

---

# Convergence of $\log(1/\epsilon)$ for Gradient-Based Algorithms in Zero-Sum Games without the Condition Number: A Smoothed Analysis

---

**Ioannis Anagnostides**  
Carnegie Mellon University  
ianagnos@cs.cmu.edu

**Tuomas Sandholm**  
Carnegie Mellon University  
Strategic Machine, Inc.  
Strategy Robot, Inc.  
Optimized Markets, Inc.  
sandholm@cs.cmu.edu

## Abstract

Gradient-based algorithms have shown great promise in solving large (two-player) zero-sum games. However, their success has been mostly confined to the low-precision regime since the number of iterations grows polynomially in  $1/\epsilon$ , where  $\epsilon > 0$  is the duality gap. While it has been well-documented that linear convergence—an iteration complexity scaling as  $\log(1/\epsilon)$ —can be attained even with gradient-based algorithms, that comes at the cost of introducing a dependency on certain condition number-like quantities which can be exponentially large in the description of the game.

To address this shortcoming, we examine the iteration complexity of several gradient-based algorithms in the celebrated framework of *smoothed analysis*, and we show that they have *polynomial smoothed complexity*, in that their number of iterations grows as a polynomial in the dimensions of the game,  $\log(1/\epsilon)$ , and  $1/\sigma$ , where  $\sigma$  measures the magnitude of the smoothing perturbation. Our result applies to optimistic gradient and extra-gradient descent/ascent, as well as a certain iterative variant of Nesterov’s smoothing technique. From a technical standpoint, the proof proceeds by characterizing and performing a smoothed analysis of a certain *error bound*, the key ingredient driving linear convergence in zero-sum games. En route, our characterization also makes a natural connection between the convergence rate of such algorithms and perturbation-stability properties of the equilibrium, which is of interest beyond the model of smoothed complexity.

## 1 Introduction

We consider the fundamental problem of computing an *equilibrium* strategy for a (two-player) zero-sum game

$$\min_{\mathbf{x} \in \Delta^n} \max_{\mathbf{y} \in \Delta^m} \langle \mathbf{x}, \mathbf{A} \mathbf{y} \rangle, \quad (1)$$

where  $\Delta^{d+1} := \{\mathbf{x} \in \mathbb{R}_{\geq 0}^{d+1} : \mathbf{x}^\top \mathbf{1}_{d+1} = 1\}$  represents the  $d$ -dimensional probability simplex and  $\mathbf{A} \in \mathbb{R}^{n \times m}$  is the payoff matrix of the game. Tracing all the way back to Von Neumann’s celebrated minimax theorem [von Neumann, 1928], zero-sum games played a pivotal role in the early development of game theory [von Neumann and Morgenstern, 1947] and the crystallization of linear programming duality [Dantzig, 1951]. Indeed, in light of the equivalence between zero-sum games and linear programming [Adler, 2013, von Stengel, 2023, Brooks and Reny, 2023], many central optimization problems can be cast as (1).

State of the art algorithms for solving zero-sum games can be coarsely classified based on the desired accuracy of a feasible solution  $(\mathbf{x}, \mathbf{y})$ , measured in terms of the *duality gap*

$$\Phi(\mathbf{x}, \mathbf{y}) := \max_{\mathbf{y}' \in \Delta^m} \langle \mathbf{x}, \mathbf{A}\mathbf{y}' \rangle - \min_{\mathbf{x}' \in \Delta^n} \langle \mathbf{x}', \mathbf{A}\mathbf{y} \rangle. \quad (2)$$

In the so-called low-precision regime, where one is content with a crude solution  $(\mathbf{x}^*, \mathbf{y}^*)$  such that  $\Phi(\mathbf{x}^*, \mathbf{y}^*) =: \epsilon \gg 0$ , the best available algorithms typically revolve around the framework of *regret minimization*, both in practice [Farina et al., 2021, Brown and Sandholm, 2019, Zinkevich et al., 2007, Tang et al., 2023] and in theory [Carmon et al., 2020, 2019, 2024, Grigoriadis and Khachiyan, 1995, Clarkson et al., 2012, Alacaoglu and Malitsky, 2022]—in conjunction with other techniques to speed up the per-iteration complexity, such as variance reduction, data structure design, and sparsification [Zhang and Sandholm, 2020, Farina and Sandholm, 2022]. Such algorithms have been central to landmark results in practical computation of equilibrium strategies even in enormous games [Brown and Sandholm, 2018, Bowling et al., 2015, Moravčík et al., 2017, Perolat et al., 2022].

The high-precision regime, where  $\epsilon \ll \frac{1}{\text{poly}(nm)}$ , has turned out to be more elusive, with current LP-based techniques struggling to scale favorably in large instances. This deficiency can be in part attributed to the relatively high per-iteration complexity of LP-based approaches, such as interior-point methods or the ellipsoid algorithm, as well as their intense memory requirements. A promising antidote is to instead rely on iterative gradient-based methods that have a minimal per-iteration cost. Indeed, in a line of work pioneered by Tseng [1995], it is known well-documented that *linear convergence*—an iteration complexity scaling only as  $\log(1/\epsilon)$ —can be achieved even with such methods [Tseng, 1995, Gilpin et al., 2012, Wei et al., 2021, Applegate et al., 2023, Fercoq, 2023]. There is, however, a major caveat to those results: the number of iterations no longer grows polynomially with the dimensions of the game  $n$  and  $m$ , but instead depends on certain condition number-like quantities that could be exponentially large in the description of the problem; it is thus unclear how to interpret those results from a computational standpoint.

To address those shortcomings, in this paper we work in the celebrated framework of *smoothed analysis* pioneered by Spielman and Teng [2004]. Namely, our goal is to characterize the iteration complexity of certain gradient-based algorithms in zero-sum games when the payoff matrix  $\mathbf{A}$  is subjected to small but random perturbations, as formally introduced below.

**Definition 1.1** (Zero-sum games under Gaussian perturbations). Let  $\bar{\mathbf{A}} \in [-1, 1]^{n \times m}$ . We assume that the payoff matrix is given by  $\mathbf{A} := \bar{\mathbf{A}} + \mathbf{G}$ , where each entry of  $\mathbf{G}$  is an independent (univariate) Gaussian random variable with zero mean and variance  $\sigma^2 \leq 1$ .

Randomness here is only injected into the payoff matrix and not the set of constraints (that is, the probability simplex), which is the natural model; after applying the perturbation, the problem should still be a zero-sum game in the form of (1). Under this model, we investigate the convergence of the following gradient-based algorithms.<sup>1</sup> (Their formal description is given later in [Appendix B](#).)

1. *optimistic gradient descent/ascent (OGDA)* [Popov, 1980];
2. *optimistic multiplicative weights update (OMWU)* [Syrgkanis et al., 2015, Chiang et al., 2012, Rakhlin and Sridharan, 2013];
3. *extra-gradient descent/ascent (EGDA)* [Korpelevich, 1976]; and
4. an iterative variant of Nesterov’s *smoothing technique (IterSmooth)* [Gilpin et al., 2012, Nesterov, 2005].

Smoothed complexity allows interpolating between worst-case analysis—when the variance of the noise  $\sigma^2$  is negligible—and average-case analysis—when the noise dominates over the underlying input. An average-case analysis is often unreliable since—as Edelman [1993] convincingly argued—a fully random matrix does not necessarily capture typical instances encountered in practice. Spielman and Teng [2004] put forward the framework of smoothed analysis as an attempt to explain the performance of algorithms in realistic scenarios; to understand how brittle worst-case instances really are. They famously proved that the simplex algorithm, under a certain pivoting rule, enjoys *polynomial smoothed complexity*, meaning that its running time is bounded by some polynomial in the

<sup>1</sup>The vanilla gradient descent/ascent algorithm does not even converge (in a last-iterate sense) in zero-sum games (e.g., [Mertikopoulos et al., 2018]), which is why our analysis revolves around certain variants thereof. It is worth noting that regret minimization techniques provide guarantees concerning the average iterates, a distinction blurred in our introduction.

size of the input and  $1/\sigma$ . Smoothed analysis is by now a well-accepted algorithmic framework with a tremendous impact in the analysis of algorithms. We also argue that it is particularly well-motivated from a game-theoretic perspective: there is often misspecification or noise when modeling a game, so smoothed analysis offers a compelling way of bypassing pathological instances that are perhaps artificial in the first place.

Nevertheless, we are not aware of any prior work operating in the smoothed complexity model per [Definition 1.1](#) in the context of zero-sum games. To clarify this point, it is important to stress here that although zero-sum games can be immediately reduced to linear programs, that reduction is less clear in the smoothed complexity model. In particular, one set of constraints in the induced linear program takes the form  $\mathbf{A}\mathbf{y} \leq v\mathbf{1}_n =: \mathbf{b}$ , where  $\mathbf{1}_n \in \mathbb{R}^n$  is the all-ones vector. According to the usual model of smoothed complexity in the context of linear programs, randomness has to be injected into both  $\mathbf{A}$  and  $\mathbf{b}$ , but that clearly disturbs the validity of the equivalence. More broadly, reductions in the smoothed complexity model are quite delicate [[Bläser and Manthey, 2015](#)]; as a further example, even reductions involving solely linear transformations can break in the smoothed complexity model since independence—a crucial assumption in this framework—is not guaranteed to carry over. Relatedly, one interesting direction arising from the work of [Spielman and Teng \[2003\]](#) is to perform smoothed analysis in linear programs which are guaranteed to be feasible and bounded, no matter the perturbation; zero-sum games under [Definition 1.1](#) constitute such a class. Besides the point above, different algorithms designed for the same problem can have entirely different properties, not least in terms of their smoothed complexity. The class of algorithms we consider in this paper is quite distinct from the ones shown to have polynomial smoothed complexity in the context of linear programs (described further in [Appendix A](#)). In many ways, gradient-based methods are simpler and more natural, which partly justifies their tremendous practical use. As a result, understanding their smoothed complexity is an important question.

## 1.1 Our results

Our main contribution is to show that, with the exception of [OMWU](#), the other gradient-based algorithms mentioned above ([Items 1, 3 and 4](#)) have polynomial smoothed complexity with high probability—that is to say, with probability at least  $1 - \frac{1}{\text{poly}(nm)}$ .

**Theorem 1.2.** *With high probability over the randomness of  $\mathbf{A} \in \mathbb{R}^{n \times m}$  ([Definition 1.1](#)), [OGDA](#), [EGDA](#) and [IterSmooth](#) converge to an  $\epsilon$ -equilibrium after  $\text{poly}(n, m, 1/\sigma) \cdot \log(1/\epsilon)$  iterations.*

The main takeaway of this result is that, modulo pathological instances, certain gradient-based algorithms are reliable solvers in zero-sum games even in the high-precision regime. Similarly to earlier endeavors in the context of linear programs [[Spielman and Teng, 2004](#), [Blum and Dunagan, 2002](#)], a dependency of  $\text{poly}(1/\sigma)$  (as in [Theorem 1.2](#)) is what we should expect; the one exception is the class of interior-point methods whose running time grows as  $\log(1/\sigma)$ , but those algorithms are (weakly) polynomial even in the worst case. We further remark that the polynomial dependency on  $n$  and  $m$  in [Theorem 1.2](#) can almost certainly be improved, and we made no effort to optimize it.

Regarding [OMWU](#), which is not covered by [Theorem 1.2](#), we also obtain a significant improvement in the iteration complexity compared to the worst-case analysis of [Wei et al. \[2021\]](#), but our bound is still not polynomial. As we explain further in [Appendix C.3](#), the main difficulty pertaining to [OMWU](#) is that the analysis of [Wei et al. \[2021\]](#) gives (at best) an exponential bound *no matter the geometry of the problem*. With that mind, our result is essentially the best one could hope for without refining the worst-case analysis of [OMWU](#), which is not within our scope here. We anticipate that our characterization herein will prove useful in conjunction with future developments in the worst-case complexity of [OMWU](#), as well as in the analysis of other iterative methods.

**The error bound** The central ingredient that enables gradient-based algorithms to exhibit linear convergence is a certain *error bound*, given below as [Definition 1.3](#). For compactness in our notation, we let  $\mathcal{X} := \Delta^n$  and  $\mathcal{Y} := \Delta^m$ . We then let  $\mathbf{z} := (\mathbf{x}, \mathbf{y})$ ,  $\mathcal{Z} := \mathcal{X} \times \mathcal{Y}$ , and  $\mathcal{Z}^* := \mathcal{X}^* \times \mathcal{Y}^*$ , where  $\mathcal{X}^*$  and  $\mathcal{Y}^*$  represent the (convex) set of equilibria for Player  $x$  and Player  $y$ , respectively.

**Definition 1.3** (Error bound). Let  $\Phi(\mathbf{z})$  denote the duality gap as introduced in [\(2\)](#). We say that the zero-sum game [\(1\)](#) satisfies an *error bound* with modulus  $\kappa \in \mathbb{R}_{>0}$  if

$$\Phi(\mathbf{z}) \geq \kappa \|\mathbf{z} - \Pi_{\mathcal{Z}^*}(\mathbf{z})\| \quad \forall \mathbf{z} \in \mathcal{Z}. \quad (3)$$

Above,  $\Pi_{\mathcal{Z}^*}(\cdot)$  denotes the (Euclidean) projection operator; the set of games with a unique equilibrium has measure one, so we can safely replace  $\Pi_{\mathcal{Z}^*}(z)$  by the unique equilibrium  $z^* \in \mathcal{Z}^*$ . It has been known at least since the work of Tseng [1995] that affine variational inequalities indeed satisfy (3). Nevertheless, it should come to no surprise that, even in  $3 \times 3$  games,  $\kappa$  can be arbitrarily small (Proposition 3.1), which in turn means that, linear convergence notwithstanding, the number of iterations prescribed by an analysis revolving around (3) can be arbitrarily large. In fact, with the exception of OMWU, which is to be discussed further below, Definition 1.3 suffices to establish linear convergence (essentially) based on existing results.<sup>2</sup> Our main result pertaining to Definition 1.3 is that the modulus  $\kappa$  is likely to be polynomial in the smoothed complexity model:

**Theorem 1.4.** *With high probability over the randomness of  $\mathbf{A}$  (Definition 1.1), the error bound per Definition 1.3 is satisfied for any sufficiently small  $\kappa \geq \text{poly}(\sigma, 1/(nm))$ .*

To establish this result, the first step is to lower bound  $\kappa$  in terms of certain natural geometric features of the problem (Theorem 3.6), which is discussed further in Section 3.1. Establishing Theorem 1.4 then reduces to analyzing each of those quantities under Definition 1.1. It turns out that bounding those quantities also suffices for characterizing OMWU, whose existing analysis due to Wei et al. [2021] involves some further ingredients besides the error bound of Definition 1.3.

**Further implications** Our characterization of the error bound given in Theorem 3.6 has some further important implications. First, a well-known vexing issue regarding computing equilibria even in zero-sum games is that a solution with small duality gap can still be relatively far from the equilibrium in the geometric sense, a phenomenon further exacerbated in multi-player games [Etessami and Yannakakis, 2007]. Therefore, results providing guarantees in terms of the duality gap are not particularly informative when it comes to computing strategies close to the equilibrium in a geometric sense. At the same time, there are ample reasons why the latter guarantee is more appealing [Etessami and Yannakakis, 2007]. Theorem 1.4 implies that such concerns can be alleviated in the smoothed complexity model:

**Corollary 1.5.** *With high probability over the randomness of  $\mathbf{A}$  (Definition 1.1), any point  $z \in \mathcal{Z}$  with  $\Phi(z) \leq \epsilon$  satisfies  $\|z - z^*\| \leq \epsilon \cdot \text{poly}(n, m, 1/\sigma)$ .*

Beyond smoothed analysis, Theorem 3.6 applies to any non-degenerate game (Definition 3.2), and can be thereby used to parameterize the rate of convergence of gradient-based algorithms based on natural and interpretable game-theoretic quantities of the underlying game, which has eluded prior work. In particular, we make a natural connection between the complexity of gradient-based algorithms and *perturbation stability* properties of the equilibrium. In light of misspecifications which are often present in game-theoretic modeling, focusing on games with perturbation-stable equilibria is well-motivated and has already received ample of interest in prior work [Balcan and Braverman, 2017, Awasthi et al., 2010]; more broadly, perturbation stability is a common assumption in the analysis of algorithms beyond the worst-case model [Makarychev and Makarychev, 2021]. There are different natural ways of defining perturbation-stable games; here, we assume that any perturbation with magnitude below  $\delta > 0$ , in that  $\|\mathbf{A}' - \mathbf{A}\|_2 \leq \delta$ , maintains the support of the equilibrium and the non-degeneracy of the game; we call such games  $\delta$ -*support-stable* (Definition 4.1). In this context, we show the following result.

**Corollary 1.6.** *For any  $\delta$ -support-stable zero-sum game, OGDA, EGDA and IterSmooth converge to an  $\epsilon$ -equilibrium after  $\text{poly}(n, m, 1/\delta) \cdot \log(1/\epsilon)$  iterations.*

That is, games in which  $\delta$  is not too close to 0 are more amenable to gradient-based algorithms, which is a quite natural connection. Corollary 1.6 is shown by relating each of the quantities involved in Theorem 3.6 to parameter  $\delta$  defined above.

## 2 Notation

Before we proceed with our technical content, we first take the opportunity to streamline our notation; further background on smoothed analysis and a description of the algorithms referred to earlier (Items 1 to 4) is given later in Appendix B, as it is not important for the purpose of the main body.

---

<sup>2</sup>Definition 1.3 also readily establishes linear convergence for other compelling primal-dual algorithms, as shown recently by Applegate et al. [2023]; in that paper, the error bound was referred to as “sharpness,” a terminology employed in other papers as well (e.g., [Zarifis et al., 2024]).

We use boldface letters, such as  $\mathbf{x}, \mathbf{y}, \mathbf{b}, \mathbf{c}$ , to represent vectors in a Euclidean space. For a vector  $\mathbf{x} \in \mathbb{R}^n$ , we access its  $i$ th coordinate via a subscript, namely  $\mathbf{x}_i$ . Superscripts (together with parentheses) are typically reserved for the (discrete) time index. We denote by  $\|\mathbf{x}\|$  the Euclidean norm,  $\|\mathbf{x}\| := \sqrt{\sum_{i=1}^n x_i^2}$ , the  $\ell_\infty$  norm by  $\|\mathbf{x}\|_\infty := \max_{1 \leq i \leq n} |\mathbf{x}_i|$ , and the  $\ell_1$  norm by  $\|\mathbf{x}\|_1 := \sum_{i=1}^n |\mathbf{x}_i|$ . For  $\mathbf{x}, \mathbf{x}' \in \mathbb{R}^n$ , we let  $\text{dist}(\mathbf{x}, \mathbf{x}') := \|\mathbf{x} - \mathbf{x}'\|$ .  $\text{span}(\cdot)$  represents the linear space spanned by a given set of vectors. For  $\mathbf{x} \in \mathbb{R}^n$  and a subset  $B \subseteq [n]$ , we denote by  $\mathbf{x}_B \in \mathbb{R}^B$  the subvector of  $\mathbf{x}$  induced by  $B$ . We let  $\mathbf{1}_n \in \mathbb{R}^n$  be the all-ones vector of dimension  $n$ ; we will typically omit the subscript when it is clear from the context. For vectors  $\mathbf{x} \in \mathbb{R}^n$  and  $\mathbf{y} \in \mathbb{R}^m$ , we write  $(\mathbf{x}, \mathbf{y}) \in \mathbb{R}^{n+m}$  to denote their concatenation. Throughout this paper, we use  $\mathbf{x}$  and  $\mathbf{y}$  to denote the strategy of Player  $x$  and Player  $y$ , respectively.

To represent matrices, we use boldface capital letter, such as  $\mathbf{A}, \mathbf{Q}$ . It will sometimes be convenient to use  $\mathbf{A}^b \in \mathbb{R}^{nm}$  to represent a vectorization of  $\mathbf{A} \in \mathbb{R}^{n \times m}$ . We overload notation by letting  $\|\mathbf{A}\|$  be the spectral norm of  $\mathbf{A}$ . For a matrix  $\mathbf{A} \in \mathbb{R}^{n \times m}$  and subsets  $B \subseteq [n], N \subseteq [m]$ , we denote by  $\mathbf{A}_{B,N} \in \mathbb{R}^{B \times N}$  the submatrix of  $\mathbf{A}$  induced by  $B$  and  $N$ .  $\mathbf{A}_{i,:}$  and  $\mathbf{A}_{:,j}$  represent the  $i$ th row and  $j$ th column of  $\mathbf{A}$ , respectively. The singular values of a matrix  $\mathbf{M} \in \mathbb{R}^{d \times d}$  are denoted by  $\sigma_1(\mathbf{M}) \geq \sigma_2(\mathbf{M}) \geq \dots \geq \sigma_d(\mathbf{M}) \geq 0$  (not to be confused with our notation for the variance  $\sigma^2$ ). To be more explicit, we may also use  $\sigma_{\max}(\mathbf{M}) := \sigma_1(\mathbf{M})$  and  $\sigma_{\min}(\mathbf{M}) := \sigma_d(\mathbf{M})$ .

### 3 Smoothed analysis of the error bound

In this section, we perform a smoothed analysis of the error bound—as introduced earlier in [Definition 1.3](#)—in (two-player) zero-sum games. It is first instructive to point out why smoothed analysis is useful in the first place: the modulus  $\kappa$  can be arbitrarily close to 0 even when  $n = m = 3$  (that is,  $3 \times 3$  games); this is detrimental as the iteration complexity of algorithms such as OGDA grows as a polynomial in  $1/\kappa$ .

**Proposition 3.1.** *There exists a  $3 \times 3$  zero-sum game such that  $\kappa$  per [Definition 1.3](#) is arbitrarily close to 0.*

In proof, it is enough to consider the ill-conditioned diagonal matrix

$$\mathbf{A} = \begin{pmatrix} \gamma & 0 & 0 \\ 0 & 2\gamma & 0 \\ 0 & 0 & 1 \end{pmatrix}, \quad (4)$$

where  $0 < \gamma \ll 1$ . The (unique) equilibrium of (4) reads  $\mathbf{x}^* = \mathbf{y}^* = \frac{1}{3+2\gamma}(2, 1, 2\gamma) \in \Delta^3$ . Now, considering  $\mathbf{x} = (1, 0, 0)$  and  $\mathbf{y} = (0, 0, 1)$ , for the duality gap we have  $\Phi(\mathbf{x}, \mathbf{y}) = \gamma$ , while the distance of  $(\mathbf{x}, \mathbf{y})$  from the optimal solution  $(\mathbf{x}^*, \mathbf{y}^*)$  is at least  $\frac{3}{3+2\gamma}$ . In turn, by [Definition 1.3](#), this means that  $\kappa \leq 2\gamma$ . So, [Proposition 3.1](#) follows by taking  $\gamma \rightarrow 0$ .<sup>3</sup>

[Proposition 3.1](#) exposes one type of pathology that can decelerate gradient-based algorithms, which is evidently related to the poor spectral properties of the payoff matrix. This intuition is quite helpful when equilibria are fully supported—as is the case in (4)—but has to be significantly refined more broadly, as we formalize in the sequel.

To sidestep such pathological examples, we thus turn to the smoothed analysis framework of [Definition 1.1](#).

#### 3.1 Overview

The most natural approach to analyze the error bound in the smoothed complexity model is to rely on an existing (worst-case) analysis proving that a positive  $\kappa$  exists, and then attempt to refine that analysis. Yet, at least based on such prior results we are aware of, that turns out to be challenging. As an example, let us consider the recent analysis of [Wei et al. \[2021\]](#). As we explain in more detail in [Appendix C.3](#), [Wei et al. \[2021\]](#) relate the modulus  $\kappa$  of the error bound to the (inverse of the) norm of a solution to a certain feasible linear program; the existence of a legitimate  $\kappa > 0$  then follows readily from feasibility. Now, this reduction seems quite promising: [Renegar \[1994\]](#) has shown that the

<sup>3</sup>If we want to specify the game with a (finite) number of  $L$  bits, [Proposition 3.1](#) tells us that the modulus  $\kappa$  can be exponentially small in  $L$ .

norm of a solution to a linear program can be bounded in terms of its *condition number*—the distance to infeasibility in our case, and [Dunagan et al. \[2011\]](#) later proved that the condition number of linear programs is polynomial in the smoothed complexity model. Nevertheless, there are some difficulties in materializing that argument. First, the induced linear program involves terms depending on both the payoff matrix and the geometry of the constraints (the probability simplex in our case). Consequently, the analysis of [Dunagan et al. \[2011\]](#) does not carry over since randomness is only injected into the payoff matrix. The second and more important obstacle is that the induced linear program depends on the optimal solution, which in turn depends on the randomness of the payoff matrix; this significantly entangles the underlying distribution. As there are exponentially many possible configurations, we cannot afford to argue about each one separately and then apply the union bound. This difficulty is in fact known to be the crux in performing smoothed analysis [[Spielman and Teng, 2004](#)].<sup>4</sup>

To address those challenges, we provide a new characterization of the error bound in terms of some natural quantities of the underlying game ([Theorem 3.6](#)), which in some sense capture the difficulty of the problem. We are then able to use a technique due to [Spielman and Teng \[2004\]](#), exposed in [Section 3.3](#), to bound the probability that each of the involved quantities is close to 0 ([Propositions 3.8 to 3.10](#)), even though the underlying distribution is quite convoluted. The resulting analysis follows the one given by [Spielman and Teng \[2003\]](#) in the context of termination of linear programs, but still has to account for a number of structural differences.

