
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 ω 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

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

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

Can we solve for Nash equilibria via unbiased stochastic optimization?

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

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

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

2 Preliminaries

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

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

As an abuse of notation, let the atomic action $a_i = e_i$ also denote the m_i -dimensional “one-hot” vector 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/ $\text{BR}_k \leftarrow \text{aBR}_k = \Pi_{\Delta}(x_k + \eta \nabla_{x_k} u_k(\mathbf{x}))$	Π_{Δ} 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 Δ^{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 ϵ .
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 BR_i if and only if it has zero projected-gradient¹ norm:*

¹Not 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 η_k : $aBR_k = x_k + \eta_k \Pi_{T\Delta}(\nabla_{x_k}^k)$.*

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 Π_{Δ} 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	$n \prod_{i=1}^n m_i$	$2nm$	$2n$

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

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

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 ϵ 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:

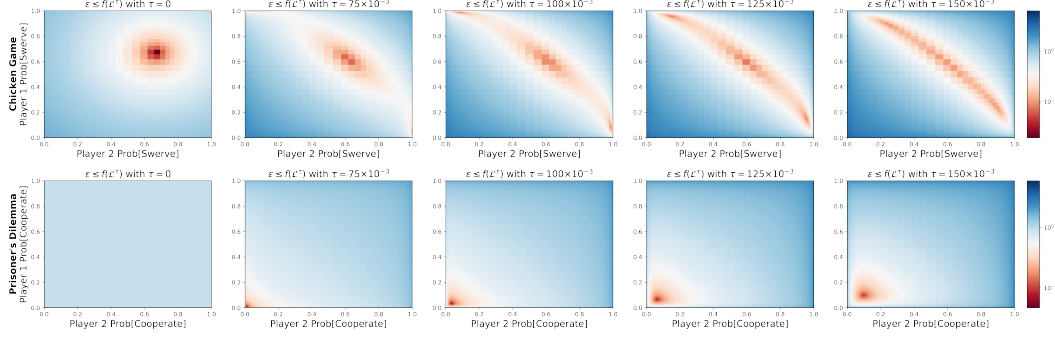


Figure 1: Upper Bound ($\epsilon \leq f(\mathcal{L}^\tau)$) Heatmap Visualization. The first row examines the loss landscape for the classic anti-coordination game of Chicken (Nash equilibria: $(0, 1)$, $(1, 0)$, $(2/3, 1/3)$) while the second row examines the Prisoner’s dilemma (Unique Nash equilibrium: $(0, 0)$). 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.

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

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

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

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

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

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

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

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

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

5 Analysis

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

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

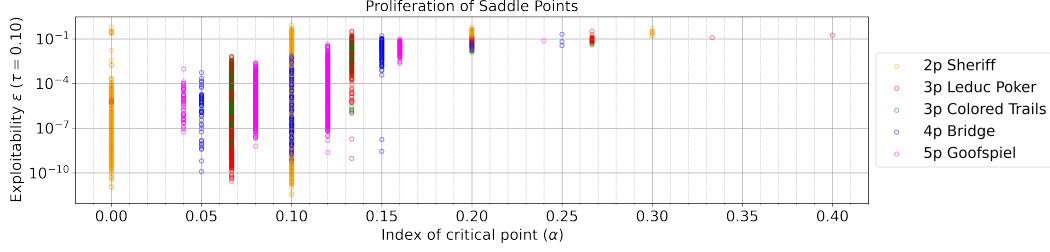


Figure 2: We reapply the analysis of [12], originally designed to understand the success of SGD in deep learning, to “slices” of several popular extensive form games. To construct a slice (or *meta-game*), we randomly sample 6 deterministic policies and then consider the corresponding n -player, 6-action 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}^τ 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 α 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 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}^τ 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}^τ '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.

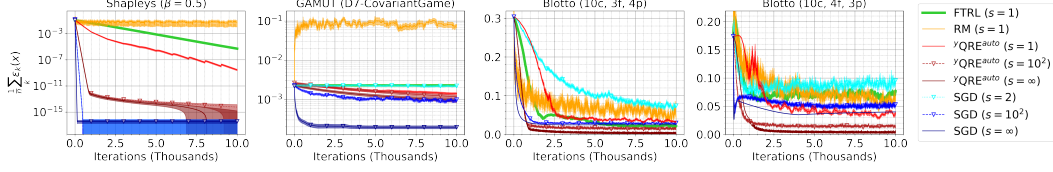


Figure 3: Comparison of SGD on $\mathcal{L}^{\tau=0}$ against baselines on four games evaluated in [15]. From left to right: 2-player, 3-action, nonsymmetric; 6-player, 5-action, nonsymmetric; 4-player, 66-action, symmetric; 3-player, 286-action, symmetric. SGD struggles at saddle points in Blotto.

