measure is the Lebesgue measure μ on \mathbb{R}^n . If d is the metric induced by the Euclidean norm, then $Q = n$ and the bound (12) is tight with $c_0 = C_0 = 1$, a dimensional constant. However, for general sets $S \subset \mathbb{R}^n$, μ need not be locally Q -regular. A sufficient condition for local regularity of μ is that S is v -uniformly fat (Krichene et al., 2015). Assumption 2. The measure μ is r_0 -locally Q -regular on (S, d) . Under Assumption 2, $B(s, r) \cap S \neq \emptyset$ for all $s \in S$ and $r \leq r_0$, hence we may hope for a bound on $\inf_{x \in B(s, r)} h(x)$ uniform in s . To obtain explicit convergence rates, we have to consider a more specific class of regularizers. 3.2

Explicit Rates for f -Divergences on $L_p(S)$

We consider a particular class of regularizers called f -divergences or Csiszar divergences (Csiszar, 1967). Following Audibert et al. (2014), we define η -potentials and the associated f -divergence. Definition 5. Let $\eta \in [0, 1]$ and a η -potential if $\lim_{z \rightarrow \infty} \eta(z) = 0$, $\lim_{z \rightarrow -\infty} \eta(z) = 1$ and $\eta'(0) = 1$. Associated to η is the convex function $R_\eta : [0, \infty) \rightarrow \mathbb{R}$ defined by $R_\eta(x) = \int_0^x \eta(z) dz$ and the f_η -divergence, defined by $h_\eta(x) = \int_S f_\eta(x(s)) d\mu(s) + \eta(X(x))$, where $\eta(X(x))$ is the indicator function of X (i.e. $\eta(X(x)) = 0$ if $x \notin X$ and $\eta(X(x)) = 1$ if $x \in X$). A remarkable fact is that for regularizers based on η -potentials, the DA update (7) can be computed efficiently. More precisely, it can be shown (see Proposition 3 in Krichene (2015)) that the maximizer in this case has a simple expression in terms of the dual problem, and the problem of computing $\pi_{t+1} = \text{Dh}_\eta(\pi_t, u)$ reduces to computing a scalar dual variable λ_t . 5

Proposition 1. Suppose that $\mu(S) = 1$, and that Assumption 2 holds with constants $r_0 > 0$ and $0 < c_0 \leq C_0 < \infty$. Under the Assumptions of Theorem 3, with $h = h_\eta$ the regularizer associated to an η -potential η , we have that, for any positive sequence $(\lambda_t)_{t \geq 1}$ with $\lambda_t \geq r_0$, t

$$R_t \min(C_0 \int_S \eta(X) d\mu, \mu(S)) \leq \sum_{t=1}^T \lambda_t \left(\int_S \eta(X) d\mu + \eta(X) \right) + k \lambda_t \leq k \lambda_t$$
 (13) $t \geq 1$. For particular choices of the sequences $(\lambda_t)_{t \geq 1}$ and $(\lambda_t)_{t \geq 1}$, we can derive explicit regret rates. 3.3

Analysis for Entropy Dual Averaging (The Generalized Hedge Algorithm) Rx Taking $\eta(z) = e^{-z}$, we have that $f_\eta(x) = 1 - e^{-x}$ and hence the regularizer is $R_\eta(x) = \int_0^x (1 - e^{-z}) dz = x - e^{-x}$. Then $\text{Dh}_\eta(\pi)(s) = \log \pi(s)$. This corresponds to a generalized Hedge algorithm (Arora et al., 2012; Krichene et al., 2015) or the entropic barrier of Bubeck and Eldan (2014) for Euclidean spaces. The regularizer h_η can be shown to be essentially strongly

convex with modulus $\psi(r) = 12r^2$. Corollary 3. Suppose that $\psi(S) = 1$, that ψ is r_0 -locally Q -regular with constants c_0, C_0 , that $k_t \leq M$ for all t , and that $\psi(r) = C_0 r^2$ for $0 \leq r \leq 1$ (that is, the rewards are ψ -Hölder continuous). Then, under Entropy Dual Averaging, choosing $\eta_t = \eta \log t/t$ with

