

---

# Approximating Nash Equilibria in Normal-Form Games via Unbiased Stochastic Optimization

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 We propose the first, to our knowledge, loss function for approximate Nash equi-  
2 libria of normal-form games that is amenable to unbiased Monte Carlo estimation.  
3 This construction allows us to deploy standard non-convex stochastic optimiza-  
4 tion techniques for approximating Nash equilibria, resulting in novel algorithms  
5 with provable guarantees. We complement our theoretical analysis with exper-  
6 iments demonstrating that stochastic gradient descent can outperform previous  
7 state-of-the-art approaches.

## 8 1 Introduction

9 Nash equilibrium famously encodes stable behavioral outcomes in multi-agent systems and is arguably  
10 the most influential solution concept in game theory. Formally speaking, if  $n$  players independently  
11 choose  $n$ , possibly mixed, strategies ( $x_i$  for  $i \in [n]$ ) and their joint strategy ( $\mathbf{x} = \prod_i x_i$ ) constitutes a  
12 *Nash equilibrium*, then no player has any incentive to unilaterally deviate from their strategy. This  
13 concept has sparked extensive research in various fields, ranging from economics [30] to machine  
14 learning [16], and has even inspired behavioral theory generalizations such as quantal response  
15 equilibria which allow for more realistic models of boundedly rational agents [28].

16 Unfortunately, when considering Nash equilibria beyond the special case of the 2-player, zero-sum  
17 scenario, two significant challenges arise. First, it becomes unclear how a group of  $n$  independent  
18 players would collectively identify a Nash equilibrium when multiple equilibria are possible, giving  
19 rise to the *equilibrium selection* problem [18]. Secondly, even approximating a single Nash equilib-  
20 rium is known to be computationally intractable and specifically PPAD-complete [11]. Combining  
21 both problems together, e.g., testing for the existence of equilibria with welfare greater than some  
22 fixed threshold is NP-hard and it is in fact even hard to approximate (i.e., finding a Nash equilibrium  
23 with welfare greater than  $\omega$  for any  $\omega > 0$ , even when the best equilibrium has welfare  $1 - \omega$ ) [2].

24 From a machine learning (ML) practitioner’s perspective, however, such computational complexity  
25 results hardly give pause for thought as collectively we have become all too familiar with the  
26 unreasonable effectiveness of ML heuristics in circumventing such obstacles. Famously, non-convex  
27 optimization is NP-hard, even if the goal is to compute a local minimizer [31], however, stochastic  
28 gradient descent (and variants thereof) succeed in training models with billions of parameters [7].

29 Unfortunately, computational techniques for Nash equilibrium have so far not achieved anywhere  
30 near the same level of success. In contrast, most modern Nash equilibrium solvers for  $n$ -player,  
31  $m$ -action, general-sum, normal-form games (NFGs) are practically restricted to a handful of players  
32 and/or actions per player except in special cases (e.g., symmetric [38] or mean-field games [34]). This  
33 is partially due to the fact that an NFG is represented by a tensor with an exponential  $nm^n$  entries;  
34 even *reading* this description into memory can be computationally prohibitive. More to the point, any

35 computational technique that presumes *exact* computation of the *expectation* of any function sampled  
 36 according to  $\mathbf{x}$  similarly does not have any hope of scaling beyond small instances.

37 This inefficiency arguably lies at the core of the differential success between ML optimization and  
 38 equilibrium computation. For example, numerous techniques exist that reduce the problem of Nash  
 39 equilibrium computation to finding the minimum of the expectation of a random variable (see related  
 40 work section). Unfortunately, unlike the source of randomness in ML applications where batch  
 41 learning suffices to easily produce unbiased estimators, these techniques do not extend easily to game  
 42 theory which incorporates non-linear functions such as maximum, best-response amongst others.  
 43 This raises our motivating goal:

### Can we solve for Nash equilibria via unbiased stochastic optimization?

44 **Our results.** Following in the successful steps of the interplay between ML and stochastic optimiza-  
 45 tion, we reformulate the approximation of Nash equilibria in an NFG as a stochastic non-convex  
 46 optimization problem admitting unbiased Monte-Carlo estimation. This enables the use of powerful  
 47 solvers and advances in parallel computing to efficiently enumerate Nash equilibria for  $n$ -player,  
 48 general-sum games. Furthermore, this re-casting allows practitioners to incorporate other desirable  
 49 objectives into the problem such as “find an approximate Nash equilibrium with welfare above  $\omega$ ”  
 50 or “find an approximate Nash equilibrium nearest the current observed joint strategy” resolving the  
 51 equilibrium selection problem in effectively ad-hoc and application tailored manner. Concretely, we  
 52 make the following contributions by producing:

- 53 • A loss function  $\mathcal{L}(\mathbf{x})$  1) whose global minima coincide with interior Nash equilibria in normal  
 54 form games, 2) admits unbiased Monte-Carlo estimation, and 3) is Lipschitz and bounded.
- 55 • A loss function  $\mathcal{L}^\tau(\mathbf{x})$  1) whose global minima coincide with logit equilibria (QREs) in normal  
 56 form games, 2) admits unbiased Monte-Carlo estimation, and 3) is Lipschitz and bounded.
- 57 • An efficient randomized algorithm for approximating Nash equilibria in a novel class of games. The  
 58 algorithm emerges by employing a recent  $\mathcal{X}$ -armed bandit approach to  $\mathcal{L}^\tau(\mathbf{x})$  and connecting its  
 59 stochastic optimization guarantees to approximate Nash guarantees. For large games, this enables  
 60 approximating equilibria *faster* than the game can even be read into memory.
- 61 • An empirical comparison of stochastic gradient descent against state-of-the-art baselines for  
 62 approximating NEs in large games. In some games, vanilla SGD actually improves upon previous  
 63 state-of-the-art; in others, SGD is slowed by saddle points, a familiar challenge in deep learning [12].

64 Overall, this perspective showcases a promising new route to approximating equilibria at scale in  
 65 practice. We conclude the paper with discussion for future work.

## 66 2 Preliminaries

67 In an  $n$ -player, normal-form game, each player  $i \in \{1, \dots, n\}$  has a strategy set  $\mathcal{A}_i =$   
 68  $\{a_{i1}, \dots, a_{im_i}\}$  consisting of  $m_i$  pure strategies. These strategies can be naturally indexed, so  
 69 we redefine  $\mathcal{A}_i = \{1, \dots, m_i\}$  as an abuse of notation. Each player  $i$  also has a utility function,  
 70  $u_i : \mathcal{A} = \prod_i \mathcal{A}_i \rightarrow [0, 1]$ , (equiv. “payoff tensor”) that maps joint actions to payoffs in the unit-  
 71 interval. Note that equilibria are invariant to payoff shift and scale [27] so we are effectively assuming  
 72 we know bounds on possible payoffs. We denote the average cardinality of the players’ action sets  
 73 by  $\bar{m} = \frac{1}{n} \sum_k m_k$  and maximum by  $m^* = \max_k m_k$ . Player  $i$  may play a mixed strategy by  
 74 sampling from a distribution over their pure strategies. Let player  $i$ ’s mixed strategy be represented  
 75 by a vector  $x_i \in \Delta^{m_i-1}$  where  $\Delta^{m_i-1}$  is the  $(m_i - 1)$ -dimensional probability simplex embedded  
 76 in  $\mathbb{R}^{m_i}$ . Each function  $u_i$  is then extended to this domain so that  $u_i(\mathbf{x}) = \sum_{\mathbf{a} \in \mathcal{A}} u_i(\mathbf{a}) \prod_j x_{ja_j}$   
 77 where  $\mathbf{x} = (x_1, \dots, x_n)$  and  $a_j \in \mathcal{A}_j$  denotes player  $j$ ’s component of the joint action  $\mathbf{a} \in \mathcal{A}$ . For  
 78 convenience, let  $x_{-i}$  denote all components of  $\mathbf{x}$  belonging to players other than player  $i$ .

79 The joint strategy  $\mathbf{x} \in \prod_i \Delta^{m_i-1}$  is a Nash equilibrium if and only if, for all  $i \in \{1, \dots, n\}$ ,  
 80  $u_i(z_i, x_{-i}) \leq u_i(\mathbf{x})$  for all  $z_i \in \Delta^{m_i-1}$ , i.e., no player has any incentive to unilaterally deviate from  
 81  $\mathbf{x}$ . Nash is typically relaxed with  $\epsilon$ -Nash, our focus:  $u_i(z_i, x_{-i}) \leq u_i(\mathbf{x}) + \epsilon$  for all  $z_i \in \Delta^{m_i-1}$ .

82 As an abuse of notation, let the atomic action  $a_i = e_i$  also denote the  $m_i$ -dimensional “one-hot” vector  
 83 with all zeros aside from a 1 at index  $a_i$ ; its use should be clear from the context. We also introduce

Loss	Function	Obstacle
Exploitability	$\max_k \epsilon_k(\mathbf{x})$	max of r.v.
Nikaido-Isoda (NI)	$\sum_k \epsilon_k(\mathbf{x})$	max of r.v.
Fully-Diff. Exp	$\sum_k \sum_{a_k \in \mathcal{A}_k} [\max(0, u_k(a_k, x_{-i}) - u_k(\mathbf{x}))]^2$	max of r.v.
Gradient-based NI	NI w/ $\text{BR}_k \leftarrow \text{aBR}_k = \Pi_{\Delta}(x_k + \eta \nabla_{x_k} u_k(\mathbf{x}))$	$\Pi_{\Delta}$ of r.v.
Unconstrained	Loss + Simplex Deviation Penalty	sampling from $x_i \in \mathbb{R}^{m_k}$

Table 1: Previous loss functions for NFGs and their obstacles to unbiased estimation.

84  $\nabla_{x_i}^i$  as player  $i$ 's utility gradient. And for convenience, denote by  $H_{il}^i = \mathbb{E}_{x_{-il}}[u_i(a_i, a_l, x_{-il})]$  the  
85 bimatrix game approximation [20] between players  $i$  and  $l$  with all other players marginalized out;  
86  $x_{-il}$  denotes all strategies belonging to players other than  $i$  and  $l$  and  $u_i(a_i, a_l, x_{-il})$  separates out  $l$ 's  
87 strategy  $x_l$  from the rest of the players  $x_{-i}$ . Similarly, denote by  $T_{ilq}^i = \mathbb{E}_{x_{-ilq}}[u_i(a_i, a_l, a_q, x_{-ilq})]$   
88 the 3-player tensor approximation to the game. Note player  $i$ 's utility can now be written succinctly  
89 as  $u_i(x_i, x_{-i}) = x_i^\top \nabla_{x_i}^i = x_i^\top H_{il}^i x_l = x_i^\top T_{ilq}^i x_l x_q$  for any  $l, q$  where we use Einstein notation for  
90 tensor arithmetic. For convenience, define  $\text{diag}(z)$  as the function that places a vector  $z$  on the  
91 diagonal of a square matrix, and  $\text{diag3} : z \in \mathbb{R}^d \rightarrow \mathbb{R}^{d \times d \times d}$  as a 3-tensor of shape  $(d, d, d)$  where  
92  $\text{diag3}(z)_{iii} = z_i$ . Following convention from differential geometry, let  $T_v \mathcal{M}$  be the tangent space  
93 of a manifold  $\mathcal{M}$  at  $v$ . For the interior of the  $d$ -action simplex  $\Delta^{d-1}$ , the tangent space is the same at  
94 every point, so we drop the  $v$  subscript, i.e.,  $T\Delta^{d-1}$ . We denote the projection of a vector  $z \in \mathbb{R}^d$   
95 onto this tangent space as  $\Pi_{T\Delta^{d-1}}(z) = z - \frac{1}{d} \mathbf{1}^\top z$ . We drop  $d$  when the dimensionality is clear  
96 from the context. Finally, let  $\mathcal{U}(S)$  denote a discrete uniform distribution over elements from set  $S$ .

### 97 3 Related Work

98 Representing the problem of computing a Nash equilibrium as an optimization problem is not new. A  
99 variety of loss functions and pseudo-distance functions have been proposed. Most of them measure  
100 some function of how much each player can exploit the joint strategy by unilaterally deviating:

$$\epsilon_k(\mathbf{x}) \stackrel{\text{def}}{=} u_k(\text{BR}_k, x_{-k}) - u_k(\mathbf{x}) \text{ where } \text{BR}_k \in \arg \max_z u_k(z, x_{-k}). \quad (1)$$

101 As argued in the introduction, we believe it is important to be able to subsample payoff tensors of  
102 normal-form games in order to scale to large instances. As Nash equilibria can consist of mixed  
103 strategies, it is advantageous to be able to sample from an equilibrium to estimate its exploitability  $\epsilon$ .  
104 However none of these losses is amenable to unbiased estimation under sampled play. Each of the  
105 functions currently explored in the literature is biased under sampled play either because 1) a random  
106 variable appears as the argument of a complex, nonlinear (non-polynomial) function or because 2) how  
107 to sample play is unclear. Exploitability, Nikaido-Isoda (NI) [32] (also known by NashConv [21] and  
108 ADI [15]), as well as fully-differentiable options ([36], p. 106, Eqn 4.31) introduce bias when a max  
109 over payoffs is estimated using samples from  $\mathbf{x}$ . Gradient-based NI [35] requires projecting the result  
110 of a gradient-ascent step onto the simplex; for the same reason as the max, this is prohibitive because  
111 it is a nonlinear operation which introduces bias. Lastly, unconstrained optimization approaches ([36],  
112 p. 106) that instead penalize deviation from the simplex lose the ability to sample from strategies  
113 when iterates are no longer proper distributions. Table 1 summarizes these complications.

## 114 4 Nash Equilibrium as Stochastic Optimization

115 We will now develop our proposed loss function which is amenable to unbiased estimation. Our key  
116 technical insight is to pay special attention to the geometry of the simplex. To our knowledge, prior  
117 works have failed to recognize the role of the tangent space  $T\Delta$ . Proofs are in the appendix.

### 118 4.1 Stationarity on the Simplex Interior

119 **Lemma 1.** *Assuming player  $i$ 's utility,  $u_i(x_i, x_{-i})$ , is concave in its own strategy  $x_i$ , a strategy in  
120 the interior of the simplex is a best response  $\text{BR}_i$  if and only if it has zero projected-gradient<sup>1</sup> norm:*

<sup>1</sup>Not to be confused with the nonlinear (i.e., introduces bias) projected gradient operator introduced in [19].

$$BR_i \in \left( \text{int}\Delta \cap \arg \max_z u_i(z, x_{-i}) - u_i(x_i, x_{-i}) \right) \iff (BR_i \in \text{int}\Delta) \wedge (\|\Pi_{T\Delta}[\nabla_{BR_i}^i]\| = 0). \quad (2)$$

121 In NFGs, each player’s utility is linear in  $x_i$ , thereby satisfying the concavity condition of Lemma 1.

## 122 4.2 Projected Gradient Norm as Loss

123 An equivalent description of a Nash equilibrium is a joint strategy  $\mathbf{x}$  where every player’s strategy is  
 124 a best response to the equilibrium (i.e.,  $x_i = BR_i$  so that  $\epsilon_i(\mathbf{x}) = 0$ ). Lemma 1 states that any interior  
 125 best response has zero projected-gradient norm, which inspires the following loss function

$$\mathcal{L}(\mathbf{x}) = \sum_k \eta_k \|\Pi_{T\Delta}(\nabla_{x_k}^k)\|^2 \quad (3)$$

126 where  $\eta_k > 0$  represent scalar weights, or equivalently, step sizes to be explained next.

127 **Proposition 1.** *The loss  $\mathcal{L}$  is equivalent to NashConv, but where player  $k$ ’s best response is approxi-*  
 128 *mated by a single step of projected-gradient ascent with step size  $\eta_k$ :  $\mathbf{a}BR_k = x_k + \eta_k \Pi_{T\Delta}(\nabla_{x_k}^k)$ .*

129 This connection was already pointed out in prior work for unconstrained problems [15, 35], but this  
 130 result is the first for strategies constrained to the simplex.

## 131 4.3 Connection to True Exploitability

132 In general, we can bound exploitability in terms of the projected-gradient norm as long as each  
 133 player’s utility is concave (this result extends beyond gradients to subgradients of non-smooth  
 134 functions).

135 **Lemma 2.** *The amount a player can gain by exploiting a joint strategy  $\mathbf{x}$  is upper bounded by a*  
 136 *quantity proportional to the norm of the projected-gradient:*

$$\epsilon_k(\mathbf{x}) \leq \sqrt{2} \|\Pi_{T\Delta}(\nabla_{x_k}^k)\|. \quad (4)$$

137 This bound is not tight on the boundary of the simplex, which can be seen clearly by considering  $x_k$   
 138 to be part of a pure strategy equilibrium. In that case, this analysis assumes  $x_k$  can be improved upon  
 139 by a projected-gradient ascent step (via the equivalence pointed out in Proposition 1). However, that  
 140 is false because the probability of a pure strategy cannot be increased beyond 1. We mention this to  
 141 provide further intuition for why  $\mathcal{L}(\mathbf{x})$  is only valid for interior equilibria.

142 Note that  $\|\Pi_{T\Delta}(\nabla_{x_k}^k)\| \leq \|\nabla_{x_k}^k\|$  because  $\Pi_{T\Delta}$  is a projection. Therefore, this improves the naive  
 143 bounds on exploitability and distance to best responses given using the “raw” gradient  $\nabla_{x_k}^k$ .

144 **Lemma 3.** *The exploitability of a joint strategy  $\mathbf{x}$ , is upper bounded by a function of  $\mathcal{L}(\mathbf{x})$ :*

$$\epsilon \leq \sqrt{\frac{2n}{\min_k \eta_k}} \sqrt{\mathcal{L}(\mathbf{x})} \stackrel{\text{def}}{=} f(\mathcal{L}). \quad (5)$$

## 145 4.4 Unbiased Estimation

146 As discussed in Section 3, a primary obstacle to unbiased estimation of  $\mathcal{L}(\mathbf{x})$  is the presence of  
 147 complex, nonlinear functions of random variables, with the projection of a point onto the simplex  
 148 being one such example (see  $\Pi_{\Delta}$  in Table 1). However,  $\Pi_{T\Delta}$ , *the projection onto the tangent space*  
 149 *of the simplex, is linear!* This is the key that allows us to design an unbiased estimator (Lemma 5).

150 Our proposed loss requires computing the squared norm of the *expected value* of the gradient  
 151 under the players’ mixed strategies, i.e., the  $l$ -th entry of player  $k$ ’s gradient equals  $\nabla_{x_{kl}}^k =$   
 152  $\mathbb{E}_{a_{-k} \sim x_{-k}} u_k(a_{kl}, a_{-k})$ . By analogy, consider a random variable  $Y$ . In general,  $\mathbb{E}[Y]^2 \neq \mathbb{E}[Y^2]$ .  
 153 This means that we cannot just sample projected-gradients and then compute their average norm to  
 154 estimate our loss. However, consider taking two independent samples from two corresponding identi-  
 155 cally distributed, independent random variables  $Y^{(1)}$  and  $Y^{(2)}$ . Then  $\mathbb{E}[Y^{(1)}]^2 = \mathbb{E}[Y^{(1)}]\mathbb{E}[Y^{(2)}] =$

	Exact	Sample Others	Sample All
Estimator of $\nabla_{x_k}^{k(p)}$	$u_k(a_{kl}, x_{-k})$	$u_k(a_{kl}, a_{-k} \sim x_{-k})$	$m_k u_k(a_{kl} \sim \mathcal{U}(\mathcal{A}_k), a_{-k} \sim x_{-k}) e_l$
$\hat{\nabla}_{x_k}^{k(p)}$ Bounds	$[0, 1]$	$[0, 1]$	$[0, m_k]$
$\hat{\nabla}_{x_k}^{k(p)}$ Query Cost	$\prod_{i=1}^n m_i$	$m_k$	1
$\mathcal{L}$ Bounds	$\pm \frac{1}{4} \sum_k \eta_k m_k$	$\pm \frac{1}{4} \sum_k \eta_k m_k$	$\pm \frac{1}{4} \sum_k \eta_k m_k^3$
$\mathcal{L}$ Query Cost	$n \prod_{i=1}^n m_i$	$2n\bar{m}$	$2n$

Table 2: Examples and Properties of Unbiased Estimators of Loss and Player Gradients ( $\hat{\nabla}_{x_k}^{k(p)}$ ).

156  $\mathbb{E}[Y^{(1)}Y^{(2)}]$  by properties of expected value over products of independent random variables. This is  
157 a common technique to construct unbiased estimates of expectations over polynomial functions of  
158 random variables. Proceeding in this way, define  $\nabla_{x_k}^{k(1)}$  as a random variable distributed according to  
159 the distribution induced by all other players' mixed strategies ( $j \neq k$ ). Let  $\nabla_{x_k}^{k(2)}$  be independent and  
160 distributed identically to  $\nabla_{x_k}^{k(1)}$ . Then

$$\mathcal{L}(\mathbf{x}) = \mathbb{E}\left[\sum_k \eta_k \underbrace{\left(\hat{\nabla}_{x_k}^{k(1)} - \frac{\mathbf{1}}{m_k} (\mathbf{1}^\top \hat{\nabla}_{x_k}^{k(1)}) \mathbf{1}\right)^\top}_{\text{projected-gradient 1}} \underbrace{\left(\hat{\nabla}_{x_k}^{k(2)} - \frac{\mathbf{1}}{m_k} (\mathbf{1}^\top \hat{\nabla}_{x_k}^{k(2)}) \mathbf{1}\right)}_{\text{projected-gradient 2}}\right] \quad (6)$$

161 where  $\hat{\nabla}_{x_k}^{k(p)}$  is an unbiased estimator of player  $k$ 's gradient. This unbiased estimator can be con-  
162 structed in several ways. The most expensive, an exact estimator, is constructed by marginalizing  
163 player  $k$ 's payoff tensor over all other players' strategies. However, a cheaper estimate can be obtained  
164 at the expense of higher variance by approximating this marginalization with a Monte Carlo estimate  
165 of the expectation. Specifically, if we sample a single action for each of the remaining players, we  
166 can construct an unbiased estimate of player  $k$ 's gradient by considering the payoff of each of its  
167 actions against the sampled background strategy. Lastly, we can consider constructing a Monte Carlo  
168 estimate of player  $k$ 's gradient by sampling only a single action from player  $k$  to represent their entire  
169 gradient. Each of these approaches is outlined in Table 2 along with the query complexity [3] of  
170 computing the estimator and bounds on the values it can take (derived via Lemma 19).

171 We can extend Lemma 3 to one that holds under  $T$  samples with probability  $1 - \delta$  by applying, for  
172 example, a Hoeffding bound:  $\epsilon \leq f(\hat{\mathcal{L}}(\mathbf{x}) + \mathcal{O}(\sqrt{\frac{1}{T} \ln(1/\delta)}))$ .

## 173 4.5 Interior Equilibria

174 We discussed earlier that  $\mathcal{L}(\mathbf{x})$  captures interior equilibria. But some games may only have *pure*  
175 equilibria. We show how to circumvent this shortcoming by considering quantal response equilibria  
176 (QREs), specifically, logit equilibria. By adding an entropy bonus to each player's utility, we can

- 177 • guarantee **all** equilibria are interior,
- 178 • still obtain unbiased estimates of our loss,
- 179 • maintain an upper bound on the exploitability  $\epsilon$  of any approximate equilibrium in the  
180 original game (i.e., the game without an entropy bonus).

181 Define  $u_k^\tau(\mathbf{x}) = u_k(\mathbf{x}) + \tau S(x_k)$  where the Shannon entropy  $S(x_k) = -\sum_l x_{kl} \ln(x_{kl})$  is a 1-  
182 strongly concave function with respect to the 1-norm [6]. Also define  $\mathcal{L}^\tau(\mathbf{x})$  as before except where  
183  $\nabla_{x_k}^k$  is replaced with  $\nabla_{x_k}^{k\tau} = \nabla_{x_k} u_k^\tau(\mathbf{x})$ , i.e., the gradient of player  $k$ 's utility *with* the entropy bonus.

184 It is well known that Nash equilibria of entropy-regularized games satisfy the conditions for logit  
185 equilibria [23], which are solutions to the fixed point equation  $x_k = \text{softmax}(\frac{\nabla_{x_k}^k}{\tau})$ . The appearance  
186 of the `softmax` makes clear that all probabilities have positive mass at positive temperature.