In what follows, we structure our argument as follows. First, in [Section 3.2](#), we relate the modulus  $\kappa$  to some natural quantities capturing key geometric features of the problem. [Section 3.3](#) then proceed by analyzing those quantities in the smoothed analysis framework.

### 3.2 Characterization of the error bound

Our first goal is to characterize the error bound in terms of certain natural quantities, which will then enable us to provide polynomial error bounds in the smoothed complexity model. Our only assumption here is that the zero-sum game is *non-degenerate*, in the sense of [Definition 3.2](#) below; this can always be met with the addition of an arbitrarily small amount of noise ([Lemma C.1](#)). As such, our characterization here has an interest beyond the smoothed analysis framework, casting the error bound in terms of more interpretable game-theoretic quantities; for example, a concrete implication is given in [Section 4](#).

Let us denote by  $v$  the *value* of game (1), that is,

$$v = \min_{\mathbf{x} \in \mathcal{X}} \max_{\mathbf{y} \in \mathcal{Y}} \langle \mathbf{x}, \mathbf{A} \mathbf{y} \rangle = \max_{\mathbf{y} \in \mathcal{Y}} \min_{\mathbf{x} \in \mathcal{X}} \langle \mathbf{x}, \mathbf{A} \mathbf{y} \rangle,$$

which is a consequence of the minimax theorem [[von Neumann, 1928](#)]. We are now ready to state the formal definition of a non-degenerate game.

**Definition 3.2** (Non-degenerate game). A zero-sum game described with a payoff matrix  $\mathbf{A}$  and value  $v$  is said to be *non-degenerate* if it admits a unique equilibrium  $(\mathbf{x}^*, \mathbf{y}^*) \in \mathcal{Z}$ , and  $\mathbf{x}^*$  and  $\mathbf{y}^*$  make tight exactly  $n$  of the inequalities  $\{\mathbf{x}_i \geq 0\}_{i \in [n]} \cup \{\langle \mathbf{x}, \mathbf{A}_{:,j} \rangle \leq v\}_{j \in [m]}$  and  $m$  of the inequalities  $\{\mathbf{y}_j \geq 0\}_{j \in [m]} \cup \{\langle \mathbf{y}, \mathbf{A}_{i,:} \rangle \geq v\}_{i \in [n]}$ , respectively.

In the sequel, we will make constant use of the fact that the set of degenerate games has measure zero under the law induced by [Definition 1.1](#) ([Lemma C.1](#)).

In this context, we let  $B(\mathbf{x}^*) := \{i \in [n] : \mathbf{x}_i^* > 0\}$  denote the *support* of  $\mathbf{x}^*$  (corresponding to Player  $x$ ), and similarly  $N(\mathbf{y}^*) := \{j \in [m] : \mathbf{y}_j^* > 0\}$  for the support of Player  $y$ . The strict complementarity theorem [[Ye, 2011](#)] tells us that  $B$  indexes exactly the set of tight inequalities  $\{\langle \mathbf{y}, \mathbf{A}_{i,:} \rangle \geq v\}_{i \in [n]}$ , and symmetrically,  $N$  indexes exactly the set of tight inequalities  $\{\langle \mathbf{x}, \mathbf{A}_{:,j} \rangle \leq v\}_{j \in [m]}$ . In particular, this implies that  $|B| = |N|$  with probability 1. It will also be convenient to define  $\bar{B} := [n] \setminus B$  and  $\bar{N} := [m] \setminus N$ .

Now, at a high level, one can split solving a zero-sum game into two subproblems: i) identifying the support of the equilibrium, and ii) solving the induced *linear system* to specify the exact probabilities

<sup>4</sup>This is not a concern in the *unconstrained* setting, where  $\mathcal{X} = \mathbb{R}^n$  and  $\mathcal{Y} = \mathbb{R}^m$ , in which a polynomial smoothed complexity follows readily from existing results relating the convergence of OGDA or EGDA to the condition number of the payoff matrix  $\mathbf{A}$  (e.g., [[Mokhtari et al., 2020](#), [Li et al., 2023](#), [Azizian et al., 2020](#)]), which in turn is well-known to be polynomial in the smoothed complexity model [[Spielman and Teng, 2004](#)].

within the support. It will be helpful to have that viewpoint in mind in the upcoming analysis, and in particular in the proof of [Theorem 3.6](#). Roughly speaking, thinking of  $\kappa$  as a measure of the problem's difficulty, we will relate  $\kappa$  to i) the difficulty of identifying the support of the equilibrium, and ii) the difficulty of solving the induced linear system. To be clear, those two subproblems are only helpful for the purpose of the analysis, and they are certainly intertwined when using algorithms such as OGDA.

Staying on the latter task, we will make use of a certain transformation so as to eliminate one of the redundant variables. Namely, for any  $\widehat{\mathbf{x}}_B \in \Delta(B)$  and  $\widehat{\mathbf{y}}_N \in \Delta(N)$ , let us select a fixed pair of coordinates  $(i, j) \in B \times N$  (for example, the ones with the smallest index). Using the fact that  $\langle \widehat{\mathbf{x}}_B, \mathbf{1} \rangle = 1$  and  $\langle \widehat{\mathbf{y}}_N, \mathbf{1} \rangle = 1$ , we can eliminate  $\widehat{\mathbf{x}}_i$  and  $\widehat{\mathbf{y}}_j$  by writing

$$\langle \widehat{\mathbf{x}}_B, \mathbf{A}_{B,N} \widehat{\mathbf{y}}_N \rangle = \langle \widetilde{\mathbf{x}}, \mathbf{Q} \widetilde{\mathbf{y}} \rangle - \langle \widetilde{\mathbf{x}}, \mathbf{c} \rangle - \langle \widetilde{\mathbf{y}}, \mathbf{b} \rangle + d, \quad (5)$$

where  $\widetilde{\mathbf{x}} \in \mathbb{R}_{\geq 0}^{\widetilde{B}}$ ,  $\widetilde{\mathbf{y}} \in \mathbb{R}_{\geq 0}^{\widetilde{N}}$  (for  $\widetilde{B} := B \setminus \{i\}$  and  $\widetilde{N} := N \setminus \{j\}$ ) coincide with  $\widehat{\mathbf{x}}_B$  and  $\widehat{\mathbf{y}}_N$  on all coordinates in  $\widetilde{B}$  and  $\widetilde{N}$ , respectively, and  $\mathbf{A}_{B,N}^b = \mathbf{T}(\mathbf{Q}^b, \mathbf{b}, \mathbf{c}, d)$  for a (non-singular) linear transformation  $\mathbf{T} \in \mathbb{R}^{(BN) \times (BN)}$ . (We spell out the exact definition of  $\mathbf{T}$  later in [Appendix C.1](#), as it is not important for our purposes here; it follows by simply writing  $\widehat{\mathbf{x}}_i = 1 - \langle \widetilde{\mathbf{x}}, \mathbf{1} \rangle$  and  $\widehat{\mathbf{y}}_j = 1 - \langle \widetilde{\mathbf{y}}, \mathbf{1} \rangle$ .) The point of transformation (5) is that, by eliminating one of the redundant variables, there is a convenient characterization of the equilibrium ([Claim C.3](#)); namely,  $\mathbf{Q} \mathbf{y}^* = \mathbf{c}$  and  $\mathbf{Q}^\top \mathbf{x}^* = \mathbf{b}$ .

We are now ready to introduce the key quantities upon which our characterization relies on. It turns out that those are analogous to the ones considered by [Spielman and Teng \[2003\]](#) in the context of analyzing the termination of linear programs; this is not coincidental, as our analysis was especially targeted to do so.

**Definition 3.3.** Let  $\mathbf{A}$  be the payoff matrix of a non-degenerate game,  $(\mathbf{x}^*, \mathbf{y}^*) \in \mathcal{Z}$  the unique equilibrium, and  $B \subseteq [n], N \subseteq [m]$  the support of  $\mathbf{x}^*$  and  $\mathbf{y}^*$  respectively. We introduce the following quantities.

1.  $\alpha_P(\mathbf{A}) := \min_{i \in B} (\mathbf{x}_i^*)$  and  $\alpha_D(\mathbf{A}) := \min_{j \in N} (\mathbf{y}_j^*)$ ;
2.  $\beta_P(\mathbf{A}) := \min_{j \in N} (v - \langle \mathbf{x}_B^*, \mathbf{A}_{B,j} \rangle)$  and  $\beta_D(\mathbf{A}) := \min_{i \in B} (\langle \mathbf{A}_{i,N}, \mathbf{y}_N^* \rangle - v)$ ; and
3.  $\gamma_P(\mathbf{A}) := \min_{j \in \widetilde{N}} \text{dist}(\mathbf{Q}_{:,j}, \text{span}(\mathbf{Q}_{:, \widetilde{N} - j}))$  and  $\gamma_D(\mathbf{A}) := \min_{i \in \widetilde{B}} \text{dist}(\mathbf{Q}_{i,:}, \text{span}(\mathbf{Q}_{\widetilde{B} - i, :}))$ , where we use the shorthand notation  $\widetilde{B} - i := \widetilde{B} \setminus \{i\}$  ( $\widetilde{N} - j := \widetilde{N} \setminus \{j\}$ ), and  $\mathbf{Q} = \mathbf{Q}(\mathbf{A})$  is defined in (5).

(Above, we adopt the convention that if a minimization problem is with respect to an empty set, the minimum is to be evaluated as 1.)

[Item 3](#) above will enable us to control the norm of solutions to any linear system induced by  $\mathbf{Q}$ , as we explain in the sequel. Our proof will actually rely on a slightly different matrix, which we call  $\overline{\mathbf{Q}}$ ; the lemma below relates the geometry of  $\overline{\mathbf{ Q}}$  to  $\mathbf{Q}$ , and reassures us that the condition number of  $\overline{\mathbf{ Q}}$  cannot be far from that of  $\mathbf{Q}$  so long as  $1 - \sum_{j \in \widetilde{N}} \mathbf{y}_j^* \geq \alpha_D(\mathbf{A})$  (by [Item 1](#)) is not too close to 0. (A symmetric statement holds when focusing on Player  $y$ .)

**Lemma 3.4.** Let  $\mathbf{c} = \mathbf{Q} \mathbf{y}^* = \sum_{j \in \widetilde{N}} \mathbf{y}_j^* \mathbf{Q}_{:,j}$ , and suppose that  $\overline{\mathbf{Q}} \in \mathbb{R}^{\widetilde{B} \times \widetilde{N}}$  is such that its  $j$ th column is equal to  $\mathbf{Q}_{:,j} - \mathbf{c}$ . Then,

$$\min_{j \in \widetilde{N}} \text{dist}(\mathbf{Q}_{:,j}, \text{span}(\mathbf{Q}_{:, \widetilde{N} - j})) \leq \left( 1 + \frac{|\widetilde{N}|}{1 - \sum_{j \in \widetilde{N}} \mathbf{y}_j^*} \right) \min_{j \in \widetilde{N}} \text{dist}(\overline{\mathbf{Q}}_{:,j}, \text{span}(\overline{\mathbf{Q}}_{:, \widetilde{N} - j})).$$

Next, we recall a fairly standard bound relating the magnitude of a solution to a linear system  $\widetilde{\mathbf{x}} = \mathbf{M}\mathbf{p}$  with the smallest singular value of a full-rank matrix  $\mathbf{M}$ .

**Lemma 3.5.** Let  $\mathbf{M} \in \mathbb{R}^{d \times d}$  be a full-rank matrix. For any  $\widetilde{\mathbf{x}} \in \mathbb{R}^d$  there is  $\mathbf{p} \in \mathbb{R}^d$  with  $\|\mathbf{p}\| \leq \frac{1}{\sigma_{\min}(\mathbf{M})} \|\widetilde{\mathbf{x}}\|$  such that

$$\widetilde{\mathbf{x}} = \mathbf{M}\mathbf{p} = \sum_{j=1}^d \mathbf{p}_j \mathbf{M}_{:,j}.$$

Moreover, to connect [Lemma 3.5](#) with  $\gamma_P(\mathbf{A})$ , we observe that the smallest singular value can also be lower bounded in terms of the smallest distance of a column from the linear space spanned by the rest of the columns—which now matches the expression of [Item 3](#) we saw earlier. In particular, we will make use of the so-called negative second moment identity [Tao et al., 2010] ([Proposition C.4](#)), which implies that

$$\sigma_{\min}(\bar{\mathbf{Q}}) \geq \sqrt{\frac{1}{\sum_{j \in \tilde{N}} \text{dist}^{-2}(\bar{\mathbf{Q}}_{:,j}, \text{span}(\bar{\mathbf{Q}}_{:,\tilde{N}-j}))}} \geq \frac{1}{\sqrt{|\tilde{N}|}} \min_{j \in \tilde{N}} \text{dist}(\bar{\mathbf{Q}}_{:,j}, \text{span}(\bar{\mathbf{Q}}_{:,\tilde{N}-j})). \quad (6)$$

[Proposition C.4](#) also implies that  $\gamma_D(\mathbf{A}) \geq \frac{1}{\sqrt{|B|}} \gamma_P(\mathbf{A})$ , and so it will suffice to lower bound  $\gamma_P(\mathbf{A})$  in the sequel. We are now ready to proceed with the main result of this subsection. Below, we use the notation “ $\gtrsim$ ” to suppress lower-order terms and absolute constants.

**Theorem 3.6.** *Let  $\mathbf{A}$  be a non-degenerate payoff matrix, and suppose that  $(\alpha_P(\mathbf{A}), \alpha_D(\mathbf{A}))$ ,  $(\beta_P(\mathbf{A}), \beta_D(\mathbf{A}))$  and  $(\gamma_P(\mathbf{A}), \gamma_D(\mathbf{A}))$  are as in [Definition 3.3](#). Then, the error bound ([Definition 1.3](#)) is satisfied for any sufficiently small modulus*

$$\kappa \gtrsim \frac{1}{\|\mathbf{A}^b\|_\infty} \frac{1}{\min(n, m)^3} \min \{(\alpha_D(\mathbf{A}))^2 \beta_D(\mathbf{A}) \gamma_P(\mathbf{A}), (\alpha_P(\mathbf{A}))^2 \beta_P(\mathbf{A}) \gamma_D(\mathbf{A})\}.$$

It is enough to explain how to lower bound  $\kappa > 0$  such that  $\max_{\mathbf{y}' \in \mathcal{Y}} \langle \mathbf{x}, \mathbf{A} \mathbf{y}' \rangle - v \geq \kappa \|\mathbf{x} - \Pi_{\mathcal{X}^*}(\mathbf{x})\| = \kappa \|\mathbf{x} - \mathbf{x}^*\|$  for any  $\mathbf{x} \in \mathcal{X}$ . In a nutshell, our argument is divided based on the magnitude  $\lambda := \|\mathbf{x}_B\|$ , which can be thought of as a measure of closeness from the support of the equilibrium. When  $\lambda \ll 1$ , which means that  $\mathbf{x}$  is still far from the support of the equilibrium,  $\max_{\mathbf{y}' \in \mathcal{Y}} \langle \mathbf{x}, \mathbf{A} \mathbf{y}' \rangle - v$  is governed by  $\beta_D(\mathbf{A})$ . In the contrary case, our basic strategy revolves around showing that the error bound can be treated as in the unconstrained case, which would then relate the modulus  $\kappa$  to the smallest singular value of the underlying matrix (essentially by [Lemma 3.5](#))—and subsequently to  $\gamma_P(\mathbf{A})$  due to (6). Indeed, this turns out to be possible by working with matrix  $\bar{\mathbf{Q}}$ , as defined earlier in [Lemma 3.4](#). We defer the precise argument to [Appendix C.1](#).

### 3.3 Smoothed analysis

Having established [Theorem 3.6](#), our next step is to show that each of the quantities introduced in [Definition 3.3](#) is unlikely to be too close to 0 in the smoothed complexity model, which would then imply [Theorem 1.4](#). The main difficulty lies in the fact that each configuration that may arise depends on the support of the equilibrium, which in turn depends on the underlying randomization of  $\mathbf{A}$ , thereby significantly complicating the underlying distribution. Further, one cannot afford to argue about each configuration separately and then apply the union bound as there are too many possible configurations. To tackle this challenge, we follow the approach put forward by [Spielman and Teng \[2003\]](#).

In particular, given that all quantities of interest in [Theorem 3.6](#) depend on the support of the equilibrium, it is natural to proceed by partitioning the probability space over all possible supports, and then bound the worst possible one—that is, the one maximizing the probability we want to minimize. In doing so, the challenge is that one has to condition on the equilibrium having a given support (formally justified by [Proposition C.5](#)). To argue about the induced probability density function upon such a conditioning, it is convenient to perform a change of variables from  $\mathbf{A}$  to a new set of variables that now contains the equilibrium  $(\mathbf{x}^*, \mathbf{y}^*)$  ([Lemma C.6](#)). The basic idea here is that since the event we condition on concerns the equilibrium, it is helpful to have that equilibrium being part of our set of variables. The induced probability density function is now quite complicated, but can still be analyzed using the following lemma.

**Lemma 3.7** ([Spielman and Teng, 2003](#)). *Let  $\rho$  be the probability density function of a random variable  $X$ . If there exist  $\delta > 0$  and  $c \in (0, 1]$  such that*

$$0 \leq t \leq t' \leq \delta \implies \frac{\rho(t')}{\rho(t)} \geq c, \quad (7)$$

then

$$\mathbb{P}[X \leq \epsilon \mid X \geq 0] \leq \frac{\epsilon}{c\delta}.$$

In words, random variables whose density is smooth—in the sense of (7)—are unlikely to be too close to 0. Gaussian random variables certainly have that property (Lemma C.8), but it is not confined to the Gaussian law; the analysis of [Spielman and Teng \[2003\]](#)—and subsequently our result—is not tailored to the Gaussian case.

We are now ready to state our main results in the smoothed complexity model; the proofs are deferred to [Appendix C.2](#). We commence with  $\beta_P(\mathbf{A})$ , which is the easiest to analyze. In particular, the following result is a consequence of an anti-concentration bound with respect to a conditional Gaussian random variable (Lemma C.7).

**Proposition 3.8.** *Let  $\beta_P(\mathbf{A})$  be defined as in [Item 2](#). For any  $\epsilon \geq 0$ ,*

$$\mathbb{P}_{\mathbf{A}} \left[ \beta_P(\mathbf{A}) \leq \frac{\epsilon}{5\|\mathbf{A}^b\|_\infty} \right] \leq \epsilon \frac{e \min(n, m)^2}{\sigma^2}.$$

The analysis of  $\gamma_P(\mathbf{A})$  is more challenging, and makes crucial use of [Lemma 3.7](#). As we alluded to earlier, a key step is to change variables from  $\mathbf{A}_{B,N}$  to  $(\mathbf{Q}, \mathbf{b}, \mathbf{c}, \cdot)$ —in accordance with (5)—and then to  $(\mathbf{Q}, \mathbf{x}^*, \mathbf{y}^*, \cdot)$  based on  $\mathbf{Q}\tilde{\mathbf{y}}^* = \mathbf{c}$ ,  $\mathbf{Q}^\top \tilde{\mathbf{x}}^* = \mathbf{b}$ . It is important to note that  $\mathbf{Q}$  no longer contains independent random variables even though  $\mathbf{A}_{B,N}$  is (by [Definition 1.1](#)); this stems from the presence of a redundant variable in  $\mathbf{x}_B^*$  (since  $\langle \mathbf{x}_B^*, \mathbf{1} \rangle = 1$ ). Nevertheless, we can still overcome this issue using [Lemma 3.7](#), leading to the following bound.

**Proposition 3.9.** *Let  $\gamma_P(\mathbf{A})$  be defined as in [Item 3](#). For any  $\epsilon \geq 0$ ,*

$$\mathbb{P}_{\mathbf{A}} \left[ \gamma_P(\mathbf{A}) \leq \frac{\epsilon}{4 \max_{j \in \tilde{N}} \|\mathbf{Q}_{\cdot,j}\| + 20\|\mathbf{A}^b\|_\infty + 3} \right] \leq \epsilon \frac{4e \min(n, m)^3}{\sigma^2}.$$

Similar reasoning, albeit with some further complications, provides a bound for  $\alpha_P(\mathbf{A})$ , which is given below.

**Proposition 3.10.** *Let  $\alpha_P(\mathbf{A})$  be defined as in [Item 1](#). For any  $\epsilon \geq 0$ ,*

$$\mathbb{P}_{\mathbf{A}} \left[ \alpha_P(\mathbf{A}) \leq \frac{\epsilon}{25(\|\mathbf{A}^b\|_\infty + 1)^2} \right] \leq \epsilon \frac{8e^2 mn \min(n, m)}{\sigma^2}.$$

Armed with [Propositions 3.8 to 3.10](#) and [Theorem 3.6](#), we can establish [Theorem 1.2](#) by suitably leveraging existing results, as we formalize in [Appendix C.3](#).

## 4 Parameterized results for perturbation-stable games

Another important implication of our characterization in [Theorem 3.6](#) is that it enables connecting the convergence rate of gradient-based algorithms to natural and interpretable game-theoretic quantities. In particular, here we highlight a connection with perturbation-stable games, in the following formal sense.

**Definition 4.1** (Perturbation-stable games). Let  $\mathbf{A}$  be the payoff matrix of a non-degenerate game. We say that the game is  $\delta$ -support-stable, with  $\delta > 0$ , if for any  $\mathbf{A}'$  with  $\|\mathbf{A} - \mathbf{A}'\| \leq \delta$  it holds that  $\mathbf{A}'$  is a non-degenerate game whose equilibrium has the same support as  $\mathbf{A}$ .

Perhaps the simplest example of a support-stable game with a favorable parameter  $\delta > 0$  arises when  $\mathbf{A}$  is the  $2 \times 2$  identity matrix. Indeed, as long as the perturbation parameter  $\delta$  remains below a certain absolute constant, the perturbed game still admits a unique full-support equilibrium. To see this, suppose for the sake of contradiction that the perturbed game has an equilibrium such that Player  $x$  plays one of the two actions with probability 1. Player  $y$  would then obtain a utility of at least  $1 - O(\delta)$ . But the value of the original game was  $1/2$ , which in turn implies that the value of the perturbed game is  $1/2 \pm \Theta(\delta)$ ; for a sufficiently small  $\delta$  this leads to a contradiction. Similar reasoning applies with respect to Player  $y$ . (The previous argument carries over more broadly to diagonally dominant  $2 \times 2$  payoff matrices.)

As we have highlighted already, games with perturbation-stable equilibria—albeit under different notions of stability—have already received attention in the literature [[Balcan and Braverman, 2017](#), [Awasthi et al., 2010](#)] (cf. [Cohen \[1986\]](#)), and are part of a broader trend in the analysis of algorithms beyond the worst case (for further background, we refer to the excellent book edited by [Roughgarden \[2021\]](#)). Our goal here is to make the following natural connection.

**Theorem 4.2.** Any  $\delta$ -support-stable game (per Definition 4.1) satisfies the error bound for any sufficiently small modulus

$$\kappa \geq \text{poly} \left( \frac{1}{n}, \frac{1}{m}, \delta \right).$$

By virtue of our discussion in Appendix C.3, Theorem 4.2 immediately implies Corollary 1.6. Indeed, we observe that all parameters involved in Theorem 3.6 can be lower bounded in terms of the stability parameter of Definition 4.1, as we formalize in Appendix C.4.

## 5 Conclusions and future research

In conclusion, we performed the first smoothed analysis with respect to a number of well-studied gradient-based algorithms in zero-sum games. In particular, we showed that OGDA, EGDA and IterSmooth all enjoy polynomial smoothed complexity, meaning that their iteration complexity grows as a polynomial in the dimensions of the game,  $1/\sigma$ , and  $\log(1/\epsilon)$ ; for OMWU, our analysis reveals a significant improvement over the worst-case bound due to Wei et al. [2021], but it still remains superpolynomial. We also made a connection between the rate of convergence of the above algorithms and a natural perturbation-stability property of the equilibrium, which is interesting beyond the model of smoothed complexity.

