

# Perturbing Best Responses in Zero-Sum Games

Adam Dziwoki, Rostislav Horčík

Czech Technical University in Prague  
dziwoada@fel.cvut.cz, xhorcik@fel.cvut.cz

## Abstract

This paper examines the impact of perturbations on best-response-based algorithms that approximate Nash equilibria in zero-sum games, specifically Double Oracle and Fictitious Play. More precisely, we assume that the oracle computing the best responses perturbs the utilities before selecting the best response. We show that using such an oracle reduces the number of iterations for both algorithms. For some cases, suitable perturbations ensure the expected number of iterations is logarithmic. Although the utility perturbation is computationally demanding as it requires iterating through all pure strategies, we demonstrate that one can efficiently perturb the utilities in games where pure strategies have further inner structure.

## Code —

<https://github.com/geoborek/perturbing-best-responses>

## Extended version — <https://arxiv.org/abs/2511.12523>

## Introduction

Computing Nash equilibria (NE) in two-player zero-sum games with huge strategy spaces is a computationally demanding problem. Among algorithms approximating NE in such games, a substantial role is played by those based on *best-response oracles* (BROs) due to their ability to consider only a subspace of the strategy spaces. Prominent examples of such algorithms are *Fictitious Play* (FP) introduced in (Brown 1951) and *Double Oracle* (DO) (McMahan, Gordon, and Blum 2003). The latter served as a basis for algorithms leveraging deep reinforcement learning to approximate best responses, such as *Policy Space Response Oracles* (Lanctot et al. 2017; McAleer et al. 2022; Bighashdel et al. 2024).

It is known that each two-player zero-sum game has an approximated  $\varepsilon$ -NE of logarithmic size in the number of pure strategies  $n$  for given  $\varepsilon > 0$  (Althöfer 1994; Lipton and Young 1994). On the other hand, it is not difficult to construct games where FP and DO need at least  $n$  iterations to find an  $\varepsilon$ -NE. This raises a natural question whether there is a BRO-based algorithm computing  $\varepsilon$ -NE in a logarithmic number of iterations. (Hazan and Koren 2016) answered

this question negatively as they proved that any (randomized) BRO-based algorithm needs at least  $\Omega(\sqrt{n}/\log^3 n)$  iterations. At the same time, they introduced quite a complicated BRO-based algorithm providing a quadratic speed up  $O(\sqrt{n}/\varepsilon^2)$  (up to a poly-logarithmic factor). Consequently, it follows that one needs to make BRO more powerful to achieve logarithmic number of iterations.

One way to make BRO stronger is by adding perturbations. Such an oracle perturbs the utilities before computing a best response for a given mixed strategy. We call such an oracle *perturbed BRO* (PBRO). A PBRO-variant of FP known as *Stochastic Fictitious Play* (SFP) was introduced in (Fudenberg and Kreps 1993). Its convergence to NE in zero-sum games was proven in (Hofbauer and Sandholm 2002). In this paper, we prove that SFP achieves the logarithmic complexity in the number of pure strategies  $n$  in expectation (note that PBRO makes SFP a randomized algorithm). Analogously, we define a PBRO-based variant of DO called *Stochastic Double Oracle* (SDO). Although its complexity in terms of  $n$  is still an open problem, we prove that perturbations ensure logarithmic behavior in expectation for some examples where DO needs  $O(n)$  iterations, namely for examples introduced in (Zhang and Sandholm 2024). We also tested experimentally on other games that perturbations reduce the number of iterations. However, it turned out that the perturbations do not accelerate the convergence in random games. Implementing perturbations into a BRO for normal-form games requires iterating through all pure strategies to perturb the utility, which is not computationally efficient. However, for games whose strategy spaces have an inner structure like *partially-observable stochastic game* (POSG), one can efficiently perturb only the rewards for transitions or terminal states. Even though such perturbations do not precisely correspond to the perturbations in normal-form games, we experimentally show that they are able to reduce the number of iterations. More precisely, we demonstrate that on POSGs from (Zhang and Sandholm 2024) and a path-planning game where one player looks for the shortest path in a grid, while the other player might choose any edge and multiply its cost by a fixed coefficient.

## Related Work

The convergence of FP was studied in several papers focusing mainly on the convergence rate w.r.t.  $\varepsilon$  instead of the

size of the game. FP converges to NE in zero-sum games, as proved in (Robinson 1951). Its convergence rate is one of the oldest open problems in game theory, known as Karlin’s conjecture (Karlin 1959). However, there are several partial results (Abernethy, Lai, and Wibisono 2021; Daskalakis and Pan 2014). Recently, *Anticipatory Fictitious Play* (AFP) was introduced (Cloud, Wang, and Kerr 2023). AFP still needs  $O(n)$  iterations in the worst case, but provides often faster convergence. Moreover, it calls the BRO four times in an iteration, whereas FP only calls the BRO two times. It is known that perturbations are related to the regularization (Hofbauer and Sandholm 2002; Abernethy et al. 2014). A regularized versions of AFP (PU and OMWU) were investigated in (Cen, Wei, and Chi 2024). Both variants need at most  $O(\log n)$  iterations, but they maintain probability distributions over all pure strategies, making them unsuitable for large games.

It is easy to see that DO needs  $\Theta(n)$  iterations. However, to our knowledge, there is no deeper theoretical study of its convergence rate. Exponential lower bounds for DO applied to POSGs were obtained in (Zhang and Sandholm 2024). Further variants of DO were investigated in (McAleer et al. 2022) but without any convergence guarantees.

## Background

This section introduces our notation and the necessary background for the paper. The set of natural numbers  $\{1, \dots, n\}$  is denoted  $[n]$ . The real-valued vectors are denoted by bold lowercase letters, e.g.  $\mathbf{p} \in \mathbb{R}^n$ . Its  $i$ -th component is denoted  $p_i$ . The projection into the  $i$ -th component is the map  $\pi_i(\mathbf{p}) = p_i$ . The vectors from the standard basis of  $\mathbb{R}^n$  are denoted  $\mathbf{e}_1, \dots, \mathbf{e}_n$ . Thus  $\pi_i(\mathbf{e}_j) = 1$  iff  $i = j$  and  $\pi_i(\mathbf{e}_j) = 0$  otherwise. Similarly, bold uppercase letters denote matrices, e.g.  $\mathbf{M} \in \mathbb{R}^{m \times n}$ . The entry of  $\mathbf{M}$  in a row  $i \in [m]$  and a column  $j \in [n]$  is denoted  $M(i, j) = \mathbf{e}_i^\top \mathbf{M} \mathbf{e}_j$ . The symbol  $\Delta_n$  stands for the simplex of probability distributions over  $[n]$ . We identified the members of  $\Delta_n$  with vectors  $\mathbf{p} \in [0, 1]^n$  such that  $\sum_{i=1}^n p_i = 1$ . Given sets of indexes  $R \subseteq [m]$  and  $C \subseteq [n]$ ,  $\mathbf{M}[R, C]$  is the submatrix of  $\mathbf{M}$  having only rows with indexes in  $R$  and columns with indexes in  $C$ .

