

596 A Omitted Proofs Section 2

597 **Proposition A.1.** *Given any convex-concave min-max game with dependent strategy sets (X, Y, f, g) ,*
 598 *a Stackelberg equilibrium always exists.*

599 *Proof of Proposition A.1* By Berge's maximum theorem [5], the outer player's value function
 600 $V(x) = \max_{y \in Y: g(x, y) \geq 0} f(x, y)$ is continuous, and the inner solution correspondence $Y^*(x) =$
 601 $\arg \max_{y \in \mathcal{Y}(x)} f(x, y)$ is non-empty, for all $x \in X$. Since V is continuous and X is compact and
 602 non-empty, by the extreme value theorem [50], there exists a minimizer x^* of V . Hence $(x^*, y^*(x^*))$
 603 where $y^*(x^*) \in Y^*(x^*)$ is well-defined and is a Stackelberg equilibrium of (X, Y, f, g) . \square

604 B Envelope Theorem

605 Danskin's theorem [15] offers insights into optimization problems of the form:

$$\max_{y \in Y} f(x, y), \quad (4)$$

606 where $Y \subset \mathbb{R}^m$ is compact and non-empty. Among other things, Danskin's theorem allows us to
 607 compute the gradient of the objective function of this optimization problem with respect to x .

608 **Theorem B.1** (Danskin's Theorem). *Consider Equation (4). Suppose that Y is convex and that*
 609 *f is concave in y . Let $V(x) = \max_{y \in Y} f(x, y)$ and $Y^*(x) = \arg \max_{y \in Y} f(x, y)$. Then, V is*
 610 *differentiable at \hat{x} if $Y^*(\hat{x})$ is a singleton. Additionally, the gradient at \hat{x} is given by $V'(\hat{x}) =$*
 611 *$\nabla_x f(\hat{x}, y^*(\hat{x}))$, where $y^*(\hat{x}) \in Y^*(\hat{x})$.*

612 Unfortunately, Danskin's theorem does not hold when Y is replaced by even a non-empty compact-
 613 valued correspondence $\mathcal{Y} : X \rightrightarrows Y$, in which case the inner problem becomes $\max_{y \in \mathcal{Y}(x)} f(x, y)$.

614 **Example B.2** (Danskin's theorem does not apply to min-max games with dependent strategy sets).
 615 *Consider the optimization problem:*

$$\max_{y \in \mathbb{R}: y+x \geq 0} -y^2 + y + 2x + 2 \quad (5)$$

616 *The solution function $y^*(x) = \arg \max_{y \in \mathbb{R}: y+x \geq 0} -y^2 + y + 2x + 2$ for this problem is well defined*
 617 *since the solution is singleton-valued and is given by:*

$$y^*(x) = \begin{cases} 1/2 & \text{if } x \geq -1/2 \\ -x & \text{if } x < -1/2 \end{cases} \quad (6)$$

618 *The value function $V(x) = \max_{y \in \mathbb{R}: y+x \geq 0} -y^2 + y + 2x + 2$ is given by:*

$$V(x) = f(x, y^*(x)) \quad (7)$$

$$= -y^*(x)^2 + y^*(x) + 2x + 2 \quad (8)$$

$$= \begin{cases} -1/4 + 1/2 + 2x + 2 & \text{if } x \geq -1/2 \\ -x^2 - x + 2x + 2 & \text{if } x < -1/2 \end{cases} \quad (9)$$

$$= \begin{cases} 9/4 + 2x & \text{if } x \geq -1/2 \\ -x^2 + x + 2 & \text{if } x < -1/2 \end{cases} \quad (10)$$

619 *The derivative of the value function is given by:*

$$\frac{\partial V}{\partial x} = \begin{cases} 2 & \text{if } x \geq -1/2 \\ 1 - 2x & \text{if } x < -1/2 \end{cases} \quad (11)$$

620 *However, the derivative predicted by Danskin's theorem is 2, for all x . Hence, Danskin's theorem does not*
 621 *hold when the constraints are parameterized, i.e., when the problem is of the form $\min_{y \in \mathcal{Y}(x)} f(x, y)$*
 622 *rather than $\min_{y \in Y} f(x, y)$ where $X \subset \mathbb{R}^n$, $Y \subset \mathbb{R}^m$, and $\mathcal{Y} : X \rightrightarrows Y$.*

623 **N.B.** *For simplicity, we do not assume the constraint set is compact in this example; however, the*
 624 *conclusion still applies, since compactness of the constraint set is used to guarantee existence of a*
 625 *solution for all x , but as a solution to this particular problem exists we can do away with the assumption.*

626 An answer to Danskin's theorem not holding when the constraints are parameterized can be found
 627 in the mathematical economics literature. In particular the following theorem due to Milgrom and
 628 Segal [41] generalizes Danskin's theorem (Theorem B.1).

629 **Theorem B.3** (Envelope Theorem [41]). *Consider the maximization problem*

$$V(\mathbf{x}) = \max_{\mathbf{y} \in \mathbb{R}^m} f(\mathbf{x}, \mathbf{y}) \text{ subject to } g_k(\mathbf{x}, \mathbf{y}) \geq 0 \text{ for all } k = 1, \dots, K . \quad (12)$$

630 Define the solution correspondence $Y^*(\mathbf{x}) = \arg \max_{\mathbf{y} \in \mathbb{R}^m: g(\mathbf{x}, \mathbf{y}) \geq 0} f(\mathbf{x}, \mathbf{y})$. Now suppose that
 631 Assumption 3.1 holds. Then, the value function V is absolutely continuous, and at any point $\hat{\mathbf{x}}$ where
 632 it is differentiable:

$$\nabla_{\mathbf{x}} V(\hat{\mathbf{x}}) = \nabla_{\mathbf{x}} L(\mathbf{y}^*(\hat{\mathbf{x}}), \boldsymbol{\lambda}(\hat{\mathbf{x}}, \mathbf{y}^*(\hat{\mathbf{x}})), \hat{\mathbf{x}}) = \nabla_{\mathbf{x}} f(\hat{\mathbf{x}}, \mathbf{y}^*(\hat{\mathbf{x}})) + \sum_{k=1}^K \lambda_k(\hat{\mathbf{x}}, \mathbf{y}^*(\hat{\mathbf{x}})) \nabla_{\mathbf{x}} g_k(\hat{\mathbf{x}}, \mathbf{y}^*(\hat{\mathbf{x}})) , \quad (13)$$

633 where $\boldsymbol{\lambda}(\hat{\mathbf{x}}, \mathbf{y}^*(\hat{\mathbf{x}})) = (\lambda_1(\hat{\mathbf{x}}, \mathbf{y}^*(\hat{\mathbf{x}})), \dots, \lambda_K(\hat{\mathbf{x}}, \mathbf{y}^*(\hat{\mathbf{x}})))^T \in \Lambda(\hat{\mathbf{x}}, \mathbf{y}^*(\hat{\mathbf{x}}))$ are the Lagrange multi-
 634 pliers associated with $\mathbf{y}^*(\hat{\mathbf{x}}) \in Y^*(\hat{\mathbf{x}})$.