A number of interesting avenues for future research remain open. First, is it the case that OMWU has polynomial smoothed complexity or is there an inherent separation with the other algorithms we studied? Answering this question in the positive would necessitate significantly improving the worst-case analysis of OMWU due to Wei et al. [2021] (cf. Cai et al. [2024] for a recent development concerning the last-iterate convergence of OMWU). Beyond OMWU, our results could also prove useful for establishing polynomial bounds for other natural dynamics in the smoothed analysis framework. Moreover, our characterization of the error bound in Theorem 3.6 assumes that the game is non-degenerate. This is an innocuous assumption in the smoothed complexity model, as it holds with probability 1, but nevertheless it would be interesting to generalize it to any game. Doing so could shed some light into whether Theorem 4.2 holds with respect to other, perhaps more natural notions of perturbation stability beyond Definition 4.1. It would also be interesting to investigate other models of smoothed complexity that account for dependencies between the entries of the payoff matrix [Bhaskara et al., 2024]. Moreover, our focus has been on zero-sum games under simplex constraints, but we suspect that more general positive results should be attainable under polyhedral constraint sets; perhaps the most notable such candidate is the class of *extensive-form games* [Romanovskii, 1962, von Stengel, 1996]. Even beyond (two-player) zero-sum games, Theorem 1.2 could apply to (multi-player) *polymatrix* zero-sum games [Cai et al., 2016]. It is less clear whether the model of smoothed complexity can be informative when it comes to convergence to *coarse correlated equilibria* in multi-player games.

## Acknowledgments

We are grateful to the anonymous reviewers at NeurIPS for their helpful feedback. The first author is indebted to Ioannis Panageas for many insightful discussions. This material is based on work supported by the Vannevar Bush Faculty Fellowship ONR N00014-23-1-2876, National Science Foundation grants RI-2312342 and RI-1901403, ARO award W911NF2210266, and NIH award A240108S001.

## References

Kenshi Abe, Kaito Ariu, Mitsuki Sakamoto, Kentaro Toyoshima, and Atsushi Iwasaki. Last-iterate convergence with full and noisy feedback in two-player zero-sum games. In *International Conference on Artificial Intelligence and Statistics (AISTATS)*, 2023.

Ilan Adler. The equivalence of linear programs and zero-sum games. *Int. J. Game Theory*, 42(1): 165–177, 2013.

Ahmet Alacaoglu and Yura Malitsky. Stochastic variance reduction for variational inequality methods. In *Conference on Learning Theory (COLT)*, 2022.

Kimon Antonakopoulos, Elena Veronica Belmega, and Panayotis Mertikopoulos. Adaptive extra-gradient methods for min-max optimization and games. In *International Conference on Learning Representations (ICLR)*, 2021.

David L. Applegate, Oliver Hinder, Haihao Lu, and Miles Lubin. Faster first-order primal-dual methods for linear programming using restarts and sharpness. *Mathematical Programming*, 201(1):133–184, 2023.

Pranjal Awasthi, Maria-Florina Balcan, Avrim Blum, Or Sheffet, and Santosh S. Vempala. On nash-equilibria of approximation-stable games. In *International Symposium on Algorithmic Game Theory (SAGT)*, 2010.

Waïss Azizian, Damien Scieur, Ioannis Mitliagkas, Simon Lacoste-Julien, and Gauthier Gidel. Accelerating smooth games by manipulating spectral shapes. In *International Conference on Artificial Intelligence and Statistics (2020)*, Proceedings of Machine Learning Research, 2020.

Maria-Florina Balcan and Mark Braverman. Nash equilibria in perturbation-stable games. *Theory Comput.*, 13(1):1–31, 2017.

Aditya Bhaskara, Eric Evert, Vaidehi Srinivas, and Aravindan Vijayaraghavan. New tools for smoothed analysis: Least singular value bounds for random matrices with dependent entries. In *Proceedings of the Annual Symposium on Theory of Computing (STOC)*, 2024.

Markus Bläser and Bodo Manthey. Smoothed complexity theory. *ACM Trans. Comput. Theory*, 7(2):6:1–6:21, 2015.

Avrim Blum and John Dunagan. Smoothed analysis of the perceptron algorithm for linear programming. In *Annual ACM-SIAM Symposium on Discrete Algorithms (SODA)*, 2002.

Shant Boodaghians, Joshua Brakensiek, Samuel B. Hopkins, and Aviad Rubinstein. Smoothed complexity of 2-player nash equilibria. In *Proceedings of the Annual Symposium on Foundations of Computer Science (FOCS)*, 2020.

Michael Bowling, Neil Burch, Michael Johanson, and Oskari Tammelin. Heads-up limit hold’em poker is solved. *Science*, 347(6218):145–149, 2015.

Benjamin Brooks and Philip J. Reny. A canonical game—75 years in the making—showing the equivalence of matrix games and linear programming. *Economic Theory Bulletin*, 2023.

Noam Brown and Tuomas Sandholm. Superhuman AI for heads-up no-limit poker: Libratus beats top professionals. *Science*, 359(6374):418–424, 2018.

Noam Brown and Tuomas Sandholm. Solving imperfect-information games via discounted regret minimization. In *Conference on Artificial Intelligence (AAAI)*, 2019.

Luciana S. Buriol, Marcus Ritt, Félix Carvalho Rodrigues, and Guido Schäfer. On the smoothed price of anarchy of the traffic assignment problem. In *Workshop on Algorithmic Approaches for Transportation Modeling, Optimization, and Systems (ATMOS)*, 2011.

Yang Cai, Ozan Candogan, Constantinos Daskalakis, and Christos H. Papadimitriou. Zero-sum polymatrix games: A generalization of minimax. *Mathematics of Operations Research*, 41(2):648–655, 2016.

Yang Cai, Argyris Oikonomou, and Weiqiang Zheng. Finite-time last-iterate convergence for learning in multi-player games. In *Proceedings of the Annual Conference on Neural Information Processing Systems (NeurIPS)*, 2022.

Yang Cai, Gabriele Farina, Julien Grand-Clément, Christian Kroer, Chung-Wei Lee, Haipeng Luo, and Weiqiang Zheng. Fast last-iterate convergence of learning in games requires forgetful algorithms. In *Proceedings of the Annual Conference on Neural Information Processing Systems (NeurIPS)*, 2024.

Yair Carmon, Yujia Jin, Aaron Sidford, and Kevin Tian. Variance reduction for matrix games. In *Proceedings of the Annual Conference on Neural Information Processing Systems (NeurIPS)*, 2019.

Yair Carmon, Yujia Jin, Aaron Sidford, and Kevin Tian. Coordinate methods for matrix games. In *Proceedings of the Annual Symposium on Foundations of Computer Science (FOCS)*, 2020.

Yair Carmon, Arun Jambulapati, Yujia Jin, and Aaron Sidford. A whole new ball game: A primal accelerated method for matrix games and minimizing the maximum of smooth functions. In *Annual ACM-SIAM Symposium on Discrete Algorithms (SODA)*, 2024.

Xi Chen, Xiaotie Deng, and Shang-Hua Teng. Settling the complexity of computing two-player Nash equilibria. *Journal of the ACM*, 2009.

Xi Chen, Chenghao Guo, Emmanouil-Vasileios Vlatakis-Gkaragkounis, Mihalis Yannakakis, and Xinzhong Zhang. Smoothed complexity of local max-cut and binary max-csp. In *Proceedings of the Annual Symposium on Theory of Computing (STOC)*, 2020.

Xi Chen, Chenghao Guo, Emmanouil-Vasileios Vlatakis-Gkaragkounis, and Mihalis Yannakakis. Smoothed complexity of SWAP in local graph partitioning. *Annual ACM-SIAM Symposium on Discrete Algorithms (SODA)*, 2024.

Chao-Kai Chiang, Tianbao Yang, Chia-Jung Lee, Mehrdad Mahdavi, Chi-Jen Lu, Rong Jin, and Shenghuo Zhu. Online optimization with gradual variations. In *Conference on Learning Theory (COLT)*, 2012.

Miranda Christ and Mihalis Yannakakis. The smoothed complexity of policy iteration for markov decision processes. In *Proceedings of the Annual Symposium on Theory of Computing (STOC)*, 2023.

Kenneth L. Clarkson, Elad Hazan, and David P. Woodruff. Sublinear optimization for machine learning. *Journal of the ACM*, 59(5):23:1–23:49, 2012.

Joel E. Cohen. Perturbation theory of completely mixed matrix games. *Linear Algebra and its Applications*, 79:153–162, 1986.

Johanne Cohen, Amélie Héliou, and Panayotis Mertikopoulos. Hedging under uncertainty: Regret minimization meets exponentially fast convergence. In *International Symposium on Algorithmic Game Theory (SAGT)*, 2017.

Michael B. Cohen, Yin Tat Lee, and Zhao Song. Solving linear programs in the current matrix multiplication time. *Journal of the ACM*, 68(1):3:1–3:39, 2021.

Leonardo Cunha, Gauthier Gidel, Fabian Pedregosa, Damien Scieur, and Courtney Paquette. Only tails matter: Average-case universality and robustness in the convex regime. In *International Conference on Machine Learning (ICML)*, 2022.

George Dantzig. A proof of the equivalence of the programming problem and the game problem. In Tjalling Koopmans, editor, *Activity Analysis of Production and Allocation*, pages 330–335. John Wiley & Sons, 1951.

Constantinos Daskalakis and Ioannis Panageas. Last-iterate convergence: Zero-sum games and constrained min-max optimization. In *Innovations in Theoretical Computer Science Conference (ITCS)*, 2019.

Constantinos Daskalakis, Alan Deckelbaum, and Anthony Kim. Near-optimal no-regret algorithms for zero-sum games. *Games and Economic Behavior*, 92:327–348, 2015.

Constantinos Daskalakis, Noah Golowich, Nika Haghtalab, and Abhishek Shetty. Smooth nash equilibria: Algorithms and complexity. In *Innovations in Theoretical Computer (ITCS)*, 2024.

Asen L Dontchev and R Tyrrell Rockafellar. *Implicit functions and solution mappings: A view from variational analysis*, volume 616. Springer, 2009.

John Dunagan, Daniel A. Spielman, and Shang-Hua Teng. Smoothed analysis of condition numbers and complexity implications for linear programming. *Mathematical Programming*, 126(2):315–350, 2011.

Alan Edelman. Eigenvalue roulette and random test matrices. *Linear Algebra for Large Scale and Real-Time Applications*, pages 365–368, 1993.

Kousha Etessami and Mihalis Yannakakis. On the complexity of Nash equilibria and other fixed points (extended abstract). In *Proceedings of the Annual Symposium on Foundations of Computer Science (FOCS)*, 2007.

Gabriele Farina and Tuomas Sandholm. Fast payoff matrix sparsification techniques for structured extensive-form games. In *Conference on Artificial Intelligence (AAAI)*, 2022.

Gabriele Farina, Christian Kroer, and Tuomas Sandholm. Faster game solving via predictive blackwell approachability: Connecting regret matching and mirror descent. In *Conference on Artificial Intelligence (AAAI)*, 2021.

Olivier Fercoq. Quadratic error bound of the smoothed gap and the restarted averaged primal-dual hybrid gradient, 2023.

Nicola Gatti, Marco Rocco, and Tuomas Sandholm. Strong Nash equilibrium is in smoothed P. In *Conference on Artificial Intelligence (AAAI)*, 2013. Late-breaking paper track.

Yiannis Giannakopoulos. A smoothed FPTAS for equilibria in congestion games. *CoRR*, abs/2306.10600, 2023.

Yiannis Giannakopoulos, Alexander Grosz, and Themistoklis Melissourgos. On the smoothed complexity of combinatorial local search. *CoRR*, abs/2211.07547, 2022.

Angeliki Giannou, Emmanouil-Vasileios Vlatakis-Gkaragkounis, and Panayotis Mertikopoulos. On the rate of convergence of regularized learning in games: From bandits and uncertainty to optimism and beyond. In *Proceedings of the Annual Conference on Neural Information Processing Systems (NeurIPS)*, 2021.

Andrew Gilpin, Javier Peña, and Tuomas Sandholm. First-order algorithm with  $\mathcal{O}(\ln(1/\epsilon))$  convergence for  $\epsilon$ -equilibrium in two-person zero-sum games. *Mathematical Programming*, 133(1–2): 279–298, 2012.

Noah Golowich, Sarath Pattathil, and Constantinos Daskalakis. Tight last-iterate convergence rates for no-regret learning in multi-player games. In *Proceedings of the Annual Conference on Neural Information Processing Systems (NeurIPS)*, 2020a.

Noah Golowich, Sarath Pattathil, Constantinos Daskalakis, and Asuman E. Ozdaglar. Last iterate is slower than averaged iterate in smooth convex-concave saddle point problems. In *Conference on Learning Theory (COLT)*, 2020b.

Eduard Gorbunov, Adrien Taylor, and Gauthier Gidel. Last-iterate convergence of optimistic gradient method for monotone variational inequalities. In *Proceedings of the Annual Conference on Neural Information Processing Systems (NeurIPS)*, 2022.

Michael D. Grigoriadis and Leonid G. Khachiyan. A sublinear-time randomized approximation algorithm for matrix games. *Operations Research Letters*, 18(2):53–58, 1995.

Nika Haghtalab, Michael I. Jordan, and Eric Zhao. On-demand sampling: Learning optimally from multiple distributions. In *Proceedings of the Annual Conference on Neural Information Processing Systems (NeurIPS)*, 2022.

Nika Haghtalab, Michael I. Jordan, and Eric Zhao. A unifying perspective on multi-calibration: Unleashing game dynamics for multi-objective learning. In *Proceedings of the Annual Conference on Neural Information Processing Systems (NeurIPS)*, 2023.

Yu-Guan Hsieh, Franck Iutzeler, Jérôme Malick, and Panayotis Mertikopoulos. On the convergence of single-call stochastic extra-gradient methods. In *Proceedings of the Annual Conference on Neural Information Processing Systems (NeurIPS)*, 2019.

Sophie Huiberts, Yin Tat Lee, and Xinzhi Zhang. Upper and lower bounds on the smoothed complexity of the simplex method. In *Proceedings of the Annual Symposium on Theory of Computing (STOC)*, 2023.

Galina M Korpelevich. The extragradient method for finding saddle points and other problems. *Matecon*, 12:747–756, 1976.

Chung-Wei Lee, Christian Kroer, and Haipeng Luo. Last-iterate convergence in extensive-form games. In *Proceedings of the Annual Conference on Neural Information Processing Systems (NeurIPS)*, 2021.

Chris Junchi Li, Huizhuo Yuan, Gauthier Gidel, Quanquan Gu, and Michael I. Jordan. Nesterov meets optimism: Rate-optimal separable minimax optimization. In *International Conference on Machine Learning (ICML)*, 2023.

Pouria Mahdavina, Yuyang Deng, Haochuan Li, and Mehrdad Mahdavi. Tight analysis of extra-gradient and optimistic gradient methods for nonconvex minimax problems. In *Proceedings of the Annual Conference on Neural Information Processing Systems (NeurIPS)*, 2022.

Arnab Maiti, Kevin G. Jamieson, and Lillian J. Ratliff. Instance-dependent sample complexity bounds for zero-sum matrix games. In *International Conference on Artificial Intelligence and Statistics (AISTATS)*, 2023.

Konstantin Makarychev and Yury Makarychev. *Perturbation Resilience*, page 95–119. Cambridge University Press, 2021.

Panayotis Mertikopoulos, Christos H. Papadimitriou, and Georgios Piliouras. Cycles in adversarial regularized learning. In *Annual ACM-SIAM Symposium on Discrete Algorithms (SODA)*, 2018.

Panayotis Mertikopoulos, Bruno Lecouat, Houssam Zenati, Chuan-Sheng Foo, Vijay Chandrasekhar, and Georgios Piliouras. Optimistic mirror descent in saddle-point problems: Going the extra (gradient) mile. In *International Conference on Learning Representations (ICLR)*, 2019.

Aryan Mokhtari, Asuman E. Ozdaglar, and Sarath Pattathil. A unified analysis of extra-gradient and optimistic gradient methods for saddle point problems: Proximal point approach. In *International Conference on Artificial Intelligence and Statistics (AISTATS)*, 2020.

Matej Moravčík, Martin Schmid, Neil Burch, Viliam Lisý, Dustin Morrill, Nolan Bard, Trevor Davis, Kevin Waugh, Michael Johanson, and Michael Bowling. Deepstack: Expert-level artificial intelligence in heads-up no-limit poker. *Science*, 356(6337):508–513, 2017.

Yurii Nesterov. Smooth minimization of non-smooth functions. *Mathematical Programming*, 103, 2005.

Courtney Paquette, Bart van Merriënboer, Elliot Paquette, and Fabian Pedregosa. Halting time is predictable for large models: A universality property and average-case analysis. *Found. Comput. Math.*, 23(2):597–673, 2023.

Julien Perolat, Bart De Vylder, Daniel Hennes, Eugene Tarassov, Florian Strub, Vincent de Boer, Paul Muller, Jerome T. Connor, Neil Burch, Thomas Anthony, Stephen McAleer, Romuald Elie, Sarah H. Cen, Zhe Wang, Audrunas Gruslys, Aleksandra Malysheva, Mina Khan, Sherjil Ozair, Finbarr Timbers, Toby Pohlen, Tom Eccles, Mark Rowland, Marc Lanctot, Jean-Baptiste Lespiau, Bilal Piot, Shayegan Omidshafiei, Edward Lockhart, Laurent Sifre, Nathalie Beauguerlange, Remi Munos, David Silver, Satinder Singh, Demis Hassabis, and Karl Tuyls. Mastering the game of stratego with model-free multiagent reinforcement learning. *Science*, 378(6623):990–996, 2022.

L.D. Popov. A modification to the Arrow-Hurwicz method for search of saddle-points. *Mathematical Notes of the Academy of Sciences of the USSR*, 28(5):845–848, 1980.

Alexander Rakhlin and Karthik Sridharan. Online learning with predictable sequences. In *Conference on Learning Theory*, pages 993–1019, 2013.

James Renegar. Incorporating condition measures into the complexity theory of linear programming. *SIAM Journal on Optimization*, 5(3):506–524, 1995.

Jarnes Renegar. Some perturbation theory for linear programming. *Mathematical Programming*, 65: 73–91, 1994.

Ralph Tyrell Rockafellar. *Convex Analysis*. Princeton university press, 2015.

I. Romanovskii. Reduction of a game with complete memory to a matrix game. *Soviet Mathematics*, 3, 1962.

Tim Roughgarden. *Beyond the Worst-Case Analysis of Algorithms*. Cambridge University Press, 2021.

Aviad Rubinstein. Settling the complexity of computing approximate two-player nash equilibria. In Irit Dinur, editor, *Proceedings of the Annual Symposium on Foundations of Computer Science (FOCS)*, 2016.

Damien Scieur and Fabian Pedregosa. Universal average-case optimality of polyak momentum. In *International Conference on Machine Learning (ICML)*, 2020.

Zhuoqing Song, Jason D. Lee, and Zhuoran Yang. Can we find nash equilibria at a linear rate in markov games? In *International Conference on Learning Representations (ICLR)*, 2023.

Daniel A. Spielman and Shang-Hua Teng. Smoothed analysis of termination of linear programming algorithms. *Math. Program.*, 97(1-2):375–404, 2003.

Daniel A. Spielman and Shang-Hua Teng. Smoothed analysis of algorithms: Why the simplex algorithm usually takes polynomial time. *Journal of the ACM*, 51(3):385–463, 2004.

Daniel A. Spielman and Shang-Hua Teng. Smoothed analysis: an attempt to explain the behavior of algorithms in practice. *Commun. ACM*, 52(10):76–84, 2009.

Vasilis Syrgkanis, Alekh Agarwal, Haipeng Luo, and Robert E Schapire. Fast convergence of regularized learning in games. In *Advances in Neural Information Processing Systems*, 2015.

Xiaohang Tang, Le Cong Dinh, Stephen Marcus McAleer, and Yaodong Yang. Regret-minimizing double oracle for extensive-form games. In *International Conference on Machine Learning (ICML)*, Proceedings of Machine Learning Research, 2023.

Terence Tao. *Topics in random matrix theory*, volume 132. American Mathematical Society, 2023.

Terence Tao, Van Vu, and Manjunath Krishnapur. Random matrices: Universality of ESDs and the circular law. *The Annals of Probability*, 38(5):2023 – 2065, 2010.

Paul Tseng. On linear convergence of iterative methods for the variational inequality problem. *Journal of Computational and Applied Mathematics*, 60(1):237–252, 1995.

Eric van Damme. *Stability and perfection of Nash equilibria*, volume 339. Springer, 1991.

Jan van den Brand, Yin Tat Lee, Yang P. Liu, Thatchaphol Saranurak, Aaron Sidford, Zhao Song, and Di Wang. Minimum cost flows, mdps, and  $\ell_1$ -regression in nearly linear time for dense instances. In *Proceedings of the Annual Symposium on Theory of Computing (STOC)*, 2021.

Daniil Vankov, Angelia Nedić, and Lalitha Sankar. Last iterate convergence of popov method for non-monotone stochastic variational inequalities, 2023.

John von Neumann. Zur Theorie der Gesellschaftsspiele. *Mathematische Annalen*, 100:295–320, 1928.

John von Neumann and Oskar Morgenstern. *Theory of Games and Economic Behavior*. Princeton University Press, 1947.

Bernhard von Stengel. Efficient computation of behavior strategies. *Games and Economic Behavior*, 14(2):220–246, 1996.

Bernhard von Stengel. Zero-sum games and linear programming duality. *Mathematics of Operations Research*, 2023.

Chen-Yu Wei, Chung-Wei Lee, Mengxiao Zhang, and Haipeng Luo. Linear last-iterate convergence in constrained saddle-point optimization. In *International Conference on Learning Representations (ICLR)*, 2021.

Yinyu Ye. *Interior point algorithms: theory and analysis*. John Wiley & Sons, 2011.

Nikos Zarifis, Puqian Wang, Ilias Diakonikolas, and Jelena Diakonikolas. Robustly learning single-index models via alignment sharpness. In *International Conference on Machine Learning (ICML)*, 2024.

Brian Hu Zhang and Tuomas Sandholm. Sparsified linear programming for zero-sum equilibrium finding. In *International Conference on Machine Learning (ICML)*, 2020.

Martin Zinkevich, Michael Bowling, Michael Johanson, and Carmelo Piccione. Regret minimization in games with incomplete information. In *Proceedings of the Annual Conference on Neural Information Processing Systems (NIPS)*, 2007.

## A Further related work

Besides the pioneering work of [Spielman and Teng \[2004\]](#), which revolved around the simplex algorithm, other prominent algorithms for solving linear programs have also been investigated through the lens of smoothed complexity. [Blum and Dunagan \[2002\]](#) showed that perceptron, a popular algorithm in machine learning, also enjoys a polynomial smoothed complexity (with high probability) for solving linear programming feasibility problems, which can also capture general linear programs via a binary search procedure. Further, [Dunagan et al. \[2011\]](#) performed a smoothed analysis of interior-point methods by relying on an earlier characterization due to [Renegar \[1995\]](#).

Beyond linear programming and (two-player) zero-sum games, there has been a considerable interest in understanding the smoothed complexity of Nash equilibria in general-sum games, but the outlook that has emerged from this endeavor is rather bleak [[Chen et al., 2009](#), [Boodaghians et al., 2020](#), [Rubinstein, 2016](#)]. On a more positive note, [Daskalakis et al. \[2024\]](#) recently considered a more permissive solution concept they refer to as a *smooth Nash equilibrium*; the basic idea of their relaxation is that instead of considering best-response deviations, they restrict to deviations that do not assign too much probability mass on any pure strategy, as controlled by a certain parameter. For a certain regime of that parameter, they obtained positive results, bypassing the intractability of the usual Nash equilibrium. Considering smooth Nash equilibria could also be fruitful in the context of zero-sum games. In particular, we surmise that, if one is content with convergence to smooth Nash equilibria, the error bound could exhibit more favorable properties. Smoothed analysis has also been applied to more structured classes of games, such as congestion or potential games [[Giannakopoulos, 2023](#), [Giannakopoulos et al., 2022](#), [Chen et al., 2020](#)], as well as other important problems in game theory [[Gatti et al., 2013](#), [Buriol et al., 2011](#)]. Other notable developments in a broader context were covered in an older survey by [Spielman and Teng \[2009\]](#); for more recent developments, we point to, for example, [Christ and Yannakakis \[2023\]](#), [Chen et al. \[2024\]](#), [Huiberts et al. \[2023\]](#), and the many references therein.