$\eta = 1/2 + 1/Q + 1/C_0 + M + 2\epsilon$ and $\epsilon \geq 0$, we have that $2c_0 \log(c_0 + r_t) \leq R_t + 2C_0 \log t + Q + 2M \log(c_0 + 1) + C_0$ (14) ≤ 0 $t \leq c_0 2^t$ whenever $\log t/t \leq r_0 + \epsilon$. One can now further optimize over the choice of η to obtain the best constant in the bound. Note also that the case $\epsilon = 1$ corresponds to Lipschitz continuity. 3.4

A General Lower Bound

Theorem 4. Let (S, d) be compact, suppose that Assumption 2 holds, and let $w : R \rightarrow R$ be any function with modulus of continuity $\psi \in Z$ such that $\|w(d(s, s_0))\|_q \leq M$ for some $s_0 \in S$ for which there exists $s \in S$ with $d(s, s_0) = DS$. Then for any online algorithm, there exist a sequence $(u_t)_{t=1}^\infty$ of reward vectors $u_t \in X$ with $\|u_t\|_q \leq M$ and modulus of continuity $\psi \in Z$ such that $w(DS) \leq R_t \leq t$, (15) 2.2 Maximizing the constant in (15) is of interest in order to benchmark the bound against the upper bounds obtained in the previous sections. This problem is however quite challenging, and we will defer this analysis to future work. For Hölder-continuous functions, we have the following result: Proposition 2. In the setting of Theorem 4, suppose that $\psi(S) = 1$ and that $\psi(r) = C_0 r^2$ for some $0 \leq r \leq 1$. Then

$1/\eta \min C_0 DS^2, M \leq R_t \leq t$. (16) 2.2 Observe that, up to a $\log t$ factor, the asymptotic rate of this general lower bound for any online algorithm matches that of the upper bound (14) of Entropy Dual Averaging.

4

Learning in Continuous Two-Player Zero-Sum Games

Consider a two-player zero sum game $G = (S_1, S_2, u)$, in which the strategy spaces S_1 and S_2 of player 1 and 2, respectively, are Hausdorff spaces, and $u : S_1 \times S_2 \rightarrow R$ is the payoff function of player 1 (as G is zero-sum, the payoff function of player 2 is $-u$). For each i , denote by $P_i := P(S_i)$ the set of Borel probability measures on S_i . Denote $S := S_1 \times S_2$ and $P := P_1 \times P_2$. For a (joint) mixed strategy $x \in P$, we define the natural extension $u : P \rightarrow R$ by $u(x) := E_x[u] = \int u(s_1, s_2) dx(s_1, s_2)$, which is the expected payoff of player 1 under x . S 6

A continuous zero-sum game G is said to have value V if $\psi(x_1, x_2) = V$. $\psi(x_1, x_2) = \inf \sup u \sup \inf u$ $x_1 \in P_1$ $x_2 \in P_2$

(17)

$x_2 \in P_2$ $x_1 \in P_1$

The elements $x_1 \in P_1$ $x_2 \in P_2$ at which (17) holds are the (mixed) Nash Equilibria of G . We denote the set of Nash equilibria of G by $N(G)$. In the case of finite games, it is well known that every two-player zero-sum game has a value. This is not true in general for continuous games, and additional conditions on strategy sets and payoffs are required, see e.g. (Glicksberg, 1950). 4.1

Repeated Play

We consider repeated play of the continuous two-player zero-sum game. Given a game G and a sequence of plays $(s_1^t)_{t=1}^\infty$ and $(s_2^t)_{t=1}^\infty$, we say that

player i has sublinear (realized) regret if

$$\limsup_{t \rightarrow \infty} \sup_{s_i \in S_i} \sum_{\tau=1}^t u_i(s_i, s_{-i}^\tau) - \max_{s_i \in S_i} \sum_{\tau=1}^t u_i(s_i, s_{-i}^\tau) = o(t) \quad (18)$$

where we use s_{-i} to denote the other player.

