

# Optimality and Stability in Non-Convex Smooth Games

Guojun Zhang

GUOJUN.ZHANG@UWATERLOO.CA

Pascal Poupart

PPOUPART@UWATERLOO.CA

Yaoliang Yu

YAOLIANG.YU@UWATERLOO.CA

*School of Computer Science*

*University of Waterloo*

*Vector Institute*

**Editor:** Simon Lacoste-Julien

## Abstract

Convergence to a saddle point for convex-concave functions has been studied for decades, while recent years has seen a surge of interest in *non-convex* (zero-sum) smooth games, motivated by their recent wide applications. It remains an intriguing research challenge how local optimal points are defined and which algorithm can converge to such points. An interesting concept is known as the local minimax point (Jin et al., 2020), which strongly correlates with the widely-known gradient descent ascent algorithm. This paper aims to provide a comprehensive analysis of local minimax points, such as their relation with other solution concepts and their optimality conditions. We find that local saddle points can be regarded as a special type of local minimax points, called *uniformly local minimax points*, under mild continuity assumptions. In (non-convex) quadratic games, we show that local minimax points are (in some sense) equivalent to global minimax points. Finally, we study the stability of gradient algorithms near local minimax points. Although gradient algorithms can converge to local/global minimax points in the non-degenerate case, they would often fail in general cases. This implies the necessity of either novel algorithms or concepts beyond saddle points and minimax points in non-convex smooth games.

**Keywords:** non-convex, minimax points, local optimality, stability, smooth games

## 1. Introduction

The existence of a saddle point in convex-concave minimax optimization follows from the celebrated minimax theorem (e.g. von Neumann, 1928; Sion et al., 1958) and numerical algorithms for finding it have a long history in optimization (e.g. Dem'yanov and Malozemov, 1974; Nemirovsky and Yudin, 1983; Zhang et al., 2019; Lin et al., 2020). Recent success in generative adversarial networks (GANs) (Goodfellow et al., 2014; Heusel et al., 2017), adversarial training (Madry et al., 2018) and reinforcement learning (Sutton et al., 1998) has lead to new challenges (Razaviyayn et al., 2020) for *non-convex non-concave* (NCNC) minimax optimization, a.k.a. NCNC zero-sum games. In such a formulation, we are given a non-convex non-concave bi-variate function  $f(\mathbf{x}, \mathbf{y})$ . One player chooses  $\mathbf{x}$  to minimize  $f(\mathbf{x}, \mathbf{y})$ , and another player chooses  $\mathbf{y}$  to maximize  $f(\mathbf{x}, \mathbf{y})$  (see detailed settings in Section 2). Since non-convex minimax optimization include non-convex minimization as a special case, one cannot hope to find a global optimal solution efficiently. Therefore, we need to look for local

optimal solutions as surrogates. The fundamental gap between the theory for *convex-concave games* and applications using *non-convex non-concave games* raises an important question:

*What is a reasonable definition, in terms of both computational and theoretical convenience, of a local optimal point in non-convex (two-player, zero-sum) games?*

Unlike conventional minimization problems where local optimal solutions are well-defined, for non-convex games a satisfying definition is still under debate. Daskalakis and Panageas (2018) used a local version of saddle points to define local optimality. They studied the local convergence behavior of gradient descent ascent (GDA) (Arrow et al., 1958) and optimistic gradient descent (OGD) (Popov, 1980; Daskalakis et al., 2018). Following this work, an important step was made by Jin et al. (2020), who proposed a new definition of local optimality called local minimax points, compared them with local saddle points, and showed that they are equivalent to the stable solutions of GDA (in some sense). As GDA is widely used in practice, such as for adversarial training (Madry et al., 2018) and for GANs, an enhanced understanding of local minimax points is needed from both theory and application perspectives.

Our work is based on Jin et al. (2020) and we aim to discuss the consequences and implications of their local minimax points to a greater extent. We believe this somewhat pedagogical study can help readers better understand local optimality in non-convex zero-sum games. Specifically, we aim to address the following questions:

- What is the relation between local saddle and local minimax points? Jin et al. (2020) showed that every local saddle point is local minimax, but is there a deeper connection? In Prop. 3.7, we show that local saddle points are a special category of local minimax points called *uniformly* local minimax points, under mild continuity assumptions.
- How can we interpret local minimax points? We give a simplified and unified approach that recovers and extends existing notions of “local mini-maximality,” from the perspective of *infinitesimal robustness* (Hampel, 1974). Local minimax points are understood as the min-player doing infinitesimal robust optimization and the max-player following the strategy of the min-player (Section 3.1).
- One of the benefits of local minimax points is that they are stationary points. Based on the interpretation using infinitesimal robustness, we go one step further and propose a new type of local optimal solutions, called *local robust points* (Def. F.1), which are still stationary points, but strictly include local minimax points as a special case. This new solution concept opens up the possibility to explore solutions in games that are not sequential, in contrast to the sequential Stackelberg games studied in Jin et al. (2020).
- How do we identify local optimal solutions based on derivatives of the function? We analyze natural properties of local minimax points, including first- and second-order optimality conditions. These conditions extend the optimality conditions in Jin et al. (2020) to cases where the domains are constrained and where the Hessian for the max-player is not invertible.
- What is the connection between local and global optimal solutions? We analyze convex-concave games (Theorem 3.10) and non-convex quadratic games (see below), and point out their difference from general non-convex games.

- Is a gradient algorithm stable at a certain local optimal solution? Under suitable conditions, Jin et al. (2020) showed the equivalence between the stable solutions of GDA and local minimax points when the Hessian for the max-player is invertible. We extend this study by analyzing the stability of several other popular gradient algorithms for min-max games and study if they converge to local optimal solutions (see below), even when the Hessian for the max-player is not invertible. Such study provides us with new insights for designing algorithms for minimax points.

As a case study, we thoroughly characterize unconstrained quadratic games, which are potentially non-convex (Daskalakis and Panageas, 2018; Jin et al., 2020; Ibrahim et al., 2020; Wang et al., 2020). On the one hand, quadratic games could help us understand *local convergence* of various gradient algorithms even on NCNC games. On the other hand, w.r.t. the existence and equivalence of global and local versions of minimax points and saddle points, properties for quadratic games are not usually true for general NCNC games. For quadratic games:

- whenever both global (local) minimax and maximin points exist, global (local) saddle points must exist (Corollary 4.6; Example 2.6, Example 4.10);
- global minimax points exist iff local minimax points exist (Theorem 4.4; Example 4.9);
- being stationary and global minimax is equivalent to being local minimax (Theorem 4.4; Example 4.8).

The exact statements formalized as theorems and the corresponding NCNC counterexamples are listed in the parentheses above. Hence, we should be careful when using unconstrained quadratic games as a typical representative in the NCNC setting, especially w.r.t. the optimality properties.

Since our unified definitions of local optimal points are all stationary points, a natural followup question is whether there exist gradient algorithms that can converge to them. In Section 5 we discuss extra-gradient algorithms (Korpelevich, 1976; Popov, 1980; Hsieh et al., 2019). By analyzing the spectrum of the Jacobian, we characterize the stable sets of hyperparameters, which yields insights on how to find local optimal points:

- EG/OGD always locally converge to any non-degenerate local saddle points, and having larger extra-gradient steps increases the local stability;
- for convergence to local minimax points, it is necessary to use two different step sizes and one step size cannot be arbitrarily small;
- for convergence to local robust points, it is more appropriate to use OGD than EG as there are cases where OGD converges, but EG does not.

For one-dimensional quadratic games, we establish the equivalence between local robust points and the stable solution of OGD, extending Jin et al. (2020) for local minimax points.

We delay most proofs to the appendices to keep the main text concise. To help readers navigate the results, we add a title for each definition, theorem, proposition, corollary, remark and example. We also provide a table for easier navigation on the next page.

**Notation:** In this paper we will use several conventions to denote optimality. To distinguish the concepts clearly, we use  $\mathbf{z}_\star = (\mathbf{x}_\star, \mathbf{y}_\star)$  for global/local saddle points;  $\mathbf{z}^\star = (\mathbf{x}^\star, \mathbf{y}^\star)$  for global/local minimax points;  $\mathbf{z}_\star = (\mathbf{x}_\star, \mathbf{y}_\star)$  for global/local maximin points and  $\mathbf{z}^\star = (\mathbf{x}^\star, \mathbf{y}^\star)$

	<b>Statement</b>	<b>Reference</b>
Definitions	global/local saddle point	Definitions 2.1, 3.2
	global/local envelope function	Definitions 2.2, 3.1
	global/local minimax (maximin) point	Definitions 2.3, 3.3
	local robust point (LRP)	Definition F.1
Global results	global saddle = global minimax + global maximin	Theorem 2.5
	both global minimax and maximin points exist, but there is no global saddle point	Example 2.6
	instability of GDA	Example 2.7
	global minimax points exist; no global maximin or global saddle points	Example 2.8
Local minimax	optimality condition when $\partial_{\mathbf{y}\mathbf{y}}^2 f$ is invertible	Theorem 3.4
	equivalence with Jin et al. (2020)	Props. 3.6, 3.9
	local saddle $\approx$ uniformly local minimax	Prop. 3.7, Example 3.8
	stationary and/or global minimax $\neq$ local minimax	Examples 3.11, 4.8
	first-order sufficient condition and examples	Thm 3.14, Example 3.15
	second-order sufficient condition and examples	Thms 3.22, 3.23, Cor 3.24, 3.25, Examples 3.26, 3.27
	necessary conditions and related examples	Thm 3.17, Cor. 3.21 Examples 3.18, 3.19, 3.20
	local minimax exists, no global minimax	Example 4.9
	local minimax & maximin exist, no local saddle	Example 4.10
Convex-concave	local minimax = stationary $\implies$ global minimax	Theorem 3.10
	local minimax = local saddle = LRP	Corollary 3.13
Quadratics	optimality conditions	Thm 4.1, Remark 4.2
	quadratic games can be non-convex	Example 4.3
	stationary + global minimax = local minimax	Theorem 4.4
	bilinear games	Corollary 4.5
	minimax + maximin = saddle	Corollary 4.6
	non-uniformly minimax in quadratic games	Remark 4.7
Stability	equivalence between past-extra gradient and OGD	Lemma 5.1
	stability criteria of EG/OGD	Theorem 5.2
	more aggressive extra-gradient steps, more stable	Theorem 5.3
	EG/OGD are more stable than GDA	Corollary 5.4
	local stability at local saddle points	Lemma 5.5, Theorem 5.6
	local stability at strict local minimax points	Lemma 5.7, Theorem 5.8
	local stability of gradient algorithms at general local minimax points	Proposition 5.9

for local robust points (Appendix F). In Section 5 we also use  $\mathbf{z}^* = (\mathbf{x}^*, \mathbf{y}^*)$  for general stationary points. When two different notions of optimality appear (such as in the proof of Prop. 3.7), we choose the notation based on which notion comes first.

## 2. Global optimal points

We focus on a *two-player zero-sum smooth game* with a payoff function  $f : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$  that is sufficiently many times differentiable depending on the context. We consider  $\mathcal{X} \subset \mathbb{R}^n$  and  $\mathcal{Y} \subset \mathbb{R}^m$  to be non-empty subsets of Euclidean spaces and will add additional assumptions (convexity, closedness) when necessary. The min-player selects a strategy  $\mathbf{x} \in \mathcal{X}$  while the max-player selects a strategy  $\mathbf{y} \in \mathcal{Y}$ , after which the min-player receives utility  $-f(\mathbf{x}, \mathbf{y})$  and the max-player receives  $f(\mathbf{x}, \mathbf{y})$ . In our setting the min-player aims to minimize  $f(\cdot, \mathbf{y})$  given (an estimate of) the max-player’s strategy  $\mathbf{y}$  and conversely the max-player tries to maximize  $f(\mathbf{x}, \cdot)$  given (an estimate of) the min-player’s strategy  $\mathbf{x}$ . In general,  $f$  is not convex in  $\mathbf{x}$  and not concave in  $\mathbf{y}$  (NCNC), which has become extremely popular in machine learning (ML) recently, due to the rise of deep models. For instance, in generative adversarial networks (Goodfellow et al., 2014),  $\mathbf{x}$  models the parameter of a generator while  $\mathbf{y}$  models that of a discriminator. In adversarial training (Madry et al., 2018),  $\mathbf{x}$  is the robust model that we aim to train while  $\mathbf{y}$  represents possible adversarial attacks. In those examples (and many others), the function  $f$  of interest is NCNC. A major challenge is to define proper notions of optimality (stationarity) and to understand the limiting behaviour of popular algorithms that are currently used by practitioners.

In the convex setting, the following solution concept is well-known:

**Definition 2.1 (global saddle)** *We call  $(\mathbf{x}_*, \mathbf{y}_*) \in \mathcal{X} \times \mathcal{Y}$  global saddle if for all  $\mathbf{x} \in \mathcal{X}$  and  $\mathbf{y} \in \mathcal{Y}$ :*

$$f(\mathbf{x}_*, \mathbf{y}) \leq f(\mathbf{x}_*, \mathbf{y}_*) \leq f(\mathbf{x}, \mathbf{y}_*). \quad (2.1)$$

*In other words, we have simultaneously:*

$$\mathbf{x}_* \in \underset{\mathbf{x} \in \mathcal{X}}{\operatorname{argmin}} f(\mathbf{x}, \mathbf{y}_*), \quad \mathbf{y}_* \in \underset{\mathbf{y} \in \mathcal{Y}}{\operatorname{argmax}} f(\mathbf{x}_*, \mathbf{y}). \quad (2.2)$$

Global saddle points correspond to Nash equilibria (Nash, 1950), where each player knows the opponent’s strategy exactly and aims to maximize the gain, but has no incentive to deviate from his/her current strategy.

We may also encounter a scenario where the players move in sequence, and we need the following definitions:

**Definition 2.2 (global envelope function)** *The upper and lower envelope functions are defined respectively as:*

$$\bar{f}(\mathbf{x}) := \sup_{\mathbf{y} \in \mathcal{Y}} f(\mathbf{x}, \mathbf{y}), \quad \underline{f}(\mathbf{y}) := \inf_{\mathbf{x} \in \mathcal{X}} f(\mathbf{x}, \mathbf{y}). \quad (2.3)$$

For envelope functions, we allow  $\bar{f}$  to take value  $+\infty$  and  $\underline{f}$  to take value  $-\infty$ . In Definition 2.2, the min-player for  $\mathbf{x}$  moves first and knows nothing about the max-player for

$\mathbf{y}$ . A natural strategy is to minimize the worst-case payoff, i.e., the upper envelope function  $\bar{f}(\mathbf{x})$ , which is typically non-convex and non-smooth (even when  $f$  is itself smooth):

$$\min_{\mathbf{x} \in \mathcal{X}} \bar{f}(\mathbf{x}). \quad (2.4)$$

On the other hand, the max-player simply maximizes  $f(\mathbf{x}, \cdot)$  given any  $\mathbf{x}$ . This leads immediately to the following solution concept:

**Definition 2.3 (global minimax and maximin)**  $(\mathbf{x}^*, \mathbf{y}^*) \in \mathcal{X} \times \mathcal{Y}$  is *global minimax* if

$$\textcircled{1} \mathbf{x}^* \in \operatorname{argmin}_{\mathbf{x} \in \mathcal{X}} \bar{f}(\mathbf{x}), \quad \textcircled{2} \mathbf{y}^* = \mathbf{y}^*(\mathbf{x}^*) \in \operatorname{argmax}_{\mathbf{y} \in \mathcal{Y}} f(\mathbf{x}^*, \mathbf{y}). \quad (2.5)$$

In other words, for all  $\mathbf{x} \in \mathcal{X}$  and  $\mathbf{y} \in \mathcal{Y}$ :

$$f(\mathbf{x}^*, \mathbf{y}) \leq f(\mathbf{x}^*, \mathbf{y}^*) = \bar{f}(\mathbf{x}^*) \leq \bar{f}(\mathbf{x}). \quad (2.6)$$

Similarly, we call  $(\mathbf{x}_*, \mathbf{y}_*) \in \mathcal{X} \times \mathcal{Y}$  *global maximin* if

$$\textcircled{1} \mathbf{y}_* \in \operatorname{argmax}_{\mathbf{y} \in \mathcal{Y}} \underline{f}(\mathbf{y}), \quad \textcircled{2} \mathbf{x}_* = \mathbf{x}_*(\mathbf{y}_*) \in \operatorname{argmin}_{\mathbf{x} \in \mathcal{X}} f(\mathbf{x}, \mathbf{y}_*). \quad (2.7)$$

In other words, for all  $\mathbf{y} \in \mathcal{Y}$  and  $\mathbf{x} \in \mathcal{X}$ :

$$\underline{f}(\mathbf{y}) \leq \underline{f}(\mathbf{y}_*) = f(\mathbf{x}_*, \mathbf{y}_*) \leq f(\mathbf{x}, \mathbf{y}_*). \quad (2.8)$$

The concept of global minimax points is used widely in machine learning. For example, in the formulation of GAN (Goodfellow et al., 2014), we first find the optimal parameters of the discriminator,  $\boldsymbol{\theta}_D$ , based on the parameters of the generator  $\boldsymbol{\theta}_G$ , and then optimize over  $\boldsymbol{\theta}_G$ . In other words, the optimal solution  $(\boldsymbol{\theta}_G^*, \boldsymbol{\theta}_D^*)$  is a global minimax point (see the definition of  $V$  in Goodfellow et al. (2014)):

$$V(\boldsymbol{\theta}_G^*, \boldsymbol{\theta}_D) \leq V(\boldsymbol{\theta}_G^*, \boldsymbol{\theta}_D^*), \quad \max_{\boldsymbol{\theta}_D} V(\boldsymbol{\theta}_G, \boldsymbol{\theta}_D) \geq \max_{\boldsymbol{\theta}_D} V(\boldsymbol{\theta}_G^*, \boldsymbol{\theta}_D), \quad \forall \boldsymbol{\theta}_G, \boldsymbol{\theta}_D. \quad (2.9)$$

In the distributional robustness formulation (Sinha et al., 2018), we find the global minimax point  $(\boldsymbol{\theta}^*, P^*)$ , where  $\boldsymbol{\theta}^*$  is the best model parameter and  $P^*$  is the worst adversarial distribution, such that:

$$\mathbb{E}_P[\ell(\boldsymbol{\theta}^*; Z)] \leq \mathbb{E}_{P^*}[\ell(\boldsymbol{\theta}^*; Z)], \quad \sup_{P \in \mathcal{P}} \mathbb{E}_P[\ell(\boldsymbol{\theta}; Z)] \geq \sup_{P \in \mathcal{P}} \mathbb{E}_P[\ell(\boldsymbol{\theta}^*; Z)], \quad \forall \boldsymbol{\theta} \in \Theta, P \in \mathcal{P}. \quad (2.10)$$

Since we use neural networks in these applications, the payoff function is non-convex non-concave, and thus a saddle point may not always exist.

**Remark 2.4 (difficulty of finding global minimax)** *Although the notion of global minimax is well-defined, it suffers from some major issues once we enter the NCNC world:*

- *We are not aware of an efficient algorithm (Murty and Kabadi, 1987) for finding a global minimizer  $\mathbf{x}^*$  of the non-convex function  $\bar{f}$ . This can be mitigated by contending with a local minimizer or even stationary point.*

- Given  $\mathbf{x}^*$ , it is NP-hard to find a global maximizer  $\mathbf{y}^*$  of the non-concave function  $f(\mathbf{x}^*, \mathbf{y})$ . While it is tempting to relax again to a local solution, this will unfortunately affect our notion of optimality for  $\mathbf{x}^*$  in the first place. We will return to this issue in the next section.
- The envelope function  $\bar{f}$  is not smooth even when  $f$  is. Although we can turn to non-smooth optimization techniques, it will be inevitably slow to optimize  $\bar{f}$ .

If we define the “mirror” function  $\lambda(\mathbf{y}, \mathbf{x}) = f(\mathbf{x}, \mathbf{y})$ , then  $(\mathbf{x}_*, \mathbf{y}_*)$  is global maximin for  $f$  iff  $(\mathbf{y}_*, \mathbf{x}_*)$  is global minimax for  $-\lambda$ . For this reason, we will limit our discussion mainly to minimax. Definition 2.3 arises in the optimization literature as well since it can be treated as a global solution to the minimax optimization problem:

$$\min_{\mathbf{x} \in \mathcal{X}} \max_{\mathbf{y} \in \mathcal{Y}} f(\mathbf{x}, \mathbf{y}).$$

We note that the ordering of  $\mathbf{x}$  and  $\mathbf{y}$ , i.e. which player moves first, matters: for instance, to get a global minimax pair  $(\mathbf{x}^*, \mathbf{y}^*)$ , we must first find  $\mathbf{x}^*$  and then conditioned on  $\mathbf{x}^*$  we find the “certificate”  $\mathbf{y}^*$ . In game-theoretic terms, this is also known as a Stackelberg game (von Stackelberg, 1934), where  $\mathbf{x}$  is the leader while  $\mathbf{y}$  is the follower.

It is well-known that weak duality, namely the inequality

$$\max_{\mathbf{y} \in \mathcal{Y}} \underline{f}(\mathbf{y}) \leq \min_{\mathbf{x} \in \mathcal{X}} \bar{f}(\mathbf{x}) \tag{2.11}$$

always holds. Strong duality, namely when equality is attained in (2.11), holds only under stringent conditions. The following theorem easily follows from the definitions:

**Theorem 2.5 (e.g. Facchinei and Pang 2007, Theorem 1.4.1)** *For any function  $f$ , the pair  $(\mathbf{x}_*, \mathbf{y}_*) \in \mathcal{X} \times \mathcal{Y}$  is global saddle iff it is both global minimax and global maximin iff strong duality holds and*

$$\mathbf{x}_* \in \operatorname{argmin}_{\mathbf{x} \in \mathcal{X}} \bar{f}(\mathbf{x}), \quad \mathbf{y}_* \in \operatorname{argmax}_{\mathbf{y} \in \mathcal{Y}} \underline{f}(\mathbf{y}). \tag{2.12}$$

Let us give some examples to digest the definitions. In general, it is possible to find a game where both global maximin and minimax points exist, but there is no saddle point:

**Example 2.6 (both global minimax and maximin points exist; no saddle point)**

*Consider the bivariate function*

$$f(x, y) = x^4/4 - x^2/2 + xy \tag{2.13}$$

*defined on  $\mathbb{R} \times \mathbb{R}$ . Global minimax points are clearly  $\{0\} \times \mathbb{R}$  with value 0. On the other hand, global maximin points are  $(\pm 1, 0)$  with value  $-1/4$ . Indeed,*

$$\max_y \min_x x^4/4 - x^2/2 + xy \leq \max_y \min_x x^4/4 - x^2/2 \leq -\frac{1}{4}, \tag{2.14}$$

*with equality attained at  $(\pm 1, 0)$ . The failure of strong duality proves the non-existence of saddle points (Theorem 2.5).*

Note that given a global saddle pair  $(\mathbf{x}_*, \mathbf{y}_*)$ ,  $\mathbf{y}_* \in \mathcal{Y}_* := \operatorname{argmax}_{\mathbf{y} \in \mathcal{Y}} f(\mathbf{x}_*, \mathbf{y})$  but not every certificate  $\bar{\mathbf{y}} \in \mathcal{Y}_*$  forms a global saddle pair with  $\mathbf{x}_*$ . This is known as “instability,” which is the reason underlying the non-convergence of the gradient descent ascent (GDA) algorithm (Golshtein, 1972; Nemirovsky and Yudin, 1983).

**Example 2.7 (instability of GDA)** *Consider the bilinear (hence convex-concave) function*

$$f(x, y) = xy$$

*defined on  $\mathbb{R} \times \mathbb{R}$ . It is easy to verify that global minimax points are precisely the set  $\{0\} \times \mathbb{R}$  while global maximin points are  $\mathbb{R} \times \{0\}$ . Taking the intersection we have the unique global saddle point  $(0, 0)$ . This bilinear function is unstable, since given  $x^* = 0$ , not every global minimax certificate (namely the entire  $\mathbb{R}$ ) forms a global saddle point with  $x^*$ . The last iterates of GDA do not converge to the unique global saddle point for this function with any (constant or not) step size, provided that it is not initialized at the saddle point (Nemirovsky and Yudin, 1983, p. 211).*

Another interesting example consists of quadratic games, which we completely classify in Section 4. Below we give a one-dimensional example where there is no global maximin or saddle point, but global minimax points exist.

**Example 2.8 (global minimax points exist; no global maximin or saddle points)** *Let  $f(x, y) = ax^2 + by^2 + cxy$  with  $a < 0$ ,  $b < 0$  and  $c^2 \geq ab$ . According to the characterization in Theorem 4.1,  $f$  only admits global minimax points. Note that for quadratic games, the existence of both global minimax and maximin points implies the existence of a saddle point, in sharp contrast with Example 2.6.*

From the example above, we see that even for simple quadratic games, saddle points may not exist. In fact, unconstrained quadratic games are often given as typical examples for NCNC minimax optimization (Daskalakis and Panageas, 2018; Jin et al., 2020; Ibrahim et al., 2020; Wang et al., 2020). Locally, they can also be regarded as second-order approximations of a smooth function, and thus seem to be good representatives of NCNC games. However, we will show in Section 4 that they are quite special in many aspects.

### 3. Local optimal points

In this section, we study definitions of local optimal points based on envelope functions. Compared to global optimal points, for local versions, we assume that we only have access to local information of  $f$ , i.e., given a point  $(\mathbf{x}, \mathbf{y})$ , we only know  $f$  over a neighborhood  $\mathcal{N}(\mathbf{x}) \times \mathcal{N}(\mathbf{y})$ . Therefore, each player can only evaluate its current strategy by comparing with other strategies in the current neighborhood, corresponding to the notion of a local minimum (maximum). This can be achieved with the following local envelope functions. In the definition below, we denote

$$\mathcal{N}(\mathbf{y}^*, \epsilon) := \{\mathbf{y} \in \mathcal{Y} : \|\mathbf{y} - \mathbf{y}^*\| \leq \epsilon\}, \tag{3.1}$$

as the intersection of  $\mathcal{Y}$  with a ball of radius  $\epsilon$  surrounding  $\mathbf{y}^*$  in  $\mathbb{R}^m$ , and similarly for  $\mathcal{N}(\mathbf{x}^*, \epsilon)$ . Of course, the exact form of the ball depends on the norm we choose.



**Definition 3.1 (local envelope function)** Fix a reference point  $\mathbf{y}^* \in \mathcal{Y}$  and radius  $\epsilon \geq 0$ , we localize the envelope function:

$$\bar{f}_\epsilon(\mathbf{x}) = \bar{f}_{\epsilon, \mathbf{y}^*}(\mathbf{x}) := \max_{\mathbf{y} \in \mathcal{N}(\mathbf{y}^*, \epsilon)} f(\mathbf{x}, \mathbf{y}). \quad (3.2)$$

The definition for  $\underline{f}_\epsilon(\mathbf{y}) = \underline{f}_{\epsilon, \mathbf{x}^*}(\mathbf{y})$  is similar if we fix some  $\mathbf{x}^* \in \mathcal{X}$ .

In § 3.1 we propose a unified framework for local optimality and then study the differential optimality conditions in § 3.2.

### 3.1 Definitions of local optimality

In this subsection, we start from the simplest definition of local optimality – local saddle points, and then relax the constraints on the players to obtain the more general local minimax points (Jin et al., 2020). It is also possible to extend local minimax points further to local robust points (LRPs), which we delay to Appendix F.

In the NCNC setting, it is natural to consider local versions of saddle points (see Definition 2.1) by localizing around neighborhoods  $\mathcal{N}(\mathbf{x}_*, \epsilon)$  and  $\mathcal{N}(\mathbf{y}_*, \epsilon)$ . Below, when we mention the local envelope functions  $\bar{f}_\epsilon(\mathbf{x})$  and  $\underline{f}_\epsilon(\mathbf{y})$  (see Definition 3.1) the centers and the neighborhoods are often omitted since they are clear from the context.

**Definition 3.2 (local saddle)** We call the pair  $(\mathbf{x}_*, \mathbf{y}_*) \in \mathcal{X} \times \mathcal{Y}$  local saddle if there exists  $\epsilon > 0$ , such that for all  $\mathbf{x} \in \mathcal{N}(\mathbf{x}_*, \epsilon)$  and  $\mathbf{y} \in \mathcal{N}(\mathbf{y}_*, \epsilon)$ ,  $f(\mathbf{x}_*, \mathbf{y}) \leq f(\mathbf{x}_*, \mathbf{y}_*) \leq f(\mathbf{x}, \mathbf{y}_*)$ . In other words,

- Fixing  $\mathbf{x}_*$ , then  $\mathbf{y}_*$  is a local maximizer of  $\underline{f}_{0, \mathbf{x}_*}(\mathbf{y}) = f(\mathbf{x}_*, \mathbf{y})$ ;
- Fixing  $\mathbf{y}_*$ , then  $\mathbf{x}_*$  is a local minimizer of  $\bar{f}_{0, \mathbf{y}_*}(\mathbf{x}) = f(\mathbf{x}, \mathbf{y}_*)$ .

In the above definition, each player contends with the local optimality of its strategy by comparing with other strategies in a neighborhood. For local saddle points, we can WLOG choose the Euclidean norm  $\|\cdot\|_2$  in the neighborhood definition (see (3.1)).

We can now generalize the definition above. One player may not be aware of the exact strategy of the opponent, and thus doing robust optimization, given a certain range of the opponent’s strategy. If  $\mathbf{x}$  is doing (a sequence of) local robust optimization and  $\mathbf{y}$  is doing usual optimization given the strategy of  $\mathbf{x}$ , we have the following definition:

**Definition 3.3 (local minimax)** We call  $(\mathbf{x}^*, \mathbf{y}^*) \in \mathcal{X} \times \mathcal{Y}$  a local minimax point if

- Fixing  $\mathbf{x}^*$ , then  $\mathbf{y}^*$  is a local maximizer of  $\underline{f}_{0, \mathbf{x}^*}(\mathbf{y}) = f(\mathbf{x}^*, \mathbf{y})$ ;
- Fixing  $\mathbf{y}^*$ , then  $\mathbf{x}^*$  is a local minimizer of  $\bar{f}_{\epsilon_n, \mathbf{y}^*}(\mathbf{x})$  for all  $\epsilon_n$  in some sequence  $0 < \epsilon_n \rightarrow 0$ .

Furthermore, if there is a neighborhood  $\mathcal{N}$  of  $\mathbf{x}^*$  such that for all  $\epsilon_n$  in the sequence,  $\mathbf{x}^*$  is a local minimizer of  $\bar{f}_{\epsilon_n}$  on  $\mathcal{N}$ , then we call  $(\mathbf{x}^*, \mathbf{y}^*)$  uniformly local minimax.

In the definition above, we also proposed uniformly local minimax points. By uniformity we mean that the neighborhood  $\mathcal{N}$  does not depend on the element  $\epsilon_n$  in the sequence. We will show a close relation between local saddle points and uniformly local minimax points in Proposition 3.7.

Definition 3.3 reveals the asymmetric position between the two players for  $\mathbf{x}$  and  $\mathbf{y}$ :  $\mathbf{y}$  needs only be a local certificate to testify the local optimality of  $\mathbf{x}$ , but  $\mathbf{x}$  minimizes the envelope function  $\bar{f}_\epsilon(\mathbf{x})$ , the worst-case payoff, simultaneously for a sequence of  $\epsilon_n \rightarrow 0$ . By switching the role of  $\mathbf{x}$  and  $\mathbf{y}$  we obtain a similar notion of local maximin. When both players satisfy this stringent condition, we obtain a new optimality notion that we term as local robust points (Appendix F).

In Proposition 3.6 we will see that Definition 3.3 has a seemingly stronger but equivalent form. To digest the somewhat complicated definition, we mention the following interpretation (e.g. Wang et al., 2020):

**Theorem 3.4 (sufficient and necessary condition of local minimax when  $\partial_{\mathbf{y}\mathbf{y}}^2 f$  is invertible)** *Let  $\mathcal{X} = \mathbb{R}^n, \mathcal{Y} = \mathbb{R}^m$  and  $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$  be twice continuously differentiable. Suppose  $\partial_{\mathbf{y}\mathbf{y}}^2 f(\mathbf{x}^*, \mathbf{y}^*)$  is invertible (i.e. non-degenerate), then  $(\mathbf{x}^*, \mathbf{y}^*)$  is local minimax iff*

- $\partial_{\mathbf{y}} f(\mathbf{x}^*, \mathbf{y}^*) = \mathbf{0}$ ,  $\partial_{\mathbf{y}\mathbf{y}}^2 f(\mathbf{x}^*, \mathbf{y}^*) \prec \mathbf{0}$ , and
- $\mathbf{x}^*$  is a local minimizer of the total function  $f(\mathbf{x}, \mathbf{y}(\mathbf{x}))$  where  $\mathbf{y}$  is defined implicitly near  $\mathbf{x}^*$  through the non-linear equation

$$\partial_{\mathbf{y}} f(\mathbf{x}, \mathbf{y}) = \mathbf{0}. \tag{3.3}$$

We emphasize that, unlike the definition in Jin et al. (2020), we do not allow  $\epsilon_n$  to take 0 in Definition 3.3 for two reasons: (a) This allows us to better separate local saddle from local minimax; (b) It is unnecessary to have  $\epsilon_n = 0$ , as we will see in Proposition 3.9.

We now show how to simplify Definition 3.3, starting with the following key lemma:

**Lemma 3.5** *Suppose  $\mathbf{y}^*$  maximizes  $f(\mathbf{x}^*, \mathbf{y})$  over some neighborhood  $\mathcal{N}(\mathbf{y}^*, \epsilon_0)$ . If  $\mathbf{x}^*$  is a local minimizer of  $\bar{f}_{\epsilon, \mathbf{y}^*}$ , for some  $0 \leq \epsilon \leq \epsilon_0$ , then it remains a local minimizer (even over the same local neighborhood) of  $\bar{f}_{\mathcal{N}}(\mathbf{x}) := \max_{\mathbf{y} \in \mathcal{N}} f(\mathbf{x}, \mathbf{y})$  for any  $\mathcal{N}(\mathbf{y}^*, \epsilon) \subseteq \mathcal{N} \subseteq \mathcal{N}(\mathbf{y}^*, \epsilon_0)$ .*

Note that in the lemma above we allow  $\epsilon = 0$ . Lemma 3.5 reveals a key property of the local minimax point in Definition 3.3: the norm in the neighborhood definition (see (3.1)) is immaterial (since we can shrink the neighborhood using Lemma 3.5 without impairing local minimaximality). In other words, the definition of local minimax points is topological and it does not depend on the norm we actually choose. Using Lemma 3.5 we can “strengthen” the notion of local minimax even more. In particular, if Definition 3.3 holds for one diminishing sequence such that  $\epsilon_0 \geq \epsilon_n \downarrow 0$  then it automatically holds for *all* sequences that satisfy this same condition. We can even extend the sequence to an interval of  $\epsilon$ ’s:

**Proposition 3.6 (equivalent definition of local minimax)** *The pair  $(\mathbf{x}^*, \mathbf{y}^*) \in \mathcal{X} \times \mathcal{Y}$  is a local minimax point iff*

- Fixing  $\mathbf{x}^*$ , then  $\mathbf{y}^*$  is a local maximizer of  $\underline{f}_{0, \mathbf{x}^*}(\mathbf{y}) = f(\mathbf{x}^*, \mathbf{y})$ ;

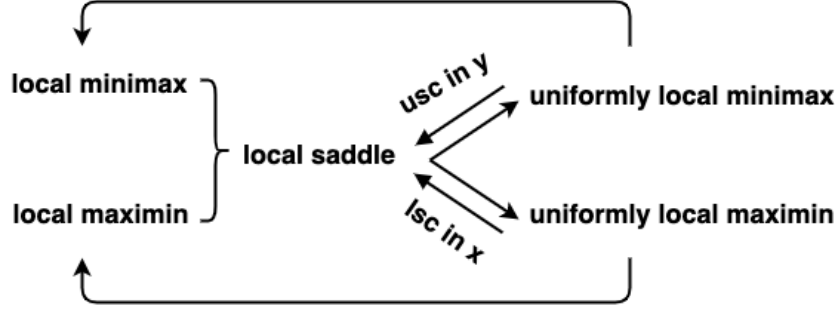


Figure 1 The relationship among different notions of local optimality. usc: upper semi-continuity and lsc: lower semi-continuity. The arrow and the bracket signs mean “to imply.” For example, a uniformly local minimax point is *bona fide* local minimax, and if a point is both local minimax and local maximin, it is local saddle.

- Fixing  $\mathbf{y}^*$ , then  $\mathbf{x}^*$  is a local minimizer of  $\bar{f}_{\epsilon, \mathbf{y}^*}(\mathbf{x})$  for all  $\epsilon \in (0, \epsilon_0]$  with some  $\epsilon_0 > 0$ .

From Definition 3.3, every uniformly local minimax point is local minimax. In fact, much more can be said between uniformly local minimax and local saddle:

**Proposition 3.7 (local saddle and uniformly local minimax)** *Every local saddle point is uniformly local minimax. If for any  $\mathbf{x} \in \mathcal{X}$ ,  $f(\mathbf{x}, \cdot)$  is upper semi-continuous, then every uniformly local minimax point is local saddle.*

Thus, for upper semi-continuous functions (in  $\mathbf{y}$ ), surprisingly, local saddle points coincide with uniformly local minimax points. We cannot drop the semi-continuity assumption:

**Example 3.8 (uniformly local minimax does not imply local saddle without semi-continuity)** *Fix any  $\mathbf{y}^* \in \mathcal{Y}$  and consider the lower semi-continuous function*

$$f(x, y) = \begin{cases} -x^2, & y = y^* \\ x^2, & y \neq y^* \end{cases}, \quad \text{with } \bar{f}_{\epsilon, \mathbf{y}^*}(x) = \begin{cases} -x^2, & \epsilon = 0 \\ x^2, & \epsilon \neq 0 \end{cases}. \quad (3.4)$$

$(0, y^*)$  is uniformly local minimax but not local saddle.

Figure 1 shows the relation between local saddle and (uniformly) local minimax (maximin) points. Finally, we prove that our Definition 3.3 coincides with the seemingly different one in Definition 14 of Jin et al. (2020). Effectively, we manage to remove the continuity assumption in Lemma 16 of Jin et al. (2020) (cf. Proposition 3.6).

**Proposition 3.9 (equivalence with Jin et al. (2020))** *The pair  $(\mathbf{x}^*, \mathbf{y}^*)$  is local minimax w.r.t. function  $f$  iff there exists  $\delta_0 > 0$  and a non-negative function  $h$  satisfying  $h(\delta) \rightarrow 0$  as  $\delta \rightarrow 0$ , such that for any  $\delta \in (0, \delta_0]$  and any  $(\mathbf{x}, \mathbf{y}) \in \mathcal{N}(\mathbf{x}^*, \delta) \times \mathcal{N}(\mathbf{y}^*, \delta)$  we have*

$$f(\mathbf{x}^*, \mathbf{y}) \leq f(\mathbf{x}^*, \mathbf{y}^*) \leq \left[ \max_{\mathbf{y}' \in \mathcal{N}(\mathbf{y}^*, h(\delta))} f(\mathbf{x}, \mathbf{y}') \right] =: \bar{f}_{h(\delta)}(\mathbf{x}). \quad (3.5)$$

From this equivalence, we can also derive that every local saddle point is local minimax (Jin et al., 2020, Proposition 17). However, our Proposition 3.7 gives a more detailed depiction of local saddle points. For functions that are convex in  $\mathbf{x}$  and concave in  $\mathbf{y}$ , we naturally expect that local optimality is somehow equivalent to global optimality:

**Theorem 3.10 (local and global minimax points in the convex-concave case)**

*Let the function  $f(\mathbf{x}, \mathbf{y})$  be convex in  $\mathbf{x}$  and concave in  $\mathbf{y}$ . Then, an interior point  $(\mathbf{x}, \mathbf{y})$  is local minimax iff it is stationary, i.e.,  $\partial_{\mathbf{x}}f(\mathbf{x}, \mathbf{y}) = \mathbf{0}$  and  $\partial_{\mathbf{y}}f(\mathbf{x}, \mathbf{y}) = \mathbf{0}$  iff it is saddle. In particular, local minimax implies global minimax.*

However, non-stationary global minimax points cannot be local minimax, see Example 2.7 and Theorem 3.12 (below). Even with stationarity, the convex-concave assumption in Theorem 3.10 cannot be appreciably weakened, as illustrated in the following example:

**Example 3.11 (stationary global minimax points are not local minimax in the non-convex case)** *Let  $f(x, y) = x^3y$  be non-convex in  $x$  but linear in  $y$ . The point  $(x^*, y^*) = (0, 1)$  is clearly stationary and global minimax. We verify that*

$$\bar{f}_\epsilon(x) = \begin{cases} (1 + \epsilon)x^3, & x \geq 0 \\ (1 - \epsilon)x^3, & x \leq 0 \end{cases}, \quad (3.6)$$

*hence  $x^* = 0$  is not a local minimizer of  $\bar{f}_\epsilon$  (for any  $\epsilon < 1$ ) and  $(0, 1)$  is not local minimax. This counterexample is constructed by performing the  $\mathcal{C}^1$  homeomorphic transformation  $(x, y) \mapsto (x^3, y)$  of the bilinear game  $b(x, y) = xy$ . We can verify that (separate) homeomorphisms transform local/global minimax points accordingly. However,  $\mathcal{C}^1$  homeomorphisms can turn non-stationary points into stationary (which is not possible in presence of convexity since stationarity equates minimality which is preserved under homeomorphisms).*

Nevertheless, for quadratic games, we can remove the convexity-concavity assumption, as will be shown in Theorem 4.1 below.