Average-case analysis has also been a popular topic in the optimization literature [[Cunha et al., 2022](#), [Paquette et al., 2023](#), [Scieur and Pedregosa, 2020](#)], and so it is worth relating our results to that line of work. In particular, let us focus on the recent work of [Cunha et al. \[2022\]](#). First, that paper targets a certain class of convex quadratic problems, whereas we examine zero-sum games. They also operate under a different perturbation model, deriving a parametrization based on the concentration of the eigenvalues of a certain matrix. Further, without strong convexity, [Cunha et al. \[2022\]](#) establish a complexity scaling with  $\text{poly}(1/\epsilon)$ , while here we target the  $\log(1/\epsilon)$  regime. We finally remark that the techniques employed are also quite different. In particular, [Cunha et al. \[2022, Problem 2.1\]](#) posit that the optimal solution does not depend on the underlying randomization. In contrast, as we have already highlighted, the fact that the equilibrium is a function of the randomization constitutes the main technical crux in our setting. At the same time, [Cunha et al. \[2022\]](#) encountered several challenges not present in our setting, so overall those results are complementary.

Beyond smoothed complexity, understanding the last-iterate convergence of gradient-based methods such as OGDA and EGDA has received tremendous interest in the literature; *e.g.*, [[Golowich et al., 2020a](#), [Cai et al., 2022](#), [Gorbunov et al., 2022](#), [Vankov et al., 2023](#), [Golowich et al., 2020b](#), [Mahdavinia et al., 2022](#), [Antonakopoulos et al., 2021](#), [Mertikopoulos et al., 2019](#), [Abe et al., 2023](#)]. It is worth noting that linear convergence has also been documented for the more challenging class of extensive-form games [[Lee et al., 2021](#)], as well as Markov games [[Song et al., 2023](#)]. Nevertheless, there are lower bounds precluding linear convergence beyond affine variational inequalities [[Golowich et al., 2020a](#), [Wei et al., 2021](#)]. We also refer to the works of [Cohen et al. \[2017\]](#) and [Giannou et al. \[2021\]](#) for further characterizations of the convergence rate of no-regret dynamics in multi-player games.

Contrary to the above line of work, which focuses on last-iterate convergence, the most common approach to solving zero-sum games revolves around regret minimization whereby optimality guarantees concern the average strategies. Learning in such settings has been a popular research topic as it captures many central problems; two notable recent applications are learning from multiple distributions [[Haghtalab et al., 2022](#)] and multi-calibration [[Haghtalab et al., 2023](#)]. Yet, there are at least three limitations of the no-regret framework worth highlighting here. The first one, which has been stressed extensively already, is that the number of iterations must grow at least as  $\Omega(1/\epsilon)$  when one insists on taking (uniform) averages [[Daskalakis et al., 2015](#)]. The second and more nuanced caveat is that the no-regret framework does not provide instance-based guarantees based on natural game-theoretic parameters of the problem (see, for example, the discussion of [Maiti](#)

et al. [2023]). Building on earlier work [Wei et al., 2021, Tseng, 1995], some of our results here attempt to address this shortcoming by coming up with a more interpretable parameterization of the iteration complexity of algorithms such as OGDA. The final limitation is that, convergence to the set of equilibria notwithstanding, no-regret guarantees provide no information regarding properties of the equilibrium reached. Although not an issue in non-degenerate zero-sum games, equilibrium selection still remains a central problem. Earlier results [Wei et al., 2021, Tseng, 1995] provide an interesting characterization for the last iterate of OGDA and EGDA by showing that the limit point is the projection of the initial point to the set of equilibria.

Finally, it is worth pointing out the best available theoretical guarantees for solving zero-sum games. Assuming that each entry of  $\mathbf{A}$  has absolute value bounded by 1, (1) can be solved in  $\tilde{O}(\max\{n, m\}^\omega)$  [Cohen et al., 2021] or  $\tilde{O}(nm + \min\{n, m\}^{5/2})$  [van den Brand et al., 2021]. Here,  $\omega$  is the exponent of matrix multiplication and  $\tilde{O}$  suppresses polylogarithmic factors in  $n$  and  $m$ . The complexity we obtain for algorithms such as OGDA is not competitive even though we work in the more benign smoothed complexity model; we reiterate that we did not attempt to optimize the polynomial factors in terms of  $n$  and  $m$ , and those can almost certainly be improved. On the other hand, there are two main aspects in which algorithms such as OGDA are more appealing in terms of their scalability: the per-iteration complexity and the memory requirements. An algorithm such as OGDA requires a single matrix-vector product in each iteration, which can be implemented in linear time for sparse matrices, and has a limited memory footprint. In contrast, implementing interior-point methods in large games can be prohibitive.

## B Preliminaries

In this section, we introduce some further background on smoothed complexity and define the algorithms cited earlier (Items 1 to 4).

**Further notation** For a random variable  $X$ , we denote by  $\mathbb{E}[X]$  its expectation and by  $\mathbb{V}[X]$  its variance, under the assumption that both are finite. For a sequence of random variables  $X_1, \dots, X_d$  and scalars  $\alpha_1, \dots, \alpha_d \in \mathbb{R}$ , linearity of expectation yields that  $\mathbb{E}[\alpha_1 X_1 + \dots + \alpha_d X_d] = \alpha_1 \mathbb{E}[X_1] + \dots + \alpha_d \mathbb{E}[X_d]$ . Assuming independence, it also holds that  $\mathbb{V}[\alpha_1 X_1 + \dots + \alpha_d X_d] = (\alpha_1)^2 \mathbb{V}[X_1] + \dots + (\alpha_d)^2 \mathbb{V}[X_d]$ . We will also use the fact that a linear combination of independent Gaussian random variables is also Gaussian. More broadly, linear combinations can be understood through a convolution in the space of probability density functions, which means that smoothness (in the sense of Lemma C.7) is preserved in a certain regime.

### B.1 Smoothed complexity

To fully specify Definition 1.1, we first recall that a (univariate) Gaussian random variable with zero mean and variance  $\sigma^2$  admits a probability density function of the form

$$\mu : t \mapsto \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{t^2}{2\sigma^2}\right).$$

The law of such a Gaussian random variable will be denoted by  $\mathcal{N}(0, \sigma^2)$ . In the original work of Spielman and Teng [2004], smoothed complexity was defined as the expected running time (or some other cost function) of some algorithm over the perturbed input. More precisely, let  $\mathcal{A}$  be an algorithm whose inputs can be expressed as vectors in  $\mathbb{R}^d$ , and let  $T_{\mathcal{A}}(\mathcal{I})$  be the running time of algorithm  $\mathcal{A}$  on input  $\mathcal{I} \in \mathbb{R}^d$ . Then, the *smoothed complexity* of  $\mathcal{A}$  is

$$\mathcal{C}_{\mathcal{A}}(d, \sigma) := \max_{\mathcal{I} \in \mathbb{R}^d} \mathbb{E}_{\mathbf{g} \sim \mathcal{N}(\mathbf{0}_d, \sigma^2 \mathbf{I}_{d \times d})} [T_{\mathcal{A}}(\mathcal{I} + \|\mathcal{I}\|\mathbf{g})].$$

As pointed out by Spielman and Teng [2003], one does not need to limit smoothed analysis to measure the expected running time, and high probability guarantees are also quite natural; see, for example, the smoothed analysis of the perceptron algorithm due to Blum and Dunagan [2002]. Our main result also provides a guarantee with high probability; it is not clear whether the expected running time can also be bounded by  $\text{poly}(n, m, 1/\sigma)$ , which is left for future work.

### B.2 Algorithms

Next, we specify the algorithms we consider in this work.

**Optimistic gradient descent/ascent** Originally proposed by Popov [1980], optimistic gradient descent/ascent (OGDA)—and variants thereof [Hsieh et al., 2019]—has been recently revived in the online learning literature commencing from the pioneering works of Rakhlin and Sridharan [2013] and Chiang et al. [2012]. If we denote for compactness  $F(\mathbf{z}) := (\mathbf{A}\mathbf{y}, -\mathbf{A}^\top \mathbf{x})$ , OGDA can be expressed as follows for  $t \in \mathbb{N}(\{1, 2, \dots, \})$ .

$$\begin{aligned}\mathbf{z}^{(t)} &:= \Pi_{\mathcal{Z}}(\hat{\mathbf{z}}^{(t)} - \eta F(\mathbf{z}^{(t-1)})), \\ \hat{\mathbf{z}}^{(t+1)} &:= \Pi_{\mathcal{Z}}(\hat{\mathbf{z}}^{(t)} - \eta F(\mathbf{z}^{(t)})).\end{aligned}\tag{OGDA}$$

Here,  $\eta > 0$  is the *learning rate*;  $\Pi_{\mathcal{Z}}(\cdot)$  denotes the (Euclidean) projection operator on set  $\mathcal{Z} := \mathcal{X} \times \mathcal{Y}$ ; and  $\mathbf{z}^{(0)} = \hat{\mathbf{z}}^{(1)} \in \mathcal{Z}$  is the initialization. That is, players simultaneously update their strategies through optimistic gradient steps. Given that  $\mathcal{X}$  and  $\mathcal{Y}$  are probability simplexes, each projection can be computed exactly in nearly linear time. The key reference point for OGDA in affine variational inequalities is the work of Wei et al. [2021] who established linear convergence using the notion of *metric subregularity* (Definition C.9), which is strongly related to Definition 1.3; we discuss their approach later in Appendix C.3.

**Optimistic multiplicative weights update** Deriving from the same class of online learning algorithms as OGDA, optimistic multiplicative weights (OMWU) is the incarnation of *optimistic mirror descent* with an entropic regularizer, namely

$$\begin{aligned}\mathbf{x}^{(t)} &\propto \mathbf{x}^{(t-1)} \circ \exp\left(-2\eta \mathbf{A}\mathbf{y}^{(t-1)} + \eta \mathbf{A}\mathbf{y}^{(t-2)}\right), \\ \mathbf{y}^{(t)} &\propto \mathbf{y}^{(t-1)} \circ \exp\left(2\eta \mathbf{A}^\top \mathbf{x}^{(t-1)} - \eta \mathbf{A}^\top \mathbf{x}^{(t-2)}\right)\end{aligned}\tag{OMWU}$$

for  $t \in \mathbb{N}$ .<sup>5</sup> Above,  $\circ$  denotes the component-wise product; the exponential mapping  $\exp(\cdot)$  is also to be applied component-wise; and  $\mathbf{z}^{(-1)} := \mathbf{z}^{(0)} := (\frac{1}{n}\mathbf{1}_n, \frac{1}{m}\mathbf{1}_m)$ . Daskalakis and Panageas [2019] first proved that OMWU exhibits asymptotic (last-iterate) convergence, and Wei et al. [2021] later established linear convergence.

*Remark B.1.* It is important to note here that the exponential map of OMWU can produce iterates with an arbitrarily large number of bits. Nevertheless, it is not hard to show that the analysis of Wei et al. [2021] carries over when the iterates are truncated up to a certain length of the most significant bits, and so we will not dwell further on this issue here.

**Extra-gradient descent/ascent** The extra-gradient method of Korpelevich [1976] is quite similar to OGDA, namely

$$\begin{aligned}\hat{\mathbf{z}}^{(t)} &:= \Pi_{\mathcal{Z}}(\mathbf{z}^{(t)} - \eta F(\mathbf{z}^{(t)}), \\ \mathbf{z}^{(t+1)} &:= \Pi_{\mathcal{Z}}(\mathbf{z}^{(t)} - \eta F(\hat{\mathbf{z}}^{(t)}))\end{aligned}\tag{EGDA}$$

for  $t \in \mathbb{N}$ . Unlike OGDA, one caveat is that it requires two gradient evaluations per each iteration  $t$ . EGDA is also less suited to use in an online environment: it requires more feedback than what is provided in the online learning setting, and in fact, even legitimate variants of EGDA can still incur substantial regret [Golowich et al., 2020a]. Tseng [1995] first established that EGDA exhibits linear convergence for problems such as (1), discussed further in Appendix C.3.

**Iterative smoothing** This is a refinement of Nesterov’s classical smoothing technique [Nesterov, 2005] due to Gilpin et al. [2012]. Let us first recall the vanilla version of Nesterov, which we refer to as *Smoothing*( $\mathbf{A}, \mathbf{z}^{(0)}, \epsilon$ ):

1. Initialize  $\eta := \frac{\epsilon}{D_{\mathcal{Z}}}$  and  $\hat{\mathbf{z}}^{(0)} := \mathbf{z}^{(0)}$ , where  $D_{\mathcal{Z}}$  is the  $\ell_2$  diameter of  $\mathcal{Z}$ .
2. For  $t = 0, 1, \dots$ 
  - (a)  $\mathbf{u}^{(t)} := \frac{2}{2+t}\hat{\mathbf{z}}^{(t)} + \frac{t}{t+2}\mathbf{z}^{(t)}$ .
  - (b) 
$$\mathbf{z}^{(t+1)} := \arg \min_{\mathbf{z} \in \mathcal{Z}} \left\{ \langle \nabla F_{\eta}(\mathbf{u}^{(t)}), \mathbf{z} - \mathbf{u}^{(t)} \rangle + \frac{L^2}{2\eta} \|\mathbf{z} - \mathbf{u}^{(t)}\|^2 \right\},$$

<sup>5</sup>OMWU is oftentimes expressed via the (optimistic) mirror descent viewpoint, but the form we provide here is easily seen to be equivalent.

where  $F_\eta(\mathbf{z}) := \max_{\widehat{\mathbf{z}} \in \mathcal{Z}} \{ \langle F(\mathbf{z}), \mathbf{z} - \widehat{\mathbf{z}} \rangle - \frac{\eta}{2} \|\mathbf{z} - \widehat{\mathbf{z}}\|^2 \}$  and  $L$  is a suitable matrix norm.

- (c) If  $\Phi(\mathbf{z}^{(t+1)}) < \epsilon$ , **return**.
- (d)

$$\widehat{\mathbf{z}}^{(t+1)} := \arg \min_{\widehat{\mathbf{z}} \in \mathcal{Z}} \left\{ \sum_{\tau=0}^t \frac{\tau+1}{2} \langle \nabla F_\eta(\mathbf{u}^{(\tau)}), \widehat{\mathbf{z}} - \mathbf{u}^{(\tau)} \rangle + \frac{L^2}{2\eta} \|\widehat{\mathbf{z}} - \mathbf{z}^{(0)}\|^2 \right\}.$$

In this context, `IterSmooth(A, z(0), ρ, ε)` is simple refinement of `Smoothing`, which nonetheless attains linear convergence [Gilpin et al., 2012].

1. Let  $\epsilon^{(0)} = F(\mathbf{z}^{(0)})$ .
2. For  $t = 0, 1, \dots$ 
  - (a)  $\epsilon^{(t+1)} := \frac{\epsilon^{(t)}}{\rho}$ .
  - (b)  $\mathbf{z}^{(t+1)} := \text{Smoothing}(\mathbf{A}, \mathbf{z}^{(t)}, \epsilon^{(t+1)})$ .
  - (c) If  $\Phi(\mathbf{z}^{(t+1)}) < \epsilon$ , **return**.

## C Omitted proofs

We dedicate this section to the proofs omitted earlier from the main body.

### C.1 Proofs from Section 3.2

We first point out that degenerates games have measure zero (cf. Spielman and Teng [2003, Proposition 5.1]).

**Lemma C.1.** *For a Gaussian distributed payoff matrix  $\mathbf{A}$  per Definition 1.1, the game is non-degenerate (Definition 3.2) with probability 1 (almost surely).*

Indeed, the set of games with a non unique equilibrium has measure zero [van Damme, 1991, Theorem 3.5.1]. Regarding the characterization in terms of the number of tight inequalities of the corresponding (primal and dual) linear programs, gathered in Definition 3.2, we note that if  $n + 1$  of the inequalities were tight at  $\mathbf{x}^*$ , that would induce a feasible linear system of  $n$  equalities (by eliminating  $v$ ) in  $n - 1$  variables (by eliminating one of the redundant variables); such degeneracies have measure zero, and there are only finitely many possible such degeneracies, leading to Lemma C.1. As a result, in the smoothed complexity model, we can safely assume that the game is non-degenerate.

Now, as we alluded to earlier, establishing Definition 1.3 reduces to showing that for any points  $\mathbf{x} \in \mathcal{X}$  and  $\mathbf{y} \in \mathcal{Y}$ ,

$$\max_{\mathbf{y}' \in \mathcal{Y}} \langle \mathbf{x}, \mathbf{A} \mathbf{y}' \rangle - v \geq \kappa \|\mathbf{x} - \Pi_{\mathcal{X}^*}(\mathbf{x})\| = \kappa \|\mathbf{x} - \mathbf{x}^*\|, \quad (8)$$

$$v - \min_{\mathbf{x}' \in \mathcal{X}} \langle \mathbf{x}', \mathbf{A} \mathbf{y} \rangle \geq \kappa \|\mathbf{y} - \Pi_{\mathcal{Y}^*}(\mathbf{y})\| = \kappa \|\mathbf{y} - \mathbf{y}^*\|. \quad (9)$$

(Definition 1.3 then indeed follows from the obvious fact  $\|\mathbf{x} - \mathbf{x}^*\| + \|\mathbf{y} - \mathbf{y}^*\| \geq \|\mathbf{z} - \mathbf{z}^*\|$ .) Accordingly, our proof of Theorem 3.6 below will focus on lower bounding  $\kappa$  so that (8) holds, and (9) can then be treated similarly.

Before we proceed, let us make some observations regarding transformation (5) we saw earlier. First, one can understand the transformation  $\mathbf{A}_{B,N}^b = \mathbf{T}(\mathbf{Q}^b, \mathbf{b}, \mathbf{c}, d)$  through the equations

$$d = \mathbf{A}_{i,j}; \mathbf{b}_{j'} = -\mathbf{A}_{i,j'} + \mathbf{A}_{i,j}; \mathbf{c}_{i'} = -\mathbf{A}_{i',j} + \mathbf{A}_{i,j}; \mathbf{Q}_{i',j'} = \mathbf{A}_{i',j'} - \mathbf{A}_{i,j'} - \mathbf{A}_{i',j} + \mathbf{A}_{i,j} \quad (10)$$

for all  $(i', j') \in \widetilde{B} \times \widetilde{N}$ . This can easily be derived from (5) by using the fact that  $\widehat{\mathbf{x}}_B = (\widetilde{\mathbf{x}}, 1 - \mathbf{1}^\top \widetilde{\mathbf{x}})$  and  $\widehat{\mathbf{y}}_N = (\widetilde{\mathbf{y}}, 1 - \mathbf{1}^\top \widetilde{\mathbf{y}})$ . From (10), we see that there is a permutation of the rows of  $\mathbf{T}$  that is upper triangular, with every entry being either 1 or  $-1$ . This implies that  $|\det(\mathbf{T})| = 1$ . With a slight abuse of notation, we will write  $\mathbf{T}_{i,j}$  (as opposed to  $\mathbf{T}_{(i,j),:}$ ) to access the  $(i, j)$  row of  $\mathbf{T}$ , so that  $\mathbf{A}_{i,j} = \langle \mathbf{T}_{i,j}, (\mathbf{Q}^b, \mathbf{b}, \mathbf{c}, d) \rangle$ . From (10), we also see that  $\mathbf{T}_{i,j}$  contains at most 4 non-zero entries. In turn, this implies that  $\|\mathbf{T}_{i,j}\| \leq 2$  and  $\|\mathbf{T}_{i,j}\|_1 \leq 4$ . We gather the above observations in the claim below, which will be used in the sequel.

**Claim C.2.** For the (linear) transformation  $\mathbf{T} \in \mathbb{R}^{(BN) \times (BN)}$  given in (10), it holds that  $|\det(\mathbf{T})| = 1$ . Further,  $\|\mathbf{T}_{i,j}\| \leq 2$  and  $\|\mathbf{T}_{i,j}\|_1 \leq 4$  for all  $(i, j) \in B \times N$ .

The point of transformation (5) is that, as we claimed earlier, the spectral properties of matrix  $\mathbf{Q}$  (as opposed to  $\mathbf{A}_{B,N}$ , which is a natural candidate) suffice to capture the difficulty of addressing the second subproblem identified in Section 3.2. In addition, there is a straightforward but convenient characterization of the equilibrium  $(\tilde{\mathbf{x}}^*, \tilde{\mathbf{y}}^*)$  in terms of the transformed game in (5), as stated below.

**Claim C.3.** It holds that  $\mathbf{Q}\tilde{\mathbf{y}}^* = \mathbf{c}$  and  $\mathbf{Q}^\top\tilde{\mathbf{x}}^* = \mathbf{b}$ .

*Proof.* It is clear that the vector  $\mathbf{Q}\tilde{\mathbf{y}}^* - \mathbf{c}$  must have the same value in every coordinate since  $\tilde{\mathbf{x}}^*$  is fully supported and a best response (by assumption). If that entry was positive, then  $\tilde{\mathbf{x}}^*$  would not be a best response since Player  $x$  could profit from removing all the probability mass (which is possible since  $\sum_{i \in \tilde{B}} \tilde{\mathbf{x}}_i^* > 0$ ). If there was a negative entry, Player  $x$  would profit from increasing its probability mass (which is possible since  $\sum_{i \in \tilde{B}} \tilde{\mathbf{x}}_i^* < 1$ ). Similar reasoning yields  $\mathbf{Q}^\top\tilde{\mathbf{x}}^* = \mathbf{b}$ .  $\square$

Having made the above observations, we next prove some lemmas claimed earlier in Section 3.2 which will be used for the proof of Theorem 3.6. First, we give the proof of Lemma 3.4.

**Lemma 3.4.** Let  $\mathbf{c} = \mathbf{Q}\tilde{\mathbf{y}}^* = \sum_{j \in \tilde{N}} \tilde{\mathbf{y}}_j^* \mathbf{Q}_{:,j}$ , and suppose that  $\bar{\mathbf{Q}} \in \mathbb{R}^{\tilde{B} \times \tilde{N}}$  is such that its  $j$ th column is equal to  $\mathbf{Q}_{:,j} - \mathbf{c}$ . Then,

$$\min_{j \in \tilde{N}} \text{dist}(\mathbf{Q}_{:,j}, \text{span}(\mathbf{Q}_{:, \tilde{N} - j})) \leq \left(1 + \frac{|\tilde{N}|}{1 - \sum_{j \in \tilde{N}} \tilde{\mathbf{y}}_j^*}\right) \min_{j \in \tilde{N}} \text{dist}(\bar{\mathbf{Q}}_{:,j}, \text{span}(\bar{\mathbf{Q}}_{:, \tilde{N} - j})).$$

*Proof.* Let  $\tilde{N} \ni j' \in \arg \min_{j \in \tilde{N}} \text{dist}(\bar{\mathbf{Q}}_{:,j}, \text{span}(\bar{\mathbf{Q}}_{:, \tilde{N} - j}))$ . By definition, there is  $\rho \in \mathbb{R}^{\tilde{N} - j'}$  and  $r \in \mathbb{R}^{\tilde{N}}$  with  $\|r\| = 1$  such that

$$\bar{\mathbf{Q}}_{:,j} := - \sum_{j \in \tilde{N} - j'} \tilde{\mathbf{y}}_j^* \mathbf{Q}_{:,j} + (1 - \tilde{\mathbf{y}}_{j'}^*) \mathbf{Q}_{:,j'} = \sum_{j \in \tilde{N} - j'} \rho_j (\mathbf{Q}_{:,j} - \mathbf{c}) + \epsilon r,$$

where  $\epsilon := \min_{j \in \tilde{N}} \text{dist}(\bar{\mathbf{Q}}_{:,j}, \text{span}(\bar{\mathbf{Q}}_{:, \tilde{N} - j}))$ . Rearranging, we have

$$\mathbf{Q}_{:,j'} \overbrace{\left(1 - \tilde{\mathbf{y}}_{j'}^* + \tilde{\mathbf{y}}_{j'}^* \sum_{j \in \tilde{N} - j'} \rho_j\right)}^{\phi_{j'}} + \sum_{j \in \tilde{N} - j'} \mathbf{Q}_{:,j} \overbrace{\left(-\tilde{\mathbf{y}}_j^* - \rho_j + \tilde{\mathbf{y}}_j^* \sum_{j'' \in \tilde{N} - j'} \rho_{j''}\right)}^{\phi_j} = \epsilon r. \quad (11)$$

Now, let us suppose that all coefficients above are such that  $|\phi_j| \leq \epsilon' := \frac{1 - \sum_{j \in \tilde{N}} \tilde{\mathbf{y}}_j^*}{1 - \sum_{j \in \tilde{N}} \tilde{\mathbf{y}}_j^* + |\tilde{N}|}$  for all  $j \in \tilde{N}$ . Then,  $\sum_{j \in \tilde{N}} \phi_j = \pm |\tilde{N}| \epsilon'$  since  $|\sum_{j \in \tilde{N}} \phi_j| \leq \sum_{j \in \tilde{N}} |\phi_j| \leq \epsilon |\tilde{N}|$ , where for convenience we used the notation  $\sum_{j \in \tilde{N}} \phi_j = \pm |\tilde{N}| \epsilon' \iff -|\tilde{N}| \epsilon' \leq \sum_{j \in \tilde{N}} \phi_j \leq |\tilde{N}| \epsilon'$ . Thus, by definition of  $\phi_j$ ,

$$\left(1 - \sum_{j \in \tilde{N}} \tilde{\mathbf{y}}_j^*\right) \left(\sum_{j \in \tilde{N} - j'} \rho_j\right) = \left(1 - \sum_{j \in \tilde{N}} \tilde{\mathbf{y}}_j^*\right) \pm \epsilon' |\tilde{N}|.$$

Since  $0 < 1 - \sum_{j \in \tilde{N}} \tilde{\mathbf{y}}_j^*$ , we have

$$\left(\sum_{j \in \tilde{N} - j'} \rho_j\right) = 1 \pm \epsilon' \frac{|\tilde{N}|}{1 - \sum_{j \in \tilde{N}} \tilde{\mathbf{y}}_j^*}.$$

Thus,

$$\phi_{j'} = 1 - \tilde{\mathbf{y}}_{j'}^* + \tilde{\mathbf{y}}_{j'}^* \sum_{j \in \tilde{N} - j'} \rho_j = 1 \pm \epsilon' \frac{|\tilde{N}|}{1 - \sum_{j \in \tilde{N}} \tilde{\mathbf{y}}_j^*} > \epsilon'$$

since  $\epsilon' \leq \frac{1 - \sum_{j \in \tilde{N}} \mathbf{y}_j^*}{1 - \sum_{j \in \tilde{N}} \mathbf{y}_j^* + |\tilde{N}|}$ . The last displayed inequality contradicts our earlier assumption that  $|\phi_{j'}| \leq \epsilon'$ . As a result, we conclude that at least one coefficient  $\phi_j$  has an absolute value at least  $\epsilon'$ . Dividing (11) by that coefficient, we get

$$\min_{j \in \tilde{N}} \text{dist}(\mathbf{Q}_{:,j}, \text{span}(\mathbf{Q}_{:,\tilde{N}-j})) \leq \frac{\epsilon}{\epsilon'} \leq \left(1 + \frac{|\tilde{N}|}{1 - \sum_{j \in \tilde{N}} \mathbf{y}_j^*}\right) \min_{j \in \tilde{N}} \text{dist}(\bar{\mathbf{Q}}_{:,j}, \text{span}(\bar{\mathbf{Q}}_{:,\tilde{N}-j})).$$

This completes the proof.  $\square$

We continue with the proof of [Lemma 3.5](#).