This paper focuses on 2-player finite zero-sum games that can be identified with matrix games if we index the sets of pure strategies by natural numbers. A *matrix game* is given by a matrix  $\mathbf{M} \in \mathbb{R}^{m \times n}$ . The row indexes  $[m]$  correspond to the set of pure strategies of the *row player*. Analogously, the column indexes  $[n]$  correspond to the pure strategies of the *column player*. The entry  $M(i, j)$  denotes the game’s outcome if the row player chooses the  $i$ -th strategy and the column player the  $j$ -th strategy. A *mixed strategy* for the row player is a probability distribution  $\mathbf{p} \in \Delta_m$  and analogously for the column player. The expected outcome of the game when the row (resp. column) player plays  $\mathbf{p} \in \Delta_m$  (resp.  $\mathbf{q} \in \Delta_n$ ) can be expressed as  $M(\mathbf{p}, \mathbf{q}) = \mathbf{p}^\top \cdot \mathbf{M} \cdot \mathbf{q}$ . We assume that the entries in  $\mathbf{M}$  represent losses (resp. rewards) for the row (resp. column) player. Thus, the row player chooses her strategy to minimize the expected outcome, whereas the column player wants to maximize it.

The FP and DO algorithms rely on a best-response oracle for each player. Given a mixed strategy  $\mathbf{q} \in \Delta_n$ , the oracle for the row player computes the index of the best response and the corresponding value, namely

$$\text{BR}_r(\mathbf{q}) \in \operatorname{argmin}_{i \in [m]} \pi_i(\mathbf{M}\mathbf{q}), \quad \text{BRVal}_r(\mathbf{q}) = \min_{i \in [m]} \pi_i(\mathbf{M}\mathbf{q})$$

Analogously for  $\mathbf{p} \in \Delta_m$ , the column player’s oracle computes

$$\text{BR}_c(\mathbf{p}) \in \operatorname{argmax}_{j \in [n]} \pi_j(\mathbf{p}^\top \mathbf{M}), \quad \text{BRVal}_c(\mathbf{p}) = \max_{j \in [n]} \pi_j(\mathbf{p}^\top \mathbf{M})$$

Note that  $\text{BR}_r(\mathbf{q})$  and  $\text{BR}_c(\mathbf{p})$  are pure strategies. In general, there might be several best responses for a given mixed strategy, and any of them can be returned by the oracle.

Given  $\varepsilon \geq 0$ , a pair of strategies  $\langle \mathbf{p}^*, \mathbf{q}^* \rangle$  is said to be  $\varepsilon$ -*Nash equilibrium* ( $\varepsilon$ -NE) of the game  $\mathbf{M}$  if the following inequalities hold:

$$M(\mathbf{p}^*, \mathbf{q}^*) - \text{BRVal}_r(\mathbf{q}^*) \leq \varepsilon, \quad \text{BRVal}_c(\mathbf{p}^*) - M(\mathbf{p}^*, \mathbf{q}^*) \leq \varepsilon$$

A *Nash equilibrium* (NE) of the game  $\mathbf{M}$  is 0-NE. Given a pair of strategies  $\langle \mathbf{p}, \mathbf{q} \rangle$ , a sufficient condition for the pair to form  $\varepsilon$ -NE is

$$\text{BRVal}_c(\mathbf{p}) - \text{BRVal}_r(\mathbf{q}) \leq \varepsilon \quad (1)$$

that we use as a termination condition for FP and DO.

This paper investigates the influence of perturbations before selecting the best response. Let  $\mathbf{u}$  and  $\mathbf{v}$  be random vectors with respective dimensions  $m, n$  whose components are i.i.d. random variables coming from a given distribution. *Perturbed best responses* are defined as follows:

$$\begin{aligned} \widetilde{\text{BR}}_r(\mathbf{q}) &\in \operatorname{argmin}_{i \in [m]} \pi_i(\mathbf{M}\mathbf{q} - \mathbf{u}), \\ \widetilde{\text{BR}}_c(\mathbf{p}) &\in \operatorname{argmax}_{j \in [n]} \pi_j(\mathbf{p}^\top \mathbf{M} + \mathbf{v}) \end{aligned}$$

When we use the perturbed best responses in FP or DO, we assume that the probability distribution for the perturbations vectors  $\mathbf{u}, \mathbf{v}$  is fixed.

This paper considers two common distributions *uniform*  $U(a, b)$  on the interval  $[a, b]$  and *Gumbel*  $G(\mu, \beta)$ . The Gumbel distribution  $G(\mu, \beta)$  is a continuous probability distribution with a CDF given by  $F(x) = e^{-e^{-(x-\mu)/\beta}}$  where  $\mu$  is a location and  $\beta > 0$  a scale. The Gumbel distribution is tightly related to the well-known function softmax, see e.g. (Goodfellow, Bengio, and Courville 2016), mapping  $\mathbf{x} \in \mathbb{R}^n$  to a probability distribution  $\text{softmax}(\mathbf{x}) \in \Delta_n$  defined by

$$\pi_i(\text{softmax}(\mathbf{x})) = \frac{e^{x_i}}{\sum_{j=1}^n e^{x_j}}$$

The following lemma is due to (Gumbel 1954) and is called *Gumbel-max trick* in the machine learning community; for details see (Train 2009, Chapter 3) or (Franke and Degen 2023). Given  $\mathbf{x} \in \mathbb{R}^n$ , the Gumbel-max trick allows to sample from  $\text{softmax}(\mathbf{x})$  by taking  $\operatorname{argmax}$  of the perturbed values  $x_i$ .

**Lemma 1.** *Let  $\mathbf{x} \in \mathbb{R}^n$  and  $i^* \in \{1, \dots, n\}$ . Let  $\mathbf{z} \in \mathbb{R}^n$  be a random vector whose components  $z_i$  be drawn i.i.d. from  $G(0, \beta)$ . Then the probability*

$$P(\operatorname{argmax}_{i \in [n]} \pi_i(\mathbf{x} + \mathbf{z}) = i^*) = \pi_{i^*}(\text{softmax}(\mathbf{x}/\beta))$$

---

**Algorithm 1: Stochastic Fictitious Play**


---

**Require:** Initial strategies  $k, l$ ;  $\varepsilon > 0$ ; a probability distribution for  $\widetilde{\text{BR}}_r$  and  $\widetilde{\text{BR}}_c$

**Ensure:**  $\varepsilon$ -NE  $\langle \mathbf{p}^*, \mathbf{q}^* \rangle$

```

 $t \leftarrow 1$ 
 $\mathbf{p} \leftarrow \mathbf{e}_k, \mathbf{q} \leftarrow \mathbf{e}_l$ 
 $lb \leftarrow \text{BRVal}_r(\mathbf{q}), ub \leftarrow \text{BRVal}_c(\mathbf{p})$ 
while  $ub - lb > t\varepsilon$  do
   $t \leftarrow t + 1$ 
   $i \leftarrow \widetilde{\text{BR}}_r(\mathbf{q}), j \leftarrow \widetilde{\text{BR}}_c(\mathbf{p})$ 
   $\mathbf{p} \leftarrow \mathbf{p} + \mathbf{e}_i, \mathbf{q} \leftarrow \mathbf{q} + \mathbf{e}_j$ 
   $lb \leftarrow \text{BRVal}_r(\mathbf{q}), ub \leftarrow \text{BRVal}_c(\mathbf{p})$ 
end while
return  $\langle \mathbf{p}^*, \mathbf{q}^* \rangle = \langle \mathbf{p}/t, \mathbf{q}/t \rangle$ 

```

---

Dually, it follows from Lemma 1 that

$$P(\text{argmin}_{i \in [n]} \pi_i(\mathbf{x} - \mathbf{z}) = i^*) = \pi_{i^*}(\text{softmax}(-\mathbf{x}/\beta))$$

as  $\text{argmin}_i \pi_i(\mathbf{x} - \mathbf{z}) = i^*$  iff  $\text{argmax}_i \pi_i(-\mathbf{x} + \mathbf{z}) = i^*$ . Consequently, if we perturb the utilities with the Gumbel distribution  $G(0, \beta)$ , the perturbed best-response  $\widetilde{\text{BR}}_r(\mathbf{q})$  is sampled from  $\text{softmax}(\mathbf{M}\mathbf{q}/\beta)$  and analogously  $\widetilde{\text{BR}}_c(\mathbf{p})$  from  $\text{softmax}(-\mathbf{p}^\top \mathbf{M}/\beta)$ .