### 3.2 Optimality conditions

Optimality conditions are an indispensable part of optimization (Bertsekas, 1997) since they help us identify local optimal points and design new algorithms. In this section, we provide first- and second-order necessary and sufficient conditions for local minimax (maximin) points. Our results extend existing ones in Jin et al. (2020). We assume  $\mathcal{X}$  and  $\mathcal{Y}$  are closed<sup>1</sup> and thus  $\mathcal{N}(\mathbf{y}^*, \epsilon)$  and  $\mathcal{N}(\mathbf{x}^*, \epsilon)$  are compact. We build on some classical results in non-smooth analysis, for which we provide a self-contained review in Appendix A, including the definition of the directional derivative  $D\bar{f}_\epsilon(\mathbf{x}; \mathbf{t})$  of an envelope function  $\bar{f}_\epsilon$  at  $\mathbf{x}$  along direction  $\mathbf{t}$ :

$$D\bar{f}_\epsilon(\mathbf{x}; \mathbf{t}) = \lim_{\alpha \rightarrow 0^+} \frac{\bar{f}_\epsilon(\mathbf{x} + \alpha\mathbf{t}) - \bar{f}_\epsilon(\mathbf{x})}{\alpha}. \quad (3.7)$$

Specifically, if  $f$  and  $\partial_{\mathbf{x}}f$  are jointly continuous (continuous w.r.t.  $(\mathbf{x}, \mathbf{y})$ ), then the directional derivative  $D\bar{f}_\epsilon(\mathbf{x}; \mathbf{t})$  always exist (Theorem A.9). In the following subsections,  $f \in \mathcal{C}^p$  means that  $f$  is  $p^{\text{th}}$  continuously differentiable.

---

1. Of course they are contained in bigger open sets where derivatives of  $f$  are well defined.

## 3.2.1 FIRST-ORDER NECESSARY CONDITIONS

**Theorem 3.12 (first-order necessary, local minimax)** *Let  $f \in \mathcal{C}^1$ . At a local minimax point  $(\mathbf{x}^*, \mathbf{y}^*)$ , we have:*

$$\partial_{\mathbf{x}} f(\mathbf{x}^*, \mathbf{y}^*)^\top \bar{\mathbf{t}} \geq 0 \geq \partial_{\mathbf{y}} f(\mathbf{x}^*, \mathbf{y}^*)^\top \mathbf{t}, \quad (3.8)$$

for any directions  $\bar{\mathbf{t}} \in \mathbf{K}_d(\mathcal{X}, \mathbf{x}^*)$ ,  $\mathbf{t} \in \mathbf{K}_d(\mathcal{Y}, \mathbf{y}^*)$ , where the cone

$$\mathbf{K}_d(\mathcal{X}, \mathbf{x}) := \liminf_{\alpha \rightarrow 0^+} \frac{\mathcal{X} - \mathbf{x}}{\alpha} := \{\mathbf{t} : \forall \{\alpha_k\} \rightarrow 0^+ \exists \{\alpha_{k_i}\} \rightarrow 0^+, \{\mathbf{t}_{k_i}\} \rightarrow \mathbf{t},$$

such that  $\mathbf{x} + \alpha_{k_i} \mathbf{t}_{k_i} \in \mathcal{X}\}$

and  $\mathbf{K}_d(\mathcal{Y}, \mathbf{y})$  is defined similarly.

**Proof** This result follows from its more general version for local robust points, Theorem F.6. ■

In the theorem above,  $\mathbf{K}_d(\mathcal{X}, \mathbf{x})$  is known as the derivable cone (Rockafellar and Wets, 2009, p. 198), which may strictly include the feasible tangent cone. When the set  $\mathcal{X}$  is closed and convex, the two coincide (Hiriart-Urruty and Lemaréchal, 2004, p. 65):

$$\mathbf{K}_d(\mathcal{X}, \mathbf{x}) = \overline{\text{cone}(\mathcal{X} - \mathbf{x})} := \text{cl}(\mathbf{t} \in \mathbb{R}^n : \mathbf{t} = \alpha(\mathbf{y} - \mathbf{x}), \mathbf{y} \in \mathcal{X}, \alpha \geq 0), \quad (3.9)$$

with  $\text{cl}$  denoting the closure of a set. We can derive a similar reduction when  $\mathcal{Y}$  is closed and convex. If both  $\mathcal{X}$  and  $\mathcal{Y}$  are closed and convex, then (3.8) reduces to:

$$\partial_{\mathbf{x}} f(\mathbf{x}^*, \mathbf{y}^*)^\top (\mathbf{x} - \mathbf{x}^*) \geq 0 \geq \partial_{\mathbf{y}} f(\mathbf{x}^*, \mathbf{y}^*)^\top (\mathbf{y} - \mathbf{y}^*), \quad \text{for any } \mathbf{x} \in \mathcal{X}, \mathbf{y} \in \mathcal{Y}. \quad (3.10)$$

This can be regarded as a bi-variate version of first-order (necessary) optimality condition for a local minimum (Bertsekas, 1997, Prop. 2.1.2). Solutions that satisfy such condition are often called stationary points. It extends the result in Jin et al. (2020) to the constrained case. Specifically, if  $(\mathbf{x}^*, \mathbf{y}^*)$  is in the interior of  $\mathcal{X} \times \mathcal{Y}$ , which always holds when  $\mathcal{X} = \mathbb{R}^n$  and  $\mathcal{Y} = \mathbb{R}^m$ , then Theorem 3.12 simplifies to

$$\partial_{\mathbf{x}} f(\mathbf{x}^*, \mathbf{y}^*) = \mathbf{0}, \quad \partial_{\mathbf{y}} f(\mathbf{x}^*, \mathbf{y}^*) = \mathbf{0}, \quad (3.11)$$

agreeing with Jin et al. (2020). Moreover, Theorem F.6 in Appendix F shows that there is an even broader class of local optimal points named local robust points (LRPs) that has the same necessary conditions, (3.8), (3.10) and (3.11), as local saddle points (e.g. Barazandeh and Razaviyayn, 2020, Definition 2) and local minimax points. It also implies that in the convex-concave case, all local notions of optimality agree:

**Corollary 3.13 (local optimal solutions in the convex-concave case)** *Let  $\mathcal{X}$  and  $\mathcal{Y}$  be convex and the function  $f(\mathbf{x}, \mathbf{y})$  be convex in  $\mathbf{x}$  and concave in  $\mathbf{y}$ . A point is local (global) saddle iff it is local minimax (maximin) iff it is an LRP.*

This corollary does not hold in the non-convex setting, see Examples 4.3 and F.3.

## 3.2.2 FIRST-ORDER SUFFICIENT CONDITIONS

Let us define the *active set* of the *zeroth* order (by “zeroth” we mean that only the function values are involved):

$$\mathcal{Y}_0(\mathbf{x}^*; \epsilon) = \{\mathbf{y} \in \mathcal{N}(\mathbf{y}^*, \epsilon) : \bar{f}_\epsilon(\mathbf{x}^*) = f(\mathbf{x}^*, \mathbf{y})\}. \quad (3.12)$$

We derive the first-order sufficient conditions for local minimax points (which follow from the sufficient condition in Theorem A.5 and Danskin’s theorem in Theorem A.9):

**Theorem 3.14 (first-order sufficient condition, local minimax)** *Assume  $\partial_{\mathbf{x}}f(\mathbf{x}, \mathbf{y})$  is continuous. If  $f(\mathbf{x}^*, \cdot)$  is maximized at  $\mathbf{y}^*$  over a neighborhood around  $\mathbf{y}^*$ , and there exists  $\epsilon_0 > 0$  such that for any  $\epsilon \in (0, \epsilon_0)$ ,*

$$\mathbf{0} \neq \mathbf{t} \in \mathcal{K}_c(\mathcal{X}, \mathbf{x}^*) \implies D\bar{f}_\epsilon(\mathbf{x}^*; \mathbf{t}) = \max_{\mathbf{y} \in \mathcal{Y}_0(\mathbf{x}^*; \epsilon)} \partial_{\mathbf{x}}f(\mathbf{x}^*, \mathbf{y})^\top \mathbf{t} > 0, \quad (3.13)$$

where the contingent cone is defined as:

$$\mathcal{K}_c(\mathcal{X}, \mathbf{x}) := \limsup_{\alpha \rightarrow 0^+} \frac{\mathcal{X} - \mathbf{x}}{\alpha} := \{\mathbf{t} : \exists \{\alpha_k\} \rightarrow 0^+, \{\mathbf{t}_k\} \rightarrow \mathbf{t}, \text{ such that } \mathbf{x} + \alpha_k \mathbf{t}_k \in \mathcal{X}\},$$

then  $(\mathbf{x}^*, \mathbf{y}^*)$  is a local minimax point.

In the case when  $\mathcal{X}$  is a convex set.  $\mathcal{K}_c(\mathcal{X}, \mathbf{x})$  reduces to the usual cone of feasible directions:

$$\mathcal{K}_c(\mathcal{X}, \mathbf{x}) = \overline{\text{cone}(\mathcal{X} - \mathbf{x})} := \text{cl}(\mathbf{t} \in \mathbb{R}^n : \mathbf{t} = \alpha(\mathbf{y} - \mathbf{x}), \mathbf{y} \in \mathcal{X}, \alpha \geq 0). \quad (3.14)$$

If furthermore  $\text{cone}(\mathcal{X} - \mathbf{x})$  is closed, (3.13) becomes:

$$\max_{\mathbf{y} \in \mathcal{Y}_0(\mathbf{x}^*; \epsilon)} \partial_{\mathbf{x}}f(\mathbf{x}^*, \mathbf{y})^\top (\mathbf{x} - \mathbf{x}^*) > 0, \forall \mathbf{x}^* \neq \mathbf{x} \in \mathcal{X}. \quad (3.15)$$

Let us demonstrate the first order condition with the following example:

**Example 3.15 (application of the first-order sufficient condition of local minimax points)** *Suppose  $f(x, y) = xy$  is bilinear. At  $(x^*, y^*) = (0, 0)$ , we have:*

$$\bar{f}_\epsilon(x^*) = f(x^*, y) = 0, \forall y \in \mathbb{R}. \quad (3.16)$$

Therefore, according to (3.12),  $\mathcal{Y}_0(\mathbf{x}^*; \epsilon) = \mathcal{N}(y^*, \epsilon)$ . Also,  $\partial_x f(x^*, y) = y$  and

$$D\bar{f}_\epsilon(x^*; x - x^*) = \max_{\mathcal{N}(y^*, \epsilon)} y(x - x^*) = \epsilon|x| > 0, \forall x \neq x^*. \quad (3.17)$$

According to Theorem 3.14,  $(x^*, y^*)$  is a local minimax point.

## 3.2.3 SECOND-ORDER NECESSARY CONDITIONS

We now turn to the second-order necessary condition of local minimax points. We sometimes use  $\partial_{\mathbf{xx}}^2 f$  as a shorthand for the second-order derivative  $\partial_{\mathbf{xx}}^2 f(\mathbf{x}^*, \mathbf{y}^*)$ , and similarly for other second-order partial derivatives. For a local minimax point  $(\mathbf{x}^*, \mathbf{y}^*)$ ,  $\mathbf{y}^*$  maximizes  $f(\mathbf{x}^*, \cdot)$  locally, and thus we have the property that  $\bar{f}_\epsilon(\mathbf{x}^*) = f(\mathbf{x}^*, \mathbf{y}^*)$  for any small  $\epsilon$ , from which we can make significant simplifications. The following technical lemma, when combined with the necessity condition in Theorem A.3, allows us to classify the directions:

**Lemma 3.16 (directional derivatives for different  $\bar{f}_\epsilon$ )** *Suppose  $f$  and  $\partial_{\mathbf{x}} f$  are jointly continuous and thus the directional derivative (3.7) exists. If  $\mathbf{y}^*$  is a local maximizer of  $f(\mathbf{x}^*, \cdot)$  over a neighborhood  $\mathcal{N}(\mathbf{y}^*, \epsilon_0)$ , then for any  $0 \leq \epsilon_1 \leq \epsilon_2 \leq \epsilon_0$ ,  $\mathcal{Y}_0(\mathbf{x}^*; \epsilon_1) \subseteq \mathcal{Y}_0(\mathbf{x}^*; \epsilon_2)$  and for each  $\mathbf{t} \in \mathcal{K}_d(\mathcal{X}, \mathbf{x}^*)$ ,  $D\bar{f}_{\epsilon_2}(\mathbf{x}^*; \mathbf{t}) \geq D\bar{f}_{\epsilon_1}(\mathbf{x}^*; \mathbf{t})$ .*

Indeed, for a local minimax point  $(\mathbf{x}^*, \mathbf{y}^*)$  and any direction  $\mathbf{t} \in \mathcal{K}_d(\mathcal{X}, \mathbf{x}^*)$ , we know from the necessity condition in Theorem A.3 that  $D\bar{f}_\epsilon(\mathbf{x}^*; \mathbf{t}) \geq 0$  for all small  $\epsilon$ , which, combined with Lemma 3.16 above, leaves us with two possibilities:

1.  $D\bar{f}_\epsilon(\mathbf{x}^*; \mathbf{t}) > 0$  for all  $\epsilon > 0$  smaller than some  $\epsilon_0(\mathbf{t})$ ;
2.  $D\bar{f}_\epsilon(\mathbf{x}^*; \mathbf{t}) = 0$  for all  $\epsilon > 0$  smaller than some  $\epsilon_0(\mathbf{t})$ .

We call the direction  $\mathbf{t}$  a *critical direction* in the second case above. With this distinction among directions, we derive the second-order necessary condition for local minimax points:

**Theorem 3.17 (second-order necessary condition, local minimax)** *Suppose  $f, \partial_{\mathbf{x}} f$  and  $\partial_{\mathbf{xx}}^2 f$  are all (jointly) continuous. If  $(\mathbf{x}^*, \mathbf{y}^*)$  is a local minimax point, then for each direction  $\mathbf{t} \in \mathcal{K}_d(\mathcal{X}, \mathbf{x}^*)$ , one of the following holds:*

1.  $D\bar{f}_\epsilon(\mathbf{x}^*; \mathbf{t}) > 0$  for all  $\epsilon > 0$  smaller than some  $\epsilon_0(t)$ ;
2.  $D\bar{f}_\epsilon(\mathbf{x}^*; \mathbf{t}) = 0$  for all  $\epsilon > 0$  smaller than some  $\epsilon_0(\mathbf{t})$  (i.e.  $\mathbf{t}$  is critical), in which case we further have

$$\mathbf{t}^\top \partial_{\mathbf{xx}}^2 f(\mathbf{x}^*, \mathbf{y}^*) \mathbf{t} + \frac{1}{2} \limsup_{\mathbf{z} \rightarrow \mathbf{y}^*} \left[ \max\{\partial_{\mathbf{x}} f(\mathbf{x}^*, \mathbf{z})^\top \mathbf{t}, 0\}^2 (f(\mathbf{x}^*, \mathbf{y}^*) - f(\mathbf{x}^*, \mathbf{z}))^\dagger \right] \geq 0, \quad (3.18)$$

where  $t^\dagger = 1/t$  if  $t \neq 0$  and 0 otherwise.

The important point to take from Theorem 3.17 is that we should test the second-order condition (3.18) only for critical directions, and the second-order derivatives of  $f$  may not fully capture the second-order derivatives of the envelope function  $\bar{f}_\epsilon$ , which can be clearly demonstrated from the following examples:

**Example 3.18 (the importance of critical directions)** *Let*

$$f(x, y) = -x^2 + xy^3$$

*be defined over  $\mathcal{X} = \mathcal{Y} = \mathbb{R}$  and consider the local minimax point  $(x^*, y^*) = (0, 0)$ . Indeed, for any  $\epsilon > 0$ ,  $x^*$  is a local minimizer of  $\bar{f}_\epsilon(x) = |x|\epsilon^3 - x^2$ . However,  $\partial_{xx}^2 f = -2$  while  $f(x^*, y^*) = f(x^*, z) = 0$  for any  $z$ . Thus, the second-order condition (3.18) fails at the directions  $t = \pm 1$ . However, there is no contradiction since these directions are not critical. Indeed, using Theorem A.9 we can verify that  $D\bar{f}_\epsilon(x^*; \pm 1) = \epsilon^3 > 0$ .*

**Example 3.19 (the importance of critical directions under multiple dimensions)**

Let

$$f(\mathbf{x}, \mathbf{y}) = -x_2^2 + x_2 y_2^3 - (y_1 + y_2)^2 + 2x_1(y_1 + y_2)$$

be defined over  $\mathcal{X} = \mathcal{Y} = \mathbb{R}^2$  and consider the local minimax point  $(\mathbf{x}^*, \mathbf{y}^*) = (\mathbf{0}, \mathbf{0})$ : Indeed,  $f(\mathbf{x}^*, \cdot)$  is clearly maximized locally at  $\mathbf{y}^* = \mathbf{0}$  and upon choosing  $y_1 = x_1 - \text{sgn}(x_2)\epsilon/2$ ,  $y_2 = \text{sgn}(x_2)\epsilon/2$  and considering  $|x_1| < \epsilon/2$  and  $|x_2| < (\epsilon/2)^3$ , we have

$$\|\mathbf{y} - \mathbf{x}\|_\infty \leq \epsilon/2 + (\epsilon/2)^3, \bar{f}_\epsilon(\mathbf{x}) \geq f(\mathbf{x}, \mathbf{y}) = x_1^2 + |x_2|(\epsilon/2)^3 - x_2^2 \geq 0 = \bar{f}_\epsilon(\mathbf{x}^*), \quad (3.19)$$

where we choose WLOG the  $\ell_\infty$  norm in our neighborhood definition (3.1). The second-order derivatives are:

$$\partial_{\mathbf{y}\mathbf{x}}^2 f = \begin{bmatrix} 2 & 0 \\ 2 & 0 \end{bmatrix}, \partial_{\mathbf{y}\mathbf{y}}^2 f = \begin{bmatrix} -2 & -2 \\ -2 & -2 \end{bmatrix}, \partial_{\mathbf{x}\mathbf{x}}^2 f = \begin{bmatrix} 0 & 0 \\ 0 & -2 \end{bmatrix}. \quad (3.20)$$

We have  $\mathcal{Y}_0(\mathbf{x}^*; \epsilon) = \{\mathbf{y} \in \mathcal{N}_\infty(\mathbf{x}^*, \epsilon) : y_1 + y_2 = 0\}$  and for any direction  $\mathbf{t}$ ,

$$\text{D}\bar{f}_\epsilon(\mathbf{x}^*; \mathbf{t}) = \max_{\mathbf{y} \in \mathcal{Y}_0(\mathbf{x}^*; \epsilon)} \mathbf{t}^\top \partial_{\mathbf{x}} f(\mathbf{x}^*, \mathbf{y}) = \epsilon^3 |t_2| \geq 0. \quad (3.21)$$

It follows that the critical directions satisfy  $t_2 = 0$ . Take a non-critical direction  $\mathbf{t} = (1, 3)$ , we easily verify that  $(\partial_{\mathbf{y}\mathbf{x}}^2 f)\mathbf{t} = (2, 2)$  lies in the range space of  $\partial_{\mathbf{y}\mathbf{y}}^2 f$ . However,

$$\begin{aligned} & \limsup_{\mathbf{z} \rightarrow \mathbf{y}^*} \left[ \max\{\partial_{\mathbf{x}} f(\mathbf{x}^*, \mathbf{z})^\top \mathbf{t}, 0\}^2 (f(\mathbf{x}^*, \mathbf{y}^*) - f(\mathbf{x}^*, \mathbf{z}))^\dagger \right] \\ &= \limsup_{\mathbf{z} \rightarrow \mathbf{0}, z_1 + z_2 \neq 0} \frac{[2(z_1 + z_2) + 3z_2^3]_+^2}{(z_1 + z_2)^2} = 4, \end{aligned} \quad (3.22)$$

so that the second-order condition in (3.18), which in this case coincides with

$$\mathbf{t}^\top (\partial_{\mathbf{x}\mathbf{x}}^2 f - \partial_{\mathbf{x}\mathbf{y}}^2 f (\partial_{\mathbf{y}\mathbf{y}}^2 f)^\dagger \partial_{\mathbf{y}\mathbf{x}}^2 f) \mathbf{t},$$

does not hold ( $-18 + 2 = -16 \not\geq 0$ ). Nevertheless, along a critical direction  $\mathbf{t}$  (where  $t_2 = 0$ ):

$$\mathbf{t}^\top \partial_{\mathbf{x}\mathbf{x}}^2 f(\mathbf{x}^*, \mathbf{y}^*) \mathbf{t} = 0, f(\mathbf{x}^*, \mathbf{z}) = -(z_1 + z_2)^2, \partial_{\mathbf{x}} f(\mathbf{x}^*, \mathbf{z})^\top \mathbf{t} = 2t_1(z_1 + z_2), \quad (3.23)$$

and thus the left-hand side of (3.18) simplifies to  $2t_1^2 \geq 0$ . In other words, the second-order condition indeed holds for critical directions.

**Example 3.20 (high order derivatives might be involved in Theorem 3.17)** The second term in (3.18) may involve higher-order information of  $f$ , rather than the standard second-order optimality condition for e.g. the minimizer of a smooth function. The higher-order term comes from the difference of function values. Let  $f(x, y) = -x^2 - y^4 + 4xy^2$  and consider the local minimax point  $(x^*, y^*) = (0, 0)$ . We have  $\mathcal{Y}_0(x^*; \epsilon) = \{y^*\}$  hence every direction is critical. In the direction  $t = 1$ , the l.h.s. of (3.18) becomes

$$-2 + \max\{4z^2 t, 0\}^2 / (2z^4) = 6 > 0.$$

Under the condition that  $\partial_{\mathbf{y}\mathbf{y}}^2 f$  is invertible, we recover the following result from Jin et al. (2020):



**Corollary 3.21 (second-order necessary condition, invertible)** *Let  $f \in \mathcal{C}^2$ . At a local minimax point  $(\mathbf{x}^*, \mathbf{y}^*)$  in the interior of  $\mathcal{X} \times \mathcal{Y}$ , if  $\partial_{\mathbf{y}\mathbf{y}}^2 f$  is invertible, then*

$$\partial_{\mathbf{y}\mathbf{y}}^2 f \prec \mathbf{0} \text{ and } \partial_{\mathbf{x}\mathbf{x}}^2 f - \partial_{\mathbf{x}\mathbf{y}}^2 f (\partial_{\mathbf{y}\mathbf{y}}^2 f)^{-1} \partial_{\mathbf{y}\mathbf{x}}^2 f \succeq \mathbf{0}. \quad (3.24)$$

**Proof** It is easy to prove  $\partial_{\mathbf{y}\mathbf{y}}^2 f \preceq \mathbf{0}$  and since  $\partial_{\mathbf{y}\mathbf{y}}^2 f$  is invertible, we have  $\partial_{\mathbf{y}\mathbf{y}}^2 f \prec \mathbf{0}$ . By expanding  $f(\mathbf{x}^*, \mathbf{z})$  to the second order, the second term in (3.18) becomes:

$$\limsup_{\mathbf{z} \rightarrow \mathbf{y}^*} \frac{\max\{(\mathbf{z} - \mathbf{y}^*)^\top (\partial_{\mathbf{y}\mathbf{x}}^2 f) \mathbf{t}, 0\}^2}{(\mathbf{z} - \mathbf{y}^*)^\top (-\partial_{\mathbf{y}\mathbf{y}}^2 f) (\mathbf{z} - \mathbf{y}^*)}. \quad (3.25)$$

With a change of variables  $\mathbf{z} - \mathbf{y}^* = (-\partial_{\mathbf{y}\mathbf{y}}^2 f)^{-1/2} (\mathbf{w} - \mathbf{y}^*)$  and using Cauchy–Schwarz inequality, we obtain  $-\mathbf{t}^\top \partial_{\mathbf{x}\mathbf{y}}^2 f (\partial_{\mathbf{y}\mathbf{y}}^2 f)^{-1} (\partial_{\mathbf{y}\mathbf{x}}^2 f) \mathbf{t}$ . It follows that  $\partial_{\mathbf{x}\mathbf{x}}^2 f - \partial_{\mathbf{x}\mathbf{y}}^2 f (\partial_{\mathbf{y}\mathbf{y}}^2 f)^{-1} \partial_{\mathbf{y}\mathbf{x}}^2 f \succeq \mathbf{0}$ . ■

Finally, we can compare our second-order necessary condition with Proposition 19 of Jin et al. (2020), which applies to quadratic functions (cf. Remark 4.2). The difference is that Proposition 19 of Jin et al. (2020) did not take the critical directions and higher-order derivatives into consideration, as demonstrated by Examples 3.18 and 3.20.

#### 3.2.4 SECOND-ORDER SUFFICIENT CONDITIONS

We introduce two second-order sufficient conditions for local minimax points, with the help of results from non-smooth optimization literature (Seeger, 1988; Kawasaki, 1992). Our results extend Jin et al. (2020) to the case when  $\partial_{\mathbf{y}\mathbf{y}}^2 f$  is not invertible, which may happen in real applications.

In the following theorem, we define  $x_+ = \max\{x, 0\}$  and the first order activation set:

$$\mathcal{Y}_1(\mathbf{x}^*; \epsilon; \mathbf{t}) = \{\mathbf{y} \in \mathcal{Y}_0(\mathbf{x}^*, \epsilon) : \text{D}\bar{f}_\epsilon(\mathbf{x}^*; \mathbf{t}) = \partial_{\mathbf{x}} f(\mathbf{x}^*, \mathbf{y})^\top \mathbf{t}\}. \quad (3.26)$$

**Theorem 3.22 (second-order sufficient condition, local minimax)** *Assume  $\mathcal{X} = \mathbb{R}^n$  and  $\mathcal{Y}$  is convex and  $f, \partial_{\mathbf{x}} f, \partial_{\mathbf{x}\mathbf{x}}^2 f$  are (jointly) continuous. At a stationary point  $(\mathbf{x}^*, \mathbf{y}^*)$ , if there exists  $\epsilon_0 > 0$  such that:*

- $f(\mathbf{x}^*, \cdot)$  is maximized at  $\mathbf{y}^*$  on  $\mathcal{N}(\mathbf{y}^*, \epsilon_0)$ ;
- along each critical direction  $\mathbf{t} \neq \mathbf{0}$ :

$$\mathbf{t}^\top \partial_{\mathbf{x}\mathbf{x}}^2 f(\mathbf{x}^*, \mathbf{y}^*) \mathbf{t} + \frac{1}{2} \limsup_{\mathbf{z} \rightarrow \mathbf{y}^*} \left( ((\partial_{\mathbf{x}} f(\mathbf{x}^*, \mathbf{z})^\top \mathbf{t})_+)^2 (f(\mathbf{x}^*, \mathbf{y}^*) - f(\mathbf{x}^*, \mathbf{z}))^\dagger \right) > 0, \quad (3.27)$$

and in any direction  $\mathbf{d} \in \mathbb{R}^m$ , there exist  $\alpha, \beta \neq 0$  and  $p, q > 0$  such that for every  $\mathbf{y} \in \mathcal{Y}_1(\mathbf{x}^*; \epsilon_0; \mathbf{t})$ , the following Taylor expansion holds:

$$f(\mathbf{x}^*, \mathbf{y} + \delta \mathbf{d}) = f(\mathbf{x}^*, \mathbf{y}) + \alpha \delta^p + o(\delta^p), \quad \partial_{\mathbf{x}} f(\mathbf{x}^*, \mathbf{y} + \delta \mathbf{d})^\top \mathbf{t} = \beta \delta^q + o(\delta^q), \quad (3.28)$$

then  $(\mathbf{x}^*, \mathbf{y}^*)$  is a local minimax point.

Note that in the statement above, the variables  $\alpha, \beta$  and  $p, q$  may depend on the direction  $\mathbf{d}$ . If  $f \in \mathcal{C}^\infty$  is smooth and both  $f(\mathbf{x}^*, \cdot)$  and  $\partial_{\mathbf{x}} f(\mathbf{x}^*, \cdot)^\top \mathbf{t}$  have non-zero Taylor expansions, then (3.28) is always true for every  $\mathbf{y} \in \mathcal{Y}_1(\mathbf{x}^*; \epsilon_0; \mathbf{t})$ . Here by ‘‘critical direction’’ we mean that  $\text{D}\bar{f}_\epsilon(\mathbf{x}^*; \mathbf{t}) = 0$  for some  $\epsilon_0 > 0$  and any  $\epsilon \in [0, \epsilon_0]$ , as discussed in Section 3.2.3. Another second-order sufficient condition for  $f \in \mathcal{C}^2$  is:

**Theorem 3.23 (second-order sufficient condition, local minimax)** *Assume  $f \in \mathcal{C}^2$  and let  $\mathcal{X}$  be convex. Suppose  $\mathbf{y}^*$  is a local maximizer of  $f(\mathbf{x}^*, \cdot)$  and that  $(\mathbf{x}^*, \mathbf{y}^*)$  is an interior stationary point. If there is  $\epsilon_0 > 0$  and for any  $\epsilon \in (0, \epsilon_0]$ , there exist  $R, r > 0$  such that for any feasible direction  $\|\mathbf{t}\| = 1$  that satisfies  $0 \leq \text{D}\bar{f}_\epsilon(\mathbf{x}^*; \mathbf{t}) \leq r$ , we have*

$$\begin{aligned} \max_{\mathbf{y} \in \mathcal{Y}_0(\mathbf{x}^*; \epsilon)} \max_{\substack{\mathbf{v} \in \mathcal{V}(\mathbf{x}^*, \mathbf{y}; \mathbf{t}) \\ \|\mathbf{v}\| \leq R}} \max_{\substack{\mathbf{w} \in \mathbf{K}_d(\Omega, \mathbf{y}; \mathbf{v}) \\ \|\mathbf{w}\| \leq R}} & \left\langle \begin{bmatrix} \partial_{\mathbf{xx}}^2 f(\mathbf{x}^*, \mathbf{y}) & \partial_{\mathbf{xy}}^2 f(\mathbf{x}^*, \mathbf{y}) \\ \partial_{\mathbf{yx}}^2 f(\mathbf{x}^*, \mathbf{y}) & \partial_{\mathbf{yy}}^2 f(\mathbf{x}^*, \mathbf{y}) \end{bmatrix} \begin{pmatrix} \mathbf{t} \\ \mathbf{v} \end{pmatrix}, \begin{pmatrix} \mathbf{t} \\ \mathbf{v} \end{pmatrix} \right\rangle + \\ & + \langle \partial_{\mathbf{y}} f(\mathbf{x}^*, \mathbf{y}), \mathbf{w} \rangle > \mathbf{0}, \end{aligned} \quad (3.29)$$

then this point is local minimax, where  $\mathcal{V}(\mathbf{x}, \mathbf{y}; \mathbf{t}) := \{\mathbf{v} \in \mathbf{K}_d(\Omega, \mathbf{y}) : \text{D}\bar{f}_\epsilon(\mathbf{x}; \mathbf{t}) = \partial_{\mathbf{x}} f(\mathbf{x}, \mathbf{y})^\top \mathbf{t} + \partial_{\mathbf{y}} f(\mathbf{x}, \mathbf{y})^\top \mathbf{v}\}$ ,  $\Omega := \mathcal{N}(\mathbf{y}^*, \epsilon)$  and

$$\begin{aligned} \mathbf{K}_d(\Omega, \mathbf{y}; \mathbf{v}) := \liminf_{t \rightarrow 0^+} \frac{\Omega - \mathbf{y} - t\mathbf{v}}{t^2/2} & := \{\mathbf{g} : \forall \{t_k\} \downarrow 0 \exists \{t_{k_i}\} \downarrow 0, \{\mathbf{g}_{k_i}\} \rightarrow \mathbf{g}, \\ & \mathbf{y} + t_{k_i} \mathbf{v} + t_{k_i}^2 \mathbf{g}_{k_i} / 2 \in \Omega\}. \end{aligned} \quad (3.30)$$

The definition of feasible directions for convex sets can be found in e.g. Hiriart-Urruty and Lemaréchal (2013). We used the convention that maximizing over an empty set yields  $-\infty$ . Specifically, if there exists  $\mathbf{y} \in \mathcal{Y}_0(\mathbf{x}^*, \epsilon)$  such that it is in the interior of  $\mathcal{Y}$ , Theorem 3.23 can be simplified as:

**Corollary 3.24 (second-order sufficient condition, interior version)** *Assume  $f \in \mathcal{C}^2$  and let  $\mathcal{X}$  be convex. Suppose  $\mathbf{y}^*$  is a local maximizer of  $f(\mathbf{x}^*, \cdot)$  and that  $(\mathbf{x}^*, \mathbf{y}^*)$  is an interior stationary point. If there is  $\epsilon_0 > 0$  such that  $\mathcal{N}(\mathbf{y}^*, \epsilon_0) \subset \mathcal{Y} \subset \mathbb{R}^m$ , and for any  $\epsilon \in (0, \epsilon_0)$ , there exist  $R, r > 0$  such that for any feasible direction  $\|\mathbf{t}\| = 1$  that satisfies  $0 \leq \text{D}\bar{f}_\epsilon(\mathbf{x}^*; \mathbf{t}) \leq r$ , we have:*

$$\max_{\mathbf{y} \in \mathcal{Y}_0(\mathbf{x}^*; \epsilon)} \max_{\substack{\mathbf{v} \in \mathcal{V}(\mathbf{x}^*, \mathbf{y}; \mathbf{t}) \\ \|\mathbf{v}\| \leq R}} \max_{\|\mathbf{w}\| \leq R} \left\langle \begin{bmatrix} \partial_{\mathbf{xx}}^2 f(\mathbf{x}^*, \mathbf{y}) & \partial_{\mathbf{xy}}^2 f(\mathbf{x}^*, \mathbf{y}) \\ \partial_{\mathbf{yx}}^2 f(\mathbf{x}^*, \mathbf{y}) & \partial_{\mathbf{yy}}^2 f(\mathbf{x}^*, \mathbf{y}) \end{bmatrix} \begin{pmatrix} \mathbf{t} \\ \mathbf{v} \end{pmatrix}, \begin{pmatrix} \mathbf{t} \\ \mathbf{v} \end{pmatrix} \right\rangle + \langle \partial_{\mathbf{y}} f(\mathbf{x}^*, \mathbf{y}), \mathbf{w} \rangle > \mathbf{0}, \quad (3.31)$$

then this point is local minimax, where  $\mathcal{V}(\mathbf{x}, \mathbf{y}; \mathbf{t}) := \{\mathbf{v} \in \mathbb{R}^m : \text{D}\bar{f}_\epsilon(\mathbf{x}; \mathbf{t}) = \partial_{\mathbf{x}} f(\mathbf{x}, \mathbf{y})^\top \mathbf{t} + \partial_{\mathbf{y}} f(\mathbf{x}, \mathbf{y})^\top \mathbf{v}\}$ .

**Proof** If  $\mathbf{y} \in \mathcal{N}(\mathbf{y}^*, \epsilon)$ , then we have  $\mathbf{K}_d(\mathbf{y}) = \mathbf{K}_d(\Omega, \mathbf{y}; \mathbf{v}) = \mathbb{R}^m$ . ■

In the special case when  $\partial_{\mathbf{yy}}^2 f(\mathbf{x}^*, \mathbf{y}^*) \prec \mathbf{0}$ , we have the following corollary. This special type of local minimax points that satisfy (3.32) are also known as *strict local minimax points* (Jin et al., 2020).

**Corollary 3.25 (second-order sufficient condition, invertible, Jin et al. (2020))**  
 Let  $f$  be twice continuously differentiable. At an interior stationary point  $(\mathbf{x}^*, \mathbf{y}^*) \in \mathcal{X} \times \mathcal{Y}$ , if

$$\partial_{\mathbf{y}\mathbf{y}}^2 f \prec \mathbf{0} \text{ and } \partial_{\mathbf{x}\mathbf{x}}^2 f - \partial_{\mathbf{x}\mathbf{y}}^2 f (\partial_{\mathbf{y}\mathbf{y}}^2 f)^{-1} \partial_{\mathbf{y}\mathbf{x}}^2 f \succ \mathbf{0}, \quad (3.32)$$

then  $(\mathbf{x}^*, \mathbf{y}^*)$  is a local minimax point.

**Proof** The active set  $\mathcal{Y}_0(\mathbf{x}^*; \epsilon) = \{\mathbf{y}^*\}$  is a singleton. From Danskin's theorem (Theorem A.9) all directions are critical. The l.h.s. of (3.29) becomes  $\mathbf{t}^\top (\partial_{\mathbf{x}\mathbf{x}}^2 f - \partial_{\mathbf{x}\mathbf{y}}^2 f (\partial_{\mathbf{y}\mathbf{y}}^2 f)^{-1} \partial_{\mathbf{y}\mathbf{x}}^2 f) \mathbf{t}$  if we choose  $R = \|(\partial_{\mathbf{y}\mathbf{y}}^2 f)^{-1} \partial_{\mathbf{y}\mathbf{x}}^2 f\|$ .  $\blacksquare$

However, Corollary 3.25 does not fully cover Theorem 3.23 when  $\partial_{\mathbf{y}\mathbf{y}}^2$  is not invertible:

**Example 3.26 (Theorem 3.23 strictly includes Corollary 3.25)** Take

$$f(x, y) = xy^2 + x^2$$

and a stationary point  $(x^*, y^*) = (0, 0)$ .  $D\bar{f}_\epsilon(x^*; t) = \epsilon^2$  if  $t = 1$  and  $D\bar{f}_\epsilon(x^*; t) = 0$  if  $t = -1$ . Take  $r = \epsilon^2/2$ . Along the critical direction  $t = -1$ , the l.h.s. of (3.29) becomes  $2 > 0$ , since  $\partial_y f(x^*, y) = 0$ , and  $\mathcal{V}(x^*, y; t) = \emptyset$  if  $y \neq 0$  and  $\mathbb{R}$  if  $y = 0$ . So,  $(0, 0)$  is local minimax from Theorem 3.23. Note that Theorem 3.22 does not apply since  $f(x^*, y)$  does not have a non-zero Taylor expansion.

We also give an example when Theorem 3.23 is not applicable but Theorem 3.22 is:

**Example 3.27 (application of Theorem 3.22 where Theorem 3.23 cannot be applied)** Take

$$f(x, y) = xy^3 - y^6$$

and a stationary point  $(x^*, y^*) = (0, 0)$ . Fixing  $x^* = 0$ ,  $f(x^*, \cdot)$  is maximized at 0, and for any  $t \neq 0$ ,  $D\bar{f}_\epsilon(x^*; t) = \max_{y^6=0} y^3 t = 0$ . Since  $\partial_x f(x^*, z) = z^3 t$  and  $f(x^*, y^*) - f(x^*, z) = z^6$ , the l.h.s. of (3.27) is  $t^2/2 > 0$ . Moreover,  $\mathcal{Y}_1(x^*; \epsilon_0; t) = \{y^*\}$  for any  $\epsilon_0 > 0$ , and

$$f(x^*, y^* + \delta d) = -\delta^6 d^6, \partial_x f(x^*, y^* + \delta d)^\top t = \delta^3 d^3 t.$$

So,  $(0, 0)$  is a local minimax point. Note that Theorem 3.23 does not apply since  $\mathcal{Y}_0(x^*; \epsilon) = \{0\}$  and all second-order derivatives are zero.

#### 4. Quadratic games: A case study

In this section we study quadratic games with the following form:

$$q(\mathbf{x}, \mathbf{y}) = \frac{1}{2} \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \\ 1 \end{bmatrix}^\top \begin{bmatrix} \mathbf{A} & \mathbf{C} & \mathbf{a} \\ \mathbf{C}^\top & \mathbf{B} & \mathbf{b} \\ \mathbf{a}^\top & \mathbf{b}^\top & c \end{bmatrix} \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \\ 1 \end{bmatrix}, \quad (4.1)$$

where  $\mathbf{x} \in \mathcal{X} = \mathbb{R}^n$  and  $\mathbf{y} \in \mathcal{Y} = \mathbb{R}^m$ . In particular, a game is *bilinear* if  $\mathbf{A}, \mathbf{B}$  vanish and *homogeneous* if  $\mathbf{a}, \mathbf{b}$  vanish. Since quadratic games are continuous, local saddle points are the same as uniformly local minimax points (see Proposition 3.7).