**Lemma 3.5.** *Let  $\mathbf{M} \in \mathbb{R}^{d \times d}$  be a full-rank matrix. For any  $\tilde{\mathbf{x}} \in \mathbb{R}^d$  there is  $\mathbf{p} \in \mathbb{R}^d$  with  $\|\mathbf{p}\| \leq \frac{1}{\sigma_{\min}(\mathbf{M})} \|\tilde{\mathbf{x}}\|$  such that*

$$\tilde{\mathbf{x}} = \mathbf{M}\mathbf{p} = \sum_{j=1}^d \mathbf{p}_j \mathbf{M}_{:,j}.$$

*Proof.* Let  $\mathbf{M} = \mathbf{U}\Sigma\mathbf{V}^\top$  be a singular value decomposition (SVD) of  $\mathbf{Q}$ , where  $\mathbf{U}$  and  $\mathbf{V}$  are orthonormal. Then, given that  $\mathbf{Q}$  is invertible (by assumption),

$$\mathbf{p} = \mathbf{V}\Sigma^{-1}\mathbf{U}^\top \tilde{\mathbf{x}},$$

where  $\Sigma^{-1} = \text{diag}(\sigma_{\min}^{-1}, \dots, \sigma_{\max}^{-1})$ . (Here,  $\sigma_{\max}$  and  $\sigma_{\min}$  are the maximum and minimum singular values of  $\mathbf{M}$ , respectively.) Thus,  $\|\mathbf{p}\| \leq \|\mathbf{V}\| \|\Sigma^{-1}\| \|\mathbf{U}^\top\| \|\tilde{\mathbf{x}}\| \leq \frac{1}{\sigma_{\min}(\mathbf{Q})} \|\tilde{\mathbf{x}}\|$ , where we used the fact that the spectral norm of any orthonormal matrix is 1 and the spectral norm of any diagonal matrix is its maximum entry in absolute value.  $\square$

We next state the negative second moment identity that connects the smallest singular values in terms of a certain geometric property of the matrix (namely, [Item 3](#)) (see also [\[Tao, 2023\]](#) for further background).

**Proposition C.4** (Negative second moment identity [\[Tao et al., 2010\]](#)). *Let  $\mathbf{M} \in \mathbb{R}^{d \times d}$  be an invertible matrix. Then,*

$$\sum_{r=1}^d \frac{1}{\sigma_r^2(\mathbf{M})} = \sum_{r=1}^d \frac{1}{\text{dist}^2(\mathbf{M}_{r,:}, H_{-r,:})} = \sum_{r=1}^d \frac{1}{\text{dist}^2(\mathbf{M}_{:,r}, H_{:,r})}, \quad (12)$$

where  $H_{-r,:} := \text{span}(\mathbf{M}_{1,:}, \dots, \mathbf{M}_{r-1,:}, \mathbf{M}_{r+1,:}, \dots, \mathbf{M}_{d,:})$ .

One can readily prove this identity by equivalently expressing the negative second moment  $\text{tr}((\mathbf{M}^{-1})^\top \mathbf{M}^{-1})$  as either  $\sum_{r=1}^d \sigma_r^2(\mathbf{M}^{-1}) = \sum_{r=1}^d \sigma_r^{-2}(\mathbf{M})$  or  $\sum_{r=1}^d \|\mathbf{M}_{:,r}^{-1}\|^2$ , leading to the first identity in (12). The second one follows from the fact that the singular values of  $\mathbf{M}^\top$  coincide with the singular values of  $\mathbf{M}$ .

We are now ready to prove [Theorem 3.6](#), restated below.

**Theorem 3.6.** *Let  $\mathbf{A}$  be a non-degenerate payoff matrix, and suppose that  $(\alpha_P(\mathbf{A}), \alpha_D(\mathbf{A}))$ ,  $(\beta_P(\mathbf{A}), \beta_D(\mathbf{A}))$  and  $(\gamma_P(\mathbf{A}), \gamma_D(\mathbf{A}))$  are as in [Definition 3.3](#). Then, the error bound ([Definition 1.3](#)) is satisfied for any sufficiently small modulus*

$$\kappa \gtrsim \frac{1}{\|\mathbf{A}^b\|_\infty \min(n, m)^3} \min \{(\alpha_D(\mathbf{A}))^2 \beta_D(\mathbf{A}) \gamma_P(\mathbf{A}), (\alpha_P(\mathbf{A}))^2 \beta_P(\mathbf{A}) \gamma_D(\mathbf{A})\}.$$

*Proof.* We lower bound  $\kappa$  so that (8) holds; bound (9) will then be treated in a symmetric fashion, and [Definition 1.3](#) will follow.

Let us fix any point  $\mathbf{x} \in \mathcal{X}$ . We can write  $\mathbf{x}$  as  $\lambda \hat{\mathbf{x}}_B + (1 - \lambda) \hat{\mathbf{x}}_{\bar{B}}$  for some  $\lambda \in [0, 1]$  such that  $\hat{\mathbf{x}}_B \in \mathcal{X}$  and all coordinates of  $\hat{\mathbf{x}}_B$  in  $\bar{B}$  are zero, and  $\hat{\mathbf{x}}_{\bar{B}} \in \mathcal{X}$  and all coordinates of  $\hat{\mathbf{x}}_{\bar{B}}$  in  $B$  are zero. For notational convenience, we define

$$P(\mathbf{A}) := \frac{1}{2|N| \sqrt{|B|}} \sigma_{\min}(\bar{\mathbf{Q}}) \left(1 + \frac{1}{\alpha_D(\mathbf{A})}\right)^{-1}. \quad (13)$$

We consider the following two cases.

**Case I:**  $\lambda P(\mathbf{A})\|\widehat{\mathbf{x}}_B - \mathbf{x}_B^*\| \geq 4(1-\lambda)\|\mathbf{A}^\flat\|_\infty$ . If  $\widehat{\mathbf{x}}_B = \mathbf{x}_B^*$ , it follows that  $\mathbf{x} = \mathbf{x}^*$  (since  $\lambda = 1$ ), and the conclusion trivially follows. We can thus assume that  $\widehat{\mathbf{x}}_B \neq \mathbf{x}_B^*$ . In this case, it follows that  $\widetilde{B} \neq \emptyset$ , and we proceed as follows.

$$\max_{\mathbf{y}' \in \mathcal{Y}} \langle \mathbf{x}, \mathbf{A}\mathbf{y}' \rangle - v \geq \lambda \max_{j \in N} \langle \widehat{\mathbf{x}}_B - \mathbf{x}_B^*, \mathbf{A}_{B,j} \rangle + (1-\lambda) \left( \langle \mathbf{x}_{\widetilde{B}}, \mathbf{A}_{\widetilde{B},j} \rangle - v \right) \quad (14)$$

$$\geq \lambda \max_{j \in N} \langle \widehat{\mathbf{x}}_B - \mathbf{x}_B^*, \mathbf{A}_{B,j} \rangle - 2(1-\lambda)\|\mathbf{A}^\flat\|_\infty, \quad (15)$$

where (14) follows from the definition  $\mathbf{x} := \lambda \widehat{\mathbf{x}}_B + (1-\lambda) \widehat{\mathbf{x}}_{\widetilde{B}}$  and the fact that  $v = \langle \mathbf{x}_B^*, \mathbf{A}_{B,j} \rangle$  for all  $j \in N$ ; and (15) uses definition of  $\|\mathbf{A}^\flat\|_\infty$  to lower bound the second term in (14). Continuing from (15), we can use the transformation defined in (5) to get

$$\max_{j \in N} \langle \widehat{\mathbf{x}}_B - \mathbf{x}_B^*, \mathbf{A}_{B,j} \rangle = \max_{j \in N} \langle \widetilde{\mathbf{x}} - \widetilde{\mathbf{x}}^*, \mathbf{Q}_{:,j} - \mathbf{c} \rangle, \quad (16)$$

where, with an abuse of notation, the convention above is that  $\mathbf{Q}_{:,j} = \mathbf{0}$  if  $j \neq \widetilde{N}$ . For convenience, let us define  $\chi_j := \langle \widetilde{\mathbf{x}} - \widetilde{\mathbf{x}}^*, \mathbf{Q}_{:,j} - \mathbf{c} \rangle$  for all  $j \in N$ . Our goal is to lower bound  $\max_{j \in N} \chi_j$ . To that end, we first observe that, by the fact that  $\mathbf{Q}\widetilde{\mathbf{y}}^* = \mathbf{c}$  (Claim C.3),

$$\begin{aligned} 0 &= \langle \widetilde{\mathbf{x}} - \widetilde{\mathbf{x}}^*, \mathbf{Q}\widetilde{\mathbf{y}}^* - \mathbf{c} \rangle = \sum_{j \in \widetilde{N}} \widetilde{\mathbf{y}}_j^* \langle \widetilde{\mathbf{x}} - \widetilde{\mathbf{x}}^*, \mathbf{Q}_{:,j} \rangle - \langle \widetilde{\mathbf{x}} - \widetilde{\mathbf{x}}^*, \mathbf{c} \rangle \\ &= \sum_{j \in \widetilde{N}} \widetilde{\mathbf{y}}_j^* \langle \widetilde{\mathbf{x}} - \widetilde{\mathbf{x}}^*, \mathbf{Q}_{:,j} - \mathbf{c} \rangle + \left( 1 - \sum_{j \in \widetilde{N}} \widetilde{\mathbf{y}}_j^* \right) \langle \widetilde{\mathbf{x}} - \widetilde{\mathbf{x}}^*, -\mathbf{c} \rangle. \end{aligned}$$

In other words,

$$\sum_{j \in N} \mathbf{y}_j^* \chi_j = 0,$$

which in turn implies that

$$\begin{aligned} \sum_{j \in N} \max(0, \chi_j) &\geq \sum_{j \in N} \mathbf{y}_j^* \max(0, \chi_j) = - \sum_{j \in N} \mathbf{y}_j^* \min(0, \chi_j) \\ &\geq -\alpha_D(\mathbf{A}) \sum_{j \in N} \min(0, \chi_j), \end{aligned} \quad (17)$$

where we made use of the obvious identity  $t = \max(0, t) + \min(0, t)$  for all  $t \in \mathbb{R}$ , as well as the definition of  $\alpha_D(\mathbf{A})$  (Item 1). We let  $\mathbf{p} \in \mathbb{R}^{\widetilde{N}}$  be the (unique) solution to the linear system

$$\widetilde{\mathbf{x}} - \widetilde{\mathbf{x}}^* = \mathbf{Q}\mathbf{p} = \sum_{j \in \widetilde{N}} (\mathbf{Q}_{:,j} - \mathbf{c})\mathbf{p}_j,$$

and  $\mathbf{p}_j = 0$  for  $j \in N \setminus \widetilde{N}$ . By Lemma 3.5, we know that  $\|\mathbf{p}\| \leq (\sigma_{\min}(\mathbf{Q}))^{-1} \|\widetilde{\mathbf{x}} - \widetilde{\mathbf{x}}^*\|$ . Then, we have

$$\sum_{j \in N} \chi_j \mathbf{p}_j = \sum_{j \in \widetilde{N}} \chi_j \mathbf{p}_j = \left\langle \widetilde{\mathbf{x}} - \widetilde{\mathbf{x}}^*, \sum_{j \in \widetilde{N}} (\mathbf{Q}_{:,j} - \mathbf{c})\mathbf{p}_j \right\rangle = \|\widetilde{\mathbf{x}} - \widetilde{\mathbf{x}}^*\|^2. \quad (18)$$

Moreover,

$$\begin{aligned} \sum_{j \in N} \chi_j \mathbf{p}_j &= \sum_{j \in N} \mathbf{p}_j \max(0, \chi_j) + \sum_{j \in N} \mathbf{p}_j \min(0, \chi_j) \\ &\leq \sum_{j \in N} \max(0, \mathbf{p}_j) \max(0, \chi_j) + \sum_{j \in N} \min(0, \mathbf{p}_j) \min(0, \chi_j) \end{aligned} \quad (19)$$

$$\leq \|\mathbf{p}\|_\infty \sum_{j \in N} \max(0, \chi_j) - \|\mathbf{p}\|_\infty \sum_{j \in N} \min(0, \chi_j) \quad (20)$$

$$\leq \|\mathbf{p}\|_\infty \left( 1 + \frac{1}{\alpha_D(\mathbf{A})} \right) \sum_{j \in N} \max(0, \chi_j) \quad (21)$$

$$\leq \frac{1}{\sigma_{\min}(\mathbf{Q})} \left( 1 + \frac{1}{\alpha_D(\mathbf{A})} \right) |N| \max_{j \in N} \chi_j \|\widetilde{\mathbf{x}} - \widetilde{\mathbf{x}}^*\|, \quad (22)$$

where (19) follows from the fact that  $\mathbf{p}_j \max(0, \chi_j) \leq \max(0, \mathbf{p}_j) \max(0, \chi_j)$  (by nonnegativity of  $\max(0, \chi_j)$ ) and  $\mathbf{p}_j \min(0, \chi_j) \leq \min(0, \mathbf{p}_j) \min(0, \chi_j)$  (by nonpositivity of  $\min(0, \chi_j)$ ); (20) uses that  $\min(0, \mathbf{p}_j) \geq -|\mathbf{p}_j| \geq -\|\mathbf{p}\|_\infty$ , which gives  $\min(0, \mathbf{p}_j) \min(0, \chi_j) \leq -\|\mathbf{p}\|_\infty \min(0, \chi_j)$ ; (21) follows from (17); and (22) uses the bound  $\|\mathbf{p}\|_2 \leq (\sigma_{\min}(\bar{\mathbf{Q}}))^{-1} \|\tilde{\mathbf{x}} - \tilde{\mathbf{x}}^*\|$  (Lemma 3.5). Combining (18) and (22),

$$\max_{j \in \tilde{N}} \langle \hat{\mathbf{x}}_B - \mathbf{x}_B^*, \mathbf{A}_{B,j} \rangle \geq \frac{1}{|N|} \sigma_{\min}(\bar{\mathbf{Q}}) \left( 1 + \frac{1}{\alpha_D(\mathbf{A})} \right)^{-1} \|\tilde{\mathbf{x}} - \tilde{\mathbf{x}}^*\| \quad (23)$$

$$\geq \frac{1}{2|N|\sqrt{|B|}} \sigma_{\min}(\bar{\mathbf{Q}}) \left( 1 + \frac{1}{\alpha_D(\mathbf{A})} \right)^{-1} \|\hat{\mathbf{x}}_B - \mathbf{x}_B^*\|, \quad (24)$$

where (23) uses the definition of  $\chi_j$  and the assumption that  $\tilde{\mathbf{x}} \neq \tilde{\mathbf{x}}^*$  (equivalently,  $\mathbf{x}_B^* \neq \hat{\mathbf{x}}_B$ ), and (24) follows from the bound

$$\|\hat{\mathbf{x}}_B - \mathbf{x}_B^*\| \leq \|\hat{\mathbf{x}}_B - \mathbf{x}_B^*\|_1 \leq \sum_{i \in \tilde{B}} |\mathbf{x}_i - \mathbf{x}_i^*| + \left| \sum_{i \in \tilde{B}} (\mathbf{x}_i - \mathbf{x}_i^*) \right| \leq 2\|\tilde{\mathbf{x}} - \tilde{\mathbf{x}}^*\|_1 \leq 2\sqrt{|B|} \|\tilde{\mathbf{x}} - \tilde{\mathbf{x}}^*\|.$$

Returning to (15), we have

$$\begin{aligned} \max_{\mathbf{y}' \in \mathcal{Y}} \langle \mathbf{x}, \mathbf{A}\mathbf{y}' \rangle - v &\geq \lambda \frac{1}{2|N|\sqrt{|B|}} \sigma_{\min}(\bar{\mathbf{Q}}) \left( 1 + \frac{1}{\alpha_D(\mathbf{A})} \right)^{-1} \|\hat{\mathbf{x}}_B - \mathbf{x}_B^*\| - 2(1-\lambda) \|\mathbf{A}^\flat\|_\infty \\ &= \lambda P(\mathbf{A}) \|\hat{\mathbf{x}}_B - \mathbf{x}_B^*\| - 2(1-\lambda) \|\mathbf{A}^\flat\|_\infty, \end{aligned} \quad (25)$$

where the equality above follows from the definition of  $P(\mathbf{A})$  in (13). Next, we bound

$$\begin{aligned} \|\mathbf{x} - \mathbf{x}^*\|^2 &= \|\lambda \hat{\mathbf{x}}_B - \mathbf{x}_B^*\|^2 + (1-\lambda)^2 \|\hat{\mathbf{x}}_{\bar{B}}\|^2 \\ &= \|\lambda(\hat{\mathbf{x}}_B - \mathbf{x}_B^*) - (1-\lambda)\mathbf{x}_B^*\|^2 + (1-\lambda)^2 \|\hat{\mathbf{x}}_{\bar{B}}\|^2 \\ &\leq 2\lambda^2 \|\hat{\mathbf{x}}_B - \mathbf{x}_B^*\|^2 + 2(1-\lambda)^2 \|\mathbf{x}_B^*\|^2 + (1-\lambda)^2 \|\hat{\mathbf{x}}_{\bar{B}}\|^2 \\ &\leq 2\lambda^2 \|\hat{\mathbf{x}}_B - \mathbf{x}_B^*\|^2 + 3(1-\lambda)^2, \end{aligned} \quad (26)$$

where (26) uses triangle inequality with respect to  $\|\cdot\|$  along with the inequality  $(t_1 + t_2)^2 \leq 2t_1^2 + 2t_2^2$ , and (27) uses that  $\|\mathbf{x}_B^*\|, \|\hat{\mathbf{x}}_{\bar{B}}\| \leq 1$ . Since we are assuming that  $\lambda P(\mathbf{A}) \|\hat{\mathbf{x}}_B - \mathbf{x}_B^*\| \geq 4(1-\lambda) \|\mathbf{A}^\flat\|_\infty$ , (27) in turn implies that

$$\begin{aligned} \|\mathbf{x} - \mathbf{x}^*\|^2 &\leq 2\lambda^2 \|\hat{\mathbf{x}}_B - \mathbf{x}_B^*\|^2 + \lambda^2 \left( \frac{P(\mathbf{A})}{\|\mathbf{A}^\flat\|_\infty} \right)^2 \|\hat{\mathbf{x}}_B - \mathbf{x}_B^*\|^2 \\ &= \lambda^2 \left( 2 + \left( \frac{P(\mathbf{A})}{\|\mathbf{A}^\flat\|_\infty} \right)^2 \right) \|\hat{\mathbf{x}}_B - \mathbf{x}_B^*\|^2. \end{aligned} \quad (28)$$

Combining (25) and (28) with the assumption that  $\lambda P(\mathbf{A}) \|\hat{\mathbf{x}}_B - \mathbf{x}_B^*\| \geq 4(1-\lambda) \|\mathbf{A}^\flat\|_\infty$ ,

$$\begin{aligned} \max_{\mathbf{y}' \in \mathcal{Y}} \langle \mathbf{x}, \mathbf{A}\mathbf{y}' \rangle - v &\geq \frac{\lambda}{2} P(\mathbf{A}) \|\hat{\mathbf{x}}_B - \mathbf{x}_B^*\| \\ &\geq \frac{1}{2} P(\mathbf{A}) \left( 2 + \left( \frac{P(\mathbf{A})}{\|\mathbf{A}^\flat\|_\infty} \right)^2 \right)^{-2} \|\mathbf{x} - \mathbf{x}^*\| \geq \kappa(\mathbf{A}) \|\mathbf{x} - \mathbf{x}^*\|. \end{aligned}$$

It is easy to see that  $P(\mathbf{A})/\|\mathbf{A}^\flat\|_\infty$  is upper bounded by an absolute constant, and so we have

$$\begin{aligned} \max_{\mathbf{y}' \in \mathcal{Y}} \langle \mathbf{x}, \mathbf{A}\mathbf{y}' \rangle - v &\gtrsim P(\mathbf{A}) \|\mathbf{x} - \mathbf{x}^*\| = \frac{1}{2|N|\sqrt{|B|}} \sigma_{\min}(\bar{\mathbf{Q}}) \left( 1 + \frac{1}{\alpha_D(\mathbf{A})} \right)^{-1} \|\mathbf{x} - \mathbf{x}^*\| \\ &\gtrsim \frac{1}{|B|^3} (\alpha_D(\mathbf{A}))^2 \gamma_P(\mathbf{A}) \|\mathbf{x} - \mathbf{x}^*\|. \end{aligned}$$

Above, the last bound uses the fact that

$$\begin{aligned} \sigma_{\min}(\bar{\mathbf{Q}}) &\geq \frac{1}{\sqrt{|\tilde{B}|}} \min_{j \in \tilde{N}} \text{dist}(\bar{\mathbf{Q}}_{:,j}, \text{span}(\bar{\mathbf{Q}}_{:, \tilde{N}-j})) \\ &\geq \frac{1}{|\tilde{B}|^{3/2}} \min_{j \in \tilde{N}} \text{dist}(\mathbf{Q}_{:,j}, \text{span}(\mathbf{Q}_{:, \tilde{N}-j})) \alpha_D(\mathbf{A}) = \frac{1}{|\tilde{B}|^{3/2}} \gamma_P(\mathbf{A}) \alpha_D(\mathbf{A}), \end{aligned}$$

where the first inequality uses (6), while the second one is a consequence of Lemma 3.4.