187 Recall that in order to construct an unbiased estimate of our loss, we simply needed to construct  
188 unbiased estimates of player gradients. The introduction of the entropy term to player  $k$ 's utility is  
189 special in that it depends entirely on known quantities, i.e., the player's own mixed strategy. We  
190 can directly and deterministically compute  $\tau \frac{dS}{dx_k} = -\tau(\ln(x_k) + 1)$  and add this to our estimator of  
191  $\nabla_{x_k}^{k(p)}$ :  $\hat{\nabla}_{x_k}^{k\tau(p)} = \hat{\nabla}_{x_k}^{k(p)} + \tau \frac{dS}{dx_k}$ . Consider our refined loss function with changes in **blue**:

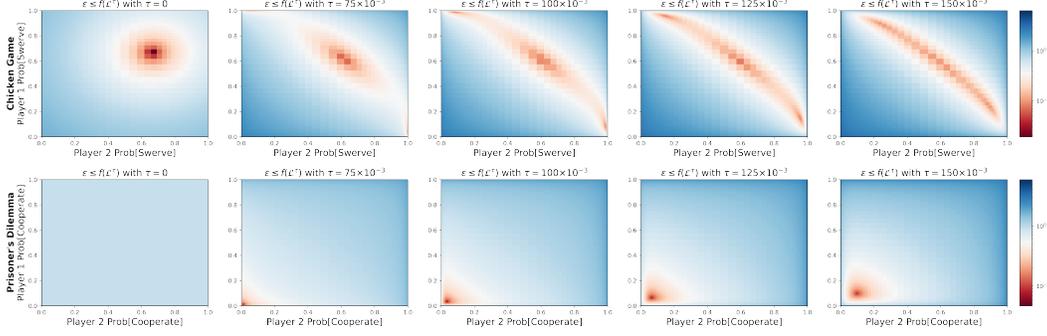


Figure 1: Upper Bound ( $\epsilon \leq f(\mathcal{L}^\tau)$ ) Heatmap Visualization. The first row examines the loss landscape for the classic anti-coordination game of Chicken (Nash equilibria:  $(0, 1)$ ,  $(1, 0)$ ,  $(2/3, 1/3)$ ) while the second row examines the Prisoner’s dilemma (Unique Nash equilibrium:  $(0, 0)$ ). Temperature increases for each plot moving to the right. For high temperatures, interior (fully-mixed) strategies are incentivized while for lower temperatures, nearly pure strategies can achieve minimum exploitability. For zero temperature, pure strategy equilibria (e.g., defect-defect) are not captured by the loss as illustrated by the bottom-left Prisoner’s Dilemma plot with a constant loss surface.

$$\mathcal{L}^\tau(\mathbf{x}) = \sum_k \eta_k \|\Pi_{T\Delta}(\nabla_{x_k}^{k\tau})\|^2. \quad (7)$$

192 As mentioned above, the utilities with entropy bonuses are still concave, therefore, a similar bound  
 193 to Lemma 2 applies. We use this to prove the QRE counterpart to Lemma 3 where  $\epsilon_{QRE}$  is the  
 194 exploitability of an approximate equilibrium in a game with entropy bonuses.

195 **Lemma 4.** *The entropy regularized exploitability,  $\epsilon_{QRE}$ , of a joint strategy  $\mathbf{x}$ , is upper bounded as:*

$$\epsilon_{QRE} \leq \sqrt{\frac{2n}{\min_k \eta_k}} \sqrt{\mathcal{L}^\tau(\mathbf{x})} \stackrel{\text{def}}{=} f(\mathcal{L}^\tau). \quad (8)$$

196 Lastly, we establish a connection between quantal response equilibria and Nash equilibria that allows  
 197 us to approximate Nash equilibria in the original game via minimizing our modified loss  $\mathcal{L}^\tau(\mathbf{x})$ .

198 **Lemma 14** ( $\mathcal{L}^\tau$  Scores Nash Equilibria). *Let  $\mathcal{L}^\tau(\mathbf{x})$  be our proposed entropy regularized loss  
 199 function with payoffs bounded in  $[0, 1]$  and  $\mathbf{x}$  be an approximate QRE. Then it holds that*

$$\epsilon \leq n\tau(W(1/e) + \frac{\bar{m} - 2}{e}) + 2\sqrt{\frac{n \max_k m_k}{\min_k \eta_k}} \sqrt{\mathcal{L}^\tau(\mathbf{x})} \quad (9)$$

200 where  $W$  is the Lambert function:  $W(1/e) = W(\exp(-1)) \approx 0.278$ .

201 This upper bound is plotted as a heatmap for familiar games in Figure 1. Notice how pure equilibria  
 202 are not visible as minima for zero temperature, but appear for slightly warmer temperatures.

## 203 5 Analysis

204 In the preceding section we established a loss function that upper bounds the exploitability of an  
 205 approximate equilibrium. In addition, the zeros of this loss function have a one-to-one correspondence  
 206 with quantal response equilibria (which approximate Nash equilibria at low temperature).

207 Here, we derive properties that suggest it is “easy” to optimize. While this function is generally  
 208 non-convex and may suffer from a proliferation of saddle points and local maxima (Figure 2), it is  
 209 Lipschitz continuous (over a subset of the interior) and bounded. These are two commonly made  
 210 assumptions in the literature on non-convex optimization, which we leverage in Section 6. In addition,  
 211 we can derive its gradient, its Hessian, and characterize its behavior around global minima.

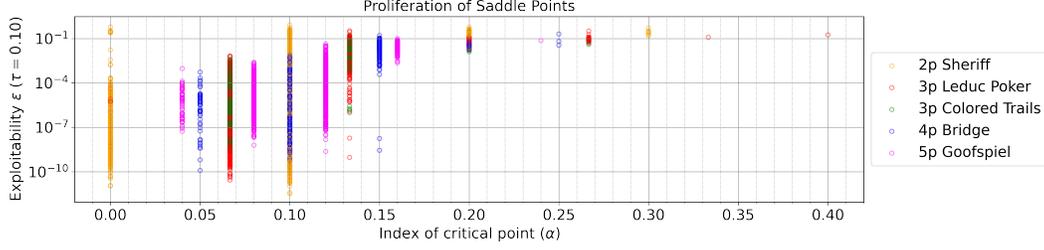


Figure 2: We reapply the analysis of [12], originally designed to understand the success of SGD in deep learning, to “slices” of several popular extensive form games. To construct a slice (or *meta-game*), we randomly sample 6 deterministic policies and then consider the corresponding  $n$ -player, 6-action normal-form game at  $\tau = 0.1$  (with payoffs normalized to  $[0, 1]$ ). The index of a critical point  $\mathbf{x}_c$  ( $\nabla_{\mathbf{x}} \mathcal{L}^\tau(\mathbf{x}_c) = \mathbf{0}$ ) indicates the fraction of negative eigenvalues in the Hessian of  $\mathcal{L}^\tau$  at  $\mathbf{x}_c$ ;  $\alpha = 0$  indicates a local minimum, 1 a maximum, else a saddle point. We see a positive correlation between exploitability and  $\alpha$  indicating a lower prevalence of local minima at high exploitability.

212 **Lemma 15.** *The gradient of  $\mathcal{L}^\tau(\mathbf{x})$  with respect to player  $l$ 's strategy  $x_l$  is*

$$\nabla_{x_l} \mathcal{L}^\tau(\mathbf{x}) = 2 \sum_k \eta_k B_{kl}^\top \Pi_{T\Delta}(\nabla_{x_k}^{k\tau}) \quad (10)$$

213 where  $B_{ll} = -\tau [I - \frac{1}{m_l} \mathbf{1}\mathbf{1}^\top] \text{diag}(\frac{1}{x_l})$  and  $B_{kl} = [I - \frac{1}{m_k} \mathbf{1}\mathbf{1}^\top] H_{kl}^k$  for  $k \neq l$ .

214 **Lemma 17.** *The Hessian of  $\mathcal{L}^\tau(\mathbf{x})$  can be written*

$$\text{Hess}(\mathcal{L}^\tau) = 2[\tilde{B}^\top \tilde{B} + T \Pi_{T\Delta}(\tilde{\nabla}^\tau)] \quad (11)$$

215 where  $\tilde{B}_{kl} = \sqrt{\eta_k} B_{kl}$ ,  $\Pi_{T\Delta}(\tilde{\nabla}^\tau) = [\eta_1 \Pi_{T\Delta}(\nabla_{x_1}^{1\tau}), \dots, \eta_n \Pi_{T\Delta}(\nabla_{x_n}^{n\tau})]$ , and we augment  $T$  (the  
216 3-player approximation to the game,  $T_{lqk}^l$ ) so that  $T_{lll}^l = \tau \text{diag}(\frac{1}{x_l^2})$ .

217 At an equilibrium, the latter term disappears because  $\Pi_{T\Delta}(\nabla_{x_k}^{k\tau}) = \mathbf{0}$  for all  $k$  (Lemma 1). If  $\mathcal{X}$   
218 was  $\mathbb{R}^{n\bar{m}}$ , then we could simply check if  $\tilde{B}$  is full-rank to determine if  $\text{Hess} \succ 0$ . However,  $\mathcal{X}$  is a  
219 simplex product, and we only care about curvature in directions toward which we can update our  
220 equilibrium. Toward that end, define  $M$  to be the  $n(\bar{m} + 1) \times n\bar{m}$  matrix that stacks  $\tilde{B}$  on top of a  
221 repeated identity matrix that encodes orthogonality to the simplex:

$$M(\mathbf{x}) = \begin{bmatrix} -\tau \sqrt{\eta_1} \Pi_{T\Delta}(\frac{1}{x_1}) & \sqrt{\eta_1} \Pi_{T\Delta}(H_{12}^1) & \dots & \sqrt{\eta_1} \Pi_{T\Delta}(H_{1n}^1) \\ \vdots & \vdots & \vdots & \vdots \\ \sqrt{\eta_n} \Pi_{T\Delta}(H_{n1}^n) & \dots & \sqrt{\eta_n} \Pi_{T\Delta}(H_{n,n-1}^n) & -\tau \sqrt{\eta_n} \Pi_{T\Delta}(\frac{1}{x_n}) \\ \mathbf{1}_1^\top & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \vdots \\ 0 & \dots & 0 & \mathbf{1}_n^\top \end{bmatrix} \quad (12)$$

222 where  $\Pi_{T\Delta}(z \in \mathbb{R}^{a \times b}) = [I_a - \frac{1}{a} \mathbf{1}_a \mathbf{1}_a^\top] z$  subtracts the mean from each column of  $z$  and  $\frac{1}{x_i}$  is  
223 shorthand for  $\text{diag}(\frac{1}{x_i})$ . If  $M(\mathbf{x})z = \mathbf{0}$  for a nonzero vector  $z \in \mathbb{R}^{n\bar{m}}$ , this implies there exists a  $z$   
224 that 1) is orthogonal to the ones vectors of each simplex (i.e., is a valid equilibrium update direction)  
225 and 2) achieves zero curvature in the direction  $z$ , i.e.,  $z^\top (\tilde{B}^\top \tilde{B})z = z^\top (\text{Hess})z = 0$ , and so  $\text{Hess}$   
226 is not positive definite. Conversely, if  $M(\mathbf{x})$  is of rank  $n\bar{m}$  for a quantal response equilibrium  $\mathbf{x}$ , then  
227 the Hessian of  $\mathcal{L}^\tau$  at  $\mathbf{x}$  in the tangent space of the simplex product ( $\mathcal{X} = \prod_i \mathcal{X}_i$ ) is positive definite.  
228 In this case, we call  $\mathbf{x}$  *well-isolated* because it implies it is not connected to any other equilibria.

229 By analyzing the rank of  $M$ , we can confirm that many classical matrix games including Rock-  
230 Paper-Scissors, Chicken, Matching Pennies, and Shapley's game all induce strongly convex  $\mathcal{L}^\tau$ 's at  
231 zero temperature (i.e., they have unique mixed Nash equilibria). In contrast, a game like Prisoner's  
232 Dilemma has a unique pure strategy that will not be captured by our loss at zero temperature.

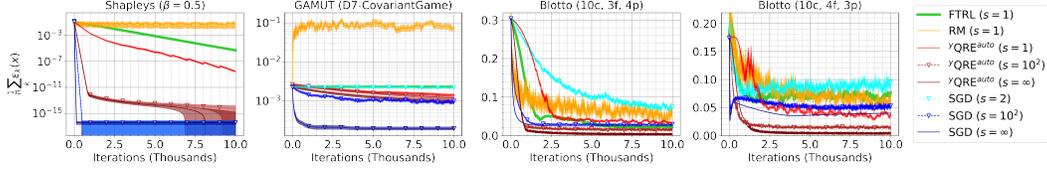


Figure 3: Comparison of SGD on  $\mathcal{L}^{\tau=0}$  against baselines on four games evaluated in [15]. From left to right: 2-player, 3-action, nonsymmetric; 6-player, 5-action, nonsymmetric; 4-player, 66-action, symmetric; 3-player, 286-action, symmetric. SGD struggles at saddle points in Blotto.

## 233 6 Algorithms

234 We have formally transformed the approximation of Nash equilibria in NFGs into a **stochastic**  
 235 optimization problem. To our knowledge, this is the first such formulation that allows one-shot  
 236 unbiased Monte-Carlo estimation which is critical to introduce the use of powerful algorithms capable  
 237 of solving high dimensional optimization problems. We explore two off-the-shelf approaches.

238 Stochastic gradient descent is the workhorse of high-dimensional stochastic optimization. It comes  
 239 with guaranteed convergence to stationary points [10], however, it may converge to local, rather than  
 240 global minima. It also enjoys implicit gradient regularization [4], seeking “flat” minima and performs  
 241 approximate Bayesian inference [26]. Despite the lack of global convergence guarantee, in the next  
 242 section, we find it performs well empirically in games previously examined by the literature.

243 We explore one other algorithmic approach to non-convex optimization based on minimizing regret,  
 244 which enjoys finite time convergence rates.  $\mathcal{X}$ -armed bandits [8] systematically explore the space of  
 245 solutions by refining a mesh over the joint strategy space, trading off exploration versus exploitation  
 246 of promising regions.<sup>2</sup> Several approaches exist [5, 37] with open source implementations (e.g., [24]).

### 247 6.1 High Probability, Polynomial Convergence Rates

248 We use a recent  $\mathcal{X}$ -armed bandit approach called BLiN [14] to establish a high probability  $\tilde{\mathcal{O}}(T^{-1/4})$   
 249 convergence rate to Nash equilibria in  $n$ -player, general-sum games under mild assumptions. The  
 250 quality of this approximation improves as  $\tau \rightarrow 0$ , at the same time increasing the constant on the  
 251 convergence rate via the Lipschitz constant  $\sqrt{\hat{L}}$  defined below. For clarity, we assume users provide  
 252 a temperature in the form  $\tau = \frac{1}{\ln(1/p)}$  with  $p \in (0, 1)$  which ensures all equilibria have probability  
 253 mass greater than  $\frac{p}{m^*}$  for all actions (Lemma 9). Lower  $p$  corresponds with lower temperature.

254 The following convergence rate depends on bounds on the exploitability in terms of the loss  
 255 (Lemma 14), bounds on the magnitude of estimates of the loss (Lemma 8), Lipschitz bounds on the  
 256 infinity norm of the gradient (Corollary 2), and the number of distinct strategies ( $n\bar{m} = \sum_k m_k$ ).

257 **Theorem 1** (BLiN PAC Rate). *Assume  $\eta_k = \eta = 2/\hat{L}$ ,  $\tau = \frac{1}{\ln(1/p)}$ , and a previously pulled arm is*  
 258 *returned uniformly at random (i.e.,  $t \sim U([T])$ ). Then for any  $w > 0$*

$$\epsilon_t \leq w \left[ \frac{n}{\ln(1/p)} \left( W(1/e) + \frac{\bar{m} - 2}{e} \right) + 4 \left( 1 + (4c^2)^{1/3} \right) \sqrt{nm^* \hat{L}} \left( \frac{\ln T}{T} \right)^{\frac{1}{2(d_z + 2)}} \right] \quad (13)$$

259 *with probability  $(1 - w^{-1})(1 - 2T^{-2})$  where  $W$  is the Lambert function ( $W(1/e) \approx 0.278$ ),*

260  *$m^* = \max_k m_k$ ,  $c \leq \frac{1}{4} \frac{n\bar{m}}{\hat{L}} \left( \frac{\ln(m^*)}{\ln(1/p)} + 2 \right)^2 \leq \frac{1}{4} \left( \frac{\ln(m^*)}{\ln(1/p)} + 2 \right)$  upper bounds the range of stochastic*

261 *estimates of  $\mathcal{L}^\tau$  (see Lemma 8), and  $\hat{L} = \left( \frac{\ln(m^*)}{\ln(1/p)} + 2 \right) \left( \frac{m^{*2}}{p \ln(1/p)} + n\bar{m} \right)$  (see Corollary 2).*

262 This result depends on the *near-optimality* [37] or *zooming-dimension*  $d_z = n\bar{m} \left( \frac{\alpha_{hi} - \alpha_{lo}}{\alpha_{lo} \alpha_{hi}} \right) \in [0, \infty)$   
 263 (Theorem 2) where  $\alpha_{lo}$  and  $\alpha_{hi}$  denote the degree of the polynomials that lower and upper bound the  
 264 function  $\mathcal{L}^\tau \circ s$  locally around an equilibrium. For example, in the case where the Hessian is positive  
 265 definite,  $\alpha_{lo} = \alpha_{hi} = 2$  and  $d_z = 0$ . Here,  $s : [0, 1]^{n(\bar{m}-1)} \rightarrow \prod_i \Delta^{m_i-1}$  is any function that maps  
 266 from the unit hypercube to a product of simplices; we analyze two such maps in the appendix.

<sup>2</sup>Zhou et al. [39] developed a similar approach but only for pure Nash equilibria.

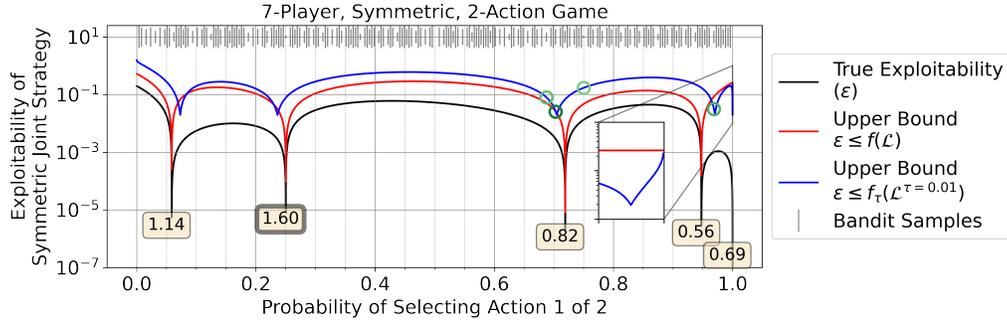


Figure 4: Bandit-based (BLiN) Nash solver applied to an artificial 7-player, symmetric, 2-action game. We search for a symmetric equilibrium, which is represented succinctly as the probability of selecting action 1. The plot shows the true exploitability  $\epsilon$  of all symmetric strategies in black and indicates there exist potentially 5 NEs (the dips in the curve). Upper bounds on our unregularized loss  $\mathcal{L}$  capture 4 of these equilibria, missing only the pure NE on the right. By considering our regularized loss,  $\mathcal{L}^\tau$ , we are able to capture this pure NE (see zoomed inset). The bandit algorithm selects strategies to evaluate, using 10 Monte-Carlo samples for each evaluation (arm pull) of  $\mathcal{L}^\tau$ . These samples are displayed as vertical bars above with the height of the vertical bar representing additional arm pulls. The best arms throughout search are denoted by green circles (darker indicates later in the search). The boxed numbers near equilibria display the welfare of the strategy.

267 Note that Theorem 1 implies that for games whose corresponding  $\mathcal{L}^\tau$  has zooming dimension  $d_z = 0$ ,  
 268 NEs can be approximated with high probability in polynomial time. This general property is difficult  
 269 to translate concisely into game theory parlance. For this reason, we present the following more  
 270 interpretable corollary which applies to a more restricted class of games.

271 **Corollary 1.** Consider the class of NFGs with at least one QRE( $\tau$ ) whose local polymatrix approx-  
 272 imation indicates it is isolated (i.e.,  $M$  from equation (12) is rank- $n\bar{m}$  implies  $\text{Hess} \succ 0$  implies  
 273  $d_z = n\bar{m}(\frac{2-2}{4}) = 0$ ). Then by Theorem 1, BLiN is a fully polynomial-time randomized approximation  
 274 scheme (FPRAS) for QREs and is a PRAS for NEs of games in this class.

275 To convey the impact of stochastic optimization guarantees more concretely, assume we are given  
 276 that an interior well-isolated NE exists. Then for a 20-player, 50-action game, it is  $1000\times$  cheaper to  
 277 compute a  $1/100$ -NE with probability 95% than it is to just list the  $nm^n$  payoffs that define the game.

## 278 6.2 Empirical Evaluation

279 Figure 3 shows SGD is competitive with scalable techniques to approximating NEs. Shapley’s game  
 280 induces a strongly convex  $\mathcal{L}$  (see Section 5) leading to SGD’s strong performance. Blotto shows  
 281 signs of convergence to low, but nonzero  $\epsilon$ , demonstrating the challenges of local minima.

282 We demonstrate BLiN (applied to  $\mathcal{L}^\tau$ ) on a 7-player, symmetric, 2-action game. Figure 4 shows the  
 283 bandit algorithm discovers two equilibria, settling on one near  $\mathbf{x} = [0.7, 0.3] \times 7$  with a wider basin  
 284 of attraction (and higher welfare). In theory, BLiN can enumerate all NEs as  $T \rightarrow \infty$ .

## 285 7 Conclusion

286 In this work, we proposed a stochastic loss for approximate Nash equilibria in normal-form games.  
 287 An unbiased loss estimator of Nash equilibria is the “key” to the stochastic optimization “door”  
 288 which holds a wealth of research innovations uncovered over several decades. Thus, it allows the  
 289 development of new algorithmic techniques for computing equilibria. We consider bandit and vanilla  
 290 SGD methods in this work, but these are only two of the many options now at our disposal (e.g,  
 291 adaptive methods [1], Gaussian processes [9], evolutionary algorithms [17], etc.). Such approaches as  
 292 well as generalizations of these techniques to imperfect-information games are promising directions  
 293 for future work. Similarly to how deep learning research first balked at and then marched on to train  
 294 neural networks via NP-hard non-convex optimization, we hope computational game theory can  
 295 march ahead to make useful equilibrium predictions of large multiplayer systems.

## References

- 296
- 297 [1] K. Antonakopoulos, P. Mertikopoulos, G. Piliouras, and X. Wang. Adagrad avoids saddle points.  
298 In *International Conference on Machine Learning*, pages 731–771. PMLR, 2022.
- 299 [2] P. Austrin, M. Braverman, and E. Chlamtáč. Inapproximability of NP-complete variants of Nash  
300 equilibrium. In *Approximation, Randomization, and Combinatorial Optimization. Algorithms  
301 and Techniques: 14th International Workshop, APPROX 2011, and 15th International Workshop,  
302 RANDOM 2011, Princeton, NJ, USA, August 17-19, 2011. Proceedings*, pages 13–25. Springer,  
303 2011.
- 304 [3] Y. Babichenko. Query complexity of approximate Nash equilibria. *Journal of the ACM (JACM)*,  
305 63(4):36:1–36:24, 2016.
- 306 [4] D. Barrett and B. Dherin. Implicit gradient regularization. In *International Conference on  
307 Learning Representations*, 2020.
- 308 [5] P. L. Bartlett, V. Gabillon, and M. Valko. A simple parameter-free and adaptive approach to  
309 optimization under a minimal local smoothness assumption. In *Algorithmic Learning Theory*,  
310 pages 184–206. PMLR, 2019.
- 311 [6] A. Beck and M. Teboulle. Mirror descent and nonlinear projected subgradient methods for  
312 convex optimization. *Operations Research Letters*, 31(3):167–175, 2003.
- 313 [7] T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam,  
314 G. Sastry, A. Askell, et al. Language models are few-shot learners. *Advances in neural  
315 information processing systems*, 33:1877–1901, 2020.
- 316 [8] S. Bubeck, R. Munos, G. Stoltz, and C. Szepesvári.  $\mathcal{X}$ -armed bandits. *Journal of Machine  
317 Learning Research*, 12(5), 2011.
- 318 [9] D. Calandriello, L. Carratino, A. Lazaric, M. Valko, and L. Rosasco. Scaling gaussian process  
319 optimization by evaluating a few unique candidates multiple times. In *International Conference  
320 on Machine Learning*, pages 2523–2541. PMLR, 2022.
- 321 [10] A. Cutkosky, H. Mehta, and F. Orabona. Optimal stochastic non-smooth non-convex optimiza-  
322 tion through online-to-non-convex conversion. *arXiv preprint arXiv:2302.03775*, 2023.
- 323 [11] C. Daskalakis, P. W. Goldberg, and C. H. Papadimitriou. The complexity of computing a Nash  
324 equilibrium. *Communications of the ACM*, 52(2):89–97, 2009.
- 325 [12] Y. N. Dauphin, R. Pascanu, C. Gulcehre, K. Cho, S. Ganguli, and Y. Bengio. Identifying and  
326 attacking the saddle point problem in high-dimensional non-convex optimization. *Advances in  
327 neural information processing systems*, 27, 2014.
- 328 [13] A. Deligkas, J. Fearnley, A. Hollender, and T. Melissourgos. Pure-circuit: Strong inapproxima-  
329 bility for PPAD. In *2022 IEEE 63rd Annual Symposium on Foundations of Computer Science  
330 (FOCS)*, pages 159–170. IEEE, 2022.
- 331 [14] Y. Feng, T. Wang, et al. Lipschitz bandits with batched feedback. *Advances in Neural  
332 Information Processing Systems*, 35:19836–19848, 2022.
- 333 [15] I. Gemp, R. Savani, M. Lanctot, Y. Bachrach, T. Anthony, R. Everett, A. Tacchetti, T. Eccles,  
334 and J. Kramár. Sample-based approximation of Nash in large many-player games via gradient  
335 descent. In *Proceedings of the 21st International Conference on Autonomous Agents and  
336 Multiagent Systems*, pages 507–515, 2022.
- 337 [16] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and  
338 Y. Bengio. Generative adversarial nets. *Advances in Neural Information Processing Systems*,  
339 27, 2014.
- 340 [17] N. Hansen, S. D. Müller, and P. Koumoutsakos. Reducing the time complexity of the de-  
341 randomization evolution strategy with covariance matrix adaptation (CMA-ES). *Evolutionary  
342 computation*, 11(1):1–18, 2003.