Our first result completely characterizes stationary, global minimax and local minimax points for homogeneous quadratic games:

**Theorem 4.1 (sufficient and necessary conditions for optimality in quadratic games)** For (homogeneous) unconstrained quadratic games, a pair  $(\mathbf{x}, \mathbf{y})$  is

- stationary iff

$$\begin{bmatrix} \mathbf{A} & \mathbf{C} \\ \mathbf{C}^\top & \mathbf{B} \end{bmatrix} \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} = \mathbf{0}; \quad (4.2)$$

- global minimax iff  $\mathbf{B} \preceq \mathbf{0}$ ,  $\mathbf{P}_\mathbf{L}^\perp(\mathbf{A} - \mathbf{C}\mathbf{B}^\dagger\mathbf{C}^\top)\mathbf{P}_\mathbf{L}^\perp \succeq \mathbf{0}$  where  $\mathbf{L} = \mathbf{C}\mathbf{P}_\mathbf{B}^\perp$ , and

$$\begin{bmatrix} \mathbf{P}_\mathbf{L}^\perp & \\ & \mathbf{I} \end{bmatrix} \begin{bmatrix} \mathbf{A} & \mathbf{C} \\ \mathbf{C}^\top & \mathbf{B} \end{bmatrix} \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} = \mathbf{0}; \quad (4.3)$$

(Recall that  $\mathbf{P}_\mathbf{L}^\perp = \mathbf{I} - \mathbf{L}\mathbf{L}^\dagger$  is the orthogonal projection onto the null space of  $\mathbf{L}^\top$ .)

- local minimax iff  $\mathbf{B} \preceq \mathbf{0}$ ,  $\mathbf{P}_\mathbf{L}^\perp(\mathbf{A} - \mathbf{C}\mathbf{B}^\dagger\mathbf{C}^\top)\mathbf{P}_\mathbf{L}^\perp \succeq \mathbf{0}$ , and stationary (i.e. (4.2) holds). In particular, local minimax points are always global minimax.

Comparing Theorem 4.1 with Theorem 3.10, we find that in both cases, local minimax points are global minimax, which is not true in general (Example 4.9). This shows that there exists some “hidden convexity” in quadratic games when local/global minimax points exist: fixing any  $\mathbf{x}$ ,  $q(\mathbf{x}, \cdot)$  is concave in  $\mathbf{y}$ ;  $\bar{q}(\mathbf{x})$  is convex in  $\mathbf{x}$  (see (C.4)).

**Remark 4.2 (application of Theorem 3.17 in quadratic games)** We could also use Theorem 3.17 to obtain the necessary condition of local minimax points for quadratic games. First write

$$f(\mathbf{x}^*, \mathbf{y}^*) - f(\mathbf{x}^*, \mathbf{y}) = -\mathbf{y}^\top \mathbf{B}\mathbf{y}/2 \text{ and } -\partial_{\mathbf{x}}f(\mathbf{x}^*, \mathbf{y})^\top \mathbf{t} = -\mathbf{y}^\top \mathbf{C}^\top \mathbf{t}$$

and  $D\bar{f}_\epsilon(\mathbf{x}^*; \mathbf{t}) \geq \delta \|\mathbf{P}_\mathbf{B}^\perp \mathbf{C}^\top \mathbf{t}\|$  for some  $\delta > 0$ . The critical directions are  $\mathbf{t} \in \mathcal{N}(\mathbf{P}_\mathbf{B}^\perp \mathbf{C}^\top)$ . If  $\mathbf{B}\mathbf{C}^\top = \mathbf{0}$ , then  $\partial_{\mathbf{x}}f(\mathbf{x}^*, \mathbf{y})^\top \mathbf{t} = 0$  for any  $\mathbf{y}$  and thus the second term in (3.18) is zero. So, we have  $\mathbf{P}_\mathbf{L}^\perp \mathbf{A} \mathbf{P}_\mathbf{L}^\perp \succeq \mathbf{0}$  with  $\mathbf{L} = \mathbf{C}\mathbf{P}_\mathbf{B}^\perp$ . Otherwise, take critical directions  $\mathbf{t}$  such that  $\mathbf{t} \in \mathcal{N}(\mathbf{P}_\mathbf{B}^\perp \mathbf{C}^\top)$ . The second term in (3.18) becomes  $-\mathbf{t}^\top \mathbf{C}\mathbf{B}^\dagger \mathbf{C}^\top \mathbf{t}$  (using Cauchy–Schwarz). Combining with the case  $\mathbf{B}\mathbf{C}^\top = \mathbf{0}$ , we have  $\mathbf{P}_\mathbf{L}^\perp(\mathbf{A} - \mathbf{C}\mathbf{B}^\dagger\mathbf{C}^\top)\mathbf{P}_\mathbf{L}^\perp \succeq \mathbf{0}$ .

We remark that the last claim of Theorem 4.1 does not follow from Theorem 3.10:

**Example 4.3 (quadratic games can be non-convex)** Let  $A = -1, C = 1, B = 0, a = b = 0$ . Then, from Theorem 4.1  $(x, y) = (0, 0)$  is local and global minimax. However,  $q(x, y) = -\frac{1}{2}x^2 + xy$  is clearly non-convex in  $x$  (although  $\bar{q}$  is convex). Also,  $(0, 0)$  is not local saddle since  $q(x, 0) \geq q(0, 0)$  does not hold.

**Theorem 4.4 (equivalence between global and local minimax in quadratic games)**

An unconstrained quadratic game admits a global minimax point iff it admits a local minimax point iff

$$\mathbf{B} \preceq \mathbf{0}, \quad \mathbf{P}_\mathbf{L}^\perp(\mathbf{A} - \mathbf{C}\mathbf{B}^\dagger\mathbf{C}^\top)\mathbf{P}_\mathbf{L}^\perp \succeq \mathbf{0}, \text{ and } \begin{bmatrix} \mathbf{a} \\ \mathbf{b} \end{bmatrix} \in \mathcal{R} \left( \begin{bmatrix} \mathbf{A} & \mathbf{C} \\ \mathbf{C}^\top & \mathbf{B} \end{bmatrix} \right). \quad (4.4)$$

For such quadratic games, local minimax points are exactly the same as stationary global minimax points.

In this theorem we used  $\mathcal{R}(\cdot)$  to denote the range of a matrix. It is clear that stationary points, global minimax points, and local minimax points are characterized in the same way as in Theorem 4.1: we need only replace  $\mathbf{0}$  on the right-hands of (4.2) and (4.3) with the vector  $[\mathbf{a}; \mathbf{b}]$ . These points always form an affine subspace for quadratic games.

Theorem 4.4 allows us to completely classify (unconstrained) quadratic games:

- there are no stationary points (hence no local or global minimax points);
- there exist stationary points but no global or local minimax point;
- there exist local minimax points which coincide with global minimax points;
- there exist local minimax points which are strictly contained in global minimax points.

Clearly, for homogeneous (unconstrained) quadratic games, stationary points always exist hence only the last three cases can happen. For (non-trivial) bilinear games, only the last case can happen:

**Corollary 4.5 (bilinear games)** *For (homogeneous) unconstrained bilinear games ( $\mathbf{A} = \mathbf{0}, \mathbf{B} = \mathbf{0}, \mathbf{C} \neq \mathbf{0}, \mathbf{a} = \mathbf{0}, \mathbf{b} = \mathbf{0}$ ), global minimax points are  $\text{null}(\mathbf{C}^\top) \times \mathbb{R}^n$  while local minimax points (i.e. stationary points) are  $\text{null}(\mathbf{C}^\top) \times \text{null}(\mathbf{C})$ .*

It is thus clear that even in bilinear games, there exist global minimax points that are not local minimax. From Theorem 4.4, we can derive that:

**Corollary 4.6 (saddle points in quadratic games)** *For (unconstrained) quadratic games, the following statements are equivalent:*

1. *Local saddle points exist.*
2. *Local maximin and minimax points exist.*
3. *Global saddle points exist.*
4. *Global maximin and minimax points exist.*
5.  $\mathbf{A} \succeq \mathbf{0} \succeq \mathbf{B}$ , and

$$\begin{bmatrix} \mathbf{a} \\ \mathbf{b} \end{bmatrix} \in \mathcal{R} \left( \begin{bmatrix} \mathbf{A} & \mathbf{C} \\ \mathbf{C}^\top & \mathbf{B} \end{bmatrix} \right). \tag{4.5}$$

6. *stationary points exist and they are all local (global) saddle.*

Note that we used  $\mathcal{R}(\cdot)$  to denote the range of a matrix. We remark that Corollary 4.6 does not follow from typical minimax theorems (such as Sion's) since our domain is unbounded and we do not assume convexity-concavity from the outset. Thus, Corollary 4.6 reveals strong duality under weaker assumptions than the usual convexity-concavity. This is in stark contrast with generic NCNC games (see Example 2.6).

**Remark 4.7 (non-uniformly local minimax in quadratic games)** *Since quadratic functions are continuous (and thus upper semi-continuous), from Proposition 3.7 we know that local saddle points are equivalent to uniformly minimax points. By comparing Corollary 4.6 and Theorem 4.4, whenever  $\mathbf{A} \succeq \mathbf{0} \succeq \mathbf{B}$  and (4.5) holds, local saddle points and thus uniformly local minimax points exist. However, if (4.4) holds but  $\mathbf{A} \succeq \mathbf{0}$  does not hold, local saddle points/uniformly local minimax points do not exist from Corollary 4.6, but local minimax points still exist from Theorem 4.4 which are hence non-uniform. We can see it more clearly from Example 4.3. One can compute  $\bar{q}_\epsilon(x) = \epsilon|x| - \frac{1}{2}x^2$ , and obtain that  $\bar{q}_\epsilon(x) \geq \bar{q}_\epsilon(0) = 0$  iff  $|x| \leq 2\epsilon$ . According to Definition 3.3 the point  $(0, 0)$  is non-uniformly local minimax.*

Corollary 4.6 reveals some fundamental and surprising properties of quadratic games. On the one hand, quadratic games consist of an important theoretical tool for understanding general smooth NCNC games (through local Taylor expansion) (e.g. Daskalakis and Panageas, 2018; Jin et al., 2020; Ibrahim et al., 2020; Wang et al., 2020); see also Section 5 below. On the other hand, they are really special and many of their unique properties do not carry over to general smooth NCNC games, as we demonstrate in the following examples:

**Example 4.8 (stationary/global minimax points exist, no local minimax points)** *For general NCNC games, the existence of a global minimax point may not imply the existence of local minimax points. Indeed, consider*

$$f(x, y) = -y^4/4 + y^2/2 - xy, \quad x \in \mathbb{R}, \quad y \in \mathbb{R}. \quad (4.6)$$

*We claim  $(\pm 1, 0)$  are the only global minimax points. Indeed,*

$$\bar{f}(x) = \max_y -y^4/4 + y^2/2 - xy = \max_{y \geq 0} -y^4/4 + y^2/2 + |x|y \geq \max_{y \geq 0} -y^4/4 + y^2/2 = 1/4.$$

*Clearly, the inequality is attained only at  $x_* = 0$  and  $y_* = \pm 1$ . Its only stationary point is  $(x, y) = (0, 0)$ . However,  $\partial_{yy}^2 f(0, 0) = 1$  hence  $y = 0$  cannot be a local maximizer of  $f(0, \cdot)$ .*

*Note that in this example the global minimax points are not stationary. For an example where a stationary and global minimax point exists with no local minimax point, please refer to Example 3.11.*

**Example 4.9 (local minimax exists, no global minimax)** *This is possible even for separable functions, such as  $f(x, y) = x^3 - x - y^2$  defined on  $\mathbb{R} \times \mathbb{R}$ . Clearly, it has a local minimax point at  $(1/\sqrt{3}, 0)$  but no global minimax points exist.*

**Example 4.10 (local minimax and local maximin points exist; no local saddle)** *We can also construct an example when both local minimax and local maximin points exist but there is no local saddle point. Take  $f_1(x, y) = g(x, y)h(x, y)$ , where*

$$g(x, y) = xy - x^2, \quad \text{and} \quad h(x, y) = \exp\left(-\frac{1}{1-x^2}\right) \mathbf{1}_{|x|<1} \exp\left(-\frac{1}{1-y^2}\right) \mathbf{1}_{|y|<1}$$

*is a bump function that smoothly interpolates between the unit box and the outside. By numerically computing the stationary points and checking the second-order conditions, we found there is no such a point where  $\partial_{xx}^2 f_1 \geq 0$  and  $\partial_{yy}^2 f_1 \leq 0$  in the open box  $\mathbb{B}_1 = \{(x, y) :$*

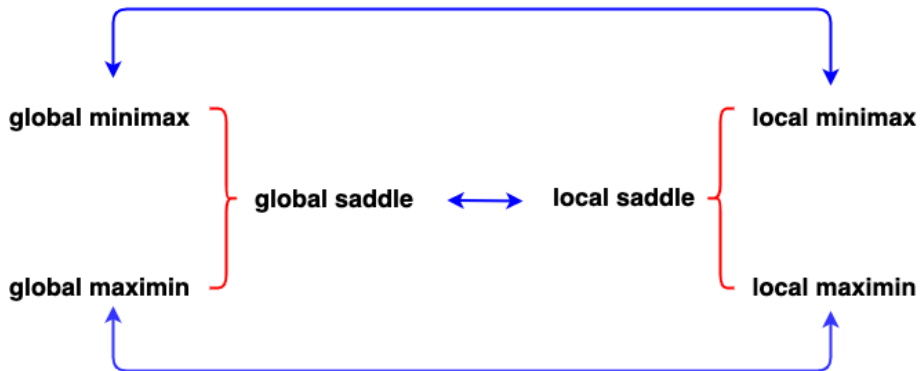


Figure 2 The relation among definitions in quadratic games.  $A \longleftrightarrow B$  means  $A$  exists iff  $B$  exists. The brackets also show the existence relation. For example, global saddle points exist iff both global minimax and maximin points exist.

$|x| < 1, |y| < 1\}$ . In other words, local saddle points do not exist. There is a local minimax point  $(0, 0)$  since

$$\bar{f}_\epsilon(x) \geq (\epsilon|x| - x^2) \exp(-1/(1 - x^2)) \exp(-1/(1 - \epsilon^2)) \geq 0$$

when  $|x| \leq \epsilon$  and  $\epsilon^2 < 1$ . Similarly we can construct  $f_2(x, y) = -g(y - 10, x - 10)h(x - 10, y - 10)$  where there is a local maximin point but no local saddle point in the open box  $\mathbb{B}_2 = \{(x, y) : |x - 10| < 1, |y - 10| < 1\}$ . Therefore,  $f(x, y) = f_1(x, y) + f_2(x, y)$  has both local minimax and local maximin points, but there is no local saddle point on  $\mathbb{B}_1 \cup \mathbb{B}_2$ .

Some special properties for quadratic games in this subsection are illustrated in Figure 2.

## 5. Stability of gradient algorithms near local optimal points

In this section, we assume that  $\mathcal{X} = \mathbb{R}^n$ ,  $\mathcal{Y} = \mathbb{R}^m$  and that  $f$  is twice continuously differentiable ( $f \in \mathcal{C}^2$ ). From (3.11) we know that local minimax points are stationary points, and thus fixed points of gradient algorithms. We focus on *local linear convergence* around stationary points using spectral analysis. Spectral analysis of a matrix  $A$  mainly involves two types of quantities: the spectrum of  $A$ ,  $\text{Sp}(A) := \{\lambda : \lambda \text{ is an eigenvalue of } A\}$ , as well as the spectral radius,  $\rho(A) := \max_{\lambda \in \text{Sp}(A)} |\lambda|$ . An iterative algorithm is *exponentially stable* if the Jacobian matrix of its update function has a spectral radius of less than one, which guarantees local linear convergence (Polyak, 1987). A more rigorous definition uses the Hartman–Grobman theorem (e.g. Katok and Hasselblatt, 1995). Below when we refer to convergence, we always mean local linear convergence.

To obtain convergence near local minimax points, we consider two-time-scale (2TS)<sup>2</sup> gradient algorithms, as applied to GANs by Heusel et al. (2017). Also, Jin et al. (2020) proved the “equivalence” between the stable points of 2TS-GDA and strict local minimax points. The intuition is that 2TS algorithms help the convergence by taking a much larger

2. This terminology comes from analogy with the continuous training dynamics. In our paper we simply mean choosing two different step sizes.

step w.r.t. the variable  $\mathbf{y}$ . We denote  $\mathbf{z}_t = (\mathbf{x}_t, \mathbf{y}_t)$  and define the vector field for the gradient update

$$\mathbf{v}(\mathbf{z}) = (-\alpha_1 \partial_{\mathbf{x}} f(\mathbf{z}), \alpha_2 \partial_{\mathbf{y}} f(\mathbf{z})).$$

Local stability results can be obtained by analyzing the Jacobian of  $\mathbf{v}(\mathbf{z})$  at a stationary point  $(\mathbf{x}^*, \mathbf{y}^*)$ :

$$\mathbf{H}_{\alpha_1, \alpha_2} = \mathbf{H}_{\alpha_1, \alpha_2}(f) := \begin{bmatrix} -\alpha_1 \partial_{\mathbf{x}\mathbf{x}}^2 f & -\alpha_1 \partial_{\mathbf{x}\mathbf{y}}^2 f \\ \alpha_2 \partial_{\mathbf{y}\mathbf{x}}^2 f & \alpha_2 \partial_{\mathbf{y}\mathbf{y}}^2 f \end{bmatrix}. \quad (5.1)$$

Define  $\alpha_2 = \gamma \alpha_1$ , and  $\mathbf{H}_{\alpha_1, \alpha_2} = \alpha_1 \mathbf{H}_{1, \gamma}$ . Note that  $\mathbf{H}_{\alpha_1, \alpha_2}(f)$  may not be symmetric, hence its spectrum lies on the complex plane. We also define  $\mathbf{H} := \mathbf{H}_{\alpha, \alpha} / \alpha$  which is independent of  $\alpha$ . To characterize the stable set of an algorithm, we ask the following question:

Given hyper-parameters  $\{\mu_i\}_{i=0}^k$  (e.g. step size, momentum coefficient) of an algorithm  $\mathbf{A}$ , what exactly is the geometric characterization on the spectrum of  $\mathbf{H}_{\alpha_1, \alpha_2}$  such that  $\mathbf{A}$  is exponentially stable at  $\mathbf{z}^*$ ?

Similar questions have been asked in Niethammer and Varga (1983) for problems of linear equations, where the Jacobian is a constant matrix. Such geometric characterizations allow us to analyze the convergence near local saddle and local minimax points.

Even with two-time-scale modification, GDA (even with momentum) does not converge near local saddle points for bilinear games (Zhang and Yu, 2020). Therefore, we will focus on extra gradient methods in this work. For completeness, thorough treatment of GDA, heavy ball (HB) and Nesterov’s momentum (NAG) is included in Appendix D. Note that second- and zeroth-order algorithms (Zhang et al., 2021; Liu et al., 2020) have also been considered very recently for minimax problems but they are beyond the scope of our work.

Note that in this section we are mostly considering one type of algorithmic modification in sequential games using two-time-scale (except in Proposition 5.9). For non-convex sequential smooth games, it is possible to use alternating updates in algorithms as studied in e.g. Zhang and Yu (2020) for bilinear games. We leave such systematic study to future work.

### 5.1 Stable sets of Extra-gradient (EG) and Optimistic gradient descent (OGD)

We consider the generalized extra-gradient method  $\text{EG}(\alpha_1, \alpha_2, \beta)$  (Korpelevich, 1976) (the original version has  $\beta = 1$ ):

$$\mathbf{z}_{t+1} = \mathbf{z}_t + \mathbf{v}(\mathbf{z}_{t+1/2}) / \beta, \quad \mathbf{z}_{t+1/2} = \mathbf{z}_t + \mathbf{v}(\mathbf{z}_t). \quad (5.2)$$

and the generalized optimistic gradient descent (Peng et al., 2020) (denoted as  $\text{OGD}(k, \alpha_1, \alpha_2)$ ):

$$\mathbf{z}_{t+1} = \mathbf{z}_t + k \mathbf{v}(\mathbf{z}_t) - \mathbf{v}(\mathbf{z}_{t-1}). \quad (5.3)$$

In (5.2), we call the first equation to be the extra-gradient step and the second equation to be the gradient step. EG was recently studied in e.g. Mertikopoulos et al. (2019) for special NCNC games, and in Azizian et al. (2020a,b) for convex-concave settings using spectral analysis. OGD was originally proposed in Popov (1980) as the past extra-gradient method, and was recently studied in the GAN literature (e.g. Daskalakis et al., 2018). Hsieh et al. (2019); Mokhtari et al. (2019) showed a close connection between EG and OGD:



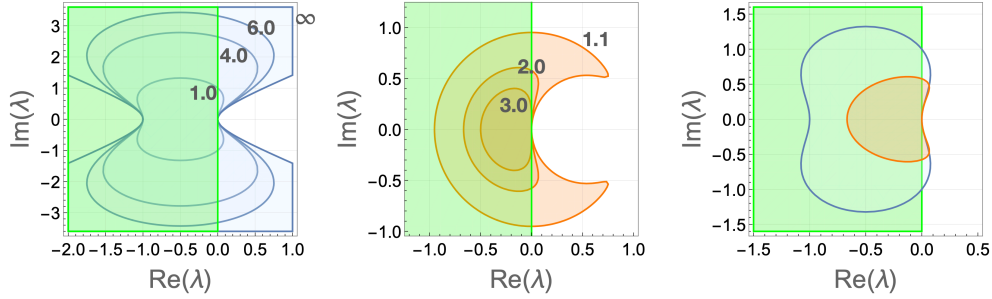


Figure 3 The blue/orange regions are where EG/OGD are exponentially stable. The green region represents where the eigenvalues of  $\text{Sp}(\mathbf{H}_{\alpha_1, \alpha_2})$  at local saddle points may occur. **(left)**  $\text{EG}(\alpha_1, \alpha_2, \beta)$  with  $\beta \in \{1.0, 4.0, 6.0, \infty\}$ ; **(middle)**  $\text{OGD}(k, \alpha_1, \alpha_2)$  with  $k \in \{1 + 1/10, 1 + 1/1, 1 + 1/0.5\}$ . **(right)** Comparison between  $\text{EG}(\alpha_1, \alpha_2, 1)$  (**blue**) and  $\text{OGD}(2, \alpha_1, \alpha_2)$  (**orange**). Best viewed in color.

**Lemma 5.1 (equivalence between past extra-gradient and OGD)** *The past extra-gradient method*

$$\mathbf{z}_{t+1} = \mathbf{z}_t + \mathbf{v}(\mathbf{z}_{t+1/2})/\beta, \quad \mathbf{z}_{t+1/2} = \mathbf{z}_t + \mathbf{v}(\mathbf{z}_{t-1/2}) \quad (5.4)$$

can be rewritten as  $\mathbf{z}'_{t+1} = \mathbf{z}'_t + k\mathbf{v}(\mathbf{z}'_t) - \mathbf{v}(\mathbf{z}'_{t-1})$  with  $k = 1 + 1/\beta$  and  $\mathbf{z}'_t = \mathbf{z}_{t-1/2}$ .

Due to this correspondence, we will only consider OGD with  $k > 1$ . We now characterize the stable sets of EG and OGD, or the *necessary and sufficient conditions* for local convergence (see the proof in Appendix E):

**Theorem 5.2 (stability of EG/OGD)** *At  $(\mathbf{x}^*, \mathbf{y}^*)$ ,  $\text{EG}(\alpha_1, \alpha_2, \beta)$  is exponentially stable iff for any  $\lambda \in \text{Sp}(\mathbf{H}_{\alpha_1, \alpha_2})$ ,  $|1 + \lambda/\beta + \lambda^2/\beta| < 1$ .  $\text{OGD}(k, \alpha_1, \alpha_2)$  is exponentially stable iff for any  $\lambda \in \text{Sp}(\mathbf{H}_{\alpha_1, \alpha_2})$ ,  $|\lambda| < 1$  and  $|\lambda|^2(k - 3 + (k + 1)|\lambda|^2) < 2\Re(\lambda)(k|\lambda|^2 - 1)$ .*

In this theorem,  $\Re(\cdot)$  represents the real part of a complex number. From this theorem, we can plot the stable region of EG and OGD with the original parameters ( $\beta = 1$  and  $k = 2$ ), and find that EG and OGD are indeed similar, as shown on the right of Figure 3. For EG, we note that Azizian et al. (2020b) used the spectral shapes of the support of  $\text{Sp}(\mathbf{H}_{\alpha_1, \alpha_2})$  to give upper and lower bounds of the convergence rates of EG, but our results are orthogonal to it since we do not assume a geometric shape of the support of  $\text{Sp}(\mathbf{H}_{\alpha_1, \alpha_2})$ .

When  $\beta \rightarrow \infty$ , we have  $k \rightarrow 1_+$ , and the step size of extra-gradient step is much larger than the step size of the gradient step. A similar conclusion can be found in Theorem 4.1 of Zhang and Yu (2020),<sup>3</sup> which states that for bilinear games, taking very small gradient steps and very large extra-gradient steps gives the best convergence rate among all hyper-parameter choices of gradient and extra-gradient steps.

Moreover, we show that larger  $\beta$  increases the local stability as well (see also Prop. 1', Hsieh et al. (2020) for a similar conclusion in saddle point problems, where  $\beta$  corresponds to  $\gamma_t/\eta_t$ ). The proof of the following theorem can be found in Appendix E:

3. Note that the exact definitions of  $\beta$  are different. Suppose the gradient step sizes are  $\alpha_1 = \alpha_2 = \alpha$ , and the extra-gradient step sizes are  $\gamma_1 = \gamma_2 = \gamma$ . Our definition gives  $\beta = \alpha/\gamma$  while Zhang and Yu (2020) gives  $\beta = \alpha\gamma$ .

**Theorem 5.3 (more aggressive extra-gradient steps, more stable)** *For  $\beta_1 > \beta_2 > 1$ , whenever  $\text{EG}(\alpha_1, \alpha_2, \beta_2)$  is exponentially stable at  $(\mathbf{x}^*, \mathbf{y}^*)$ ,  $\text{EG}(\alpha_1, \alpha_2, \beta_1)$  is exponentially stable at  $(\mathbf{x}^*, \mathbf{y}^*)$  as well. For  $k_1 > k_2 > 1$ , whenever  $\text{OGD}(k_1, \alpha_1, \alpha_2)$  is exponentially stable at  $(\mathbf{x}^*, \mathbf{y}^*)$ ,  $\text{OGD}(k_2, \alpha_1, \alpha_2)$  is exponentially stable at  $(\mathbf{x}^*, \mathbf{y}^*)$  as well.*

In the limit when  $\beta \rightarrow \infty$ , the stable region is  $\Re(\lambda + \lambda^2) < 0$  whose boundary is a hyperbola. Similarly, when  $k \rightarrow 1_+$ , OGD has the largest convergence region:  $\{\lambda \in \mathbb{C} : |\lambda| < 1, |\lambda - 1/2| > 1/2\}$ . Figure 3 visualizes the stable sets of EG/OGD. Their convergence regions strictly include that of GDA, and thus these algorithms are more stable:

**Corollary 5.4 (EG/OGD are more stable than GDA)** *When the step sizes  $\alpha_1, \alpha_2$  are small enough, whenever GDA converges, EG and OGD converge as well.*

The formal version of Corollary 5.4 can be found in Corollary E.1.

## 5.2 Local convergence to local optimal points

After characterizing the stable sets of EG and OGD, we move on to see the spectral behavior of local optimal points. For local saddle points, the spectrum of  $\mathbf{H}_{\alpha_1, \alpha_2}$  is on the left closed half plane. However, the spectrum of local minimax points (and thus LRPs, see Appendix F) can be quite arbitrary. With these results we can study how gradient algorithms (GDA with momentum, EG/OGD) converge to local optimal points.

### 5.2.1 LOCAL SADDLE POINTS

Even though the matrix  $\mathbf{H}_{\alpha_1, \alpha_2}(f)$  is not symmetric, it is still negative semi-definite near local saddle points.<sup>4</sup> Therefore, we can prove that its spectrum lies on the left (closed) complex plane:

**Lemma 5.5 (local saddle)** *Suppose  $\alpha_1, \alpha_2 > 0$  are fixed. For  $f \in \mathcal{C}^2$ , at a local saddle point, for all  $\lambda \in \text{Sp}(\mathbf{H}_{\alpha_1, \alpha_2}(f))$ , we have  $\Re(\lambda) \leq 0$ . For all  $z \in \mathbb{C}$  with  $\Re(z) \leq 0$ , there exists a quadratic function  $q$  and a local saddle point  $(\mathbf{x}^*, \mathbf{y}^*)$  such that  $z \in \text{Sp}(\mathbf{H}_{\alpha_1, \alpha_2}(q))$ . For bilinear functions, at a local saddle point we have  $\Re(\lambda) = 0$  for all  $\lambda \in \text{Sp}(\mathbf{H}_{\alpha_1, \alpha_2})$ .*

This result is a slight extension of Lemma 2.4 in Daskalakis and Panageas (2018). Combined with Lemma 5.5, we can show that EG converges around any local saddle point where the Jacobian  $\mathbf{H}(f)$  is non-singular, and a similar result holds for OGD if  $k$  is in a certain range:

**Theorem 5.6 (stability of EG/OGD at local saddle points)**  *$\text{EG}(\alpha, \alpha, 1)$  is exponentially stable at any local saddle point if at such a point,  $0 < |\lambda| < 1/\alpha$  for every  $\lambda \in \text{Sp}(\mathbf{H})$ .  $\text{OGD}(k, \alpha, \alpha)$  is exponentially stable at any local saddle point if  $1 < k \leq 2$  and  $0 < |\lambda| < 1/(k\alpha)$  for every  $\lambda \in \text{Sp}(\mathbf{H})$ . If  $k \geq 3$ ,  $\text{OGD}(k, \alpha_1, \alpha_2)$  is not exponentially stable for bilinear games.*

Given a fixed non-singular Jacobian matrix, we can always choose  $\alpha$  to be small enough, such that  $0 < |\lambda| < 1/\alpha$  (or  $0 < |\lambda| < 1/(k\alpha)$ ) for any  $\lambda \in \text{Sp}(\mathbf{H})$ . Therefore, EG and OGD always locally converge to any local saddle point as long as  $\mathbf{H}(f)$  is non-singular.

---

4. A real  $n \times n$  matrix  $\mathbf{A}$  is negative semi-definite if for any  $\mathbf{x} \in \mathbb{R}^n$ ,  $\mathbf{x}^\top \mathbf{A} \mathbf{x} \leq 0$ , i.e.  $\mathbf{A} + \mathbf{A}^\top$  is symmetric and negative semi-definite.

## 5.2.2 LOCAL MINIMAX POINTS

Now we study how gradient algorithms converge to local minimax points. We do not have the results in Theorem 5.6, since different from local saddle points, the spectrum of the Jacobian  $\mathbf{H}_{\alpha_1, \alpha_2}(f)$  is quite arbitrary:

**Lemma 5.7 (spectrum of local minimax can be arbitrary)** *Given  $\alpha_1, \alpha_2 > 0$ , for any  $z \in \mathbb{C}$ , there exists a quadratic function  $q$  and a local minimax point  $(\mathbf{x}^*, \mathbf{y}^*)$  where  $z \in \text{Sp}(\mathbf{H}_{\alpha_1, \alpha_2}(q))$ .*

This result shows that local minimax points are a more general class than the class of local stable stationary points (LSSPs) as studied recently in Berard et al. (2020), in terms of zero-sum games, since LSSPs are defined such that  $\Re(\lambda) < 0$  for any  $\lambda \in \text{Sp}(\mathbf{H}_{\alpha, \alpha})$  and  $\alpha > 0$  (note the slight change of signs due to the difference of notations). Under certain assumptions, 2TS gradient algorithms can converge to local minimax points. The following result slightly extends Jin et al. (2020) where only GDA is analyzed:

**Theorem 5.8 (stability of EG/OGD at strict local minimax points)** *Assume at a stationary point  $(\mathbf{x}^*, \mathbf{y}^*)$ ,*

$$\partial_{\mathbf{y}\mathbf{y}}^2 f \prec \mathbf{0} \text{ and } \partial_{\mathbf{x}\mathbf{x}}^2 f - \partial_{\mathbf{x}\mathbf{y}}^2 f (\partial_{\mathbf{y}\mathbf{y}}^2 f)^{-1} \partial_{\mathbf{y}\mathbf{x}}^2 f \succ \mathbf{0}. \quad (5.5)$$

*Then there exist  $\gamma_0 > 0$  and  $\alpha_0 > 0$  such that for any  $\gamma > \gamma_0, 0 < \alpha_2 < \alpha_0$  and  $\alpha_1 = \alpha_2/\gamma$ , EG and OGD (with  $k > 1$ ) are exponentially stable.*

In fact, the theorem above can be extended to momentum methods as well (see Appendix D). As we have seen in Corollary 3.25, (5.5) is sufficient for being local minimax (see also Fiez et al. (2019); Wang et al. (2020); Zhang et al. (2021) for applications in GANs). However, without assumption (5.5) (see also Jin et al. (2020, Theorem 28) for GDA), convergence is more difficult:

**Proposition 5.9 (stability of gradient algorithms at general local minimax points)**

*There exists a quadratic function (e.g.,  $q(x, y) = -x^2 + xy$ ) and a global (thus local, from Theorem 4.4) minimax point  $\mathbf{z}^* = (\mathbf{x}^*, \mathbf{y}^*)$  where*

- *GDA (with momentum or alternating updates) does not converge to  $\mathbf{z}^*$ , for any hyper-parameter choice.*
- *If  $\alpha_1 = \alpha_2$ , or  $\alpha_2 \rightarrow 0$ , EG/OGD do not converge to  $\mathbf{z}^*$ . Otherwise there exist hyper-parameter choices such that EG/OGD converge to  $\mathbf{z}^*$ .*
- *Alternating OGD does not converge to  $\mathbf{z}^*$  given  $\alpha_2 \rightarrow 0$ .*

The exact forms of alternating updates can be found in Zhang and Yu (2020) which we have also included in the proof of Proposition 5.9. It basically says that we update  $\mathbf{x}$  and  $\mathbf{y}$  one after the other rather than simultaneously. Proposition 5.9 extends Jin et al. (2020) by studying the degenerate case of  $\partial_{\mathbf{y}\mathbf{y}}^2 f$  and gradient algorithms other than GDA. The implication is two-fold:

- On the algorithmic aspect, we may not always rely on the usual ODE analysis (Mescheder et al., 2017; Mertikopoulos et al., 2018; Fiez et al., 2019) when trying to find global/local minimax points, as such analysis relies on approximating gradient algorithms with their continuous versions, by taking the step sizes to be arbitrarily small. For EG/OGD, the step size of the follower ( $\alpha_2$ ) has to be large while the step size of the leader can be arbitrarily small, reflecting the asymmetric position of players in Stackelberg games (Jin et al., 2020).
- We may also need new solution concepts in addition to global/local minimax points in machine learning applications (e.g. Farnia and Ozdaglar, 2020; Schaefer et al., 2020), even though many machine learning applications, including GANs (Goodfellow et al., 2014) and adversarial training (Madry et al., 2018) are essentially based on the notion of global minimax points. This is because when applying standard gradient-based algorithms to do a local search in machine learning applications, we cannot always expect the final solutions found by the algorithms to cover all global/local minimax points.

## 6. Conclusion

The aim of this work is to provide a comprehensive study of the recently proposed local minimax points (Jin et al., 2020). We discussed the relations between local saddle and local minimax points, between local and global minimax points, and interpreted local minimax points based on infinitesimal robustness. This new interpretation allows us to further generalize local minimax points such that they are still stationary (Theorem F.6). We presented the first- and second-order optimality conditions of these local optimal solutions, which extend Jin et al. (2020) to the constrained and degenerate cases. Specifically, in (potentially non-convex) quadratic games, local minimax points are (in some sense) equivalent to global minimax points. We also studied the stability of popular gradient algorithms near local optimal solutions, which provides insights for the design of algorithms to find minimax points.

The implication of this work is two-fold: **(a)** we may need new algorithms for smooth games, since we have shown in Proposition 5.9 that our common intuition might fail w.r.t. the convergence to a local and global minimax point; **(b)** we need to think about new solution concepts other than global/local minimax points. As many theoretical works aim to go beyond the definition of Nash equilibria (a.k.a. saddle points) such as Jin et al. (2020); Farnia and Ozdaglar (2020); Berard et al. (2020), to name a few, we may need to take one step further, beyond the definition of Stackelberg equilibria (a.k.a. minimax points), as also pointed out in Schaefer et al. (2020). Our new definition of local robust points sheds some light on going beyond Stackelberg games (Appendix F).

## Acknowledgments

We thank NSERC, the Canada CIFAR AI Chairs Program, Borealis AI and the Waterloo-Huawei Joint Innovation Lab for financial support. GZ is also supported by David R.

Cheriton scholarship and Vector research grant. We thank Chi Jin and Oliver Schulte for useful discussion.

## Appendix A. Non-smooth analysis: A short detour

We give a short detour on some classical optimality conditions in non-smooth optimization. These results will be used in Section 3 to yield necessary and sufficient conditions for local optimality in zero-sum two-player games, since the optimality conditions for local optimal points can be reduced to those for the envelope functions, which are in general non-smooth. A more thorough version of this appendix can be found in Zhang et al. (2020).

Let  $h$  be a function defined on some set  $\mathcal{X} \subseteq \mathbb{R}^m$ . Its upper and lower (Dini) directional derivatives are defined as:

$$D^+h(\mathbf{x}; \mathbf{d}) := \limsup_{t \rightarrow 0^+} \frac{h(\mathbf{x} + t\mathbf{d}) - h(\mathbf{x})}{t}, \quad D_+h(\mathbf{x}; \mathbf{d}) := \liminf_{t \rightarrow 0^+} \frac{h(\mathbf{x} + t\mathbf{d}) - h(\mathbf{x})}{t}. \quad (\text{A.1})$$

When the two limits coincide, we use the notation  $Dh(\mathbf{x}; \mathbf{d})$  and call the function  $h$  directionally differentiable (at  $\mathbf{x}$  along direction  $\mathbf{d}$ ). We can similarly define the upper and lower second-order directional derivatives<sup>5</sup> according to Ben-Tal and Zowe (1982):

$$Hh(\mathbf{x}; \mathbf{d}, \mathbf{g}) = \limsup_{t \rightarrow 0^+} \frac{h(\mathbf{x} + t\mathbf{d} + t^2\mathbf{g}/2) - h(\mathbf{x}) - t \cdot Dh(\mathbf{x}; \mathbf{d})}{t^2/2}, \quad (\text{A.2})$$

$$H_+h(\mathbf{x}; \mathbf{d}, \mathbf{g}) = \liminf_{t \rightarrow 0^+} \frac{h(\mathbf{x} + t\mathbf{d} + t^2\mathbf{g}/2) - h(\mathbf{x}) - t \cdot Dh(\mathbf{x}; \mathbf{d})}{t^2/2}. \quad (\text{A.3})$$

Similarly, when the two limits coincide we use the simplified notation  $Hh(\mathbf{x}; \mathbf{d}, \mathbf{g})$  and call  $h$  twice directionally differentiable (at  $\mathbf{x}$  along parabolic  $(\mathbf{d}, \mathbf{g})$ ). Note that, when  $\mathbf{d} = \mathbf{0}$ , we recover the directional derivative:

$$Hh(\mathbf{x}; \mathbf{0}, \mathbf{g}) = H_+h(\mathbf{x}; \mathbf{0}, \mathbf{g}) = Dh(\mathbf{x}; \mathbf{g}), \quad (\text{A.4})$$

while if  $\mathbf{g} = \mathbf{0}$ ,

$$Hh(\mathbf{x}; \mathbf{d}) := Hh(\mathbf{x}; \mathbf{d}, \mathbf{0}), \quad H_+h(\mathbf{x}; \mathbf{d}) := H_+h(\mathbf{x}; \mathbf{d}, \mathbf{0}), \quad Hh(\mathbf{x}; \mathbf{d}) := Hh(\mathbf{x}; \mathbf{d}, \mathbf{0}) \quad (\text{A.5})$$

reduces to the second-order directional derivatives of Dem'yanov (1973). The advantage of the definition of Ben-Tal and Zowe (1982) is evidenced in the following chain rule:

**Theorem A.1 (Ben-Tal and Zowe 1982)** *Let  $h : \mathbb{R}^m \rightarrow \mathbb{R}$  be locally Lipschitz and  $k : \mathbb{R}^n \rightarrow \mathbb{R}^m$  be (twice) directionally differentiable. Then,*

$$D^+(h \circ k)(\mathbf{x}; \mathbf{d}) = D^+h(k(\mathbf{x}); Dk(\mathbf{x}; \mathbf{d})), \quad (\text{A.6})$$

$$H(h \circ k)(\mathbf{x}; \mathbf{d}, \mathbf{g}) = Hh(k(\mathbf{x}); Dk(\mathbf{x}; \mathbf{d}), Hk(\mathbf{x}; \mathbf{d}, \mathbf{g})). \quad (\text{A.7})$$

(The same result holds for the lower derivatives, and hence the derivatives when they exist.)