To prove our result on SFP, we need to recall the *randomized exponentially weighted forecaster* (REWF), introduced in (Cesa-Bianchi and Lugosi 2006, Chapter 4). Let  $\mathbf{M} \in [0, 1]^{m \times n}$  be a  $m \times n$ -matrix viewed as a loss function  $[m] \times [n] \rightarrow [0, 1]$ . Consider a “sort of” game between a player and an opponent with  $T \geq 1$  rounds. In the round  $t$ , the player chooses an action  $i_t \in [m]$  and the opponent an action  $j_t \in [n]$ , making the player to suffer the loss  $\mathbf{M}(i_t, j_t)$ . The player’s goal is to minimize the cumulative loss  $\sum_{t=1}^T \mathbf{M}(i_t, j_t)$ . REWF is a randomized strategy for the player based on the opponent’s previously selected actions. In the round  $t$ , REWF samples  $i_t$  based on  $j_1, \dots, j_{t-1}$  from the distribution  $\text{softmax}(-\eta \sum_{s=1}^{t-1} \mathbf{M}e_{j_s})$ , where  $\eta > 0$  is a parameter. Following REWF ensures an upper bound on the regret, i.e., the difference between the cumulative loss and the best fixed player’s action  $i$ ; it immediately follows from (Cesa-Bianchi and Lugosi 2006, Theorem 2.2 and Lemma 4.1).

**Corollary 2.** Let  $\mathbf{M} \in [0, 1]^{m \times n}$ ,  $T \geq 1$ ,  $\eta > 0$ , and  $\delta \in (0, 1)$ . REWF satisfies, with probability  $1 - \delta$ ,

$$\sum_{t=1}^T \mathbf{M}(i_t, j_t) - \min_{i \in [m]} \sum_{t=1}^T \mathbf{M}(i, j_t) \leq \frac{\ln m}{\eta} + \frac{T\eta}{8} + \sqrt{\frac{T}{2} \ln \frac{1}{\delta}}$$

## Stochastic Fictitious Play

*Stochastic Fictitious Play* (SFP) is a randomized version of *Fictitious Play* (FP). Its pseudocode is shown in Algorithm 1. We expand it with the termination condition (1) so that it stops after reaching  $\varepsilon$ -NE. Note that the best-response values  $\text{BRVal}_r(\mathbf{q})$ ,  $\text{BRVal}_c(\mathbf{p})$  are not perturbed, hence the termination condition is not affected by the perturbations.

The algorithm maintains two vectors  $\mathbf{p}$  and  $\mathbf{q}$  representing multisets of already played pure strategies. In each iteration, SFP finds perturbed best responses w.r.t.  $\mathbf{p}$  and  $\mathbf{q}$  respectively<sup>1</sup>. When the termination condition is satisfied, SFP returns the average of all the played strategies so that the resulting vectors  $\mathbf{p}^*, \mathbf{q}^*$  form probability distributions. Note that even though Algorithm 1 operates with vectors  $\mathbf{e}_i \in \mathbb{R}^m$ ,  $\mathbf{e}_j \in \mathbb{R}^n$ , SFP can be implemented efficiently without considering all the dimensions by storing only the multisets of pure strategies played by the players up to the current iteration.

Using the Gumbel-max trick (see Lemma 1), we show that SFP combined with Gumbel perturbations corresponds to the computation when both players apply the randomized exponentially weighted forecaster (REWF) against each other. W.l.o.g. we assume that the matrix game  $\mathbf{M} \in \mathbb{R}^{m \times n}$  is normalized into  $[0, 1]$  and  $m \leq n$ .

**Theorem 3.** Let  $\mathbf{M} \in [0, 1]^{m \times n}$  be a matrix game such that  $m \leq n$  and  $\varepsilon > 0$ . An  $\varepsilon$ -NE can be computed by SFP with perturbations sampled from  $G(0, \beta)$  in  $O(\frac{\log n}{\varepsilon^2})$  expected iterations, where  $\beta = \frac{2 + \sqrt{2 \ln n}}{\varepsilon \sqrt{8 \ln n}}$ .

*Proof.* We first prove that SFP finds  $\varepsilon$ -NE in  $O(\frac{\log n}{\varepsilon^2})$  iterations with probability at least  $1/2$ . We let SFP to iterate for  $T = \left(\frac{2 + \sqrt{2 \ln n}}{\varepsilon}\right)^2 \in O(\frac{\log n}{\varepsilon^2})$  iterations.

Let  $\mathbf{p}_t, \mathbf{q}_t$  denote the vectors  $\mathbf{p}, \mathbf{q}$  at the iteration  $t \leq T$ . Further, let  $i_t = \widetilde{\text{BR}}_r(\mathbf{q}_{t-1})$  and  $j_t = \widetilde{\text{BR}}_c(\mathbf{p}_{t-1})$ . At the iteration  $t$ , we have  $\mathbf{q}_{t-1} = \sum_{s=1}^{t-1} \mathbf{e}_{j_s}$ . Thus  $\mathbf{M}\mathbf{q}_{t-1} = \sum_{s=1}^{t-1} \mathbf{M}e_{j_s}$ . By Lemma 1,  $i_t$  is sampled from the distribution  $\text{softmax}(-\eta \sum_{s=1}^{t-1} \mathbf{M}e_{j_s})$  for  $\eta = 1/\beta$ . By our choice of  $\beta$  and  $T$ , we have  $\eta = \sqrt{8 \ln n / T}$ . Consequently, Corollary 2 for  $\delta = 1/4$  implies with probability at least  $3/4$ , using also  $m \leq n$  and  $\ln 4 \leq 2$ ,

$$\begin{aligned} \sum_{t=1}^T \mathbf{M}(i_t, j_t) - \min_{i \in [m]} \sum_{t=1}^T \mathbf{M}(i, j_t) &\leq \frac{\ln m}{\eta} + \frac{T\eta}{8} + \sqrt{\frac{T}{2} \ln 4} \\ &\leq \frac{\ln n}{\eta} + \frac{T\eta}{8} + \sqrt{T} = \sqrt{T} \left(1 + \sqrt{\frac{\ln n}{2}}\right) \end{aligned}$$

Applying Corollary 2 to the matrix  $\mathbf{1} - \mathbf{M}^\top$ , we can derive an analogous bound for the column player with probability at least  $3/4$

$$\max_{j \in [n]} \sum_{t=1}^T \mathbf{M}(i_t, j) - \sum_{t=1}^T \mathbf{M}(i_t, j_t) \leq \sqrt{T} \left(1 + \sqrt{\frac{\ln n}{2}}\right)$$

Consequently, both inequalities hold simultaneously with probability at least  $1/2$  using the well-known lower bound  $P(A \cap B) \geq P(A) + P(B) - 1$  on the probability of the intersection of two events  $A, B$ . Summing the above inequalities, we get

$$\max_{j \in [n]} \sum_{t=1}^T \mathbf{M}(i_t, j) - \min_{i \in [n]} \sum_{t=1}^T \mathbf{M}(i, j_t) \leq \sqrt{T}(2 + \sqrt{2 \ln n})$$

---

<sup>1</sup>Formally,  $\mathbf{p}, \mathbf{q}$  are not probability distributions. We assume w.l.o.g. that the best-response oracle works for them as well.