- 343 [18] J. C. Harsanyi, R. Selten, et al. A general theory of equilibrium selection in games. *MIT Press*  
344 *Books*, 1, 1988.
- 345 [19] E. Hazan, K. Singh, and C. Zhang. Efficient regret minimization in non-convex games. In  
346 *International Conference on Machine Learning*, pages 1433–1441. PMLR, 2017.
- 347 [20] E. Janovskaja. Equilibrium points in polymatrix games. *Lithuanian Mathematical Journal*, 8  
348 (2):381–384, 1968.
- 349 [21] M. Lanctot, V. Zambaldi, A. Gruslys, A. Lazaridou, K. Tuyls, J. Pérolat, D. Silver, and  
350 T. Graepel. A unified game-theoretic approach to multiagent reinforcement learning. In  
351 *Advances in Neural Information Processing Systems*, pages 4190–4203, 2017.
- 352 [22] M. Lanctot, E. Lockhart, J.-B. Lespiau, V. Zambaldi, S. Upadhyay, J. Pérolat, S. Sriniva-  
353 sathan, F. Timbers, K. Tuyls, S. Omidshafiei, D. Hennes, D. Morrill, P. Muller, T. Ewalds,  
354 R. Faulkner, J. Kramár, B. D. Vylder, B. Saeta, J. Bradbury, D. Ding, S. Borgeaud, M. Lai,  
355 J. Schrittwieser, T. Anthony, E. Hughes, I. Danihelka, and J. Ryan-Davis. OpenSpiel:  
356 A framework for reinforcement learning in games. *CoRR*, abs/1908.09453, 2019. URL  
357 <http://arxiv.org/abs/1908.09453>.
- 358 [23] S. Leonardos, G. Piliouras, and K. Spendlove. Exploration-exploitation in multi-agent com-  
359 petition: convergence with bounded rationality. *Advances in Neural Information Processing*  
360 *Systems*, 34:26318–26331, 2021.
- 361 [24] W. Li, H. Li, J. Honorio, and Q. Song. Pyxab – a python library for  $\mathcal{X}$ -armed bandit and online  
362 blackbox optimization algorithms, 2023. URL <https://arxiv.org/abs/2303.04030>.
- 363 [25] C. K. Ling, F. Fang, and J. Z. Kolter. What game are we playing? end-to-end learning in normal  
364 and extensive form games. *arXiv preprint arXiv:1805.02777*, 2018.
- 365 [26] S. Mandt, M. D. Hoffman, and D. M. Blei. Stochastic gradient descent as approximate bayesian  
366 inference. *Journal of Machine Learning Research*, 18:1–35, 2017.
- 367 [27] L. Marris, I. Gemp, and G. Piliouras. Equilibrium-invariant embedding, metric space, and  
368 fundamental set of  $2 \times 2$  normal-form games. *arXiv preprint arXiv:2304.09978*, 2023.
- 369 [28] R. D. McKelvey and T. R. Palfrey. Quantal response equilibria for normal form games. *Games*  
370 *and Economic Behavior*, 10(1):6–38, 1995.
- 371 [29] D. Milec, J. Černý, V. Lisý, and B. An. Complexity and algorithms for exploiting quantal  
372 opponents in large two-player games. *Proceedings of the AAAI Conference on Artificial*  
373 *Intelligence*, 35(6):5575–5583, 2021.
- 374 [30] P. R. Milgrom and R. J. Weber. A theory of auctions and competitive bidding. *Econometrica:*  
375 *Journal of the Econometric Society*, pages 1089–1122, 1982.
- 376 [31] K. G. Murty and S. N. Kabadi. Some NP-complete problems in quadratic and nonlinear  
377 programming. Technical report, 1985.
- 378 [32] H. Nikaidô and K. Isoda. Note on non-cooperative convex games. *Pacific Journal of Mathemat-*  
379 *ics*, 5(1):807815, 1955.
- 380 [33] E. Nudelman, J. Wortman, Y. Shoham, and K. Leyton-Brown. Run the GAMUT: A comprehen-  
381 sive approach to evaluating game-theoretic algorithms. In *AAMAS*, volume 4, pages 880–887,  
382 2004.
- 383 [34] J. Pérolat, S. Perrin, R. Elie, M. Laurière, G. Piliouras, M. Geist, K. Tuyls, and O. Pietquin.  
384 Scaling mean field games by online mirror descent. In *Proceedings of the 21st International*  
385 *Conference on Autonomous Agents and Multiagent Systems*, 2022.
- 386 [35] A. Raghunathan, A. Cherian, and D. Jha. Game theoretic optimization via gradient-based  
387 Nikaido-Isoda function. In *International Conference on Machine Learning*, pages 5291–5300.  
388 PMLR, 2019.

- 389 [36] Y. Shoham and K. Leyton-Brown. *Multiagent systems: Algorithmic, game-theoretic, and logical*  
390 *foundations*. Cambridge University Press, 2008.
- 391 [37] M. Valko, A. Carpentier, and R. Munos. Stochastic simultaneous optimistic optimization. In  
392 *International Conference on Machine Learning*, pages 19–27. PMLR, 2013.
- 393 [38] B. Wiedenbeck and E. Brinkman. Data structures for deviation payoffs. In *Proceedings of the*  
394 *22nd International Conference on Autonomous Agents and Multiagent Systems*, 2023.
- 395 [39] Y. Zhou, J. Li, and J. Zhu. Identify the Nash equilibrium in static games with random payoffs.  
396 In *International Conference on Machine Learning*, pages 4160–4169. PMLR, 2017.

397	<b>Appendix: Approximating Nash Equilibria in Normal-Form Games via</b>	
398	<b>Unbiased Stochastic Optimization</b>	
399	<b>A Loss: Connection to Exploitability, Unbiased Estimation, and Upper Bounds</b>	<b>14</b>
400	A.1 KKT Conditions Imply Fixed Point Sufficiency . . . . .	14
401	A.2 Norm of Projected-Gradient and Equivalence to NFG Exploitability with Approximate Best Responses . . . . .	15
402		
403	A.3 Connection to True Exploitability . . . . .	15
404	A.4 Unbiased Estimation . . . . .	17
405	A.5 Bound on Loss . . . . .	18
406	<b>B QREs Approximate NEs at Low Temperature</b>	<b>19</b>
407	<b>C Gradient of Loss</b>	<b>22</b>
408	C.1 Unbiased Estimation . . . . .	23
409	C.2 Bound on Gradient / Lipschitz Property . . . . .	23
410	<b>D Hessian of Loss</b>	<b>24</b>
411	<b>E Regret Bounds</b>	<b>26</b>
412	<b>F Complexity</b>	<b>27</b>
413	F.1 Polymatrix Games . . . . .	27
414	F.2 Normal-Form Games . . . . .	28
415	<b>G Helpful Lemmas and Propositions</b>	<b>28</b>
416	<b>H Maps from Hypercube to Simplex Product</b>	<b>30</b>
417	H.1 Hessian of Bandit Reward Function . . . . .	30
418	<b>I A2: Bounded Diameters and Well-shaped Cells</b>	<b>33</b>
419	I.1 $L_2$ -Norm . . . . .	33
420	I.2 $L_\infty$ -Norm . . . . .	34
421	I.3 Near Optimality Dimension . . . . .	35
422	<b>J D-BLiN</b>	<b>39</b>
423	<b>K Experimental Setup and Details</b>	<b>43</b>
424	K.1 Loss Visualization and Rank Test . . . . .	43
425	K.2 Saddle Point Analysis . . . . .	44
426	K.3 SGD on Classical Games . . . . .	45
427	K.4 BLiN on Artificial Game . . . . .	45

428 **A Loss: Connection to Exploitability, Unbiased Estimation, and Upper**  
429 **Bounds**

430 **A.1 KKT Conditions Imply Fixed Point Sufficiency**

431 Consider the following constrained optimization problem:

$$\max_{\mathbf{x} \in \mathbb{R}^d} f(\mathbf{x}) \quad (14)$$

$$s.t. g_i(\mathbf{x}) \leq 0 \quad \forall i \quad (15)$$

$$h_j(\mathbf{x}) = 0 \quad \forall j \quad (16)$$

432 where  $f$  is concave and  $g_i$  and  $h_j$  represent inequality and equality constraints respectively. If  $g_i$   
433 and  $h_j$  are affine functions, then any maximizer  $\mathbf{x}^*$  of  $f$  must satisfy the following KKT conditions  
434 (necessary and sufficient):

- 435 • Stationarity:  $0 \in \partial f(\mathbf{x}^*) - \sum_j \lambda_j \partial h_j(\mathbf{x}^*) - \sum_i \mu_i \partial g_i(\mathbf{x}^*)$
- 436 • Primal feasibility:  $h_j(\mathbf{x}^*) = 0$  for all  $j$  and  $g_i(\mathbf{x}^*) \leq 0$  for all  $i$
- 437 • Dual feasibility:  $\mu_i \geq 0$  for all  $i$
- 438 • Complementary slackness:  $\mu_i g_i(\mathbf{x}^*) = 0$  for all  $i$ .

439 **Lemma 1.** *Assuming player  $i$ 's utility,  $u_i(x_i, x_{-i})$ , is concave in its own strategy  $x_i$ , any best*  
440 *response in the interior of the simplex has zero projected-gradient norm:*

$$z^* \in \left( \text{int}\Delta \cup \arg \max_z u_i(z, x_{-i}) - u_i(x_i, x_{-i}) \right) \iff (z^* \in \text{int}\Delta) \wedge (\|\Pi_\Delta[\nabla_{z^*}^i]\| = 0). \quad (17)$$

441 *Proof.* Consider the problem of formally computing  $\text{exp}_i(\mathbf{x}) = \max_{z \in \text{int}\Delta} u_i(z, x_{-i}) - u_i(x_i, x_{-i})$ :

$$\max_{z \in \mathbb{R}^d} u_i(z, x_{-i}) - u_i(x_i, x_{-i}) \quad (18)$$

$$s.t. -z_i + x_{\min} \leq 0 \quad \forall i \quad (19)$$

$$1 - \sum_i z_i = 0. \quad (20)$$

442 where  $x_{\min} > 0$  is some constant that captures our given assumption that the solution  $z^*$  lies in  
443 the interior of the simplex. Note that the objective is linear (concave) in  $z$  and the constraints are  
444 affine, therefore the KKT conditions are necessary and sufficient for optimality. Mapping the KKT  
445 conditions onto this problem yields the following:

- 446 • Stationarity:  $0 \in \partial u_i(z^*, x_{-i}) + \lambda \mathbf{1} + \sum_i \mu_i e_i$
- 447 • Primal feasibility:  $\sum_i z_i^* = 1$  and  $z_i^* \geq x_{\min}$  for all  $i$
- 448 • Dual feasibility:  $\mu_i \geq 0$  for all  $i$
- 449 • Complementary slackness:  $\mu_i z_i^* = 0$  for all  $i$ .

450 For any point  $z \in \text{int}\Delta$ , primal feasibility will be satisfied for some  $x_{\min} > 0$ . This implies each  $z_j$   
451 is strictly positive. By complementary slackness and dual feasibility, each  $\mu_i$  must be identically zero.  
452 This implies the stationarity condition can be simplified to  $0 \in \partial u_i(z^*, x_{-i}) + \lambda \mathbf{1}$ . Rearranging  
453 terms we find that for any  $z^*$ , there exists a  $\lambda$  such that

$$\partial u_i(z^*, x_{-i}) \in \lambda \mathbf{1}. \quad (21)$$

454 Equivalently,  $\partial u_i(z^*, x_{-i}) \propto \mathbf{1}$  at  $z^* \in \text{int}\Delta$ . Any vector proportional to the ones vector has zero  
455 projected-gradient norm, completing the claim.  $\square$

456 **A.2 Norm of Projected-Gradient and Equivalence to NFG Exploitability with Approximate**  
457 **Best Responses**

458 **Proposition 1.** *The loss  $\mathcal{L}$  is equivalent to NashConv, but where player  $k$ 's best response is approxi-*  
459 *mated by a single step of projected-gradient ascent with step size  $\eta_k$ :  $\mathbf{aBR}_k = x_k + \eta_k \Pi_\Delta[\nabla_{x_k}^k]$ .*

460 *Proof.* Define an approximate best response as the result of a player adjusting their strategy via a  
461 projected-gradient ascent step, i.e.,  $\mathbf{aBR}_k = x_k + \eta_k \Pi_\Delta[\nabla_{x_k}^k]$  for player  $k$ .

462 In a normal form game, player  $k$ 's utility at this new strategy is  $u_k(\mathbf{aBR}_k, x_{-k}) = (\nabla_{x_k}^k)^\top (x_k +$   
463  $\eta_k \Pi_\Delta[\nabla_{x_k}^k]) = u_k(\mathbf{x}) + \eta_k (\nabla_{x_k}^k)^\top \Pi_\Delta[\nabla_{x_k}^k]$ .

464 Therefore, the amount player  $k$  gains by playing  $\mathbf{aBR}$  is

$$\hat{\epsilon}_k(\mathbf{x}) = u_k(\mathbf{aBR}_k, x_{-k}) - u_k(\mathbf{x}) \quad (22)$$

$$= \eta_k (\nabla_{x_k}^k)^\top \Pi_\Delta[\nabla_{x_k}^k] \quad (23)$$

$$= \eta_k \left( \nabla_{x_k}^k - \frac{1}{m_k} (\mathbf{1}^\top \nabla_{x_k}^k) \mathbf{1} \right)^\top \Pi_\Delta[\nabla_{x_k}^k] \quad (24)$$

$$= \eta_k \|\Pi_\Delta[\nabla_{x_k}^k]\|^2 \quad (25)$$

465 where the third equality follows from the fact that the projected-gradient,  $\Pi_\Delta[\nabla_{x_k}^k]$ , is orthogonal to  
466 the ones vector.  $\square$

467 **A.3 Connection to True Exploitability**

468 **Lemma 2.** *The amount a player can gain by deviating is upper bounded by a quantity proportional*  
469 *to the norm of the projected-gradient:*

$$\epsilon_k(\mathbf{x}) \leq \sqrt{2} \|\Pi_\Delta(\nabla_{x_k}^k)\|. \quad (26)$$

470 *Proof.* Let  $z$  be any point on the simplex. Then

$$u_k(z, x_{-k}) - u_k(\mathbf{x}) \leq (\nabla_{x_k}^k)^\top (z - x_k) \quad (27)$$

$$= (\nabla_{x_k}^k)^\top (z - x_k) - \frac{1}{m_k} (\mathbf{1}^\top \nabla_{x_k}^k) \overbrace{\mathbf{1}^\top (z - x_k)}^{1-1=0} \quad (28)$$

$$= (\Pi_\Delta[\nabla_{x_k}^k])^\top (z - x_k) \quad (29)$$

$$\leq \sqrt{2} \|\Pi_\Delta(\nabla_{x_k}^k)\|. \quad (30)$$

471  $\square$

472 Continuing, we can prove a bound on NashConv in terms of projected-gradient loss:

473 **Lemma 3.** *The exploitability,  $\epsilon$ , of a joint strategy  $\mathbf{x}$ , is upper bounded as a function of our proposed*  
474 *loss:*

$$\epsilon \leq \sqrt{\frac{2n}{\min_k \eta_k}} \sqrt{\mathcal{L}(\mathbf{x})}. \quad (31)$$

*Proof.*

$$\epsilon = \max_k \max_z u_k(z, x_{-k}) - u_k(\mathbf{x}) \quad (32)$$

$$\leq \sum_k \max_z u_k(z, x_{-k}) - u_k(\mathbf{x}) \quad (33)$$

$$\leq \sum_k \sqrt{2} \|\Pi_\Delta(\nabla_{x_k}^k)\|_2 \quad (34)$$

$$= \sqrt{2} \left\| \|\Pi_\Delta(\nabla_{x_1}^1)\|_2, \dots, \sqrt{2} \|\Pi_\Delta(\nabla_{x_n}^n)\|_2 \right\|_1 \quad (35)$$

$$\leq \sqrt{2n} \left\| \|\Pi_\Delta(\nabla_{x_1}^1)\|_2, \dots, \|\Pi_\Delta(\nabla_{x_n}^n)\|_2 \right\|_2 \quad (36)$$

$$= \sqrt{2n} \sqrt{\sum_k \|\Pi_\Delta(\nabla_{x_k}^k)\|_2^2} \quad (37)$$

$$\leq \sqrt{2n} \sqrt{\sum_k \left(\frac{1}{\eta_k}\right) \eta_k \|\Pi_\Delta(\nabla_{x_k}^k)\|_2^2} \quad (38)$$

$$\leq \sqrt{\frac{2n}{\min_k \eta_k}} \sqrt{\sum_k \eta_k \|\Pi_\Delta(\nabla_{x_k}^k)\|_2^2} \quad (39)$$

$$= \sqrt{\frac{2n}{\min_k \eta_k}} \sqrt{\mathcal{L}(\mathbf{x})} \quad (40)$$

475

□

476 **Lemma 4.** *The entropy regularized exploitability,  $\epsilon_{QRE}$ , of a joint strategy  $\mathbf{x}$ , is upper bounded as a*  
 477 *function of our proposed loss:*

$$\epsilon_{QRE} \leq \sqrt{\frac{2n}{\min_k \eta_k}} \sqrt{\mathcal{L}^\tau(\mathbf{x})}. \quad (41)$$

478 *Proof.* Recall that  $u_k^\tau(x_k, x_{-k})$  is also concave with respect to  $x_k$ . Then

$$\epsilon_{QRE} = \max_k \max_z u_k^\tau(z, x_{-k}) - u_k^\tau(\mathbf{x}) \quad (42)$$

$$\leq \sum_k \max_z u_k^\tau(z, x_{-k}) - u_k^\tau(\mathbf{x}) \quad (43)$$

$$\leq \sum_k \sqrt{2} \|\Pi_\Delta(\nabla_{x_k}^{k\tau})\|_2 \quad (44)$$

$$= \sqrt{2} \left\| \|\Pi_\Delta(\nabla_{x_1}^{1\tau})\|_2, \dots, \sqrt{2} \|\Pi_\Delta(\nabla_{x_n}^{n\tau})\|_2 \right\|_1 \quad (45)$$

$$\leq \sqrt{2n} \left\| \|\Pi_\Delta(\nabla_{x_1}^{1\tau})\|_2, \dots, \|\Pi_\Delta(\nabla_{x_n}^{n\tau})\|_2 \right\|_2 \quad (46)$$

$$= \sqrt{2n} \sqrt{\sum_k \|\Pi_\Delta(\nabla_{x_k}^{k\tau})\|_2^2} \quad (47)$$

$$\leq \sqrt{2n} \sqrt{\sum_k \left(\frac{1}{\eta_k}\right) \eta_k \|\Pi_\Delta(\nabla_{x_k}^{k\tau})\|_2^2} \quad (48)$$

$$\leq \sqrt{\frac{2n}{\min_k \eta_k}} \sqrt{\sum_k \eta_k \|\Pi_\Delta(\nabla_{x_k}^{k\tau})\|_2^2} \quad (49)$$

$$= \sqrt{\frac{2n}{\min_k \eta_k}} \sqrt{\mathcal{L}^\tau(\mathbf{x})} \quad (50)$$

479

□

480 **A.4 Unbiased Estimation**

481 **Lemma 5.** *An unbiased estimate of  $\mathcal{L}(\mathbf{x})$  can be obtained by drawing two samples (pure strategies)*  
 482 *from each players' mixed strategy and observing payoffs.*

483 *Proof.* Define  $\nabla_{x_k}^k$  as the random variable distributed according to the distribution induced by all  
 484 players' mixed strategies. Let  $\nabla_{x_k}^{k(1)}$  and  $\nabla_{x_k}^{k(2)}$  represent two other independent random variables,  
 485 distributed identically to  $\nabla_{x_k}^k$ . Then

$$\mathbb{E}_{a_k \sim x_k \forall k}[\mathcal{L}(\mathbf{x})] = \mathbb{E}_{a_k \sim x_k \forall k} \left[ \sum_k \eta_k ( \|\nabla_{x_k}^k\|^2 - \frac{1}{m_k} (\mathbf{1}^\top \nabla_{x_k}^k)^2 ) \right] \quad (51)$$

$$= \sum_k \eta_k ( \mathbb{E}_{a_k \sim x_k \forall k} [\|\nabla_{x_k}^k\|^2] - \frac{1}{m_k} \mathbb{E}_{a_k \sim x_k \forall k} [(\mathbf{1}^\top \nabla_{x_k}^k)^2] ) \quad (52)$$

$$= \sum_k \eta_k ( \mathbb{E}_{a_k \sim x_k \forall k} [\sum_l (\nabla_{x_{kl}}^k)^2] - \frac{1}{m_k} \mathbb{E}_{a_k \sim x_k \forall k} [(\sum_l \nabla_{x_{kl}}^k)^2] ) \quad (53)$$

$$= \sum_k \eta_k ( \sum_l \mathbb{E}_{a_k \sim x_k \forall k} [(\nabla_{x_{kl}}^k)^2] - \frac{1}{m_k} \mathbb{E}_{a_k \sim x_k \forall k} [(\sum_l \nabla_{x_{kl}}^k)^2] ) \quad (54)$$

$$= \sum_k \eta_k \left( \sum_l \mathbb{E}_{a_k \sim x_k \forall k} [\nabla_{x_{kl}}^{k(1)}] \mathbb{E}_{a_k \sim x_k \forall k} [\nabla_{x_{kl}}^{k(2)}] \right) \quad (55)$$

$$- \frac{1}{m_k} \mathbb{E}_{a_k \sim x_k \forall k} \left[ \sum_l \nabla_{x_{kl}}^{k(1)} \right] \mathbb{E}_{a_k \sim x_k \forall k} \left[ \sum_l \nabla_{x_{kl}}^{k(2)} \right] \quad (56)$$

$$= \sum_k \eta_k \left( \sum_l \mathbb{E}_{a_j \sim x_j \forall j \neq k} [\nabla_{x_{kl}}^{k(1)}] \mathbb{E}_{a_j \sim x_j \forall j \neq k} [\nabla_{x_{kl}}^{k(2)}] \right) \quad (57)$$

$$- \frac{1}{m_k} \mathbb{E}_{a_j \sim x_j \forall j \neq k} \left[ \sum_l \nabla_{x_{kl}}^{k(1)} \right] \mathbb{E}_{a_j \sim x_j \forall j \neq k} \left[ \sum_l \nabla_{x_{kl}}^{k(2)} \right] \quad (58)$$

$$= \sum_k \eta_k \left( [\hat{\nabla}_{x_k}^{k(1)}]^\top \hat{\nabla}_{x_k}^{k(2)} - \frac{1}{m_k} (\mathbf{1}^\top \hat{\nabla}_{x_k}^{k(1)}) (\mathbf{1}^\top \hat{\nabla}_{x_k}^{k(2)}) \right) \quad (59)$$

$$= \sum_k \eta_k \underbrace{\left( \hat{\nabla}_{x_k}^{k(1)} - \frac{1}{m_k} (\mathbf{1}^\top \hat{\nabla}_{x_k}^{k(1)}) \mathbf{1} \right)^\top}_{\text{appx. br gradient}} \underbrace{\hat{\nabla}_{x_k}^{k(2)}}_{\text{exp. payoffs}} \quad (60)$$

486 where  $\hat{\nabla}_{x_k}^{k(p)}$  is an unbiased estimator of player  $k$ 's gradient.

487 □

488 **Lemma 6.** *The loss formed as the sum of the squared norms of the projected-gradients,  $\mathcal{L}^\tau$ , can be*  
 489 *decomposed into three terms as follows:*

$$\mathcal{L}^\tau(\mathbf{x}) = \underbrace{\sum_k \eta_k x_q^\top B_{kq}^\top B_{kq} x_q}_{(A)} + 2 \underbrace{\sum_k \eta_k E_k^\top B_{kq} x_q}_{(B)} + \underbrace{\sum_k \eta_k E_k^\top E_k}_{(C)} \quad (61)$$

490 where  $q$  is any player other than  $k$ .

491 *Proof.* Let  $S^\tau = -\tau \sum_l x_{kl} \log(x_{kl})$  so that  $\frac{\partial S^\tau}{\partial x_k} = -\tau(\ln(x_k) + \mathbf{1})$ . Note that  $\Pi_{T\Delta}[\frac{\partial S^\tau}{\partial x_k}] =$   
 492  $-\tau \ln(x_k)$ .

$$\mathcal{L}^\tau(\mathbf{x}) = \sum_k \eta_k (\Pi_{T\Delta}[\nabla_{x_k}^k])^\top \Pi_{T\Delta}[\nabla_{x_k}^k] \quad (62)$$

$$= \sum_k \eta_k [H_{kq}^k x_q + \frac{\partial S^\tau}{\partial x_k}]^\top [I - \frac{1}{m_k} \mathbf{1}\mathbf{1}^\top] [I - \frac{1}{m_k} \mathbf{1}\mathbf{1}^\top] [H_{kq}^k x_q + \frac{\partial S^\tau}{\partial x_k}] \quad (63)$$

$$= \sum_k \eta_k \left( x_q^\top [H_{kq}^k]^\top [I - \frac{1}{m_k} \mathbf{1}\mathbf{1}^\top]^2 [H_{kq}^k] x_q + 2 [\frac{\partial S^\tau}{\partial x_k}]^\top [I - \frac{1}{m_k} \mathbf{1}\mathbf{1}^\top]^2 [H_{kq}^k x_q] \right. \quad (64)$$

$$\left. + [\frac{\partial S^\tau}{\partial x_k}]^\top [I - \frac{1}{m_k} \mathbf{1}\mathbf{1}^\top]^2 [\frac{\partial S^\tau}{\partial x_k}] \right) \quad (65)$$

$$= \underbrace{\sum_k \eta_k x_q^\top B_{kq}^\top B_{kq} x_q}_{(A)} + 2 \underbrace{\sum_k \eta_k E_k^\top B_{kq} x_q}_{(B)} + \underbrace{\sum_k \eta_k E_k^\top E_k}_{(C)} \quad (66)$$

493 where  $B_{kq} = [I - \frac{1}{m_k} \mathbf{1}\mathbf{1}^\top] H_{kq}^k$  and  $E_k = [I - \frac{1}{m_k} \mathbf{1}\mathbf{1}^\top] [\frac{\partial S^\tau}{\partial x_k}] = -\tau \ln(x_k)$ .