5. A popular directional derivative in non-smooth analysis, due to Clarke (1990), is to replace  $h(\mathbf{x} + t\mathbf{d})$  with  $h(\mathbf{y} + t\mathbf{d})$  for some sequence  $\mathbf{y} \rightarrow \mathbf{x}$ . The second-order counterpart appeared in Cominetti and Correa (1990). For our purpose here, the classical Dini definitions suffice.

In contrast, the definition of Dem'yanov (1973) fails to satisfy the chain rule above. Indeed, if  $h$  is differentiable, then

$$Dh(\mathbf{x}; \mathbf{d}) = \langle \nabla h(\mathbf{x}), \mathbf{d} \rangle \quad (\text{A.8})$$

while if  $h$  is twice differentiable, then

$$Hh(\mathbf{x}; \mathbf{d}, \mathbf{g}) = Dh(\mathbf{x}; \mathbf{g}) + Hh(\mathbf{x}; \mathbf{d}) = \langle \nabla h(\mathbf{x}), \mathbf{g} \rangle + \langle \mathbf{d}, \nabla^2 h(\mathbf{x}) \mathbf{d} \rangle, \quad (\text{A.9})$$

where  $\nabla h$  and  $\nabla^2 h$  are the gradient and Hessian of  $h$ , respectively. (A slightly more general setting is discussed in Seeger 1988, Proposition 1.1.) The following properties of the directional derivatives are clear:

**Theorem A.2** *For any  $\lambda \geq 0$  we have*

$$Dh(\mathbf{x}; \lambda \mathbf{d}) = \lambda \cdot Dh(\mathbf{x}; \mathbf{d}), \quad (\text{A.10})$$

$$Hh(\mathbf{x}; \lambda \mathbf{d}, \lambda^2 \mathbf{g}) = \lambda^2 \cdot Hh(\mathbf{x}; \mathbf{d}, \mathbf{g}) \quad (\text{A.11})$$

*If  $h$  is locally Lipschitz around  $\mathbf{x}$ , then  $Dh(\mathbf{x}; \cdot)$  and  $Hh(\mathbf{x}; \mathbf{d}, \cdot)$  are Lipschitz continuous. (Similar results hold for the upper and lower derivatives.)*

### A.1 Necessary conditions

Consider the non-smooth optimization problem

$$\min_{\mathbf{x} \in \mathcal{X} \subseteq \mathbb{R}^m} h(\mathbf{x}). \quad (\text{A.12})$$

We define three tangent cones of the (closed) constraint set  $\mathcal{X}$ :

$$\mathbf{K}_f(\mathcal{X}, \mathbf{x}) := \{\mathbf{d} : \forall \{t_k\} \rightarrow 0^+ \exists \{t_{k_i}\} \rightarrow 0^+, \mathbf{x} + t_{k_i} \mathbf{d} \in \mathcal{X}\} \subseteq \text{cone}(\mathcal{X} - \mathbf{x}) \quad (\text{A.13})$$

$$\mathbf{K}_d(\mathcal{X}, \mathbf{x}) := \liminf_{t \rightarrow 0^+} \frac{\mathcal{X} - \mathbf{x}}{t} := \{\mathbf{d} : \forall \{t_k\} \rightarrow 0^+ \exists \{t_{k_i}\} \rightarrow 0^+, \{\mathbf{d}_{k_i}\} \rightarrow \mathbf{d}, \mathbf{x} + t_{k_i} \mathbf{d}_{k_i} \in \mathcal{X}\} \quad (\text{A.14})$$

$$\mathbf{K}_c(\mathcal{X}, \mathbf{x}) := \limsup_{t \rightarrow 0^+} \frac{\mathcal{X} - \mathbf{x}}{t} := \{\mathbf{d} : \exists \{t_k\} \rightarrow 0^+, \{\mathbf{d}_k\} \rightarrow \mathbf{d}, \mathbf{x} + t_k \mathbf{d}_k \in \mathcal{X}\}. \quad (\text{A.15})$$

Obviously, the (feasible) cone  $\mathbf{K}_f$  is contained in the (derivable) cone  $\mathbf{K}_d$ , which is itself contained in the (contingent) cone  $\mathbf{K}_c$ .  $\mathbf{K}_d$  and  $\mathbf{K}_c$  are always closed while  $\mathbf{K}_f$  may not be so (even when  $\mathcal{X}$  is closed). On the other hand, if  $\mathcal{X}$  is convex (and  $\mathbf{x} \in \mathcal{X}$ ), then all three tangent cones are convex,  $\mathbf{K}_f = \text{cone}(\mathcal{X} - \mathbf{x})$  and  $\mathbf{K}_d = \mathbf{K}_c = \overline{\mathbf{K}_f}$ . Note that for all tangent cones, we have

$$\forall \mathbf{x} \notin \bar{\mathcal{X}}, \mathbf{K}(\mathcal{X}, \mathbf{x}) = \emptyset, \text{ and } \forall \mathbf{x} \in \mathcal{X}^\circ, \mathbf{K}(\mathcal{X}, \mathbf{x}) = \mathbb{R}^m, \quad (\text{A.16})$$

where  $\bar{\mathcal{X}}$  and  $\mathcal{X}^\circ$  denote the closure and interior of  $\mathcal{X}$ , respectively. The following necessary condition is well-known:

**Theorem A.3 (first-order necessary condition, e.g. Dem'yanov (1966))** *Let  $\mathbf{x}^*$  be a local minimizer of  $h$  over  $\mathcal{X}$ . Then,*

$$\forall \mathbf{d} \in \mathbf{K}_f(\mathcal{X}, \mathbf{x}^*), \quad \mathbf{D}_+ h(\mathbf{x}^*; \mathbf{d}) \geq 0. \quad (\text{A.17})$$

*The converse is also true if  $h$  and  $\mathcal{X}$  are both convex around  $\mathbf{x}^*$ . If  $h$  is locally Lipschitz, then*

$$\forall \mathbf{d} \in \mathbf{K}_d(\mathcal{X}, \mathbf{x}^*), \quad \mathbf{D}_+ h(\mathbf{x}^*; \mathbf{d}) \geq 0. \quad (\text{A.18})$$

**Proof** We first prove the converse part. Suppose to the contrary there exists  $\mathbf{x}$  around  $\mathbf{x}^*$  so that  $h(\mathbf{x}) < h(\mathbf{x}^*)$ . Then,  $\mathbf{d} = \mathbf{x} - \mathbf{x}^* \in \mathbf{K}_f(\mathcal{X}, \mathbf{x}^*)$  and we have

$$\mathbf{D}_+ h(\mathbf{x}^*; \mathbf{d}) = \liminf_{t \rightarrow 0^+} \frac{h((1-t)\mathbf{x}^* + t\mathbf{x}) - h(\mathbf{x}^*)}{t} \leq h(\mathbf{x}) - h(\mathbf{x}^*) < 0, \quad (\text{A.19})$$

which is a contradiction.

To see the claim when  $h$  is locally Lipschitz, note that  $\mathbf{d} \in \mathbf{K}_d(\mathcal{X}, \mathbf{x}^*)$  implies for any  $\{t_k\} \rightarrow 0$  there exist  $\{t_{k_i}\} \rightarrow 0^+$  and  $\{\mathbf{d}_{k_i}\} \rightarrow \mathbf{d}$  such that  $\mathbf{x}^* + t_{k_i} \mathbf{d}_{k_i} \in \mathcal{X}$ . For sufficiently large  $k_i$  we have  $h(\mathbf{x}^* + t_{k_i} \mathbf{d}_{k_i}) \geq h(\mathbf{x}^*)$  since  $\mathbf{x}^*$  by assumption is a local minimizer. Thus,

$$\liminf_{t \rightarrow 0^+} \frac{h(\mathbf{x}^* + t\mathbf{d}) - h(\mathbf{x}^*)}{t} := \lim_{t_k \rightarrow 0^+} \frac{h(\mathbf{x}^* + t_k \mathbf{d}) - h(\mathbf{x}^*)}{t_k} \quad (\text{A.20})$$

$$\begin{aligned} &\geq \limsup_{t_{k_i} \rightarrow 0^+} \frac{h(\mathbf{x}^* + t_{k_i} \mathbf{d}_{k_i}) - h(\mathbf{x}^*)}{t_{k_i}} \\ &\quad - \limsup_{t_{k_i} \rightarrow 0^+} \frac{h(\mathbf{x}^* + t_{k_i} \mathbf{d}) - h(\mathbf{x}^* + t_{k_i} \mathbf{d}_{k_i})}{t_{k_i}} \end{aligned} \quad (\text{A.21})$$

$$\geq 0 - 0 = 0. \quad (\text{A.22})$$

The proof for a general function  $h$  is similar. ■

To derive second-order conditions, we define similarly the second-order tangent cones:

$$\mathbf{K}_f(\mathcal{X}, \mathbf{x}; \mathbf{d}) := \{\mathbf{g} : \forall \{t_k\} \downarrow 0 \exists \{t_{k_i}\} \downarrow 0, \mathbf{x} + t_{k_i} \mathbf{d} + t_{k_i}^2 \mathbf{g}/2 \in \mathcal{X}\}, \quad (\text{A.23})$$

$$\begin{aligned} \mathbf{K}_d(\mathcal{X}, \mathbf{x}; \mathbf{d}) &:= \liminf_{t \rightarrow 0^+} \frac{\mathcal{X} - \mathbf{x} - t\mathbf{d}}{t^2/2} \\ &:= \{\mathbf{g} : \forall \{t_k\} \downarrow 0 \exists \{t_{k_i}\} \downarrow 0, \{\mathbf{g}_{k_i}\} \rightarrow \mathbf{g}, \mathbf{x} + t_{k_i} \mathbf{d} + t_{k_i}^2 \mathbf{g}_{k_i}/2 \in \mathcal{X}\}. \end{aligned} \quad (\text{A.24})$$

The proof of the following result is completely similar to that of Theorem A.3:

**Theorem A.4 (second-order necessary condition, e.g. Ben-Tal and Zowe 1985)**

*Let  $h$  be directionally differentiable and  $\mathbf{x}^*$  be a local minimizer of  $h$  over  $\mathcal{X}$ . Then,*

$$\forall \mathbf{d} \in \mathbf{K}_f(\mathcal{X}, \mathbf{x}^*), \forall \mathbf{g} \in \mathbf{K}_f(\mathcal{X}, \mathbf{x}^*; \mathbf{d}), \quad \mathbf{D}h(\mathbf{x}^*; \mathbf{d}) = 0 \implies \mathbf{H}_+ h(\mathbf{x}^*; \mathbf{d}, \mathbf{g}) \geq 0. \quad (\text{A.25})$$

*If  $h$  is locally Lipschitz, then*

$$\forall \mathbf{d} \in \mathbf{K}_d(\mathcal{X}, \mathbf{x}^*), \forall \mathbf{g} \in \mathbf{K}_d(\mathcal{X}, \mathbf{x}^*; \mathbf{d}), \quad \mathbf{D}h(\mathbf{x}^*; \mathbf{d}) = 0 \implies \mathbf{H}_+ h(\mathbf{x}^*; \mathbf{d}, \mathbf{g}) \geq 0. \quad (\text{A.26})$$

## A.2 Sufficient conditions

We give sufficient conditions for a non-smooth function to attain an isolated minimum.

**Theorem A.5 (first-order, e.g. Dem'yanov 1970; Ben-Tal and Zowe 1985)** *Let  $h$  be locally Lipschitz. If*

$$\forall \mathbf{0} \neq \mathbf{d} \in \mathcal{K}_c(\mathcal{X}, \mathbf{x}^*), \quad D_+h(\mathbf{x}^*; \mathbf{d}) > 0, \quad (\text{A.27})$$

then  $\mathbf{x}^*$  is an isolated local minimum of  $h$  over  $\mathcal{X}$ .

**Proof** Suppose to the contrary there exists a sequence  $\mathbf{x}_k \in \mathcal{X}$  converging to  $\mathbf{x}^*$  so that  $h(\mathbf{x}_k) \leq h(\mathbf{x}^*)$ . Let  $t_k := \|\mathbf{x}_k - \mathbf{x}^*\|$  and  $\mathbf{d}_k := (\mathbf{x}_k - \mathbf{x}^*)/\|\mathbf{x}_k - \mathbf{x}^*\|$ . By passing to a subsequence we may assume  $\mathbf{d}_k \rightarrow \mathbf{d} \neq \mathbf{0}$ , where clearly  $\mathbf{d} \in \mathcal{K}_c(\mathcal{X}, \mathbf{x}^*)$  since  $\mathbf{x}^* + t_k \mathbf{d}_k = \mathbf{x}_k \in \mathcal{X}$ . But then

$$D_+h(\mathbf{x}^*; \mathbf{d}) \leq \liminf_{t_k \rightarrow 0^+} \frac{h(\mathbf{x}^* + t_k \mathbf{d}) - h(\mathbf{x}^*)}{t_k} \quad (\text{A.28})$$

$$\leq \liminf_{t_k \rightarrow 0^+} \frac{h(\mathbf{x}^* + t_k \mathbf{d}_k) - h(\mathbf{x}^*)}{t_k} + \limsup_{t_k \rightarrow 0^+} \frac{h(\mathbf{x}^* + t_k \mathbf{d}) - h(\mathbf{x}^* + t_k \mathbf{d}_k)}{t_k} \quad (\text{A.29})$$

$$\leq 0 + 0 = 0, \quad (\text{A.30})$$

arriving at a contradiction. ■

Note that when  $\mathcal{X}$  is convex, we may replace  $\mathcal{K}_c = \overline{\mathcal{K}_f}$  with  $\mathcal{K}_f$  (recall the Lipschitz continuity in Theorem A.2).

**Theorem A.6 (second-order, e.g. Dem'yanov 1970)** *Let  $h$  be locally Lipschitz and directional differentiable, and  $\mathcal{X}$  be convex. If*

1.  $\forall \mathbf{d} \in \mathcal{K}_f(\mathcal{X}, \mathbf{x}^*), \quad Dh(\mathbf{x}^*; \mathbf{d}) \geq 0,$
2.  $\exists \gamma > 0$  such that for all  $\mathbf{d} \in \mathcal{K}_f(\mathcal{X}, \mathbf{x}^*), \|\mathbf{d}\| = 1, Dh(\mathbf{x}^*; \mathbf{d}) \in [0, \gamma]$  we have for all small  $t$  and uniformly on bounded sets in  $\mathbf{d}$ :

$$\frac{h(\mathbf{x}^* + t\mathbf{d}) - h(\mathbf{x}^*) - tDh(\mathbf{x}^*; \mathbf{d})}{t^2/2} \geq A_h(\mathbf{x}^*; \mathbf{d}) > 0, \quad (\text{A.31})$$

then  $\mathbf{x}^*$  is an isolated local minimum of  $h$  over  $\mathcal{X}$ .

**Proof** Let  $\mathbf{x} \in \mathcal{X}$  and  $\mathbf{x} \neq \mathbf{x}^*$ , then  $\mathbf{d} := (\mathbf{x} - \mathbf{x}^*)/\|\mathbf{x} - \mathbf{x}^*\| \in \mathcal{K}_f(\mathcal{X}, \mathbf{x}^*)$  (since  $\mathcal{X}$  is convex). Suppose  $Dh(\mathbf{x}^*; \mathbf{d}) \geq \gamma > 0$ , then

$$h(\mathbf{x}^* + t\mathbf{d}) = h(\mathbf{x}^*) + tDh(\mathbf{x}^*; \mathbf{d}) + o(t) \geq h(\mathbf{x}^*) + \gamma t + o(t) > h(\mathbf{x}^*) + \gamma t/2, \quad (\text{A.32})$$

for sufficiently small  $t \leq t_{\mathbf{d}}$ . Since the function  $\mathbf{d} \mapsto h(\mathbf{x}^* + t\mathbf{d})$  is locally Lipschitz, we may choose a non-empty open subset from each set  $\{\mathbf{v} : \forall t \in (0, t_{\mathbf{d}}], h(\mathbf{x}^* + t\mathbf{v}) > h(\mathbf{x}^*)\}$ . Hence, using a standard compactness argument, we know for all small positive  $t$ ,

$$\mathbf{d} \in \mathcal{K}_f(\mathcal{X}, \mathbf{x}^*), \|\mathbf{d}\| = 1, Dh(\mathbf{x}^*; \mathbf{d}) \geq \gamma \implies h(\mathbf{x}^* + t\mathbf{d}) > h(\mathbf{x}^*). \quad (\text{A.33})$$



Suppose instead  $Dh(\mathbf{x}^*, \mathbf{d}) \in [0, \gamma]$ , then for all small positive  $t$  and uniformly in  $\mathbf{d}$  we have

$$h(\mathbf{x}^* + t\mathbf{d}) \geq h(\mathbf{x}^*) + tDh(\mathbf{x}^*; \mathbf{d}) + \frac{1}{2}t^2A_h(\mathbf{x}^*; \mathbf{d}) \quad (\text{A.34})$$

$$\geq h(\mathbf{x}^*) + \frac{1}{2}t^2A_h(\mathbf{x}^*; \mathbf{d}) \quad (\text{A.35})$$

$$> h(\mathbf{x}^*). \quad (\text{A.36})$$

Finally, combining the above two cases completes the proof.  $\blacksquare$

We make a few remarks regarding Theorem A.6:

- In general we cannot let  $\gamma = 0$  (for an explicit counterexample, see Dem'yanov 1970). This is one of the subtleties to work with directional derivatives: even when  $Dh(\mathbf{x}^*; \mathbf{d})$  vanishes for some direction  $\mathbf{d}$  we may still have  $Dh(\mathbf{x}^*; \mathbf{d})$  approaching 0 for other directions, but with  $\gamma = 0$  we will not know how  $A_h(\mathbf{x}^*; \mathbf{d})$  behaves (e.g. negative) along the latter directions.
- It is clear that  $H_+h \geq A_h$ . In some cases it is easier to verify the uniformity (along directions) in (A.31) if we relax the lower 2nd-order directional derivative  $H_+h$  to some convenient function  $A_h$ . See Theorem A.11 for an example.
- If  $\mathcal{X} = \mathbb{R}^m$  and  $h$  is Fréchet differentiable with locally Lipschitz gradient  $\nabla h$  around  $\mathbf{x}^*$ , then we can verify the uniformity in (A.31) as follows. Note first that we have  $\nabla h(\mathbf{x}^*) = \mathbf{0}$  from the necessary condition. Second, for all small  $t$  we have

$$\frac{h(\mathbf{x}^* + t\bar{\mathbf{d}}) - h(\mathbf{x}^*)}{t^2/2} = \frac{h(\mathbf{x}^* + t\mathbf{d} + t(\bar{\mathbf{d}} - \mathbf{d})) - h(\mathbf{x}^*)}{t^2/2} \quad (\text{A.37})$$

$$= \frac{h(\mathbf{x}^* + t\mathbf{d}) - h(\mathbf{x}^*) + t \langle \nabla h(\mathbf{x}^* + \theta t\mathbf{d}) - \nabla h(\mathbf{x}^*), \bar{\mathbf{d}} - \mathbf{d} \rangle}{t^2/2} \quad (\text{A.38})$$

$$\geq \frac{h(\mathbf{x}^* + t\mathbf{d}) - h(\mathbf{x}^*)}{t^2/2} - 2L\|\mathbf{d}\|\|\bar{\mathbf{d}} - \mathbf{d}\|, \quad (\text{A.39})$$

where  $\theta \in [0, 1]$  and  $L$  is the local Lipschitz constant of  $\nabla h$ . Thus, if  $\frac{h(\mathbf{x}^* + t\mathbf{d}) - h(\mathbf{x}^*)}{t^2/2} > 0$  then for all nearby  $\bar{\mathbf{d}}$  we also have  $\frac{h(\mathbf{x}^* + t\bar{\mathbf{d}}) - h(\mathbf{x}^*)}{t^2/2} > 0$ . In this case we may let  $A_h = H_+h$  and recover (Ben-Tal and Zowe, 1985, Theorem 3.2).

Another result that directly uses the second-order derivative is:

**Theorem A.7 (second-order sufficient condition, e.g. Dem'yanov and Malozemov 1974)** *Suppose  $h$  is uniformly first-order and second-order directional differentiable (at  $\mathbf{x}^*$ ) and  $\mathcal{X}$  is convex. If there exist  $r, q > 0$  such that for all normalized feasible direction  $\mathbf{t}$ ,  $Dh(\mathbf{x}^*; \mathbf{t}) \geq 0$ , and*

$$0 \leq Dh(\mathbf{x}^*; \mathbf{t}) < r \implies Hh(\mathbf{x}^*; \mathbf{t}) \geq q > 0, \quad (\text{A.40})$$

*then  $\mathbf{x}^*$  is an isolated local minimum.*

**Proof** If  $Dh(\mathbf{x}^*; \mathbf{t}) \geq r$ , it reduces to the proof of Thm. A.5. Otherwise, (A.40) holds, and

$$h(\mathbf{x}^* + \alpha \mathbf{t}) = h(\mathbf{x}^*) + \alpha Dh(\mathbf{x}^*; \mathbf{t}) + \frac{\alpha^2}{2} Hh(\mathbf{x}^*; \mathbf{t}) + o(\alpha^2; \mathbf{t}). \quad (\text{A.41})$$

Since  $h$  is uniformly second-order directional differentiable in any direction  $\mathbf{t}$ , there exist  $0 < \alpha_1 < \alpha_0$  such that for any  $0 < \alpha < \alpha_1$  and for any  $\|\mathbf{t}\| = 1$ ,  $o(\alpha^2; \mathbf{t}) \geq -q\alpha^2/4$ . Therefore, for any  $\mathbf{x} \in \mathcal{N}(\mathbf{x}^*, \alpha_1)$  not equal to  $\mathbf{x}^*$ , we can take  $\mathbf{t} = (\mathbf{x} - \mathbf{x}^*)/\|\mathbf{x} - \mathbf{x}^*\|$  (which is feasible from convexity of  $\mathcal{X}$ ),  $\alpha = \|\mathbf{x} - \mathbf{x}^*\|$  and obtain:

$$h(\mathbf{x}) = h(\mathbf{x}^* + \alpha \mathbf{t}) \geq h(\mathbf{x}^*) + \alpha^2 q/4 > h(\mathbf{x}^*). \quad (\text{A.42})$$

■

In the theorem above, we are considering “approximately” critical directions, rather than only the second-order derivatives along the critical directions. The following example demonstrates this point, as inspired by Ben-Tal and Zowe (1985, Example 2.1):

**Example A.8** *We cannot take  $r = 0$  in (3.29). Consider  $f((x_1, x_2), y) = (2x_1 + x_1^2 + x_2^2)y + x_1^3$  and  $(\mathbf{x}^*, y^*) = (\mathbf{0}, 0)$ .  $\bar{f}_\epsilon(x_1, x_2) = \epsilon|2x_1 + x_1^2 + x_2^2| + x_1^3$  and it is uniformly twice directional differentiable. We can evaluate  $D\bar{f}_\epsilon((0, 0); (t_1, t_2)) = 2\epsilon|t_1|$  and*

$$H\bar{f}_\epsilon((0, 0); (t_1, t_2)) = \begin{cases} 2\epsilon(t_1^2 + t_2^2) & t_1 > 0, \\ 2\epsilon t_2^2 & t_1 = 0, \\ -2\epsilon(t_1^2 + t_2^2) & t_1 < 0. \end{cases}$$

The critical directions are  $(0, t_2)$  along which  $H\bar{f}_\epsilon(\mathbf{0}, \mathbf{t}) = 2\epsilon t_2^2 > 0$ . However,

$$\bar{f}_\epsilon((0, 0), (x_1, \sqrt{-2x_1 - x_1^2})) = x_1^3 < 0$$

if  $-2 \leq x_1 \leq 0$ .

### A.3 Envelope function

Our main interest in this work is the envelope function:

$$\bar{f}(\mathbf{x}) := \max_{\mathbf{y} \in \mathcal{Y}} f(\mathbf{x}, \mathbf{y}) \quad (\text{A.43})$$

where  $\mathcal{Y}$  is some compact topological Hausdorff space<sup>6</sup>. It is easy to verify:

- If  $f : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$  is (jointly) continuous, then so is  $\bar{f}$  (in  $\mathbf{x}$ ).
- If also  $\partial_{\mathbf{x}} f : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$  is (jointly) continuous, then  $\bar{f}$  is locally Lipschitz.

The envelope function turns out to be directionally differentiable:

6. Results in this section can be extended to the more general case where the constraint set  $\mathcal{Y}$  depends on  $\mathbf{x}$  (in some semicontinuous manner); see Seeger (1988) for an excellent treatment. For our purpose here it suffices to consider a constant  $\mathcal{Y}$ .

**Theorem A.9 (e.g. Danskin 1966; Dem'yanov 1966)** *Let  $f$  and  $\partial_{\mathbf{x}}f$  be (jointly) continuous. Then, the envelope function  $\bar{f}$  is directionally differentiable:*

$$D\bar{f}(\mathbf{x}; \mathbf{d}) = \max_{\mathbf{y} \in \mathcal{Y}_0(\mathbf{x})} \langle \partial_{\mathbf{x}}f(\mathbf{x}, \mathbf{y}), \mathbf{d} \rangle, \text{ where } \mathcal{Y}_0(\mathbf{x}) := \{\mathbf{y} \in \mathcal{Y} : \bar{f}(\mathbf{x}) = f(\mathbf{x}; \mathbf{y})\}. \quad (\text{A.44})$$

Clearly,  $D\bar{f}(\mathbf{x}; \cdot)$  is Lipschitz continuous.

The following theorem explains the necessity of the function  $A_h$  in Theorem A.6:

**Theorem A.10 (Seeger 1988; Dem'yanov 1970)** *Let  $f$  and  $\partial_{\mathbf{x}}f$  be continuous. Then,*

$$D\bar{f}(\mathbf{x}; \mathbf{d}) = \max_{\mathbf{y} \in \mathcal{Y}_0(\mathbf{x})} \langle \partial_{\mathbf{x}}f(\mathbf{x}, \mathbf{y}), \mathbf{d} \rangle, \quad \mathcal{Y}_0(\mathbf{x}) := \{\mathbf{y} \in \mathcal{Y} : \bar{f}(\mathbf{x}) = f(\mathbf{x}, \mathbf{y})\} \quad (\text{A.45})$$

$$H_+\bar{f}(\mathbf{x}; \mathbf{d}, \mathbf{g}) \geq \max_{\mathbf{y} \in \mathcal{Y}_1(\mathbf{x}; \mathbf{d})} H_+f(\mathbf{x}, \mathbf{y}; \mathbf{d}, \mathbf{g}), \quad \mathcal{Y}_1(\mathbf{x}; \mathbf{d}) := \{\mathbf{y} \in \mathcal{Y}_0(\mathbf{x}) : D\bar{f}(\mathbf{x}; \mathbf{d}) = \langle \partial_{\mathbf{x}}f(\mathbf{x}, \mathbf{y}), \mathbf{d} \rangle\}. \quad (\text{A.46})$$

If  $\partial_{\mathbf{xx}}^2 f$  is also (jointly) continuous, then

$$A_{\bar{f}}(\mathbf{x}; \mathbf{d}) := \max_{\mathbf{y} \in \mathcal{Y}_1(\mathbf{x}; \mathbf{d})} \langle \partial_{\mathbf{xx}}^2 f(\mathbf{x}, \mathbf{y}) \mathbf{d}, \mathbf{d} \rangle \quad (\text{A.47})$$

satisfies the uniformity condition in Theorem A.6.

**Proof** We need only prove the last claim. Indeed

$$\begin{aligned} \frac{\bar{f}(\mathbf{x} + t\mathbf{d}) - \bar{f}(\mathbf{x}) - tD\bar{f}(\mathbf{x}; \mathbf{d})}{t^2/2} &\geq \max_{\mathbf{y} \in \mathcal{Y}_1(\mathbf{x}; \mathbf{d})} \frac{f(\mathbf{x} + t\mathbf{d}, \mathbf{y}) - f(\mathbf{x}, \mathbf{y}) - t \langle \partial_{\mathbf{x}}f(\mathbf{x}, \mathbf{y}), \mathbf{d} \rangle}{t^2/2} \\ &= \max_{\mathbf{y} \in \mathcal{Y}_1(\mathbf{x}; \mathbf{d})} \langle \partial_{\mathbf{xx}}^2 f(\mathbf{x} + t\theta(\mathbf{y}, \mathbf{d}) \cdot \mathbf{d}, \mathbf{y}) \mathbf{d}, \mathbf{d} \rangle. \end{aligned} \quad (\text{A.48})$$

Since  $\partial_{\mathbf{xx}}^2 f$  is continuous (hence uniformly continuous over compact sets), the right-hand side converges to  $A_{\bar{f}}(\mathbf{x}; \mathbf{d})$  uniformly on bounded sets in  $\mathbf{d}$  as  $t$  goes to 0.  $\blacksquare$

When  $\mathcal{Y}$  has limit points, proving  $A_{\bar{f}}(\mathbf{x}; \mathbf{d}) = H\bar{f}(\mathbf{x}; \mathbf{d})$  may be difficult (even with additional regularity conditions). Nevertheless, we can still apply the sufficient condition in Theorem A.6.

Seeger (1988) pointed out the following equivalence:

$$D\bar{f}(\mathbf{x}; \mathbf{d}) = \max_{\mathbf{y} \in \mathcal{Y}_0(\mathbf{x})} Df(\mathbf{x}, \mathbf{y}; \mathbf{d}) = \max_{\mathbf{y} \in \mathcal{Y}_0(\mathbf{x})} \sup_{\mathbf{v} \in K_{\mathbf{d}}(\mathcal{Y}, \mathbf{y})} Df(\mathbf{x}, \mathbf{y}; (\mathbf{d}, \mathbf{v})), \quad (\text{A.49})$$

where the first two directional derivatives are taken wrt  $\mathbf{x}$  only while the last directional derivative is joint wrt  $(\mathbf{x}, \mathbf{y})$ . Indeed, when  $f$  is (jointly) continuously differentiable,  $Df(\mathbf{x}, \mathbf{y}; (\mathbf{d}, \mathbf{v})) = \langle \partial_{\mathbf{x}}f(\mathbf{x}, \mathbf{y}), \mathbf{d} \rangle + \langle \partial_{\mathbf{y}}f(\mathbf{x}, \mathbf{y}), \mathbf{v} \rangle$ . However, since  $\mathbf{y} \in \mathcal{Y}_0(\mathbf{x})$ , we know from the necessary condition in Theorem A.3 that  $\langle \partial_{\mathbf{y}}f(\mathbf{x}, \mathbf{y}), \mathbf{v} \rangle \leq 0$  for all  $\mathbf{v} \in K_{\mathbf{d}}(\mathcal{Y}, \mathbf{y})$ . Surprisingly, the second-order counterparts are no longer equivalent:

**Theorem A.11 (Seeger 1988)** *Let  $f : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$  be continuously differentiable. Then,*

$$H_+\bar{f}(\mathbf{x}; \mathbf{d}, \mathbf{g}) \geq \max_{\mathbf{y} \in \mathcal{Y}_0(\mathbf{x})} \sup_{\mathbf{v} \in \mathcal{V}(\mathbf{x}, \mathbf{y}; \mathbf{d})} \sup_{\mathbf{w} \in K_{\mathbf{d}}(\mathcal{Y}, \mathbf{y})} H_+f(\mathbf{x}, \mathbf{y}; (\mathbf{d}, \mathbf{v}), (\mathbf{g}, \mathbf{w})), \quad (\text{A.50})$$

where  $\mathcal{Y}_0(\mathbf{x}) = \{\mathbf{y} \in \mathcal{Y} : \bar{f}(\mathbf{x}) = f(\mathbf{x}, \mathbf{y})\}$  and  $\mathcal{V}(\mathbf{x}, \mathbf{y}; \mathbf{d}) := \{\mathbf{v} \in \mathcal{K}_d(\mathcal{Y}, \mathbf{y}) : D\bar{f}(\mathbf{x}; \mathbf{d}) = Df(\mathbf{x}, \mathbf{y}; (\mathbf{d}, \mathbf{v}))\}$ .

If the second-order derivative of  $f$  is also (jointly) continuous, then

$$\begin{aligned} A_{\bar{f}}(\mathbf{x}; \mathbf{d}) := & \max_{\mathbf{y} \in \mathcal{Y}_0(\mathbf{x})} \sup_{\mathbf{v} \in \mathcal{V}(\mathbf{x}, \mathbf{y}; \mathbf{d})} \sup_{\mathbf{w} \in \mathcal{K}_d(\mathcal{Y}, \mathbf{y}; \mathbf{v})} \left\langle \begin{bmatrix} \partial_{\mathbf{x}\mathbf{x}}^2 f(\mathbf{x}, \mathbf{y}) & \partial_{\mathbf{x}\mathbf{y}}^2 f(\mathbf{x}, \mathbf{y}) \\ \partial_{\mathbf{y}\mathbf{x}}^2 f(\mathbf{x}, \mathbf{y}) & \partial_{\mathbf{y}\mathbf{y}}^2 f(\mathbf{x}, \mathbf{y}) \end{bmatrix} \begin{pmatrix} \mathbf{d} \\ \mathbf{v} \end{pmatrix}, \begin{pmatrix} \mathbf{d} \\ \mathbf{v} \end{pmatrix} \right\rangle + \\ & + \langle \partial_{\mathbf{y}} f(\mathbf{x}, \mathbf{y}), \mathbf{w} \rangle \end{aligned} \quad (\text{A.51})$$

satisfies the uniformity condition in Theorem A.6, provided that the directions  $\mathbf{d}, \mathbf{v}$  and  $\mathbf{w}$  are bounded.

**Proof** We assume  $\mathcal{K}_d(\mathcal{Y}, \mathbf{y}; \mathbf{v})$  is not empty for otherwise the theorem is vacuous. For any  $\mathbf{w} \in \mathcal{K}_d(\mathcal{Y}, \mathbf{y}; \mathbf{v})$  we know for any sequence  $t_k \downarrow 0$  there exist a subsequence  $t_{k_i} \downarrow 0$  and  $\mathbf{w}_{k_i} \rightarrow \mathbf{w}$  such that  $\mathbf{y} + t_{k_i} \mathbf{v} + t_{k_i}^2 \mathbf{w}_{k_i} \in \mathcal{Y}$ . Thus, fix any  $\mathbf{y} \in \mathcal{Y}_0(\mathbf{x})$ ,  $\mathbf{v} \in \mathcal{V}(\mathbf{x}, \mathbf{y}; \mathbf{d})$  and  $\mathbf{w} \in \mathcal{K}_d(\mathcal{Y}, \mathbf{y}; \mathbf{v})$ , we know (after passing to a subsequence if necessary)

$$\frac{\bar{f}(\mathbf{x} + t_k \mathbf{d} + t_k^2 \mathbf{g}/2) - \bar{f}(\mathbf{x}) - t_k D\bar{f}(\mathbf{x}; \mathbf{d})}{t_k^2/2} \quad (\text{A.52})$$

$$\geq \frac{f(\mathbf{x} + t_k \mathbf{d} + t_k^2 \mathbf{g}/2, \mathbf{y} + t_k \mathbf{v} + t_k^2 \mathbf{w}_{k_i}/2) - f(\mathbf{x}, \mathbf{y}) - t_k Df(\mathbf{x}, \mathbf{y}; (\mathbf{d}, \mathbf{v}))}{t_k^2/2} \quad (\text{A.53})$$

$$\geq \frac{f(\mathbf{x} + t_k \mathbf{d} + t_k^2 \mathbf{g}/2, \mathbf{y} + t_k \mathbf{v} + t_k^2 \mathbf{w}/2) - f(\mathbf{x}, \mathbf{y}) - t_k Df(\mathbf{x}, \mathbf{y}; (\mathbf{d}, \mathbf{v}))}{t_k^2/2} + \quad (\text{A.54})$$

$$+ \frac{f(\mathbf{x} + t_k \mathbf{d} + t_k^2 \mathbf{g}/2, \mathbf{y} + t_k \mathbf{v} + t_k^2 \mathbf{w}_{k_i}/2) - f(\mathbf{x} + t_k \mathbf{d} + t_k^2 \mathbf{g}/2, \mathbf{y} + t_k \mathbf{v} + t_k^2 \mathbf{w}/2)}{t_k^2/2} \quad (\text{A.55})$$

$$= H_+ f(\mathbf{x}, \mathbf{y}; (\mathbf{d}, \mathbf{v}), (\mathbf{g}, \mathbf{w})) + o(t_k), \quad (\text{A.56})$$

where the small order term  $o(t_k)$  is independent of  $\mathbf{d}, \mathbf{v}$  and  $\mathbf{w}$  if they are bounded.  $\blacksquare$

By setting  $\mathbf{y} \in \mathcal{Y}_1(\mathbf{x}; \mathbf{d})$ ,  $\mathbf{v} = \mathbf{w} = \mathbf{0}$ , we see that the lower bounds in Theorem A.11 are always shaper than the ones in Theorem A.10. However, note that Theorem A.10 only requires  $\mathcal{Y}$  to be any compact topological space while Theorem A.11 only applies when  $\mathcal{Y}$  is a compact set of some finite dimensional vector space.

**Example A.12 (Seeger 1988)** Let  $\mathcal{Y} = \mathbb{R}^m$  and  $f(\mathbf{x}, \mathbf{y}) = \begin{pmatrix} \mathbf{x} \\ \mathbf{y} \end{pmatrix}^\top \left\{ \frac{1}{2} \begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{B}^\top & \mathbf{C} \end{bmatrix} \begin{pmatrix} \mathbf{x} \\ \mathbf{y} \end{pmatrix} + \begin{pmatrix} \mathbf{p} \\ \mathbf{q} \end{pmatrix} \right\}$ . Assume  $\mathbf{C} \prec \mathbf{0}$ . Then,  $\mathcal{Y}_0(\mathbf{x})$  is a singleton,  $\mathcal{Y}_1 = \mathbb{R}^m$ , and WLOG  $\mathbf{w} = \mathbf{0}$ . Therefore,

$$A_{\bar{f}}(\mathbf{x}; \mathbf{d}) = \mathbf{d}^\top (\mathbf{A} - \mathbf{B}\mathbf{C}^{-1}\mathbf{B}^\top) \mathbf{d}, \quad (\text{A.57})$$

whence  $(\mathbf{x}, \mathbf{y}) = \begin{bmatrix} \mathbf{A} & \mathbf{B} \\ \mathbf{B}^\top & \mathbf{C} \end{bmatrix}^{-1} \begin{pmatrix} \mathbf{p} \\ \mathbf{q} \end{pmatrix}$  is a (unique) global saddle point if  $\mathbf{C} \prec \mathbf{0}$  and  $\mathbf{A} - \mathbf{B}\mathbf{C}^{-1}\mathbf{B}^\top \succ \mathbf{0}$ .

However, if we apply Theorem A.10 we can only conclude that

$$A_{\bar{f}}(\mathbf{x}; \mathbf{d}) = \mathbf{d}^\top \mathbf{A} \mathbf{d}, \quad (\text{A.58})$$

which is clearly a looser lower bound (recall that  $\mathbf{C} \prec \mathbf{0}$ ).