635 C Omitted Proofs Section 3

636 *Proof of Theorem 3.2.* Let $V(\mathbf{x}) = \max_{\mathbf{y} \in Y: g(\mathbf{x}, \mathbf{y}) \geq 0} f(\mathbf{x}, \mathbf{y})$. Reformulating the problem as a La-
 637 grangian saddle point problem, for all $\hat{\mathbf{x}} \in X$, it holds that:

$$V(\hat{\mathbf{x}}) = \max_{\mathbf{y} \in Y: g(\hat{\mathbf{x}}, \mathbf{y}) \geq 0} f(\hat{\mathbf{x}}, \mathbf{y}) \quad (14)$$

$$= \max_{\mathbf{y} \in Y} \min_{\boldsymbol{\lambda} \in \mathbb{R}_{++}^K} \left\{ f(\hat{\mathbf{x}}, \mathbf{y}) + \sum_{k=1}^K \lambda_k g_k(\hat{\mathbf{x}}, \mathbf{y}) \right\} \quad (15)$$

638 Since an interior point exists by the assumptions, the Karush-Kuhn-Tucker Theorem [36] applies,
 639 so for all $\hat{\mathbf{x}} \in X$, there exists $\boldsymbol{\lambda} \in \mathbb{R}^K$ that solves the above optimization problem Equation (15).

640 Let $Y^*(\hat{\mathbf{x}}) = \arg \max_{\mathbf{y} \in Y: g(\hat{\mathbf{x}}, \mathbf{y}) \geq 0} f(\hat{\mathbf{x}}, \mathbf{y})$ and $\Lambda(\hat{\mathbf{x}}, \mathbf{y}) =$
 641 $\arg \min_{\boldsymbol{\lambda} \in \mathbb{R}_{++}^K} \left\{ f(\hat{\mathbf{x}}, \mathbf{y}) + \sum_{k=1}^K \lambda_k g_k(\hat{\mathbf{x}}, \mathbf{y}) \right\}$. We can then re-express the value function
 642 as:

$$V(\hat{\mathbf{x}}) = f(\hat{\mathbf{x}}, \mathbf{y}^*(\hat{\mathbf{x}})) + \sum_{k=1}^K \lambda_k(\hat{\mathbf{x}}, \mathbf{y}^*(\hat{\mathbf{x}})) g_k(\hat{\mathbf{x}}, \mathbf{y}^*(\hat{\mathbf{x}})), \quad \forall \mathbf{y}^*(\hat{\mathbf{x}}) \in Y^*(\hat{\mathbf{x}}), \lambda_k(\hat{\mathbf{x}}, \mathbf{y}^*(\hat{\mathbf{x}})) \in \Lambda(\hat{\mathbf{x}}, \mathbf{y}^*(\hat{\mathbf{x}})) .$$

643 Alternatively, we can take the maximum over $\boldsymbol{\lambda}$'s and \mathbf{y} 's to obtain:

$$V(\hat{\mathbf{x}}) = \max_{\mathbf{y}^*(\hat{\mathbf{x}}) \in Y^*(\hat{\mathbf{x}})} \max_{\lambda_k(\hat{\mathbf{x}}, \mathbf{y}^*(\hat{\mathbf{x}})) \in \Lambda(\hat{\mathbf{x}}, \mathbf{y}^*(\hat{\mathbf{x}}))} \left\{ f(\hat{\mathbf{x}}, \mathbf{y}) + \sum_{k=1}^K \lambda_k(\hat{\mathbf{x}}, \mathbf{y}^*(\hat{\mathbf{x}})) g_k(\hat{\mathbf{x}}, \mathbf{y}) \right\} .$$

644 Note that for fixed $\mathbf{y}^*(\hat{\mathbf{x}}) \in Y^*(\hat{\mathbf{x}})$ and corresponding $\lambda_k(\hat{\mathbf{x}}, \mathbf{y}^*(\hat{\mathbf{x}})) \in \Lambda(\hat{\mathbf{x}}, \mathbf{y}^*(\hat{\mathbf{x}}))$,
 645 $f(\hat{\mathbf{x}}, \mathbf{y}) + \sum_{k=1}^K \lambda_k(\hat{\mathbf{x}}, \mathbf{y}^*(\hat{\mathbf{x}})) g_k(\hat{\mathbf{x}}, \mathbf{y})$ is differentiable, since f, g_1, \dots, g_K are differentiable.

646 Additionally, recall the pointwise maximum subdifferential property, i.e., if $f(\mathbf{x}) = \max_{\alpha \in \mathcal{A}} f_{\alpha}(\mathbf{x})$
 647 for a family of functions $\{f_{\alpha}\}_{\alpha \in \mathcal{A}}$, then $\partial_{\mathbf{x}} f(\mathbf{a}) = \text{conv} \left(\bigcup_{\alpha \in \mathcal{A}} \{\partial_{\mathbf{x}} f_{\alpha}(\mathbf{a}) \mid f_{\alpha}(\mathbf{a}) = f(\mathbf{a})\} \right)$

648 (see, for example, [7]), which then gives:

$$\partial_{\mathbf{x}} V(\hat{\mathbf{x}}) = \partial_{\mathbf{x}} \left(\max_{\mathbf{y}^*(\hat{\mathbf{x}}) \in Y^*(\hat{\mathbf{x}})} \max_{\lambda_k(\hat{\mathbf{x}}, \mathbf{y}^*(\hat{\mathbf{x}})) \in \Lambda(\hat{\mathbf{x}}, \mathbf{y}^*(\hat{\mathbf{x}}))} \left\{ f(\hat{\mathbf{x}}, \mathbf{y}^*(\hat{\mathbf{x}})) + \sum_{k=1}^K \lambda_k(\hat{\mathbf{x}}, \mathbf{y}^*(\hat{\mathbf{x}})) g_k(\hat{\mathbf{x}}, \mathbf{y}^*(\hat{\mathbf{x}})) \right\} \right) \quad (16)$$