494

□

## 495 A.5 Bound on Loss

496 By equation (51), we can also rewrite this loss as a weighted sum of 2-norms,  $\mathcal{L}(\mathbf{x}) = \sum_k \eta_k \|\nabla_{x_k}^k -$   
 497  $\mu_k\|_2^2$  where  $\mu_k = \frac{1}{m_k} (\mathbf{1}^\top \nabla_{x_k}^k) \in [0, 1]$  for brevity. This will allow us to more easily analyze our  
 498 loss.

499 **Lemma 7.** Assume payoffs are bounded by 1, then setting  $\eta_k \leq \frac{4}{nm_k}$  or  $\eta_k \leq \frac{4}{n\bar{m}}$  or  $\sum_k \eta_k \leq \frac{4}{n}$   
 500 ensures  $0 \leq \mathcal{L}(x) \leq 1$  for all  $x \in \mathcal{X}$ .

*Proof.*

$$0 \leq \mathcal{L}(\mathbf{x}) = \sum_k \eta_k \|\nabla_{x_k}^k - \mu_k\|_2^2 \quad (67)$$

$$= \sum_k \eta_k m_k \left[ \frac{1}{m_k} \sum_l (\nabla_{x_{kl}}^k - \mu_k)^2 \right] \quad (68)$$

$$= \sum_k \eta_k m_k \text{Var}[\nabla_{x_k}^k] \quad (69)$$

$$\leq \frac{1}{4} \sum_k \eta_k m_k \quad (70)$$

$$\leq \frac{1}{4} (\max_k \eta_k) (\sum_k m_k) \quad (71)$$

$$= \frac{1}{4} (\max_k \eta_k) n\bar{m} \leq 1 \quad (72)$$

$$\implies (\max_k \eta_k) \leq \frac{4}{n\bar{m}}. \quad (73)$$

501

□

502 The  $k$ th element of the sum in the loss does not depend on agent  $k$ 's strategy. We will rewrite the loss  
 503 to make its dependence on all other players' strategies more obvious ( $l, q \neq k$  below).

$$\mathcal{L}(\mathbf{x}) = \sum_k \eta_k ([H_{kl}^k x_l]^\top [H_{kq}^k x_q] - \frac{1}{m_k} (\mathbf{1}^\top [H_{kl}^k x_l]) (\mathbf{1}^\top [H_{kq}^k x_q])) \quad (74)$$

$$= \sum_k \eta_k ([H_{kl}^k x_l]^\top [H_{kq}^k x_q] - \frac{1}{m_k} [H_{kl}^k x_l]^\top \mathbf{1} \mathbf{1}^\top [H_{kq}^k x_q]) \quad (75)$$

$$= \sum_k \eta_k [H_{kl}^k x_l]^\top [I - \frac{1}{m_k} \mathbf{1} \mathbf{1}^\top] [H_{kq}^k x_q] \quad (76)$$

$$= \sum_k \eta_k x_l^\top [H_{kl}^k]^\top [I - \frac{1}{m_k} \mathbf{1} \mathbf{1}^\top] [H_{kq}^k] x_q \quad (77)$$

$$= \sum_k \eta_k x_q^\top [H_{kq}^k]^\top [I - \frac{1}{m_k} \mathbf{1} \mathbf{1}^\top] [H_{kq}^k] x_q \quad \text{isolate dep. on } q \quad (78)$$

$$= \sum_k \eta_k x_q^\top [H_{qk}^k] [I - \frac{1}{m_k} \mathbf{1} \mathbf{1}^\top] [H_{kq}^k] x_q \quad (79)$$

$$= \sum_k \eta_k x_q^\top A_{qkq} x_q. \quad (80)$$

504 where  $A_{qkq} = [H_{qk}^k] [I - \frac{1}{m_k} \mathbf{1} \mathbf{1}^\top] [H_{kq}^k]$  does not depend on  $x_k$ .

505 Note this means we can also write  $\mathcal{L}(\mathbf{x}) = \sum_k \eta_k x_l^\top A_{lkq} x_q$  for any  $l, q \neq k$ .

506 **Lemma 8.** Assume payoffs are bounded in  $[0, 1]$ , then

$$|\mathcal{L}^\tau(\mathbf{x})| \leq \frac{1}{4} (\max_k \eta_k) n \bar{m} \left( \frac{\ln(m^*)}{\ln(1/p)} + 2 \right)^2 \quad (81)$$

507 for any  $\mathbf{x}$  such that  $x_{kl} \geq \frac{p}{m^*} \forall k, l$ .

508 *Proof.* Starting from the definition of  $\mathcal{L}^\tau$  and applying Lemma 19 along with intermediate results  
 509 from Lemma 16, we find

$$|\mathcal{L}^\tau(\mathbf{x})| = \left| \sum_k \eta_k \|\Pi_{T\Delta}(\nabla_{x_k}^k)\|^2 \right| \quad (82)$$

$$\leq \frac{1}{4} \sum_k \eta_k m_k \left( \tau \ln\left(\frac{1}{x_{\min}}\right) + 1 \right)^2 \quad (83)$$

$$= \frac{1}{4} \sum_k \eta_k m_k \left( \frac{1}{\ln(1/p)} \ln\left(\frac{m^*}{p}\right) + 1 \right)^2 \quad (84)$$

$$= \frac{1}{4} \sum_k \eta_k m_k \left( \frac{\ln(m^*)}{\ln(1/p)} + 2 \right)^2 \quad (85)$$

$$= \frac{1}{4} (\max_k \eta_k) n \bar{m} \left( \frac{\ln(m^*)}{\ln(1/p)} + 2 \right)^2. \quad (86)$$

510 □

## 511 B QREs Approximate NEs at Low Temperature

512 **Lemma 9.** Setting  $\tau = \ln(1/p)^{-1}$  with  $p \in [0, 1)$  ensures that all QREs contain probabilities greater  
 513 than  $\frac{p}{\max_k m_k}$ .

514 *Proof.* What must  $\tau$  be to ensure  $x_{lp} \geq x_{\min}$  for any  $l, p$ ? We can check the case where  $\nabla = e_i$ . Let  
 515  $m^* = \max_k m_k$ . Then

$$x_{\min} = \min_k \min_{\nabla_{x_k}^k} \min_l [\text{softmax}(\frac{\nabla_{x_k}^k}{\tau})]_l \quad (87)$$

$$= \frac{e^0}{(m^* - 1)e^{\frac{1}{\tau}} + e^0} \quad (88)$$

$$= \frac{1}{(m^* - 1)e^{\frac{1}{\tau}} + 1} \quad (89)$$

$$\implies e^{\frac{1}{\tau}} = \frac{1}{m^* - 1} \left( \frac{1}{x_{\min}} - 1 \right) \quad (90)$$

$$\implies \tau = \frac{1}{\ln\left(\frac{1}{m^* - 1} \left( \frac{1}{x_{\min}} - 1 \right)\right)}. \quad (91)$$

516 If  $x_{\min} = \frac{p}{m^*}$  with  $p \in [0, 1]$ , then

$$\tau^* = \frac{1}{\ln\left(\frac{1}{m^* - 1} \left( \frac{1}{x_{\min}} - 1 \right)\right)} \quad (92)$$

$$= \frac{1}{\ln\left(\frac{1}{m^* - 1} \left( \frac{m^*}{p} - 1 \right)\right)} \quad (93)$$

$$= \frac{1}{\ln\left(\frac{m^* - p}{m^* - 1} \frac{1}{p}\right)} \quad (94)$$

$$\leq \frac{1}{\ln\left(\frac{1}{p}\right)}. \quad (95)$$

517 This implies if we set  $\tau = \ln(1/p)^{-1}$ , then we are guaranteed that all QREs contain probabilities  
 518 greater than  $x_{\min} = \frac{p}{\max_k m_k}$ .  $\square$

519 **Lemma 10** (Repeated from Lemma 1 of [29]). Let  $\nabla_{x_k}^k$  be player  $k$ 's gradient ( $m_k \geq 2$ ) with payoffs  
 520 bounded in  $[0, 1]$  and  $x$  be a QRE at temperature  $\tau$ . Then it holds that

$$u_k(BR_k, x_{-k}) - u_k(x) = \max(\nabla_{x_k}^k) - (\nabla_{x_k}^k)^\top \text{softmax}\left(\frac{\nabla_{x_k}^k}{\tau}\right) \leq \tau(W(1/e) + \frac{m_k - 2}{e}) \quad (96)$$

521 where  $W$  is the Lambert function ( $W(1/e) \approx 0.278$ ).

522 **Lemma 11** (Slightly modified from Proposition 5.1a of [6]). Let  $\psi_e(x_k) = \sum_l x_{kl} \ln(x_{kl})$  if  
 523  $x_k \in \Delta^{m_k - 1}$  else  $+\infty$ . Then  $\psi_e(x_k)$  is 1-strongly convex over  $\text{int}\Delta^{m_k - 1}$  w.r.t. the  $\|\cdot\|_1$  and  $\|\cdot\|_2$   
 524 norms, i.e.,

$$\langle \nabla \psi_e(x) - \nabla \psi_e(y), x - y \rangle \geq \|x - y\|_1^2 \geq \|x - y\|_2^2 \quad (97)$$

$$\implies \psi_e(y) \geq \psi_e(x) + \nabla \psi_e(x)^\top (y - x) + \frac{1}{2} \|y - x\|_2^2. \quad (98)$$

525 for all  $x, y \in \text{int}\Delta^{m_k - 1}$ .

526 **Lemma 12.** Let  $l(x|x_k) = \langle \nabla f_k(x_k), x \rangle + \frac{1}{t_k} B_{\psi_e}(x, x_k)$  where  $t_k > 0$ ,  $f_k(z) = -\epsilon_k(z) =$   
 527  $-[u_k(z, x_{-i}) + S^\tau(z) - u_k(x) - S^\tau(x_k)]$ , and  $B_{\psi_e}(x, x_k) = \psi_e(x) - \psi_e(y) - \langle x - y, \nabla \psi_e(y) \rangle$   
 528 with  $\psi_e$  defined in Lemma 11. Finally, let  $x_{k+1} = \arg \min_{x \in \text{int}\Delta} l(x|x_k)$ . Then

$$\|x_k - x_{k+1}\| \leq 2 \|\Pi_\Delta(\nabla_{x_k}^{k\tau})\|. \quad (99)$$

529 *Proof.* Plugging  $\psi_e(x_k) = \sum_l x_{kl} \ln(x_{kl}) = -S(x_k)$  on  $\text{int}\Delta$  into  $B_{\psi_e}(x, x_k)$ , we find

$$B_{\psi_e}(x, x_k) = S(x_k) - S(x) - \langle \ln(x_k) + 1, x - x_k \rangle \quad (100)$$

$$= S(x_k) - S(x) - \langle \ln(x_k), x - x_k \rangle \quad (101)$$

530 for all  $x, x_k \in \text{int}\Delta$ . Note that  $-S(x)$  is 1-strongly convex on  $\text{int}\Delta$ , therefore,  $B_{\psi_e}(x, x_k)$  is also  
 531 1-strongly convex in  $x$ . Continuing, this also implies  $l(x|x_k)$  is 1-strongly convex.

532 Let  $x_{k+1} = \arg \min_{x \in \text{int}\Delta} l(x|x_k)$  and note that  $\nabla f_k(x_k) = -\nabla \epsilon_k = -\nabla_{x_k}^{k\tau}$ . Strong convexity of  
 533  $l$  implies

$$l(x_{k+1}) \geq l(x_k) + \nabla_x l(x_k)^\top (x_{k+1} - x_k) + \frac{1}{2} \|x_{k+1} - x_k\|_2^2 \quad (102)$$

$$\implies \|x_k - x_{k+1}\|_2^2 \leq 2 \left[ \underbrace{l(x_{k+1}) - l(x_k)}_{\leq 0} + \nabla_x l(x_k)^\top (x_k - x_{k+1}) \right] \quad (103)$$

$$\leq 2 \nabla_x l(x_k)^\top (x_k - x_{k+1}) = 2(\nabla f_k(x_k) + \frac{1}{t_k} [\ln(x_k) + \mathbf{1} - \ln(x_k)])^\top (x_k - x_{k+1}) \quad (104)$$

$$= 2 \nabla f_k(x_k)^\top (x_k - x_{k+1}) = 2(\nabla_{x_k}^{k\tau})^\top (x_{k+1} - x_k) = 2(\Pi_\Delta(\nabla_{x_k}^{k\tau}))^\top (x_{k+1} - x_k) \quad (105)$$

$$\leq 2 \|\Pi_\Delta(\nabla_{x_k}^{k\tau})\| \|x_k - x_{k+1}\|. \quad (106)$$

534 Rearranging the inequality achieves the desired result.  $\square$

535 **Lemma 13.** [Low Temperature Approximate QREs are Approximate Nash Equilibria] Let  $\nabla_{x_k}^{k\tau}$  be  
 536 player  $k$ 's entropy regularized gradient with payoffs bounded in  $[0, 1]$  and  $\mathbf{x}$  be an approximate QRE.  
 537 Then it holds that

$$u_k(\text{BR}_k, x_{-k}) - u_k(\mathbf{x}) \leq \tau(W(1/e) + \frac{m_k - 2}{e}) + 2\sqrt{m_k} \|\Pi_\Delta(\nabla_{x_k}^{k\tau})\| \quad (107)$$

538 where  $W$  is the Lambert function ( $W(1/e) \approx 0.278$ ).

539 *Proof.* First note that  $x_k = \text{softmax}(\ln(x_k))$  for  $x_k \in \text{int}\Delta$ . Recall that the `softmax` is invariant  
 540 to constant offsets to its argument, i.e.,  $\text{softmax}(z + c\mathbf{1}) = \text{softmax}(z)$  for any  $c \in \mathbb{R}$ . Then

$$\text{softmax}\left(\frac{\nabla_{x_k}^k}{\tau}\right) = \text{softmax}\left(\ln(x_k) - \frac{1}{\tau}[-\nabla_{x_k}^k + \tau \ln(x_k)]\right) \quad (108)$$

$$= \text{softmax}\left(\ln(x_k) - \frac{1}{\tau}[-\nabla_{x_k}^k + \tau \ln(x_k) + \tau \mathbf{1}]\right) \quad (109)$$

$$= \text{softmax}\left(\ln(x_k) - \frac{1}{\tau} \nabla f_k(x_k)\right) \quad (110)$$

$$= \arg \min_{x \in \text{int}\Delta} l(x|x_k) \text{ with } t_k = \frac{1}{\tau} \quad (111)$$

$$= x_k^* \quad (112)$$

541 where the closed-form solution to the minimization problem as a `softmax` formula comes from  
 542 inspecting the Entropic Descent Algorithm (EDA) of [6].

543 Then, beginning with the definition of exploitability, we find

$$u_k(\text{BR}_k, x_{-k}) - u_k(\mathbf{x}) = u_k(\text{BR}_k, x_{-k}) - (\nabla_{x_k}^k)^\top x_k \quad (113)$$

$$= u_k(\text{BR}_k, x_{-k}) - (\nabla_{x_k}^k)^\top \text{softmax}\left(\frac{\nabla_{x_k}^k}{\tau}\right) - (\nabla_{x_k}^k)^\top (x_k - \text{softmax}\left(\frac{\nabla_{x_k}^k}{\tau}\right)) \quad (114)$$

$$\leq \tau(W(1/e) + \frac{m_k - 2}{e}) + \|\nabla_{x_k}^k\| \cdot \|x_k - \text{softmax}\left(\frac{\nabla_{x_k}^k}{\tau}\right)\| \quad (115)$$

$$= \tau(W(1/e) + \frac{m_k - 2}{e}) + \|\nabla_{x_k}^k\| \cdot \|x_k - x_k^*\| \quad (116)$$

$$\leq \tau(W(1/e) + \frac{m_k - 2}{e}) + 2\sqrt{m_k} \|\Pi_\Delta(\nabla_{x_k}^{k\tau})\|. \quad (117)$$

544  $\square$

545 **Lemma 14.** [ $\mathcal{L}^\tau$  Scores Nash Equilibria] Let  $\mathcal{L}^\tau(\mathbf{x})$  be our proposed entropy regularized loss  
 546 function with payoffs bounded in  $[0, 1]$  and  $\mathbf{x}$  be an approximate QRE. Then it holds that

$$\epsilon \leq n\tau(W(1/e) + \frac{\bar{m} - 2}{e}) + 2\sqrt{\frac{n \max_k m_k}{\min_k \eta_k}} \sqrt{\mathcal{L}^\tau(\mathbf{x})} \quad (118)$$

547 where  $W$  is the Lambert function ( $W(1/e) \approx 0.278$ ).

548 *Proof.* Beginning with the definition of exploitability and applying Lemma 13, we find

$$\epsilon = \max_k u_k(\text{BR}_k, x_{-k}) - u_k(\mathbf{x}) \quad (119)$$

$$\leq \sum_k u_k(\text{BR}_k, x_{-k}) - u_k(\mathbf{x}) \quad (120)$$

$$\leq \sum_k \left[ \tau(W(1/e) + \frac{m_k - 2}{e}) + 2\sqrt{m_k} \|\Pi_\Delta(\nabla_{x_k}^{k\tau})\| \right] \quad (121)$$

$$= n\tau(W(1/e) + \frac{\bar{m} - 2}{e}) + 2 \sum_k \sqrt{m_k} \|\Pi_\Delta(\nabla_{x_k}^{k\tau})\| \quad (122)$$

$$\leq n\tau(W(1/e) + \frac{\bar{m} - 2}{e}) + 2\sqrt{\max_k m_k} \sum_k \|\Pi_\Delta(\nabla_{x_k}^{k\tau})\| \quad (123)$$

$$\leq n\tau(W(1/e) + \frac{\bar{m} - 2}{e}) + 2\sqrt{\frac{n \max_k m_k}{\min_k \eta_k}} \sqrt{\mathcal{L}^\tau(\mathbf{x})}. \quad (124)$$

549 where the last inequality follows from the same steps outlined in Lemma 3, which established the  
 550 relationship between  $\mathcal{L}(\mathbf{x})$  and  $\epsilon$ .

551 □

## 552 C Gradient of Loss

553 **Lemma 15.** The gradient of  $\mathcal{L}^\tau(\mathbf{x})$  with respect to player  $l$ 's strategy  $x_l$  is

$$\nabla_{x_l} \mathcal{L}(\mathbf{x}) = 2 \sum_k \eta_k B_{kl}^\top \Pi_{T\Delta}(\nabla_{x_k}^{k\tau}) \quad (125)$$

554 where  $B_{ll} = -\tau[I - \frac{1}{m_l} \mathbf{1}\mathbf{1}^\top] \text{diag}(\frac{1}{x_l})$  and  $B_{kl} = [I - \frac{1}{m_k} \mathbf{1}\mathbf{1}^\top] H_{kl}^k$  for  $k \neq l$ .

555 *Proof.* Recall from Lemma 6 that the loss can be decomposed as  $\mathcal{L}^\tau(\mathbf{x}) = (A) + (B) + (C)$ .

556 Then

$$D_{x_l}[(A)] = D_{x_l} \left[ \sum_k \eta_k x_q^\top B_{kq}^\top B_{kq} x_q \right] = 2 \sum_{k \neq l} \eta_k B_{kl}^\top B_{kl} x_l \quad (126)$$

557 where  $q \neq k$  and  $B_{kq} = [I - \frac{1}{m_k} \mathbf{1}\mathbf{1}^\top] [H_{kq}^k]$  does not depend on  $x_k$ .

558 Also, letting  $B_{ll} = -\tau[I - \frac{1}{m_l} \mathbf{1}\mathbf{1}^\top] \text{diag}(\frac{1}{x_l})$ ,

$$D_{x_l}[(B)] = D_{x_l} \left[ -2\tau \sum_k \eta_k \ln(x_k)^\top B_{kq} x_q \right] \quad (127)$$

$$= -2\tau \left[ \eta_l D_{x_l} [\ln(x_l)^\top B_{lq} x_q] + \sum_{k \neq l} \eta_k D_{x_l} [\ln(x_k)^\top B_{kl} x_l] \right] \quad (128)$$

$$= -2\tau \left[ \eta_l \text{diag}(\frac{1}{x_l}) B_{lq} x_q + \sum_{k \neq l} \eta_k B_{kl}^\top \ln(x_k) \right] \quad (129)$$

$$= -2\tau \left[ \eta_l \left( [I - \frac{1}{m_l} \mathbf{1}\mathbf{1}^\top] \text{diag}(\frac{1}{x_l}) \right)^\top \Pi_{T\Delta}(\nabla^l) + \sum_{k \neq l} \eta_k B_{kl}^\top \ln(x_k) \right] \quad (130)$$

$$= 2 \left[ \eta_l B_{ll}^\top \Pi_{T\Delta}(\nabla^l) - \tau \sum_{k \neq l} \eta_k B_{kl}^\top \ln(x_k) \right]. \quad (131)$$

559 And

$$D_{x_l}[(C)] = D_{x_l} \left[ \sum_k \eta_k \tau^2 \ln(x_k)^\top \left[ I - \frac{1}{m_k} \mathbf{1}\mathbf{1}^\top \right] \ln(x_k) \right] \quad (132)$$

$$= 2\tau^2 \left[ \eta_l \text{diag}\left(\frac{1}{x_l}\right) \left[ I - \frac{1}{m_l} \mathbf{1}\mathbf{1}^\top \right] \ln(x_l) \right] \quad (133)$$

$$= -2\tau \eta_l \left( \left[ I - \frac{1}{m_l} \mathbf{1}\mathbf{1}^\top \right] \text{diag}\left(\frac{1}{x_l}\right) \right)^\top \Pi_{T\Delta}(-\tau \ln(x_l)) \quad (134)$$

$$= 2\eta_l B_{ll}^\top \Pi_{T\Delta}(-\tau \ln(x_l)). \quad (135)$$

560 Putting these together, we find

$$\nabla_{x_l} \mathcal{L}(\mathbf{x}) = 2 \sum_{k \neq l} \eta_k B_{kl}^\top (B_{kl} x_l - \tau \ln(x_k)) + 2\eta_l B_{ll}^\top \left[ \Pi_{T\Delta}(\nabla^l) + \Pi_{T\Delta}(-\tau \ln(x_l)) \right] \quad (136)$$

$$= 2\eta_l B_{ll}^\top \Pi_{T\Delta}(\nabla_{x_k}^{k\tau}) + 2 \sum_{k \neq l} \eta_k B_{kl}^\top \Pi_{T\Delta}(\nabla_{x_k}^{k\tau}) \quad (137)$$

$$= 2 \sum_k \eta_k B_{kl}^\top \Pi_{T\Delta}(\nabla_{x_k}^{k\tau}). \quad (138)$$

561

□

## 562 C.1 Unbiased Estimation

563 In order to construct an unbiased estimate of  $A_{lk}$ , we will need to form two independent unbiased  
 564 estimates of  $H_{kl}^k$ . Recall that  $H_{kl}^k$  is simply the expected bimatrix game between players  $k$  and  $l$   
 565 when all other players sample their actions according to their current strategies.

## 566 C.2 Bound on Gradient / Lipschitz Property

567 **Lemma 16.** *Assume payoffs are upper bounded by 1, then the infinity norm of the gradient is bounded*  
 568 *as*

$$\|\nabla_{\mathbf{x}} \mathcal{L}^\tau(\mathbf{x})\|_\infty \leq \frac{1}{2} (\max_k \eta_k) \left( \tau \ln\left(\frac{1}{x_{\min}}\right) + 1 \right) \left[ \tau m^* \left( \frac{1}{x_{\min}} - 1 \right) + n\bar{m} \right]. \quad (139)$$

569 *Proof.* Recall from Lemma 15 that the gradient of  $\mathcal{L}(\mathbf{x})$  with respect to player  $l$ 's strategy  $x_l$  is

$$\nabla_{x_l} \mathcal{L}(\mathbf{x}) = 2 \sum_k \eta_k B_{kl}^\top \Pi_{T\Delta}(\nabla_{x_k}^{k\tau}) \quad (140)$$

570 where  $B_{ll} = -\tau \left[ I - \frac{1}{m_l} \mathbf{1}\mathbf{1}^\top \right] \text{diag}\left(\frac{1}{x_l}\right)$  and  $B_{kl} = \left[ I - \frac{1}{m_k} \mathbf{1}\mathbf{1}^\top \right] H_{kl}^k$  for  $k \neq l$ .

571 For payoffs in  $[0, 1]$ , the entries in  $\nabla_{x_k}^{k\tau} = \nabla_{x_k}^k - \tau \ln(x_k)$  are bounded within  $[0, \tau \ln(\frac{1}{x_{\min}}) + 1]$   
 572 with a range  $\tau \ln(\frac{1}{x_{\min}}) + 1$ . Similarly, the entries in  $-\tau \text{diag}\left(\frac{1}{x_l}\right)$  are bounded within  $[-\tau \frac{1}{x_{\min}}, -\tau]$   
 573 with a range of  $\tau \left( \frac{1}{x_{\min}} - 1 \right)$ .