In principle, one should use the lower second-order directional derivative)  $H_+(\mathbf{x}^*; \mathbf{d}, \mathbf{g}) \geq 0$  for a stronger necessary condition. However, to our knowledge, we do not have an appropriate formula for it. We therefore look into *upper* second-order derivatives instead for which Kawasaki (1988) showed a result. From this result, we are able to introduce the second-order necessary conditions for  $\mathbf{x}^*$  being a local minimizer of  $\bar{f}(\mathbf{x})$ :

**Theorem A.13 (Kawasaki 1988)** *Let  $f$  be twice (jointly) continuously differentiable. Then,*

$$H\bar{f}(\mathbf{x}; \mathbf{d}, \mathbf{g}) = \max_{\mathbf{y} \in \mathcal{Y}_1(\mathbf{x}, \mathbf{d})} \langle \partial_{\mathbf{x}} f(\mathbf{x}, \mathbf{y}), \mathbf{g} \rangle + \langle \partial_{\mathbf{xx}}^2 f(\mathbf{x}, \mathbf{y}) \mathbf{d}, \mathbf{d} \rangle + \limsup_{\mathbf{z} \rightarrow \mathbf{y}} \frac{1}{2} v_-^2(\mathbf{z}; \mathbf{d}) u^\dagger(\mathbf{z}), \quad (\text{A.59})$$

$$\text{where } (t)_- = \min\{t, 0\}, \quad t^\dagger = \begin{cases} 1/t, & t \neq 0 \\ 0, & t = 0 \end{cases}, \text{ and}$$

$$u(\mathbf{y}) := \bar{f}(\mathbf{x}) - f(\mathbf{x}, \mathbf{y}) \geq 0, \quad v(\mathbf{y}; \mathbf{d}) := D\bar{f}(\mathbf{x}; \mathbf{d}) - Df(\mathbf{x}, \mathbf{y}; \mathbf{d}). \quad (\text{A.60})$$

**Proof** We give a direct (and arguably simpler) proof of this result. Denote

$$\Delta(t) := \frac{\bar{f}(\mathbf{x} + t\mathbf{d} + t^2\mathbf{g}/2) - \bar{f}(\mathbf{x}) - tD\bar{f}(\mathbf{x}; \mathbf{d})}{t^2/2}. \quad (\text{A.61})$$

Using the definitions of  $u$  and  $v$  we have

$$\Delta(t) = \frac{\bar{f}(\mathbf{x} + t\mathbf{d} + t^2\mathbf{g}) - f(\mathbf{x}, \mathbf{z}) - tDf(\mathbf{x}, \mathbf{z}; \mathbf{d}) - u(\mathbf{z}) - tv(\mathbf{z}; \mathbf{d})}{t^2/2}, \quad (\text{A.62})$$

which holds for any  $\mathbf{z} \in \mathcal{Y}$ . Let us first choose  $\mathbf{z} = \mathbf{z}_t \in \mathcal{Y}_0(\mathbf{x} + t\mathbf{d} + t^2\mathbf{g})$ :

$$\Delta(t) = \frac{f(\mathbf{x} + t\mathbf{d} + t^2\mathbf{g}, \mathbf{z}_t) - f(\mathbf{x}, \mathbf{z}_t) - tDf(\mathbf{x}, \mathbf{z}_t; \mathbf{d})}{t^2/2} - \frac{u(\mathbf{z}_t) + tv(\mathbf{z}_t; \mathbf{d})}{t^2/2}. \quad (\text{A.63})$$

Let  $\mathbf{y} \in \mathcal{Y}_0(\mathbf{x})$  be a limit point of  $\mathbf{z}_t$ . Suppose  $\mathbf{y} \in \mathcal{Y}_0(\mathbf{x}) \setminus \mathcal{Y}_1(\mathbf{x}; \mathbf{d})$ . Then, for small  $t$  we have (in the corresponding subsequence)  $v(\mathbf{z}_t; \mathbf{d}) \approx v(\mathbf{y}; \mathbf{d}) > 0$  hence  $\liminf_t \Delta(t) = H_+\bar{f}(\mathbf{x}; \mathbf{d}, \mathbf{g}) = -\infty$ , contradicting Theorem A.10. Thus,  $\mathbf{y} \in \mathcal{Y}_1(\mathbf{x}; \mathbf{d})$ . Optimizing  $t$  for the second term we obtain

$$\Delta(t) \leq \frac{f(\mathbf{x} + t\mathbf{d} + t^2\mathbf{g}, \mathbf{z}_t) - f(\mathbf{x}, \mathbf{z}_t) - tDf(\mathbf{x}, \mathbf{z}_t; \mathbf{d})}{t^2/2} + \frac{1}{2} v_-^2(\mathbf{z}_t; \mathbf{d}) u^\dagger(\mathbf{z}_t), \quad (\text{A.64})$$

where we used the fact that if  $u(\mathbf{z}_t) = 0$  then  $v(\mathbf{z}_t; \mathbf{d}) \geq 0$  (see Theorem A.9). Taking limits on both sides proves the  $\leq$  part in (A.59).

For the converse, let  $\mathbf{y} \in \mathcal{Y}_1(\mathbf{x}; \mathbf{d})$  and  $\mathbf{z}_k \rightarrow \mathbf{y}$  attain the maximum and limsup in (A.59), respectively. We need only consider  $\lim_{\mathbf{z}_k \rightarrow \mathbf{y}} \frac{1}{2} v_-^2(\mathbf{z}_k; \mathbf{d}) u^\dagger(\mathbf{z}_k) > 0$ , for otherwise the  $\geq$  part in (A.59) would already follow from Theorem A.10. We obviously have  $u(\mathbf{z}_k) > 0$  and  $v(\mathbf{z}_k; \mathbf{d}) < 0$  for sufficiently large  $t$ . Since  $u(\mathbf{z}_k) \rightarrow u(\mathbf{y}) = 0$  we also have  $v(\mathbf{z}_k; \mathbf{d}) \rightarrow v(\mathbf{y}; \mathbf{d}) = 0$ . We claim that (after passing to a subsequence if necessary)  $\lim_k u(\mathbf{z}_k)/v(\mathbf{z}_k; \mathbf{d}) = 0$ ,

for otherwise  $\lim v^2(\mathbf{z}_k; \mathbf{d})/u(\mathbf{z}_k) = 0$ , contradicting to its strict positivity. Now, setting  $t_k = -2u(\mathbf{z}_k)/v(\mathbf{z}_k; \mathbf{d})$  we have (for large  $k$ ):

$$\Delta(t_k) \geq \frac{f(\mathbf{x} + t_k \mathbf{d} + t_k^2 \mathbf{g}, \mathbf{z}_k) - f(\mathbf{x}, \mathbf{z}_k) - t_k \mathbf{D}f(\mathbf{x}, \mathbf{z}_k; \mathbf{d}) - u(\mathbf{z}_k) - t_k v(\mathbf{z}_k; \mathbf{d})}{t_k^2/2} \quad (\text{A.65})$$

$$= \frac{f(\mathbf{x} + t_k \mathbf{d} + t_k^2 \mathbf{g}, \mathbf{z}_k) - f(\mathbf{x}, \mathbf{z}_k) - t_k \mathbf{D}f(\mathbf{x}, \mathbf{z}_k; \mathbf{d})}{t_k^2/2} + \frac{1}{2} v_-^2(\mathbf{z}_k; \mathbf{d}) u^\dagger(\mathbf{z}_k). \quad (\text{A.66})$$

Taking limits on both sides we obtain the  $\geq$  part in (A.59).  $\blacksquare$

For later convenience, we remind that

$$\mathcal{Y}_0(\mathbf{x}) = \{\mathbf{y} : u(\mathbf{y}) = 0\}, \quad \mathcal{Y}_1(\mathbf{x}; \mathbf{d}) = \{\mathbf{y} : u(\mathbf{y}) = v(\mathbf{y}; \mathbf{d}) = 0\}. \quad (\text{A.67})$$

and denote  $\bar{E}(\mathbf{y}; \mathbf{t}) = \limsup_{\mathbf{z} \rightarrow \mathbf{y}} \frac{1}{2} v_-^2(\mathbf{z}; \mathbf{d}) u^\dagger(\mathbf{z})$ .

With Carathéodory's theorem for convex hulls, one can obtain from (A.59) the following necessary condition for envelope functions:

**Theorem A.14 (Kawasaki (1991))** *Assume  $f \in \mathcal{C}^2$  and  $\mathcal{X} = \mathbb{R}^n$ . If  $\mathbf{x}^*$  is a local minimum of  $\bar{f}(\mathbf{x})$ , then for each  $\mathbf{d} \in \mathbb{R}^n$  satisfying  $\mathbf{D}\bar{f}(\mathbf{x}^*; \mathbf{d}) = 0$ , there exist at most  $n + 1$  points  $\mathbf{y}_1, \dots, \mathbf{y}_{n+1} \in \mathcal{Y}_1(\mathbf{x}^*; \mathbf{d})$  and  $\lambda_1, \dots, \lambda_n \geq 0$  not all zero, such that:*

$$\sum_{i=1}^a \lambda_i \partial_{\mathbf{x}} f(\mathbf{x}^*, \mathbf{y}_i) = \mathbf{0}, \quad \sum_{i=1}^a \lambda_i \left( \mathbf{d}^\top \partial_{\mathbf{xx}}^2 f(\mathbf{x}^*, \mathbf{y}_i) \mathbf{d} + \bar{E}(\mathbf{y}_i; \mathbf{d}) \right) \geq 0. \quad (\text{A.68})$$

**Proof** We borrow the result from Kawasaki (1991). In order to write down the second-order derivative formula in Kawasaki (1988), we define

$$Y_0(\mathbf{t}) := \{\mathbf{y} \in \mathcal{Y} : \text{there exists a sequence } \{\mathbf{z}_k\} \rightarrow \mathbf{y}, u(\mathbf{z}_k) > 0 \text{ and } v(\mathbf{z}_k; \mathbf{t})/u(\mathbf{z}_k) \rightarrow -\infty\},$$

and the following upper semi-continuous function (Kawasaki, 1988):

$$\bar{E}'(\mathbf{y}; \mathbf{t}) = \begin{cases} \sup_{\{\mathbf{z}_k\} \rightarrow \mathbf{y}} \limsup_k v(\mathbf{z}_k; \mathbf{t})^2 / (2u(\mathbf{z}_k)) & \mathbf{y} \in Y_0(\mathbf{t}) \text{ and } \{\mathbf{z}_k\} \text{ is in } Y_0(\mathbf{t}), \\ 0 & u(\mathbf{y}) = v(\mathbf{y}; \mathbf{t}) = 0 \text{ \& } \mathbf{y} \notin Y_0(\mathbf{t}) \\ -\infty & \text{otherwise.} \end{cases} \quad (\text{A.69})$$

As shown in Kawasaki (1991),  $u(\mathbf{y}) = v(\mathbf{y}; \mathbf{t}) = 0$  whenever  $\mathbf{y} \in Y_0(\mathbf{t})$ . We simplify the definition above:

**Lemma A.15** *Denoting  $x_- := \min\{x, 0\}$ ,  $x^\dagger = 1/x$  if  $x \neq 0$  and  $x^\dagger = 0$  otherwise, then for any  $u(\mathbf{y}) = v(\mathbf{y}; \mathbf{t}) = 0$ ,*

$$\bar{E}(\mathbf{y}; \mathbf{t}) = \limsup_{\mathbf{z}_k \rightarrow \mathbf{y}} v_-(\mathbf{z}_k; \mathbf{t})^2 u^\dagger(\mathbf{z}_k) / 2. \quad (\text{A.70})$$

**Proof** It suffices to consider those sequences  $\{\mathbf{z}_k\} \subset \mathcal{Y}$  such that  $u(\mathbf{z}_k) \geq 0$ . We want to prove that  $\bar{E}(\mathbf{y}; \mathbf{t}) = \bar{E}'(\mathbf{y}; \mathbf{t})$ . We first prove  $\bar{E}(\mathbf{y}; \mathbf{t}) \geq \bar{E}'(\mathbf{y}; \mathbf{t})$ . If  $\mathbf{y} \in Y_0(\mathbf{t})$ , then for any  $\delta > 0$ , there exists a sequence  $\{\mathbf{z}_k\}$  such that

$$\limsup_k v(\mathbf{z}_k; \mathbf{t})^2 / (2u(\mathbf{z}_k)) \geq \bar{E}'(\mathbf{y}; \mathbf{t}) - \delta,$$

$u(\mathbf{z}_k) > 0$  and  $v(\mathbf{z}_k; \mathbf{t})/u(\mathbf{z}_k) \rightarrow -\infty$ . For large enough  $m$ ,  $v(\mathbf{z}_k; \mathbf{t}) < 0$ , and thus we take the same sequence in (A.70) to obtain  $\bar{E}(\mathbf{y}; \mathbf{t}) \geq \bar{E}'(\mathbf{y}; \mathbf{t}) - \delta$ . Since the above holds for any  $\delta > 0$ , we have  $\bar{E}(\mathbf{y}; \mathbf{t}) \geq \bar{E}'(\mathbf{y}; \mathbf{t})$ . If  $\mathbf{y} \notin Y_0(\mathbf{t})$ , then  $\bar{E}(\mathbf{y}; \mathbf{t}) \geq 0 = \bar{E}'(\mathbf{y}; \mathbf{t})$ .

Now let us prove that  $\bar{E}(\mathbf{y}; \mathbf{t}) \leq \bar{E}'(\mathbf{y}; \mathbf{t})$ . Assume for any  $\delta > 0$ ,  $\{\mathbf{z}_k\}$  is the sequence such that

$$\limsup_k v_-(\mathbf{z}_k; \mathbf{t})^2 u^\dagger(\mathbf{z}_k)/2 \geq \bar{E}(\mathbf{y}; \mathbf{t}) - \delta.$$

If  $u(\mathbf{z}_k) > 0$  or  $v(\mathbf{z}_k; \mathbf{t}) < 0$  for finite number of  $m$ , then  $\bar{E}(\mathbf{y}; \mathbf{t}) = 0 \leq \bar{E}'(\mathbf{y}; \mathbf{t})$ . Assume WLOG now that for any  $m$ ,  $u(\mathbf{z}_k) > 0$  and  $v(\mathbf{z}_k; \mathbf{t}) < 0$ , if  $v(\mathbf{z}_k; \mathbf{t})/u(\mathbf{z}_k)$  is bounded, then since  $v(\mathbf{y}; \mathbf{t}) = 0$ ,  $\bar{E}(\mathbf{y}; \mathbf{t}) = 0 \leq \bar{E}'(\mathbf{y}; \mathbf{t})$ . So we can assume further that  $v(\mathbf{z}_k; \mathbf{t})/u(\mathbf{z}_k) \rightarrow -\infty$ . Using the same sequence in (A.69), we know  $\bar{E}'(\mathbf{y}; \mathbf{t}) \geq \bar{E}(\mathbf{y}; \mathbf{t}) - \delta$  for any  $\delta > 0$ , and thus  $\bar{E}'(\mathbf{y}; \mathbf{t}) \geq \bar{E}(\mathbf{y}; \mathbf{t})$ .  $\blacksquare$

Moreover, the following assumption guarantees the existence of  $\mathbf{H}\bar{f}(\mathbf{x}; \mathbf{d}, \mathbf{g})$  from which we can get second-order sufficient conditions:

**Assumption A.16 (Kawasaki (1992))** *For each  $\mathbf{y} \in \mathcal{Y}_1(\mathbf{x}^*; \mathbf{t})$  with  $\mathbf{t} \neq \mathbf{0}$  and  $\mathbf{D}\bar{f}(\mathbf{x}^*; \mathbf{t}) = 0$ , and for each non-zero  $\mathbf{d} \in \mathbb{R}^m$ , there exist  $\alpha, \beta \neq 0$  and  $p, q > 0$  such that the following approximation holds:*

$$u(\mathbf{y} + \delta \mathbf{d}) = \alpha \delta^p + o(\delta^p), \quad v(\mathbf{y} + \delta \mathbf{d}; \mathbf{t}) = \beta \delta^q + o(\delta^q), \quad (\text{A.71})$$

whenever  $\mathbf{y} + \delta \mathbf{d} \in \mathcal{N}(\mathbf{y}^*, \epsilon)$  and  $\delta > 0$ . Note that

$$u(\mathbf{y}) := \bar{f}(\mathbf{x}^*) - f(\mathbf{x}^*, \mathbf{y}), \quad v(\mathbf{y}; \mathbf{d}) := \mathbf{D}\bar{f}(\mathbf{x}^*; \mathbf{d}) - \mathbf{D}f(\mathbf{x}^*, \mathbf{y}; \mathbf{d}).$$

**Theorem A.17 (second-order sufficient condition, Kawasaki (1992))** *Assume Assumption A.16 holds at  $\mathbf{x}^*$ . Let  $\mathcal{X} = \mathbb{R}^n$  and  $\mathcal{Y}$  be convex.  $\mathbf{x}^*$  is an isolated local minimum of  $\bar{f}(\mathbf{x})$  if for any  $\mathbf{d} \in \mathbb{R}^n$ ,  $\mathbf{D}\bar{f}(\mathbf{x}^*; \mathbf{d}) > 0$ , or  $\mathbf{D}\bar{f}(\mathbf{x}^*; \mathbf{d}) = 0$ ,  $\mathbf{d} \neq \mathbf{0}$  and there exist  $a \geq 1$  points  $\mathbf{y}_1, \dots, \mathbf{y}_a \in \mathcal{Y}_1(\mathbf{x}^*; \mathbf{d})$  and  $\lambda_1, \dots, \lambda_a > 0$  such that:*

$$\sum_{i=1}^a \lambda_i \partial_{\mathbf{x}} f(\mathbf{x}^*, \mathbf{y}_i) = \mathbf{0}, \quad \sum_{i=1}^a \lambda_i \left( \mathbf{d}^\top \partial_{\mathbf{x}\mathbf{x}}^2 f(\mathbf{x}^*, \mathbf{y}_i) \mathbf{d} + \bar{E}(\mathbf{y}_i; \mathbf{d}) \right) > 0. \quad (\text{A.72})$$

## Appendix B. Proofs in Section 3

**Theorem 3.4 (sufficient and necessary condition of local minimax when  $\partial_{\mathbf{y}\mathbf{y}}^2 f$  is invertible)** *Let  $\mathcal{X} = \mathbb{R}^n$ ,  $\mathcal{Y} = \mathbb{R}^m$  and  $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$  be twice continuously differentiable. Suppose  $\partial_{\mathbf{y}\mathbf{y}}^2 f(\mathbf{x}^*, \mathbf{y}^*)$  is invertible (i.e. non-degenerate), then  $(\mathbf{x}^*, \mathbf{y}^*)$  is local minimax iff*

- $\partial_{\mathbf{y}} f(\mathbf{x}^*, \mathbf{y}^*) = \mathbf{0}$ ,  $\partial_{\mathbf{y}\mathbf{y}}^2 f(\mathbf{x}^*, \mathbf{y}^*) \prec \mathbf{0}$ , and
- $\mathbf{x}^*$  is a local minimizer of the total function  $f(\mathbf{x}, \mathbf{y}(\mathbf{x}))$  where  $\mathbf{y}$  is defined implicitly near  $\mathbf{x}^*$  through the non-linear equation

$$\partial_{\mathbf{y}} f(\mathbf{x}, \mathbf{y}) = \mathbf{0}. \quad (3.3)$$

**Proof** Given that  $\partial_{\mathbf{y}\mathbf{y}}^2 f(\mathbf{x}^*, \mathbf{y}^*)$  is invertible, the first condition is clearly equivalent to  $\mathbf{y}^*$  being a local maximizer of  $f(\mathbf{x}^*, \cdot)$ . Consider the non-linear equation (3.3), whose solution is determined by the implicit function theorem as a continuously differentiable function  $\mathbf{y}(\mathbf{x})$  defined near  $\mathbf{x}^*$ . Fix any  $\epsilon$ . Since  $\mathbf{y}(\mathbf{x}^*) = \mathbf{y}^*$ , shrinking the neighbourhood around  $\mathbf{x}^*$  if necessary we may assume  $\mathbf{y}(\mathbf{x}) \in \mathcal{N}(\mathbf{y}^*, \epsilon)$  so that  $\bar{f}_\epsilon(\mathbf{x}) = f(\mathbf{x}, \mathbf{y}(\mathbf{x}))$ . Thus, if  $(\mathbf{x}^*, \mathbf{y}^*)$  is local minimax, then for  $\mathbf{x}$  near  $\mathbf{x}^*$ :

$$f(\mathbf{x}^*, \mathbf{y}(\mathbf{x}^*)) = f(\mathbf{x}^*, \mathbf{y}^*) = \bar{f}_\epsilon(\mathbf{x}^*) \leq \bar{f}_\epsilon(\mathbf{x}) = f(\mathbf{x}, \mathbf{y}(\mathbf{x})), \quad (\text{B.1})$$

so,  $\mathbf{x}^*$  is a local minimizer of the total function. Reversing the argument proves the converse. ■

**Lemma 3.5** *Suppose  $\mathbf{y}^*$  maximizes  $f(\mathbf{x}^*, \mathbf{y})$  over some neighborhood  $\mathcal{N}(\mathbf{y}^*, \epsilon_0)$ . If  $\mathbf{x}^*$  is a local minimizer of  $\bar{f}_{\epsilon, \mathbf{y}^*}$ , for some  $0 \leq \epsilon \leq \epsilon_0$ , then it remains a local minimizer (even over the same local neighborhood) of  $\bar{f}_{\mathcal{N}}(\mathbf{x}) := \max_{\mathbf{y} \in \mathcal{N}} f(\mathbf{x}, \mathbf{y})$  for any  $\mathcal{N}(\mathbf{y}^*, \epsilon) \subseteq \mathcal{N} \subseteq \mathcal{N}(\mathbf{y}^*, \epsilon_0)$ .*

**Proof** We first note that since  $\mathbf{y}^*$  maximizes  $f(\mathbf{x}^*, \mathbf{y})$  over  $\mathcal{N}(\mathbf{y}^*, \epsilon_0)$ , we clearly have for all  $\mathbf{y}^* \in \mathcal{N} \subseteq \mathcal{N}(\mathbf{y}^*, \epsilon_0)$ :

$$\bar{f}_{\mathcal{N}}(\mathbf{x}^*) = f(\mathbf{x}^*, \mathbf{y}^*). \quad (\text{B.2})$$

Moreover, for any  $\mathcal{N} \supseteq \mathcal{N}(\mathbf{y}^*, \epsilon)$  and any  $\mathbf{x} \in \mathcal{X}$ :

$$\bar{f}_{\mathcal{N}}(\mathbf{x}) \geq \bar{f}_{\epsilon, \mathbf{y}^*}(\mathbf{x}) =: \bar{f}_\epsilon(\mathbf{x}). \quad (\text{B.3})$$

Since  $\mathbf{x}^*$  is a local minimizer of  $\bar{f}_\epsilon$ , say over the neighborhood  $\mathcal{M}$ , we have for all  $\mathbf{x} \in \mathcal{M}$  and  $\mathcal{N}(\mathbf{y}^*, \epsilon) \subseteq \mathcal{N} \subseteq \mathcal{N}(\mathbf{y}^*, \epsilon_0)$ :

$$\bar{f}_{\mathcal{N}}(\mathbf{x}) \geq \bar{f}_\epsilon(\mathbf{x}) \geq \bar{f}_\epsilon(\mathbf{x}^*) = f(\mathbf{x}^*, \mathbf{y}^*) = \bar{f}_{\mathcal{N}}(\mathbf{x}^*), \quad (\text{B.4})$$

i.e.,  $\mathbf{x}^*$  is a local minimizer of  $\bar{f}_{\mathcal{N}}(\mathbf{x})$  over the same local neighborhood  $\mathcal{M}$ . ■

**Proposition 3.6 (equivalent definition of local minimax)** *The pair  $(\mathbf{x}^*, \mathbf{y}^*) \in \mathcal{X} \times \mathcal{Y}$  is a local minimax point iff*

- *Fixing  $\mathbf{x}^*$ , then  $\mathbf{y}^*$  is a local maximizer of  $\underline{f}_{0, \mathbf{x}^*}(\mathbf{y}) = f(\mathbf{x}^*, \mathbf{y})$ ;*
- *Fixing  $\mathbf{y}^*$ , then  $\mathbf{x}^*$  is a local minimizer of  $\bar{f}_{\epsilon, \mathbf{y}^*}(\mathbf{x})$  for all  $\epsilon \in (0, \epsilon_0]$  with some  $\epsilon_0 > 0$ .*

**Proof** We need only prove if  $(\mathbf{x}^*, \mathbf{y}^*)$  is local minimax according to Definition 3.3, then there exists some  $\epsilon_0 > 0$  such that  $\mathbf{x}^*$  is a local minimizer of  $\bar{f}_\epsilon(\mathbf{x})$  for all  $\epsilon \in (0, \epsilon_0]$ . Indeed, from Definition 3.3 we know  $f(\mathbf{x}^*, \mathbf{y})$  is maximized at  $\mathbf{y}^*$  over some neighborhood  $\mathcal{N}(\mathbf{y}^*, \epsilon_0)$  for some  $\epsilon_0 > 0$ . For any  $0 < \epsilon \leq \epsilon_0$ , one can find  $0 < \epsilon_n < \epsilon$  since the promised sequence  $\epsilon_n \rightarrow 0$ . By definition  $\mathbf{x}^*$  is a local minimizer for  $\bar{f}_{\epsilon_n}$ , hence by Lemma 3.5 it remains a local minimizer for  $\bar{f}_\epsilon$ . ■



**Proposition 3.7 (local saddle and uniformly local minimax)** *Every local saddle point is uniformly local minimax. If for any  $\mathbf{x} \in \mathcal{X}$ ,  $f(\mathbf{x}, \cdot)$  is upper semi-continuous, then every uniformly local minimax point is local saddle.*

**Proof** Let  $(\mathbf{x}_*, \mathbf{y}_*)$  be local saddle, i.e.,  $\mathbf{y}_*$  maximizes  $f(\mathbf{x}_*, \cdot)$  over the neighborhood  $\mathcal{N}(\mathbf{y}_*, \epsilon)$  and  $\mathbf{x}_*$  minimizes  $\bar{f}_{0, \mathbf{y}_*} = f(\cdot, \mathbf{y}_*)$  over the neighborhood  $\mathcal{N}(\mathbf{x}_*, \epsilon)$ . We fix the neighborhood  $\mathcal{N}(\mathbf{x}_*) = \mathcal{N}(\mathbf{x}_*, \epsilon)$  and choose any sequence  $\{\epsilon_n\} \subset (0, \epsilon]$ . Applying Lemma 3.5 we know  $\mathbf{x}_*$  remains a minimum for all  $\bar{f}_{\epsilon_n}$  over the (fixed) neighborhood  $\mathcal{N}(\mathbf{x}_*)$ . Thus,  $(\mathbf{x}_*, \mathbf{y}_*)$  is uniformly local minimax.

Conversely, let  $f$  be upper semi-continuous (in  $\mathbf{y}$  for any  $\mathbf{x}$ ) and  $(\mathbf{x}^*, \mathbf{y}^*)$  uniformly local minimax over the fixed neighborhood  $\mathcal{N}(\mathbf{x}^*)$ . By definition  $\mathbf{y}^*$  maximizes  $f(\mathbf{x}^*, \cdot)$  over some neighborhood  $\mathcal{N}(\mathbf{y}^*, \epsilon_0)$ , and  $\mathbf{x}^*$  minimizes all  $\bar{f}_{\epsilon_n}$  over the fixed neighborhood  $\mathcal{N}(\mathbf{x}^*)$ , where the positive sequence  $\epsilon_n \rightarrow 0$ . Fix any  $\mathbf{x} \in \mathcal{N}(\mathbf{x}^*)$ . Since  $f(\mathbf{x}, \cdot)$  is upper semi-continuous at  $\mathbf{y}^*$ , we have for any  $\delta > 0$ , there exists  $\epsilon_n \in (0, \epsilon_0]$  such that:

$$f(\mathbf{x}^*, \mathbf{y}^*) = \bar{f}_{\epsilon_n}(\mathbf{x}^*) \leq \bar{f}_{\epsilon_n}(\mathbf{x}) \leq f(\mathbf{x}, \mathbf{y}^*) + \delta. \quad (\text{B.5})$$

Letting  $\delta \rightarrow 0$  we know  $f(\mathbf{x}, \mathbf{y}^*) \geq f(\mathbf{x}^*, \mathbf{y}^*)$  for any  $\mathbf{x} \in \mathcal{N}(\mathbf{x}^*)$ . ■

**Proposition 3.9 (equivalence with Jin et al. (2020))** *The pair  $(\mathbf{x}^*, \mathbf{y}^*)$  is local minimax w.r.t. function  $f$  iff there exists  $\delta_0 > 0$  and a non-negative function  $h$  satisfying  $h(\delta) \rightarrow 0$  as  $\delta \rightarrow 0$ , such that for any  $\delta \in (0, \delta_0]$  and any  $(\mathbf{x}, \mathbf{y}) \in \mathcal{N}(\mathbf{x}^*, \delta) \times \mathcal{N}(\mathbf{y}^*, \delta)$  we have*

$$f(\mathbf{x}^*, \mathbf{y}) \leq f(\mathbf{x}^*, \mathbf{y}^*) \leq \left[ \max_{\mathbf{y}' \in \mathcal{N}(\mathbf{y}^*, h(\delta))} f(\mathbf{x}, \mathbf{y}') \right] =: \bar{f}_{h(\delta)}(\mathbf{x}). \quad (3.5)$$

**Proof** ( $\Leftarrow$ ) Suppose  $(\mathbf{x}^*, \mathbf{y}^*)$  satisfies (3.5). Then clearly,  $\mathbf{y}^*$  maximizes  $f(\mathbf{x}^*, \cdot)$  over the neighborhood  $\mathcal{N}(\mathbf{x}^*, \delta_0)$ . Take an arbitrary positive sequence  $\{\delta_n\}$  with  $\delta_n \rightarrow 0$  and let  $\epsilon_n = \sup_{m \geq n} h(\delta_m)$ . Since  $h(\delta) \rightarrow 0$  as  $\delta \rightarrow 0$ , we may assume WLOG that  $\epsilon_n$  is well-defined and bounded from above. If  $h(\delta_n) = 0$  for some  $n$  then  $(\mathbf{x}^*, \mathbf{y}^*)$  is local saddle and hence local minimax thanks to Proposition 3.7. Otherwise we have  $\epsilon_n > 0$  for all  $n$  and  $\epsilon_n \rightarrow 0$  since  $\lim_{\delta \rightarrow 0} h(\delta) = 0$ . WLOG we assume  $\epsilon_1 \leq \delta_0$  (for otherwise we may discard the head of the sequence  $\{\epsilon_n\}$ ). From (3.5) we know for any  $\mathbf{x} \in \mathcal{N}(\mathbf{x}^*, \delta_n)$ :

$$\bar{f}_{h(\delta_n)}(\mathbf{x}) \geq f(\mathbf{x}^*, \mathbf{y}^*) = \bar{f}_{h(\delta_n)}(\mathbf{x}^*), \quad (\text{B.6})$$

since  $h(\delta_n) \leq \epsilon_1 \leq \delta_0$  and  $\mathbf{y}^*$  maximizes  $f(\mathbf{x}^*, \mathbf{y})$  over  $\mathcal{N}(\mathbf{x}^*, \delta_0)$ . Therefore,  $\mathbf{x}^*$  is a local minimizer of  $\bar{f}_{h(\delta_n)}$  hence also of  $\bar{f}_{\epsilon_n}$  thanks to Lemma 3.5.

( $\Rightarrow$ ) Suppose  $(\mathbf{x}^*, \mathbf{y}^*)$  is local minimax (see Definition 3.3). Then,  $\mathbf{y}^*$  maximizes  $f(\mathbf{x}^*, \cdot)$  over some neighborhood  $\mathcal{N}(\mathbf{y}^*, \epsilon_0)$  where  $\epsilon_0 > 0$ . Since  $\mathbf{x}^*$  is a local minimizer of  $\bar{f}_{\epsilon_n}$ , it minimizes  $\bar{f}_{\epsilon_n}$  over some neighborhood  $\mathcal{N}(\mathbf{x}^*, \delta'_n)$  with  $\delta'_n > 0$ . From  $\{\delta'_n\}$  we construct another positive sequence  $\{\delta_n\}$  where  $\delta_0 = \min\{\delta'_1, 1, \epsilon_0\} > 0$  and

$$\delta_n = \min\{\delta'_n, \delta_{n-1}, 1/n\}, \quad n = 1, 2, \dots, \quad (\text{B.7})$$

which is diminishing by construction. Define  $h(\delta) = \epsilon_n$  if  $\delta_{n+1} < \delta \leq \delta_n$ . Since  $\epsilon_n \rightarrow 0$ ,  $\lim_{\delta \rightarrow 0} h(\delta) = 0$ . WLOG we assume  $\epsilon_1 \leq \epsilon_0$  and by definition  $\delta_0 \leq \epsilon_0$ . For any  $\delta \in (0, \delta_0]$  there exists some  $n$  such that  $\delta \in (\delta_{n+1}, \delta_n]$ . Thus, for any  $(\mathbf{x}, \mathbf{y}) \in \mathcal{N}(\mathbf{x}^*, \delta'_n) \times \mathcal{N}(\mathbf{y}^*, \epsilon_0)$ :

$$\bar{f}_{h(\delta)}(\mathbf{x}) = \bar{f}_{\epsilon_n}(\mathbf{x}) \geq \bar{f}_{\epsilon_n}(\mathbf{x}^*) = f(\mathbf{x}^*, \mathbf{y}^*) \geq f(\mathbf{x}^*, \mathbf{y}). \quad (\text{B.8})$$

Since  $\delta \leq \delta_n \leq \delta'_n$  and  $\delta \leq \epsilon_0$ , the above still holds over the smaller neighborhood  $\mathcal{N}(\mathbf{x}^*, \delta) \times \mathcal{N}(\mathbf{y}^*, \delta)$ , which is exactly (3.5).  $\blacksquare$

**Theorem 3.10 (local and global minimax points in the convex-concave case)** *Let the function  $f(\mathbf{x}, \mathbf{y})$  be convex in  $\mathbf{x}$  and concave in  $\mathbf{y}$ . Then, an interior point  $(\mathbf{x}, \mathbf{y})$  is local minimax iff it is stationary, i.e.,  $\partial_{\mathbf{x}} f(\mathbf{x}, \mathbf{y}) = \mathbf{0}$  and  $\partial_{\mathbf{y}} f(\mathbf{x}, \mathbf{y}) = \mathbf{0}$  iff it is saddle. In particular, local minimax implies global minimax.*

**Proof** Suppose  $(\mathbf{x}^*, \mathbf{y}^*)$  is stationary. For any small  $\epsilon > 0$ ,

$$\bar{f}_{\epsilon}(\mathbf{x}) = \max_{\mathbf{y} \in \mathcal{N}(\mathbf{y}^*, \epsilon)} f(\mathbf{x}, \mathbf{y}) \quad (\text{B.9})$$

is convex by assumption. To see that  $\mathbf{x}^*$  is a local (hence global) minimizer of  $\bar{f}_{\epsilon}$ , we need only verify that  $\mathbf{0} \in \partial \bar{f}_{\epsilon}(\mathbf{x}^*)$ . Since  $\mathbf{y}^*$  maximizes  $f(\mathbf{x}^*, \cdot)$  by assumption, we know from Danskin's theorem that  $\partial \bar{f}_{\epsilon}(\mathbf{x}^*) \supseteq \partial f(\mathbf{x}^*, \mathbf{y}^*) \ni \mathbf{0}$  since  $(\mathbf{x}^*, \mathbf{y}^*)$  is stationary.

Now suppose  $(\mathbf{x}^*, \mathbf{y}^*)$  is local minimax. Then,  $\mathbf{y}^*$  is a local hence global maximizer of  $f(\mathbf{x}^*, \cdot)$ . Also,  $\mathbf{x}^*$  is a local hence global minimizer of  $\bar{f}_{\epsilon}$ . Thus,

$$\bar{f}(\mathbf{x}) \geq \bar{f}_{\epsilon}(\mathbf{x}) \geq \bar{f}_{\epsilon}(\mathbf{x}^*) = f(\mathbf{x}^*, \mathbf{y}^*) = \bar{f}(\mathbf{x}^*), \quad (\text{B.10})$$

i.e.,  $\mathbf{x}^*$  is a global minimizer of  $\bar{f}$ .  $\blacksquare$

**Corollary 3.13 (local optimal solutions in the convex-concave case)** *Let  $\mathcal{X}$  and  $\mathcal{Y}$  be convex and the function  $f(\mathbf{x}, \mathbf{y})$  be convex in  $\mathbf{x}$  and concave in  $\mathbf{y}$ . A point is local (global) saddle iff it is local minimax (maximin) iff it is an LRP.*

**Proof** For convex-concave functions being local saddle is equivalent to satisfying (3.8). We also know from Proposition 3.7 that every local saddle point is local minimax (maximin) and from Definition F.1 that every local minimax point is an LRP.  $\blacksquare$

**Lemma 3.16 (directional derivatives for different  $\bar{f}_{\epsilon}$ )** *Suppose  $f$  and  $\partial_{\mathbf{x}} f$  are jointly continuous and thus the directional derivative (3.7) exists. If  $\mathbf{y}^*$  is a local maximizer of  $f(\mathbf{x}^*, \cdot)$  over a neighborhood  $\mathcal{N}(\mathbf{y}^*, \epsilon_0)$ , then for any  $0 \leq \epsilon_1 \leq \epsilon_2 \leq \epsilon_0$ ,  $\mathcal{Y}_0(\mathbf{x}^*; \epsilon_1) \subseteq \mathcal{Y}_0(\mathbf{x}^*; \epsilon_2)$  and for each  $\mathbf{t} \in \mathbf{K}_d(\mathcal{X}, \mathbf{x}^*)$ ,  $\mathbf{D} \bar{f}_{\epsilon_2}(\mathbf{x}^*; \mathbf{t}) \geq \mathbf{D} \bar{f}_{\epsilon_1}(\mathbf{x}^*; \mathbf{t})$ .*

**Proof** Clearly,  $\bar{f}_{\epsilon}(\mathbf{x}^*) = f(\mathbf{x}^*, \mathbf{y}^*)$  for any  $\epsilon \in [0, \epsilon_0]$  and  $\mathbf{y} \in \mathcal{N}(\mathbf{y}^*, \epsilon_1)$  implies  $\mathbf{y} \in \mathcal{N}(\mathbf{y}^*, \epsilon_2)$  for any  $\epsilon_1 \leq \epsilon_2$ , whence follows  $\mathcal{Y}_0(\mathbf{x}^*; \epsilon_1) \subseteq \mathcal{Y}_0(\mathbf{x}^*; \epsilon_2)$ . Using Danskin's theorem in Theorem A.9 we thus have  $\mathbf{D} \bar{f}_{\epsilon_2}(\mathbf{x}^*; \mathbf{t}) \geq \mathbf{D} \bar{f}_{\epsilon_1}(\mathbf{x}^*; \mathbf{t})$ .  $\blacksquare$

**Theorem 3.17 (second-order necessary condition, local minimax)** *Suppose  $f, \partial_{\mathbf{x}}f$  and  $\partial_{\mathbf{xx}}^2 f$  are all (jointly) continuous. If  $(\mathbf{x}^*, \mathbf{y}^*)$  is a local minimax point, then for each direction  $\mathbf{t} \in \mathcal{K}_{\mathbf{d}}(\mathcal{X}, \mathbf{x}^*)$ , one of the following holds:*

1.  $D\bar{f}_{\epsilon}(\mathbf{x}^*; \mathbf{t}) > 0$  for all  $\epsilon > 0$  smaller than some  $\epsilon_0(\mathbf{t})$ ;
2.  $D\bar{f}_{\epsilon}(\mathbf{x}^*; \mathbf{t}) = 0$  for all  $\epsilon > 0$  smaller than some  $\epsilon_0(\mathbf{t})$  (i.e.  $\mathbf{t}$  is critical), in which case we further have

$$\mathbf{t}^{\top} \partial_{\mathbf{xx}}^2 f(\mathbf{x}^*, \mathbf{y}^*) \mathbf{t} + \frac{1}{2} \limsup_{\mathbf{z} \rightarrow \mathbf{y}^*} \left[ \max\{\partial_{\mathbf{x}}f(\mathbf{x}^*, \mathbf{z})^{\top} \mathbf{t}, 0\}^2 (f(\mathbf{x}^*, \mathbf{y}^*) - f(\mathbf{x}^*, \mathbf{z}))^{\dagger} \right] \geq 0, \quad (3.18)$$

where  $t^{\dagger} = 1/t$  if  $t \neq 0$  and 0 otherwise.

**Proof** We know  $\bar{f}_{\epsilon}$  is locally Lipschitz since  $\partial_{\mathbf{x}}f$  is continuous, and there exists  $\epsilon_0 > 0$  such that  $\bar{f}_{\epsilon}(\mathbf{x}^*) = f(\mathbf{x}^*, \mathbf{y}^*)$  for any  $0 < \epsilon < \epsilon_0$ . The rest of the claim can be readily derived from Theorem A.4 and Theorem A.13, by taking  $\epsilon \rightarrow 0$  and noting that the upper directional derivative is by definition larger than the lower directional derivative.  $\blacksquare$

**Theorem 3.22 (second-order sufficient condition, local minimax)** *Assume  $\mathcal{X} = \mathbb{R}^n$  and  $\mathcal{Y}$  is convex and  $f, \partial_{\mathbf{x}}f, \partial_{\mathbf{xx}}^2 f$  are (jointly) continuous. At a stationary point  $(\mathbf{x}^*, \mathbf{y}^*)$ , if there exists  $\epsilon_0 > 0$  such that:*

- $f(\mathbf{x}^*, \cdot)$  is maximized at  $\mathbf{y}^*$  on  $\mathcal{N}(\mathbf{y}^*, \epsilon_0)$ ;
- along each critical direction  $\mathbf{t} \neq \mathbf{0}$ :

$$\mathbf{t}^{\top} \partial_{\mathbf{xx}}^2 f(\mathbf{x}^*, \mathbf{y}^*) \mathbf{t} + \frac{1}{2} \limsup_{\mathbf{z} \rightarrow \mathbf{y}^*} \left( ((\partial_{\mathbf{x}}f(\mathbf{x}^*, \mathbf{z})^{\top} \mathbf{t})_+)^2 (f(\mathbf{x}^*, \mathbf{y}^*) - f(\mathbf{x}^*, \mathbf{z}))^{\dagger} \right) > 0, \quad (3.27)$$

and in any direction  $\mathbf{d} \in \mathbb{R}^m$ , there exist  $\alpha, \beta \neq 0$  and  $p, q > 0$  such that for every  $\mathbf{y} \in \mathcal{Y}_1(\mathbf{x}^*; \epsilon_0; \mathbf{t})$ , the following Taylor expansion holds:

$$f(\mathbf{x}^*, \mathbf{y} + \delta \mathbf{d}) = f(\mathbf{x}^*, \mathbf{y}) + \alpha \delta^p + o(\delta^p), \quad \partial_{\mathbf{x}}f(\mathbf{x}^*, \mathbf{y} + \delta \mathbf{d})^{\top} \mathbf{t} = \beta \delta^q + o(\delta^q), \quad (3.28)$$

then  $(\mathbf{x}^*, \mathbf{y}^*)$  is a local minimax point.