**Case II:**  $\lambda P(\mathbf{A})\|\widehat{\mathbf{x}}_B - \mathbf{x}_B^*\| < 4(1 - \lambda)\|\mathbf{A}^\flat\|_\infty$ . This case can only arise when  $\overline{B} \neq \emptyset$  (for otherwise  $\lambda = 1$ ). Then, we bound

$$\begin{aligned} \max_{\mathbf{y}' \in \mathcal{Y}} \langle \mathbf{x}, \mathbf{A} \mathbf{y}' \rangle - v &\geq \langle \mathbf{x}, \mathbf{A} \mathbf{y}^* \rangle - v \\ &\geq \lambda(\langle \widehat{\mathbf{x}}_B - \mathbf{x}_B^*, \mathbf{A}_{B,N} \mathbf{y}_N^* \rangle) + (1 - \lambda)(\langle \widehat{\mathbf{x}}_{\overline{B}}, \mathbf{A}_{\overline{B},N} \mathbf{y}_N^* - v \rangle) \\ &\geq (1 - \lambda)\beta_D(\mathbf{A}), \end{aligned} \quad (29)$$

by definition of  $\beta_D(\mathbf{A})$  (Item 2) and the fact that  $\langle \widehat{\mathbf{x}}_B - \mathbf{x}_B^*, \mathbf{A}_{B,N} \mathbf{y}_N^* \rangle = v \langle \widehat{\mathbf{x}}_B - \mathbf{x}_B^*, \mathbf{1} \rangle = 0$ . Moreover, by (27) together with the assumption that  $\lambda P(\mathbf{A})\|\widehat{\mathbf{x}}_B - \mathbf{x}_B^*\| < 4(1 - \lambda)\|\mathbf{A}^\flat\|_\infty$ ,

$$\|\mathbf{x} - \mathbf{x}^*\|^2 \leq 32 \left( \frac{\|\mathbf{A}^\flat\|_\infty}{P(\mathbf{A})} \right)^2 (1 - \lambda)^2 + 3(1 - \lambda)^2 = \left( 32 \left( \frac{\|\mathbf{A}^\flat\|_\infty}{P(\mathbf{A})} \right)^2 + 3 \right) (1 - \lambda)^2.$$

Combining with (29) yields

$$\begin{aligned} \max_{\mathbf{y}' \in \mathcal{Y}} \langle \mathbf{x}, \mathbf{A} \mathbf{y}' \rangle - v &\geq \left( 32 \left( \frac{\|\mathbf{A}^\flat\|_\infty}{P(\mathbf{A})} \right)^2 + 3 \right)^{-2} \beta_D(\mathbf{A}) \|\mathbf{x} - \mathbf{x}^*\| \\ &\gtrsim \frac{P(\mathbf{A})}{\|\mathbf{A}^\flat\|_\infty} \beta_D(\mathbf{A}) \|\mathbf{x} - \mathbf{x}^*\| \\ &\gtrsim \frac{1}{\|\mathbf{A}^\flat\|_\infty} \frac{1}{|B|^3} \alpha_D(\mathbf{A})^2 \beta_D(\mathbf{A}) \gamma_P(\mathbf{A}) \|\mathbf{x} - \mathbf{x}^*\|. \end{aligned}$$

□

## C.2 Proofs from Section 3.3

We continue with the proofs from Section 3.3. As we have noted already, given that all quantities of interest in Definition 3.3 depend on the support of the equilibrium, it is natural to proceed by partitioning the probability space over all possible such configurations. To do so, we will use the following simple fact [Spielman and Teng, 2003, Proposition 8.1].

**Proposition C.5** (Spielman and Teng, 2003). *Let  $X$  and  $Y$  be random variables distributed according to an integrable density function. For any event  $\mathcal{E}(X, Y)$ ,*

$$\mathbb{P}_{X,Y}[\mathcal{E}(X, Y)] \leq \max_y \mathbb{P}_{X,Y}[\mathcal{E}(X, Y) \mid Y = y] =: \max_Y \mathbb{P}_{X,Y}[\mathcal{E}(X, Y) \mid Y].$$

In our application, we want to condition on the event that  $B$  is the support of  $\mathbf{x}^*$  and  $N$  is the support of  $\mathbf{y}^*$ . For convenience, we let  $\text{Type}_{B,N}(\mathbf{A})$  denote the indicator random variable representing whether  $B$  and  $N$  indeed index the positive coordinates of the equilibrium; that is,  $\text{Type}_{B,N}(\mathbf{A}) := \mathbb{1}\{B = \{i \in [n] : \mathbf{x}_i^*(\mathbf{A}) > 0\} \wedge N = \{j \in [m] : \mathbf{y}_j^*(\mathbf{A}) > 0\}\}$ . Unlike general linear programs, which can be infeasible or unbounded, the linear program induced by a zero-sum game is guaranteed to be primal and dual feasible, no matter the perturbation (under Definition 1.1). We will thus only have to condition on events in which  $B$  and  $N$  are both nonempty. To be able to control the probability density function upon conditioning on  $\text{Type}_{B,N}(\mathbf{A})$ , it will be convenient to perform a certain change of variables, which is described next.

**Change of variables** Let us denote by  $\mathbf{A}_{\overline{B},\overline{N}}$  the entries of  $\mathbf{A}$  excluding those in  $\mathbf{A}_{B,N}$ . We first perform a change of variables from  $\mathbf{A}_{\overline{B},\overline{N}}, \mathbf{A}_{B,N}$  to  $\mathbf{A}_{\overline{B},\overline{N}}, \mathbf{Q}, \mathbf{c}, \mathbf{b}, d$ , which uses the linear transformation  $\mathbf{T}$  associated with (5). With this new set of variables at hand, we can conveniently express  $\mathbf{Q}\tilde{\mathbf{y}}^* = \mathbf{c}$  and  $\mathbf{Q}^\top \tilde{\mathbf{x}}^* = \mathbf{b}$  (Claim C.3). Accordingly, we next perform a change of variables from  $\mathbf{A}_{\overline{B},\overline{N}}, \mathbf{Q}, \mathbf{c}, \mathbf{b}, d$  to  $\mathbf{A}_{\overline{B},\overline{N}}, \mathbf{Q}, \mathbf{x}^*, \mathbf{y}^*, v$ . When performing those change of variables one has to account for the transformed probability density function. This can be understood as follows. The probability of an event  $\mathcal{E}(\mathbf{A})$  can be expressed as

$$\int_{\mathbf{A}} \mathcal{E}(\mathbf{A}) \mu_{\mathbf{A}}(\mathbf{A}) d\mathbf{A}.$$

The integral above can be cast in terms of a new set of variables  $\mathbf{B}$  by computing the corresponding Jacobian, assuming that it is non-singular. We will make use of this fact in the sequel. The following lemma gathers some of the above observations regarding the change of variables.

**Lemma C.6** (Change of variables). *Let  $\mathcal{E}(\mathbf{A})$  be any event that depends on the randomness of  $\mathbf{A}$ . Then,*

$$\begin{aligned}\mathbb{P}_{\mathbf{A}}[\mathcal{E}(\mathbf{A})] &\leq \max_{B,N} \mathbb{P}_{\mathbf{A}}[\mathcal{E}(\mathbf{A}) \mid \text{Type}_{B,N}(\mathbf{A})] \\ &= \max_{B,N} \mathbb{P}_{\mathbf{A}_{\bar{B},\bar{N}}, \mathbf{Q}, \mathbf{x}^*, \mathbf{y}^*, v} [\mathcal{E}(\mathbf{A}) \mid \mathbf{A}_{\bar{B},N} \mathbf{y}_N^* \geq v\mathbf{1} \text{ and } \mathbf{A}_{\bar{N},B}^\top \mathbf{x}_B^* \leq v\mathbf{1}].\end{aligned}$$

Indeed, the first inequality above is a consequence of [Proposition C.5](#). The equality then follows from noting that, when

$$\mathbf{c} = \mathbf{Q}\tilde{\mathbf{y}}^*, \mathbf{b} = \mathbf{Q}^\top \tilde{\mathbf{x}}^*, v = d - \langle \tilde{\mathbf{x}}^*, \mathbf{Q}\tilde{\mathbf{y}}^* \rangle \iff \mathbf{A}_{B,N} \mathbf{y}^* = v\mathbf{1}, \mathbf{A}_{N,B}^\top \mathbf{x}^* = v\mathbf{1},$$

the event  $\text{Type}_{B,N}(\mathbf{A})$  can be equivalently expressed as  $\mathbf{A}_{\bar{B},N} \mathbf{y}_N^* \geq v\mathbf{1}$  and  $\mathbf{A}_{\bar{N},B}^\top \mathbf{x}_B^* \leq v\mathbf{1}$ .

We first bound the probability that  $\beta_P(\mathbf{A}) := \min_{j \in \bar{N}} (v - \langle \mathbf{x}_B^*, \mathbf{A}_{B,j} \rangle)$  is close to 0; the proof for  $\beta_D(\mathbf{A})$  is then symmetric. The key ingredient is the following anti-concentration lemma pertaining to a conditional Gaussian distribution [[Spielman and Teng, 2003, Lemma 8.3](#)].

**Lemma C.7** ([Spielman and Teng, 2003](#)). *Let  $g$  be a Gaussian random variable of variance  $\sigma^2$  and mean of absolute value at most 1. For  $\epsilon \geq 0$ ,  $\tau \geq 1$  and  $t \leq \tau$ ,*

$$\mathbb{P}[g \leq t + \epsilon \mid g \geq t] \leq \frac{\epsilon\tau}{\sigma^2} e^{\frac{\epsilon(\tau+3)}{\sigma^2}}.$$

**Proposition 3.8.** *Let  $\beta_P(\mathbf{A})$  be defined as in [Item 2](#). For any  $\epsilon \geq 0$ ,*

$$\mathbb{P}_{\mathbf{A}} \left[ \beta_P(\mathbf{A}) \leq \frac{\epsilon}{5\|\mathbf{A}^\dagger\|_\infty} \right] \leq \epsilon \frac{e \min(n, m)^2}{\sigma^2}.$$

*Proof.* By [Lemma C.6](#), it suffices to bound

$$\max_{B,N} \mathbb{P}_{\mathbf{A}_{\bar{B},\bar{N}}, \mathbf{Q}, \mathbf{x}^*, \mathbf{y}^*, v} [\beta_P(\mathbf{A}) \leq \epsilon' \mid \mathbf{A}_{\bar{B},N} \mathbf{y}_N^* \geq v\mathbf{1} \text{ and } \mathbf{A}_{\bar{N},B}^\top \mathbf{x}_B^* \leq v\mathbf{1}].$$

By [Proposition C.5](#), it suffices to prove that for all  $B, N, \mathbf{A}_{\bar{B},N}, \mathbf{A}_{\bar{B},\bar{N}}, \mathbf{Q}, \mathbf{x}^*, \mathbf{y}^*, v$  satisfying  $\mathbf{A}_{\bar{B},N} \mathbf{y}_N^* \geq v\mathbf{1}$ ,

$$\mathbb{P}_{\mathbf{A}_{\bar{B},\bar{N}}} [\exists j \in \bar{N} : v - \langle \mathbf{x}_B^*, \mathbf{A}_{B,j} \rangle \leq \epsilon' \mid \forall j \in N : v - \langle \mathbf{x}_B^*, \mathbf{A}_{B,j} \rangle \geq 0] \quad (30)$$

$$\leq \sum_{j \in \bar{N}} \mathbb{P}_{\mathbf{A}_{B,j}} [v - \langle \mathbf{x}_B^*, \mathbf{A}_{B,j} \rangle \leq \epsilon' \mid \forall j \in N : v - \langle \mathbf{x}_B^*, \mathbf{A}_{B,j} \rangle \geq 0] \quad (31)$$

$$\leq \sum_{j \in \bar{N}} \mathbb{P}_{\mathbf{A}_{B,j}} [v - \langle \mathbf{x}_B^*, \mathbf{A}_{B,j} \rangle \leq \epsilon' \mid v - \langle \mathbf{x}_B^*, \mathbf{A}_{B,j} \rangle \geq 0] \quad (32)$$

$$= \sum_{j \in \bar{N}} \mathbb{P}_{\mathbf{g}_j} [\mathbf{g}_j \leq \epsilon' - v \mid \mathbf{g}_j \geq -v]. \quad (33)$$

where in (30) the distribution of  $\mathbf{A}_{\bar{B},\bar{N}}$  after conditioning on  $\mathbf{A}_{\bar{B},N}, \mathbf{A}_{\bar{B},\bar{N}}, \mathbf{Q}, \mathbf{x}^*, \mathbf{y}^*, v$  remains the same, which is a consequence of independence per [Definition 1.1](#); (31) is an application of the union bound; (32) uses the fact that the events  $\{v - \langle \mathbf{x}_B^*, \mathbf{A}_{B,j} \rangle \geq 0\}_{j \in N}$  are pairwise independent; and (33) defines  $\mathbf{g}_j := -\langle \mathbf{x}_B^*, \mathbf{A}_{B,j} \rangle$ , which is a Gaussian random variable with expectation  $|\mathbb{E}[\mathbf{g}_j]| \leq \max_{i \in B} |\mathbf{A}_{i,j}|$  and variance  $\mathbb{V}[\mathbf{g}_j] = \sum_{i \in B} (\mathbf{x}_i^*)^2 \mathbb{V}[\mathbf{A}_{i,j}] = \sigma^2 \sum_{i \in B} (\mathbf{x}_i^*)^2$  (by independence). In particular, by Cauchy-Schwarz,  $\mathbb{V}[\mathbf{g}_j] \geq \frac{1}{|B|} \sigma^2$ . Further, by [Lemma C.7](#) (for  $\tau = \max(1, |v|/|\mathbb{E}[\mathbf{g}_j]|)$ ), we have

$$\begin{aligned}\mathbb{P}_{\mathbf{g}_j} [\mathbf{g}_j \leq \epsilon' - v \mid \mathbf{g}_j \geq -v] &\leq \epsilon' \frac{\max(|v|, |\mathbb{E}[\mathbf{g}_j]|)}{\mathbb{V}[\mathbf{g}_j]} e^{\frac{\max(4|\mathbb{E}[\mathbf{g}_j]|, 3|\mathbb{E}[\mathbf{g}_j]| + |v|)}{\mathbb{V}[\mathbf{g}_j]}} \\ &\leq \epsilon' \frac{\min(n, m) \max(|v|, |\mathbb{E}[\mathbf{g}_j]|)}{\sigma^2} e^{\frac{\min(n, m) \max(4|\mathbb{E}[\mathbf{g}_j]|, 3|\mathbb{E}[\mathbf{g}_j]| + |v|)}{\sigma^2}}\end{aligned}$$

for any  $\epsilon' \geq 0$  and  $j \in \bar{N}$ , where we note that we applied [Lemma C.7](#) for  $\mathbf{g}_j / |\mathbb{E}[\mathbf{g}_j]|$  (since the absolute value of the mean has to be at most 1), which has variance  $\mathbb{V}[\mathbf{g}_j] / (\mathbb{E}[\mathbf{g}_j])^2$ . So, setting  $\epsilon := \epsilon'(|v| + 4|\mathbb{E}[\mathbf{g}_j]|)$ ,

$$\begin{aligned} \mathbb{P}_{\mathbf{g}_j} \left[ \mathbf{g}_j \leq \frac{\epsilon}{|v| + 4 \max_{i \in B} |\mathbf{A}_{i,j}|} - v \mid \mathbf{g}_j \geq -v \right] &\leq \mathbb{P}_{\mathbf{g}_j} \left[ \mathbf{g}_j \leq \frac{\epsilon}{|v| + 4|\mathbb{E}[\mathbf{g}_j]|} - v \mid \mathbf{g}_j \geq -v \right] \\ &\leq \epsilon \frac{\min(n, m)}{\sigma^2} e^{\frac{\min(n, m)}{\sigma^2}}. \end{aligned} \quad (34)$$

Now, when  $\epsilon \frac{\min(n, m)}{\sigma^2} > 1$  the proposition is vacuously true, while in the contrary case the claim follows from (34) and (33).  $\square$

Next, we proceed with the bound on  $\gamma_P(\mathbf{A})$ . The key ingredient is the observation that a random variable with a slowly changing density function cannot be too concentrated on any any interval ([Lemma 3.7](#) due to [Spielman and Teng \[2003, Lemma 8.2\]](#); we restate it below for convenience). Gaussian random variables have this property, as pointed out by [Spielman and Teng \[2003, Lemma 8.1\]](#).

**Lemma C.8** ([Spielman and Teng, 2003](#)). *Let  $\mu$  be the probability density function of a Gaussian random variable in  $\mathbb{R}^d$  of variance  $\sigma^2$  centered at a point of norm at most 1. If  $\text{dist}(\mathbf{r}, \mathbf{r}') \leq \epsilon \leq 1$ , then*

$$\frac{\mu(\mathbf{r}')}{\mu(\mathbf{r})} \geq e^{-\frac{\epsilon(\|\mathbf{r}\|+2)}{\sigma^2}}.$$

**Lemma 3.7** ([Spielman and Teng, 2003](#)). *Let  $\rho$  be the probability density function of a random variable  $X$ . If there exist  $\delta > 0$  and  $c \in (0, 1]$  such that*

$$0 \leq t \leq t' \leq \delta \implies \frac{\rho(t')}{\rho(t)} \geq c, \quad (7)$$

then

$$\mathbb{P}[X \leq \epsilon \mid X \geq 0] \leq \frac{\epsilon}{c\delta}.$$

**Proposition 3.9.** *Let  $\gamma_P(\mathbf{A})$  be defined as in [Item 3](#). For any  $\epsilon \geq 0$ ,*

$$\mathbf{A} \left[ \gamma_P(\mathbf{A}) \leq \frac{\epsilon}{4 \max_{j \in \bar{N}} \|\mathbf{Q}_{:,j}\| + 20 \|\mathbf{A}^b\|_\infty + 3} \right] \leq \epsilon \frac{4e \min(n, m)^3}{\sigma^2}.$$

*Proof.* Let  $\mu_{\mathbf{A}}(\mathbf{A})$  be the probability density function of  $\mathbf{A}$ , which, by independence ([Definition 1.1](#)), can be expressed as  $\prod_{i \in [n], j \in [m]} \mu_{\mathbf{A}_{i,j}}$ , where  $\mu_{\mathbf{A}_{i,j}}$  is a Gaussian random variable. We first perform a change of variables from  $\mathbf{A}_{\bar{B}, \bar{N}}$ ,  $\mathbf{A}_{B, \bar{N}}$  to  $\mathbf{A}_{\bar{B}, \bar{N}}$ ,  $\mathbf{Q}$ ,  $\mathbf{b}$ ,  $\mathbf{c}$ ,  $\mathbf{d}$ , in accordance with (5); this can be understood through the (non-singular; [Claim C.2](#)) linear transformation  $\mathbf{A}_{B, \bar{N}}^b = \mathbf{T}(\mathbf{Q}^b, \mathbf{b}, \mathbf{c}, \mathbf{d})$ . To express the density in the new variables, we first note that the Jacobian of the change of variables is  $|\det(\mathbf{T})| = 1$  ([Claim C.2](#)), and so the density on  $\mathbf{Q}$ ,  $\mathbf{b}$ ,  $\mathbf{c}$ ,  $\mathbf{d}$  can be expressed as  $\mu_{\mathbf{A}_{B, \bar{N}}}(\mathbf{T}(\mathbf{Q}^b, \mathbf{b}, \mathbf{c}, \mathbf{d})) \mu_{\mathbf{A}_{\bar{B}, \bar{N}}}(\mathbf{A}_{\bar{B}, \bar{N}})$ .

Next, we perform a change of variables from  $\mathbf{A}_{\bar{B}, \bar{N}}$ ,  $\mathbf{Q}$ ,  $\mathbf{b}$ ,  $\mathbf{c}$ ,  $\mathbf{d}$  to  $\mathbf{A}_{\bar{B}, \bar{N}}$ ,  $\mathbf{Q}$ ,  $\tilde{\mathbf{x}}^*$ ,  $\tilde{\mathbf{y}}^*$ ,  $v$  according to the transformations  $\mathbf{Q}\tilde{\mathbf{y}}^* = \mathbf{c}$ ;  $\mathbf{Q}^\top \tilde{\mathbf{x}}^* = \mathbf{b}$ ; and  $v = d - \langle \tilde{\mathbf{x}}^*, \mathbf{Q}\tilde{\mathbf{y}}^* \rangle$ . It is easy to see that the Jacobian of the change of variables is  $|\det(\mathbf{T})| = 1$  ([Claim C.2](#)), and so the density on  $\mathbf{Q}$ ,  $\tilde{\mathbf{x}}^*$ ,  $\tilde{\mathbf{y}}^*$ ,  $v$  reads

$$\left| \det \left( \frac{\partial(\mathbf{A}_{\bar{B}, \bar{N}}, \mathbf{Q}, \mathbf{b}, \mathbf{c}, \mathbf{d})}{\partial(\mathbf{A}_{\bar{B}, \bar{N}}, \mathbf{Q}, \tilde{\mathbf{x}}^*, \tilde{\mathbf{y}}^*, v)} \right) \right| = \left| \det \left( \frac{\partial(\mathbf{b}, \mathbf{c}, \mathbf{d})}{\partial(\tilde{\mathbf{x}}^*, \tilde{\mathbf{y}}^*, v)} \right) \right| = \det(\mathbf{Q})^2.$$

So, the density on  $\mathbf{A}_{\bar{B}, \bar{N}}$ ,  $\mathbf{Q}$ ,  $\tilde{\mathbf{x}}^*$ ,  $\tilde{\mathbf{y}}^*$ ,  $v$  reads

$$\mu_{\mathbf{A}_{B, \bar{N}}}(\mathbf{T}(\mathbf{Q}^b, \mathbf{Q}^\top \tilde{\mathbf{x}}^*, \mathbf{Q}\tilde{\mathbf{y}}^*, v + \langle \tilde{\mathbf{x}}^*, \mathbf{Q}\tilde{\mathbf{y}}^* \rangle)) \mu_{\mathbf{A}_{\bar{B}, \bar{N}}}(\mathbf{A}_{\bar{B}, \bar{N}}) \det(\mathbf{Q})^2.$$

By [Lemma C.6](#), it suffices to upper bound

$$\max_{B, \bar{N}} \mathbb{P}_{\mathbf{A}_{\bar{B}, \bar{N}}, \mathbf{Q}, \tilde{\mathbf{x}}^*, \tilde{\mathbf{y}}^*, v} [\gamma_P(\mathbf{A}) \leq \epsilon \mid \mathbf{A}_{\bar{B}, \bar{N}} \mathbf{y}_N^* \geq v \mathbf{1} \text{ and } \mathbf{A}_{\bar{N}, B}^\top \mathbf{x}_B^* \leq v \mathbf{1}].$$

Further, by [Proposition C.5](#), it is in turn enough to bound  $\mathbb{P}_{\mathbf{Q}}[\gamma_P(\mathbf{A}) \leq \epsilon]$  for all  $B, N$  (for the non-trivial case where  $\tilde{B}, \tilde{N} \neq \emptyset$ ),  $\mathbf{A}_{\tilde{B}, \tilde{N}}$ ,  $\tilde{\mathbf{x}}^*$ ,  $\tilde{\mathbf{y}}^*$ ,  $v$  such that  $\mathbf{A}_{\tilde{B}, N} \mathbf{y}_N^* \geq v \mathbf{1}$  and  $\mathbf{A}_{\tilde{N}, B}^\top \mathbf{x}_B^* \leq v \mathbf{1}$ , where the induced distribution on  $\mathbf{Q}$  is

$$\mu_{\mathbf{A}_{B, N}}(\mathbf{T}(\mathbf{Q}^b, \mathbf{Q}^\top \tilde{\mathbf{x}}^*, \mathbf{Q} \tilde{\mathbf{y}}^*, v + \langle \tilde{\mathbf{x}}^*, \mathbf{Q} \tilde{\mathbf{y}}^* \rangle)) \det(\mathbf{Q})^2.$$