$$= \text{conv} \left(\bigcup_{\mathbf{y}^*(\hat{\mathbf{x}}) \in Y^*(\hat{\mathbf{x}})} \bigcup_{\lambda_k(\hat{\mathbf{x}}, \mathbf{y}^*(\hat{\mathbf{x}})) \in \Lambda(\hat{\mathbf{x}}, \mathbf{y}^*(\hat{\mathbf{x}}))} \partial_{\mathbf{x}} \left\{ f(\hat{\mathbf{x}}, \mathbf{y}^*(\hat{\mathbf{x}})) + \sum_{k=1}^K \lambda_k(\hat{\mathbf{x}}, \mathbf{y}^*(\hat{\mathbf{x}})) g_k(\hat{\mathbf{x}}, \mathbf{y}^*(\hat{\mathbf{x}})) \right\} \right) \quad (17)$$

$$= \text{conv} \left(\bigcup_{\mathbf{y}^*(\hat{\mathbf{x}}) \in Y^*(\hat{\mathbf{x}})} \bigcup_{\lambda_k(\hat{\mathbf{x}}, \mathbf{y}^*(\hat{\mathbf{x}})) \in \Lambda(\hat{\mathbf{x}}, \mathbf{y}^*(\hat{\mathbf{x}}))} \left\{ \nabla_{\mathbf{x}} f(\hat{\mathbf{x}}, \mathbf{y}^*(\hat{\mathbf{x}})) + \sum_{k=1}^K \lambda_k(\hat{\mathbf{x}}, \mathbf{y}^*(\hat{\mathbf{x}})) \nabla_{\mathbf{x}} g_k(\hat{\mathbf{x}}, \mathbf{y}^*(\hat{\mathbf{x}})) \right\} \right) . \quad (18)$$

649 \square

650 D Algorithms

651 The algorithms studied in our paper and described in Section 3 are presented below. We note that
 652 Π_Y is the projection operator on the set Y which is defined as $\Pi_Y(\mathbf{y}) = \arg \min_{\mathbf{z} \in Y} \|\mathbf{y} - \mathbf{z}\|_2$.

Algorithm 1 Max-Oracle Gradient Descent

Inputs: $X, Y, f, g, \eta, T, \mathbf{x}^{(0)}$

Output: $(\mathbf{x}^*, \mathbf{y}^*)$

- 1: **for** $t = 1, \dots, T$ **do**
 - 2: Find $\hat{\mathbf{y}} \in Y$ such that $f(\mathbf{x}^{(t-1)}, \hat{\mathbf{y}}) \geq \max_{\mathbf{y} \in Y: g(\mathbf{x}^{(t-1)}, \mathbf{y}) \geq 0} f(\mathbf{x}^{(t-1)}, \mathbf{y}) - \delta$ and $g(\mathbf{x}^{(t-1)}, \hat{\mathbf{y}}) \geq 0$
 - 3: Set $\mathbf{y}^{(t-1)} = \hat{\mathbf{y}}$
 - 4: Set $\boldsymbol{\lambda}^{(t-1)} = \boldsymbol{\lambda}(\mathbf{x}^{(t-1)}, \mathbf{y}^{(t-1)})$
 - 5: Set $\mathbf{x}^{(t)} = \Pi_X \left(\mathbf{x}^{(t-1)} - \eta_t \left[\nabla_{\mathbf{x}} f(\mathbf{x}^{(t-1)}, \mathbf{y}^{(t-1)}) + \sum_{k=1}^K \lambda_k^{(t-1)} \nabla_{\mathbf{x}} g_k(\mathbf{x}^{(t-1)}, \mathbf{y}^{(t-1)}) \right] \right)$
 - 6: **end for**
 - 7: Find $\hat{\mathbf{y}} \in Y$ such that $f(\mathbf{x}^{(T)}, \hat{\mathbf{y}}) \geq \max_{\mathbf{y} \in Y: g(\mathbf{x}^{(T)}, \mathbf{y}) \geq 0} f(\mathbf{x}^{(T)}, \mathbf{y}) - \delta$ and $g(\mathbf{x}^{(T)}, \hat{\mathbf{y}}) \geq 0$
 - 8: $\mathbf{y}^{(T)} = \hat{\mathbf{y}}$
 - 9: **return** $(\mathbf{x}^{(T)}, \mathbf{y}^{(T)})$
-

Algorithm 2 Nested Gradient Descent

Inputs: $X, Y, f, g, \eta_x, \eta_y, T_x, T_y, \mathbf{x}^{(0)}, \mathbf{y}^{(0)}$

Output: $\mathbf{x}^*, \mathbf{y}^*$

- 1: **for** $t = 1, \dots, T_x$ **do**
 - 2: $\mathbf{y}^{(t-1)} = \mathbf{y}^{(0)}$
 - 3: **for** $s = 1, \dots, T_y$ **do**
 - 4: $\mathbf{y}^{(t-1)} = \Pi_{\{\mathbf{y} \in Y: g(\mathbf{x}^{(t-1)}, \mathbf{y}) \geq 0\}} (\mathbf{y}^{(t-1)} + \eta_{s\mathbf{y}} \nabla_{\mathbf{y}} f(\mathbf{x}^{(t-1)}, \mathbf{y}^{(t-1)}))$
 - 5: **end for**
 - 6: Set $\boldsymbol{\lambda}^{(t-1)} = \boldsymbol{\lambda}(\mathbf{x}^{(t-1)}, \mathbf{y}^{(t-1)})$
 - 7: Set $\mathbf{x}^{(t)} = \Pi_X \left(\mathbf{x}^{(t-1)} - \eta_{t\mathbf{x}} \left[\nabla_{\mathbf{x}} f(\mathbf{x}^{(t-1)}, \mathbf{y}^{(t-1)}) + \sum_{k=1}^K \lambda_k^{(t-1)} \nabla_{\mathbf{x}} g_k(\mathbf{x}^{(t-1)}, \mathbf{y}^{(t-1)}) \right] \right)$
 - 8: **end for**
 - 9: $\mathbf{y}^{(T)} = \mathbf{y}^{(0)}$
 - 10: **for** $s = 1, \dots, T_y$ **do**
 - 11: $\mathbf{y}^{(T)} = \Pi_{\{\mathbf{y} \in Y: g(\mathbf{x}^{(T)}, \mathbf{y}) \geq 0\}} (\mathbf{y}^{(T)} + \eta_{s\mathbf{y}} \nabla_{\mathbf{y}} f(\mathbf{x}^{(T)}, \mathbf{y}^{(T)}))$
 - 12: **end for**
 - 13: **return** $(\mathbf{x}^{(T)}, \mathbf{y}^{(T)})$
-