574 The infinity norm of the gradient can then be bounded as

$$\|\nabla_{\mathbf{x}} \mathcal{L}^\tau(\mathbf{x})\|_\infty = \max_l \|\nabla_{x_l} \mathcal{L}(\mathbf{x})\|_\infty \quad (141)$$

$$= \max_l \left\| 2 \sum_k \eta_k B_{kl}^\top \Pi_{T\Delta}(\nabla_{x_k}^{k\tau}) \right\|_\infty \quad (142)$$

$$\leq 2 \sum_k \eta_k \max_l \left\| B_{kl}^\top \Pi_{T\Delta}(\nabla_{x_k}^{k\tau}) \right\|_\infty \quad (143)$$

$$\leq \frac{1}{2} \sum_{k \neq l^*} \eta_k m_k (\tau \ln(\frac{1}{x_{\min}}) + 1) + \frac{1}{2} \eta_{l^*} m_{l^*} \tau (\frac{1}{x_{\min}} - 1) (\tau \ln(\frac{1}{x_{\min}}) + 1) \quad (144)$$

$$= \frac{1}{2} (\tau \ln(\frac{1}{x_{\min}}) + 1) \left[ \eta_{l^*} m_{l^*} \tau (\frac{1}{x_{\min}} - 1) + \sum_{k \neq l^*} \eta_k m_k \right] \quad (145)$$

$$\leq \frac{1}{2} (\max_k \eta_k) (\tau \ln(\frac{1}{x_{\min}}) + 1) \left[ \tau m_{l^*} (\frac{1}{x_{\min}} - 1) + \sum_{k \neq l^*} m_k \right] \quad (146)$$

$$\leq \frac{1}{2} (\max_k \eta_k) (\tau \ln(\frac{1}{x_{\min}}) + 1) \left[ \tau m^* (\frac{1}{x_{\min}} - 1) + n\bar{m} \right] \quad (147)$$

575 where the second inequality follows from Lemma 19.

576 □

577 **Corollary 2.** *If  $\tau$  is set according to Lemma 9, then the infinity norm of the gradient is bounded as*

$$\|\nabla_{\mathbf{x}} \mathcal{L}^\tau(\mathbf{x})\|_\infty \leq \frac{1}{2} (\max_k \eta_k) \left[ \frac{\ln(m^*)}{\ln(1/p)} + 2 \right] \left[ \frac{m^{*2}}{p \ln(1/p)} + n\bar{m} \right] = \frac{1}{2} (\max_k \eta_k) \hat{L} \quad (148)$$

578 where  $m^* = \max_k m_k$  and  $\hat{L}$  is defined implicitly for convenience in other derivations.

579 *Proof.* Starting with Lemma 16 and applying Lemma 9 (i.e.,  $\tau = \ln(1/p)^{-1}$  and  $x_{\min} = \frac{p}{m^*}$  where  
580  $m^* = \max_k m_k$ ), we find

$$\|\nabla_{\mathbf{x}} \mathcal{L}^\tau(\mathbf{x})\|_\infty \leq \frac{1}{2} (\max_k \eta_k) (\tau \ln(\frac{1}{x_{\min}}) + 1) \left[ \tau m^* (\frac{1}{x_{\min}} - 1) + n\bar{m} \right] \quad (149)$$

$$= \frac{1}{2} (\max_k \eta_k) \left[ \frac{\ln(m^*/p)}{\ln(1/p)} + 1 \right] \left[ \frac{m^*}{\ln(1/p)} (\frac{m^*}{p} - 1) + n\bar{m} \right] \quad (150)$$

$$\leq \frac{1}{2} (\max_k \eta_k) \left[ \frac{\ln(m^*)}{\ln(1/p)} + 2 \right] \left[ \frac{m^{*2}}{p \ln(1/p)} + n\bar{m} \right]. \quad (151)$$

581 As  $p \rightarrow 0^+$ , the norm of the gradient blows up because the gradient of Shannon entropy blows up  
582 for small probabilities. As  $p \rightarrow 1$ , the norm of the gradient blows up because we require infinite  
583 temperature  $\tau$  to guarantee all QREs are nearly uniform; recall  $\tau$  is the regularization coefficient on  
584 the entropy bonus terms which means our modified utilities blow up for large  $\tau$ . In practice, setting  $p$   
585 to  $\mathcal{O}(1)$ , e.g.,  $p = \frac{1}{10}$  is sufficient. □

## 586 D Hessian of Loss

587 We will now derive the Hessian of our loss. This will be useful in establishing properties about global  
588 minima that enable the application of tailored minimization algorithms. Let  $D_z[f(z)]$  denote the  
589 differential operator applied to (possibly multivalued) function  $f$  with respect to  $z$ . For example,  
590  $D_{x_q}[H_{lk}^k] = D_{x_q}[x_q T_{qlk}^k] = T_{qlk}^k$  where  $T_{qlk}^k$  is player  $k$ 's payoff tensor according to the three-way  
591 approximation between players  $k, l$ , and  $q$  to the game at  $\mathbf{x}$ .

592 **Lemma 17.** *The Hessian of  $\mathcal{L}^\tau(\mathbf{x})$  can be written*

$$\text{Hess}(\mathcal{L}^\tau) = 2\tilde{B}^\top \tilde{B} + T \Pi_{T\Delta}(\tilde{\nabla}^\tau) \quad (152)$$

593 where  $\tilde{B}_{kl} = \sqrt{\eta_k} B_{kl}$ ,  $\Pi_{T\Delta}(\tilde{\nabla}^\tau) = [\eta_1 \Pi_{T\Delta}(\nabla_{x_1}^{1\tau}), \dots, \eta_n \Pi_{T\Delta}(\nabla_{x_n}^{n\tau})]$ , and we augment  $T$  (the  
594 3-player tensor approximation to the game,  $T_{lqk}^k$ ) so that  $T_{ll}^l = \tau \text{diag} \mathfrak{Z}(\frac{1}{x_l^\tau})$  and otherwise 0.

595 *Proof.* Recall the gradient of our proposed loss:

$$\nabla_{x_l} \mathcal{L}(\mathbf{x}) = 2 \sum_k \eta_k B_{kl}^\top \Pi_{T\Delta}(\nabla_{x_k}^{k\tau}) \quad (153)$$

596 where  $B_{ll} = -\tau[I - \frac{1}{m_l} \mathbf{1}\mathbf{1}^\top] \text{diag}(\frac{1}{x_l})$  and  $B_{kl} = [I - \frac{1}{m_k} \mathbf{1}\mathbf{1}^\top] H_{kl}^k$  for  $k \neq l$ .

597 Consider the following Jacobians, which will play an auxiliary role in our derivation of the Hessian:

$$D_l[B_{ll}] = \tau[I - \frac{1}{m_l} \mathbf{1}\mathbf{1}^\top] \text{diag}3(\frac{1}{x_l^2}) \quad (154)$$

$$D_q[B_{ll}] = \mathbf{0} \quad (155)$$

$$D_l[B_{kl}] = \mathbf{0} \quad (156)$$

$$D_q[B_{kl}] = [I - \frac{1}{m_k} \mathbf{1}\mathbf{1}^\top] T_{klq}^k \quad (157)$$

$$D_k[\Pi_{T\Delta}(\nabla_{x_k}^{k\tau})] = [I - \frac{1}{m_k} \mathbf{1}\mathbf{1}^\top] D_k[\nabla_{x_k}^{k\tau}] \quad (158)$$

$$= [I - \frac{1}{m_k} \mathbf{1}\mathbf{1}^\top] D_k[\nabla_{x_k}^k - \tau \ln(x_k)] \quad (159)$$

$$= [I - \frac{1}{m_k} \mathbf{1}\mathbf{1}^\top] [-\tau \text{diag}(\frac{1}{x_k})] \quad (160)$$

$$= B_{kk} \quad (161)$$

$$D_l[\Pi_{T\Delta}(\nabla_{x_k}^{k\tau})] = [I - \frac{1}{m_k} \mathbf{1}\mathbf{1}^\top] D_l[\nabla_{x_k}^{k\tau}] \quad (162)$$

$$= [I - \frac{1}{m_k} \mathbf{1}\mathbf{1}^\top] D_l[\nabla_{x_k}^k - \tau \ln(x_k)] \quad (163)$$

$$= [I - \frac{1}{m_k} \mathbf{1}\mathbf{1}^\top] [H_{kl}^k] \quad (164)$$

$$= B_{kl}. \quad (165)$$

598 We can derive the diagonal blocks of the Hessian as

$$D_{ll}[\mathcal{L}(\mathbf{x})] = D_l[\nabla_{x_l} \mathcal{L}(\mathbf{x})] \quad (166)$$

$$= 2D_l[\sum_k \eta_k B_{kl}^\top \Pi_{T\Delta}(\nabla_{x_k}^{k\tau})] \quad (167)$$

$$= 2\left[\eta_l D_l[B_{ll}^\top \Pi_{T\Delta}(\nabla_{x_l}^{l\tau})] + \sum_{k \neq l} \eta_k D_l[B_{kl}^\top \Pi_{T\Delta}(\nabla_{x_k}^{k\tau})]\right] \quad (168)$$

$$= 2\left[\eta_l [D_l[B_{ll}]^\top \Pi_{T\Delta}(\nabla_{x_l}^{l\tau}) + B_{ll}^\top D_l[\Pi_{T\Delta}(\nabla_{x_l}^{l\tau})]]\right] \quad (169)$$

$$+ \sum_{k \neq l} \eta_k [D_l[B_{kl}]^\top \Pi_{T\Delta}(\nabla_{x_k}^{k\tau}) + B_{kl}^\top D_l[\Pi_{T\Delta}(\nabla_{x_k}^{k\tau})]] \quad (170)$$

$$= 2\left[\eta_l [\tau \text{diag}3(\frac{1}{x_l^2}) [I - \frac{1}{m_l} \mathbf{1}\mathbf{1}^\top] \Pi_{T\Delta}(\nabla_{x_l}^{l\tau}) + B_{ll}^\top B_{ll}] + \sum_{k \neq l} \eta_k B_{kl}^\top B_{kl}\right] \quad (171)$$

$$= 2\left[\tau \eta_l \text{diag}([\frac{1}{x_l^2}]) \odot \Pi_{T\Delta}(\nabla_{x_l}^{l\tau}) + \sum_k \eta_k B_{kl}^\top B_{kl}\right] \quad (172)$$

599 and the off-diagonal blocks as

$$D_{lq}[\mathcal{L}(\mathbf{x})] = D_q[\nabla_{x_l} \mathcal{L}(\mathbf{x})] \quad (173)$$

$$= 2D_q\left[\sum_k \eta_k B_{kl}^\top \Pi_{T\Delta}(\nabla_{x_k}^{k\tau})\right] \quad (174)$$

$$= 2\left[\eta_l D_q[B_{ll}^\top \Pi_{T\Delta}(\nabla_{x_l}^{l\tau})] + \sum_{k \neq l} \eta_k D_q[B_{kl}^\top \Pi_{T\Delta}(\nabla_{x_k}^{k\tau})]\right] \quad (175)$$

$$= 2\left[\eta_l [D_q[B_{ll}^\top]]^\top \Pi_{T\Delta}(\nabla_{x_l}^{l\tau}) + B_{ll}^\top D_q[\Pi_{T\Delta}(\nabla_{x_l}^{l\tau})]\right] \quad (176)$$

$$+ \sum_{k \neq l} \eta_k [D_q[B_{kl}^\top]]^\top \Pi_{T\Delta}(\nabla_{x_k}^{k\tau}) + B_{kl}^\top D_q[\Pi_{T\Delta}(\nabla_{x_k}^{k\tau})] \quad (177)$$

$$= 2\left[\eta_l B_{ll}^\top B_{lq} + \sum_{k \neq l} \eta_k [T_{lqk}^k [I - \frac{1}{m_k} \mathbf{1}\mathbf{1}^\top] \Pi_{T\Delta}(\nabla_{x_k}^{k\tau}) + B_{kl}^\top B_{kq}]\right] \quad (178)$$

$$= 2\left[\sum_k \eta_k B_{kl}^\top B_{kq} + \sum_{k \neq l} \eta_k T_{lqk}^k \Pi_{T\Delta}(\nabla_{x_k}^{k\tau})\right]. \quad (179)$$

600 Therefore, the Hessian can be written concisely as

$$2[\tilde{B}^\top \tilde{B} + T \Pi_{T\Delta}(\tilde{\nabla}^\tau)] \quad (180)$$

601 where  $\tilde{B}_{kl} = \sqrt{\eta_k} B_{kl}$ ,  $\Pi_{T\Delta}(\tilde{\nabla}^\tau) = [\eta_1 \Pi_{T\Delta}(\nabla_{x_1}^{1\tau}), \dots, \eta_m \Pi_{T\Delta}(\nabla_{x_m}^{m\tau})]$ , and we augment  $T$  (the  
602 3-player tensor approximation to the game,  $T_{lqk}^k$ ) so that  $T_{lll}^l = \tau \text{diag}\mathfrak{3}(\frac{1}{x_l^2})$  and otherwise 0.

603 □

## 604 E Regret Bounds

605 **Lemma 18.** [Loss Regret to Exploitability Regret] Assume exploitability of a joint strategy  $\mathbf{x}$  is upper  
606 bounded by  $f(\mathcal{L}^\tau(\mathbf{x}))$  where  $f$  is a concave function and  $\mathcal{L}^\tau$  is a loss function. Let  $\mathbf{x}_t$  be a joint  
607 strategy randomly drawn from the set of predictions made by an online learning algorithm  $\mathcal{A}$  over  $T$   
608 steps. Then the expected exploitability of  $\mathbf{x}_t$  is bounded by the average regret of  $\mathcal{A}$ :

$$\mathbb{E}[\epsilon_t] \leq f\left(\frac{1}{T} \sum_t \mathcal{L}_t\right). \quad (181)$$

*Proof.*

$$\mathbb{E}[\epsilon_t] = \mathbb{E}[f(\mathcal{L}(\mathbf{x}_t))] \quad (182)$$

$$\leq f(\mathbb{E}[\mathcal{L}(\mathbf{x}_t)]) \quad (183)$$

$$= f\left(\frac{1}{T} \sum_t \mathcal{L}(\mathbf{x}_t)\right) \quad (184)$$

609 where the second inequality follows from Jensen's inequality. □

610 **Theorem 1.** [BLiN PAC Rate] Assume  $\eta_k = \eta = 2/\hat{L}$  as defined in Lemma 2,  $\tau = \frac{1}{\ln(1/p)}$  so that all  
611 equilibria place at least  $\frac{p}{m^*}$  mass on each strategy, and a previously pulled arm is returned uniformly  
612 at random (i.e.,  $t \sim U(T)$ ). Then for any  $w > 0$ ,

$$\epsilon_t \leq w \left[ n\tau(W(1/e) + \frac{\bar{m} - 2}{e}) + 4(1 + (4c^2)^{1/3})\sqrt{nm^* \hat{L}} \left(\frac{\ln T}{T}\right)^{\frac{1}{2(d_x+2)}} \right] \quad (185)$$

613 with probability  $(1 - w^{-1})(1 - 2T^{-2})$  where  $W$  is the Lambert function ( $W(1/e) \approx 0.278$ ),

614  $m^* = \max_k m_k$ , and  $c \leq \frac{1}{4} \frac{n\bar{m}}{\hat{L}} \left(\frac{\ln(m^*)}{\ln(1/p)} + 2\right)^2$  is an upper bound on the maximum sampled value  
615 from  $\mathcal{L}^\tau$  (see Lemma 8).

616 *Proof.* Assume  $\eta_k = \eta = \frac{2}{\hat{L}}$  as defined in Lemma 2 so that  $\mathcal{L}^\tau$  is 1-Lipschitz with respect to  $\|\cdot\|_\infty$ .  
617 Also assume a previously pulled arm is returned uniformly at random. Starting with Lemma 14 and  
618 applying Corollary 9, we find

$$\mathbb{E}[\epsilon_t] \leq n\tau(W(1/e) + \frac{\bar{m} - 2}{e}) + 2\sqrt{\frac{n \max_k m_k}{\min_k \eta_k}} \sqrt{\mathcal{L}^\tau(\mathbf{x})} \quad (186)$$

$$= \frac{n}{\ln(1/p)}(W(1/e) + \frac{\bar{m} - 2}{e}) + \sqrt{2nm^*\hat{L}}\sqrt{8(1 + (4c^2)^{1/3})^2 T^{\frac{-1}{(d_z+2)}} \ln T^{\frac{1}{(d_z+2)}}} \quad (187)$$

$$= \frac{n}{\ln(1/p)}(W(1/e) + \frac{\bar{m} - 2}{e}) + 4(1 + (4c^2)^{1/3})\sqrt{nm^*\hat{L}}\left(\frac{\ln T}{T}\right)^{\frac{1}{2(d_z+2)}} \quad (188)$$

619 with probability  $1 - 2T^{-2}$  where  $W$  is the Lambert function ( $W(1/e) \approx 0.278$ ),  $m^* = \max_k m_k$ ,  
620 and  $c \leq \frac{1}{4} \frac{n\bar{m}}{\hat{L}} \left(\frac{\ln(m^*)}{\ln(1/p)} + 2\right)^2$  is an upper bound on the range of sampled values from  $\mathcal{L}^\tau$  (see  
621 Lemma 8).

622 Recall  $\hat{L} = \left[\frac{\ln(m^*)}{\ln(1/p)} + 2\right] \left[\frac{m^{*2}}{p \ln(1/p)} + n\bar{m}\right]$ . Therefore,

$$c \leq \frac{1}{4} \frac{n\bar{m}}{\hat{L}} \left(\frac{\ln(m^*)}{\ln(1/p)} + 2\right)^2 \quad (189)$$

$$= \frac{1}{4} n\bar{m} \left(\frac{\frac{\ln(m^*)}{\ln(1/p)} + 2}{\frac{m^{*2}}{p \ln(1/p)} + n\bar{m}}\right). \quad (190)$$

623 Markov's inequality then allows us to bound the pointwise exploitability of any arm returned by the  
624 algorithm as

$$\epsilon_t \leq w \left[ \frac{n}{\ln(1/p)}(W(1/e) + \frac{\bar{m} - 2}{e}) + 4(1 + (4c^2)^{1/3})\sqrt{nm^*\hat{L}}\left(\frac{\ln T}{T}\right)^{\frac{1}{2(d_z+2)}} \right] \quad (191)$$

625 with probability  $(1 - w^{-1})(1 - 2T^{-2})$  for any  $w > 0$ . □

## 626 F Complexity

### 627 F.1 Polymatrix Games

628 Interestingly, at zero temperature (where QRE = Nash),  $M$  is constant for a polymatrix game, so  
629 the rank of this matrix can be computed just once to extract information about all possible interior  
630 equilibria in the game. Furthermore, the Hessian is positive semi-definite over the entire joint strategy  
631 space, implying the loss function is convex (see Figure 1 (left) for empirical support). This indicates,  
632 by convex optimization theory, 1) all mixed Nash equilibria in polymatrix games form a convex set  
633 (i.e., they are connected) and 2) assuming mixed equilibria exist, they can be computed simply by  
634 stochastic gradient descent on  $\mathcal{L}$ . If  $M$  is rank- $n\bar{m}$ , then this interior equilibrium is unique.

635 **Complexity** Approximation of Nash equilibria in polymatrix games is known to be PPAD-hard [13].  
636 In contrast, if we restrict our class of polymatrix games to those with at least one interior Nash  
637 equilibrium, our analysis proves we can find an approximate Nash equilibrium in deterministic,  
638 polynomial time (Corollary 3). This follows directly from the fact that  $\mathcal{L}$  is convex, our decision  
639 set  $\mathcal{X} = \prod_i \mathcal{X}_i$  is convex, and convex optimization theory admits polynomial time approximation  
640 algorithms (e.g., gradient descent). We consider the assumption of the existence of an interior Nash  
641 equilibrium to be relatively mild<sup>3</sup>, so this positive complexity result is surprising.

642 Also, note that the Hessian of the loss at Nash equilibria is encoded entirely by the polymatrix  
643 approximation at the equilibrium. Therefore, approximating the Hessian of  $\mathcal{L}$  about the equilibrium  
644 (which amounts to observing near-equilibrium behavior [25]) allows one to recover this polymatrix  
645 approximation (up to constant offsets of the columns which equilibria are invariant to [27]).

<sup>3</sup>Marris et al. [27] shows 2-player, 2-action polymatrix games with interior Nash equilibria make up a non-trivial  $1/4$  of the space of possible  $2 \times 2$  games.

646 **Corollary 3** (Approximating Nash Equilibria of Polymatrix Games with Interior Equilibria). *Con-*  
647 *sider the class of polymatrix games with interior Nash equilibria. This class of games admits a fully*  
648 *polynomial time deterministic approximation scheme (FPTAS).*

649 *Proof.* Lemma 3 relates the approximation of Nash equilibria to the minimization of the loss function  
650  $\mathcal{L}(\mathbf{x})$ . By Lemma 1, this loss function attains its minimum value of zero if and only if  $\mathbf{x}$  is a  
651 Nash equilibrium. For polymatrix games, Hessian of this loss function is everywhere finite and  
652 positive definite (Lemma 17), therefore, this loss function is convex. The decision set for this  
653 minimization problem is the product space of simplices, therefore it is also convex. Given that we  
654 only consider polymatrix games with interior Nash equilibria, we know that our loss function attains  
655 a global minimum within this set. By convex optimization theory, this function can be approximately  
656 minimized in a polynomial number of steps by, for example, (projected) gradient descent. Gradient  
657 descent requires computing the gradient of the loss function at each step. From Lemma 15, we see  
658 that computing the gradient (at zero temperature) simply requires reading the polymatrix description  
659 of the game (i.e., each bi-matrix game  $H_{kl}^k$  between players), which is clearly polynomial in the size  
660 of the input (the polymatrix description). The remaining computational steps of gradient descent  
661 (e.g., projection onto simplices) are polynomial as well. In conclusion, gradient descent approximates  
662 a Nash equilibrium in polynomial number of steps (logarithmic if strongly-convex), each of which  
663 costs polynomial time, therefore the entire scheme is polynomial.  $\square$

## 664 F.2 Normal-Form Games

665 **Corollary 1.** *Consider the class of NFGs with at least one QRE( $\tau$ ) whose local polymatrix approx-*  
666 *imation indicates it is isolated (i.e.,  $M$  from equation (12) is rank- $n\bar{m}$  implies Hess  $\succ 0$  implies*  
667  *$d_z = n\bar{m}(\frac{2-2}{4}) = 0$ ). Then by Theorem 1, BLiN is a fully polynomial-time randomized approximation*  
668 *scheme (FPRAS) for QREs and is a PRAS for NEs of games in this class.*

669 *Proof.* If  $\alpha = 0$ , an  $\epsilon$ -QRE can be obtained with BLiN in a number of iterations that is polynomial in  
670 the game description length ( $nm^n$ ). The same holds for an  $\epsilon$ -NE, however, the temperature must be  
671 exponentially small to achieve a given  $\epsilon$ ; hence, we lose the *fully* qualifier. Specifically,

$$p \leq e^{-\frac{8n}{\epsilon} \left( W^{(1/\epsilon) + \frac{m-2}{e}} \right)}. \quad (192)$$

672 This, in turn, causes the Lipschitz constant  $\hat{L}$  to grow exponentially large, leading to an exponential  
673 blow up in the number of iterations required for convergence.  $\square$

### 674 F.2.1 Concrete Example

675 The end of Section 6 stated a concrete result for a 20-player, 50-action game *assuming* we are given  
676 that the game as an interior Nash equilibrium. This result requires re-deriving a rate similar to  
677 Theorem 1, but for the unregularized game.

678 For example, revisiting Corollary 2 but for zero temperature, we find  $\hat{L} = n\bar{m}$ . Let  $\eta = \frac{2}{L}$  as  
679 before. Now, consulting Table 2, we find that samples from  $\mathcal{L}$  are constrained to a range of size  
680  $c = \frac{1}{2}n\bar{m}\eta = 1$ . Applying Corollary 9 to Lemma 3, we find:

$$\epsilon_t \leq w \left[ 2\sqrt{2}(1 + 4^{1/3})n\sqrt{\bar{m}} \left( \frac{\ln T}{T} \right)^{\frac{1}{4}} \right] \quad (193)$$

681 with probability  $(1 - w^{-1})(1 - 2T^{-2}) = 0.95(1 - 2T^{-2})$ . Plugging in  $w = 20$ ,  $n = 20$ , and  $m = 50$   
682 and solving for  $T$  numerically, we find that  $T \leq 10^{28.7}$ . For such large  $T$ ,  $0.95(1 - 2T^{-2}) \approx 0.95$ .  
683 Again consulting Table 2, each call (arm pull) of BLiN costs  $2nm$ , implying a total query cost of  
684  $10^{32.0}$ . In contrast, there exist  $10^{35.2}$  scalar entries in the  $nm^n$  payoff tensor, which is a factor larger  
685 by 1000.

## 686 G Helpful Lemmas and Propositions

687 **Proposition 2.** *The matrix  $I - \frac{1}{m_k} \mathbf{1}\mathbf{1}^\top$  is a projection matrix and therefore idempotent. It is also*  
688 *symmetric, which implies it is its own square root.*

*Proof.*

$$\left[I - \frac{1}{m_k} \mathbf{1}\mathbf{1}^\top\right]^\top \left[I - \frac{1}{m_k} \mathbf{1}\mathbf{1}^\top\right] = I - \frac{2}{m_k} \mathbf{1}\mathbf{1}^\top + \frac{1}{m_k^2} \mathbf{1}(\mathbf{1}^\top \mathbf{1})\mathbf{1}^\top \quad (194)$$

$$= I - \frac{2}{m_k} \mathbf{1}\mathbf{1}^\top + \frac{1}{m_k} \mathbf{1}\mathbf{1}^\top \quad (195)$$

$$= \left[I - \frac{1}{m_k} \mathbf{1}\mathbf{1}^\top\right]. \quad (196)$$

689

□

690 **Proposition 3.** *The matrix  $I - \frac{1}{m_k} \mathbf{1}\mathbf{1}^\top$  is positive semi-definite.*

691 *Proof.* Let  $z \in \mathbb{R}^{m_k}$ . Then

$$z^\top \left[I - \frac{1}{m_k} \mathbf{1}\mathbf{1}^\top\right] z = \|z\|_2^2 - \frac{1}{m_k} \langle z, \mathbf{1} \rangle^2 \quad (197)$$

$$\geq \|z\|_2^2 - \frac{1}{m_k} \langle |z|, \mathbf{1} \rangle^2 \quad (198)$$

$$= \|z\|_2^2 - \frac{1}{m_k} \|z\|_1^2 \quad (199)$$

$$\geq \|z\|_2^2 - \|z\|_2^2 = 0 \quad \forall z \quad (200)$$

$$\implies \left[I - \frac{1}{m_k} \mathbf{1}\mathbf{1}^\top\right] \succeq 0. \quad (201)$$

692

□

693 **Proposition 4.** *The matrix  $I - \frac{1}{m_k} \mathbf{1}\mathbf{1}^\top$  has rank  $m_k - 1$  and its 1-d nullspace lies along  $\mathbf{1}_k$ .*

694 *Proof.* Note that  $\text{rank}(A+B) \leq \text{rank}(A) + \text{rank}(B)$  for matrices  $A$  and  $B$  of the same dimension.