We will prove that for any  $j \in \tilde{N}$  and  $\mathbf{Q}_{:, \tilde{N}-j}$ ,

$$\mathbb{P}_{\mathbf{Q}_{:, j}} \left[ \text{dist}(\mathbf{Q}_{:, j}, \text{span}(\mathbf{Q}_{:, \tilde{N}-j})) \leq \frac{\epsilon}{4\|\mathbf{Q}_{:, j}\| + 4|v| + 4\|\mathbf{Q}_{:, \tilde{N}-j}^b\|_\infty + 3} \right] \leq \epsilon \frac{4e \min(n, m)^2}{\sigma^2}, \quad (35)$$

and then apply a union bound over  $j \in \tilde{N}$ . Having fixed  $\mathbf{Q}_{:, \tilde{N}-j}$ , we can express  $\mathbf{Q}_{:, j}$  as  $\mathbf{q}^\parallel + t\mathbf{q}^\perp$ , where  $\mathbb{R}^{\tilde{B}} \ni \mathbf{q}^\parallel \in \text{span}(\mathbf{Q}_{:, \tilde{N}-j})$  and  $\mathbb{R}^{\tilde{B}} \ni \mathbf{q}^\perp$  is the unit vector orthogonal to  $\text{span}(\mathbf{Q}_{:, \tilde{N}-j})$ . Then,  $|t| = \text{dist}(\mathbf{Q}_{:, j}, \text{span}(\mathbf{Q}_{:, \tilde{N}-j}))$  and  $|\det(\mathbf{Q})| = tC(\mathbf{Q}_{:, \tilde{N}-j})$ , where  $C(\mathbf{Q}_{:, \tilde{N}-j})$  does not depend on  $\mathbf{Q}_{:, j}$  (this can be obtained by expressing the determinant using the formula for parallelepipeds). By symmetry, we can prove (35) by bounding the probability that  $t$  is at most  $\epsilon$  given that  $t$  is at least 0. We can thus focus on proving

$$\max_{\mathbf{q}^\parallel \in \text{span}(\mathbf{Q}_{:, \tilde{N}-j})} \frac{\mathbb{P}[t \leq \epsilon \mid t \geq 0]}{t} \leq \epsilon \frac{4e \min(n, m)^2 (4\|\mathbf{q}^\parallel\|_\infty + 4|v| + 4\|\mathbf{Q}_{:, \tilde{N}-j}^b\|_\infty + 3)}{\sigma^2}, \quad (36)$$

and then (35) follows from the fact that  $\|\mathbf{Q}_{:, j}\| \geq \|\mathbf{q}^\parallel\|$ . Now, the induced distribution on  $t$  is proportional to

$$\rho(t) := t^2 \prod_{(i, j) \in B \times N} \mu_{\mathbf{A}_{i, j}}(\langle \mathbf{T}_{i, j}, \mathbf{r}_{i, j}(t) \rangle)$$

for  $\mathbf{r}_{i, j}(t)$  defined as

$$\begin{aligned} & (\mathbf{q}^\parallel + t\mathbf{q}^\perp, \mathbf{Q}_{:, \tilde{N}-j}^b, \mathbf{Q}_{\tilde{N}-j, :}^\top, \tilde{\mathbf{x}}^*, \langle \tilde{\mathbf{x}}^*, \mathbf{q}^\parallel + t\mathbf{q}^\perp \rangle, \mathbf{Q}_{:, \tilde{N}-j} \tilde{\mathbf{y}}_{\tilde{N}-j}^* + \tilde{\mathbf{y}}_j^* \langle \mathbf{q}^\parallel + t\mathbf{q}^\perp \rangle, \\ & \quad v + \langle \tilde{\mathbf{x}}^*, \mathbf{Q}_{:, \tilde{N}-j} \tilde{\mathbf{y}}_{\tilde{N}-j}^* \rangle + \tilde{\mathbf{y}}_j^* \langle \tilde{\mathbf{x}}^*, \mathbf{q}^\parallel + t\mathbf{q}^\perp \rangle). \end{aligned}$$

We now want to apply [Lemma 3.7](#). To that end, we have

$$|\langle \mathbf{T}_{i, j}, \mathbf{r}_{i, j}(t) - \mathbf{r}_{i, j}(t') \rangle|^2 \leq \|\mathbf{T}_{i, j}\|^2 \|\mathbf{r}_{i, j}(t) - \mathbf{r}_{i, j}(t')\|^2$$

$$\leq 4(t - t')^2 \|\langle \mathbf{q}^\perp, \langle \tilde{\mathbf{x}}^*, \mathbf{q}^\perp \rangle, \tilde{\mathbf{y}}_j^* \mathbf{q}^\perp, \tilde{\mathbf{y}}_j^* \langle \tilde{\mathbf{x}}^*, \mathbf{q}^\perp \rangle \rangle\|^2 \quad (37)$$

$$\leq 16(t - t')^2, \quad (38)$$

where (37) follows from  $\|\mathbf{T}_{i, j}\|_2 \leq 2$  ([Claim C.2](#)), and (38) follows from the fact that  $\|\mathbf{q}^\perp\|, \|\tilde{\mathbf{x}}^*\|, \|\tilde{\mathbf{y}}^*\| \leq 1$ . Moreover, again by [Claim C.2](#),

$$|\langle \mathbf{T}_{i, j}, \mathbf{r}_{i, j}(t) \rangle| \leq \|\mathbf{T}_{i, j}\|_1 \|\mathbf{r}_{i, j}(t)\|_\infty \leq 4(\|\mathbf{q}^\parallel\|_\infty + |v| + \|\mathbf{Q}_{:, \tilde{N}-j}^b\|_\infty + t).$$

Let  $0 \leq t \leq t' \leq \delta \leq \frac{1}{4}$  for  $\delta = \frac{\sigma^2}{4|B||N|(4\|\mathbf{q}^\parallel\| + 4|v| + 4\|\mathbf{Q}_{:, \tilde{N}-j}^b\|_\infty + 3)}$ . [Lemma C.8](#) then implies that

$$\frac{\mu_{\mathbf{A}_{i, j}}(\langle \mathbf{T}_{i, j}, \mathbf{r}_{i, j}(t') \rangle)}{\mu_{\mathbf{A}_{i, j}}(\langle \mathbf{T}_{i, j}, \mathbf{r}_{i, j}(t) \rangle)} \geq e^{-\frac{1}{|B||N|}}.$$

Thus,

$$\frac{\rho(t')}{\rho(t)} \geq \left(\frac{t'}{t}\right)^2 \prod_{(i, j) \in B \times N} \frac{\mu_{\mathbf{A}_{i, j}}(\langle \mathbf{T}_{i, j}, \mathbf{r}_{i, j}(t') \rangle)}{\mu_{\mathbf{A}_{i, j}}(\langle \mathbf{T}_{i, j}, \mathbf{r}_{i, j}(t) \rangle)} \geq e^{-1}.$$

We conclude that (36) can be obtained from [Lemma 3.7](#), and the theorem follows.  $\square$

Finally, we bound the probability that  $\alpha_P(\mathbf{A})$  ([Item 1](#)) is close to 0;  $\alpha_D(\mathbf{A})$  can be bounded in a similar fashion.

**Proposition 3.10.** Let  $\alpha_P(\mathbf{A})$  be defined as in Item 1. For any  $\epsilon \geq 0$ ,

$$\mathbb{P}_{\mathbf{A}} \left[ \alpha_P(\mathbf{A}) \leq \frac{\epsilon}{25(\|\mathbf{A}^\flat\|_\infty + 1)^2} \right] \leq \epsilon \frac{8e^2 mn \min(n, m)}{\sigma^2}.$$

*Proof.* By Lemma C.6, it suffices to bound

$$\max_{B, N} \mathbb{P}_{\mathbf{A}_{\overline{B}, \overline{N}}, \mathbf{Q}, \mathbf{x}^*, \mathbf{y}^*, v} [\alpha_P(\mathbf{A}) \leq \epsilon \mid \mathbf{A}_{\overline{B}, N} \mathbf{y}_N^* \geq v \mathbf{1} \text{ and } \mathbf{A}_{\overline{N}, B}^\top \mathbf{x}_B^* \leq v \mathbf{1}],$$

where we recall that the induced probability density function on  $\mathbf{A}_{\overline{B}, \overline{N}}, \mathbf{Q}, \mathbf{x}^*, \mathbf{y}^*, v$  reads

$$\mu_{\mathbf{A}_{B, N}}(\mathbf{T}(\mathbf{Q}^\flat, \mathbf{Q}^\top \tilde{\mathbf{x}}^*, \mathbf{Q} \tilde{\mathbf{y}}^*, v + \langle \tilde{\mathbf{x}}^*, \mathbf{Q} \tilde{\mathbf{y}}^* \rangle)) \mu_{\mathbf{A}_{B, \overline{N}}}(\mathbf{A}_{B, \overline{N}}) \mu_{\mathbf{A}_{\overline{B}, N}}(\mathbf{A}_{\overline{B}, N}) \mu_{\mathbf{A}_{\overline{B}, \overline{N}}}(\mathbf{A}_{\overline{B}, \overline{N}}) \det(\mathbf{Q})^2.$$

We consider the non-trivial case where  $\tilde{B}, \tilde{N} \neq \emptyset$ . We will perform a further change of variables. Namely, let  $\mathbf{a} = \mathbf{A}_{\overline{N}, i}$  for  $i \in B \setminus \tilde{B}$ . We map  $\mathbf{A}_{B, \overline{N}}$  to  $\bar{\mathbf{A}}_{\tilde{B}, \overline{N}} := \mathbf{A}_{\tilde{B}, \overline{N}} - \mathbf{1} \mathbf{a}^\top$ ,  $\mathbf{a}$ , so that  $\mathbf{A}_{\overline{N}, B}^\top \mathbf{x}_B^* \leq v \mathbf{1}$  can be equivalently expressed as  $\bar{\mathbf{A}}_{\overline{N}, \tilde{B}}^\top \tilde{\mathbf{x}}^* \leq v \mathbf{1} - \mathbf{a}$ . The induced density function is now proportional to

$$\mu_{\mathbf{A}_{B, N}}(\mathbf{T}(\mathbf{Q}^\flat, \mathbf{Q}^\top \tilde{\mathbf{x}}^*, \mathbf{Q} \tilde{\mathbf{y}}^*, v + \langle \tilde{\mathbf{x}}^*, \mathbf{Q} \tilde{\mathbf{y}}^* \rangle)) \mu_{\mathbf{a}}(\mathbf{a}) \mu_{\mathbf{A}_{\tilde{B}, \overline{N}}}(\bar{\mathbf{A}}_{\tilde{B}, \overline{N}} + \mathbf{1} \mathbf{a}^\top) \nu(\cdot),$$

where  $\nu(\cdot)$  does not depend on  $\tilde{\mathbf{x}}^*$  and  $\mathbf{a}$ . By Proposition C.5, it is enough to show that for any  $B, N, \bar{\mathbf{A}}_{\tilde{B}, \overline{N}}, \mathbf{A}_{\overline{B}, N}, \mathbf{A}_{\overline{B}, \overline{N}}, \mathbf{Q}, \mathbf{y}^*, v$  satisfying  $\mathbf{A}_{\overline{B}, N} \mathbf{y}^* \geq v \mathbf{1}$ ,

$$\begin{aligned} \mathbb{P}_{\tilde{\mathbf{x}}^*, \mathbf{a}} \left[ \alpha_P \leq \frac{\epsilon}{\max((\|\mathbf{Q}^\flat\|_\infty + 1)^2, (1 + \|\bar{\mathbf{A}}_{\tilde{B}, \overline{N}}^\flat\|_\infty)(5\|\bar{\mathbf{A}}_{\tilde{B}, \overline{N}}^\flat\|_\infty + |v| + 4))} \mid \bar{\mathbf{A}}_{\overline{N}, \tilde{B}}^\top \tilde{\mathbf{x}}^* \leq v \mathbf{1} - \mathbf{a} \right] \\ \leq \epsilon \frac{8e^2 mn \min(n, m)}{\sigma^2}, \end{aligned}$$

where the induced distribution on  $\tilde{\mathbf{x}}^*$  and  $\mathbf{a}$  is proportional to

$$\mu_{\mathbf{A}_{B, N}}(\mathbf{T}(\mathbf{Q}^\flat, \mathbf{Q}^\top \tilde{\mathbf{x}}^*, \mathbf{Q} \tilde{\mathbf{y}}^*, v + \langle \tilde{\mathbf{x}}^*, \mathbf{Q} \tilde{\mathbf{y}}^* \rangle)) \mu_{\mathbf{a}}(\mathbf{a}) \mu_{\mathbf{A}_{\tilde{B}, \overline{N}}}(\bar{\mathbf{A}}_{\tilde{B}, \overline{N}} + \mathbf{1} \mathbf{a}^\top). \quad (39)$$

We see that  $\tilde{\mathbf{x}}^*$  is independent of  $\mathbf{a}$  and  $\{\mathbf{a}_j\}_{j \in \overline{N}}$  are pairwise independent. Thus, conditioning on the event  $\bar{\mathbf{A}}_{\overline{N}, \tilde{B}}^\top \tilde{\mathbf{x}}^* \leq v \mathbf{1} - \mathbf{a}$ , the induced distribution on  $\tilde{\mathbf{x}}^*$  is proportional to

$$\mu_{\mathbf{A}_{B, N}}(\mathbf{T}(\mathbf{Q}^\flat, \mathbf{Q}^\top \tilde{\mathbf{x}}^*, \mathbf{Q} \tilde{\mathbf{y}}^*, v + \langle \tilde{\mathbf{x}}^*, \mathbf{Q} \tilde{\mathbf{y}}^* \rangle)) \prod_{j \in \overline{N}} \mathbb{P}_{\mathbf{a}_j} [\langle \bar{\mathbf{A}}_{\tilde{B}, j}, \tilde{\mathbf{x}}^* \rangle \leq v - \mathbf{a}_j].$$

We can proceed by showing that for any fixed  $i \in \tilde{B}$  and  $\tilde{\mathbf{x}}_{\tilde{B} - i}^*$ ,

$$\begin{aligned} \mathbb{P}_{\tilde{\mathbf{x}}_i^*} \left[ \tilde{\mathbf{x}}_i^* \leq \frac{\epsilon}{\max((\|\mathbf{Q}^\flat\|_\infty + 1)^2, (1 + \|\bar{\mathbf{A}}_{\tilde{B}, N}^\flat\|_\infty)(5\|\bar{\mathbf{A}}_{\tilde{B}, N}^\flat\|_\infty + |v| + 4))} \mid \bar{\mathbf{A}}_{\overline{N}, \tilde{B}}^\top \tilde{\mathbf{x}}^* \leq v \mathbf{1} - \mathbf{a} \right] \\ \leq \epsilon \frac{8e^2 m \min(n, m)}{\sigma^2}, \end{aligned}$$

and then applying the union bound over all  $i \in \tilde{B}$ . Having fixed  $\tilde{\mathbf{x}}_{\tilde{B} - i}^*$ , the induced density on  $\tilde{\mathbf{x}}_i^*$ , say  $\rho(t)$ , is proportional to  $\rho_1(t) \cdot \rho_2(t)$ , where

$$\rho_1(t) := \mu_{\mathbf{A}_{B, N}}(\mathbf{T}(\mathbf{Q}^\flat, \mathbf{Q}_{:, \tilde{B} - i}^\top \tilde{\mathbf{x}}_{\tilde{B} - i}^* + t \mathbf{Q}_{:, i}^\top, \mathbf{Q} \tilde{\mathbf{y}}^*, v + \langle \tilde{\mathbf{x}}_{\tilde{B} - i}^*, \mathbf{Q}_{\tilde{B} - i, :} \tilde{\mathbf{y}}^* \rangle + t \langle \mathbf{Q}_{i, :}, \tilde{\mathbf{y}}^* \rangle))$$

and

$$\rho_2(t) := \prod_{j \in \overline{N}} \mathbb{P}_{\mathbf{a}_j} [\langle \bar{\mathbf{A}}_{j, \tilde{B} - i}, \tilde{\mathbf{x}}_{\tilde{B} - i}^* \rangle + \bar{\mathbf{A}}_{i, j} t \leq v - \mathbf{a}_j].$$

We will first apply Lemma 3.7 to bound  $\rho_1(t')/\rho_1(t)$  for  $0 \leq t \leq t' \leq \delta \leq 1$  and a sufficiently small  $\delta$ . We define

$$r_{i,j}(t) := (\mathbf{Q}^\flat, \mathbf{Q}_{:, \tilde{B} - i}^\top \tilde{\mathbf{x}}_{\tilde{B} - i}^* + t \mathbf{Q}_{:, i}^\top, \mathbf{Q} \tilde{\mathbf{y}}^*, v + \langle \tilde{\mathbf{x}}_{\tilde{B} - i}^*, \mathbf{Q}_{\tilde{B} - i, :} \tilde{\mathbf{y}}^* \rangle + t \langle \mathbf{Q}_{i, :}, \tilde{\mathbf{y}}^* \rangle),$$

so that  $\rho_1(t) = \prod_{(i,j) \in B \times N} \mu_{\mathbf{A}_{i,j}}(\langle \mathbf{T}_{i,j}, \mathbf{r}_{i,j}(t) \rangle)$ . Then, we have

$$|\langle \mathbf{T}_{i,j}, \mathbf{r}_{i,j}(t) - \mathbf{r}_{i,j}(t') \rangle| \leq 4|t - t'| \|\mathbf{Q}^b\|_\infty,$$

where we used [Claim C.2](#). Further,

$$|\langle \mathbf{T}_{i,j}, \mathbf{r}_{i,j}(t) \rangle| \leq (t+1) \|\mathbf{Q}^b\|_\infty,$$

and so [Lemma C.8](#) implies that for  $\delta \leq \frac{1}{4\|\mathbf{Q}^b\|_\infty}$ ,

$$\frac{\mu_{\mathbf{A}_{i,j}}(\langle \mathbf{T}_{i,j}, \mathbf{r}_{i,j}(t') \rangle)}{\mu_{\mathbf{A}_{i,j}}(\langle \mathbf{T}_{i,j}, \mathbf{r}_{i,j}(t) \rangle)} \geq e^{-\frac{8\delta \|\mathbf{Q}^b\|_\infty (\|\mathbf{Q}^b\|_\infty + 1)}{\sigma^2}}.$$

As a result, for  $\delta \leq \frac{\sigma^2}{8|B||N|\|\mathbf{Q}^b\|_\infty(\|\mathbf{Q}^b\|_\infty + 1)}$ ,

$$\frac{\rho_1(t')}{\rho_1(t)} = \prod_{(i,j) \in B \times N} \frac{\mu_{\mathbf{A}_{i,j}}(\langle \mathbf{T}_{i,j}, \mathbf{r}_{i,j}(t') \rangle)}{\mu_{\mathbf{A}_{i,j}}(\langle \mathbf{T}_{i,j}, \mathbf{r}_{i,j}(t) \rangle)} \geq e^{-1}.$$

Next, we focus on lower bounding  $\rho_2(t')/\rho_2(t)$ . From (39), it is not hard to see that  $\mathbf{a}_j$  is a Gaussian random variable with expectation  $|\mathbb{E}[\mathbf{a}_j]| \leq 1 + \|\bar{\mathbf{A}}_{\tilde{B}, \bar{N}}^b\|_\infty$  and variance  $\mathbb{V}[\mathbf{a}_j] \geq \frac{\sigma^2}{\min(n, m)}$ . Also,

$$\begin{aligned} \frac{\rho_2(t')}{\rho_2(t)} &= \prod_{j \in \bar{N}} \frac{\mathbb{P}_{\mathbf{a}_j}[\langle \bar{\mathbf{A}}_{\tilde{B}-i,j}, \tilde{\mathbf{x}}_{\tilde{B}-i}^* \rangle + \bar{\mathbf{A}}_{i,j}t' \leq v - \mathbf{a}_j]}{\mathbb{P}_{\mathbf{a}_j}[\langle \bar{\mathbf{A}}_{\tilde{B}-i,j}, \tilde{\mathbf{x}}_{\tilde{B}-i}^* \rangle + \bar{\mathbf{A}}_{i,j}t \leq v - \mathbf{a}_j]} \\ &\geq \prod_{j \in \bar{N}} \mathbb{P}_{\mathbf{a}_j}[\langle \bar{\mathbf{A}}_{\tilde{B}-i,j}, \tilde{\mathbf{x}}_{\tilde{B}-i}^* \rangle + \bar{\mathbf{A}}_{i,j}t' \leq v - \mathbf{a}_j \mid \langle \bar{\mathbf{A}}_{\tilde{B}-i,j}, \tilde{\mathbf{x}}_{\tilde{B}-i}^* \rangle + \bar{\mathbf{A}}_{i,j}t \leq v - \mathbf{a}_j]. \end{aligned}$$

By [Lemma C.7](#) (for  $\tau = (2\|\bar{\mathbf{A}}_{\tilde{B}, \bar{N}}^b\|_\infty + |v| + 1)/(1 + \|\bar{\mathbf{A}}_{\tilde{B}, \bar{N}}^b\|_\infty)$ ),

$$\begin{aligned} \mathbb{P}_{\mathbf{a}_j}[\langle \bar{\mathbf{A}}_{\tilde{B}-i,j}, \tilde{\mathbf{x}}_{\tilde{B}-i}^* \rangle + \bar{\mathbf{A}}_{i,j}t' \leq v - \mathbf{a}_j \mid \langle \bar{\mathbf{A}}_{\tilde{B}-i,j}, \tilde{\mathbf{x}}_{\tilde{B}-i}^* \rangle + \bar{\mathbf{A}}_{i,j}t \leq v - \mathbf{a}_j] \\ \geq 1 - \delta \frac{\min(n, m) \|\bar{\mathbf{A}}_{\tilde{B}, \bar{N}}^b\|_\infty (2\|\bar{\mathbf{A}}_{\tilde{B}, \bar{N}}^b\|_\infty + |v| + 1)}{\sigma^2} e^{\frac{\min(n, m) \|\bar{\mathbf{A}}_{\tilde{B}, \bar{N}}^b\|_\infty (5\|\bar{\mathbf{A}}_{\tilde{B}, \bar{N}}^b\|_\infty + |v| + 4)}{\sigma^2}}. \end{aligned}$$

Thus, for  $\delta \leq \frac{1}{2em} \frac{\sigma^2}{\min(n, m) \|\bar{\mathbf{A}}_{\tilde{B}, \bar{N}}^b\|_\infty (5\|\bar{\mathbf{A}}_{\tilde{B}, \bar{N}}^b\|_\infty + |v| + 4)}$ ,

$$\mathbb{P}_{\mathbf{a}_j}[\langle \bar{\mathbf{A}}_{\tilde{B}-i,j}, \tilde{\mathbf{x}}_{\tilde{B}-i}^* \rangle + \bar{\mathbf{A}}_{i,j}t' \leq v - \mathbf{a}_j \mid \langle \bar{\mathbf{A}}_{\tilde{B}-i,j}, \tilde{\mathbf{x}}_{\tilde{B}-i}^* \rangle + \bar{\mathbf{A}}_{i,j}t \leq v - \mathbf{a}_j] \geq 1 - \frac{1}{2m},$$

which in turn implies that

$$\frac{\rho_2(t')}{\rho_2(t)} \geq \left(1 - \frac{1}{2m}\right)^{\bar{N}} \geq e^{-1}.$$

We conclude that  $\frac{\rho(t')}{\rho(t)} \geq e^{-2}$ , and the proof follows from [Lemma 3.7](#) by lower bounding the value of  $\delta$ .  $\square$

Armed with [Propositions 3.8 to 3.10](#), [Theorem 1.4](#) can be obtained from [Theorem 3.6](#), in conjunction with a union bound and the fact that  $\|\mathbf{A}^b\|_\infty \leq \text{poly}(n, m)$  with high probability (by Gaussian concentration).

### C.3 Proof of Theorem 1.2

Having established [Theorem 1.4](#), here we explain how existing results imply [Theorem 1.2](#). We first focus on OGDA. We also take the opportunity to explain in more detail how [Wei et al. \[2021\]](#) established [Definition 1.3](#), which was sketched earlier in [Section 3.1](#). Our treatment of the rest of the algorithms will be more brief.

**Metric subregularity** A central ingredient in the approach of Wei et al. [2021] is what they refer to as saddle-point *metric subregularity*, stated below as [Definition C.9](#). For the sake of generality, we give the definition for a general objective function  $f : \mathcal{X} \times \mathcal{Y} \ni (\mathbf{x}, \mathbf{y}) \mapsto f(\mathbf{x}, \mathbf{y})$ , assumed to be continuously differentiable; (1) corresponds to the bilinear case  $f(\mathbf{x}, \mathbf{y}) = \langle \mathbf{x}, \mathbf{A}\mathbf{y} \rangle$ . We use again the notation  $F(\mathbf{z}) := (\nabla_{\mathbf{x}} f(\mathbf{x}, \mathbf{y}), -\nabla_{\mathbf{y}} f(\mathbf{x}, \mathbf{y}))$ , where  $\mathbb{R}^{n+m} \ni \mathbf{z} := (\mathbf{x}, \mathbf{y})$ . We also let  $L \in \mathbb{R}_{>0}$  be a Lipschitz continuity parameter for  $F$  with respect to  $\|\cdot\|$ , so that  $\|F(\mathbf{z}) - F(\mathbf{z}')\| \leq L\|\mathbf{z} - \mathbf{z}'\|$ ; in the context of (1), one can always take  $L := \|\mathbf{A}\|$ .