**Proof** It follows from Theorem A.17. From Danskin's theorem  $D\bar{f}_{\epsilon}(\mathbf{x}^*; \mathbf{t}) \geq 0$  for any small  $\epsilon > 0$ . Besides, for any small enough  $\epsilon$ , (A.72) is satisfied since  $\mathbf{y}^* \in \mathcal{Y}_1(\mathbf{x}^*; \epsilon_0; \mathbf{t})$ . Noting that  $\bar{f}_{\epsilon}(\mathbf{x}^*) = f(\mathbf{x}^*, \mathbf{y}^*) = f(\mathbf{x}^*, \mathbf{y})$  for any  $0 \leq \epsilon < \epsilon_0$  and  $\mathbf{y} \in \mathcal{Y}_1(\mathbf{x}^*; \epsilon_0; \mathbf{t})$ , (3.28) follows from Assumption A.16.  $\blacksquare$

**Theorem 3.23 (second-order sufficient condition, local minimax)** *Assume  $f \in \mathcal{C}^2$  and let  $\mathcal{X}$  be convex. Suppose  $\mathbf{y}^*$  is a local maximizer of  $f(\mathbf{x}^*, \cdot)$  and that  $(\mathbf{x}^*, \mathbf{y}^*)$  is an*

interior stationary point. If there is  $\epsilon_0 > 0$  and for any  $\epsilon \in (0, \epsilon_0]$ , there exist  $R, r > 0$  such that for any feasible direction  $\|\mathbf{t}\| = 1$  that satisfies  $0 \leq D\bar{f}_\epsilon(\mathbf{x}^*; \mathbf{t}) \leq r$ , we have

$$\begin{aligned} \max_{\mathbf{y} \in \mathcal{Y}_0(\mathbf{x}^*; \epsilon)} \max_{\substack{\mathbf{v} \in \mathcal{V}(\mathbf{x}^*, \mathbf{y}; \mathbf{t}) \\ \|\mathbf{v}\| \leq R}} \max_{\substack{\mathbf{w} \in \mathcal{K}_d(\Omega, \mathbf{y}; \mathbf{v}) \\ \|\mathbf{w}\| \leq R}} & \left\langle \begin{bmatrix} \partial_{\mathbf{x}\mathbf{x}}^2 f(\mathbf{x}^*, \mathbf{y}) & \partial_{\mathbf{x}\mathbf{y}}^2 f(\mathbf{x}^*, \mathbf{y}) \\ \partial_{\mathbf{y}\mathbf{x}}^2 f(\mathbf{x}^*, \mathbf{y}) & \partial_{\mathbf{y}\mathbf{y}}^2 f(\mathbf{x}^*, \mathbf{y}) \end{bmatrix} \begin{pmatrix} \mathbf{t} \\ \mathbf{v} \end{pmatrix}, \begin{pmatrix} \mathbf{t} \\ \mathbf{v} \end{pmatrix} \right\rangle + \\ & + \langle \partial_{\mathbf{y}} f(\mathbf{x}^*, \mathbf{y}), \mathbf{w} \rangle > \mathbf{0}, \end{aligned} \quad (3.29)$$

then this point is local minimax, where  $\mathcal{V}(\mathbf{x}, \mathbf{y}; \mathbf{t}) := \{\mathbf{v} \in \mathcal{K}_d(\Omega, \mathbf{y}) : D\bar{f}_\epsilon(\mathbf{x}; \mathbf{t}) = \partial_{\mathbf{x}} f(\mathbf{x}, \mathbf{y})^\top \mathbf{t} + \partial_{\mathbf{y}} f(\mathbf{x}, \mathbf{y})^\top \mathbf{v}\}$ ,  $\Omega := \mathcal{N}(\mathbf{y}^*, \epsilon)$  and

$$\begin{aligned} \mathcal{K}_d(\Omega, \mathbf{y}; \mathbf{v}) &:= \liminf_{t \rightarrow 0^+} \frac{\Omega - \mathbf{y} - t\mathbf{v}}{t^2/2} := \{\mathbf{g} : \forall \{t_k\} \downarrow 0 \exists \{t_{k_i}\} \downarrow 0, \{\mathbf{g}_{k_i}\} \rightarrow \mathbf{g}, \\ &\quad \mathbf{y} + t_{k_i} \mathbf{v} + t_{k_i}^2 \mathbf{g}_{k_i} / 2 \in \Omega\}. \end{aligned} \quad (3.30)$$

**Proof** Since  $\mathbf{y}^* \in \mathcal{Y}_0(\mathbf{x}^*; \epsilon)$ , from Danskin's theorem (Theorem A.9) we know that  $D\bar{f}_\epsilon(\mathbf{x}^*; \mathbf{t}) \geq 0$  for any  $\epsilon$  small enough. We then combine Theorem A.6 with Theorem A.11. Note that all the directions  $\mathbf{t}, \mathbf{v}, \mathbf{w}$  are bounded.  $\blacksquare$

## Appendix C. Proofs in Section 4

**Theorem 4.1 (sufficient and necessary conditions for optimality in quadratic games)** For (homogeneous) unconstrained quadratic games, a pair  $(\mathbf{x}, \mathbf{y})$  is

- stationary iff

$$\begin{bmatrix} \mathbf{A} & \mathbf{C} \\ \mathbf{C}^\top & \mathbf{B} \end{bmatrix} \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} = \mathbf{0}; \quad (4.2)$$

- global minimax iff  $\mathbf{B} \preceq \mathbf{0}$ ,  $\mathbf{P}_\mathbf{L}^\perp (\mathbf{A} - \mathbf{C}\mathbf{B}^\dagger \mathbf{C}^\top) \mathbf{P}_\mathbf{L}^\perp \succeq \mathbf{0}$  where  $\mathbf{L} = \mathbf{C}\mathbf{P}_\mathbf{B}^\perp$ , and

$$\begin{bmatrix} \mathbf{P}_\mathbf{L}^\perp & \\ & \mathbf{I} \end{bmatrix} \begin{bmatrix} \mathbf{A} & \mathbf{C} \\ \mathbf{C}^\top & \mathbf{B} \end{bmatrix} \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix} = \mathbf{0}; \quad (4.3)$$

(Recall that  $\mathbf{P}_\mathbf{L}^\perp = \mathbf{I} - \mathbf{L}\mathbf{L}^\dagger$  is the orthogonal projection onto the null space of  $\mathbf{L}^\top$ .)

- local minimax iff  $\mathbf{B} \preceq \mathbf{0}$ ,  $\mathbf{P}_\mathbf{L}^\perp (\mathbf{A} - \mathbf{C}\mathbf{B}^\dagger \mathbf{C}^\top) \mathbf{P}_\mathbf{L}^\perp \succeq \mathbf{0}$ , and stationary (i.e. (4.2) holds). In particular, local minimax points are always global minimax.

**Proof** The first claim follows directly from the definition of stationarity.

To prove the second claim, we note that fixing  $\mathbf{x}$ ,  $q(\mathbf{x}, \cdot)$  is clearly quadratic in  $\mathbf{y}$ . Thus, it admits a local (hence also global) maximizer  $\mathbf{y}$  iff

$$\mathbf{B} \preceq \mathbf{0}, \quad (\text{C.1})$$

$$\mathbf{C}^\top \mathbf{x} + \mathbf{B}\mathbf{y} = \mathbf{0}. \quad (\text{C.2})$$

Note that there exists some  $\mathbf{y}$  to satisfy (C.2) iff  $\mathbf{C}^\top \mathbf{x}$  belongs to the range space of  $\mathbf{B}$  iff

$$\mathbf{P}_\mathbf{B}^\perp \mathbf{C}^\top \mathbf{x} = \mathbf{0}, \text{ i.e. } \mathbf{L}^\top \mathbf{x} = \mathbf{0}, \quad (\text{C.3})$$

or equivalently  $\mathbf{x} = \mathbf{P}_\mathbf{L}^\perp \mathbf{z}$  for some  $\mathbf{z} \in \mathbb{R}^m$ . Therefore, we have the envelope function:

$$\bar{q}(\mathbf{x}) = \begin{cases} \frac{1}{2} \mathbf{x}^\top (\mathbf{A} - \mathbf{C}\mathbf{B}^\dagger \mathbf{C}^\top) \mathbf{x}, & \mathbf{L}^\top \mathbf{x} = \mathbf{0} \\ \infty, & \text{otherwise} \end{cases}. \quad (\text{C.4})$$

Thus, the quadratic function  $\bar{q}$  (when restricted to the null space of  $\mathbf{L}^\top$ ) admits a local (hence also global) minimizer iff

$$\mathbf{P}_\mathbf{L}^\perp (\mathbf{A} - \mathbf{C}\mathbf{B}^\dagger \mathbf{C}^\top) \mathbf{P}_\mathbf{L}^\perp \succeq \mathbf{0}, \quad (\text{C.5})$$

in which case the minimizer  $\mathbf{x}$  satisfies

$$\mathbf{L}^\top \mathbf{x} = \mathbf{0} = \mathbf{P}_\mathbf{L}^\perp (\mathbf{A} - \mathbf{C}\mathbf{B}^\dagger \mathbf{C}^\top) \mathbf{x}, \quad (\text{C.6})$$

whereas the maximizer  $\mathbf{y}$  satisfies (C.2). It is easy to verify that (C.6) and (C.2) are equivalent to (4.3). For the last claim, note first that we have proved in Theorem 3.12 that any local minimax point is stationary. Moreover, if  $(\mathbf{x}^*, \mathbf{y}^*)$  is local minimax, then  $\mathbf{x}^*$  locally minimizes  $\bar{q}_{\epsilon, \mathbf{y}^*}$  (for all small  $\epsilon$ ), i.e., for  $\mathbf{x}$  close to  $\mathbf{x}^*$ , we have

$$\bar{q}(\mathbf{x}) \geq \bar{q}_{\epsilon, \mathbf{y}^*}(\mathbf{x}) \geq \bar{q}_{\epsilon, \mathbf{y}^*}(\mathbf{x}^*) = q(\mathbf{x}^*, \mathbf{y}^*) = \bar{q}(\mathbf{x}^*), \quad (\text{C.7})$$

where the last equality follows since fixing  $\mathbf{x}^*, \mathbf{y}^*$  is a local hence also global maximizer of the quadratic function  $q(\mathbf{x}^*, \cdot)$ . We have shown above that any local minimizer of  $\bar{q}(\mathbf{x})$  is necessarily global. Therefore,  $(\mathbf{x}^*, \mathbf{y}^*)$  is global minimax.

Lastly, we prove the converse of the last claim. Let  $\mathbf{B} \preceq \mathbf{0}$ ,  $\mathbf{P}_\mathbf{L}^\perp (\mathbf{A} - \mathbf{C}\mathbf{B}^\dagger \mathbf{C}^\top) \mathbf{P}_\mathbf{L}^\perp \succeq \mathbf{0}$ , and  $(\mathbf{x}^*, \mathbf{y}^*)$  be stationary, i.e. they satisfy (4.2). Fixing  $\mathbf{y}^*$  we have for all small  $\epsilon > 0$ :

$$2\bar{q}_\epsilon(\mathbf{x}) = 2\bar{q}_{\epsilon, \mathbf{y}^*}(\mathbf{x}) = \max_{\|\mathbf{y} - \mathbf{y}^*\| \leq \epsilon} \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix}^\top \begin{bmatrix} \mathbf{A} & \mathbf{C} \\ \mathbf{C}^\top & \mathbf{B} \end{bmatrix} \begin{bmatrix} \mathbf{x} \\ \mathbf{y} \end{bmatrix}. \quad (\text{C.8})$$

We are left to prove  $\mathbf{x}^*$  is a local minimizer of  $\bar{q}_\epsilon$  for all small  $\epsilon$ .<sup>7</sup> Let  $c = \max\{\|\mathbf{B}^\dagger \mathbf{C}^\top\|, \|\mathbf{A} - \mathbf{C}\mathbf{B}^\dagger \mathbf{C}^\top\|\}$ . We assume first  $c > 0$  and  $\mathbf{L} \neq \mathbf{0}$ . Let  $\sigma$  be the smallest positive singular value of  $\mathbf{L} = \mathbf{C}\mathbf{P}_\mathbf{B}^\perp$ . Consider any  $\mathbf{x}$  such that  $\|\mathbf{x} - \mathbf{x}^*\| \leq \epsilon(\sigma \wedge 1)/(3c)$ . We decompose

$$\mathbf{x} - \mathbf{x}^* = \boldsymbol{\delta}_\parallel + \boldsymbol{\delta}_\perp, \text{ where } \boldsymbol{\delta}_\perp = \mathbf{P}_\mathbf{L}^\perp (\mathbf{x} - \mathbf{x}^*), \quad (\text{C.9})$$

and define

$$\mathbf{y} - \mathbf{y}^* = -\mathbf{B}^\dagger \mathbf{C}^\top (\mathbf{x} - \mathbf{x}^*) + \epsilon \mathbf{L}^\top (\mathbf{x} - \mathbf{x}^*) / (2\|\mathbf{L}^\top (\mathbf{x} - \mathbf{x}^*)\|), \quad (\text{C.10})$$

7. Unfortunately we cannot use the sufficient conditions in Section 3.2.4 since  $\mathbf{x}^*$  may not be an isolated local minimizer.

where by convention  $0/0 := 0$ . Clearly,  $\|\mathbf{y} - \mathbf{y}^*\| \leq \epsilon/3 + \epsilon/2 < \epsilon$ . Thus, using the stationarity of  $(\mathbf{x}^*, \mathbf{y}^*)$ :

$$2\bar{q}_\epsilon(\mathbf{x}) \geq 2q(\mathbf{x}, \mathbf{y}) = \begin{bmatrix} \mathbf{x} - \mathbf{x}^* \\ \mathbf{y} - \mathbf{y}^* \end{bmatrix}^\top \begin{bmatrix} \mathbf{A} & \mathbf{C} \\ \mathbf{C}^\top & \mathbf{B} \end{bmatrix} \begin{bmatrix} \mathbf{x} - \mathbf{x}^* \\ \mathbf{y} - \mathbf{y}^* \end{bmatrix} \quad (\text{C.11})$$

$$\begin{aligned} (\text{note } \mathbf{B}\mathbf{L}^\top = \mathbf{0}) &= (\mathbf{x} - \mathbf{x}^*)^\top (\mathbf{A} - \mathbf{C}\mathbf{B}^\dagger\mathbf{C}^\top)(\mathbf{x} - \mathbf{x}^*) + \epsilon\|\mathbf{L}^\top(\mathbf{x} - \mathbf{x}^*)\| \\ &= \delta_\parallel^\top (\mathbf{A} - \mathbf{C}\mathbf{B}^\dagger\mathbf{C}^\top)\delta_\parallel + 2\delta_\parallel^\top (\mathbf{A} - \mathbf{C}\mathbf{B}^\dagger\mathbf{C}^\top)\delta_\perp + \end{aligned} \quad (\text{C.12})$$

$$+ \delta_\perp^\top (\mathbf{A} - \mathbf{C}\mathbf{B}^\dagger\mathbf{C}^\top)\delta_\perp + \epsilon\|\mathbf{L}^\top\delta_\parallel\| \quad (\text{C.13})$$

$$\geq -\epsilon\sigma\|\delta_\parallel\|/3 - 2\epsilon\sigma\|\delta_\parallel\|/3 + 0 + \epsilon\sigma\|\delta_\parallel\| = 0 = 2\bar{q}_\epsilon(\mathbf{x}^*), \quad (\text{C.14})$$

where we used the fact that  $\|\delta_\parallel\| \vee \|\delta_\perp\| \leq \epsilon\sigma/(3c)$  and  $\mathbf{P}_\mathbf{L}^\perp(\mathbf{A} - \mathbf{C}\mathbf{B}^\dagger\mathbf{C}^\top)\mathbf{P}_\mathbf{L}^\perp \succeq \mathbf{0}$ . Finally, we note that if  $c = 0$ , then  $\mathbf{A} - \mathbf{C}\mathbf{B}^\dagger\mathbf{C}^\top = \mathbf{0}$  hence the proof still goes through (with  $c$  replaced by 1 say). Similarly, if  $\mathbf{L} = \mathbf{0}$ , then  $\delta_\parallel = \mathbf{0}$  hence the proof again goes through (with  $\sigma$  replaced by 1 say).  $\blacksquare$

#### Theorem 4.4 (equivalence between global and local minimax in quadratic games)

*An unconstrained quadratic game admits a global minimax point iff it admits a local minimax point iff*

$$\mathbf{B} \preceq \mathbf{0}, \quad \mathbf{P}_\mathbf{L}^\perp(\mathbf{A} - \mathbf{C}\mathbf{B}^\dagger\mathbf{C}^\top)\mathbf{P}_\mathbf{L}^\perp \succeq \mathbf{0}, \quad \text{and} \quad \begin{bmatrix} \mathbf{a} \\ \mathbf{b} \end{bmatrix} \in \mathcal{R} \left( \begin{bmatrix} \mathbf{A} & \mathbf{C} \\ \mathbf{C}^\top & \mathbf{B} \end{bmatrix} \right). \quad (4.4)$$

*For such quadratic games, local minimax points are exactly the same as stationary global minimax points.*

**Proof** If (4.4) holds, let

$$\begin{bmatrix} \mathbf{A} & \mathbf{C} \\ \mathbf{C}^\top & \mathbf{B} \end{bmatrix} \begin{bmatrix} \mathbf{x}^* \\ \mathbf{y}^* \end{bmatrix} = \begin{bmatrix} \mathbf{a} \\ \mathbf{b} \end{bmatrix}. \quad (\text{C.15})$$

Then, performing the translation  $(\mathbf{x}, \mathbf{y}) \leftarrow (\mathbf{x} - \mathbf{x}^*, \mathbf{y} - \mathbf{y}^*)$  we reduce to the homogeneous case and applying Theorem 4.1 we obtain the existence of a local (or global) minimax point. If a local minimax point exists, then stationarity yields the range condition. Performing translation and applying Theorem 4.1 again establishes all conditions in (4.4).

All we are left to prove is when a global minimax point  $(\mathbf{x}^*, \mathbf{y}^*)$  exists the range condition holds. Indeed, fixing  $\mathbf{x}^*, \mathbf{y}^*$  maximizes the quadratic  $q(\mathbf{x}^*, \cdot)$  hence from stationarity:

$$\mathbf{C}^\top \mathbf{x}^* + \mathbf{B}\mathbf{y}^* = \mathbf{b}. \quad (\text{C.16})$$

The above equation has a solution  $\mathbf{y}^*$  iff  $\mathbf{P}_\mathbf{B}^\perp \mathbf{C}^\top \mathbf{x}^* = \mathbf{P}_\mathbf{B}^\perp \mathbf{b}$ , i.e.  $\mathbf{L}^\top \mathbf{x}^* = \mathbf{P}_\mathbf{B}^\perp \mathbf{b}$  (recall that  $\mathbf{L} := \mathbf{C}\mathbf{P}_\mathbf{B}^\perp$ ). Solving  $\mathbf{y}$  and plugging back in  $q$  we obtain: for all  $\mathbf{x}$  such that  $\mathbf{L}^\top \mathbf{x} = \mathbf{P}_\mathbf{B}^\perp \mathbf{b}$ ,

$$\bar{q}(\mathbf{x}) = \frac{1}{2}\mathbf{x}^\top (\mathbf{A} - \mathbf{C}\mathbf{B}^\dagger\mathbf{C}^\top)\mathbf{x} + \mathbf{x}^\top \mathbf{C}\mathbf{B}^\dagger \mathbf{b} - \mathbf{a}^\top \mathbf{x}. \quad (\text{C.17})$$

Since  $\mathbf{x}^*$  is a global minimizer of  $\bar{q}$ , we obtain the stationarity condition:

$$\mathbf{P}_L^\perp[(\mathbf{A} - \mathbf{C}\mathbf{B}^\dagger\mathbf{C}^\top)\mathbf{x}^* + \mathbf{C}\mathbf{B}^\dagger\mathbf{b} - \mathbf{a}] = \mathbf{0}. \quad (\text{C.18})$$

Combined with (C.16) we obtain:

$$\mathbf{P}_L^\perp[\mathbf{A}\mathbf{x}^* + \mathbf{C}\mathbf{B}^\dagger\mathbf{B}\mathbf{y}^* - \mathbf{a}] = \mathbf{0} \iff \mathbf{A}\mathbf{x}^* + \mathbf{C}\mathbf{B}^\dagger\mathbf{B}\mathbf{y}^* - \mathbf{a} = \mathbf{L}\mathbf{z} = \mathbf{C}\mathbf{P}_B^\perp\mathbf{z} \text{ for some } \mathbf{z} \quad (\text{C.19})$$

$$\iff \mathbf{A}\mathbf{x}^* + \mathbf{C}(\mathbf{B}^\dagger\mathbf{B}\mathbf{y}^* + \mathbf{P}_B^\perp\mathbf{z}) = \mathbf{a} \quad (\text{C.20})$$

From (C.16) and (C.20) we deduce  $(\mathbf{x}^*, \mathbf{B}^\dagger\mathbf{B}\mathbf{y}^* + \mathbf{P}_B^\perp\mathbf{z})$  satisfies the range condition (C.15). ■

## Appendix D. Momentum algorithms

We study the effect of momentum for convergence to local saddle points, including heavy ball (Polyak, 1964) and Nesterov's momentum (Nesterov, 1983). They are similar to GDA and do not converge even for bilinear games, as proved in Zhang and Yu (2020). In the following two subsections, we study the effect of momentum for convergence to local saddle points. GDA is a special case if we take the momentum parameter  $\beta = 0$ .

Many of the proofs in this appendix and Appendix E rely on Schur's theorem:

**Theorem D.1 (Schur (1917))** *The roots of a real polynomial  $p(\lambda) = a_0\lambda^n + a_1\lambda^{n-1} + \dots + a_n$  are within the (open) unit disk of the complex plane iff  $\forall k \in \{1, 2, \dots, n\}$ ,  $\det(\mathbf{P}_k\mathbf{P}_k^\mathbf{H} - \mathbf{Q}_k^\mathbf{H}\mathbf{Q}_k) > 0$ , where  $\mathbf{P}_k, \mathbf{Q}_k$  are  $k \times k$  matrices defined as:  $[\mathbf{P}_k]_{i,j} = a_{i-j}\mathbf{1}_{i \geq j}$ ,  $[\mathbf{Q}_k]_{i,j} = a_{n-i+j}\mathbf{1}_{i \leq j}$ .*

In this theorem, we use  $A^\mathbf{H}$  to denote the Hermitian conjugate of  $A$ , and

$$\mathbf{1}_{\text{condition}} = \begin{cases} 1 & \text{if condition is true,} \\ 0 & \text{otherwise.} \end{cases} \quad (\text{D.1})$$

Schur's theorem has been applied to analyze bilinear zero-sum games to give necessary and sufficient convergence conditions (Zhang and Yu, 2020). However, in that paper only real polynomials have been studied. Here we give a corollary for complex quadratic polynomials:

**Lemma D.2 (Schur)** *For complex quadratic polynomials  $\lambda^2 + a\lambda + b$ , the exact convergence condition is:*

$$|b| < 1, (1 - |b|^2)^2 + 2\Re(a^2\bar{b}) > |a|^2(1 + |b|^2). \quad (\text{D.2})$$

**Proof** For quadratic polynomials, we compute

$$\mathbf{P}_1 = [1], \mathbf{Q}_1 = [b], \quad (\text{D.3})$$

$$\mathbf{P}_2 = \begin{bmatrix} 1 & 0 \\ a & 1 \end{bmatrix}, \mathbf{Q}_2 = \begin{bmatrix} b & a \\ 0 & b \end{bmatrix}, \quad (\text{D.4})$$

We require  $\det(\mathbf{P}_k \mathbf{P}_k^H - \mathbf{Q}_k^H \mathbf{Q}_k) =: \delta_k > 0$ , for  $k = 1, 2$ . If  $k = 1$ , we have  $1 - |b|^2 > 0$ . If  $k = 2$ , we have:

$$\mathbf{P}_k \mathbf{P}_k^H - \mathbf{Q}_k^H \mathbf{Q}_k = \begin{bmatrix} 1 - |b|^2 & \bar{a} - a\bar{b} \\ a - \bar{a}b & 1 - |b|^2 \end{bmatrix}, \quad (\text{D.5})$$

where  $\bar{a}$  means the complex conjugate. The determinant should be positive, so we have:

$$(1 - |b|^2)^2 + 2\Re(a^2\bar{b}) > |a|^2(1 + |b|^2). \quad (\text{D.6})$$

■

Some proofs in this section rely on Mathematica code, mostly with the built-in function `Reduce`. This function relies on cylindrical algebraic decomposition (Basu et al., 2005) and can be verified manually.

### D.1 Heavy ball (HB)

We study the heavy ball method  $\text{HB}(\alpha_1, \alpha_2, \beta)$  (Polyak, 1964) in the context of minimax optimization, as also studied in Gidel et al. (2019); Zhang and Yu (2020):

$$\mathbf{z}_{t+1} = \mathbf{z}_t + \mathbf{v}(\mathbf{z}_t) + \beta(\mathbf{z}_t - \mathbf{z}_{t-1}), \mathbf{v}(\mathbf{z}) = (-\alpha_1 \partial_{\mathbf{x}} f(\mathbf{z}), \alpha_2 \partial_{\mathbf{y}} f(\mathbf{z})). \quad (\text{D.7})$$

**Theorem D.3 (HB)**  $\text{HB}(\alpha_1, \alpha_2, \beta)$  is exponentially stable iff  $\forall \lambda \in \text{Sp}(\mathbf{H}_{\alpha_1, \alpha_2}), |\beta| < 1$ ,

$$2\beta\Re(\lambda^2) - 2(1 - \beta)^2(1 + \beta)\Re(\lambda) > (1 + \beta^2)|\lambda|^2.$$

**Proof** With state augmentation  $\mathbf{z}_t \rightarrow (\mathbf{z}_{t+1}, \mathbf{z}_t)$ , the Jacobian for  $\text{HB}(\alpha_1, \alpha_2, \beta)$  is:

$$\mathbf{J}_{\text{HB}}(f) = \begin{bmatrix} (1 + \beta)\mathbf{I}_{n+m} + \mathbf{H}_{\alpha_1, \alpha_2} & -\beta\mathbf{I}_{n+m} \\ \mathbf{I}_{n+m} & \mathbf{0} \end{bmatrix}, \quad (\text{D.8})$$

The spectrum can be computed as:

$$\text{Sp}(\mathbf{J}_{\text{HB}}(f)) = \{w : p(w) := (w - 1)(w - \beta) - w\lambda = 0, \lambda \in \mathbf{H}_{\alpha_1, \alpha_2}\}. \quad (\text{D.9})$$

This quadratic equation can be further expanded as:

$$w^2 - (\beta + 1 + \lambda)w + \beta = 0. \quad (\text{D.10})$$

With Lemma D.2, we obtain the necessary and sufficient conditions for which all the roots are within a unit disk:

$$|\beta| < 1, 2\beta\Re(\lambda^2) - 2(1 - \beta)^2(1 + \beta)\Re(\lambda) > (1 + \beta^2)|\lambda|^2. \quad (\text{D.11})$$

■



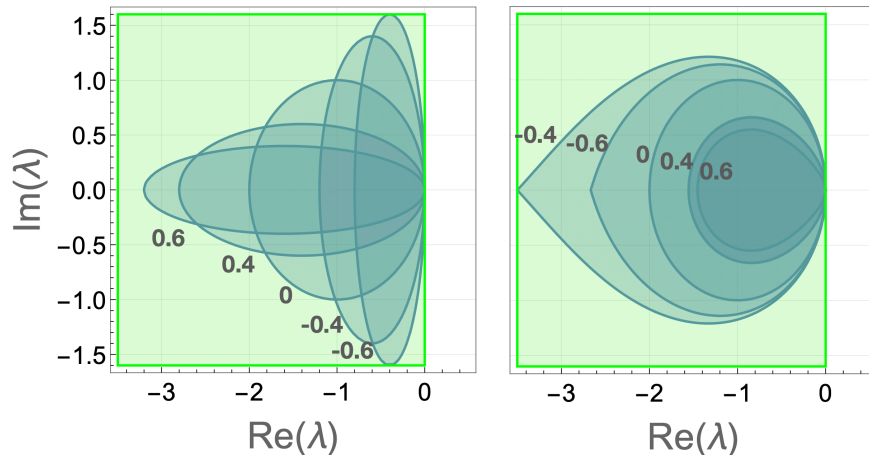


Figure 4 Convergence regions of momentum methods with different momentum parameter  $\beta$ : **(left)**  $\text{HB}(\alpha, \beta)$ ; **(right)**  $\text{NAG}(\alpha, \beta)$ . We take  $\beta = 0, \pm 0.4, \pm 0.6$  (as shown in the figure). The green region represents the one where the eigenvalues of  $\text{Sp}(\mathbf{H}_{\alpha_1, \alpha_2})$  at local saddle points may occur.

This theorem can also be derived from Euler transform as in (Niethammer and Varga, 1983, Section 6) which is used in analyzing methods for solving linear equations. The first inequality  $|\beta| < 1$  can be easily used to guide hyper-parameter tuning in practice. The second condition in fact describes an ellipsoid centered at  $(-\beta - 1, 0)$ . If we define  $\lambda = u + iv$  and  $(u, v) \in \mathbb{R}^2$ , then this condition can be simplified as:

$$\frac{(u + \beta + 1)^2}{(\beta + 1)^2} + \frac{v^2}{(\beta - 1)^2} < 1. \quad (\text{D.12})$$

As shown on the left of Figure 4, if the momentum factor  $\beta$  is positive, the ellipsoid is elongated in the horizontal direction; otherwise, it is elongated in the vertical direction. This agrees with existing results on negative momentum (Gidel et al., 2019; Zhang and Yu, 2020), where they studied bilinear games.

**Corollary D.4 (HB)** *For any  $|\beta| < 1$ ,  $\text{HB}(\alpha, \alpha, \beta)$  is exponentially stable for small enough  $\alpha$  at a local saddle point iff at such a point  $\Re(\lambda) \neq 0$  for all  $\lambda \in \text{Sp}(\mathbf{H})$ .*

**Proof** From Lemma 5.5, for any  $\lambda \in \text{Sp}(\mathbf{H})$ ,  $\Re(\lambda) \leq 0$ . If  $\Re(\lambda) \neq 0$  for all  $\lambda \in \text{Sp}(\mathbf{H})$ , then (D.12) holds for small enough  $\alpha$ . If  $\Re(\lambda) = 0$  for some  $\lambda \in \text{Sp}(\mathbf{H})$ , we cannot have (D.12). ■

## D.2 Nesterov’s accelerated gradient (NAG)

Nesterov’s accelerated gradient (Nesterov, 1983) is a variant of Polyak’s heavy ball, which achieves the optimal convergence rate for convex functions. It has been widely applied in deep learning (Sutskever et al., 2013). In Bollapragada et al. (2019), the authors analyzed the spectrum of NAG using numerical range in the context of linear regression, which is equivalent to the case when  $\text{Sp}(\mathbf{H}) \subset \mathbb{R}$  (cf. Bollapragada et al. (2019, p. 11)).

The key difference between HB and NAG is the order of momentum update and the gradient update. We study Nesterov's momentum for minimax optimization:

$$\mathbf{z}_{t+1} = \mathbf{z}'_t + \alpha \mathbf{v}(\mathbf{z}'_t), \quad \mathbf{z}'_t = \mathbf{z}_t + \beta(\mathbf{z}_t - \mathbf{z}_{t-1}), \quad (\text{D.13})$$

which we denote as  $\text{NAG}(\alpha_1, \alpha_2, \beta)$ . We have the following stability result for NAG:

**Theorem D.5 (NAG)** *NAG* $(\alpha_1, \alpha_2, \beta)$  *is exponentially stable iff for any*  $\lambda \in \text{Sp}(\mathbf{H}_{\alpha_1, \alpha_2})$ :

$$|1 + \lambda|^{-2} > 1 + 2\beta(\beta^2 - \beta - 1)\Re(\lambda) + \beta^2|\lambda|^2(1 + 2\beta), \quad |\beta| \cdot |1 + \lambda| < 1. \quad (\text{D.14})$$

**Proof** With state augmentation  $\mathbf{z}_t \rightarrow (\mathbf{z}_{t+1}, \mathbf{z}_t)$ , the Jacobian for NAG is:

$$\begin{bmatrix} (1 + \beta)(\mathbf{I}_{n+m} + \mathbf{H}_{\alpha_1, \alpha_2}) & -\beta(\mathbf{I}_{n+m} + \mathbf{H}_{\alpha_1, \alpha_2}) \\ \mathbf{I}_{n+m} & \mathbf{0} \end{bmatrix}.$$

The spectrum can be computed as:

$$\text{Sp}(\mathbf{J}(f)) = \{w : p(w) := w^2 - w(1 + \beta)(1 + \lambda) + \beta(1 + \lambda) = 0, \lambda \in \mathbf{H}_{\alpha_1, \alpha_2}\}.$$

Comparing with (D.10), we find that the two characteristic polynomials are different only by  $O(\alpha\beta)$ . With Lemma D.2, the condition for local linear convergence is:

$$|1 + \lambda|^{-2} > 1 + 2\beta(\beta^2 - \beta - 1)\Re(\lambda) + \beta^2|\lambda|^2(1 + 2\beta), \quad (\text{D.15})$$

$$|\beta| \cdot |1 + \lambda| < 1. \quad (\text{D.16})$$

■

From Figure 4, the convergence region of NAG is better conditioned than HB. However, NAG is still similar to HB and GDA in terms of the local convergence behavior:

**Corollary D.6 (NAG)** *If*  $\Re(\lambda) \geq 0$  *for some*  $\lambda \in \mathbf{H}_{\alpha_1, \alpha_2}$ , *then*  $\text{NAG}(\alpha_1, \alpha_2, \beta)$  *is not exponentially stable.*

**Proof** Take  $\lambda \in \mathbf{H}_{\alpha_1, \alpha_2}$  and assume  $\lambda = u + iv$  with  $u, v \in \mathbb{R}$ . (D.14) can be translated to the following Mathematica code:

```
Reduce[b^2 ((1 + u)^2 + v^2) < 1 && ((1 + u)^2 + v^2) (1 +
2 b (b^2 - b - 1) u + b^2 (u^2 + v^2) (1 + 2 b)) < 1 && u >= 0],
```

and the result is `False`. ■

According to Lemma 5.5,  $\text{NAG}(\alpha_1, \alpha_2, \beta)$  never converges on bilinear games. Summarizing the previous subsections, we conclude that adding momentum does not help in converging to local saddle points.

## Appendix E. Proofs in Section 5

**Lemma 5.1 (equivalence between past extra-gradient and OGD)** *The past extra-gradient method*

$$\mathbf{z}_{t+1} = \mathbf{z}_t + \mathbf{v}(\mathbf{z}_{t+1/2})/\beta, \mathbf{z}_{t+1/2} = \mathbf{z}_t + \mathbf{v}(\mathbf{z}_{t-1/2}) \quad (5.4)$$

can be rewritten as  $\mathbf{z}'_{t+1} = \mathbf{z}'_t + k\mathbf{v}(\mathbf{z}'_t) - \mathbf{v}(\mathbf{z}'_{t-1})$  with  $k = 1 + 1/\beta$  and  $\mathbf{z}'_t = \mathbf{z}_{t-1/2}$ .

**Proof** From the second equation of (5.4) we obtain

$$\begin{aligned} \mathbf{z}_{t+3/2} &= \mathbf{z}_{t+1} + \mathbf{v}(\mathbf{z}_{t+1/2}) \\ &= \mathbf{z}_t + \left(1 + \frac{1}{\beta}\right) \mathbf{v}(\mathbf{z}_{t+1/2}) + \mathbf{v}(\mathbf{z}_{t-1/2}) - \mathbf{v}(\mathbf{z}_{t-1/2}) \\ &= \mathbf{z}_{t+1/2} + \left(1 + \frac{1}{\beta}\right) \mathbf{v}(\mathbf{z}_{t+1/2}) - \mathbf{v}(\mathbf{z}_{t-1/2}). \end{aligned} \quad (E.1)$$

In the second line we used the first equation of (5.4) and in the third line we used the second equation of (5.4). ■

**Theorem 5.2 (stability of EG/OGD)** *At  $(\mathbf{x}^*, \mathbf{y}^*)$ ,  $EG(\alpha_1, \alpha_2, \beta)$  is exponentially stable iff for any  $\lambda \in Sp(\mathbf{H}_{\alpha_1, \alpha_2})$ ,  $|1 + \lambda/\beta + \lambda^2/\beta| < 1$ .  $OGD(k, \alpha_1, \alpha_2)$  is exponentially stable iff for any  $\lambda \in Sp(\mathbf{H}_{\alpha_1, \alpha_2})$ ,  $|\lambda| < 1$  and  $|\lambda|^2(k - 3 + (k + 1)|\lambda|^2) < 2\Re(\lambda)(k|\lambda|^2 - 1)$ .*

**Proof** From (5.2) the update of EG can be rewritten as  $\mathbf{z}_{t+1} = \mathbf{z}_t + \mathbf{v}(\mathbf{z}_t + \mathbf{v}(\mathbf{z}_t))/\beta$ . We compute the Jacobian matrix of this update:

$$\mathbf{J} = \mathbf{J}(f) = \mathbf{I} + \mathbf{H}_{\alpha_1, \alpha_2}/\beta + \mathbf{H}_{\alpha_1, \alpha_2}^2/\beta.$$

It then follows that  $Sp(\mathbf{J}) = 1 + Sp(\mathbf{H}_{\alpha_1, \alpha_2})/\beta + Sp(\mathbf{H}_{\alpha_1, \alpha_2})^2/\beta$ , where the operation is element-wise. Therefore,  $\rho(\mathbf{J}(f)) < 1$  iff

$$\max_{\lambda \in \mathbf{H}_{\alpha_1, \alpha_2}} |1 + \lambda/\beta + \lambda^2/\beta| < 1.$$

Similarly for OGD, the spectrum can be computed as:

$$Sp(\mathbf{J}_{OGD}) = \{x : p(x) := x^2 - (1 + k\lambda)x + \lambda = 0, \lambda \in \mathbf{H}_{\alpha_1, \alpha_2}\}. \quad (E.2)$$

With Lemma D.2, we obtain the necessary and sufficient conditions when the roots of  $p(x)$  are in the unit circle:

$$|\lambda| < 1, (k - 1)|\lambda|^2(k - 3 + (k + 1)|\lambda|^2) < 2(k - 1)\Re(\lambda)(k|\lambda|^2 - 1), \forall \lambda \in \mathbf{H}_{\alpha_1, \alpha_2}. \quad \blacksquare$$

**Theorem 5.3 (more aggressive extra-gradient steps, more stable)** For  $\beta_1 > \beta_2 > 1$ , whenever  $\text{EG}(\alpha_1, \alpha_2, \beta_2)$  is exponentially stable at  $(\mathbf{x}^*, \mathbf{y}^*)$ ,  $\text{EG}(\alpha_1, \alpha_2, \beta_1)$  is exponentially stable at  $(\mathbf{x}^*, \mathbf{y}^*)$  as well. For  $k_1 > k_2 > 1$ , whenever  $\text{OGD}(k_1, \alpha_1, \alpha_2)$  is exponentially stable at  $(\mathbf{x}^*, \mathbf{y}^*)$ ,  $\text{OGD}(k_2, \alpha_1, \alpha_2)$  is exponentially stable at  $(\mathbf{x}^*, \mathbf{y}^*)$  as well.