695 Let  $A = I - \frac{1}{m_k} \mathbf{1}\mathbf{1}^\top$  and  $B = \frac{1}{m_k} \mathbf{1}\mathbf{1}^\top$  and apply  $\text{rank}(A) \geq \text{rank}(A+B) - \text{rank}(B)$ :

$$\text{rank}\left(I - \frac{1}{m_k} \mathbf{1}\mathbf{1}^\top\right) \geq \text{rank}(I) - \text{rank}\left(\frac{1}{m_k} \mathbf{1}\mathbf{1}^\top\right) = m_k - 1. \quad (202)$$

696 We can confirm the nullspace by inspection:

$$\left[I - \frac{1}{m_k} \mathbf{1}\mathbf{1}^\top\right] \mathbf{1} = \mathbf{1} - \frac{m_k}{m_k} \mathbf{1} = 0. \quad (203)$$

697

□

698 **Lemma 19.** *The product  $A\left[I_m - \frac{1}{m} \mathbf{1}_m \mathbf{1}_m^\top\right]^p B$  for any  $p > 0$  has entries whose absolute value is*  
 699 *bounded by  $\frac{m}{4}(A_{\max} - A_{\min})(B_{\max} - B_{\min})$  where  $A_{\min}, A_{\max}, B_{\min}, B_{\max}$  represent the minima*  
 700 *and maxima of the matrices respectively.*

701 *Proof.* The matrix  $\left[I - \frac{1}{m} \mathbf{1}\mathbf{1}^\top\right]$  is idempotent so we can rewrite the product for any  $p$  as

$$A\left[I - \frac{1}{m} \mathbf{1}\mathbf{1}^\top\right]\left[I - \frac{1}{m} \mathbf{1}\mathbf{1}^\top\right]B. \quad (204)$$

702 The matrix  $\left[I - \frac{1}{m} \mathbf{1}\mathbf{1}^\top\right]$  has the property that it removes the mean from every row of a matrix when  
 703 right multiplied against it, i.e.,  $A\left[I - \frac{1}{m} \mathbf{1}\mathbf{1}^\top\right]$  removes the means from the rows of  $A$ . Similarly, left  
 704 multiplying it removes the means from the column. Let  $\tilde{A}$  and  $\tilde{B}$  represent these mean-centered results  
 705 respectively. The absolute value of the  $ij$ th entry in the resulting product can then be recognized as

$$\left|\sum_k \tilde{A}_{ik} \tilde{B}_{kj}\right| = \left|\sum_k \left(A_{ik} - \frac{1}{m} \sum_{k'} A_{ik'}\right) \left(B_{kj} - \frac{1}{m} \sum_{k'} B_{k'j}\right)\right| \quad (205)$$

$$= |m \cdot \text{Corr}(A_{i,\cdot}, B_{\cdot,j}) \cdot \sigma_{A_{i,\cdot}} \cdot \sigma_{B_{\cdot,j}}| \quad (206)$$

$$\leq m \sigma_{A_{i,\cdot}} \cdot \sigma_{B_{\cdot,j}}. \quad (207)$$

706 The variance of a bounded random variable  $X$  is upper bounded by  $\text{Var}[X] \leq \frac{1}{4}(\max_X - \min_X)^2$ .

707 Hence its standard deviation is bounded by  $\text{Std}[X] \leq \frac{1}{2}(\max_X - \min_X)$ . Plugging these bounds  
 708 for  $A$  and  $B$  into equation (207) completes the claim. □

709 **H Maps from Hypercube to Simplex Product**

710 In this section, we derive properties of a map  $s$  from the unit-hypercube to the simplex product. This  
 711 map is necessary to adapt our proposed loss  $\mathcal{L}^\tau$  to the commonly assumed setting in the  $\mathcal{X}$ -armed  
 712 bandit literature [8]. We derive relevant properties of two such maps: the `softmax` and a mapping  
 713 that interprets dimensions of the hypercube as angles on a unit-sphere that is then  $\ell_1$ -normalized.

714 **Lemma 20.** *Let  $f(x) = -\mathcal{L}(s(x))$ . Then  $\|\nabla f(x)\|_\infty \leq \|J(s(x))^\top\|_\infty \|\nabla \mathcal{L}(s(x))\|_\infty$ .*

*Proof.*

$$\|\nabla f(x)\|_\infty = \|J(s(x))^\top \nabla \mathcal{L}(s(x))\|_\infty \leq \|J(s(x))^\top\|_\infty \|\nabla \mathcal{L}(s(x))\|_\infty. \quad (208)$$

715 □

716 **Lemma 21.** *The  $\infty$ -norm of the Jacobian-transpose of a transformation  $s(x)$  applied elementwise  
 717 to a product space is bounded by the  $\infty$ -norm of the Jacobian-transpose of a single transformation  
 718 from that product space, i.e.,  $\|J(s(\mathbf{x}))^\top\|_\infty \leq \max_{x_i \in \mathcal{X}_i} \|J(s(x_i))^\top\|_\infty$  for any  $i$ .*

719 *Proof.* Let  $\mathbf{x} \in \mathcal{X} = \prod_{i=1}^n \mathcal{X}_i$ ,  $\mathcal{Z} = \prod_{i=1}^n \mathcal{Z}_i$  and  $S : \mathcal{X} \rightarrow \mathcal{Z} = [s(x_1); \dots; s(x_n)]^\top$  where ;  
 720 denotes column-wise stacking,  $x_i \in \mathcal{X}_i$ . Also,  $\mathcal{X}_i = \mathcal{X}_j$  and  $\mathcal{Z}_i = \mathcal{Z}_j$  for all  $i$  and  $j$ . Then the  
 721 Jacobian of  $S(\mathbf{x})$  is

$$J(S(\mathbf{x}))^\top = \begin{bmatrix} J(s(x_1))^\top & 0 \dots & 0 \\ 0 & J(s(x_2))^\top \dots & 0 \\ 0 & 0 \dots & 0 \\ 0 & 0 \dots & J(s(x_n))^\top \end{bmatrix}. \quad (209)$$

722 The  $\infty$ -norm of this matrix is the max 1-norm of any row. This matrix is diagonal, therefore, the  
 723  $\infty$ -norm of each elementwise Jacobian-transpose represents the max 1-norm of the rows spanned  
 724 by its block. Given that the domains, ranges, and transformations  $s$  for all blocks are the same,  
 725 their  $\infty$ -norms are also the same. The max  $\infty$  over the blocks is then equal to the  $\infty$ -norm of any  
 726 individual  $J(s(x_i))^\top$ . □

727 **H.1 Hessian of Bandit Reward Function**

728 **Lemma 22.** *Let  $s(x)$  be a function that maps the unit hypercube to the simplex product (mixed  
 729 strategy space). Then the objective function  $f(x) = -\mathcal{L}(s(x))$ . The Hessian of  $-f(x)$  at an optimum  
 730  $x^*$  in direction  $\Delta$  is  $\Delta x^\top [Ds(x)^\top H_{\mathcal{L}}(x)Ds(x)] \Big|_{x^*} \Delta x$  where  $H_{\mathcal{L}}$  is the Hessian of  $\mathcal{L}$  and  $Ds(x)$  is  
 731 the Jacobian of  $s(x)$ .*

*Proof.*

$$(D^2(\mathcal{L} \circ s)(x^*))(\Delta x, \Delta x) = \Delta x^\top \left[ \sum_i \overbrace{\partial_i \mathcal{L}(s(x))}^{=0 \text{ at } x=x^*} D^2 h_i(x) \right] \Big|_{x^*} \Delta x + \Delta x^\top [Ds(x)^\top H_{\mathcal{L}}(x)Ds(x)] \Big|_{x^*} \Delta x \quad (210)$$

$$= \Delta x^\top [Ds(x)^\top H_{\mathcal{L}}(x)Ds(x)] \Big|_{x^*} \Delta x. \quad (211)$$

732 □

733 **Lemma 23.** *Let  $s(x) : \mathcal{X} \rightarrow \prod_k \Delta^{m_k-1}$  be an injective function, i.e.,  $x \neq y \implies s(x) \neq s(y)$ .  
 734 Also let  $J = J(s(x))$  be the Jacobian of  $s$  with respect to  $x$  and  $\Delta x$  be a nonzero vector in the  
 735 tangent space of  $\mathcal{X}$ . Then*

$$J\Delta x \neq \mathbf{0}. \quad (212)$$

736 *Proof.* Recall that the  $ij$ th entry of the Jacobian represents  $\frac{\partial s_i}{\partial x_j}$  so that the  $i$ th entry of  $J\Delta x$  is

$$[J\Delta x]_i = \sum_j \frac{\partial s_i}{\partial x_j} \Delta x_j = ds_i. \quad (213)$$

737 Assume  $J\Delta x = \mathbf{0}$ . This would imply a change in  $x \in \mathcal{X}$  results in no change in  $s$  ( $ds = \mathbf{0}$ ),  
738 contradicting the fact that  $s$  is injective. Therefore, we must conclude the claim that  $J\Delta x \neq \mathbf{0}$ .  $\square$

739 **Lemma 24.** *Let  $J$  be the Jacobian of the softmax operator. Then  $\|J\|_\infty \leq 2$  and  $\|J^\top\|_\infty \leq 2$ .*

740 *Proof.* Let  $S_i$  represent the  $i$ th entry of  $S = \text{softmax}(z)$  for any  $z \in \mathbb{R}^m$ . Then the 1-norm of row  $i$   
741 is upper bounded as

$$D_j S_i = S_i(\delta_{ij} - S_j) \quad (214)$$

$$\implies \sum_j |D_j S_i| = \sum_j |S_i(\delta_{ij} - S_j)| \quad (215)$$

$$\leq \sum_j |\delta_{ij} S_i| + |S_i S_j| \quad (216)$$

$$= S_i + \sum_j S_i S_j \quad (217)$$

$$= S_i + S_i \sum_j S_j \quad (218)$$

$$= 2S_i \quad (219)$$

$$\leq 2 \forall i. \quad (220)$$

742 Also, the 1-norm of row  $j$  is upper bounded similarly as

$$(221)$$

$$\sum_i |D_j S_i| = \sum_i |S_i(\delta_{ij} - S_j)| \quad (222)$$

$$\leq \sum_i |\delta_{ij} S_i| + |S_i S_j| \quad (223)$$

$$= S_j + \sum_i S_i S_j \quad (224)$$

$$= S_j + S_j \sum_i S_i \quad (225)$$

$$= 2S_j \quad (226)$$

$$\leq 2 \forall j. \quad (227)$$

743 The  $\infty$ -norm of a matrix is the maximum 1-norm of any row. Therefore,  $\|J\|_\infty$  and  $\|J^\top\|_\infty$  are both  
744 upper bounded by 2.  $\square$

745 **Lemma 25.** *Let  $J = J(s(x))$  be the Jacobian of any composition of transformations  $s = s_t \circ \dots \circ s_1$   
746 where  $s_t(z) = [z_i / \sum_j z_j]_i$ . Then  $J\Delta x$  lies in the tangent space of the simplex.*

747 *Proof.* We aim to show  $\mathbf{1}^\top J\Delta x = \mathbf{0}$  for any  $\Delta x$  and  $x$ . By chain rule, the Jacobian of  $s$  is  
748  $J = J(s) = \prod_{t'=t}^{t'=1} J(s_{t'})$ . Therefore,  $\mathbf{1}^\top J\Delta x = \mathbf{1}^\top (\prod_{t'=t}^{t'=1} J(s_{t'}))\Delta x$ . Consider the first product:

$$\mathbf{1}^\top J(s_t) = \mathbf{0} \quad (228)$$

749 by Lemma 27. Therefore  $\mathbf{1}^\top J\Delta x = \mathbf{1}^\top J(s_t) (\prod_{t'=t-1}^{t'=1} J(s_{t'}))\Delta x = \mathbf{0}^\top (\prod_{t'=t-1}^{t'=1} J(s_{t'}))\Delta x = \mathbf{0}$ .  
750 This implies  $J\Delta x$  is orthogonal to  $\mathbf{1}$  for any  $x \in \mathcal{X}$  and  $\Delta x$ , therefore  $J\Delta x$  lies in the tangent space  
751 of the simplex for any  $x \in \mathcal{X}$  and  $\Delta x$ .  $\square$

752 For spherical coordinates,  $s(x) = n(l(c(x)))$  where  $c(x) = \pi/2x$ ,  $l(\psi)$  maps angles to the unit  
753 sphere, and  $n(z) = [z_i / \sum_j z_j]_i$ .

754 **Definition 1.** Define  $l(\psi)$  as the transformation to the unit-sphere using spherical coordinates:

$$l_1(\psi) = \cos(\psi_1) \quad (229)$$

$$l_2(\psi) = \sin(\psi_1) \cos(\psi_2) \quad (230)$$

$$l_3(\psi) = \sin(\psi_1) \sin(\psi_2) \cos(\psi_3) \quad (231)$$

$$\vdots = \vdots \quad (232)$$

$$l_{m-1}(\psi) = \sin(\psi_1) \sin(\psi_2) \dots \cos(\psi_{m-1}) \quad (233)$$

$$l_m(\psi) = \sin(\psi_1) \sin(\psi_2) \dots \sin(\psi_{m-1}). \quad (234)$$

755 **Lemma 26.** Let  $J$  be the Jacobian of the transformation to the unit-sphere using spherical coordinates,  
756 i.e.  $z = l(\psi)$  where  $\|l\|^2 = 1$  and  $\psi_i \in [0, \frac{\pi}{2}]$  represents an angle for each  $i$ . Then  $\|J\|_F \leq \sqrt{m}$ .

757 *Proof.* The Jacobian of the transformation is

$$J(l) = \begin{bmatrix} -\sin(\psi_1) & 0 & \dots & 0 \\ \cos(\psi_1) \cos(\psi_2) & -\sin(\psi_1) \sin(\psi_2) & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ \cos(\psi_1) \sin(\psi_2) \dots \cos(\psi_{m-1}) & \dots & \dots & -\sin(\psi_1) \dots \sin(\psi_{m-2}) \sin(\psi_{m-1}) \\ \cos(\psi_1) \sin(\psi_2) \dots \sin(\psi_{m-1}) & \dots & \dots & \sin(\psi_1) \dots \sin(\psi_{m-2}) \cos(\psi_{m-1}) \end{bmatrix} \quad (235)$$

758 and it square is

$$J(l) = \begin{bmatrix} t_1 & 0 & \dots & 0 \\ \cos(\psi_1)^2 \cos(\psi_2)^2 & \sin(\psi_1)^2 t_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ \cos(\psi_1)^2 \sin(\psi_2)^2 \dots \cos(\psi_{m-1})^2 & \dots & \dots & \sin(\psi_1)^2 \dots \sin(\psi_{m-2})^2 t_{m-1} \\ \cos(\psi_1)^2 \sin(\psi_2)^2 \dots \sin(\psi_{m-1})^2 & \dots & \dots & \sin(\psi_1)^2 \dots \sin(\psi_{m-2})^2 t_m \end{bmatrix} \quad (236)$$

759 where

$$\delta_{im} = 1 \text{ if } i = m, 0 \text{ else} \quad (237)$$

$$t_i = \delta_{im} \cos^2(\psi_{i-1}) + (1 - \delta_{im}) \sin^2(\psi_i) \leq 1. \quad (238)$$

760 To compute the Frobenius norm, we will need the sum of the squares of all entries. We will consider  
761 the sum of each row individually using the following auxiliary variable  $R_{i,k \leq i}$  where  $\sum_j J_{ij}^2 = R_{i,1}$   
762 and apply a recursive inequality.

$$R_{i,k \leq i} = \sum_{k'=k}^{i-1} \cos^2(\psi_{k'}) \left[ \prod_{l=k, l \neq k'}^{i-1} \sin^2(\psi_l) \right] \cos^2(\psi_i) + t_i \prod_{l=k}^{i-1} \sin^2(\psi_l) \quad (240)$$

$$= \cos^2(\psi_k) \underbrace{\left[ \prod_{l=k+1}^{i-1} \sin^2(\psi_l) \right]}_{\leq 1} \cos^2(\psi_i) \quad (241)$$

$$+ \sin^2(\psi_k) \sum_{k'=k+1}^{i-1} \cos^2(\psi_{k'}) \left[ \prod_{l=k+1, l \neq k'}^{i-1} \sin^2(\psi_l) \right] \cos^2(\psi_i) \quad (242)$$

$$+ \sin^2(\psi_k) t_i \prod_{l=k+1}^{i-1} \sin^2(\psi_l) \quad (243)$$

$$\leq \cos^2(\psi_k) \quad (244)$$

$$+ \sin^2(\psi_k) \left( \sum_{k'=k+1}^{i-1} \cos^2(\psi_{k'}) \left[ \prod_{l=k+1, l \neq k'}^{i-1} \sin^2(\psi_l) \right] \cos^2(\psi_i) + t_i \prod_{l=k+1}^{i-1} \sin^2(\psi_l) \right) \quad (245)$$

$$= \cos^2(\psi_k) + \sin^2(\psi_k) R_{i,k+1}. \quad (246)$$

763 Note then that  $R_{i,k+1} \leq 1 \implies R_{i,k} \leq 1$ . We know  $R_{i,i} = t_i \leq 1$ , therefore,  $R_{i,1} \leq 1$  by applying  
 764 the inequality recursively. Finally,  $\sum_j J_{ij}^2 = R_{i,1} \leq 1$  implies the claim  $\|J\|_F^2 = \sum_i R_{i,1} \leq m$ .  $\square$

765 **Lemma 27.** *Let  $J$  be the Jacobian of  $n(z) = z/Z$  where  $Z = \sum_k z_k$ . Then  $\mathbf{1}^\top J = \mathbf{0}^\top$ .*

766 *Proof.* The  $ij$ th entry of the Jacobian of  $n(z)$  is

$$J(n)_{ij} = \frac{1}{Z^2}(-z_i + \delta_{ij}Z). \quad (247)$$

767 Therefore  $[\mathbf{1}^\top J]_j = \sum_i J(n)_{ij} = \frac{1}{Z^2}(-Z + Z) = 0$  where  $z$  is a point on the unit-sphere in the  
 768 positive orthant.  $\square$

## 769 I A2: Bounded Diameters and Well-shaped Cells

770 We assume the feasible set is a unit-hypercube of dimensionality  $d$  where cells are evenly split along  
 771 the longest edge to give  $b$  new partitions and  $x_{h,i}$  represents the center of each cell.

772 There exists a decreasing sequence  $w(h) > 0$ , such that for any depth  $h \geq 0$  and for any cell  $\mathcal{X}_{h,i}$  of  
 773 depth  $h$ , we have  $\sup_{x \in \mathcal{X}_{h,i}} \ell(x_{h,i}, x) \leq w(h)$ . Moreover, there exists  $\nu > 0$  such that for any depth  
 774  $h \geq 0$ , any cell  $\mathcal{X}_{h,i}$  contains an  $\ell$ -ball of radius  $\nu w(h)$  centered at  $x_{h,i}$ .

$\ell(x, y)$	$c$	$\gamma$	$\nu$
$\ell(x, y) = \ x - y\ _2^\alpha$	$d^{\alpha/2} \left(\frac{b}{2}\right)^\alpha$	$b^{-\alpha/d}$	$d^{-\alpha/2} b^{-2\alpha}$
$\ell(x, y) = \ x - y\ _\infty^\alpha$	$\left(\frac{b}{2}\right)^\alpha$	$b^{-\alpha/d}$	$b^{-2\alpha}$

Table 3: Bounding Constants:  $\sup_{x \in \mathcal{X}_{h,i}} \ell(x_{h,i}, x) \leq w(h) = c\gamma^h$ .

### 775 I.1 $L_2$ -Norm

776 **Lemma 28** ( $L_2$ -Norm Bounding Ball). *Let  $\ell(x, y) = \|x - y\|_2^\alpha$ . Then  $\sup_{x \in \mathcal{X}_{h,i}} \ell(x_{h,i}, x) \leq$   
 777  $w_2(h) = c\gamma^h$  where  $c = \left(\frac{db^2}{4}\right)^{\alpha/2}$  and  $\gamma = b^{-\alpha/d}$ .*

*Proof.*

$$w(0) = \left[ \sum_{i=1}^d (1/2)^2 \right]^{\alpha/2} = \left( \frac{d}{4} \right)^{\alpha/2} \quad (248)$$

$$w(1) = \left[ (1/b \cdot 1/2)^2 + \sum_{i=2}^d (1/2)^2 \right]^{\alpha/2} = \left[ (1/b^2)(1/4) + (d-1)(1/4) \right]^{\alpha/2} \quad (249)$$

$$= \left( \frac{d-1+1/b^2}{4} \right)^{\alpha/2} \quad (250)$$

$$w(d) = \left[ \sum_{i=1}^d (1/b \cdot 1/2)^2 \right]^{\alpha/2} = \left( \frac{d}{4 \cdot b^2} \right)^{\alpha/2} \quad (251)$$

$$w(h) = \left[ r(1/b)^{2(q+1)}(1/2)^2 + \sum_{i=r}^d (1/b)^{2q}(1/2)^2 \right]^{\alpha/2} \quad (252)$$

$$= \left[ (1/b)^{2q}(1/2)^2 (r(1/b)^2 + (d-r)) \right]^{\alpha/2} \quad (253)$$

$$= \left[ (1/b^2)^q (1/4) \left( d - r \left( 1 - \frac{1}{b^2} \right) \right) \right]^{\alpha/2} \quad (254)$$

$$\leq \left[ (1/b^2)^q (1/4) d \right]^{\alpha/2} \quad (255)$$

$$\leq \left[ (1/b^2)^{h/d-1} (1/4) d \right]^{\alpha/2} \quad (256)$$

$$= \left[ (1/b^2)^{h/d} (b^2/4) d \right]^{\alpha/2} \quad (257)$$

$$= \left( \frac{db^2}{4} \right)^{\alpha/2} (1/b)^{\frac{\alpha}{d}h} \quad (258)$$

$$= c\gamma^h \quad (259)$$

778 where  $q, r = \text{divmod}(h, d) \implies q \geq h/d - 1$ ,  $c = \left( \frac{db^2}{4} \right)^{\alpha/2}$ , and  $\gamma = (1/b)^{\alpha/d} = b^{-\alpha/d}$ .  $\square$

779 **Lemma 29** ( $L_2$ -Norm Inner Ball). *Let  $\ell(x, y) = \|x - y\|_2^\alpha$ . Any cell  $\mathcal{X}_{h,i}$  contains an  $\ell$ -ball of*  
 780 *radius  $\nu w_2(h)$  where  $\nu = (db^4)^{-\alpha/2}$ .*

781 *Proof.* Any cell  $\mathcal{X}_{h,i}$  contains an  $\ell$ -ball of radius equal to its shortest axis:

$$r_{\min} = \left[ (1/4)(1/b^2)^{\lceil h/d \rceil} \right]^{\alpha/2} \quad (260)$$

$$\geq \left[ (1/4)(1/b^2)^{h/d+1} \right]^{\alpha/2} \quad (261)$$

$$= \left( \frac{1}{b^2 \cdot 4} \right)^{\alpha/2} (1/b)^{\frac{\alpha}{d}h} \quad (262)$$

$$= w(h) \cdot \left( \frac{1}{db^4} \right)^{\alpha/2}. \quad (263)$$

782

$\square$

## 783 I.2 $L_\infty$ -Norm

784 **Lemma 30** ( $L_\infty$ -Norm Bounding Ball). *Let  $\ell(x, y) = \|x - y\|_\infty^\alpha$ . Then  $\sup_{x \in \mathcal{X}_{h,i}} \ell(x_{h,i}, x) \leq$*   
 785  *$w_\infty(h) = c\gamma^h$  where  $c = \left( \frac{b}{2} \right)^\alpha$  and  $\gamma = b^{-\alpha/d}$ .*

786 *Proof.* Any cell  $\mathcal{X}_{h,i}$  is contained by an  $\ell$ -ball of radius equal to its longest axis:

$$r_{\max} = \left[ (1/4)(1/b^2)^{\lfloor h/d \rfloor} \right]^{\alpha/2} \quad (264)$$

$$\leq \left[ (1/4)(1/b^2)^{h/d-1} \right]^{\alpha/2} \quad (265)$$

$$= \left( \frac{b^2}{4} \right)^{\alpha/2} (1/b)^{\frac{\alpha}{d}h} \quad (266)$$

$$= c\gamma^h \quad (267)$$

787 where  $c = \left(\frac{b^2}{4}\right)^{\alpha/2}$ , and  $\gamma = (1/b)^{\alpha/d} = b^{-\alpha/d}$ . □

788 **Lemma 31** ( $L_\infty$ -Norm Inner Ball). *Let  $\ell(x, y) = \|x - y\|_\infty^\alpha$ . Any cell  $\mathcal{X}_{h,i}$  contains an  $\ell$ -ball of*  
 789 *radius  $\nu w_\infty(h)$  where  $\nu = b^{-2\alpha}$ .*

790 *Proof.* Any cell  $\mathcal{X}_{h,i}$  contains an  $\ell$ -ball of radius equal to its shortest axis:

$$r_{\min} = [(1/4)(1/b^2)^{\lceil h/d \rceil}]^{\alpha/2} \quad (268)$$

$$\geq [(1/4)(1/b^2)^{h/d+1}]^{\alpha/2} \quad (269)$$

$$= \left(\frac{1}{b^2 \cdot 4}\right)^{\alpha/2} (1/b)^{\frac{\alpha}{d}h} \quad (270)$$

$$= w(h) \cdot \left(\frac{1}{b^4}\right)^{\alpha/2}. \quad (271)$$

791 □

### 792 I.3 Near Optimality Dimension

793 This is written in terms of a maximizing  $f$ .

794 **Assumption 1.** *Locally around each interior  $x^*$ ,  $-f(x)$  is lower bounded by  $-f(x^*) + \sigma_- \|x -$   
 795  $x^*\|^{\alpha_{hi}}$  and upper bounded by  $-f(x^*) + \ell(x, x^*)$  where  $\ell(x, x^*) = \sigma_+ \|x - x^*\|^{\alpha_{lo}}$  with  $\alpha_{lo} \leq \alpha_{hi}$   
 796 and  $\sigma_- \leq \sigma_+$  if  $\alpha_{lo} = \alpha_{hi}$ . In other words, for all  $f(x) \geq f(x^*) - \eta$ :*

$$f(x^*) - f(x) \leq \sigma_+ \|x - x^*\|^{\alpha_{lo}} \quad (272)$$

$$f(x^*) - f(x) \geq \sigma_- \|x - x^*\|^{\alpha_{hi}} \quad (273)$$

797 where we have left the precise norm unspecified for generality.