**Definition C.9** (Metric subregularity for saddle-point problems [Wei et al., 2021]). A saddle-point problem satisfies *metric subregularity* if there exists a problem-dependent parameter  $\kappa' \in \mathbb{R}_{>0}$  such that for any  $\mathbf{z} \in \mathcal{Z}$  and  $\mathbf{z}^* := \Pi_{\mathcal{Z}^*}(\mathbf{z})$ ,

$$\sup_{\mathbf{z}' \in \mathcal{Z}} \frac{\langle F(\mathbf{z}), \mathbf{z} - \mathbf{z}' \rangle}{\|\mathbf{z} - \mathbf{z}'\|} \geq \kappa' \|\mathbf{z} - \mathbf{z}^*\|. \quad (40)$$

The nomenclature of [Definition C.9](#) can be justified by the fact that (40) is equivalent to a common type of metric subregularity [Wei et al., 2021, Appendix F]; for more background, we refer to Dontchev and Rockafellar [2009]. We further remark that Wei et al. [2021] introduced (40) in a more general form by allowing an exponent  $\beta \in \mathbb{R}_{\geq 0}$  in the right-hand side, but that additional flexibility is not relevant for our purposes.<sup>6</sup>

Now, there is an obvious connection between [Definition 1.3](#) and [Definition C.9](#) in bilinear problems with bounded domain; namely, we have

$$\sup_{\mathbf{z}' \in \mathcal{Z}} \frac{\langle F(\mathbf{z}), \mathbf{z} - \mathbf{z}' \rangle}{\|\mathbf{z} - \mathbf{z}'\|} \geq \frac{1}{2} \Phi(\mathbf{z}),$$

where we used the fact that  $\langle F(\mathbf{z}), \mathbf{z} \rangle = 0$  and  $\|\mathbf{z} - \mathbf{z}'\| \leq D_{\mathcal{Z}} = 2$ . So, [Definition 1.3](#) with respect to parameter  $\kappa$  implies [Definition C.9](#) with parameter  $\kappa' := \kappa/2$ .

**Linear convergence of OGDA** Under metric subregularity, in the sense of [Definition C.9](#), Wei et al. [2021] were able to establish that OGDA converges to the set  $\mathcal{Z}^*$  at a linear rate:

**Theorem C.10** (Wei et al., 2021). *Consider a saddle-point problem (1) satisfying metric subregularity with respect to some  $\kappa' \in \mathbb{R}_{>0}$ . For any  $\eta \leq \frac{1}{8L}$ , the iterates  $(\mathbf{z}^{(\tau)})_{1 \leq \tau \leq t}$  of OGDA satisfy*

$$\text{dist}(\mathbf{z}^{(t)}, \mathcal{Z}^*) \leq 8 \left( 1 + \frac{16\eta^2(\kappa')^2}{81} \right)^{-t/2} \text{dist}(\widehat{\mathbf{z}}^{(1)}, \mathcal{Z}^*). \quad (41)$$

As a result, [Theorem C.10](#) implies that OGDA guarantees  $\text{dist}(\mathbf{z}^{(t)}, \mathcal{Z}^*) \leq \epsilon$  so long as

$$t \geq 2 \left\lceil \frac{\log \left( \frac{8D_{\mathcal{Z}}}{\epsilon} \right)}{\log \left( 1 + \frac{(\kappa')^2}{324\|\mathbf{A}\|^2} \right)} \right\rceil. \quad (42)$$

In conjunction with [Theorem 3.6](#) and [Propositions 3.8](#) to [3.10](#), this immediately implies that OGDA has a polynomial smoothed complexity with high probability, as claimed earlier in [Theorem 1.2](#).

Before we proceed, it is instructive to explain how Wei et al. [2021] treated the error bound in bilinear problems where  $\mathcal{X}$  and  $\mathcal{Y}$  are polyhedral sets. As we explained earlier, it is enough to show that for any  $\mathbf{x} \in \mathcal{X}$  and  $\mathbf{y} \in \mathcal{Y}$ ,

$$\begin{aligned} \max_{\mathbf{y} \in \mathcal{Y}} \mathbf{x}^\top \mathbf{A}\mathbf{y} - v &\geq \kappa \|\mathbf{x} - \Pi_{\mathcal{X}^*}(\mathbf{x})\|, \\ v - \min_{\mathbf{x} \in \mathcal{X}} \mathbf{x}^\top \mathbf{A}\mathbf{y} &\geq \kappa \|\mathbf{y} - \Pi_{\mathcal{Y}^*}(\mathbf{y})\|. \end{aligned}$$

We focus on the first inequality, which is with respect to Player  $x$ . We let  $\mathcal{X} := \{\mathbf{x} \in \mathbb{R}^n : \mathbf{c}_i^\top \mathbf{x} \leq b_i \quad \forall i \in [\ell_x]\}$ , where  $\ell_x$  denotes the number of vertices of  $\mathcal{X}$ . We also let  $\mathbf{o}_j := \mathbf{A}\mathbf{y}_j$ , where  $\mathbf{y}_j$  denotes the  $j$ th vertex of  $\mathcal{Y}$ ; for simplicity, we will denote by  $k_y \in \mathbb{N}$  the number of vertices of  $\mathcal{Y}$ . We consider a fixed  $\mathbf{x} \in \mathcal{X} \setminus \mathcal{X}^*$  and  $\mathbf{x}^* = \Pi_{\mathcal{X}^*}(\mathbf{x})$ .

<sup>6</sup>Wei et al. [2021] impose (40) only for points  $\mathbf{z} \in \mathcal{Z} \setminus \mathcal{Z}^*$ , which is easily seen to be equivalent.

It is easy to see that the set of optimal strategies for Player  $x$ ,  $\mathcal{X}^* := \{\mathbf{x} \in \mathcal{X} : \max_{\mathbf{y} \in \mathcal{Y}} \langle \mathbf{x}, \mathbf{A}\mathbf{y} \rangle \leq v\}$ , can be expressed as

$$\mathcal{X}^* := \{\mathbf{x} \in \mathbb{R}^n : \mathbf{c}_i^\top \mathbf{x} \leq b_i, \mathbf{o}_j^\top \mathbf{x} \leq v \quad \forall (i, j) \in [\ell_x] \times [k_y]\}.$$

Indeed, any point  $\mathbf{y} \in \mathcal{Y}$  is a convex combination of the vertices of  $\mathcal{Y}$ , and the converse direction is also obvious. A feasibility constraint  $i \in [\ell_x]$  is said to be *tight* if  $\mathbf{c}_i^\top \mathbf{x}^* = b_i$ ; similarly, an optimality constraint  $j \in [k_y]$  is tight if  $\mathbf{o}_j^\top \mathbf{x}^* = v$ . We let  $L_x = L_x(\mathbf{x}^*) \subseteq [\ell_x]$  be the set of tight feasibility constraints and  $K_y = K_y(\mathbf{x}^*) \subseteq [k_y]$  be the set of tight optimality constraints. We can assume without any loss that  $L_x, K_y \neq \emptyset$ . It is well-known (e.g., [Rockafellar, 2015]) that the *normal cone* of  $\mathcal{X}^*$  at  $\mathbf{x}^*$  with respect to  $\mathcal{X}^*$  can be expressed as

$$N_{\mathbf{x}^*} := \left\{ \sum_{i \in L_x} p_i \mathbf{c}_i + \sum_{j \in K_y} q_j \mathbf{o}_j \quad \forall (\mathbf{p}, \mathbf{q}) \in \mathbb{R}_{\geq 0}^{L_x} \times \mathbb{R}_{\geq 0}^{K_y} \right\}.$$

Wei et al. [2021] also define  $M_{\mathbf{x}^*} \subseteq N_{\mathbf{x}^*}$  as

$$N_{\mathbf{x}^*} \cap \{\mathbf{c}_i^\top \mathbf{x} \leq 0 \quad \forall i \in L_x\}.$$

Now, the main parameter of interest that relates to [Definition 1.3](#) in the analysis of Wei et al. [2021] stems from the following quantity.

**Definition C.11.** We let  $C \in \mathbb{R}_{>0}$  be defined as the infimum over  $(0, \infty)$  so that

$$\left\{ \sum_{i \in L_x} p_i \mathbf{c}_i + \sum_{j \in K_y} q_j \mathbf{o}_j, 0 \leq p_i, q_j \leq C \right\} \supseteq M_{\mathbf{x}^*} \cap \mathcal{B}_\infty, \quad (43)$$

where  $\mathcal{B}_\infty \subseteq \mathbb{R}^n$  is the set of points with  $\ell_\infty$  norm upper bounded by 1.

By definition of  $M_{\mathbf{x}^*}$ , it is evident that there always exists a finite problem-dependent parameter  $C \in \mathbb{R}_{>0}$  such that [Definition C.11](#) is satisfied. It is then not hard to show that

$$\max_{\mathbf{y} \in \mathcal{Y}} \mathbf{x}^\top \mathbf{A}\mathbf{y} - v \geq \frac{1}{C|K_y|} \|\mathbf{x} - \Pi_{\mathcal{X}}(\mathbf{x}^*)\|.$$

Assuming that the number of vertices is polynomial in the dimensions,<sup>7</sup> this shows that [Definition C.11](#) essentially captures the complexity of satisfying [Definition 1.3](#). As we explained earlier in [Section 3.1](#), the constraint matrix of the linear program induced by [Definition C.11](#) depends both on the payoff matrix  $\mathbf{A}$  as well as the set of constraints. It is thus unclear how to use existing results in the model of smoothed complexity [Dunagan et al., 2011] to bound  $C$ . The second and more important challenge revolves around the fact that [Definition C.11](#) depends solely on the tight constraints of the optimal solution, which in turn depends on the randomness of  $\mathbf{A}$ . Under our characterization, the latter challenge was addressed earlier in [Section 3.3](#).

Continuing for OMWU, we again rely on the analysis of Wei et al. [2021], which relates the rate of convergence of OMWU to three quantities. The first one [Wei et al., 2021, Definition 3] is similar to [Definition C.9](#), but with the difference that the maximization is now constrained to be over points whose support is a subset of the support of the equilibrium; namely,

$$\kappa_x := \min_{\mathbf{x} \in \mathcal{X} \setminus \{\mathbf{x}^*\}} \max_{\mathbf{y} \in \mathcal{V}^*(\mathcal{Y})} \frac{\langle \mathbf{x} - \mathbf{x}^*, \mathbf{A}\mathbf{y} \rangle}{\|\mathbf{x} - \mathbf{x}^*\|_1}, \quad (44)$$

where  $\mathcal{V}^*(\mathcal{Y}) := \{\mathbf{y} \in \Delta^m : \text{supp}(\mathbf{y}) \subseteq \text{supp}(\mathbf{y}^*)\}$ . A symmetric definition is to be considered with respect to Player  $y$ . To connect this to (8), we note that, when  $\mathbf{y} \in \mathcal{V}^*(\mathcal{Y})$ ,  $\langle \mathbf{x}^*, \mathbf{A}\mathbf{y} \rangle = v$ . We are thus left to lower bound  $\max_{\mathbf{y}} \langle \mathbf{x}, \mathbf{A}\mathbf{y} \rangle - v$  in terms of  $\|\mathbf{x} - \mathbf{x}^*\|_1$ , but under the constraint that  $\mathbf{y} \in \mathcal{V}^*(\mathcal{Y})$ . An inspection of our proof of [Theorem 3.6](#) (and in particular the proof of (8)) reveals that its conclusion holds even when the maximization is subject to the above constraint, and so our analysis

<sup>7</sup>In fact, by virtue of Carathéodory's theorem, one can refine [Definition C.11](#) so that this holds even when the number of vertices is exponential in the dimensions. Namely, a point  $\mathbf{v} \in M_{\mathbf{x}^*} \cap \mathcal{B}_\infty$  can be written as the conical combination of at most  $n$  of the vectors describing the cone in (43), thereby maintaining feasibility. This observation can be used to refine the (worst-case) analysis of Wei et al. [2021] to, for example, extensive-form games wherein the number of vertices is typically exponential in the dimensions.

immediately lower bounds (44) as well. The second quantity introduced by Wei et al. [2021, Definition 2] corresponds exactly to **Item 2**, which was bounded in **Proposition 3.8**. The third quantity [Wei et al., 2021, Definition 4] is where the exponential overhead is introduced. Namely, the iteration complexity of OMWU in their analysis depends on  $\exp(\min(\alpha_P(\mathbf{A}), \alpha_D(\mathbf{A}))^{-1})$ , where we recall the definition in **Item 1**.<sup>8</sup> Unfortunately, for any game, it holds that  $\alpha_P(\mathbf{A}) \leq 1/n$  and  $\alpha_D(\mathbf{A}) \leq 1/m$ , and so even if the geometry of the problem is favorable, the obtained bound is exponential. (The reason the above quantity is crucial in their analysis is because it lower bounds the probability of playing any action through the trajectory of OMWU.) Nevertheless, using **Proposition 3.10**, our analysis provides instead a bound of  $\exp(\text{poly}(n, m, 1/\sigma))$  with high probability, which is still a major improvement over the worst-case bound of Wei et al. [2021], which can be doubly exponential in the number of bits  $L$  describing the game—one can easily make sure that  $\alpha_P(\mathbf{A}) \approx 1/2^L$  (**Proposition 3.1**).

Next, for EGDA, Tseng [1995] established linear convergence under the error bound

$$\text{dist}(\mathbf{z}, \mathbf{z}^*) \leq \tau \|\mathbf{z} - \Pi_{\mathcal{Z}}(\mathbf{z} - \eta F(\mathbf{z}))\|$$

for some  $\tau > 0$  and a suitable  $\eta > 0$  [Tseng, 1995, Corollary 3.3]. It is easy to make the following connection.

**Lemma C.12.** *It holds that  $\Phi(\mathbf{z}) \leq \frac{2}{\eta} \|\mathbf{z} - \Pi_{\mathcal{Z}}(\mathbf{z} - \eta F(\mathbf{z}))\|$ .*

*Proof.* Indeed, by the first-order optimality condition for the optimization problem associated with

$$\mathbf{z}' := \Pi_{\mathcal{Z}}(\mathbf{z} - \eta F(\mathbf{z})) = \arg \min_{\mathbf{z}' \in \mathcal{Z}} \{\|\mathbf{z}' - (\mathbf{z} - \eta F(\mathbf{z}))\|^2 := h(\mathbf{z}')\},$$

we get  $\langle \widehat{\mathbf{z}} - \mathbf{z}', \nabla h(\mathbf{z}') \rangle \geq 0$  for any  $\widehat{\mathbf{z}} \in \mathcal{Z}$ , or equivalently,  $\min_{\widehat{\mathbf{z}} \in \mathcal{Z}} \langle \widehat{\mathbf{z}} - \mathbf{z}', \mathbf{z}' - \mathbf{z} + \eta F(\mathbf{z}) \rangle \geq 0$ . Observing that  $\min_{\widehat{\mathbf{z}} \in \mathcal{Z}} \langle \widehat{\mathbf{z}}, F(\mathbf{z}) \rangle = -\Phi(\mathbf{z})$  and bounding

$$\langle \mathbf{z} - \mathbf{z}', \widehat{\mathbf{z}} - \mathbf{z}' \rangle \geq -\|\mathbf{z} - \mathbf{z}'\| \|\widehat{\mathbf{z}} - \mathbf{z}'\| \geq -D_{\mathcal{Z}} \|\mathbf{z} - \mathbf{z}'\| = -2 \|\mathbf{z} - \Pi_{\mathcal{Z}}(\mathbf{z} - \eta F(\mathbf{z}))\|$$

leads to the claim.  $\square$

It can thus be shown that **Definition 1.3** is again sufficient to dictate the rate of convergence of EGDA.

Finally, for **IterSmooth**, Gilpin et al. [2012] introduced a “condition measure” of the payoff matrix  $\mathbf{A}$ , which in fact corresponds precisely to **Definition 1.3**. Thus, **Theorem 1.2** with respect to **IterSmooth** follows readily from [Gilpin et al., 2012, Theorem 2].

#### C.4 Proof of Theorem 4.2

Finally, we conclude with the proof of **Theorem 4.2**, which is restated below.

**Theorem 4.2.** *Any  $\delta$ -support-stable game (per **Definition 4.1**) satisfies the error bound for any sufficiently small modulus*

$$\kappa \geq \text{poly}\left(\frac{1}{n}, \frac{1}{m}, \delta\right).$$

*Proof of Theorem 4.2.* We treat each parameter separately.

- Let us start from  $\beta_P(\mathbf{A})$  (**Item 2**). We let  $j' \in \arg \min_{j \in \overline{N}} (v - \langle \mathbf{x}_B^*, \mathbf{A}_{B,j} \rangle)$ , where we assume that  $\overline{N} \neq \emptyset$ . We consider a perturbed matrix  $\mathbf{A}'$  such that

$$\mathbf{A}'_{i,j} = \begin{cases} \mathbf{A}_{i,j} - \beta_P(\mathbf{A}) & \text{if } i \in B, j = j', \\ \mathbf{A}_{i,j} & \text{otherwise.} \end{cases}$$

<sup>8</sup>More specifically, the proof of Wei et al. [2021, Theorem 3] upper bounds the Kullback-Leibler divergence  $\text{KL}(\mathbf{z}^{(t)}, \mathbf{z}^*)$  by a quantity that is at least as large as  $\left(1 + \frac{15\eta^2 C_2}{32}\right)^{-t}$ , where  $C_2 \leq \exp(\min(\alpha_P(\mathbf{A}), \alpha_D(\mathbf{A}))^{-1})$ . Thus, to guarantee  $\text{KL}(\mathbf{z}^{(t)}, \mathbf{z}^*) \leq \epsilon$  using the analysis of Wei et al. [2021] one needs at least  $\log(1/\epsilon)/\log\left(1 + \frac{15\eta^2 C_2}{32}\right)$  iterations. When  $C_2 \ll 1$ , this grows with  $1/C_2 \geq \exp(\min(\alpha_P(\mathbf{A}), \alpha_D(\mathbf{A})))$ .

Then, the game described by  $\mathbf{A}'$  cannot be non-degenerate with the same support as  $\mathbf{A}$ . Indeed, in the contrary case it would follow that the (unique) equilibrium  $(\mathbf{x}_B^*, \mathbf{y}_N^*)$  remains the same since  $\mathbf{A}'_{B,N} = \mathbf{A}_{B,N}$ . But then,  $v - \langle \mathbf{x}_B^*, \mathbf{A}'_{B,j'} \rangle = v - \langle \mathbf{x}_B^*, \mathbf{A}_{B,j'} \rangle - \beta_P(\mathbf{A}) = 0$ , by definition of  $j'$  and  $\beta_P(\mathbf{A})$ , which is a contradiction. Further,  $\|\mathbf{A} - \mathbf{A}'\| = \beta_P(\mathbf{A})$ . In turn, this implies that  $\delta \leq \beta_P(\mathbf{A})$ . Similar reasoning yields that  $\delta \leq \beta_D(\mathbf{A})$ .

- Continuing for  $\gamma_P(\mathbf{A})$  (Item 3), we assume that  $\tilde{B}, \tilde{N} \neq \emptyset$ . We let  $\mathbf{U}\Sigma\mathbf{V}^\top$  be a singular value decomposition (SVD) of  $\mathbf{Q}$ . Then, a perturbation to  $\mathbf{Q}$  of the form  $\mathbf{U}\text{diag}(0, 0, \dots, \sigma_{\min}(\mathbf{Q}))\mathbf{V}^\top$  leads to a singular matrix  $\mathbf{Q}'$ , which cannot be the case if the perturbed game is non-degenerate with the same support. This perturbation can be cast in terms of  $\mathbf{A}'_{B,N}$  through transformation  $\mathbf{T}$  in (5). This lower bounds  $\sigma_{\min}(\mathbf{Q})$  in terms of  $\delta$ , and Proposition C.4 can in turn lower bound  $\gamma_P(\mathbf{A})$  in terms of  $\sigma_{\min}(\mathbf{Q})$ .
- Finally, we treat  $\alpha_P(\mathbf{A})$  (Item 1). The non-trivial case is again when  $\tilde{B}, \tilde{N} \neq \emptyset$ . Let  $i' \in \arg \min_{i \in B} (\mathbf{x}_i^*)$ . If  $i' \in \tilde{B}$ , we define

$$\mathbb{R}^{\tilde{B}} \ni \tilde{\mathbf{x}}'_i = \begin{cases} 0 & \text{if } i = i', \\ \mathbf{x}_i^* & \text{otherwise.} \end{cases}$$

We know that  $\mathbf{Q}^\top \tilde{\mathbf{x}}^* = \mathbf{b}$ . We then consider the perturbed vector  $\mathbf{b}' := \mathbf{Q}^\top \tilde{\mathbf{x}}'$ . If the perturbed game was non-degenerate with the same support, it would follow that  $(\tilde{\mathbf{x}}', \cdot)$  is the unique equilibrium, which is a contradiction since  $\tilde{\mathbf{x}}'_{i'} = 0$ . Further, the norm of the perturbation  $\|\mathbf{b} - \mathbf{b}'\|$  is upper bounded in terms of  $\alpha_P(\mathbf{A})$ , which can be again expressed in terms of  $\mathbf{A}_{B,N}$  through transformation (5). Similarly, if  $i' \notin \tilde{B}$ , we define

$$\mathbb{R}^{\tilde{B}} \ni \tilde{\mathbf{x}}'_i = \mathbf{x}_i^* + \frac{\alpha_P(\mathbf{A})}{|\tilde{B}|},$$

and we consider the perturbed vector  $\mathbf{b}' := \mathbf{Q}^\top \tilde{\mathbf{x}}'$ . If the perturbed game was non-degenerate with the same support, it would follow that  $(\tilde{\mathbf{x}}', \cdot)$  is the unique equilibrium, which is a contradiction since  $\sum_{i \in \tilde{B}} \tilde{\mathbf{x}}'_i = \sum_{i \in \tilde{B}} \mathbf{x}_i^* + \alpha_D(\mathbf{A}) = 1$ . The norm of the perturbation is again upper bounded in terms of  $\alpha_P(\mathbf{A})$ . Overall, we have shown that  $\delta \leq \alpha_P(\mathbf{A})\text{poly}(n, m)$ . Similar reasoning applies with respect to  $\alpha_D(\mathbf{A})$ . This completes the proof.  $\square$

## NeurIPS Paper Checklist

### 1. Claims

Question: Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope?

Answer: [\[Yes\]](#)

Justification: All claims made in the abstract and introduction are proven in [Appendices C.1](#) to [C.4](#).

Guidelines:

- The answer NA means that the abstract and introduction do not include the claims made in the paper.
- The abstract and/or introduction should clearly state the claims made, including the contributions made in the paper and important assumptions and limitations. A No or NA answer to this question will not be perceived well by the reviewers.
- The claims made should match theoretical and experimental results, and reflect how much the results can be expected to generalize to other settings.
- It is fine to include aspirational goals as motivation as long as it is clear that these goals are not attained by the paper.

### 2. Limitations

Question: Does the paper discuss the limitations of the work performed by the authors?

Answer: [\[Yes\]](#)

Justification: All limitations and assumptions are stated in [Section 1](#).

Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

### 3. Theory Assumptions and Proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [\[Yes\]](#)

Justification: The full set of assumptions and proofs are given in [Appendices C.1 to C.4](#).

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

#### 4. Experimental Result Reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [NA]

Justification: [NA]

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
  - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
  - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
  - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
  - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

#### 5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [NA]

Justification: [NA]

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

## 6. Experimental Setting/Details

Question: Does the paper specify all the training and test details (e.g., data splits, hyper-parameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [NA]

Justification: [NA]

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

## 7. Experiment Statistical Significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [NA]

Justification: [NA]

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer “Yes” if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.

- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

## 8. Experiments Compute Resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [NA]

Justification: [NA]

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

## 9. Code Of Ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines>?

Answer: [Yes]

Justification: The contribution of the paper is theoretical, and conforms in every respect with the NeurIPS Code of Ethics.

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

## 10. Broader Impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [NA]

Justification: The contribution of the paper is theoretical, and we do not foresee any societal impact.

Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.

- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

## 11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [NA]

Justification: [NA]

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

## 12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [NA]

Justification: [NA]

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, [paperswithcode.com/datasets](https://paperswithcode.com/datasets) has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.

- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

### 13. New Assets

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [NA]

Justification: [NA]

Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

### 14. Crowdsourcing and Research with Human Subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [NA]

Justification: [NA]

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.
- According to the NeurIPS Code of Ethics, workers involved in data collection, curation, or other labor should be paid at least the minimum wage in the country of the data collector.

### 15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human Subjects

Question: Does the paper describe potential risks incurred by study participants, whether such risks were disclosed to the subjects, and whether Institutional Review Board (IRB) approvals (or an equivalent approval/review based on the requirements of your country or institution) were obtained?

Answer: [NA]

Justification: [NA]

Guidelines:

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.