**Proof** Rewriting  $\lambda = x + iy$  with  $x, y \in \mathbb{R}$  for  $\lambda \in \mathbf{H}_{\alpha_1, \alpha_2}$  and using Theorem 5.2, we run the following Mathematica code ( $b_1 \equiv \beta_1$ ,  $b_2 \equiv \beta_2$ ):

```
Reduce[ForAll[{x, y, b1, b2}, ((y + 2 x y)/b2)^2 +
  (1 + (x + x^2 - y^2)/b2)^2 < 1 && b1 > b2 > 1,
  ((y + 2 x y)/b1)^2 + (1 + (x + x^2 - y^2)/b1)^2 < 1]]
```

The answer is True. For the second part, we rewrite the stability condition for OGD as:

$$k|\lambda|^2(1 + |\lambda|^2 - 2\Re(\lambda)) < 3|\lambda|^2 - |\lambda|^4 - 2\Re(\lambda). \quad (\text{E.3})$$

Since  $\Re(\lambda) \leq |\lambda|$ ,  $1 + |\lambda|^2 - 2\Re(\lambda) \geq 0$ . The left hand side increases with  $k$ .  $\blacksquare$

From Theorem D.3 and Theorem 5.2 we can easily infer the relation among the stable sets of gradient algorithms:

**Corollary E.1** Given  $|\lambda| < 1$  with  $\lambda \in \mathbf{H}_{\alpha_1, \alpha_2}$ , whenever  $\text{GDA}(\alpha_1, \alpha_2)$  converges,  $\text{EG}(\alpha_1, \alpha_2, 1)$  converges as well. Given  $|\lambda| < 1/\sqrt{3}$  with  $\lambda \in \mathbf{H}_{\alpha_1, \alpha_2}$ , whenever  $\text{GDA}(\alpha_1, \alpha_2)$  converges,  $\text{OGD}(2, \alpha_1, \alpha_2)$  converges.

**Proof** When  $\beta = 0$ , (D.11) becomes  $|1 + \lambda| < 1$ . The first part follows from:

$$|1 + \lambda| < 1 \text{ and } |\lambda| < 1 \implies |1 + \lambda + \lambda^2| < 1. \quad (\text{E.4})$$

Taking  $k = 2$ , from Theorem 5.2, the stability condition for OGD is:

$$|\lambda|^2(-1 + 3|\lambda|^2) < 2\Re(\lambda)(2|\lambda|^2 - 1). \quad (\text{E.5})$$

We want to show that for all  $|1 + \lambda| < 1$  and  $|\lambda| < 1/\sqrt{3}$ , (E.5) holds, and thus we define  $\lambda = u + iv$  ( $u, v \in \mathbb{R}$ ) and use the following Mathematica code:

```
Reduce[ForAll[{u, v}, (1 + u)^2 + v^2 < 1 && u^2 + v^2 < 1/3,
  (u^2 + v^2) (-1 + 3 (u^2 + v^2)) < 2 u (-1 + 2 (u^2 + v^2))]]
```

This result is True.  $\blacksquare$

**Lemma 5.5 (local saddle)** Suppose  $\alpha_1, \alpha_2 > 0$  are fixed. For  $f \in \mathcal{C}^2$ , at a local saddle point, for all  $\lambda \in \text{Sp}(\mathbf{H}_{\alpha_1, \alpha_2}(f))$ , we have  $\Re(\lambda) \leq 0$ . For all  $z \in \mathbb{C}$  with  $\Re(z) \leq 0$ , there exists a quadratic function  $q$  and a local saddle point  $(\mathbf{x}^*, \mathbf{y}^*)$  such that  $z \in \text{Sp}(\mathbf{H}_{\alpha_1, \alpha_2}(q))$ . For bilinear functions, at a local saddle point we have  $\Re(\lambda) = 0$  for all  $\lambda \in \text{Sp}(\mathbf{H}_{\alpha_1, \alpha_2})$ .

**Proof** The convergence analysis reduces to the spectral study of  $\mathbf{H}_{1,\gamma}$ . With the similarity transformation:

$$\mathbf{H}' = \mathbf{U}^{-1}\mathbf{H}_{1,\gamma}\mathbf{U} = \begin{bmatrix} -\partial_{\mathbf{x}\mathbf{x}}^2 f & -\sqrt{\gamma}\partial_{\mathbf{x}\mathbf{y}}^2 f \\ \sqrt{\gamma}\partial_{\mathbf{y}\mathbf{x}}^2 f & \gamma\partial_{\mathbf{y}\mathbf{y}}^2 f \end{bmatrix}, \quad \mathbf{U} = \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \mathbf{0} & \sqrt{\gamma}\mathbf{I} \end{bmatrix}, \quad (\text{E.6})$$

It suffices to study the spectrum of  $\mathbf{H}'$ . For any local saddle point  $(\mathbf{x}^*, \mathbf{y}^*)$ , we have:

$$\partial_{\mathbf{x}\mathbf{x}}^2 f(\mathbf{x}^*, \mathbf{y}^*) \succeq \mathbf{0}, \quad \partial_{\mathbf{y}\mathbf{y}}^2 f(\mathbf{x}^*, \mathbf{y}^*) \preceq \mathbf{0}. \quad (\text{E.7})$$

From this necessary condition,  $\Re(\mathbf{H}') := (\mathbf{H}' + \mathbf{H}'^\top)/2$  is negative semi-definite, and with the Ky Fan inequality (Fan (1950)) we have  $\Re(\text{Sp}(\mathbf{H}')) \prec \text{Sp}(\Re(\mathbf{H}')) \prec \mathbf{0}$ , with “ $\prec$ ” meaning majorization (Marshall et al., 1979). The second part can be proved by assuming  $z = -u + iv$  with  $u \geq 0$  and  $v \in \mathbb{R}$ . The quadratic function can be

$$q = \frac{ux^2}{2} - \frac{uy^2}{2\gamma} + \frac{v}{\sqrt{\gamma}}xy,$$

since one can verify that  $(0, 0)$  is a local saddle point where:

$$\mathbf{H}_{1,\gamma} = \begin{bmatrix} -u & -v/\sqrt{\gamma} \\ v\sqrt{\gamma} & -u \end{bmatrix}, \quad (\text{E.8})$$

whose two eigenvalues are  $z$  and  $\bar{z}$ . For bilinear games  $f = \mathbf{x}^\top \mathbf{C}\mathbf{y} + \mathbf{a}^\top \mathbf{x} + \mathbf{b}^\top \mathbf{y}$ , at any local saddle point, the Jacobian matrix of the vector field is:

$$\mathbf{H}_{1,\gamma} = \begin{bmatrix} \mathbf{0} & -\mathbf{C} \\ \gamma\mathbf{C}^\top & \mathbf{0} \end{bmatrix}. \quad (\text{E.9})$$

The eigenvalues are  $\lambda = \pm i\sqrt{\gamma}\sigma$ , with  $\sigma$  a singular value of  $\mathbf{C}$ . ■

**Theorem 5.6 (stability of EG/OGD at local saddle points)** *EG*( $\alpha, \alpha, 1$ ) is exponentially stable at any local saddle point if at such a point,  $0 < |\lambda| < 1/\alpha$  for every  $\lambda \in \text{Sp}(\mathbf{H})$ . *OGD*( $k, \alpha, \alpha$ ) is exponentially stable at any local saddle point if  $1 < k \leq 2$  and  $0 < |\lambda| < 1/(k\alpha)$  for every  $\lambda \in \text{Sp}(\mathbf{H})$ . If  $k \geq 3$ , *OGD*( $k, \alpha_1, \alpha_2$ ) is not exponentially stable for bilinear games.

**Proof** At a local saddle point, from Lemma 5.5, for any  $\lambda \in \text{Sp}(\mathbf{H})$ ,  $\Re(\lambda) \leq 0$ . The corollary follows with  $0 < |\lambda| < 1/\alpha$  for every  $\lambda \in \text{Sp}(\mathbf{H})$  and Theorem 5.2, since if  $\beta = 1$ , we can show:

$$\Re(\lambda) \leq 0 \text{ and } 0 < |\lambda| < 1 \implies |1 + \lambda + \lambda^2| < 1, \quad (\text{E.10})$$

with the following Mathematica code (rewrite  $\lambda = u + iv$  with  $u, v \in \mathbb{R}$ ):

```
Reduce[ForAll[{u, v}, u <= 0 && 0 < u^2 + v^2 < 1, (v + 2 u v)^2 + (1 + u + u^2 - v^2)^2 < 1]],
```

and the result is **True**. For OGD, if  $1 < k \leq 2$ , we use Theorem 5.2, Lemma 5.5, and the following Mathematica code (rewrite  $\lambda = u + iv$  with  $u, v \in \mathbb{R}$ ):

```
Reduce[ForAll[{u,v,k}, 0 < u^2+v^2<1/k^2 && u<=0 && 1<k<=2,
(u^2+v^2)(-3+k+(1+k)(u^2+v^2)) <2u(-1+k(u^2+v^2))]] .
```

The result is **True**. If  $k \geq 3$  and the game is bilinear, from Theorem 5.2, Theorem 5.3 and Lemma 5.5 we must have  $4|\lambda|^4 < 0$  to obtain local convergence, which is obviously false. ■

**Lemma 5.7 (spectrum of local minimax can be arbitrary)** *Given  $\alpha_1, \alpha_2 > 0$ , for any  $z \in \mathbb{C}$ , there exists a quadratic function  $q$  and a local minimax point  $(\mathbf{x}^*, \mathbf{y}^*)$  where  $z \in \text{Sp}(\mathbf{H}_{\alpha_1, \alpha_2}(q))$ .*

**Proof** Let us assume  $z = u + iv$  with  $(u, v) \in \mathbb{R}^2$ . We first construct a real polynomial:

$$(\lambda - z)(\lambda - \bar{z}) = \lambda^2 - 2u\lambda + u^2 + v^2 = 0. \quad (\text{E.11})$$

On the other hand, the characteristic polynomial of  $\mathbf{H}_{\alpha_1, \alpha_2}(q)$  with  $q(x, y) = ax^2/2 + by^2/2 + cxy$  is:

$$\lambda^2 + (\alpha_1 a - \alpha_2 b)\lambda + \alpha_1 \alpha_2 (c^2 - ab) = 0. \quad (\text{E.12})$$

Comparing (E.11) and (E.12), it suffices to require that:

$$\alpha_1 a - \alpha_2 b = -2u, \quad \alpha_1 \alpha_2 (c^2 - ab) = u^2 + v^2, \quad (\text{E.13})$$

which always has real solutions given  $(\alpha_1 > 0, \alpha_2 > 0, u, v)$ . ■

**Theorem 5.8 (stability of EG/OGD at strict local minimax points)** *Assume at a stationary point  $(\mathbf{x}^*, \mathbf{y}^*)$ ,*

$$\partial_{\mathbf{y}\mathbf{y}}^2 f \prec \mathbf{0} \text{ and } \partial_{\mathbf{x}\mathbf{x}}^2 f - \partial_{\mathbf{x}\mathbf{y}}^2 f (\partial_{\mathbf{y}\mathbf{y}}^2 f)^{-1} \partial_{\mathbf{y}\mathbf{x}}^2 f \succ \mathbf{0}. \quad (5.5)$$

*Then there exist  $\gamma_0 > 0$  and  $\alpha_0 > 0$  such that for any  $\gamma > \gamma_0, 0 < \alpha_2 < \alpha_0$  and  $\alpha_1 = \alpha_2/\gamma$ , EG and OGD (with  $k > 1$ ) are exponentially stable.*

**Proof** Assume  $\mathbf{x} \in \mathbb{R}^n$  and Using Lemma 36 of Jin et al. (2020), for any  $\delta > 0$ , there exists  $\gamma_0 > 0$ , when  $\gamma > \gamma_0$ , the eigenvalues of  $\mathbf{H}(1/\gamma, 1)$ ,  $\lambda_1, \dots, \lambda_n, \lambda_{n+1}, \dots, \lambda_{m+n}$ , are:

$$|\lambda_i + \mu_i/\gamma| < \delta/\gamma, \quad \forall i = 1, \dots, n, \quad |\lambda_{i+n} - \nu_i| < \delta, \quad \forall i = 1, \dots, m, \quad (\text{E.14})$$

where  $\mu_i \in \text{Sp}(\partial_{\mathbf{x}\mathbf{x}}^2 f - \partial_{\mathbf{x}\mathbf{y}}^2 f (\partial_{\mathbf{y}\mathbf{y}}^2 f)^{-1} \partial_{\mathbf{y}\mathbf{x}}^2 f)$  and  $\nu_i \in \text{Sp}(\partial_{\mathbf{y}\mathbf{y}}^2 f)$ . From our assumption,  $\mu_i > 0$  and  $\nu_i < 0$ . With (E.14), there exists  $\gamma_0$  such that for every  $\gamma > \gamma_0$ ,  $\Re(\lambda_i) < 0$  for all  $\lambda_i \in H(1/\gamma, 1)$ . From Theorem 5.6, EG ( $\beta = 1$ ) and OGD ( $1 < k \leq 2$ ) are exponentially stable if  $\alpha_2$  is small enough. ■

**Proposition 5.9 (stability of gradient algorithms at general local minimax points)**

There exists a quadratic function (e.g.,  $q(x, y) = -x^2 + xy$ ) and a global (thus local, from Theorem 4.4) minimax point  $\mathbf{z}^* = (\mathbf{x}^*, \mathbf{y}^*)$  where

- GDA (with momentum or alternating updates) does not converge to  $\mathbf{z}^*$ , for any hyper-parameter choice.
- If  $\alpha_1 = \alpha_2$ , or  $\alpha_2 \rightarrow 0$ , EG/OGD do not converge to  $\mathbf{z}^*$ . Otherwise there exist hyper-parameter choices such that EG/OGD converge to  $\mathbf{z}^*$ .
- Alternating OGD does not converge to  $\mathbf{z}^*$  given  $\alpha_2 \rightarrow 0$ .

**Proof** We consider  $q(x, y) := -x^2 + xy$  as the example, with  $\mathcal{X} = \mathcal{Y} = \mathbb{R}$ . From (4.1) we know that  $(0, 0)$  is a global minimax point.  $(0, 0)$  is also local minimax since it is stationary (see Theorem 4.4).  $\mathbf{H}_{1, \gamma}$  at  $(0, 0)$  is:

$$\mathbf{H}_{1, \gamma} = \begin{bmatrix} 2 & -1 \\ \gamma & 0 \end{bmatrix}. \quad (\text{E.15})$$

If  $0 < \gamma \leq 1$ , the two eigenvalues are  $1 \pm \sqrt{1 - \gamma}$  which are both real and positive. One can read from Theorem D.3 (or Figure 4) and Theorem 5.2 (or Figure 3) that GDA (with momentum) and EG/OGD do not converge to  $(0, 0)$ , locally and globally. Specifically, when  $\gamma = 1$ ,  $\alpha_1 = \alpha_2$ .

If  $\gamma > 1$ , the eigenvalues are  $\lambda_{1,2} = 1 \pm i\sqrt{\gamma - 1}$ , which have positive real parts. From Theorem D.3 (or Figure 4), GDA (with momentum) do not converge to  $(0, 0)$ . Now let us study 2TS-EG and 2TS-OGD, which corresponds to the second point of Proposition 5.9.

**2TS-EG** Taking  $\beta \rightarrow \infty$  we require that  $\Re(\lambda + \lambda^2) < 0$ , which simplifies to:

$$\alpha_1 + \alpha_1^2 - \alpha_1^2(\gamma - 1) < 0, \quad (\text{E.16})$$

and thus

$$\alpha_2 > 1 + 2\alpha_1 > 1. \quad (\text{E.17})$$

We cannot take  $\alpha_2$  to be arbitrarily small.

**2TS-OGD** For 2TS-OGD, we need  $\alpha_2$  to be  $\Omega(1)$  as well. From Theorem 5.2, we take  $k \rightarrow 1_+$  so that the convergence region is the largest:

$$|\lambda| < 1, |\lambda - 1/2| > 1/2. \quad (\text{E.18})$$

Bringing in the eigenvalues  $\alpha_1(1 \pm i\sqrt{\gamma - 1})$ , we obtain:

$$\alpha_1 < 1, 1/\alpha_1 < \gamma < 1/\alpha_1^2. \quad (\text{E.19})$$

In other words,  $1 < \alpha_2 < 1/\alpha_1$ . We could take  $\alpha_1$  infinitesimal but not  $\alpha_2$ .

**Alternating updates** Now let us study alternating updates on this example. We use the same framework as Zhang and Yu (2020). If a simultaneous algorithm takes the form of:

$$\mathbf{x}_t = T_1(\mathbf{x}_{t-1}, \mathbf{y}_{t-1}, \dots, \mathbf{x}_{t-k}, \mathbf{y}_{t-k}), \mathbf{y}_t = T_2(\mathbf{x}_{t-1}, \mathbf{y}_{t-1}, \dots, \mathbf{x}_{t-k}, \mathbf{y}_{t-k}), \quad (\text{E.20})$$

then the corresponding alternating algorithm is:

$$\mathbf{x}_t = T_1(\mathbf{x}_{t-1}, \mathbf{y}_{t-1}, \dots, \mathbf{x}_{t-k}, \mathbf{y}_{t-k}), \mathbf{y}_t = T_2(\mathbf{x}_t, \mathbf{y}_{t-1}, \dots, \mathbf{x}_{t-k+1}, \mathbf{y}_{t-k}), \quad (\text{E.21})$$

by replacing all the  $\mathbf{x}_{t-i}$  in the update function for  $\mathbf{y}_t$  to  $\mathbf{x}_{t+1-i}$ , for  $i = 1, \dots, k$ . We only study GDA and OGD in this paper for illustration purpose and other gradient algorithms follow similarly. The alternating GDA can be written as ( $\alpha_1 > 0, \alpha_2 > 0$ ):

$$\mathbf{x}_{t+1} = \mathbf{x}_t - \alpha_1 \partial_{\mathbf{x}} f(\mathbf{x}_t, \mathbf{y}_t), \mathbf{y}_{t+1} = \mathbf{y}_t + \alpha_2 \partial_{\mathbf{y}} f(\mathbf{x}_{t+1}, \mathbf{y}_t), \quad (\text{E.22})$$

and the alternating OGD can be written as (see (5.3)) ( $\alpha_1 > 0, \alpha_2 > 0, k > 1$ ):

$$\mathbf{x}_{t+1} = \mathbf{x}_t - k\alpha_1 \partial_{\mathbf{x}} f(\mathbf{x}_t, \mathbf{y}_t) + \alpha_1 \partial_{\mathbf{x}} f(\mathbf{x}_{t-1}, \mathbf{y}_{t-1}), \quad (\text{E.23})$$

$$\mathbf{y}_{t+1} = \mathbf{y}_t + k\alpha_2 \partial_{\mathbf{y}} f(\mathbf{x}_{t+1}, \mathbf{y}_t) - \alpha_2 \partial_{\mathbf{y}} f(\mathbf{x}_t, \mathbf{y}_{t-1}). \quad (\text{E.24})$$

Let us denote  $\mathbf{A} = \partial_{\mathbf{xx}}^2 f(\mathbf{x}^*, \mathbf{y}^*)$ ,  $\mathbf{B} = \partial_{\mathbf{yy}}^2 f(\mathbf{x}^*, \mathbf{y}^*)$  and  $\mathbf{C} = \partial_{\mathbf{xy}}^2 f(\mathbf{x}^*, \mathbf{y}^*)$ . Locally, we can treat the gradient algorithms as a linear dynamical system. For instance, the linear dynamical system of simultaneous GDA and simultaneous OGD can be written as:

$$\text{GDA: } \begin{pmatrix} \mathbf{x}_{t+1} - \mathbf{x}^* \\ \mathbf{y}_{t+1} - \mathbf{y}^* \end{pmatrix} = \begin{pmatrix} \mathbf{x}_t - \mathbf{x}^* \\ \mathbf{y}_t - \mathbf{y}^* \end{pmatrix} + \begin{pmatrix} -\alpha_1 \mathbf{A} & -\alpha_1 \mathbf{C} \\ \alpha_2 \mathbf{C}^\top & \alpha_2 \mathbf{B} \end{pmatrix} \begin{pmatrix} \mathbf{x}_t - \mathbf{x}^* \\ \mathbf{y}_t - \mathbf{y}^* \end{pmatrix}, \quad (\text{E.25})$$

$$\begin{aligned} \text{OGD: } \begin{pmatrix} \mathbf{x}_{t+1} - \mathbf{x}^* \\ \mathbf{y}_{t+1} - \mathbf{y}^* \end{pmatrix} &= \begin{pmatrix} \mathbf{x}_t - \mathbf{x}^* \\ \mathbf{y}_t - \mathbf{y}^* \end{pmatrix} + k \begin{pmatrix} -\alpha_1 \mathbf{A} & -\alpha_1 \mathbf{C} \\ \alpha_2 \mathbf{C}^\top & \alpha_2 \mathbf{B} \end{pmatrix} \begin{pmatrix} \mathbf{x}_t - \mathbf{x}^* \\ \mathbf{y}_t - \mathbf{y}^* \end{pmatrix} - \\ &- \begin{pmatrix} -\alpha_1 \mathbf{A} & -\alpha_1 \mathbf{C} \\ \alpha_2 \mathbf{C}^\top & \alpha_2 \mathbf{B} \end{pmatrix} \begin{pmatrix} \mathbf{x}_{t-1} - \mathbf{x}^* \\ \mathbf{y}_{t-1} - \mathbf{y}^* \end{pmatrix}. \end{aligned} \quad (\text{E.26})$$

With Theorem 2.3 from Zhang and Yu (2020), the characteristic equations for alternating GDA and alternating OGD are:

$$\text{GDA: } \det \left( (\lambda - 1) \mathbf{I} - \begin{pmatrix} -\alpha_1 \mathbf{A} & -\alpha_1 \mathbf{C} \\ \alpha_2 \lambda \mathbf{C}^\top & \alpha_2 \mathbf{B} \end{pmatrix} \right) = 0, \quad (\text{E.27})$$

$$\text{OGD: } \det \left( (\lambda - 1) \lambda \mathbf{I} - (k\lambda - 1) \begin{pmatrix} -\alpha_1 \mathbf{A} & -\alpha_1 \mathbf{C} \\ \alpha_2 \lambda \mathbf{C}^\top & \alpha_2 \mathbf{B} \end{pmatrix} \right) = 0. \quad (\text{E.28})$$

For the quadratic example  $q(x, y) = -x^2 + xy$  we are considering, we have  $\mathbf{A} = -2, \mathbf{B} = 0, \mathbf{C} = 1$ . Bringing it to (E.27), we obtain:

$$\text{GDA: } \lambda^2 + (\alpha_1 \alpha_2 - 2\alpha_1 - 2)\lambda + 2\alpha_1 + 1 = 0, \quad (\text{E.29})$$

$$\text{OGD: } \lambda^4 + (\alpha_1 \alpha_2 k^2 - 2\alpha_1 k - 2)\lambda^3 + (2\alpha_1 - 2\alpha_1 \alpha_2 k + 2\alpha_1 k + 1)\lambda^2 + (\alpha_1 \alpha_2 - 2\alpha_1)\lambda = 0. \quad (\text{E.30})$$

From Corollary 2.1 of Zhang and Yu (2020), alternating GDA is stable iff:

$$2\alpha_1 + 1 < 1, |\alpha_1 \alpha_2 - 2\alpha_1 - 2| < 2\alpha_1 + 2. \quad (\text{E.31})$$



Note that the first condition can never hold since  $\alpha_1 > 0$ . Hence, alternating GDA cannot converge to the local minimax point  $(0, 0)$  if the initialization is not at  $(0, 0)$ . For alternating OGD, the second equation of (E.29) can be simplified as  $\lambda = 0$  or:

$$\lambda^3 + (\alpha_1\alpha_2k^2 - 2\alpha_1k - 2)\lambda^2 + (2\alpha_1 - 2\alpha_1\alpha_2k + 2\alpha_1k + 1)\lambda + \alpha_1(\alpha_2 - 2) = 0. \quad (\text{E.32})$$

Using Corollary 2.1 of Zhang and Yu (2020) again we know that alternating OGD is stable iff:

$$|c| < 1, |a + c| < 1 + b, b - ac < 1 - c^2, \quad (\text{E.33})$$

where  $a = \alpha_1\alpha_2k^2 - 2\alpha_1k - 2$ ,  $b = 2\alpha_1 - 2\alpha_1\alpha_2k + 2\alpha_1k + 1$ ,  $c = \alpha_1(\alpha_2 - 2)$ . We simplify it on Mathematica:

```
Reduce[Abs[c] < 1 && Abs[a+c] < 1 + b && b - a c < 1 - c^2 && k > 1
&& \alpha_1 > 0 && \alpha_2 > 0, {\alpha_1, \alpha_2}]
```

and obtain that:

$$k > 1 \text{ and } 0 < \alpha_1 < \frac{4}{k^2 - 1} \text{ and} \\ \sqrt{\frac{-2\alpha_1 + \alpha_1^2k^2 + 1}{\alpha_1^2(k+1)^2}} + \frac{2\alpha_1 + \alpha_1k - 1}{\alpha_1(k+1)} < \alpha_2 < \frac{4\alpha_1 + 4\alpha_1k + 4}{\alpha_1 + \alpha_1k^2 + 2\alpha_1k}. \quad (\text{E.34})$$

Since  $k > 1$  and

$$\begin{aligned} \sqrt{\frac{-2\alpha_1 + \alpha_1^2k^2 + 1}{\alpha_1^2(k+1)^2}} + \frac{2\alpha_1 + \alpha_1k - 1}{\alpha_1(k+1)} &\geq \sqrt{\frac{-2\alpha_1 + \alpha_1^2 + 1}{\alpha_1^2(k+1)^2}} + \frac{2\alpha_1 + \alpha_1k - 1}{\alpha_1(k+1)} \\ &= \frac{\alpha_1k + 2\alpha_1 - 1 + |\alpha_1 - 1|}{\alpha_1(k+1)} \\ &\geq \frac{\alpha_1k + 2\alpha_1 - 1 + 1 - \alpha_1}{\alpha_1(k+1)} \\ &= 1, \end{aligned} \quad (\text{E.35})$$

we have  $\alpha_2 > 1$  for alternating updates of OGD. ■

## Appendix F. Local robust points

In this section, we summarize results about local robust points, which naturally extend local minimax points to a symmetric version. They are stationary points (Theorem F.6), but they may not correspond to solution concepts in sequential games (Example F.3). In one-dimensional case they are equivalent to the stable sets of Optimistic Gradient Descent (Proposition F.15). However, in general cases all common coordinate-independent gradient algorithms would fail to converge to some local robust point (Proposition F.16). The main results are summarized in Table 1.

	Statement	Reference
LRP	non-trivial examples	Prop. F.2, Eg. F.3
	nuances in the definition	Examples F.4, F.5
	LRPs are stationary points	Theorem F.6
	optimality conditions	Theorems F.6, F.7, F.8, F.11
	LRP in quadratic games	Theorem F.14
	equivalence with the stable set of OGD in 1D	Prop. F.15
	failure of gradient algorithms at LRP	Proposition F.16

Table 1 Results of local robust points.

### F.1 Definition of local robust points

In the definition of local minimax points,  $\mathbf{x}$  and  $\mathbf{y}$  are *asymmetric*:  $\mathbf{y}$  is the follower who knows the strategy of  $\mathbf{x}$ , but  $\mathbf{x}$  only knows a “rough” set of the strategies of  $\mathbf{y}$  and hence aims to optimize the worst-case scenario. One natural (and perhaps more realistic) generalization is to allow robust optimization for  $\mathbf{y}$  as well, so as to restore equal position for both players:

**Definition F.1 (LRP)** We call  $(\mathbf{x}^*, \mathbf{y}^*) \in \mathcal{X} \times \mathcal{Y}$  a local robust point (LRP) if

- fixing  $\mathbf{x}^*$ , there exists some sequence  $0 \leq \varepsilon_n \rightarrow 0$  such that for each  $\varepsilon_n$  in the sequence, there exists an envelope function  $\underline{f}_{\varepsilon_n, \mathbf{x}^*}(\mathbf{y})$  such that  $\mathbf{y}^*$  is a local maximizer;
- fixing  $\mathbf{y}^*$ , there exists some sequence  $0 \leq \varepsilon_n \rightarrow 0$  such that for each  $\varepsilon_n$  in the sequence, there exists an envelope functions  $\bar{f}_{\varepsilon_n, \mathbf{y}^*}(\mathbf{x})$  such that  $\mathbf{x}^*$  is a local minimizer.

In the above definition, both  $\mathbf{x}$  and  $\mathbf{y}$  are doing robust optimization:  $\bar{f}_\varepsilon(\mathbf{x})$  and  $-\underline{f}_\varepsilon(\mathbf{y})$  can be treated as the worst-case cost for each player, assuming that each one only knows an approximate strategy of the opponent ( $\mathbf{x}^*$  or  $\mathbf{y}^*$ ), up to some estimation error ( $\varepsilon$  or  $\varepsilon$ ). Since each player does not know the exact amount of perturbation, it will try to minimize a sequence of envelope functions with a series of neighborhoods that can be arbitrarily small.

LRPs are a subclass of stationary points, as we will see in Theorem F.6. The definition of LRPs includes local saddle, local minimax and local maximin points, as visualized in Figure 5. For example, if  $\{\varepsilon_n\} = \{0\}$  and  $0 < \varepsilon_n \rightarrow 0$ , then LRP reduces to local minimax points. The simplest non-trivial example for LRPs might be quadratic games. In general for one-dimensional quadratic games, it can be shown that:

**Proposition F.2 (characterization of LRPs in one-dimensional quadratic games)**

$f(x, y) = ax^2/2 + cxy + by^2/2$  has an LRP at  $(0, 0)$  iff

$$\{c = 0, a \geq 0 \geq b\} \text{ or } \{c \neq 0, c^2 \geq ab\}. \quad (\text{F.1})$$

**Proof** If  $c = 0$ ,  $f$  is separable, we obtain  $a \geq 0$  because  $x^*$  locally minimizes  $\bar{f}_\varepsilon(x)$ , and  $b \leq 0$  since  $y^*$  locally maximizes  $\bar{f}_\varepsilon(y)$ . If  $c \neq 0$ , then for small enough  $x, y$ ,

$$\bar{f}_\varepsilon(x) = \begin{cases} |cx|\varepsilon + b\varepsilon^2/2 + ax^2/2 & \text{if } b \geq 0 \\ (c^2 - ab)x^2/(-2b) & \text{if } b < 0 \end{cases}, \quad \underline{f}_\varepsilon(y) = \begin{cases} -|cy|\varepsilon + by^2/2 + a\varepsilon^2/2 & \text{if } a \leq 0 \\ -(c^2 - ab)y^2/(2a) & \text{if } a > 0 \end{cases} \quad (\text{F.2})$$

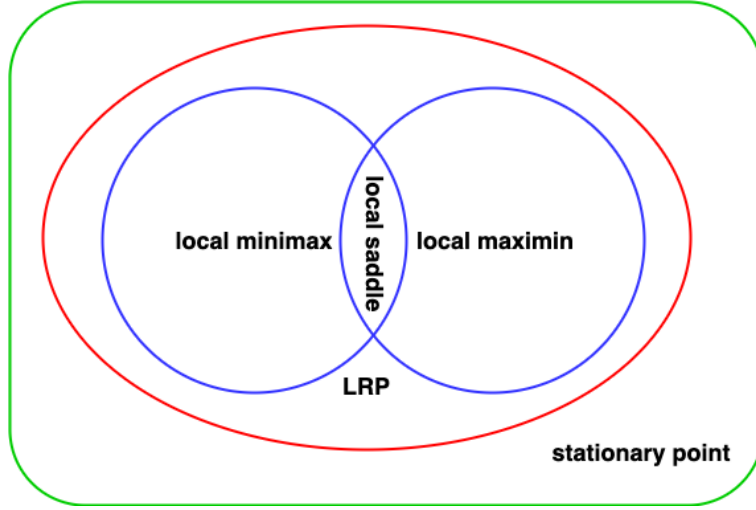


Figure 5 The relation among the sets of local saddle, local minimax and local maximin points, as well as LRPs. In the unconstrained case, they are all stationary (Theorem 3.12).

From the above, we can show that it is necessary and sufficient to have  $c^2 \geq ab$ : if  $c^2 \geq ab$ , then  $\bar{f}_\epsilon(x)$  is locally minimized at  $x = 0$  and  $\underline{f}_\epsilon(y)$  is locally maximized at  $y = 0$ ; if  $c^2 < ab$ , then  $a > 0, b > 0$ , when  $\underline{f}_\epsilon(y)$  is not locally maximized at  $y = 0$ , or  $a < 0, b < 0$ , when  $\bar{f}_\epsilon(x)$  is not locally minimized at  $x = 0$ . ■

If  $c = 0$  and  $a = -2, b = 2$ , then this quadratic function clearly does not have an LRP (but has a stationary point), which implies the non-triviality of our definition. Another interesting case is when  $a = -2, c = 1$  and  $b = 2$ :

**Example F.3 (LRPs may not be either local minimax or maximin)** Consider  $f(x, y) = -x^2 + xy + y^2$  and  $(x^*, y^*) = (0, 0)$  with the domain  $|x| \leq D, |y| \leq D$ . Straightforward calculation gives (assuming  $0 < \epsilon \leq D, 0 < \varepsilon \leq D$ ):

$$\bar{f}_\epsilon(x) = -x^2 + \epsilon|x| + \epsilon^2, \quad \underline{f}_\epsilon(y) = -\varepsilon^2 - \varepsilon|y| + y^2. \quad (\text{F.3})$$

Thus,  $f$  has an LRP at  $(0, 0)$ , which is neither local minimax or local maximin:  $f(0, y) = y^2$  is not locally maximized at  $y = 0$  and  $f(x, 0) = -x^2$  is not locally minimized at  $x = 0$ . Note that  $(0, 0)$  is not a global minimax/maximin point either. However, we have:

$$\begin{aligned} \bar{f}_D(x) &= \max_{|y| \leq D} f(x, y) = -x^2 + D|x| + D^2 \geq \bar{f}_D(0), \text{ for all } |x| \leq D \\ \underline{f}_D(y) &= \min_{|x| \leq D} f(x, y) = -D^2 - D|y| + y^2 \leq \underline{f}_D(0), \text{ for all } |y| \leq D. \end{aligned} \quad (\text{F.4})$$

So  $(0, 0)$  can be treated as some type of “global robust point”, defined as

$$\sup_{\mathbf{y} \in \mathcal{Y}} f(\mathbf{x}, \mathbf{y}) \geq \sup_{\mathbf{y} \in \mathcal{Y}} f(\mathbf{x}^*, \mathbf{y}), \text{ for any } \mathbf{x} \in \mathcal{X} \quad (\text{F.5})$$

$$\inf_{\mathbf{x} \in \mathcal{X}} f(\mathbf{x}, \mathbf{y}) \leq \inf_{\mathbf{x} \in \mathcal{X}} f(\mathbf{x}, \mathbf{y}^*), \text{ for any } \mathbf{y} \in \mathcal{Y}. \quad (\text{F.6})$$

In such a game, each player is agnostic of the opponent's strategy and only optimizing the worst case. There is no follower or leader. Such study goes beyond the regime of sequential games and we leave it to future research.

However, for LRPs, some results we derived in Section 3.1 for local minimax points cease to hold anymore. For example, for local minimax points the norm we choose in the neighborhood definition is immaterial (see Lemma 3.5), but for LRPs, that choice of the neighborhoods does matter, as can be seen from the following example:

**Example F.4 (effect of the neighborhood)** Consider the function

$$f(\mathbf{x}, \mathbf{y}) = -\mathbf{x}^\top \begin{bmatrix} 0 & 0 \\ 0 & 1 \end{bmatrix} \mathbf{x} + \mathbf{x}^\top \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix} \mathbf{y} + \mathbf{y}^\top \begin{bmatrix} 1 & 0 \\ 0 & 0 \end{bmatrix} \mathbf{y}, \quad (\text{F.7})$$

with  $\mathcal{X} = \mathcal{Y} = \mathbb{R}^2$  and  $(\mathbf{x}^*, \mathbf{y}^*) = (\mathbf{0}, \mathbf{0})$ . For the  $\ell_\infty$  normed ball  $\mathcal{N}_\infty(\mathbf{y}^*, \epsilon) = \{\mathbf{y} \in \mathbb{R}^2 : \|\mathbf{y} - \mathbf{y}^*\|_\infty \leq \epsilon\}$ ,  $\bar{f}_\epsilon(\mathbf{x}) = \epsilon^2 + \epsilon|x_1| + \epsilon|x_2| - x_2^2$  which is locally minimized at  $\mathbf{x}^*$ . However, for the Euclidean ball  $\mathcal{N}_2(\mathbf{y}^*, \epsilon) = \{\mathbf{y} \in \mathbb{R}^2 : \|\mathbf{y} - \mathbf{y}^*\|_2 \leq \epsilon\}$ ,

$$\bar{f}_\epsilon(0, x_2) = \max_{\mathbf{y} \in \mathcal{N}_2(\mathbf{y}^*, \epsilon)} x_2 y_2 + y_1^2 - x_2^2 \leq \max_{|y_2| \leq \epsilon} \epsilon^2 - y_2^2 + x_2 y_2 - x_2^2 \leq \epsilon^2 - 3x_2^2/4 < \bar{f}_\epsilon(0, 0) = \epsilon^2,$$

for any  $0 < |x_2| < 2\epsilon$ . One can show that  $(\mathbf{x}^*, \mathbf{y}^*) = (\mathbf{0}, \mathbf{0})$  is an LRP by choosing the neighborhoods of  $\mathbf{x}^*$  and  $\mathbf{y}^*$  to be  $\ell_\infty$  balls, since

$$\bar{f}_\epsilon(\mathbf{x}) = \epsilon^2 + \epsilon|x_1| + \epsilon|x_2| - x_2^2 \geq \bar{f}_\epsilon(\mathbf{0}) \text{ locally and } \underline{f}_\epsilon(\mathbf{y}) = -\epsilon^2 - \epsilon|y_1| - \epsilon|y_2| + y_1^2 \leq \underline{f}_\epsilon(\mathbf{0})$$

locally. In Appendix F.3 we will show a “meaningful” neighborhood choice for LRPs in quadratic games using the eigenspace.

In order for the class of LRPs to include the class of local minimax points, we may no longer take  $\{\epsilon_n\}$  and  $\{\varepsilon_n\}$  to be strictly positive sequences as in Def. 3.3:

**Example F.5 (The definition of LRPs need to include  $\epsilon = 0$  and  $\varepsilon = 0$ )** Take

$$f(x, y) = xy^3 - x^2/(1 + y^2)$$

and  $(x^*, y^*) = (0, 0)$ . This point is a local minimax point, since  $\underline{f}_0(y) = f(x^*, y) = 0$ , and  $\bar{f}_\epsilon(x) \geq \epsilon^3|x| - x^2/(1 + \epsilon^2) \geq 0 = \bar{f}_\epsilon(x^*)$ , given small enough  $x$ . However, for any  $\varepsilon > 0$ ,

$$\underline{f}_\varepsilon(y) = -\varepsilon|y|^3 - \varepsilon^2/(1 + y^2) \text{ and } \underline{f}_\varepsilon(y) - \underline{f}_\varepsilon(y^*) = \varepsilon y^2(\varepsilon/(1 + y^2) - |y|) > 0$$

for small enough  $y$ . Therefore, in Definition F.1 the case of  $\varepsilon = 0$  needs to be included, as otherwise  $(x^*, y^*) = (0, 0)$  does not satisfy the definition of LRPs, since for any  $\varepsilon > 0$ , the variable  $y^*$  cannot be a local maximizer of  $\underline{f}_\varepsilon$ .

## F.2 Optimality conditions for LRPs

Let us define the *active sets* of the *zeroth* order (by “zeroth” we mean that only the function values are involved):

$$\mathcal{Y}_0(\mathbf{x}^*; \epsilon) = \{\mathbf{y} \in \mathcal{N}(\mathbf{y}^*, \epsilon) : \bar{f}_\epsilon(\mathbf{x}^*) = f(\mathbf{x}^*, \mathbf{y})\}, \quad (\text{F.8})$$

$$\mathcal{X}_0(\mathbf{y}^*; \epsilon) = \{\mathbf{x} \in \mathcal{N}(\mathbf{x}^*, \epsilon) : \underline{f}_\epsilon(\mathbf{y}^*) = f(\mathbf{x}, \mathbf{y}^*)\}. \quad (\text{F.9})$$

We derive the first-order optimality conditions for LRPs.

**Theorem F.6 (first-order necessary, LRP)** *Let  $f \in \mathcal{C}^1$ . At an LRP  $(\mathbf{x}^*, \mathbf{y}^*)$ , we have:*

$$\partial_{\mathbf{x}} f(\mathbf{x}^*, \mathbf{y}^*)^\top \bar{\mathbf{t}} \geq 0 \geq \partial_{\mathbf{y}} f(\mathbf{x}^*, \mathbf{y}^*)^\top \underline{\mathbf{t}}, \quad (\text{F.10})$$

for any directions  $\bar{\mathbf{t}} \in \mathbf{K}_d(\mathcal{X}, \mathbf{x}^*)$ ,  $\underline{\mathbf{t}} \in \mathbf{K}_d(\mathcal{Y}, \mathbf{y}^*)$ , where the cone

$$\mathbf{K}_d(\mathcal{X}, \mathbf{x}) := \liminf_{\alpha \rightarrow 0^+} \frac{\mathcal{X} - \mathbf{x}}{\alpha} := \{\mathbf{t} : \forall \{\alpha_k\} \rightarrow 0^+ \exists \{\alpha_{k_i}\} \rightarrow 0^+, \{\mathbf{t}_{k_i}\} \rightarrow \mathbf{t}, \\ \text{such that } \mathbf{x} + \alpha_{k_i} \mathbf{t}_{k_i} \in \mathcal{X}\}$$

and  $\mathbf{K}_d(\mathcal{Y}, \mathbf{y})$  is defined similarly.