653 D.1 Omitted Proofs Section 3

654 **Lemma D.1** (Lipschitz Objective, Lipschitz Value Function). *Let $f : X \times Y$ be a continuous function, where $X \subset \mathbb{R}^n$, $Y \subset \mathbb{R}^m$. Suppose that $\nabla_x f$ is continuous in (x, y) , X is compact and non-empty, and $\mathcal{Y} : X \rightrightarrows Y$ is nonempty-compact-valued correspondence, then*
655 *$V(x) = \max_{y \in \mathcal{Y}(x)} f(x, y)$ is ℓ_f -Lipschitz continuous, i.e., $\|V(x_1) - V(x_2)\| \leq \ell_f \|x_1 - x_2\|$,*
656 *with $\ell_f = \max_{(\hat{x}, \hat{y}) \in X \times Y} \|\nabla_x f(\hat{x}, \hat{y})\|$.*

659 *Proof of Lemma D.1* Let $\ell_f = \max_{(\hat{x}, \hat{y}) \in X \times Y} \|\nabla_x f(\hat{x}, \hat{y})\|$. Clearly, we have $\forall x_1, x_2 \in X, y \in \mathcal{Y}(x_1) \cap \mathcal{Y}(x_2)$, $\|f(x_1, y) - f(x_2, y)\| \leq \ell_f \|x_1 - x_2\|$.

661 Fix $x_1, x_2 \in X$. Then, for all $y \in \mathcal{Y}(x_1) \cap \mathcal{Y}(x_2)$, we have:

$$f(x_1, y) \leq f(x_1, y) - f(x_2, y) + f(x_2, y) \quad (19)$$

$$\leq \ell_f \|x_1 - x_2\| + f(x_2, y) \quad (20)$$

662 Taking the max over the y 's on both sides (which is guaranteed to exist by the continuity of f , and compactness and non-emptiness of \mathcal{Y}), we obtain:

$$\max_{y \in \mathcal{Y}(x_1)} f(x_1, y) \leq \ell_f \|x_1 - x_2\| + \max_{y \in \mathcal{Y}(x_2)} f(x_2, y) \quad (21)$$

$$V(x_1) \leq \ell_f \|x_1 - x_2\| + V(x_2) \quad (22)$$

$$V(x_1) - V(x_2) \leq \ell_f \|x_1 - x_2\| \quad (23)$$

664 Since this inequality holds for arbitrary $x_1, x_2 \in X$, we also have:

$$V(x_2) - V(x_1) \leq \ell_f \|x_1 - x_2\| \quad (24)$$

665 Combining the two inequalities, we obtain

$$\|V(x_1) - V(x_2)\| \leq \ell_f \|x_1 - x_2\| \quad (25)$$

666 \square

667 *Proof of Theorem 3.3.* Note that by Theorem 3.2 we have $\nabla_x f(x^{(t-1)}, y^{(t-1)}) + \sum_{k=1}^K \lambda_k^{(t-1)} \nabla_x g_k(x^{(t-1)}, y^{(t-1)}) \in \partial_x V(x^{(t-1)}) = \partial_x \max_{y \in Y: g(x^{(t-1)}, y) \geq 0} f(x^{(t-1)}, y)$.
668 For notational clarity, let $g(t-1) = \nabla_x f(x^{(t-1)}, y^{(t-1)}) + \sum_{k=1}^K \lambda_k^{(t-1)} \nabla_x g_k(x^{(t-1)}, y^{(t-1)})$.
669 Suppose that $x^* \in \arg \min_{x \in X} \max_{y \in Y: g(x, y)} f(x, y)$

$$\|x^{(T)} - x^*\|^2 = \|\Pi_X(x^{(T-1)} - \eta_T g(T-1)) - \Pi_X(x^*)\|^2 \quad (26)$$

$$\leq \|x^{(T-1)} - \eta_T g(T-1) - x^*\|^2 \quad (27)$$

$$= \|x^{(T-1)} - x^*\|^2 - 2\eta_T \langle g(T-1), x^{(T-1)} - x^* \rangle + \eta_T^2 \|g(T-1)\|^2 \quad (28)$$

$$\leq \|x^{(T-1)} - x^*\|^2 - 2\eta_T (f(x^{(T-1)}, y^{(T-1)}) - f(x^*, y^{(T-1)})) + \eta_T^2 \|g(T-1)\|^2 \quad (29)$$

671 where the first line follows from definitions, the second from the non-expansiveness of the projection operator, the third from algebra, the fourth from the definition of subgradients, i.e.,
672 $g(t-1)^T (x^{(t-1)} - x^*) \geq f(x^{(t-1)}, y^{(t-1)}) - f(x^*, y^{(t-1)})$. Applying the inequality above
673 recursively, we obtain:

$$\|x^{(T)} - x^*\|^2 \leq \|x^{(0)} - x^*\|^2 - \sum_{t=1}^T 2\eta_t (f(x^{(t-1)}, y^{(t-1)}) - f(x^*, y^{(t-1)})) + \sum_{t=1}^T \eta_t^2 \|g(t-1)\|^2 \quad (30)$$

675 Since $\|\mathbf{x}^{(t)} - \mathbf{x}^*\| \geq 0$, we have:

$$2 \sum_{t=1}^T \eta_t \left(f(\mathbf{x}^{(t-1)}, \mathbf{y}^{(t-1)}) - f(\mathbf{x}^*, \mathbf{y}^{(t-1)}) \right) \leq \|\mathbf{x}^{(0)} - \mathbf{x}^*\|^2 + \sum_{t=1}^T \eta_t^2 \|\mathbf{g}(t-1)\|^2 \quad (31)$$

676 Let $(\mathbf{x}_{\text{best}}^{(t)}, \mathbf{y}_{\text{best}}^{(t)}) = \arg \min_{(\mathbf{x}^{(k)}, \mathbf{y}^{(k)}) : k \in [t]} f(\mathbf{x}^{(k)}, \mathbf{y}^{(k)})$, then we have we have:

$$\sum_{t=1}^T \eta_t \left(f(\mathbf{x}^{(t-1)}, \mathbf{y}^{(t-1)}) - f(\mathbf{x}^*, \mathbf{y}^{(t-1)}) \right) \geq \sum_{t=1}^T \eta_t \left(f(\mathbf{x}^{(t-1)}, \mathbf{y}^{(t-1)}) - \max_{\mathbf{y} \in Y : \mathbf{g}(\mathbf{x}^*, \mathbf{y}) \geq 0} f(\mathbf{x}^*, \mathbf{y}) \right) \quad (32)$$

$$\geq \left(\sum_{t=1}^T \eta_t \right) \min_{t \in [T]} \left(f(\mathbf{x}^{(t-1)}, \mathbf{y}^{(t-1)}) - \max_{\mathbf{y} \in Y : \mathbf{g}(\mathbf{x}^*, \mathbf{y}) \geq 0} f(\mathbf{x}^*, \mathbf{y}) \right) \quad (33)$$

$$= \left(\sum_{t=1}^T \eta_t \right) \left(f(\mathbf{x}_{\text{best}}^{(T)}, \mathbf{x}_{\text{best}}^{(T)}) - \max_{\mathbf{y} \in Y : \mathbf{g}(\mathbf{x}^*, \mathbf{y}) \geq 0} f(\mathbf{x}^*, \mathbf{y}) \right) \quad (34)$$