798 **Definition 2.**  $\mathcal{X}_\epsilon \stackrel{\text{def}}{=} \{x \in \mathcal{X} \mid f(x) \geq f(x^*) - \epsilon\}$

799 **Definition 3.**  $\mathcal{X}_\epsilon^{\text{lower}} \stackrel{\text{def}}{=} \{x \in \mathcal{X} \mid f(x^*) - \sigma_- \|x - x^*\|^{\alpha_{hi}} \geq f(x^*) - \epsilon\}$

800 **Corollary 4.**  $\mathcal{X}_\epsilon \subseteq \mathcal{X}_\epsilon^{\text{lower}}$ .

801 *Proof.* By Assumption 1,  $f(x^*) - \sigma_- \|x - x^*\|^{\alpha_{hi}} \geq f(x)$ . Therefore, any  $x \in \mathcal{X}$  that satisfies the  
 802 requirement for an element of  $\mathcal{X}_\epsilon$ ,  $f(x) \geq f(x^*) - \epsilon$ , will also satisfy the requirement for an element  
 803 of  $\mathcal{X}_\epsilon^{\text{lower}}$ . □

804 **Definition 4.** *The  $\psi$ -near optimality dimension is the smallest  $d' > 0$  such that there exists  $C > 0$   
 805 such that for any  $\epsilon > 0$ , the maximum number of disjoint  $\ell$ -balls of radius  $\psi\epsilon$  and center in  $\mathcal{X}_\epsilon$  is less  
 806 than  $C\epsilon^{-d'}$ .*

807 **Theorem 2.** *The  $\psi$ -near optimality dimension of  $f : x \in [0, 1]^d \rightarrow [-1, 1]$  under  $\ell$  is  $d' =$   
 808  $d\left(\frac{\alpha_{hi} - \alpha_{lo}}{\alpha_{lo}\alpha_{hi}}\right)$  with constant*

$$C = \max \left\{ 1, S_d^{-1} \left( r_\eta^{\frac{\alpha_{hi}}{\alpha_{lo}}} \sigma_-^{\left(\frac{\alpha_{hi} - \alpha_{lo}}{\alpha_{lo}\alpha_{hi}}\right)} \right)^{-d} \right\} \left( \frac{\sigma_+}{\psi \sigma_-^{\alpha_{lo}/\alpha_{hi}}} \right)^{d/\alpha_{lo}}. \quad (274)$$

809 *Proof.* First, let us define  $r_\eta = \left(\frac{\eta}{\sigma_-}\right)^{1/\alpha_{hi}}$  as in equation (285) which implies  $\eta = \sigma_- r_\eta^{\alpha_{hi}}$ . Then  
 810 apply Lemmas 32 ( $N_{\epsilon \leq \eta} \leq C_{\epsilon \leq \eta} \epsilon^{-d'}$ ) and 34 ( $N_{\epsilon \geq \eta} \leq C_{\epsilon \geq \eta}$ ) which bound the number of  $\ell$ -balls  
 811 required to pack  $\mathcal{X}_\epsilon$  when  $\epsilon$  is less than and greater than  $\eta$  respectively:

$$C_{\epsilon \leq \eta} = \left( \frac{\sigma_+}{\psi \sigma_-^{\alpha_{lo}/\alpha_{hi}}} \right)^{d/\alpha_{lo}} \quad (275)$$

$$d' = d \left( \frac{\alpha_{hi} - \alpha_{lo}}{\alpha_{lo}\alpha_{hi}} \right) \quad (276)$$

812 and

$$C_{\epsilon \geq \eta} = S_d^{-1} \left( \frac{\sigma_+}{\psi \eta} \right)^{d/\alpha_{1o}} \quad (277)$$

$$= S_d^{-1} \eta^{-d/\alpha_{1o}} \sigma_-^{d/\alpha_{hi}} \left( \frac{\sigma_+}{\psi \sigma_-^{\alpha_{1o}/\alpha_{hi}}} \right)^{d/\alpha_{1o}} \quad (278)$$

$$= S_d^{-1} \eta^{-d/\alpha_{1o}} \sigma_-^{d/\alpha_{hi}} C_{\epsilon \leq \eta} \quad (279)$$

$$= S_d^{-1} r_\eta^{-d\alpha_{hi}/\alpha_{1o}} \sigma_-^{-d/\alpha_{1o}} \sigma_-^{d/\alpha_{hi}} C_{\epsilon \leq \eta} \quad (280)$$

$$= S_d^{-1} r_\eta^{-d \frac{\alpha_{hi}}{\alpha_{1o}}} \sigma_-^{-d \left( \frac{\alpha_{hi} - \alpha_{1o}}{\alpha_{1o} \alpha_{hi}} \right)} C_{\epsilon \leq \eta} \quad (281)$$

$$= S_d^{-1} \left( r_\eta^{\frac{\alpha_{hi}}{\alpha_{1o}}} \sigma_-^{\left( \frac{\alpha_{hi} - \alpha_{1o}}{\alpha_{1o} \alpha_{hi}} \right)} \right)^{-d} C_{\epsilon \leq \eta} \quad (282)$$

813 where  $S_d$  is the volume constant for a  $d$ -sphere under the given norm.  $S_d^{-1}$  has been upper bounded  
 814 for the 2-norm in Lemma 33. For the  $\infty$ -norm,  $S_d^{-1} = 2^{-d}$ . We have written  $C_{\epsilon \geq \eta}$  in terms of  $C_{\epsilon \leq \eta}$   
 815 to clarify which is larger.

816 Therefore,

$$C = \max \left\{ 1, S_d^{-1} \left( r_\eta^{\frac{\alpha_{hi}}{\alpha_{1o}}} \sigma_-^{\left( \frac{\alpha_{hi} - \alpha_{1o}}{\alpha_{1o} \alpha_{hi}} \right)} \right)^{-d} \right\} C_{\epsilon \leq \eta} \quad (283)$$

$$= \max \left\{ 1, S_d^{-1} \left( r_\eta^{\frac{\alpha_{hi}}{\alpha_{1o}}} \sigma_-^{\left( \frac{\alpha_{hi} - \alpha_{1o}}{\alpha_{1o} \alpha_{hi}} \right)} \right)^{-d} \right\} \left( \frac{\sigma_+}{\psi \sigma_-^{\alpha_{1o}/\alpha_{hi}}} \right)^{d/\alpha_{1o}}. \quad (284)$$

817 Intuitively, if the radius for which the polynomial bounds hold ( $r_\eta$ ) is large and the minimum  
 818 curvature constant  $\sigma_-$  is also large, then the bound  $C_{\epsilon \leq \eta}$  holds for large deviations from optimality  $\eta$ .  
 819 The number of  $\eta$ -radius  $\ell$ -balls required to cover the remaining space,  $C_{\epsilon \geq \eta}$ , will be comparatively  
 820 small.  $\square$

821 **Corollary 5** (Zooming Dimension). *The zooming dimension of  $f : x \in [0, 1]^d \rightarrow [-1, 1]$  under*  
 822  *$\ell(x, y) = \|x - y\|_\infty$  is  $d_z = d \left( \frac{\alpha_{hi} - \alpha_{1o}}{\alpha_{1o} \alpha_{hi}} \right)$ .*

823 *Proof.* Mapping the definition of zooming dimension onto  $\psi$ -near optimality, we find  $\psi \epsilon = r/2$  and  
 824  $\epsilon = 16r$ . Then we can infer  $\psi = 1/32$ . This result only effects the constant  $C_z$ , not the zooming  
 825 dimension.

826 If  $\epsilon = 8(1 + \sqrt{c_1/c_2})r_m$ , then  $\psi = \frac{1}{16(1 + \sqrt{c_1/c_2})}$ .  $\square$

827 **Lemma 32** ( $N_{\epsilon \leq \eta} \leq C_{\epsilon \leq \eta} \epsilon^{-d'}$ ). *The number of disjoint  $\ell$ -balls that can pack into a set  $\mathcal{X}_{\epsilon \leq \eta}$ ,  $N_{\epsilon \leq \eta}$ ,*  
 828 *is upper bounded by  $C_{\epsilon \leq \eta} \epsilon^{-d'}$  where  $C_{\epsilon \leq \eta} = \left( \frac{\sigma_+}{\psi \sigma_-^{\alpha_{1o}/\alpha_{hi}}} \right)^{d/\alpha_{1o}}$  and  $d' = d \left( \frac{\alpha_{hi} - \alpha_{1o}}{\alpha_{1o} \alpha_{hi}} \right)$  and  $S_d$  is the*  
 829 *volume constant for a  $d$ -sphere under the given norm  $\|\cdot\|$ .*

830 *Proof.* The number of disjoint  $\ell$ -balls of radius  $\psi \epsilon$  and center in  $\mathcal{X}_{\epsilon \leq \eta}$  can be upper bounded as  
 831 follows.

832 Rewrite  $\mathcal{X}_\epsilon^{lower}$  by rearranging terms as

$$\mathcal{X}_\epsilon^{lower} = \{x \in \mathcal{X} \mid \|x - x^*\| \leq \left( \frac{\epsilon}{\sigma_-} \right)^{1/\alpha_{hi}} \stackrel{\text{def}}{=} r_\epsilon\} \quad (285)$$

833 and recall that from Corollary 4 that  $\mathcal{X}_\epsilon \subseteq \mathcal{X}_\epsilon^{lower}$ . Furthermore, an  $\ell$ -ball of radius  $\psi \epsilon$  implies

$$\sigma_+ \|x - y\|^{\alpha_{1o}} \leq \psi \epsilon \implies \|x - y\| \leq \left( \frac{\psi \epsilon}{\sigma_+} \right)^{1/\alpha_{1o}} \stackrel{\text{def}}{=} r_\ell. \quad (286)$$

834 The number of disjoint  $\ell$ -balls that can pack into a set  $\mathcal{X}_\epsilon$ ,  $N_{\epsilon \leq \eta}$ , is upper bounded by the ratio of the  
 835 volumes of the two sets:

$$N_{\epsilon \leq \eta} \leq \frac{\text{Vol}(\mathcal{X}_\epsilon)}{\text{Vol}(\mathcal{B}_\ell)} \quad (287)$$

$$\leq \frac{\text{Vol}(\mathcal{X}_\epsilon^{\text{lower}})}{\text{Vol}(\mathcal{B}_\ell)} \quad (288)$$

$$= \frac{S_d r_\epsilon^d}{S_d r_\ell^d} \quad (289)$$

$$\leq \frac{\left(\frac{\epsilon}{\sigma_-}\right)^{d/\alpha_{hi}}}{\left(\frac{\psi \epsilon}{\sigma_+}\right)^{d/\alpha_{lo}}} \quad (290)$$

$$= \left(\frac{\sigma_+^{1/\alpha_{lo}} \psi^{-1/\alpha_{lo}}}{\sigma_-^{1/\alpha_{hi}}}\right)^d \epsilon^{d(1/\alpha_{hi} - 1/\alpha_{lo})} \quad (291)$$

$$= \left(\frac{\sigma_+}{\psi \sigma_-^{\alpha_{lo}/\alpha_{hi}}}\right)^{d/\alpha_{lo}} \epsilon^{-d(\frac{\alpha_{hi} - \alpha_{lo}}{\alpha_{lo} \alpha_{hi}})} \quad (292)$$

$$= C_{\epsilon \leq \eta} \epsilon^{-d'} \quad (293)$$

836 where  $C_{\epsilon \leq \eta} = \left(\frac{\sigma_+}{\psi \sigma_-^{\alpha_{lo}/\alpha_{hi}}}\right)^{d/\alpha_{lo}}$  and  $d' = d(\frac{\alpha_{hi} - \alpha_{lo}}{\alpha_{lo} \alpha_{hi}})$  and  $S_d$  is the volume constant for a  $d$ -sphere

837 under the given norm  $\|\cdot\|$ , e.g.,  $S_d = 2^d$  for  $\|\cdot\|_\infty$ .  $\square$

838 **Corollary 6.** If  $\alpha_{lo} = \alpha_{hi} = \alpha$ ,

$$N_{\epsilon \leq \eta} \leq \left(\frac{\kappa}{\psi}\right)^{d/\alpha}. \quad (294)$$

839 In other words,  $N_{\epsilon \leq \eta} \leq C_{\epsilon \leq \eta} \epsilon^{-d'}$  where  $C_{\epsilon \leq \eta} = \left(\frac{\sigma_+}{\psi \sigma_-}\right)^{d/\alpha}$  and  $d' = 0$ .

840 **Corollary 7.** If Assumption 1 is given in terms of the 2-norm, these can be translated to bounds in  
 841 terms of the  $\infty$ -norm resulting in the same  $\psi$ -near optimality dimension but incurring an additional  
 842 exponential factor in the constant  $C_{\epsilon \leq \eta}^{(\infty)} \leftarrow C_{\epsilon \leq \eta}^{(2)} d^{d/2}$ .

843 *Proof.* Recall that  $\|\cdot\|_\infty \leq \|\cdot\|_2 \leq \sqrt{d} \|\cdot\|_\infty$ , therefore

$$f(x^*) - f(x) \leq \sigma_{+2} \|x - x^*\|_2^{\alpha_{lo}} \leq \sigma_{+\infty} \|x - x^*\|_\infty^{\alpha_{lo}} \quad (295)$$

$$f(x^*) - f(x) \geq \sigma_{-2} \|x - x^*\|_2^{\alpha_{hi}} \geq \sigma_{-\infty} \|x - x^*\|_\infty^{\alpha_{hi}} \quad (296)$$

844 where  $\sigma_{+\infty} = \sigma_{+2} d^{\alpha_{lo}/2}$  and  $\sigma_{-\infty} = \sigma_{-2}$ . Then

$$C_{\epsilon \leq \eta}^{(\infty)} = \left(\frac{\sigma_{+2} d^{\alpha_{lo}/2}}{\psi \sigma_-^{\alpha_{lo}/\alpha_{hi}}}\right)^{d/\alpha_{lo}} = \left(\frac{\sigma_{+2}}{\psi \sigma_-^{\alpha_{lo}/\alpha_{hi}}}\right)^{d/\alpha_{lo}} d^{d/2} = C_{\epsilon \leq \eta}^{(2)} d^{d/2}. \quad (297)$$

845  $\square$

846 Recall, these results apply when  $f(x) \geq f(x^*) - \eta$ , i.e., when  $\epsilon \leq \eta$ . Otherwise, we can upper bound  
 847 the number of  $\ell$ -balls by considering the entire set  $\mathcal{X}$  which has volume 1. First, we will bound the  
 848 constant associated with the volume of a  $d$ -sphere.

849 **Lemma 33.** The volume of a  $d$ -sphere with radius  $r$  and  $d$  even is given by  $S_d r^d$  where  $S_d^{-1} \leq$   
 850  $\sqrt{2\pi d} \left(\frac{d}{2\pi e}\right)^{d/2}$ .

851 *Proof.* First, we recall Stirling's bounds on the factorial:  $\sqrt{2\pi n} \left(\frac{n}{e}\right)^n e^{\frac{1}{12n+1}} < n! < \sqrt{2\pi n} \left(\frac{n}{e}\right)^n e^{\frac{1}{12n}}$ .  
 852 This will be useful for bounding the Gamma function:  $\Gamma(d) = (d-1)!$  for even  $d$ .

853 Given  $d$  is even, we start with the exact formula for  $S_d$ :

$$S_d^{-1} = \frac{\Gamma(d/2 + 1)}{\pi^{d/2}} \quad (298)$$

$$= \frac{(d/2)!}{\pi^{d/2}} \quad (299)$$

$$< \frac{\sqrt{2\pi(d/2)} \left(\frac{d/2}{e}\right)^{d/2} e^{\frac{1}{12(d/2)}}}{\pi^{d/2}} \quad (300)$$

$$= \frac{\pi^{1/2} d^{1/2} \left(\frac{d}{2e}\right)^{d/2} e^{\frac{1}{6d}}}{\pi^{d/2}} \quad (301)$$

$$= \frac{\pi^{1/2} d^{(d+1)/2} e^{\frac{1}{6d}}}{(2\pi e)^{d/2}} \quad (302)$$

$$\leq \sqrt{2\pi d} \left(\frac{d}{2\pi e}\right)^{d/2}. \quad (303)$$

854

□

855 **Lemma 34** ( $N_{\epsilon \geq \eta} \leq C_{\epsilon \geq \eta}$ ). *The number of disjoint  $\ell$ -balls that can pack into a set  $\mathcal{X}_{\epsilon \geq \eta}$ ,  $N_{\epsilon \geq \eta}$ , is*  
 856 *upper bounded by  $C_{\epsilon \geq \eta}$  where  $C_{\epsilon \geq \eta} = S_d^{-1} \left(\frac{\sigma_+}{\psi \eta}\right)^{d/\alpha_{1o}}$  and  $S_d$  is the volume constant for a  $d$ -sphere*  
 857 *under a given norm.*

858 *Proof.* We can upper bound the number of  $\ell$ -balls needed to pack the entire space as follows:

$$N_{\epsilon \geq \eta} \leq \frac{Vol(\mathcal{X})}{Vol(\mathcal{B}_\ell)} \quad (304)$$

$$= \frac{1}{S_d r_\ell^d} \quad (305)$$

$$\leq \frac{1}{S_d \left(\frac{\psi \eta}{\sigma_+}\right)^{d/\alpha_{1o}}} \quad (306)$$

$$= S_d^{-1} \left(\frac{\sigma_+}{\psi \eta}\right)^{d/\alpha_{1o}} \quad (307)$$

$$= C_{\epsilon \geq \eta} \quad (308)$$

859 where  $r_\ell$  was defined in equation (286).  $S_d^{-1}$  has been upper bounded for the 2-norm in Lemma 33.  
 860 For the  $\infty$ -norm,  $S_d^{-1} = 2^{-d}$ . □

861 **Corollary 8.** *If Assumption 1 is given in terms of the 2-norm, these can be translated to bounds in*  
 862 *terms of the  $\infty$ -norm resulting in the same  $\psi$ -near optimality dimension but incurring an additional*  
 863 *exponential factor in the constant  $C_{\epsilon \geq \eta}^{(\infty)} = \left(\frac{\sigma_+ 2}{2\eta^{1/\alpha_{1o}}}\right)^d C_{\epsilon \leq \eta}^{(\infty)} = \left(\frac{\sigma_+ 2}{2\eta^{1/\alpha_{1o}}}\right)^d C_{\epsilon \leq \eta}^{(2)} d^{d/2}$ .*

*Proof.*

$$C_{\epsilon \geq \eta}^{(\infty)} = 2^{-d} \left(\frac{\sigma_+ 2}{\psi \eta}\right)^{d/\alpha_{1o}} d^{d/2} \quad (309)$$

$$= 2^{-d} \eta^{-d/\alpha_{1o}} \sigma_{-2}^{d/\alpha_{hi}} \left(\frac{\sigma_+ 2}{\psi \sigma_{-2}^{\alpha_{1o}/\alpha_{hi}}}\right)^{d/\alpha_{1o}} d^{d/2} \quad (310)$$

$$= 2^{-d} \eta^{-d/\alpha_{1o}} \sigma_{-2}^{d/\alpha_{hi}} C_{\epsilon \leq \eta}^{(\infty)} \quad (311)$$

$$= \left(\frac{\sigma_+ 2}{2\eta^{1/\alpha_{1o}}}\right)^d C_{\epsilon \leq \eta}^{(\infty)}. \quad (312)$$

864

□

865 If we further assume  $\alpha = \alpha_{lo} = \alpha_{hi} = 2$ , then we can bound the number of  $\ell$ -balls required with a  
 866 constant, independent of  $\epsilon$ , as

$$C = \max\{N_{\epsilon \leq \eta}, N_{\epsilon \geq \eta}\} \quad (313)$$

$$= \max\left\{\left(\frac{\kappa}{\psi}\right)^{d/2}, \sqrt{2\pi d} \left(\frac{d\sigma_{max}}{2\pi e\psi\eta}\right)^{d/2}\right\} \quad (314)$$

$$= \beta^{d/2} \psi^{-d/2} d^{\xi/2(d+1)} \quad (315)$$

867 where  $\beta = \kappa$ ,  $\xi = 0$  for  $N_{\epsilon \leq \eta}$  and  $\beta = \frac{\sigma_{max}(2\pi)^{1/d}}{2\pi e\eta} < \frac{2\sigma_{max}}{\pi e\eta} = \frac{2\kappa}{\pi e r_\eta^2} < \frac{\kappa}{(2r_\eta)^2}$  for  $d \geq 2$ ,  $\xi = 1$  for  
 868  $N_{\epsilon \geq \eta}$ .  $N_{\epsilon \geq \eta}$  dominates for large  $d$ . The cross over occurs at

$$\left(\frac{\kappa}{\psi}\right)^{d/2} = \sqrt{2\pi d} \left(\frac{d\sigma_{max}}{2\pi e\psi\eta}\right)^{d/2} \quad (316)$$

$$\implies \frac{\kappa}{\psi} = (2\pi d)^{1/d} \left(\frac{d\sigma_{max}}{2\pi e\psi\eta}\right) \quad (317)$$

$$\implies r_\eta^2 = \frac{\eta}{\sigma_-} = (2\pi d)^{1/d} \left(\frac{d}{2\pi e}\right) = z(d). \quad (318)$$

869 where  $r_\eta$  was defined in equation (285). As  $d$  grows and  $z(d)$  exceeds  $r_\eta^2$ ,  $N_{\epsilon \geq \eta}$  begins to dominate,  
 870 therefore we will upper bound  $C$  as

$$C \leq \left(\frac{\kappa}{\psi(2r_\eta)^2}\right)^{d/2} d^{\frac{1}{2}(d+1)}. \quad (319)$$

Locality	$C$
(*) $r_\eta^2 \leq z(d)$	$N_{\epsilon \geq \eta} \leq \left(\frac{\kappa}{\psi(2r_\eta)^2}\right)^{d/2} d^{\frac{1}{2}(d+1)} = \left(\frac{3\kappa b^2}{(2r_\eta)^2}\right)^{d/2} d^{d+\frac{1}{2}}$
$r_\eta^2 > z(d)$	$N_{\epsilon \leq \eta} \leq \left(\frac{\kappa}{\psi}\right)^{d/2} = (3\kappa b^2)^{d/2} d^{d/2}$

Table 4: Bounding Constants for  $\ell(x, y) = \|x - y\|_2^2$ ,  $\psi = \nu/3 = (3b^2d)^{-1}$  and  $z(d) = (2\pi d)^{1/d} \left(\frac{d}{2\pi e}\right)$  with smoothness radius  $r_\eta$  and  $\psi = \nu/3$ . (\*) indicates the case that is more likely for difficult problems.

For convenience, we repeat the other relevant constants in Table 5.

$\ell(x, y)$	$c$	$\gamma$	$\nu$
$\ell(x, y) = \ x - y\ _2^2$	$d\left(\frac{b}{2}\right)^2$	$b^{-2/d}$	$d^{-1}b^{-2}$

Table 5: Bounding Constants

871

## 872 J D-BLiN

873 The regret bound for Doubling BLiN [14] was originally proved assuming a standard normal distribu-  
 874 tion, however, the authors state their proof can be easily adapted to any sub-Gaussian distribution,  
 875 which includes bounded random variables. This matches our setting with bounded payoffs, so we  
 876 repeat their analysis here for that setting.

877 **Definition 5** (Global Arm Accuracy).  $\mathcal{E} \stackrel{\text{def}}{=} \left\{ |\mu(x) - \hat{\mu}_m(C)| \leq r_m + \sqrt{c_1 \frac{\ln T}{n_m}}, \forall 1 \leq m \leq \right.$   
 878  $\left. B_{stop} - 1, \forall C \in \mathcal{A}_m, \forall x \in C \right\}$ .