---

**Algorithm 2: Stochastic Double Oracle**


---

**Require:** Initial strategies  $k, l; \varepsilon > 0$ ; a probability distribution for  $\widetilde{\text{BR}}_r$  and  $\widetilde{\text{BR}}_c$

**Ensure:**  $\varepsilon$ -NE  $\langle \mathbf{p}^*, \mathbf{q}^* \rangle$

$t \leftarrow 1$

$R \leftarrow \{k\}, C \leftarrow \{l\}$

$\mathbf{p} \leftarrow \mathbf{e}_k, \mathbf{q} \leftarrow \mathbf{e}_l$

$lb \leftarrow \text{BRVal}_r(\mathbf{q}), ub \leftarrow \text{BRVal}_c(\mathbf{p})$

**while**  $ub - lb > \varepsilon$  **do**

$t \leftarrow t + 1$

$\langle \mathbf{p}, \mathbf{q} \rangle \leftarrow \text{getNash}(M[R, C])$

$i \leftarrow \widetilde{\text{BR}}_r(\mathbf{q}), R \leftarrow R \cup \{i\}$

$j \leftarrow \widetilde{\text{BR}}_c(\mathbf{p}), C \leftarrow C \cup \{j\}$

$lb \leftarrow \text{BRVal}_r(\mathbf{q}), ub \leftarrow \text{BRVal}_c(\mathbf{p})$

**end while**

**return**  $\langle \mathbf{p}^*, \mathbf{q}^* \rangle = \langle \mathbf{p}, \mathbf{q} \rangle$

---

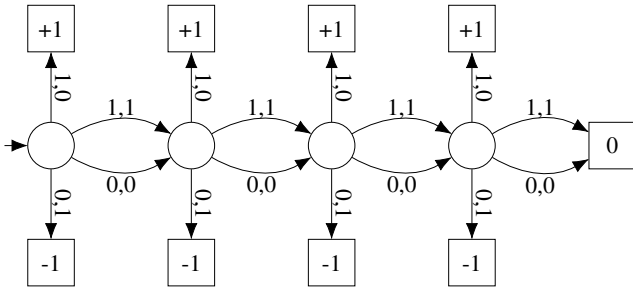


Figure 1: An example of the 4-bit stochastic game.

SFP returns the pair  $\langle \frac{1}{T} \sum_{t=1}^T \mathbf{e}_{i_t}, \frac{1}{T} \sum_{t=1}^T \mathbf{e}_{j_t} \rangle$ . Note that

$$\begin{aligned} ub &= \text{BRVal}_c \left( \frac{1}{T} \sum_{t=1}^T \mathbf{e}_{i_t} \right) = \max_{j \in [n]} \pi_j \left( \frac{1}{T} \sum_{t=1}^T \mathbf{e}_{i_t}^\top \mathbf{M} \right) \\ &= \frac{1}{T} \max_{j \in [n]} \sum_{t=1}^T M(i_t, j) \end{aligned}$$

Analogously,  $lb = \frac{1}{T} \min_{i \in [m]} \sum_{t=1}^T M(i, j_t)$ . Combining the above facts, we get  $ub - lb \leq (2 + \sqrt{2 \ln n}) / \sqrt{T} = \varepsilon$ . Thus  $\langle \frac{1}{T} \sum_{t=1}^T \mathbf{e}_{i_t}, \frac{1}{T} \sum_{t=1}^T \mathbf{e}_{j_t} \rangle$  forms  $\varepsilon$ -NE with probability at least  $1/2$ .

To obtain an algorithm with a logarithmic expected number of iterations, we execute SFP for  $T$  many iterations until we find  $\varepsilon$ -NE. It occurs in the first or second run in expectation.  $\square$

### Stochastic Double Oracle

This section introduces *Stochastic Double Oracle* (SDO), a randomized variant of Double Oracle (DO) with perturbed best responses. Its pseudocode is shown in Algorithm 2. As DO, it maintains two sets of pure strategies  $R, C$  for the row and column player, respectively. At each iteration, it computes NE for the induced subgame  $M[R, C]$ . The termination condition is the same as in SFP.

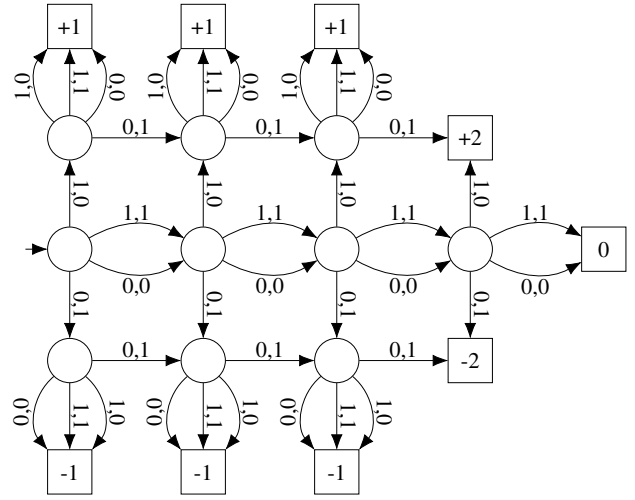


Figure 2: An example of the 4-bit POSG.

In this paper, we take the first steps in investigating convergence for SDO. We start with two matrix games introduced in (Zhang and Sandholm 2024) that are difficult to solve for DO, i.e., DO needs to iterate through all pure strategies to find  $\varepsilon$ -NE. We prove that SDO with suitable perturbations finds  $\varepsilon$ -NE in a logarithmic number of iterations in expectation for both.

**Example 1.** Let  $n \geq 1$ . In the first game, both players choose a number from  $[n]$ . The player who chooses the greater number wins. If the chosen numbers are equal, it is a draw. The corresponding square matrix  $\mathbf{L} \in \mathbb{R}^{n \times n}$  is defined as follows:

$$\mathbf{L} = \begin{pmatrix} 0 & 1 & \cdots & 1 \\ -1 & 0 & \cdots & 1 \\ \vdots & \vdots & \ddots & \vdots \\ -1 & -1 & \cdots & 0 \end{pmatrix}$$

The game  $\mathbf{L}$  has a unique pure  $\varepsilon$ -NE consisting of the last row and column for any  $0 \leq \varepsilon < 1$ . Thus finding  $\varepsilon$ -NE for  $\varepsilon < 1$  is the same as finding NE. If DO is initialized with  $k = l = 1$  and the best responses are chosen adversarially (i.e., we always take the best response with the least index), DO needs  $n$  iterations to find the NE. Note that there are several candidates when DO looks for a best response to a row or column. On the other hand, when we apply SDO to  $\mathbf{L}$ , the perturbations allow us to select a candidate uniformly randomly.