A strategy σ_i for player i is, loosely speaking, a (possibly random) mapping from past observations to its actions. Of primary interest to us are Hannan-consistent strategies: Definition 6 (Hannan, 1957). A strategy σ_i of player i is Hannan consistent if, for any sequence $(s_{-i}^\tau)_{\tau=1}^\infty$, the sequence of plays $(s^\tau)_{\tau=1}^\infty$ generated by σ_i has sublinear regret almost surely. Note that the almost sure statement in Definition 6 is with respect to the randomness in the strategy σ_i . The following result is a generalization of its counterpart for discrete games (e.g. Corollary 7.1 in (Cesa-Bianchi and Lugosi, 2006)): Proposition 3. Suppose G has value V and consider a sequence of plays $(s_1^\tau)_{\tau=1}^\infty, (s_2^\tau)_{\tau=1}^\infty$ and P_t assume that both players have sublinear realized regret. Then $\lim_{t \rightarrow \infty} \frac{1}{t} \sum_{\tau=1}^t u(s_1^\tau, s_2^\tau) = V$. As in the discrete case (Cesa-Bianchi and Lugosi, 2006), we can also say something about convergence of the empirical distributions of play to the set of Nash Equilibria. Since these distributions have finite support for every t , we can at best hope for convergence in the weak sense as follows: Theorem 5. Suppose that in a repeated two-player zero sum game $P_t G$ that has a value both players follow a Hannan-consistent strategy, and denote by $x_i^\tau = \frac{1}{t} \sum_{s=1}^\tau x_i(s)$ the marginal empirical distribution of play of player i at iteration t . Let $x_t := (\frac{1}{t} \sum_{s=1}^t x_1(s), \frac{1}{t} \sum_{s=1}^t x_2(s))$. Then $x_t \rightarrow^* N(G)$ almost surely, that is, with probability 1 the sequence $(x_t)_{t=1}^\infty$ weakly converges to the set of Nash equilibria of G . Corollary 4. If G has a unique Nash equilibrium x^* , then with probability 1, $x_t \rightarrow x^*$. 4.2

Hannan-Consistent Strategies

By Theorem 5, if each player follows a Hannan-consistent strategy, then the empirical distributions of play weakly converge to the set of Nash equilibria of the game. But do such strategies exist? Regret minimizing strategies are intuitive candidates, and the intimate connection between regret minimization and learning in games is well studied in many cases, e.g. for finite games (Cesa-Bianchi and Lugosi, 2006) or potential games (Monderer and Shapley, 1996). Using our results from Section 3, we will show that, under the appropriate assumption on the information revealed to the player, no-regret learning based on Dual Averaging leads to Hannan consistency in our setting. Specifically, suppose that after each iteration t , each player i observes a partial payoff function $u_i^\tau : S_i \rightarrow \mathbb{R}$ describing their payoff as a function of only their own action, s_i , holding the action played by the other player fixed. That is, $u_i^\tau(s_i) := u_i(s_i, s_{-i}^\tau)$ and $u_i^{\tau+1}(s_i) := u_i(s_i, s_{-i}^{\tau+1})$. Remark 2. Note that we do not assume that the players have knowledge of the joint utility function u . However, we do assume that the player has full information feedback, in the sense that they observe partial reward functions $u_i^\tau(s_i, s_{-i}^\tau)$ on their entire action set, as opposed to only observing the reward $u(s_1^\tau, s_2^\tau)$ of the action played (the latter corresponds to the bandit setting). $\mathcal{H}_i = \{u_i^\tau\}_{\tau=1}^\infty$ the sequence of partial payoff functions observed by player i . We use U_i to denote the set of all possible such histories, and define $U_{0i} := \{u_i^\tau\}_{\tau=1}^\infty$. A strategy σ_i of player i is a collection $(\sigma_i^\tau)_{\tau=1}^\infty$ of (possibly random)

mappings $\sigma_t : U_t \rightarrow S_i$, such that at iteration t , player i plays $s_{it} = \sigma_t(i)$ (U_t). We make the following assumption on the payoff function: 7