677 Hence, we get the following bound:

$$f(\mathbf{x}_{\text{best}}^{(T)}, \mathbf{y}_{\text{best}}^{(T)}) - \max_{\mathbf{y} \in Y : \mathbf{g}(\mathbf{x}^*, \mathbf{y}) \geq 0} f(\mathbf{x}^*, \mathbf{y}) \leq \frac{\|\mathbf{x}^{(0)} - \mathbf{x}^*\|^2 + \sum_{t=1}^T \eta_t^2 \|\mathbf{g}(t-1)\|^2}{2 \left(\sum_{t=1}^T \eta_t \right)} \quad (35)$$

678 Since f is ℓ_f -Lipschitz with $\ell_f = \max_{(\hat{\mathbf{x}}, \hat{\mathbf{y}}) \in X \times Y} \|\nabla_{\mathbf{x}} f(\hat{\mathbf{x}}, \hat{\mathbf{y}})\|$, then for all $k \in \mathbb{N}$ we know that

679 $\|\mathbf{g}(k-1)\| \leq \ell_f$.

$$f(\mathbf{x}_{\text{best}}^{(T)}, \mathbf{x}_{\text{best}}^{(T)}) - \max_{\mathbf{y} \in Y : \mathbf{g}(\mathbf{x}^*, \mathbf{y}) \geq 0} f(\mathbf{x}^*, \mathbf{y}) \leq \frac{\|\mathbf{x}^{(0)} - \mathbf{x}^*\|^2 + \ell_f^2 \sum_{t=1}^T \eta_t^2}{2 \left(\sum_{t=1}^T \eta_t \right)} \quad (36)$$

$$f(\mathbf{x}_{\text{best}}^{(T)}, \mathbf{x}_{\text{best}}^{(T)}) - \min_{\mathbf{x} \in X} \max_{\mathbf{y} \in Y : \mathbf{g}(\mathbf{x}, \mathbf{y}) \geq 0} f(\mathbf{x}, \mathbf{y}) \leq \frac{\|\mathbf{x}^{(0)} - \mathbf{x}^*\|^2 + \ell_f^2 \sum_{t=1}^T \eta_t^2}{2 \left(\sum_{t=1}^T \eta_t \right)} \quad (37)$$

680 Under the assumption of the theorem:

$$\sum_{k=1}^T \eta_k^2 \leq \infty \quad \sum_{k=1}^T \eta_k = \infty \quad (38)$$

681 as $t \rightarrow \infty$, $\lim_{k \rightarrow \infty} f(\mathbf{x}_{\text{best}}^{(k)}, \mathbf{y}^{(k)}) \leq \min_{\mathbf{x} \in X} \max_{\mathbf{y} \in Y : \mathbf{g}(\mathbf{x}, \mathbf{y}) \geq 0} f(\mathbf{x}, \mathbf{y})$, and since for all $k \in \mathbb{N}$,

682 $\mathbf{y}_{\text{best}}^{(k)}$ satisfies $f(\mathbf{x}_{\text{best}}^{(k)}, \mathbf{y}_{\text{best}}^{(k)}) \geq \max_{\mathbf{y} \in Y : \mathbf{g}(\mathbf{x}, \mathbf{y}) \geq 0} f(\mathbf{x}_{\text{best}}^{(k)}, \mathbf{y}) - \delta$, as the number of iterations

683 increases, the best iterate converges to a $(0, \delta)$ -Stackelberg equilibrium. Additionally, setting

684 $\eta_t = \frac{\|\mathbf{x}^{(0)} - \mathbf{x}^*\|}{\ell_f \sqrt{T}}$ for all $t = 1, \dots, T$, we get:

$$f(\mathbf{x}_{\text{best}}^{(T)}, \mathbf{x}_{\text{best}}^{(T)}) - \min_{\mathbf{x} \in X} \max_{\mathbf{y} \in Y : \mathbf{g}(\mathbf{x}, \mathbf{y}) \geq 0} f(\mathbf{x}^*, \mathbf{y}) \leq \frac{\ell_f \|\mathbf{x}^{(0)} - \mathbf{x}^*\|^2}{\sqrt{T}} \quad (39)$$

685 Hence, the best iterate converges to a (ε, δ) -Stackelberg equilibrium in $O(\varepsilon^{-2})$ iterations. \square

686 **Theorem D.2.** Suppose that Algorithm [1](#) is run on a convex-concave min-max game with de-
 687 pendent strategy sets given by (X, Y, f, \mathbf{g}) where X is convex. Suppose that Assumption [3.1](#)
 688 holds and that additionally f is μ -strongly convex in \mathbf{x} , i.e., $\forall \mathbf{x}_1, \mathbf{x}_2 \in X, \mathbf{y} \in Y, f(\mathbf{x}_1, \mathbf{y}) \geq$
 689 $f(\mathbf{x}_2, \mathbf{y}) + \langle \mathbf{g}, (\mathbf{x}_1 - \mathbf{x}_2) \rangle + \frac{\mu}{2} \|\mathbf{x}_1 - \mathbf{x}_2\|^2$ where $\mathbf{g} \in \partial_{\mathbf{x}} f(\mathbf{x}_2, \mathbf{y})$. Then, if $(\mathbf{x}_{\text{best}}^{(t)}, \mathbf{y}_{\text{best}}^{(t)}) \in$

690 $\arg \min_{(\mathbf{x}^{(k)}, \mathbf{y}^{(k)}): k \in [t]} f(\mathbf{x}^{(k)}, \mathbf{y}^{(k)})$, for $\varepsilon \in (0, 1)$, and $\eta_t = \frac{2}{\mu(t+1)}$, if we choose T large enough
 691 such that:

$$T \geq N_T(\varepsilon) \doteq O(\varepsilon^{-1})$$

692 then there exists an iteration $T^* \leq T$ such that $(\mathbf{x}_{\text{best}}^{(T^*)}, \mathbf{y}_{\text{best}}^{(T^*)})$ is an (ε, δ) -Stackelberg equilibrium.