We prove our result on SDO by applying drift theory (Kötzing and Krejca 2019). For technical reasons, we prove it w.l.o.g. for a dual game where the player who chooses the smaller number wins. Its matrix  $\mathbf{S} = \mathbf{L}^\top$  has a unique NE consisting of the first row and column.

Let  $k > 1$ . Consider the  $k$ -th row  $\mathbf{x}$  of  $\mathbf{S}$ , i.e.,  $x_i = 1$  for  $i < k$ ,  $x_k = 0$ , and  $x_i = -1$  for  $i > k$ . The following crucial lemma shows that  $\widetilde{\text{BR}}_c(\mathbf{e}_k)$  with uniform perturbations selects a best response uniformly among  $\{1, \dots, k-1\}$ . Analogous lemma can be proven for the  $k$ -th column.

**Lemma 4.** Let  $\mathbf{x} = \langle 1, \dots, 1, 0, -1, \dots, -1 \rangle \in \mathbb{R}^n$  be the  $k$ -th row of  $\mathbf{S}$ ,  $k > 1$ , and  $\mathbf{z} = \langle z_1, \dots, z_n \rangle$  a random vector whose components  $z_i \sim \mathcal{U}(-1/2, 1/2)$ . Let  $I$  be the random variable  $I = \operatorname{argmax}_{i \in [n]} \pi_i(\mathbf{x} + \mathbf{z})$ . Then  $E[I] = \frac{k}{2}$ .

*Proof.* Note that  $x_i + z_i \sim \mathcal{U}(1/2, 3/2)$  for  $i < k$ ,  $x_i + z_i \sim \mathcal{U}(-1/2, 1/2)$  for  $i = k$ , and  $x_i + z_i \sim \mathcal{U}(-3/2, -1/2)$  for  $i > k$ . Thus we have  $E[I] = \sum_{i=1}^n i \cdot P(I = i) = \sum_{i=1}^{k-1} i \cdot P(I = i)$  because  $P(I = i) = 0$  for  $i \geq k$ . Moreover,  $P(I = i) = \frac{1}{k-1}$  is the uniform distribution on  $\{1, \dots, k-1\}$  as the first  $k-1$  components of  $\mathbf{x}$  are equal. Consequently,  $E(I) = \frac{1}{k-1} \sum_{i=1}^{k-1} i = \frac{k}{2}$ .  $\square$

Next, we need to understand equilibria of the subgame  $\mathbf{S}[R, C]$  in every iteration of SDO. For notational simplicity, we index rows and columns in  $\mathbf{S}[R, C]$  with the same indexes as in the original matrix  $\mathbf{S}$ . Further, we denote the probability simplices over  $R, C$  by  $\Delta_R, \Delta_C$ , respectively.

**Lemma 5.** Let  $R, C \subseteq [n]$  and  $r = \min R$ ,  $c = \min C$ .

1. If  $r < c$ , then  $\langle \mathbf{e}_i, \mathbf{q} \rangle$  is NE in  $\mathbf{S}[R, C]$  for any  $i \in \{k \in R \mid k < c\}$  and  $\mathbf{q} \in \Delta_C$ . There are no other NEs.
2. Dually, if  $c < r$ , then  $\langle \mathbf{p}, \mathbf{e}_j \rangle$  is NE in  $\mathbf{S}[R, C]$  for any  $j \in \{k \in C \mid k < r\}$  and  $\mathbf{p} \in \Delta_R$ . There are no other NEs.
3. If  $r = c$ , then  $\langle \mathbf{e}_r, \mathbf{e}_c \rangle$  is the unique NE in  $\mathbf{S}[R, C]$ .

*Proof.* For the first item, let  $\mathbf{x}$  be an  $i$ -th row of  $\mathbf{S}$  for  $i \in \{k \in R \mid i < c\}$ . Note that  $x_j = -1$  for all  $j \in C$ . Thus the  $i$ -th row in  $\mathbf{S}[R, C]$  is equilibrial strategy for the row player. A column player can play an arbitrary strategy. The second item is proven analogously as the first one.

The last claim follows if we note that the  $r$ -th row in  $\mathbf{S}[R, C]$  is of the form  $\langle 0, -1, \dots, -1 \rangle$  and  $c$ -th column in  $\mathbf{S}[R, C]$  is of the form  $\langle 0, 1, \dots, 1 \rangle$ .  $\square$

**Theorem 6.** SDO finds NE of  $\mathbf{S}$  in  $O(\log n)$  expected number of iterations.

*Proof.* We assume  $k, l = n$ , which is the worst possible start for SDO on  $\mathbf{S}$ . Let  $R_t, C_t$  be the row and column strategies in the iteration  $t$ . Define  $r_t = \min R_t$  and  $c_t = \min C_t$ . Next, we define  $X_t = \max\{r_t, c_t\}$ . Apparently, we have  $X_{t+1} \leq X_t$  since  $r_{t+1} \leq r_t$  and  $c_{t+1} \leq c_t$ . Once  $X_t = 1$ , SDO will find the NE in the next iteration. Thus, we want to show  $E[T] \in O(\log n)$  where  $T = \inf\{t \mid X_t \leq 1\}$ . This can be proven by (Kötzing and Krejca 2019, Corollary 17), if we show that  $X_t - E[X_{t+1} \mid X_1, \dots, X_t] \geq \delta X_t$  for some  $\delta > 0$ . Regarding the sequence  $X_1, X_2, \dots$ , we will consider only subsequence  $X_1, X_3, \dots$  of odd iterations  $t$  as SDO needs often two iterations to improve strategies for both players. In other words, we will prove  $X_t - E[X_{t+2} \mid X_1, X_3, \dots, X_t] \geq \delta X_t$  for some  $\delta > 0$ . To simplify the notation, we denote the conditional expectations without the condition in the rest of proof. For instance,  $E[X_{t+2}]$  means  $E[X_{t+2} \mid X_1, X_3, \dots, X_t]$ .

Let  $t$  be an odd iteration. There are two cases. We prove the case if  $r_t \leq c_t$ . The proof for the case when  $r_t \geq c_t$  is analogous. If  $r_t \leq c_t$ , then  $\langle \mathbf{e}_i, \mathbf{q} \rangle$  is NE for  $\mathbf{S}[R_t, C_t]$  for some  $i \leq c_t$  and some  $\mathbf{q} \in \Delta_C$  by Lemma 5 (items

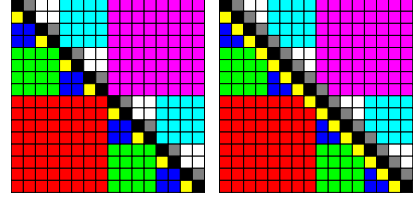


Figure 3: The clusters corresponding to the terminal states.