**Proof** Use Theorem A.3, Theorem A.9 and the assumption that  $f \in \mathcal{C}^1$ . ■

**Theorem F.7 (first-order sufficient condition, LRP)** *If  $f$  is continuously differentiable and there exist two sequences  $\epsilon_n \rightarrow 0$ ,  $\varepsilon_n \rightarrow 0$ , such that for any  $n \in \mathbb{N}^+$ :*

$$\mathbf{0} \neq \bar{\mathbf{t}} \in \mathbf{K}_c(\mathcal{X}, \mathbf{x}^*) \implies \mathbf{D}\bar{f}_{\epsilon_n}(\mathbf{x}^*; \bar{\mathbf{t}}) = \max_{\mathbf{y} \in \mathcal{Y}_0(\mathbf{x}^*; \epsilon_n)} \partial_{\mathbf{x}} f(\mathbf{x}, \mathbf{y})^\top \bar{\mathbf{t}} > 0, \quad (\text{F.11})$$

$$\mathbf{0} \neq \underline{\mathbf{t}} \in \mathbf{K}_c(\mathcal{Y}, \mathbf{y}^*) \implies \mathbf{D}\underline{f}_{\varepsilon_n}(\mathbf{y}^*; \underline{\mathbf{t}}) = \min_{\mathbf{x} \in \mathcal{X}_0(\mathbf{y}^*; \varepsilon_n)} \partial_{\mathbf{y}} f(\mathbf{x}, \mathbf{y})^\top \underline{\mathbf{t}} < 0. \quad (\text{F.12})$$

then  $(\mathbf{x}^*, \mathbf{y}^*)$  is an isolated LRP of  $f$ .

We next discuss how to obtain second-order conditions for LRPs. Recalling Definition F.1, for the second-order optimality conditions of the local maximality of min-type envelope functions  $\underline{f}_\epsilon(\mathbf{y})$ , we can simply take  $f \rightarrow -f$ ,  $\bar{f}_\epsilon(\mathbf{x}) \rightarrow -\underline{f}_\epsilon(\mathbf{y})$  and switch the roles of  $\mathbf{x}$  and  $\mathbf{y}$ . Let us define that:

$$\bar{u}_\epsilon(\mathbf{y}) := \bar{f}_\epsilon(\mathbf{x}^*) - f(\mathbf{x}^*, \mathbf{y}), \quad \bar{v}(\mathbf{y}; \mathbf{t}) = -\partial_{\mathbf{x}} f(\mathbf{x}^*, \mathbf{y})^\top \mathbf{t}, \\ \mathcal{Y}_1(\epsilon; \mathbf{t}) = \{\mathbf{y} \in \mathcal{N}(\mathbf{y}^*, \epsilon) : \bar{u}_\epsilon(\mathbf{y}) = \bar{v}(\mathbf{y}; \mathbf{t}) = 0\}, \quad (\text{F.13})$$

$$\underline{u}_\epsilon(\mathbf{x}) := f(\mathbf{x}, \mathbf{y}^*) - \underline{f}_\epsilon(\mathbf{y}^*), \quad \underline{v}(\mathbf{x}; \mathbf{t}) = \partial_{\mathbf{y}} f(\mathbf{x}, \mathbf{y}^*)^\top \mathbf{t}, \\ \mathcal{X}_1(\varepsilon; \mathbf{t}) = \{\mathbf{x} \in \mathcal{N}(\mathbf{x}^*, \varepsilon) : \underline{u}_\epsilon(\mathbf{x}) = \underline{v}(\mathbf{x}; \mathbf{t}) = 0\}, \quad (\text{F.14})$$

and

$$\bar{E}_\epsilon(\mathbf{y}; \mathbf{t}) = \limsup_{\mathbf{z} \rightarrow \mathbf{y}} \frac{1}{2} \bar{v}_-^2(\mathbf{z}; \mathbf{d}) \bar{u}_\epsilon^\dagger(\mathbf{z}), \quad \underline{E}_\varepsilon(\mathbf{x}; \mathbf{t}) = \limsup_{\mathbf{z} \rightarrow \mathbf{x}} \underline{v}_-(\mathbf{z}; \mathbf{t})^2 \underline{u}_\epsilon^\dagger(\mathbf{z})/2. \quad (\text{F.15})$$

We obtain the second-order necessary conditions for LRPs from Theorem A.14:

**Theorem F.8 (second-order necessary condition, LRP)** *If  $(\mathbf{x}^*, \mathbf{y}^*)$  is an LRP with sequence  $\{\epsilon_k\}, \{\varepsilon_k\}$ , then for any  $\epsilon_k$ , for each direction  $\bar{\mathbf{t}} \in \mathbb{R}^n$ ,  $\mathbf{D}\bar{f}_{\epsilon_k}(\mathbf{x}^*; \bar{\mathbf{t}}) > 0$ , or  $\mathbf{D}\bar{f}_{\epsilon_k}(\mathbf{x}^*; \bar{\mathbf{t}}) = 0$  and there exist at most  $n+1$  points  $\mathbf{y}_1, \dots, \mathbf{y}_{n+1} \in \mathcal{Y}_1(\epsilon_k; \bar{\mathbf{t}})$  and  $\lambda_1, \dots, \lambda_n \geq 0$  not all zero, such that:*

$$\sum_{i=1}^{n+1} \lambda_i \partial_{\mathbf{x}} f(\mathbf{x}^*, \mathbf{y}_i) = \mathbf{0}, \quad \sum_{i=1}^{n+1} \lambda_i \left( \bar{\mathbf{t}}^\top \partial_{\mathbf{xx}}^2 f(\mathbf{x}^*, \mathbf{y}_i) \bar{\mathbf{t}} + \bar{E}_{\epsilon_k}(\mathbf{y}_i, \bar{\mathbf{t}}) \right) \geq 0. \quad (\text{F.16})$$

*For each feasible direction  $\underline{\mathbf{t}} \in \mathbb{R}^m$ ,  $\mathbf{D}\underline{f}_{\varepsilon_k}(\mathbf{y}^*; \underline{\mathbf{t}}) < 0$ , or  $\mathbf{D}\underline{f}_{\varepsilon_k}(\mathbf{y}^*; \underline{\mathbf{t}}) = 0$  and there exist at most  $m+1$  points  $\mathbf{x}_1, \dots, \mathbf{x}_{m+1} \in \mathcal{X}_1(\varepsilon_k; \underline{\mathbf{t}})$  and  $\mu_1, \dots, \mu_m \geq 0$  not all zero, such that:*

$$\sum_{i=1}^{m+1} \mu_i \partial_{\mathbf{x}} f(\mathbf{x}_i, \mathbf{y}^*) = \mathbf{0}, \quad \sum_{i=1}^{m+1} \mu_i \left( \underline{\mathbf{t}}^\top \partial_{\mathbf{yy}}^2 f(\mathbf{x}_i, \mathbf{y}^*) \underline{\mathbf{t}} - \underline{E}_{\varepsilon_k}(\mathbf{x}_i, \underline{\mathbf{t}}) \right) \leq 0. \quad (\text{F.17})$$

**Remark F.9** *For LRPs we do not have the simplification as local minimax points in Theorem 3.17 since Lemma 3.16 does not necessarily hold. In fact,  $\mathbf{y}^*$  may not even be in the active set  $\mathcal{Y}_0(\mathbf{x}^*)$  (e.g. Example F.3). Comparably, for a local minimax point  $(\mathbf{x}^*, \mathbf{y}^*)$ ,  $\mathbf{y}^* \in \mathcal{Y}_0(\mathbf{x}^*)$  and  $\bar{u}_\epsilon(\mathbf{y}^*)$  is a constant for small enough  $\epsilon$ .*

It is also possible to construct second-order sufficient conditions for LRPs from Theorem A.17 and Theorem A.6. We only construct one from Theorem A.17 as the other construction is analogous. Similar to Assumption A.16, we need the following assumption:

**Assumption F.10** *For each  $\mathbf{x} \in \mathcal{X}_1(\varepsilon; \mathbf{t})$  with  $\mathbf{t} \neq \mathbf{0}$  and  $\mathbf{D}\underline{f}_\varepsilon(\mathbf{x}; \mathbf{t}) = 0$ , and for each non-zero  $\mathbf{d} \in \mathbb{R}^m$ , there exist  $\alpha, \beta \neq 0$  and  $p, q > 0$  such that the following approximation holds:*

$$\underline{u}_\varepsilon(\mathbf{x} + \delta \mathbf{d}) = \alpha \delta^p + o(\delta^p), \quad \underline{v}(\mathbf{x} + \delta \mathbf{d}; \mathbf{t}) = \beta \delta^q + o(\delta^q), \quad (\text{F.18})$$

*whenever  $\mathbf{x} + \delta \mathbf{d} \in \mathcal{N}(\mathbf{x}^*, \varepsilon)$  and  $\delta > 0$ .*

With this assumption and Assumption A.16 (with a slight change of notations) we can write down the second-order sufficient condition for LRPs, similar to Theorem F.8:

**Theorem F.11 (second-order sufficient condition, LRP)** *Assume that Assumption A.16 and Assumption F.10 hold, and let  $\mathcal{X} = \mathbb{R}^n$  and  $\mathcal{Y} = \mathbb{R}^m$ . Suppose there exists a sequence  $\{\epsilon_k\}$  such that for any  $\epsilon_k$ , for each direction  $\bar{\mathbf{t}} \in \mathbb{R}^n$ ,  $\mathbf{D}\bar{f}_{\epsilon_k}(\mathbf{x}^*; \bar{\mathbf{t}}) > 0$ , or  $\mathbf{D}\bar{f}_{\epsilon_k}(\mathbf{x}^*; \bar{\mathbf{t}}) = 0$  and there exist  $a \geq 1$  points  $\mathbf{y}_1, \dots, \mathbf{y}_a \in \mathcal{Y}_1(\epsilon_k; \bar{\mathbf{t}})$  and  $\lambda_1, \dots, \lambda_a \geq 0$  not all zero, such that:*

$$\sum_{i=1}^a \lambda_i \partial_{\mathbf{x}} f(\mathbf{x}^*, \mathbf{y}_i) = \mathbf{0}, \quad \sum_{i=1}^a \lambda_i \left( \bar{\mathbf{t}}^\top \partial_{\mathbf{xx}}^2 f(\mathbf{x}^*, \mathbf{y}_i) \bar{\mathbf{t}} + \bar{E}_{\epsilon_k}(\mathbf{y}_i, \bar{\mathbf{t}}) \right) > 0. \quad (\text{F.19})$$

*If moreover there exists a sequence  $\{\varepsilon_k\}$  such that for any  $\varepsilon_k$ , along each  $\underline{\mathbf{t}} \in \mathbb{R}^m$ ,  $\mathbf{D}\underline{f}_{\varepsilon_k}(\mathbf{y}^*; \underline{\mathbf{t}}) < 0$ , or  $\mathbf{D}\underline{f}_{\varepsilon_k}(\mathbf{y}^*; \underline{\mathbf{t}}) = 0$  and there exist  $b \geq 1$  points  $\mathbf{x}_1, \dots, \mathbf{x}_b \in \mathcal{X}_1(\varepsilon_k; \underline{\mathbf{t}})$  and  $\mu_1, \dots, \mu_b \geq 0$  not all zero, such that:*

$$\sum_{i=1}^b \mu_i \partial_{\mathbf{y}} f(\mathbf{x}_i, \mathbf{y}^*) = \mathbf{0}, \quad \sum_{i=1}^b \mu_i \left( \underline{\mathbf{t}}^\top \partial_{\mathbf{yy}}^2 f(\mathbf{x}_i, \mathbf{y}^*) \underline{\mathbf{t}} - \underline{E}_{\varepsilon_k}(\mathbf{x}_i, \underline{\mathbf{t}}) \right) < 0, \quad (\text{F.20})$$

*then  $(\mathbf{x}^*, \mathbf{y}^*)$  is an LRP.*

### F.3 Local robust points in quadratic games

In this subsection, we discuss the existence conditions for LRPs in quadratic games. Since LRPs are also stationary, we can translate the origin such that the quadratic game is homogeneous.

**Definition F.12 (positive/negative part of a symmetric matrix)** For an  $n$ -dimensional symmetric matrix  $\mathbf{A} \in \mathbb{S}^n$ , given its spectral decomposition  $\mathbf{A} = \mathbf{U}\mathbf{D}\mathbf{U}^\top$ , we define the positive part  $\mathbf{A}_p = \mathbf{U}\mathbf{D}_p\mathbf{U}^\top$ , and the negative part is  $\mathbf{A}_n = \mathbf{U}\mathbf{D}_n\mathbf{U}^\top$ , where  $[\mathbf{D}_p]_{i,j} = d_{ii}\delta_{i,j}\mathbf{1}_{d_{ii}>0}$  (resp.  $[\mathbf{D}_n]_{i,j} = d_{ii}\delta_{i,j}\mathbf{1}_{d_{ii}<0}$ ) is a diagonal matrix that takes the positive part (resp. the negative part) of  $\mathbf{D}$ .

**Definition F.13 (eigenspace neighborhood)** Given the spectral decomposition of a symmetric matrix  $\mathbf{A} = \sum_i \lambda_i \mathbf{v}_i \mathbf{v}_i^\top$ , we define the eigenspace neighborhood w.r.t.  $\mathbf{A}$  as:

$$\mathcal{N}_{\mathbf{A}}(\mathbf{x}, \epsilon) := \{\mathbf{x} + \sum_i c_i \mathbf{v}_i : |c_i| \leq \epsilon\}. \quad (\text{F.21})$$

With the decomposition of symmetric matrices and the eigenspace neighborhoods, we can derive the condition for LRPs in unconstrained quadratic games:

#### Theorem F.14 (necessary and sufficient conditions of LRPs in quadratic games)

Let us choose  $\mathcal{N}(\mathbf{y}^*, \epsilon) = \mathcal{N}_{\mathbf{B}}(\mathbf{y}^*, \epsilon)$  and  $\mathcal{N}(\mathbf{x}^*, \epsilon) = \mathcal{N}_{\mathbf{A}}(\mathbf{x}^*, \epsilon)$  for envelope functions  $\bar{f}_\epsilon(\mathbf{x})$  and  $\underline{f}_\epsilon(\mathbf{y})$  respectively. In order for  $(\mathbf{x}^*, \mathbf{y}^*) = (\mathbf{0}, \mathbf{0})$  to be an LRP for the homogeneous quadratic game, it is necessary and sufficient that:

$$\mathbf{P}_{\mathbf{L}}^\perp (\mathbf{A} - \mathbf{C}\mathbf{B}_n^\dagger \mathbf{C}^\top) \mathbf{P}_{\mathbf{L}}^\perp \succeq \mathbf{0}, \quad \mathbf{L} = \mathbf{C}\mathbf{P}_{\mathbf{B}_n}^\perp, \quad (\text{F.22})$$

$$\mathbf{P}_{\mathbf{M}}^\perp (\mathbf{B} - \mathbf{C}^\top \mathbf{A}_p^\dagger \mathbf{C}) \mathbf{P}_{\mathbf{M}}^\perp \preceq \mathbf{0}, \quad \mathbf{M} = \mathbf{C}^\top \mathbf{P}_{\mathbf{A}_p}^\perp. \quad (\text{F.23})$$

**Proof** Given the spectral decomposition  $\mathbf{B} = \sum_i b_i \mathbf{v}_i \mathbf{v}_i^\top$  and  $\mathbf{y} = \sum_i y_i \mathbf{v}_i$ , the quadratic function can be written as:

$$q(\mathbf{x}, \mathbf{y}) = \mathbf{x}^\top \mathbf{A} \mathbf{x} / 2 + \sum_i b_i y_i^2 / 2 + \sum_i y_i \mathbf{x}^\top \mathbf{C} \mathbf{v}_i. \quad (\text{F.24})$$

Maximizing over the eigenspace neighborhood of  $\mathcal{N}(\mathbf{y}^*, \epsilon)$  we obtain:

$$\bar{q}_\epsilon(\mathbf{x}) = \mathbf{x}^\top (\mathbf{A} - \mathbf{C}\mathbf{B}_n^\dagger \mathbf{C}^\top) \mathbf{x} / 2 + \sum_{i \in \mathcal{I}_+} (b_i \epsilon^2 / 2 + \epsilon |\mathbf{x}^\top \mathbf{C} \mathbf{v}_i|), \quad \mathcal{I}_+ := \{i \in [m] : b_i \geq 0\}. \quad (\text{F.25})$$

In order for  $\bar{q}_\epsilon(\mathbf{x}) \geq \bar{q}_\epsilon(\mathbf{x}^*)$ , it is necessary that for all  $\mathbf{x}$  such that  $\mathbf{v}_i^\top \mathbf{C}^\top \mathbf{x} = 0$  for  $i \in \mathcal{I}_+$ ,  $\mathbf{x}^\top (\mathbf{A} - \mathbf{C}\mathbf{B}_n^\dagger \mathbf{C}^\top) \mathbf{x} / 2 \geq 0$ . That is, for all  $\mathbf{L}^\top \mathbf{x} = \mathbf{0}$  with  $\mathbf{L} := \mathbf{C}\mathbf{P}_{\mathbf{B}_n}^\perp$ ,  $\mathbf{x}^\top (\mathbf{A} - \mathbf{C}\mathbf{B}_n^\dagger \mathbf{C}^\top) \mathbf{x} / 2 \geq 0$ , which yields (F.22). Symmetrically we obtain (F.23) for maximizing  $q_\epsilon(\mathbf{y})$ . The sufficient part is analogous to the proof of Theorem 4.1. Denote  $\boldsymbol{\eta}$  as an  $|\mathcal{I}_+|$ -dimensional vector with  $\eta_i = \mathbf{v}_i^\top \mathbf{C}^\top \mathbf{x}$  and  $i \in \mathcal{I}_+$ , then

$$\sum_{i \in \mathcal{I}_+} |\mathbf{x}^\top \mathbf{C} \mathbf{v}_i| = \|\boldsymbol{\eta}\|_1 \geq \|\boldsymbol{\eta}\|_2 = \left\| \sum_{i \in \mathcal{I}_+} (\mathbf{v}_i^\top \mathbf{C}^\top \mathbf{x}) \mathbf{v}_i \right\|_2 = \|\mathbf{L}^\top \mathbf{x}\|_2. \quad (\text{F.26})$$

The rest follows after (C.12). ■

In the special case of local minimax when  $\mathbf{B} \preceq \mathbf{0}$ , (F.22) and (F.23) reduces to (4.4).

#### F.4 Stability at local robust points

Finally, we discuss the convergence of first-order algorithms near LRPs. In Proposition F.2, we gave full characterization for LRPs in one-dimensional quadratic games. In fact, from our spectral analysis in Section 5 one can draw the following conclusion:

**Proposition F.15 (local stability at LRP)** *Suppose  $c^2 \neq ab$ . For one-dimensional homogeneous quadratic games  $q(x, y) = ax^2/2 + cxy + by^2/2$ , the stable sets of GDA (with momentum) and EG/OGD are within the set of LRPs. Moreover:*

- *There exists a quadratic game and an LRP,  $\mathbf{z}^*$ , such that no hyper-parameter choice can allow 2TS-EG to converge to  $\mathbf{z}^*$ .*
- *Whenever a LRP exists, there always exists a hyper-parameter choice  $(\alpha_1, \alpha_2, k)$  such that 2TS-OGD converges to the LRP.*

**Proof Part I** From stationarity the set of LRPs is  $\{(0, 0)\}$  if  $c^2 > ab$  and empty if  $c^2 < ab$ . The stable sets of gradient algorithms can only be empty or  $\{(0, 0)\}$ . We note that for  $q(x, y) = ax^2/2 + cxy + by^2/2$ , the characteristic polynomial of  $\mathbf{H}_{\alpha_1, \alpha_2}$  is:

$$\lambda^2 + (\alpha_1 a - \alpha_2 b)\lambda + \alpha_1 \alpha_2 (c^2 - ab) = 0. \quad (\text{F.27})$$

It is necessary that  $c^2 - ab \geq 0$  since from our spectral characterization, the two roots are either 1) both complex and are conjugate to each other; 2) both real and negative. If  $c = 0$ , we must have  $a \geq 0 \geq b$  since the two roots are both real and must be non-positive. Comparing with Proposition F.2 we have the first conclusion.

**Part II** Let us show the claim for EG. Take  $q(x, y) = -x^2 + xy + y^2/2$ . From (F.27) and Theorem 5.3, it suffices to show that:

$$p(\lambda) := \lambda^2 - (2\alpha_1 + \alpha_2)\lambda + 3\alpha_1\alpha_2 = 0 \quad (\text{F.28})$$

has no solution in the region  $\{\lambda \in \mathbb{C} : \Re(\lambda + \lambda^2) < 0\}$ . If  $(2\alpha_1 + \alpha_2)^2 \geq 12\alpha_1\alpha_2$ , it suffices to show that  $p(\lambda)$  has no root between  $-1$  and  $0$ . Otherwise, the condition  $\Re(\lambda + \lambda^2) < 0$  becomes

$$2\alpha_1 + \alpha_2 + (2\alpha_1 + \alpha_2)^2 < 6\alpha_1\alpha_2,$$

which cannot be true since  $(2\alpha_1 + \alpha_2)^2 \geq 8\alpha_1\alpha_2$  and  $\alpha_1 > 0, \alpha_2 > 0$ .

**Part III** For the claim of OGD, if  $c = 0$  then  $a > 0 > b$  and it is easy. If  $c \neq 0$ , combining (F.27) and (E.18), it suffices to show the existence of  $(\alpha_1, \alpha_2) \in \mathbb{R}_{++}$  such that

$$(\alpha_1 a - \alpha_2 b)^2 < 4\alpha_1\alpha_2(c^2 - ab) < 4, \quad \alpha_1 a - \alpha_2 b > -2\alpha_1\alpha_2(c^2 - ab), \quad (\text{F.29})$$

which, with  $\gamma = \alpha_2/\alpha_1$ , reduces to the existence of  $(\alpha_2, \gamma) \in \mathbb{R}_{++}$  such that

$$\frac{\gamma b - a}{2(c^2 - ab)} < \alpha_2, \quad \alpha_2^2 < \frac{\gamma}{c^2 - ab}, \quad (a - \gamma b)^2 < 4\gamma(c^2 - ab), \quad (\text{F.30})$$

which reduces to the existence of  $\gamma \in \mathbb{R}_{++}$  such that

$$(a - \gamma b)^2 < 4\gamma(c^2 - ab). \quad (\text{F.31})$$



this is always true no matter whether  $b = 0$  or  $b \neq 0$ . ■

This proposition shows the essential difference between EG and OGD in the convergence to LRPs. The last claim shares the same spirit with Jin et al. (2020, Theorem 28), since we can similarly write:

$$\mathcal{LRP} = 2\mathcal{T}\mathcal{S}\text{-OGD}, \quad (\text{F.32})$$

where  $\mathcal{LRP}$  is the set of LRPs and  $2\mathcal{T}\mathcal{S}\text{-OGD}$  is the set of all possible stable points of 2TS-OGD given some parameters ( $\alpha_1 > 0, \alpha_2 > 0, k > 1$ ).

However, this result does not hold in higher dimensions. We can prove the following:

**Proposition F.16 (failure of gradient algorithms at LRP)** *There exists a two-dimensional quadratic function  $q(\mathbf{x}, \mathbf{y})$  with its LRP at  $(\mathbf{0}, \mathbf{0})$ , in the same setting as Theorem F.14, such that GD (with momentum), EG or OGD cannot converge to the LRP for any hyper-parameter choice.*

**Proof** Combined with what we have in Proposition F.15 and Proposition 5.9, it suffices to prove the negative result for OGD. Since local robust points include both local minimax points and local maximin points, we construct a two-dimensional quadratic function that include both cases:

$$q(\mathbf{x}, \mathbf{y}) = -x_1^2 + x_1y_1 + x_2y_2 + y_2^2. \quad (\text{F.33})$$

Note that  $(\mathbf{0}, \mathbf{0})$  is the only stationary point. We now prove that it is also a local robust point. Writing the quadratic function in the same form as (4.1), we have:

$$\mathbf{A} = \begin{bmatrix} -2 & 0 \\ 0 & 0 \end{bmatrix}, \mathbf{B} = \begin{bmatrix} 0 & 0 \\ 0 & 2 \end{bmatrix}, \mathbf{C} = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}. \quad (\text{F.34})$$

From Definition F.12, we obtain the positive and the negative parts of  $\mathbf{A}$  and  $\mathbf{B}$ :

$$\mathbf{A}_p = \mathbf{0}, \mathbf{A}_n = \mathbf{A}, \mathbf{B}_p = \mathbf{B}, \mathbf{B}_n = \mathbf{0}, \quad (\text{F.35})$$

and thus  $\mathbf{P}_{\mathbf{B}_n}^\perp = \mathbf{P}_{\mathbf{A}_p}^\perp = \mathbf{I}$ . In (F.22) and (F.23), one can write  $\mathbf{L} = \mathbf{M} = \mathbf{I}$  and  $\mathbf{P}_{\mathbf{L}}^\perp = \mathbf{P}_{\mathbf{M}}^\perp = \mathbf{0}$ . It thus follows that (F.22) and (F.23) hold and  $(\mathbf{0}, \mathbf{0})$  is a LRP.

We now analyze the local convergence of OGD. The Jacobian of  $\mathbf{v}(\mathbf{z})$  is a constant:

$$\mathbf{H}_{\alpha_1, \alpha_2}(q) = \begin{bmatrix} -\alpha_1 \mathbf{A} & -\alpha_1 \mathbf{C} \\ \alpha_2 \mathbf{C}^\top & \alpha_2 \mathbf{B} \end{bmatrix}. \quad (\text{F.36})$$

Note that  $\mathbf{C}^\top$  and  $\mathbf{B}$  are diagonal matrices and thus they commute. So, we can compute the characteristic equation of  $\mathbf{H}_{\alpha_1, \alpha_2}(q)$  as:

$$\det((\lambda \mathbf{I} + \alpha_1 \mathbf{A})(\lambda \mathbf{I} - \alpha_2 \mathbf{B}) + \alpha_1 \alpha_2 \mathbf{C} \mathbf{C}^\top) = 0, \quad (\text{F.37})$$

from which we obtain:

$$\lambda(\lambda - 2\alpha_1) + \alpha_1 \alpha_2 = 0, \quad (\text{F.38})$$

$$\lambda(\lambda - 2\alpha_2) + \alpha_1 \alpha_2 = 0. \quad (\text{F.39})$$

For 2TS-OGD, when  $k \rightarrow 1_+$  the algorithm is the most stable (Theorem 5.3), where the condition should be (Theorem 5.2, (E.18)):

$$|\lambda| < 1, |\lambda - 1/2| > 1/2. \tag{F.40}$$

Now we separate the discussion into two cases: if  $\alpha_1 \geq \alpha_2 > 0$ , then (F.38) gives:

$$\lambda_{1,2} = \alpha_1 \pm \sqrt{\alpha_1^2 - \alpha_1 \alpha_2}, \tag{F.41}$$

and there exists a real and positive root. Similarly, if  $\alpha_2 \geq \alpha_1 > 0$ , (F.39) has a real and positive root. In either case (F.40) would be violated. ■

From the proof, we can see that the problem lies in the coordinate-independent step sizes. In fact, (F.33) could be rewritten as:

$$q(\mathbf{x}, \mathbf{y}) = q_1(x_1, y_1) + q_2(x_2, y_2), q_1(x, y) := -x^2 + xy, q_2(x, y) := xy + y^2. \tag{F.42}$$

For the function  $q_1$ ,  $(0, 0)$  is a local minimax point, and the stability constraint for 2TS-OGD is (with  $k \rightarrow 1_+$ , see (E.19)):

$$\alpha_1 < 1, 1 < \alpha_2 < 1/\alpha_1. \tag{F.43}$$

While for the function  $q_2$ ,  $(0, 0)$  is a local maximin point, and the stability constraint for 2TS-OGD is (in a similar way):

$$\alpha_2 < 1, 1 < \alpha_1 < 1/\alpha_2. \tag{F.44}$$

(F.43) and (F.44) are conflicting each other. Therefore, it tells us that coordinate-dependent step sizes might be necessary in order for stability near a LRP, such as those in Adam (Kingma and Ba, 2015), which is widely used in GAN training.

We finally mention that LRPs are a wider class that could include the stable points of gradient algorithms. For example, in the proof of Prop. 27 of Jin et al. (2020), there is a two-dimensional quadratic function that has  $(0, 0)$  as a stable solution of simultaneous GDA, but it is neither local maximin or minimax. It can be shown that it is in fact a local robust point.

## References

- K. Arrow, L. Hurwicz, and H. Uzawa. *Studies in linear and non-linear programming*. Stanford University Press, 1958.
- W. Azizian, I. Mitliagkas, S. Lacoste-Julien, and G. Gidel. A tight and unified analysis of extragradient for a whole spectrum of differentiable games. In *the 23rd International Conference on Artificial Intelligence and Statistics*, 2020a.
- W. Azizian, D. Scieur, I. Mitliagkas, S. Lacoste-Julien, and G. Gidel. Accelerating smooth games by manipulating spectral shapes. In *the 23rd International Conference on Artificial Intelligence and Statistics*, 2020b.

- B. Barazandeh and M. Razaviyayn. Solving non-convex non-differentiable min-max games using proximal gradient method. In *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 3162–3166. IEEE, 2020.
- S. Basu, R. Pollack, and M.-F. Roy. *Algorithms in real algebraic geometry*. Springer, 2005.
- A. Ben-Tal and J. Zowe. Necessary and sufficient optimality conditions for a class of nonsmooth minimization problems. *Mathematical Programming*, 24(1):70–91, 1982.
- A. Ben-Tal and J. Zowe. Directional derivatives in nonsmooth optimization. *Journal of Optimization Theory and Applications*, 47(4):483–490, 1985.
- H. Berard, G. Gidel, A. Almahairi, P. Vincent, and S. Lacoste-Julien. A closer look at the optimization landscapes of generative adversarial networks. In *International Conference on Learning Representations*, 2020. URL <https://openreview.net/forum?id=HJeVnCEkWH>.
- D. P. Bertsekas. Nonlinear programming. *Journal of the Operational Research Society*, 48(3):334–334, 1997.
- R. Bollapragada, D. Scieur, and A. d’Aspremont. Nonlinear acceleration of primal-dual algorithms. In *the 22nd International Conference on Artificial Intelligence and Statistics*, pages 739–747, 2019.
- F. H. Clarke. *Optimization and Nonsmooth Analysis*. SIAM, 1990.
- R. Cominetti and R. Correa. A Generalized Second-Order Derivative in Nonsmooth Optimization. *SIAM Journal on Control and Optimization*, 28(4):789–809, 1990.
- J. M. Danskin. The Theory of Max-Min, with Applications. *SIAM Journal on Applied Mathematics*, 14(4):641–664, 1966.
- C. Daskalakis and I. Panageas. The limit points of (optimistic) gradient descent in min-max optimization. In *Advances in Neural Information Processing Systems*, pages 9236–9246, 2018.
- C. Daskalakis, A. Ilyas, V. Syrgkanis, and H. Zeng. Training GANs with optimism. In *the 6th International Conference on Learning Representations*, 2018.
- V. F. Dem’yanov. On the solution of several minimax problems. I. *Cybernetics*, 2:47–53, 1966.
- V. F. Dem’yanov. Sufficient conditions for a local minimax. *USSR Computational Mathematics and Mathematical Physics*, 10(5):53–63, 1970.
- V. F. Dem’yanov. Second-order directional derivatives of a function of the maximum. *Cybernetics*, 9:797–800, 1973.
- V. F. Dem’yanov and V. N. Malozemov. *Introduction to Minimax*. Wiley, 1974.
- F. Facchinei and J.-S. Pang. *Finite-dimensional variational inequalities and complementarity problems*. Springer Science & Business Media, 2007.

- K. Fan. On a theorem of weyl concerning eigenvalues of linear transformations: II. *Proceedings of the National Academy of Sciences of the United States of America*, 36(1):31, 1950.
- F. Farnia and A. Ozdaglar. Do GANs always have Nash equilibria? In *International Conference on Machine Learning*, pages 3029–3039. PMLR, 2020.
- T. Fiez, B. Chasnov, and L. J. Ratliff. Convergence of learning dynamics in Stackelberg games. *arXiv*, 2019. arXiv:1906.01217.
- G. Gidel, R. A. Hemmat, M. Pezeshki, G. Huang, R. Lepriol, S. Lacoste-Julien, and I. Mitliagkas. Negative momentum for improved game dynamics. In *the 22nd International Conference on Artificial Intelligence and Statistics*, 2019.
- E. G. Golshtein. A generalized gradient method for finding saddlepoints. *Ekonomika i matematicheskie*, 8(4):36–52, 1972.
- I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio. Generative adversarial nets. In *Advances in neural information processing systems*, pages 2672–2680, 2014.
- F. R. Hampel. The influence curve and its role in robust estimation. *Journal of the american statistical association*, 69(346):383–393, 1974.
- M. Heusel, H. Ramsauer, T. Unterthiner, B. Nessler, and S. Hochreiter. GANs trained by a two time-scale update rule converge to a local Nash equilibrium. In *Advances in neural information processing systems*, pages 6626–6637, 2017.
- J.-B. Hiriart-Urruty and C. Lemaréchal. *Fundamentals of convex analysis*. Springer Science & Business Media, 2004.
- J.-B. Hiriart-Urruty and C. Lemaréchal. *Convex analysis and minimization algorithms I: Fundamentals*, volume 305. Springer, 2013.
- Y.-G. Hsieh, F. Iutzeler, J. Malick, and P. Mertikopoulos. On the convergence of single-call stochastic extra-gradient methods. In *NeurIPS*, pages 6936–6946, 2019.
- Y.-G. Hsieh, F. Iutzeler, J. Malick, and P. Mertikopoulos. Explore aggressively, update conservatively: Stochastic extragradient methods with variable stepsize scaling. In *NeurIPS 2020-34th Conference on Neural Information Processing Systems*, 2020.
- A. Ibrahim, W. Azizian, G. Gidel, and I. Mitliagkas. Linear lower bounds and conditioning of differentiable games. In *International conference on machine learning*, pages 6356–6366, 2020.
- C. Jin, P. Netrapalli, and M. Jordan. What is local optimality in nonconvex-nonconcave minimax optimization? In *International conference on machine learning*, pages 5735–5744, 2020.
- A. Katok and B. Hasselblatt. *Introduction to the modern theory of dynamical systems*, volume 54. Cambridge university press, 1995.

- H. Kawasaki. The upper and lower second order directional derivatives of a sup-type function. *Mathematical Programming*, 41(1-3):327–339, 1988.
- H. Kawasaki. Second order necessary optimality conditions for minimizing a sup-type function. *Mathematical programming*, 49(1-3):213–229, 1991.
- H. Kawasaki. Second-order necessary and sufficient optimality conditions for minimizing a sup-type function. *Applied Mathematics and Optimization*, 26(2):195–220, 1992.
- D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. In *International Conference on Learning Representations*, 2015.
- G. Korpelevich. The extragradient method for finding saddle points and other problems. *Matecon*, 12:747–756, 1976.
- T. Lin, C. Jin, and M. I. Jordan. Near-optimal algorithms for minimax optimization. In *the 33rd Conference on Learning Theory*, 2020.
- S. Liu, S. Lu, X. Chen, Y. Feng, K. Xu, A. Al-Dujaili, M. Hong, and U.-M. O’Reilly. Min-max optimization without gradients: Convergence and applications to black-box evasion and poisoning attacks. In *International conference on machine learning*, pages 2307–2318, 2020.
- A. Madry, A. Makelov, L. Schmidt, D. Tsipras, and A. Vladu. Towards deep learning models resistant to adversarial attacks. In *the 6th International Conference on Learning Representations*, 2018.
- A. W. Marshall, I. Olkin, and B. C. Arnold. *Inequalities: theory of majorization and its applications*, volume 143. Springer, 1979.
- P. Mertikopoulos, C. Papadimitriou, and G. Piliouras. Cycles in adversarial regularized learning. In *Proceedings of the Twenty-Ninth Annual ACM-SIAM Symposium on Discrete Algorithms*, pages 2703–2717, 2018.
- P. Mertikopoulos, B. Lecouat, H. Zenati, C.-S. Foo, V. Chandrasekhar, and G. Piliouras. Optimistic mirror descent in saddle-point problems: Going the extra (gradient) mile. In *the 7th International Conference on Learning Representations*, 2019.
- L. Mescheder, S. Nowozin, and A. Geiger. The numerics of GANs. In *Advances in Neural Information Processing Systems*, pages 1825–1835, 2017.
- A. Mokhtari, A. Ozdaglar, and S. Pattathil. Proximal point approximations achieving a convergence rate of  $o(1/k)$  for smooth convex-concave saddle point problems: Optimistic gradient and extra-gradient methods. *arXiv:1906.01115*, 2019.
- K. G. Murty and S. N. Kabadi. Some np-complete problems in quadratic and nonlinear programming. *Mathematical programming*, 39(2):117–129, 1987.
- J. F. Nash. Equilibrium points in  $n$ -person games. *Proceedings of the national academy of sciences*, 36(1):48–49, 1950.

- A. S. Nemirovsky and D. B. Yudin. *Problem complexity and method efficiency in optimization*. Wiley, 1983.
- Y. Nesterov. A method for unconstrained convex minimization problem with the rate of convergence  $o(1/k^2)$ . *Doklady AN USSR*, 269:543–547, 1983.
- W. Niethammer and R. S. Varga. The analysis of  $k$ -step iterative methods for linear systems from summability theory. *Numerische Mathematik*, 41(2):177–206, 1983.
- W. Peng, Y.-H. Dai, H. Zhang, and L. Cheng. Training GANs with centripetal acceleration. *Optimization Methods and Software*, 35(5):955–973, 2020.
- B. Polyak. *Introduction to Optimization*. Optimization Software Inc., 1987.
- B. T. Polyak. Some methods of speeding up the convergence of iteration methods. *USSR Computational Mathematics and Mathematical Physics*, 4(5):1–17, 1964.
- L. D. Popov. A modification of the Arrow–Hurwicz method for search of saddle points. *Mathematical Notes*, 28(5):845–848, 1980.
- M. Razaviyayn, T. Huang, S. Lu, M. Nouiehed, M. Sanjabi, and M. Hong. Nonconvex min-max optimization: Applications, challenges, and recent theoretical advances. *IEEE Signal Processing Magazine*, 37(5):55–66, 2020.
- R. T. Rockafellar and R. J.-B. Wets. *Variational analysis*, volume 317. Springer Science & Business Media, 2009.
- F. Schaefer, H. Zheng, and A. Anandkumar. Implicit competitive regularization in GANs. In *International Conference on Machine Learning*, pages 8533–8544. PMLR, 2020.
- I. Schur. Über potenzreihen, die im innern des einheitskreises beschränkt sind. *Journal für die reine und angewandte Mathematik*, 147:205–232, 1917.
- A. Seeger. Second order directional derivatives in parametric optimization problems. *Mathematics of Operations Research*, 13(1):124–139, 1988.
- A. Sinha, H. Namkoong, and J. Duchi. Certifying some distributional robustness with principled adversarial training. In *International Conference on Learning Representations*, 2018.
- M. Sion et al. On general minimax theorems. *Pacific Journal of mathematics*, 8(1):171–176, 1958.
- I. Sutskever, J. Martens, G. Dahl, and G. Hinton. On the importance of initialization and momentum in deep learning. In *International conference on machine learning*, pages 1139–1147, 2013.
- R. S. Sutton, A. G. Barto, et al. *Introduction to reinforcement learning*, volume 135. MIT press Cambridge, 1998.

- J. von Neumann. Zur theorie der gesellschaftsspiele. *Mathematische annalen*, 100(1):295–320, 1928.
- H. von Stackelberg. *Market structure and equilibrium*. Springer, 1934.
- Y. Wang, G. Zhang, and J. Ba. On solving minimax optimization locally: A follow-the-ridge approach. In *the 8th International Conference on Learning Representations*, 2020.
- G. Zhang and Y. Yu. Convergence of gradient methods on bilinear zero-sum games. In *the 8th International Conference on Learning Representations*, 2020.
- G. Zhang, P. Poupart, and Y. Yu. Optimality and stability in non-convex smooth games. arXiv:2002.11875, 2020.
- G. Zhang, K. Wu, P. Poupart, and Y. Yu. Newton-type methods for minimax optimization. In *ICML workshop on Beyond First-Order Methods in ML Systems*, 2021. arXiv:2006.14592.
- J. Zhang, M. Hong, and S. Zhang. On Lower Iteration Complexity Bounds for the Saddle Point Problems. arXiv:1912.07481, 2019.