693 *Proof of Theorem D.2* Note that by Theorem 3.2 we have $\nabla_{\mathbf{x}} f(\mathbf{x}^{(t-1)}, \mathbf{y}^{(t-1)}) +$
 694 $\sum_{k=1}^K \lambda_k^{(t-1)} \nabla_{\mathbf{x}} g_k(\mathbf{x}^{(t-1)}, \mathbf{y}^{(t-1)}) \in \partial_{\mathbf{x}} V(\mathbf{x}^{(t-1)}) = \partial_{\mathbf{x}} \max_{\mathbf{y} \in Y: \mathbf{g}(\mathbf{x}^{(t-1)}, \mathbf{y}) \geq 0} f(\mathbf{x}^{(t-1)}, \mathbf{y})$.
 695 For notational clarity, let $\mathbf{g}(t-1) = \nabla_{\mathbf{x}} f(\mathbf{x}^{(t-1)}, \mathbf{y}^{(t-1)}) + \sum_{k=1}^K \lambda_k^{(t-1)} \nabla_{\mathbf{x}} g_k(\mathbf{x}^{(t-1)}, \mathbf{y}^{(t-1)})$.
 696 Suppose that $\mathbf{x}^* \in \arg \min_{\mathbf{x} \in X} \max_{\mathbf{y} \in Y: \mathbf{g}(\mathbf{x}, \mathbf{y})} f(\mathbf{x}, \mathbf{y})$. For any $t \in \mathbb{N}$ such that $t \geq 1$, we have:

$$\|\mathbf{x}^{(t)} - \mathbf{x}^*\|^2 = \|\Pi_X(\mathbf{x}^{(t-1)} - \eta_t \mathbf{g}(t-1)) - \Pi_X(\mathbf{x}^*)\|^2 \quad (40)$$

$$\leq \|\mathbf{x}^{(t-1)} - \eta_t \mathbf{g}(t-1) - \mathbf{x}^*\|^2 \quad (41)$$

$$= \|\mathbf{x}^{(t-1)} - \mathbf{x}^*\|^2 - 2\eta_t \langle \mathbf{g}(t-1), (\mathbf{x}^{(t-1)} - \mathbf{x}^*) \rangle + \eta_t^2 \|\mathbf{g}(t-1)\|^2 \quad (42)$$

$$\leq \|\mathbf{x}^{(t-1)} - \mathbf{x}^*\|^2 - 2\eta_t \left[\frac{\mu}{2} \|\mathbf{x}^{(t-1)} - \mathbf{x}^*\|^2 + f(\mathbf{x}^{(t-1)}, \mathbf{y}^{(t)}) - f(\mathbf{x}^*, \mathbf{y}^{(t)}) \right] + \eta_t^2 \|\mathbf{g}(t-1)\|^2 \quad (43)$$

$$= \|\mathbf{x}^{(t-1)} - \mathbf{x}^*\|^2 - \eta_t \mu \|\mathbf{x}^{(t-1)} - \mathbf{x}^*\|^2 - 2\eta_t (f(\mathbf{x}^{(t-1)}, \mathbf{y}^{(t-1)}) - f(\mathbf{x}^*, \mathbf{y}^{(t-1)})) + \eta_t^2 \|\mathbf{g}(t-1)\|^2 \quad (44)$$

$$= (1 - \eta_t \mu) \|\mathbf{x}^{(t-1)} - \mathbf{x}^*\|^2 - 2\eta_t (f(\mathbf{x}^{(t-1)}, \mathbf{y}^{(t-1)}) - f(\mathbf{x}^*, \mathbf{y}^{(t-1)})) + \eta_t^2 \|\mathbf{g}(t-1)\|^2 \quad (45)$$

697 where the first line follows from definitions, the second from the non-expansiveness of the projection
 698 operator, the third from algebra, the fourth from the definition of strong convexity, i.e., $\mathbf{g}(t-1)^T (\mathbf{x}^{(t-1)} - \mathbf{x}^*) \geq \frac{\mu}{2} \|\mathbf{x}^{(t-1)} - \mathbf{x}^*\|^2 + f(\mathbf{x}^{(t-1)}, \mathbf{y}^{(t-1)}) - f(\mathbf{x}^*, \mathbf{y}^{(t-1)})$.
 699

700 Re-organizing expressions, we get:

$$f(\mathbf{x}^{(t-1)}, \mathbf{y}^{(t-1)}) - f(\mathbf{x}^*, \mathbf{y}^{(t-1)}) \leq \frac{1 - \eta_t \mu}{2\eta_t} \|\mathbf{x}^{(t-1)} - \mathbf{x}^*\|^2 - \frac{1}{2\eta_t} \|\mathbf{x}^{(t)} - \mathbf{x}^*\|^2 + \frac{\eta_t}{2} \|\mathbf{g}(t-1)\|^2 \quad (46)$$

701 Setting $\eta_t = \frac{2}{\mu(t+1)}$, we get:

$$f(\mathbf{x}^{(t-1)}, \mathbf{y}^{(t-1)}) - f(\mathbf{x}^*, \mathbf{y}^{(t-1)}) \leq \frac{\mu(t-1)}{4} \|\mathbf{x}^{(t-1)} - \mathbf{x}^*\|^2 - \frac{\mu(t+1)}{4} \|\mathbf{x}^{(t)} - \mathbf{x}^*\|^2 + \frac{1}{\mu(t+1)} \|\mathbf{g}(t-1)\|^2 \quad (47)$$

$$t (f(\mathbf{x}^{(t-1)}, \mathbf{y}^{(t-1)}) - f(\mathbf{x}^*, \mathbf{y}^{(t-1)})) \leq \frac{\mu t(t-1)}{4} \|\mathbf{x}^{(t-1)} - \mathbf{x}^*\|^2 - \frac{\mu t(t+1)}{4} \|\mathbf{x}^{(t)} - \mathbf{x}^*\|^2 + \frac{1}{\mu} \|\mathbf{g}(t-1)\|^2 \quad (48)$$

702 where the last line was obtained by multiplying by t on both sides.