1 and 3). Consequently, the best response for the column player is  $j = \widetilde{\text{BR}}_c(\mathbf{e}_i) < i \leq c_t$ . Thus  $j$  becomes the new minimum of  $C_{t+1}$ , i.e.,  $c_{t+1} = j$ . Applying Lemma 4, we get  $E[c_{t+1}] = i/2 \leq c_t/2$ .

In the iteration  $t + 1$ , we either have  $r_{t+1} \leq c_{t+1}$  or  $r_{t+1} > c_{t+1}$ . In the first case,  $X_{t+2} \leq X_{t+1} = c_{t+1}$ . Thus  $E[X_{t+2}] \leq E[c_{t+1}] \leq c_t/2 \leq \max\{r_t, c_t\}/2 = X_t/2$ . If  $r_{t+1} > c_{t+1}$ , then  $\langle \mathbf{p}, \mathbf{e}_j \rangle$  for some  $j < r_{t+1}$  and any  $\mathbf{p} \in \Delta_R$  by Lemma 5. Consequently, the row player's best response  $i = \widetilde{\text{BR}}_r(\mathbf{e}_j) < j < r_{t+1}$  becomes the minimum of  $R_{t+2}$ , i.e.,  $r_{t+2} = i$ . Applying the column analog of Lemma 4, we get  $E[r_{t+2}] = j/2 \leq r_{t+1}/2$ . As  $c_{t+2} \leq c_{t+1}$ , we have  $X_{t+2} \leq \max\{r_{t+2}, c_{t+1}\}$ . Consequently,  $E[X_{t+2}] \leq \max\{r_{t+1}/2, c_t/2\} \leq \max\{r_t/2, c_t/2\} \leq X_t/2$ .

To sum up, we have  $X_t - E[X_{t+2}] \geq X_t/2$ . (Kötzing and Krejca 2019, Corollary 17) implies  $E[T] \leq 2 \ln n$  for the subsequence of odd iterations. Thus SDO needs at most  $1 + 4 \ln n$  iterations to find the NE in expectation.  $\square$

**Example 2.** The second example from (Zhang and Sandholm 2024) is a modification of the game from Example 1 where the best responses are unique. We will again use its dual variant. Its matrix  $\mathbf{U} \in \mathbb{R}^{n \times n}$  is defined as follows:

$$\mathbf{U} = \begin{pmatrix} 0 & -2 & -1 & \cdots & -1 \\ 2 & 0 & -2 & \cdots & -1 \\ 1 & 2 & 0 & \cdots & -1 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & 1 & 1 & \cdots & 0 \end{pmatrix}$$

Again,  $\mathbf{U}$  has unique  $\varepsilon$ -NE for  $0 \leq \varepsilon < 1$  consisting of the first row and column. DO needs  $n$  iterations to find the NE, if initialized with  $k = l = n$ .

The following theorem can be proven analogously as Theorem 6. The main difference is a modification of Lemma 4; Lemma 8 below (its proof is in the extended version of the paper). It shows that the perturbed best responses w.r.t. a single row  $\mathbf{e}_k$  are not selected uniformly from  $\{1, \dots, k-1\}$  as was the case for the game  $\mathbf{S}$ . The distribution is due to the unique best responses shifted towards larger indexes. However, the expected index is less than  $3/4k$  which is sufficient to apply (Kötzing and Krejca 2019, Corollary 17).

**Theorem 7.** SDO finds NE of  $\mathbf{U}$  in  $O(\log n)$  expected number of iterations.

**Lemma 8.** Let  $\mathbf{x} = \langle 1, \dots, 1, 2, 0, -2, -1, \dots, -1 \rangle \in \mathbb{R}^n$  be the  $k$ -th row of  $\mathbf{U}$ ,  $k > 1$ , and  $\mathbf{z} = \langle z_1, \dots, z_n \rangle$  a

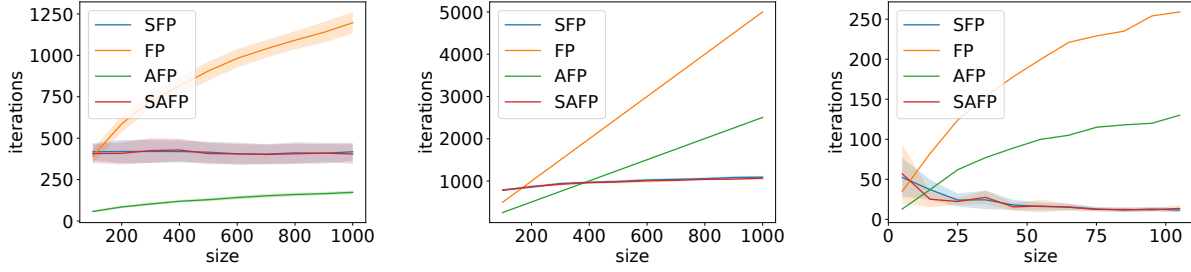


Figure 4: SFP iterations on random  $[0, 1]$ -matrix games (left),  $U^\top$  (middle), and the  $f$ -finger Morra game (right) for  $\varepsilon = 0.1$ .

random vector whose components  $z_i \sim \mathcal{U}(-1, 1)$ . Let  $I$  be the random variable  $I = \operatorname{argmax}_{i \in [n]} \pi_i(\mathbf{x} + \mathbf{z})$ . Then  $E[I] \leq 3/4k$ .

### Efficient Perturbations

Let  $M \in \mathbb{R}^{m \times n}$  be a matrix game and  $\mathbf{q} \in \Delta_n$ . Suppose we have a best-response oracle computing  $\operatorname{BR}_r$ . To implement a perturbed best response  $\widetilde{\operatorname{BR}}_r(\mathbf{q})$  using  $\operatorname{BR}_r$ , it requires perturbing every component of  $M\mathbf{q}$ . That cannot be done if SDO maintains only a fraction of  $M$ . Nevertheless, suppose  $M$  represents a game with inner structure like an extensive-form game (EFG) or a partially-observable stochastic game (POSG). In that case, we can implement an efficient PBRO using a best-response oracle and perturbation of utilities/rewards for terminal states or transitions.

Consider the (fully observable) stochastic game from Figure 1. In this game, each player chooses a sequence of four bits that corresponds to the player’s actions taken in non-terminal (circular) states. The terminal (squared) states contain their corresponding reward. This game and its generalization for any number of bits  $n$  was defined in (Zhang and Sandholm 2024). Considering the  $n$ -bit game, the equivalent matrix game is of size  $2^n \times 2^n$  and has the form of the matrix  $L$  from Example 1. To perturb best responses efficiently in this game, one can perturb only the rewards of the terminal states before applying the standard best-response oracle.

We experimentally tested that this kind of efficient perturbation speeds up the convergence of SDO for the stochastic game from Figure 1. Furthermore, we did the same for the POSG from Figure 2 also coming from (Zhang and Sandholm 2024). In this POSG, the players cannot observe the state of the game. So they again choose only four bits. The  $n$ -bit generalization of this game has an equivalent matrix game of size  $2^n \times 2^n$  and has the form of matrix  $U^\top$ ; see Example 2. The experimental results are shown in the next section. Here, we explain the details of our implementation.