879 Define:  $n_m = c_2 \frac{\ln T}{r_m^2} \implies r_m = \sqrt{c_2 \frac{\ln T}{n_m}}$ .

880 **Definition 6 (Elimination Rule).** Eliminate  $C \in \mathcal{A}_m$  if  $\hat{\mu}_m^{\max} - \hat{\mu}_m(C) \geq 2(1 + \sqrt{c_1/c_2})r_m =$   
881  $2(\sqrt{c_2} + \sqrt{c_1})\sqrt{\frac{\ln T}{n_m}}$  where  $\hat{\mu}_m^{\max} \stackrel{\text{def}}{=} \max_{C \in \mathcal{A}_m} \hat{\mu}_m(C)$ .

882 **Lemma 35.**  $Pr[\mathcal{E}] \geq 1 - 2T^{-2(c_1/c^2-1)}$ .

883 *Proof.* Assume  $y_{C,i} \in [a, b]$  with  $c = b - a$  and  $\hat{\mu}(C) = \frac{1}{n_m} \sum_{i=1}^{n_m} y_{C,i}$ . Applying a Hoeffding  
884 inequality gives

$$Pr \left[ |\hat{\mu}(C) - \mathbb{E}[\hat{\mu}(C)]| \geq \sqrt{c_1 \frac{\ln T}{n_m}} \right] \leq 2e^{-2c_1 \ln T / c^2} \quad (320)$$

$$= 2(e^{\ln T})^{-2c_1/c^2} \quad (321)$$

$$= 2T^{-2c_1/c^2} \quad \forall C. \quad (322)$$

885 By Lipschitzness of  $\mu$  we also have

$$|\mathbb{E}[\hat{\mu}(C)] - \mu(x)| \leq r_m, \quad \forall x \in C. \quad (323)$$

886 Then consider

$$\sup_{x \in C} |\mu(x) - \hat{\mu}(C)| = \sup_{x \in C} |\mu(x) - \mathbb{E}[\hat{\mu}(C)] + \mathbb{E}[\hat{\mu}(C)] - \hat{\mu}(C)| \quad (324)$$

$$\leq \sup_{x \in C} \left( |\mu(x) - \mathbb{E}[\hat{\mu}(C)]| + |\mathbb{E}[\hat{\mu}(C)] - \hat{\mu}(C)| \right) \quad (325)$$

$$= \sup_{x \in C} |\mu(x) - \mathbb{E}[\hat{\mu}(C)]| + |\mathbb{E}[\hat{\mu}(C)] - \hat{\mu}(C)| \quad (326)$$

$$\leq \sqrt{c_1 \frac{\ln T}{n_m}} + r_m \quad (327)$$

887 with probability  $1 - 2T^{-2c_1/c^2}$ . The first inequality follows by triangle inequality and the second  
888 follows from equation (323) and considering the complement of equation (322).

889 The complement of this result occurs with probability

$$Pr \left[ \sup_{x \in C} |\mu(x) - \hat{\mu}(C)| \geq r_m + \sqrt{c_1 \frac{\ln T}{n_m}} \right] \leq 2T^{-2c_1/c^2}. \quad (328)$$

890 At least 1 arm is played in each cube  $C \in \mathcal{A}_m$  for  $1 \leq m \leq B_{stop} - 1$ , therefore,  $|\mathcal{A}_m| \leq T$  must be  
891 true given the exit condition of the algorithm. In addition, assume  $B_{stop} \leq T$  ( $B_{stop}$  will be defined  
892 such that this is true). Then a union bound over all  $T^2$  events gives

$$Pr \left[ \exists m \in [1, B_{stop} - 1], C \in \mathcal{A}_m \text{ s.t. } \sup_{x \in C} |\mu(x) - \hat{\mu}(C)| \geq r_m + \sqrt{c_1 \frac{\ln T}{n_m}} \right] \quad (329)$$

$$\leq \sum_{m=1}^{B_{stop}-1} \sum_{C \in \mathcal{A}_m} Pr \left[ \sup_{x \in C} |\mu(x) - \hat{\mu}(C)| \geq r_m + \sqrt{c_1 \frac{\ln T}{n_m}} \right] \quad (330)$$

$$\leq \sum_{m=1}^{B_{stop}-1} \sum_{C \in \mathcal{A}_m} 2T^{-2c_1/c^2} \quad (331)$$

$$\leq 2T^{-2c_1/c^2} T^2. \quad (332)$$

893 Taking the complement of this event and noting that  $\sup_{x \in C} |\mu(x) - \hat{\mu}(C)| \leq r_m + \sqrt{c_1 \frac{\ln T}{n_m}} \implies$

894  $|\mu(x) - \hat{\mu}(C)| \leq r_m + \sqrt{c_1 \frac{\ln T}{n_m}} \quad \forall x \in C$  gives the desired result.  $\square$

895 **Lemma 36 (Optimal Arm Survives).** Under event  $\mathcal{E}$ , the optimal arm  $x^* = \arg \max \mu(x)$  is not  
896 eliminated after the first  $B_{stop} - 1$  batches.

897 *Proof.* Let  $C_m^*$  denote the cube containing  $x^*$  in  $\mathcal{A}_m$ . Under event  $\mathcal{E}$ , for any cube  $C \in \mathcal{A}_m$  and  
 898  $x \in C$ , the following relation shows that  $C_m^*$  avoids the elimination rule in round  $m$ :

$$\hat{\mu}(C) - \hat{\mu}(C_m^*) \leq \left( \mu(x) + r_m + \sqrt{c_1 \frac{\ln T}{n_m}} \right) + \left( -\mu(x^*) + r_m + \sqrt{c_1 \frac{\ln T}{n_m}} \right) \quad (333)$$

$$= \underbrace{(\mu(x) - \mu(x^*))}_{\leq 0} + 2r_m + 2\sqrt{c_1 \frac{\ln T}{n_m}} \quad (334)$$

$$\leq 2\sqrt{c_2 \frac{\ln T}{n_m}} + 2\sqrt{c_1 \frac{\ln T}{n_m}} \quad (335)$$

$$= 2(\sqrt{c_1} + \sqrt{c_2})\sqrt{\frac{\ln T}{n_m}} \quad (336)$$

899 where the first inequality follows from applying Lemma 35 to upper bound  $\hat{\mu}(C)$  and  $\hat{\mu}(C_m^*)$  indi-  
 900 vidually. The remaining steps use the optimality of  $x^*$ , the definition of  $r_m$ , and the elimination  
 901 rule.  $\square$

902 **Lemma 37.** *Under event  $\mathcal{E}$ , for any  $1 \leq m \leq B_{\text{stop}}$ , any  $C \in \mathcal{A}_m$  and any  $x \in C$ ,  $\Delta_x$  satisfies*

$$\Delta_x \leq 4(1 + \sqrt{c_1/c_2})r_{m-1} \quad (337)$$

903 *Proof.* For  $m = 1$ , recall that  $r_m$  is the side length of a cube  $C \in \mathcal{A}_m$ , therefore,  $\Delta_x \leq r_{m-1} \leq$   
 904  $4(1 + \sqrt{c_1/c_2})r_{m-1}$  holds directly from the Lipschitzness of  $\mu$ .

905 For  $m > 1$ , let  $C_{m-1}^* \in \mathcal{A}_{m-1}$  be the cube containing  $x^*$ . From Lemma 36, this cube has not been  
 906 eliminated under event  $\mathcal{E}$ . For any cube  $C \in \mathcal{A}_m$  and  $x \in C$ , it is clear that  $x$  is also in the parent of  
 907  $C$ , denoted  $C_{\text{par}}$  ( $x \in C \subset C_{\text{par}}$ ). Then for any  $x \in C$ , it holds that

$$\Delta_x = \mu(x^*) - \mu(x) \leq \left( \hat{\mu}_{m-1}(C_{m-1}^*) + r_{m-1} + \sqrt{c_1 \frac{\ln T}{n_{m-1}}} \right) + \left( -\hat{\mu}_{m-1}(C_{\text{par}}) + r_{m-1} + \sqrt{c_1 \frac{\ln T}{n_{m-1}}} \right) \quad (338)$$

$$= (\hat{\mu}_{m-1}(C_{m-1}^*) - \hat{\mu}_{m-1}(C_{\text{par}})) + 2(\sqrt{c_1} + \sqrt{c_2})\sqrt{\frac{\ln T}{n_{m-1}}} \quad (339)$$

$$\leq (\hat{\mu}_{m-1}^{\max} - \hat{\mu}_{m-1}(C_{\text{par}})) + 2(\sqrt{c_1} + \sqrt{c_2})\sqrt{\frac{\ln T}{n_{m-1}}} \quad (340)$$

$$\leq 2(\sqrt{c_1} + \sqrt{c_2})\sqrt{\frac{\ln T}{n_{m-1}}} + 2(\sqrt{c_1} + \sqrt{c_2})\sqrt{\frac{\ln T}{n_{m-1}}} \quad (341)$$

$$= 4(\sqrt{c_1} + \sqrt{c_2})\sqrt{\frac{\ln T}{n_{m-1}}} \quad (342)$$

$$= 4(1 + \sqrt{c_1/c_2})r_{m-1} \quad (343)$$

908 where we have applied Lemma 35 similarly as in Lemma 36 and also used the definition of  $r_{m-1}$ .  
 909 The last two inequalities use the fact that  $\hat{\mu}_{m-1}(C_{m-1}^*) \leq \hat{\mu}_{m-1}^{\max}$  and  $C_{\text{par}}$  was not eliminated.  $\square$

910 **Theorem 3.** *With probability exceeding  $1 - 2T^{-2(c_1/c^2-1)}$ , the  $T$ -step total regret  $R(T)$  of BLiN  
 911 with Doubling Edge-length Sequence (D-BLiN) [14] satisfies*

$$R(T) \leq 8(1 + \sqrt{c_1/c_2})(2c_2 + 1) \ln(T) \frac{1}{d_z+2} T^{\frac{d_z+1}{d_z+2}} \quad (344)$$

912 where  $d_z$  is the zooming dimension of the problem instance. In addition, D-BLiN only needs no more  
 913 than  $B^* = \frac{\log 2(T) - \log 2(\ln(T))}{d_z+2} + 2$  rounds of communications to achieve this regret rate.

914 *Proof.* Since  $r_m = \frac{r_{m-1}}{2} \implies r_{m-1} = 2r_m$  for the Doubling Edge-length Sequence, Lemma 37  
915 implies that every cube  $C \in \mathcal{A}_m$  is a subset of  $S(8(1 + \sqrt{c_1/c_2})r_m)$ . Thus from the definition of  
916 zooming number (Corollary 5 with appropriate condition), we have

$$|\mathcal{A}_m| \leq N_{r_m} \leq C_z r_m^{-d_z}. \quad (345)$$

917 Fix any positive number  $B$ . Also by Lemma 37, we know that any arm played after batch  $B$  incurs  
918 a regret bounded by  $8(1 + \sqrt{c_1/c_2})r_B$ , since the cubes played after batch  $B$  have edge length no  
919 larger than  $r_B$ . Then the total regret that occurs after batch  $B$  is bounded by  $8(1 + \sqrt{c_1/c_2})r_B T$   
920 (where  $T$  is an upper bound on the number of arms).

921 Thus the regret can be bounded as

$$R(T) \leq \sum_{m=1}^B \sum_{C \in \mathcal{A}_m} \sum_{i=1}^{n_m} \Delta_{x_{C,i}} + 8(1 + \sqrt{c_1/c_2})r_B T \quad (346)$$

922 where the first term bounds the regret in the first  $B$  batches of D-BLiN, and the second term bounds  
923 the regret after the first  $B$  batches. If the algorithm stops at batch  $\tilde{B} < B$ , we define  $\mathcal{A}_m =$  for any  
924  $\tilde{B} < m \leq B$  and inequality equation (346) still holds.

925 By Lemma 37, we have  $\Delta_{x_{C,i}} \leq 8(1 + \sqrt{c_1/c_2})r_m$  for all  $C \in \mathcal{A}_m$ . We can thus bound equa-  
926 tion (346) by

$$R(T) \leq \sum_{m=1}^B |\mathcal{A}_m| \cdot n_m \cdot 8(1 + \sqrt{c_1/c_2})r_m + 8(1 + \sqrt{c_1/c_2})r_B T \quad (347)$$

$$\leq \sum_{m=1}^B N_{r_m} \cdot n_m \cdot 8(1 + \sqrt{c_1/c_2})r_m + 8(1 + \sqrt{c_1/c_2})r_B T \quad (348)$$

$$= \sum_{m=1}^B N_{r_m} \cdot c_2 \frac{\ln T}{r_m^2} \cdot 8(1 + \sqrt{c_1/c_2})r_m + 8(1 + \sqrt{c_1/c_2})r_B T \quad (349)$$

$$= \sum_{m=1}^B N_{r_m} \cdot \frac{\ln T}{r_m} \cdot 8c_2(1 + \sqrt{c_1/c_2}) + 8(1 + \sqrt{c_1/c_2})r_B T \quad (350)$$

927 where equation (348) uses equation (345), and equation (349) uses equality  $n_m = c_2 \frac{\ln T}{r_m^2}$ . Since  
928  $r_m = 2^{-m+1}$  and  $N_{r_m} \leq r_m^{-d_z} \leq 2^{(m-1)d_z}$ , we have

$$R(T) \leq \sum_{m=1}^B 2^{(m-1)d_z} \cdot \frac{\ln T}{2^{-m+1}} \cdot 8c_2(1 + \sqrt{c_1/c_2}) + 8(1 + \sqrt{c_1/c_2})2^{-B+1}T \quad (351)$$

$$= 8(1 + \sqrt{c_1/c_2}) \left[ c_2 \ln T \sum_{m=1}^B 2^{(m-1)(d_z+1)} + 2^{-B+1}T \right]. \quad (352)$$

929 Continuing we find

$$R(T) \leq 8(1 + \sqrt{c_1/c_2}) \left[ c_2 \ln T \sum_{m=1}^B 2^{(m-1)(d_z+1)} + 2^{-B+1}T \right] \quad (353)$$

$$= 8(1 + \sqrt{c_1/c_2}) \left[ c_2 \ln T \sum_{m=1}^B (2^{d_z+1})^{m-1} + 2^{-B+1}T \right] \quad (354)$$

$$= 8(1 + \sqrt{c_1/c_2}) \left[ c_2 \ln T \sum_{m=0}^{B-1} (2^{d_z+1})^m + 2^{-B+1}T \right] \quad (355)$$

$$= 8(1 + \sqrt{c_1/c_2}) \left[ c_2 \ln T \left( \frac{2^{B(d_z+1)} - 1}{2^{d_z+1} - 1} \right) + 2^{-B+1}T \right] \text{ via geometric series} \quad (356)$$

$$\leq 8(1 + \sqrt{c_1/c_2}) \left[ c_2 \ln T \left( \frac{2^{B(d_z+1)}}{2^{d_z+1} - 1} \right) + 2^{-B+1}T \right] \quad (357)$$

$$\leq 8(1 + \sqrt{c_1/c_2}) \left[ c_2 \ln T \left( 2 \cdot \frac{2^{B(d_z+1)}}{2^{d_z+1}} \right) + 2^{-B+1}T \right] \quad (358)$$

$$= 8(1 + \sqrt{c_1/c_2}) \left[ 2c_2 2^{(B-1)(d_z+1)} \ln T + 2^{-(B-1)}T \right]. \quad (359)$$

930 This inequality holds for any positive  $B$ . By choosing  $B^* = 1 + \frac{\log_2(\frac{T}{\ln T})}{d_z+2}$ , we have

$$R(T) \leq 8(1 + \sqrt{c_1/c_2}) \left[ 2c_2 \left( \frac{T}{\ln T} \right)^{\frac{(d_z+1)}{(d_z+2)}} \ln T + \left( \frac{\ln T}{T} \right)^{\frac{1}{(d_z+2)}} T \right] \quad (360)$$

$$= 8(1 + \sqrt{c_1/c_2}) \left[ 2c_2 T^{\frac{(d_z+1)}{(d_z+2)}} \ln T^{1 - \frac{(d_z+1)}{(d_z+2)}} + T^{1 - \frac{1}{(d_z+2)}} \ln T^{\frac{1}{(d_z+2)}} \right] \quad (361)$$

$$= 8(1 + \sqrt{c_1/c_2}) \left[ 2c_2 T^{\frac{(d_z+1)}{(d_z+2)}} \ln T^{\frac{1}{(d_z+2)}} + T^{\frac{(d_z+1)}{(d_z+2)}} \ln T^{\frac{1}{(d_z+2)}} \right] \quad (362)$$

$$= 8(1 + \sqrt{c_1/c_2}) (2c_2 + 1) T^{\frac{(d_z+1)}{(d_z+2)}} \ln T^{\frac{1}{(d_z+2)}}. \quad (363)$$

931

□

932 **Corollary 9.** Setting  $c_1 = 2c^2$  and  $c_2 = \left(\frac{c^2}{2}\right)^{1/3}$  simplifies Theorem 3 such that

$$R(T) \leq 8(1 + (4c^2)^{1/3})^2 T^{\frac{(d_z+1)}{(d_z+2)}} \ln T^{\frac{1}{(d_z+2)}}. \quad (364)$$

933 with probability  $1 - 2T^{-2}$ .

934 *Proof.*

□

## 935 K Experimental Setup and Details

936 Here we provide further details on the experiments.

### 937 K.1 Loss Visualization and Rank Test

938 Figure 1 and claims made in Section 5 analyze several classical matrix games. We report the payoff  
 939 matrices in standard row-player / column-player payoff form below. All games are then shifted and  
 940 scaled so payoffs lie in  $[0, 1]$  (i.e., first by subtracting the minimum and then scaling by the max).

941 RPS:

$$\begin{bmatrix} 0/0 & -1/1 & 1/-1 \\ 1/-1 & 0/0 & -1/1 \\ -1/1 & 1/-1 & 0/0 \end{bmatrix}. \quad (365)$$

942 Chicken:

$$\begin{bmatrix} 0/0 & -1/1 \\ 1/-1 & -3/-3 \end{bmatrix}. \quad (366)$$

943 Matching Pennies:

$$\begin{bmatrix} 1/-1 & -1/1 \\ -1/1 & 1/-1 \end{bmatrix}. \quad (367)$$

944 Modified-Shapleys:

$$\begin{bmatrix} 1/-0.5 & 0/1 & 0.5/0 \\ 0.5/0 & 1/-0.5 & 0/1 \\ 0/1 & 0.5/0 & 1/-0.5 \end{bmatrix}. \quad (368)$$

945 Prisoner's Dilemma:

$$\begin{bmatrix} -1/-1 & -3/0 \\ 0/-3 & -2/-2 \end{bmatrix}. \quad (369)$$

### 946 K.1.1 NashConv is Biased

947 We use Chicken to demonstrate the effect of sampled play on the bias of the popular NashConv loss.  
 948 NashConv is unable to capture the interior Nash equilibrium because of its high bias. In contrast, our proposed loss  $\mathcal{L}^\tau$  is guaranteed to capture all equilibria at low temperature  $\tau$ .

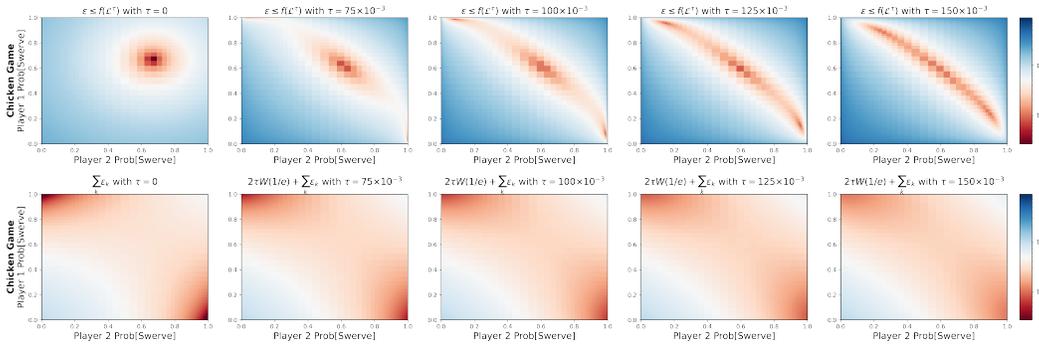


Figure 5: Effect of Sampled Play on a Biased Loss. The first row displays the expected upper bound guaranteed by our proposed loss  $\mathcal{L}^\tau$  (also displayed in Figure 1). The second row displays the expectation of NashConv under sampled play, i.e.,  $\sum_k \epsilon_k$  where  $\epsilon_k = \mathbb{E}_{a_{-k} \sim x_{-k}} [\max_{a_k} u_k^\tau(\mathbf{a})] - \mathbb{E}_{a \sim x} [u_k^\tau(\mathbf{a})]$ . To be consistent, we add the offset  $n\tau W(1/e) + \sum_k \epsilon_k$  to NashConv per Lemma 14, which relates the exploitability at positive temperature to that at zero temperature. The resulting loss surface clearly shows NashConv fails to recognize the interior Nash equilibrium due to its inherent bias. NashConv succeeds in finding pure equilibria because sampling from a pure joint equilibrium is a deterministic process (no noise means no bias).

949

### 950 K.2 Saddle Point Analysis

951 To generate Figure 2, we follow a procedure similar to the study of MNIST in [12] (Section 3 of  
 952 Supp.). Their recommended procedure searches for critical points in two ways. The first repeats a  
 953 randomized, iterative optimization process 20 times. They then sample one these 20 trials at random,  
 954 select a random point along the descent trajectory, and search for a critical point (using Newton's  
 955 method) nearby. They repeat this sampling process 100 times. The second approach randomly selects

956 a feasible point in the decision set and searches for a critical point nearby (again using Newton’s  
 957 method). They also perform this 100 times.

958 Our protocol differs from theirs slightly in a few respects. One, we use SGD, rather than the saddle-  
 959 free Newton algorithm to trace out an initial descent trajectory. Two, we do not add noise to strategies  
 960 along the descent trajectory prior to looking for critical points. Lastly, we use different experimental  
 961 hyperparameters. We run SGD for 1000 iterations rather than 20 epochs and rerun SGD 100 times  
 962 rather than 20. We sample 1000 points for each of the two approaches for finding critical points.

### 963 K.3 SGD on Classical Games

964 The games examined in Figure 3 were all taken from [15]. Each is available via open source  
 965 implementations in OpenSpiel [22] or GAMUT [33].

966 We compare against several other baselines, replicating the experiments in [15]. RM indicates  
 967 regret-matching and FTRL indicates follow-the-regularized-leader. These are, arguably, the two most  
 968 popular scalable stochastic algorithms for approximating Nash equilibria.  $y\text{QRE}^{auto}$  is a stochastic  
 969 algorithm developed in [15].

970 For each of the experiments, we sweep over learning rates in log-space from  $10^{-3}$  to  $10^2$  in increments  
 971 of 1. We also consider whether to run SGD with the projected-gradient and whether to constrain  
 972 iterates to the simplex via Euclidean projection or entropic mirror descent [6]. We then presented the  
 973 results of the best performing hyperparameters. This was the same approach taken in [15].

974 **Saddle Points in Blotto** To confirm the existence of saddle points, we computed the Hessian of  
 975  $\mathcal{L}(x_{10k})$  for SGD ( $s = \infty$ ), deflated the matrix by removing from its eigenvectors all directions  
 976 orthogonal to the simplex, and then computed its top- $(n\bar{m} - n)$  eigenvalues. We do this because  
 977 there always exists a  $n$ -dimensional nullspace of the Hessian at zero temperature that lies outside the  
 978 tangent space of the simplex, and we only care about curvature within the tangent space. Specifically,  
 979 at an equilibrium  $x$ , if we compute  $z^\top \text{Hess}(\mathcal{L})z$  where  $z$  is formed as a linear combination of the  
 980 vectors  $\{[x_1, 0, \dots, 0]^\top, \dots, [0, \dots, x_n]^\top\}$ , then each block  $\tilde{B}_{kl}$  is identically zero at an equilibrium:  
 981  $\tilde{B}_{kl}x_l = \sqrt{\eta_k}[I - \frac{1}{m_k}\mathbf{1}\mathbf{1}^\top]H_{kl}^kx_l = \sqrt{\eta_k}\Pi_{T\Delta}(\nabla_{x_k}^k) = 0$ . By Lemma 17, this implies there is zero  
 982 curvature of the loss in the direction  $z$ :  $z^\top \text{Hess}(\mathcal{L})z = 0$ .

### 983 K.4 BLiN on Artificial Game

984 To construct the 7-player, 2-action, symmetric, artificial game in Figure 4, we used the following  
 985 coefficients (discovered by trial-and-error):

$$\begin{bmatrix} 0.09906873 & 0 & 0.23116037 & 0 & 0.62743528 & 0 & 0.19813746 \\ 0 & 0.33022909 & 0 & 0.03302291 & 0 & 0.62743528 & 0 \end{bmatrix}. \quad (370)$$

986 The first row indicates the payoffs received when player  $i$  plays action 0 and the background  
 987 population plays any of the possible joint actions (number of combinations with replacement). For  
 988 example, the first column indicates the payoff when all background players play action 0. The second  
 989 column indicates all background players play action 0 except for one which plays action 1, and so on.  
 990 The last column indicates all background players play action 1. These  $2n$  scalars uniquely define the  
 991 payoffs of a symmetric game.

992 Given that this game only has two actions, we represent a mixed strategy by a single scalar  $p \in [0, 1]$ ,  
 993 i.e., the probability of the first action. Furthermore, this game is symmetric and we seek a symmetric  
 994 equilibrium, so we can represent a full Nash equilibrium by this single scalar  $p$ . This reduces our  
 995 search space from  $7 \times 2 = 14$  variables to 1 variable (and obviates any need for a map  $s$  from the  
 996 unit hypercube to the simplex—see Lemma 24).