703 Summing up across all iterations on both sides:

$$\sum_{t=0}^T t \left(f(\mathbf{x}^{(t)}, \mathbf{y}^{(t)}) - f(\mathbf{x}^*, \mathbf{y}^{(t)}) \right) \leq \sum_{t=0}^T \frac{\mu t(t-1)}{4} \|\mathbf{x}^{(t-1)} - \mathbf{x}^*\|^2 - \sum_{t=0}^T \frac{\mu t(t+1)}{4} \|\mathbf{x}^{(t)} - \mathbf{x}^*\|^2 + \sum_{t=0}^T \frac{1}{\mu} \|\mathbf{g}(t-1)\|^2 \quad (49)$$

$$= \sum_{t=0}^T \frac{\mu t(t-1)}{4} \|\mathbf{x}^{(t-1)} - \mathbf{x}^*\|^2 - \sum_{t=1}^{T+1} \frac{\mu(t-1)t}{4} \|\mathbf{x}^{(t-1)} - \mathbf{x}^*\|^2 + \sum_{t=0}^T \frac{1}{\mu} \|\mathbf{g}(t-1)\|^2 \quad (50)$$

$$= -\frac{\mu t(t+1)}{4} \|\mathbf{x}^{(t)} - \mathbf{x}^*\|^2 + \sum_{t=0}^T \frac{1}{\mu} \|\mathbf{g}(t-1)\|^2 \quad (51)$$

$$\leq \sum_{t=0}^T \frac{1}{\mu} \|\mathbf{g}(t-1)\|^2 \quad (52)$$

$$\leq \frac{T}{\mu} \ell_f \quad (53)$$

704 where the last line was obtained by noticing that f is ℓ_f -Lipschitz with $\ell_f =$
 705 $\max_{(\hat{\mathbf{x}}, \hat{\mathbf{y}}) \in X \times Y} \|\nabla_{\mathbf{x}} f(\hat{\mathbf{x}}, \hat{\mathbf{y}})\|$, which implies that for all $k \in \mathbb{N}$ we know that $\|\mathbf{g}(k-1)\| \leq \ell_f$.

706 Let $(\mathbf{x}_{\text{best}}^{(t)}, \mathbf{y}_{\text{best}}^{(t)}) = \arg \min_{(\mathbf{x}^{(k)}, \mathbf{y}^{(k)}) : k \in [t]} f(\mathbf{x}^{(k)}, \mathbf{y}^{(k)})$:

$$\sum_{t=0}^T t \left(f(\mathbf{x}^{(t)}, \mathbf{y}^{(t)}) - f(\mathbf{x}^*, \mathbf{y}^{(t)}) \right) \leq \frac{T}{\mu} \ell_f \quad (54)$$

$$\sum_{t=0}^T t \left(f(\mathbf{x}^{(t)}, \mathbf{y}^{(t)}) - \max_{\mathbf{y} \in Y : \mathbf{g}(\mathbf{x}^*, \mathbf{y}) \geq \mathbf{0}} f(\mathbf{x}^*, \mathbf{y}) \right) \leq \frac{T}{\mu} \ell_f \quad (55)$$

$$\left(\sum_{t=0}^T t \right) \min_{t \in [T]} \left(f(\mathbf{x}^{(t)}, \mathbf{y}^{(t)}) - \max_{\mathbf{y} \in Y : \mathbf{g}(\mathbf{x}^*, \mathbf{y}) \geq \mathbf{0}} f(\mathbf{x}^*, \mathbf{y}) \right) \leq \frac{T}{\mu} \ell_f \quad (56)$$

$$f(\mathbf{x}_{\text{best}}^{(T)}, \mathbf{y}_{\text{best}}^{(T)}) - \max_{\mathbf{y} \in Y : \mathbf{g}(\mathbf{x}^*, \mathbf{y}) \geq \mathbf{0}} f(\mathbf{x}^*, \mathbf{y}) \leq \frac{\ell_f}{\mu(T+1)} \quad (57)$$

$$(58)$$

707 That is, as the number of iterations increases, the best iterate converges to a $(0, \delta)$ -Stackelberg
 708 equilibrium. Additionally, the best iterate converges to a (ε, δ) -Stackelberg equilibrium in $O(\varepsilon^{-1})$
 709 iterations. \square

710 We present the following theorem which proves one of the cases given in Theorem 3.4. The proof
 711 for the other cases is the same as the proof below. We note that gradient ascent converges in
 712 $O(\varepsilon^{-1})$ iterations to a ε -maximum for a Lipschitz smooth objective, and in $O(\log(\varepsilon))$ iterations to
 713 a ε -maximum for a Lipschitz smooth and strongly concave objective [6].

714 **Theorem D.3.** Suppose that Algorithm 2 is run on a convex-concave min-max game with de-
 715 pendent strategy sets given by (X, Y, f, \mathbf{g}) where X, Y are convex. Suppose that Assumption 3.1
 716 holds and f is $\ell_{\nabla f}$ -smooth, i.e., $\forall (\mathbf{x}_1, \mathbf{y}_1), (\mathbf{x}_2, \mathbf{y}_2) \in X \times Y, \|\nabla f(\mathbf{x}_1, \mathbf{y}_1) - \nabla f(\mathbf{x}_2, \mathbf{y}_2)\| \leq$
 717 $\ell_{\nabla f} \|\mathbf{x}_1, \mathbf{y}_1) - (\mathbf{x}_2, \mathbf{y}_2)\|$.

718 Let $(\mathbf{x}_{\text{best}}^{(t)}, \mathbf{y}_{\text{best}}^{(t)}) \in \arg \min_{(\mathbf{x}^{(k)}, \mathbf{y}^{(k)}) : k \in [t]} f(\mathbf{x}^{(k)}, \mathbf{y}^{(k)})$. For $\varepsilon \in (0, 1)$, if we choose $T_{\mathbf{x}}$ and $T_{\mathbf{y}}$
 719 large enough such that:

$$T_{\mathbf{x}} \geq N_{T_{\mathbf{x}}}(\varepsilon) := O(\varepsilon^{-2}) \quad (59)$$

$$T_{\mathbf{y}} \geq N_{T_{\mathbf{y}}}(\varepsilon) := O(\varepsilon^{-1}) \quad (60)$$

720 then there exists an iteration $T^* \leq T_{\mathbf{x}} T_{\mathbf{y}} = O(\varepsilon^{-3})$ such that $(\mathbf{x}_{\text{best}}^{(T^*)}, \mathbf{y}_{\text{best}}^{(T^*)})$ is an $(\varepsilon, \varepsilon)$ -Stackelberg
 721 equilibrium.

722 *Proof of Theorem 3.3.* Since f is ℓ_f -smooth, it is well known that the inner gradient descent pro-
 723 cedure will compute an ε -maximum of $f(\mathbf{x}^{(t)}, \cdot)$ for each iterate $\mathbf{x}^{(t)}$ in $O(\varepsilon^{-2})$ iterations [6].
 724 Combining the iteration complexity of the outer and inner loops using this result and Theorem 3.3,
 725 we obtain an iteration complexity of $O(\varepsilon^{-2})O(\varepsilon^{-1}) = O(\varepsilon^{-3})$. \square

726 E An Economic Application: Details

727 E.1 Experimental details

728 E.1.1 General Experiment Setup

729 Our experimental goals were two-folds. First, we wanted to understand the convergence complexity
 730 of our algorithms for different Fisher markets under which the objective function Equation (3)
 731 satisfies different smoothness properties. Secondly, we wanted to understand the approximate
 732 optimal $\mathbf{y}^{(t)}$ computed by the max-oracle in Algorithm 1 or by the inner loop in Algorithm 2 affected
 733 the preciseness of Stackelberg equilibrium outputed.