One possibility to test our hypothesis for the game in Figure 1 would be implementing the best response oracle as an MDP solver, for instance using the value iteration algorithm (VI). However, VI would eliminate the existence of multiple best responses as it computes the optimal decision for each non-terminal state. For example, if the second player chooses 0 as her first bit, any sequence starting with 1 is a

best response for the first player. However, VI would return 1, 1, 1, 1 as the only best response.

Thus, we tested the hypothesis on the corresponding matrix game  $L$ . Each matrix entry corresponds to a terminal state. However, this correspondence is not one-to-one, as several entries correspond to the same terminal state. So, it induces a clustering of entries in  $L$  according to the terminal states. The clustering for  $L$  is shown in Figure 3 (left). Each color represents a single cluster.

Now, perturbing the reward of a terminal state is the same as perturbing its corresponding cluster by the same random value. Let  $K$  be the number of clusters and  $\mathbf{z}$  a  $K$ -dimensional random vector whose components  $z_i$  are i.i.d. random variables. For each cluster  $k \in [K]$ , we denote the  $\{0, 1\}$ -matrix  $B_k \in \mathbb{R}^{m \times n}$  that masks the cluster  $k$ , i.e.,  $B_k(i, j) = 1$  iff the  $i, j$ -entry belongs to the cluster  $k$  and  $B_k(i, j) = 0$  otherwise. Finally, for a column player’s mixed strategy  $\mathbf{q} \in \Delta_n$ , we compute the perturbed best-response as follows:

$$\widetilde{\operatorname{BR}}_r(\mathbf{q}) = \operatorname{argmin}_{i \in [m]} \pi_i \left( \left( L + \sum_{k=1}^K z_k B_k \right) \mathbf{q} \right) \quad (2)$$

The best response for the column player is defined analogously. The same clustering can be applied also for the POSG in Figure 2. The corresponding clustering is shown in Figure 3 (right).

### Experiments

All our experiments were implemented in the programming language Python and experimentally evaluated on the CPU AMD Ryzen 7 PRO 7840U. All randomized algorithms were repeated ten times to report the mean and the standard deviation. The initialization of SFP and SDO was deterministic, with the algorithm starting at the worst possible indexes. For tie-breaking best responses, the one with the least index was taken. To generate pseudorandom values, we used the numpy library with the initial seed 1.

SFP was tested in combination with Gumbel perturbations  $G(0, \beta)$  where  $\beta$  was set according to Theorem 3. We ran SFP with  $\varepsilon = 0.1$  on random square  $[0, 1]$ -matrix games, the matrix games  $L$  (Example 1) and  $U^\top$  (Example 2), and  $f$ -finger Morra game (Good 1965) all normalized to the interval  $[0, 1]$ ; Figure 4. We compared SFP with FP, AFP, and

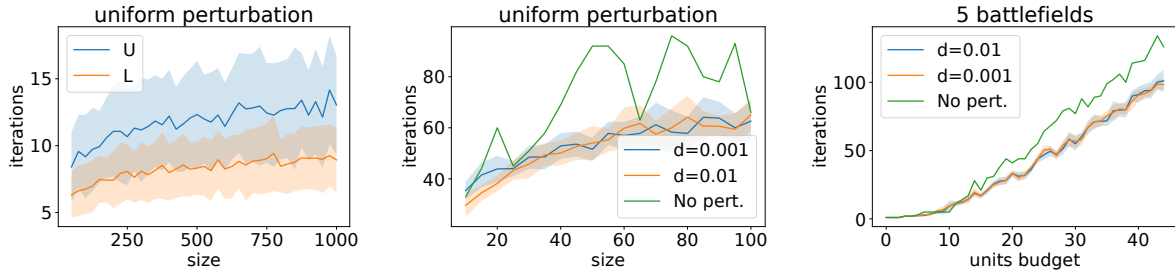


Figure 5: SDO iterations on  $L$  and  $U^\top$  with perturbations from  $\mathcal{U}(-1, 1)$  (left),  $f$ -finger Morra game (middle), and Colonel Blotto (right) with perturbations from  $\mathcal{U}(-d, d)$  for  $\varepsilon = 0.1$ .

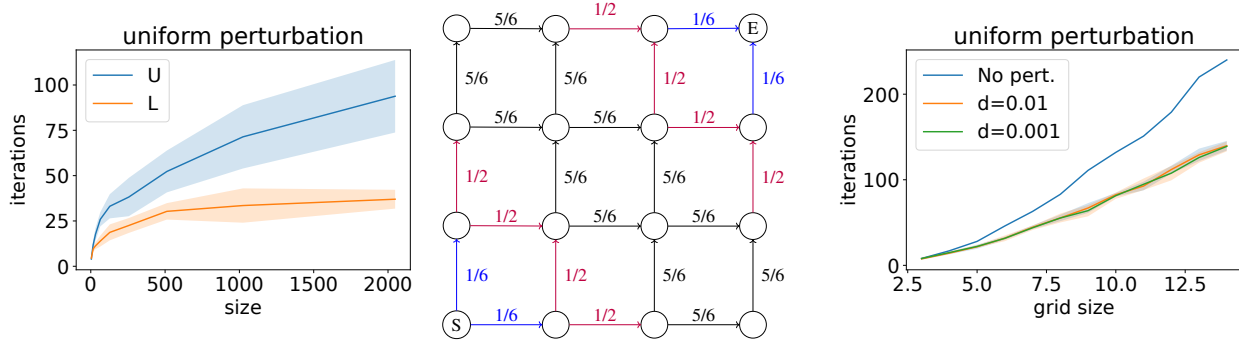


Figure 6: SDO iterations for  $L$  and  $U^\top$  with efficient perturbations from  $\mathcal{U}(-1, 1)$  (left),  $4 \times 4$ -grid with the transition costs for the path-planning game (middle), and the corresponding SDO iterations with uniform perturbations  $\mathcal{U}(-d, d)$  (right).

AFP with perturbed best responses (SAFP). The left graph shows the number of iterations on random  $n \times n$ -matrix games. In the experiment, we generated 100 random games for each  $n$ . AFP outperforms all the other methods; SFP and SAFP behave similarly on random games. The middle graph presents the number of iterations for the game  $U^\top$ . An analogous graph for  $L$  looks almost identical, so we omit it. FP and AFP show their the worst-case  $O(n)$  complexity. SFP and SAFP again have a similar performance. Finally, the right graph shows the number of iterations for the  $f$ -finger Morra game. The horizontal axis represents the number of fingers  $f$ . The size of the corresponding matrix is  $f^2 \times f^2$ . Interestingly, SFP and SAFP are able to quickly find an  $\varepsilon$ -NE for large values of  $f$ . However, this is not the case for FP and AFP.

We did similar experiments for DO and SDO combined with the uniform perturbations. The SDO iterations for matrix games  $L$ ,  $U^\top$ ,  $f$ -finger Morra game, and the Colonel Blotto game with five battlefields for different numbers of units are shown in Figure 5. We omit the DO iterations in the left graph as the DO needs  $n$  iterations for the size  $n$ . In all these cases, SDO outperforms DO. However, we noticed this is not the case for random games where perturbations do not provide a faster convergence or even degrades the convergence rate if the perturbations are too large. This outcome is somewhat expected, considering that the best responses in random games are inherently stochastic.