Assumption 3. The payoff function u is uniformly continuous in s_i with modulus of continuity independent of s_{-i} for $i = 1, 2$. That is, for each i there exists $\eta : \mathbb{R}^+ \rightarrow \mathbb{R}^+$ such that $|u(s, s_i) - u(s, s_i')| \leq \eta(|s_i - s_i'|)$ for all $s_i, s_i' \in S_i$. It is easy to see that Assumption 3 implies that the game has a value (see supplementary material). It also makes our setting compatible with that of Section 3. Suppose now that each player randomizes their play according to the sequence of probability distributions on S_i generated by DA with regularizer h_i . That is, suppose that each σ_{ti} is a random variable with the following distribution: $P_{ti}(\sigma_{ti}) = \frac{1}{Z_{ti}} \exp(-\eta(\sigma_{ti}))$. (19) Theorem 6. Suppose that player i uses strategy σ_i according to (19), and that the DA algorithm ensures sublinear regret (i.e. $\limsup_{t \rightarrow \infty} R_t/t = 0$). Then σ_i is Hannan-consistent. Corollary 5. If both players use strategies according to (19) with the respective Dual Averaging ensuring that $\limsup_{t \rightarrow \infty} R_t/t = 0$, then with probability 1 the sequence $(\sigma_{xt})_{t=1}^\infty$ of empirical distributions of play weakly converges to the set of Nash equilibria of G . Example Consider a zero-sum game G_1 between two players on the unit interval with payoff function $u(s_1, s_2) = s_1 s_2 - a_1 s_1 - a_2 s_2$, where $a_1 = e/2$ and $a_2 = e/1$. It is easy to verify that the pair

(σ_1^*, σ_2^*) is a mixed-strategy Nash equilibrium of G_1 . For sequences $(s_t^*)_{t=1}^\infty$, $(s_t^*)_{t=1}^\infty$ and $(s_t^*)_{t=1}^\infty$, the cumulative payoff functions for fixed action $s \in [0, 1]$ are given, respectively, by

$U_1(s) = \sum_{t=1}^T s_2^* - a_1 \sum_{t=1}^T s_1^* - a_2 \sum_{t=1}^T s_2^* - U_2(s) = a_2 \sum_{t=1}^T s_1^* - a_1 \sum_{t=1}^T s_2^* - U_1(s)$ If each player i uses the Generalized Hedge Algorithm with learning rates $(\eta_t)_{t=1}^\infty$, their strategy in period t is to sample from the distribution $\sigma_t(s) \propto \exp(\eta_t U_t(s))$, where $\eta_{t1} = \eta_t (\eta_t - a_1 s_2^* - a_1 s_1^*)$ and $\eta_{t2} = \eta_t (a_2 s_1^* - \eta_t - a_1 s_1^*)$. Interestingly, in this case the sum of the opponent's past plays is a sufficient statistic, in the sense that it completely determines the mixed strategy at time t . 2.5 2.0 1.5 1.0 0.5 0.0 2.5 2.0 1.5 1.0 0.5 0.0 0.0

player 1, t=5000
player 1, t=50000
x1 (s)
player 1, t=500000
x1 (s)
player 2, t=5000
x1 (s)
player 2, t=50000
player 2, t=500000
x2 (s)
0.2
0.4
0.6
0.8
x2 (s)
1.0
0.0

0.2
0.4
0.6
0.8
x2 (s)
1.0
0.0
0.2
0.4
0.6
0.8
1.0

Figure 1: Normalized histograms of the empirical distributions of play in G (100 bins) Figure 1 shows normalized histograms of the empirical distributions of play at different iterations t . As t grows the histograms approach the equilibrium densities x_1 and x_2 , respectively. However, this does not mean that the individual strategies x_{it} converge. Indeed, Figure 2 shows the θ_{it} oscillating around the equilibrium parameters θ_1 and θ_2 , respectively, even for very large t . We do, however, observe that the time-averaged parameters $\bar{\theta}_t$ converge to the equilibrium values θ_1 and θ_2 .