6 Algorithms

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

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

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

6.1 High Probability, Polynomial Convergence Rates

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

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

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

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

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

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

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

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

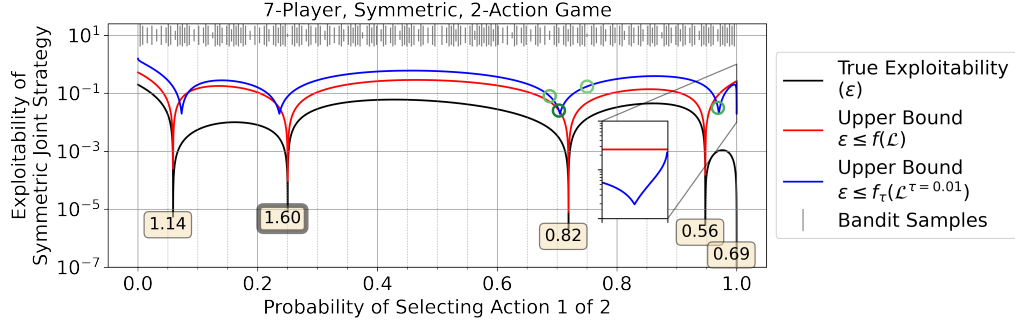


Figure 4: Bandit-based (BLiN) Nash solver applied to an artificial 7-player, symmetric, 2-action game. We search for a symmetric equilibrium, which is represented succinctly as the probability of selecting action 1. The plot shows the true exploitability ϵ of all symmetric strategies in black and indicates 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}^τ , 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}^τ . 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}^τ 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(τ) 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 ϵ , demonstrating the challenges of local minima.

282 We demonstrate BLiN (applied to \mathcal{L}^τ) 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

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

- [18] J. C. Harsanyi, R. Selten, et al. A general theory of equilibrium selection in games. *MIT Press Books*, 1, 1988.
- [19] E. Hazan, K. Singh, and C. Zhang. Efficient regret minimization in non-convex games. In *International Conference on Machine Learning*, pages 1433–1441. PMLR, 2017.
- [20] E. Janovskaja. Equilibrium points in polymatrix games. *Lithuanian Mathematical Journal*, 8 (2):381–384, 1968.
- [21] M. Lanctot, V. Zambaldi, A. Gruslys, A. Lazaridou, K. Tuyls, J. Pérolat, D. Silver, and T. Graepel. A unified game-theoretic approach to multiagent reinforcement learning. In *Advances in Neural Information Processing Systems*, pages 4190–4203, 2017.
- [22] M. Lanctot, E. Lockhart, J.-B. Lespiau, V. Zambaldi, S. Upadhyay, J. Pérolat, S. Srinivasan, F. Timbers, K. Tuyls, S. Omidshafiei, D. Hennes, D. Morrill, P. Muller, T. Ewalds, R. Faulkner, J. Kramár, B. D. Vylder, B. Saeta, J. Bradbury, D. Ding, S. Borgeaud, M. Lai, J. Schrittwieser, T. Anthony, E. Hughes, I. Danihelka, and J. Ryan-Davis. OpenSpiel: A framework for reinforcement learning in games. *CoRR*, abs/1908.09453, 2019. URL <http://arxiv.org/abs/1908.09453>.
- [23] S. Leonardos, G. Piliouras, and K. Spendlove. Exploration-exploitation in multi-agent competition: convergence with bounded rationality. *Advances in Neural Information Processing Systems*, 34:26318–26331, 2021.
- [24] W. Li, H. Li, J. Honorio, and Q. Song. Pyxab – a python library for \mathcal{X} -armed bandit and online blackbox optimization algorithms, 2023. URL <https://arxiv.org/abs/2303.04030>.
- [25] C. K. Ling, F. Fang, and J. Z. Kolter. What game are we playing? end-to-end learning in normal and extensive form games. *arXiv preprint arXiv:1805.02777*, 2018.
- [26] S. Mandt, M. D. Hoffman, and D. M. Blei. Stochastic gradient descent as approximate bayesian inference. *Journal of Machine Learning Research*, 18:1–35, 2017.
- [27] L. Marris, I. Gemp, and G. Piliouras. Equilibrium-invariant embedding, metric space, and fundamental set of 2x2 normal-form games. *arXiv preprint arXiv:2304.09978*, 2023.
- [28] R. D. McKelvey and T. R. Palfrey. Quantal response equilibria for normal form games. *Games and Economic Behavior*, 10(1):6–38, 1995.
- [29] D. Milec, J. Černý, V. Lisý, and B. An. Complexity and algorithms for exploiting quantal opponents in large two-player games. *Proceedings of the AAAI Conference on Artificial Intelligence*, 35(6):5575–5583, 2021.
- [30] P. R. Milgrom and R. J. Weber. A theory of auctions and competitive bidding. *Econometrica: Journal of the Econometric Society*, pages 1089–1122, 1982.
- [31] K. G. Murty and S. N. Kabadi. Some NP-complete problems in quadratic and nonlinear programming. Technical report, 1985.
- [32] H. Nikaidô and K. Isoda. Note on non-cooperative convex games. *Pacific Journal of Mathematics*, 5(1):807815, 1955.
- [33] E. Nudelman, J. Wortman, Y. Shoham, and K. Leyton-Brown. Run the GAMUT: A comprehensive approach to evaluating game-theoretic algorithms. In *AAMAS*, volume 4, pages 880–887, 2004.
- [34] J. Pérolat, S. Perrin, R. Elie, M. Laurière, G. Piliouras, M. Geist, K. Tuyls, and O. Pietquin. Scaling mean field games by online mirror descent. In *Proceedings of the 21st International Conference on Autonomous Agents and Multiagent Systems*, 2022.
- [35] A. Raghunathan, A. Cherian, and D. Jha. Game theoretic optimization via gradient-based Nikaido-Isoda function. In *International Conference on Machine Learning*, pages 5291–5300. 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.