Further, we tested the efficient perturbations of clusters,

see Equation (2), on the matrix games  $L$  and  $U^\top$  of size  $2^n \times 2^n$  corresponding to stochastic games from Figure 1 and 2. The results are shown in Figure 6 (left). Although SDO with efficient perturbations requires more iterations in comparison with Figure 5 (left), still the obtained results are much better than the linear complexity of DO.

To test the efficient perturbations further, we defined a path-planning game on an  $n \times n$ -grid. An example of such a grid for  $n = 4$  is shown in Figure 6 (middle). The colors and numbers denote the costs for particular transitions. The path-planning player looks for the shortest path in the grid starting in the node S and finishing in the node E. The other player selects a single edge and multiplies its cost by a given coefficient; 10 in our experiments. We implemented the perturbed best-response oracle by perturbing the transition costs every time the oracle is called. Figure 6 (right) show that uniform perturbations reduce the number of iterations with increasing grid size.

## Conclusions

To summarize, perturbing best responses improves the convergence of FP and DO. We proved that SFP has a logarithmic complexity in the number of pure strategies  $n$ . Although DO needs  $\Theta(n)$  many iterations to find  $\varepsilon$ -NE in the worst case, it is still possible that SDO finds an  $\varepsilon$ -NE in  $O(\log n)$  expected number of iterations. We leave this open problem for future research.

## Acknowledgments

The authors were supported by the Czech Science Foundation grant—no. 24-12046S.

## References

- Abernethy, J. D.; Lai, K. A.; and Wibisono, A. 2021. Fast Convergence of Fictitious Play for Diagonal Payoff Matrices. In Marx, D., ed., *Proceedings of the 2021 ACM-SIAM Symposium on Discrete Algorithms, SODA 2021, Virtual Conference, January 10 - 13, 2021*, 1387–1404. SIAM.
- Abernethy, J. D.; Lee, C.; Sinha, A.; and Tewari, A. 2014. Online Linear Optimization via Smoothing. In Balcan, M.; Feldman, V.; and Szepesvári, C., eds., *Proceedings of The 27th Conference on Learning Theory, COLT 2014, Barcelona, Spain, June 13-15, 2014*, volume 35 of *JMLR Workshop and Conference Proceedings*, 807–823. JMLR.org.
- Althöfer, I. 1994. On sparse approximations to randomized strategies and convex combinations. *Linear Algebra and its Applications*, 199: 339–355.
- Bighashdel, A.; Wang, Y.; McAleer, S.; Savani, R.; and Oliehoek, F. A. 2024. Policy Space Response Oracles: A Survey. In *Proceedings of the Thirty-Third International Joint Conference on Artificial Intelligence, IJCAI 2024, Jeju, South Korea, August 3-9, 2024*, 7951–7961. ijcai.org.
- Brown, G. W. 1951. Iterative Solution of Games by Fictitious Play. In Koopmans, T. C., ed., *Activity Analysis of Production and Allocation*. New York: Wiley.
- Cen, S.; Wei, Y.; and Chi, Y. 2024. Fast Policy Extragradient Methods for Competitive Games with Entropy Regularization. *J. Mach. Learn. Res.*, 25: 4:1–4:48.
- Cesa-Bianchi, N.; and Lugosi, G. 2006. *Prediction, learning, and games*. Cambridge University Press.
- Cloud, A.; Wang, A.; and Kerr, W. 2023. Anticipatory Fictitious Play. In *Proceedings of the Thirty-Second International Joint Conference on Artificial Intelligence, IJCAI 2023, 19th-25th August 2023, Macao, SAR, China*, 73–81. ijcai.org.
- Daskalakis, C.; and Pan, Q. 2014. A Counter-example to Karlin’s Strong Conjecture for Fictitious Play. In *55th IEEE Annual Symposium on Foundations of Computer Science, FOCS 2014, Philadelphia, PA, USA, October 18-21, 2014*, 11–20. IEEE Computer Society.
- Franke, M.; and Degen, J. 2023. The softmax function: Properties, motivation, and interpretation. <https://doi.org/10.31234/osf.io/vsw47>.
- Fudenberg, D.; and Kreps, D. M. 1993. Learning Mixed Equilibria. *Games and Economic Behavior*, 5(3): 320–367.
- Good, R. A. 1965. f-Finger Morra. *SIAM Review*, 7(1): 81–87.
- Goodfellow, I.; Bengio, Y.; and Courville, A. 2016. *Deep Learning*. MIT Press.
- Gumbel, E. J. 1954. *Statistical Theory of Extreme Values and Some Practical Applications: A Series of Lectures*, volume 33 of *Applied mathematics series*. U.S. Government Printing Office.
- Hazan, E.; and Koren, T. 2016. The computational power of optimization in online learning. In Wicks, D.; and Mansour, Y., eds., *Proceedings of the 48th Annual ACM SIGACT Symposium on Theory of Computing, STOC 2016, Cambridge, MA, USA, June 18-21, 2016*, 128–141. ACM.
- Hofbauer, J.; and Sandholm, W. H. 2002. On the Global Convergence of Stochastic Fictitious Play. *Econometrica*, 70(6): 2265–2294.
- Karlin, S. 1959. *Mathematical methods and theory in games, programming, and economics*. Reading, U. S.
- Kötzing, T.; and Krejca, M. S. 2019. First-hitting times under drift. *Theor. Comput. Sci.*, 796: 51–69.
- Lanctot, M.; Zambaldi, V. F.; Gruslys, A.; Lazaridou, A.; Tuyls, K.; Pérolat, J.; Silver, D.; and Graepel, T. 2017. A Unified Game-Theoretic Approach to Multiagent Reinforcement Learning. In Guyon, I.; von Luxburg, U.; Bengio, S.; Wallach, H. M.; Fergus, R.; Vishwanathan, S. V. N.; and Garnett, R., eds., *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*, 4190–4203.
- Lipton, R. J.; and Young, N. E. 1994. Simple strategies for large zero-sum games with applications to complexity theory. In Leighton, F. T.; and Goodrich, M. T., eds., *Proceedings of the Twenty-Sixth Annual ACM Symposium on Theory of Computing, 23-25 May 1994, Montréal, Québec, Canada*, 734–740. ACM.
- McAleer, S.; Lanier, J.; Wang, K.; Baldi, P.; Fox, R.; and Sandholm, T. 2022. Self-Play PSRO: Toward Optimal Populations in Two-Player Zero-Sum Games. arXiv:2207.06541.
- McMahan, H. B.; Gordon, G. J.; and Blum, A. 2003. Planning in the Presence of Cost Functions Controlled by an Adversary. In Fawcett, T.; and Mishra, N., eds., *Machine Learning, Proceedings of the Twentieth International Conference (ICML 2003), August 21-24, 2003, Washington, DC, USA*, 536–543. AAAI Press.
- Robinson, J. 1951. An Iterative Method of Solving a Game. *Annals of Mathematics*, 54(2): 296–301.
- Train, K. E. 2009. *Discrete Choice Methods with Simulation*. Cambridge University Press.
- Zhang, B. H.; and Sandholm, T. 2024. Exponential Lower Bounds on the Double Oracle Algorithm in Zero-Sum Games. In *Proceedings of the Thirty-Third International Joint Conference on Artificial Intelligence, IJCAI 2024, Jeju, South Korea, August 3-9, 2024*, 3032–3039. ijcai.org.