θ_{1t} vs t
 θ_{2t} vs t
1 0 θ_1 100
101
102
103
104

Figure 2: Evolution of parameters θ_{it} and $\bar{\theta}_t :=$

105
Plot of θ_{it}
 $\theta_{it} = 1$
106
 θ_{it} in G1

In the supplementary material we provide additional numerical examples, including one that illustrates how our algorithms can be utilized as a tool to compute approximate Nash equilibria in continuous zero-sum games on non-convex domains. 8

2 References

Sanjeev Arora, Elad Hazan, and Satyen Kale. The multiplicative weights update method: a metaalgorithm and applications. *Theory of Computing*, 8(1):121?164, 2012. Jean-Yves Audibert, S?bastien Bubeck, and G?bor Lugosi. Regret in online combinatorial optimization. *Mathematics of Operations Research*, 39(1):31?45, 2014. S. Bubeck and R. Eldan. The entropic barrier: a simple and optimal

universal self-concordant barrier. ArXiv e-prints, December 2014. S?bastien Bubeck and Nicol? Cesa-Bianchi. Regret analysis of stochastic and nonstochastic multiarmed bandit problems. *Foundations and Trends in Machine Learning*, 5(1):1?122, 2012. Nicol? Cesa-Bianchi and Gabor Lugosi. *Prediction, Learning, and Games*. Cambridge UP, 2006. Thomas M. Cover. Universal portfolios. *Mathematical Finance*, 1(1):1?29, 1991. Imre Csisz?r. Information-type measures of difference of probability distributions and indirect observations. *Studia Scientiarum Mathematicarum Hungarica*, 2:299?318, 1967. Irving L. Glicksberg. Minimax theorem for upper and lower semicontinuous payoffs. Research Memorandum RM-478, The RAND Corporation, Oct 1950. James Hannan. Approximation to Bayes risk in repeated play. In *Contributions to the Theory of Games*, vol III of *Annals of Mathematics Studies* 39. Princeton University Press, 1957. Sergiu Hart and Andreu Mas-Colell. A general class of adaptive strategies. *Journal of Economic Theory*, 98(1):26 ? 54, 2001. Elad Hazan, Amit Agarwal, and Satyen Kale. Logarithmic regret algorithms for online convex optimization. *Machine Learning*, 69(2-3):169?192, 2007. Juha Heinonen., Pekka Koskela, Nageswari Shanmugalingam, and Jeremy T. Tyson. *Sobolev Spaces on Metric Measure Spaces: An Approach Based on Upper Gradients*. New Mathematical Monographs. Cambridge University Press, 2015. Walid Krichene. Dual averaging on compactly-supported distributions and application to no-regret learning on a continuum. CoRR, abs/1504.07720, 2015. Walid Krichene, Maximilian Balandat, Claire Tomlin, and Alexandre Bayen. The Hedge Algorithm on a Continuum. In 32nd International Conference on Machine Learning, pages 824?832, 2015. Joon Kwon and Panayotis Mertikopoulos. A continuous-time approach to online optimization. ArXiv e-prints, January 2014. Ehud Lehrer. Approachability in infinite dimensional spaces. *International Journal of Game Theory*, 31(2):253?268, 2003. Dov Monderer and Lloyd S. Shapley. Potential games. *Games and Economic Behavior*, 14(1):124 ? 143, 1996. Yurii Nesterov. Primal-dual subgradient methods for convex problems. *Mathematical Programming*, 120(1):221?259, 2009. Shai Shalev-Shwartz. Online learning and online convex optimization. *Foundations and Trends in Machine Learning*, 4(2):107?194, 2012. Nati Srebro, Karthik Sridharan, and Ambuj Tewari. On the universality of online mirror descent. In *Advances in Neural Information Processing Systems 24 (NIPS)*, pages 2645?2653. 2011. Karthik Sridharan and Ambuj Tewari. Convex games in banach spaces. In *COLT 2010 - The 23rd Conference on Learning Theory*,., pages 1?13, Haifa, Israel, June 2010. Thomas Str?mberg. Duality between Fr?chet differentiability and strong convexity. *Positivity*, 15(3): 527?536, 2011. Lin Xiao. Dual averaging methods for regularized stochastic learning and online optimization. *J. Mach. Learn. Res.*, 11:2543?2596, December 2010. 9