734 To answer these questions, we have collected data on the prices and allocations computed by
 735 Algorithm 1 with an exact max-oracle on each iteration and by Algorithm 2 on each iteration of
 736 the outer loop algorithms by running them on randomly initialized markets. We have initialized
 737 500 different linear, Cobb-Douglas, Leontief Fisher markets with 5 buyers and 8 goods. For each of
 738 these markets, we have run Algorithm 1 and Algorithm 2 twice, one time with high starting prices
 739 and one time with low starting prices to understand the impact of initialization conditions on the
 740 algorithm. We have run Algorithm 1 and Algorithm 2's outer loop for 500, 300, and 700 iterations
 741 for linear, Cobb-Douglas, and Leontief Fisher markets respectively.

742 We have opted for a learning rate of 5 for both algorithms after manual hyper-parameter tuning
 743 and picked a decay rate of $t^{-1/2}$ for the learning rate based on our theory. For each run of the
 744 algorithm, we then computed the objective functions value for the iterates calculated by the
 745 algorithm throughout it to obtain Figure 1. Finally, to understand how much precision was lost
 746 in the accuracy of the Stackelberg equilibrium outputed by Algorithm 2 from not being able to
 747 compute a maximum of $f(\mathbf{x}, \cdot)$ for given $\mathbf{x} \in X$, we have run a first order James' test to see if the
 748 equilibrium strategies outputed by Algorithm 1 and Algorithm 2 were statistically distinguishable.

749 E.1.2 Computational Requirements, Packages, and Algorithmic Details

750 The experiments were run on MacOS machine with 8GB ram and an apple M1 chip and experiments
 751 took about 2 hours to run. Only CPU resources were used.

752 We have run our experiments in Python 3.7 [64] and have used the NumPy [28], Pandas [60], and
 753 CVXPY [19]. The data from our experiments can be found on our code repository as well (<https://anonymous.4open.science/r/min-max-fisher-CEFA/>). Figure 1 was graphed via Matplotlib
 754 [31]. To run the first order James test, we transfer the data generated by our Python code to an
 755 R script [51], which we manipulate using the Tidyverse environment [67], and finally obtain the
 756 desired p-values via the STests package in R [30].

758 **Licensing** R as a package is licensed under GPL-2 | GPL-3. Python software and documentation
 759 are licensed under the PSF License Agreement. Numpy is distributed under a liberal BSD license.
 760 Pandas is distributed under a new BSD license. Matplotlib only uses BSD compatible code, and
 761 its license is based on the PSF license. CVXPY is licensed under an APACHE license. Tidyverse is
 762 distributed under an MIT license.

763 For our execution of algorithm Algorithm 1 for linear, Cobb-Douglas and Leontief Fisher markets,
 764 we used an exact Max-Oracle since the demand has a closed form solution for these markets [25].
 765 As the computational overhead of the projection operation in the inner loop of Algorithm 2 can be
 766 high for most projection methods, we have opted to use CVXPY first order for the inner loop of
 767 Algorithm 2. In particular, we have opted for the ECOS solver and in case if any runtime exception
 768 occurred. Note that these solvers compute ε -optimal points as a result we believe that they present
 769 an accurate view of how Algorithm 2 would behave.

Algorithm 3 δ -Approximate Tâtonnement for Fisher Markets**Inputs:** $U, b, \eta, T, \mathbf{p}^{(0)}, \delta$ **Output:** $(\mathbf{X}^*, \mathbf{p}^*)$

```

1: for  $t = 1, \dots, T$  do
2:   For all  $i \in [n]$ , find  $\mathbf{x}_i^{(t)}$  s.t.  $u_i(\mathbf{x}_i^{(t)}) \geq \max_{\mathbf{x}_i: \mathbf{x}_i \cdot \mathbf{p}^{(t-1)} \leq b_i} u_i(\mathbf{x}_i) - \delta$  and  $\mathbf{x}_i^{(t)} \cdot \mathbf{p}^{(t-1)} \leq b_i$ 
3:   Set  $\mathbf{p}^{(t)} = \max \left\{ \mathbf{p}^{(t-1)} - \eta_t (1 - \sum_{i \in [n]} \mathbf{x}_i^{(t)}), 0 \right\}$ 
4: end for
5: return  $(\mathbf{X}^{(T)}, \mathbf{p}^{(T)})$ 

```

Algorithm 4 δ -Approximate Nested Tâtonnement for Fisher Markets**Inputs:** $U, b, \eta, T_p, T_X, \mathbf{p}^{(0)}$ **Output:** $(\mathbf{X}^*, \mathbf{p}^*)$

```

1: for  $t = 1, \dots, T_p$  do
2:   for  $s = 1, \dots, T_X$  do
3:     For all  $i \in [n]$ ,  $\mathbf{x}_i^{(t)} = \Pi_{\{\mathbf{x}: \mathbf{x} \cdot \mathbf{p}^{(t-1)} \leq b_i\}} \left( \mathbf{x}_i^{(t)} + \frac{b_i}{u_i(\mathbf{x}_i^{(t)})} \nabla_{\mathbf{x}_i} u_i(\mathbf{x}_i^{(t)}) \right)$ 
4:   end for
5:   Set  $\mathbf{p}^{(t)} = \max \left\{ \mathbf{p}^{(t-1)} - \eta_t (1 - \sum_{i \in [n]} \mathbf{x}_i^{(t)}), 0 \right\}$ 
6: end for
7: return  $(\mathbf{X}^{(T)}, \mathbf{p}^{(T)})$ 

```

771 **F Additional Related Work**

772 Much progress has been made recently in solving min-max games with independent strategy sets,
 773 both in the convex-concave case and in non-convex-concave case. For the former case, when
 774 f is μ_x -strongly-convex- μ_y -strongly-concave, Tseng [63], Yuri Nesterov [69], and Gidel et al.
 775 [24] proposed variational inequality methods and Mokhtari, Ozdaglar, and Pattathil [42] gradient-
 776 descent-ascent (GDA)-based methods that compute a solution in $\tilde{O}(\mu_y + \mu_x)$ iterations.. These
 777 upper bounds were recently complemented by the lower bound of $\tilde{\Omega}(\sqrt{\mu_y \mu_x})$, shown by Ibrahim
 778 et al. [32] and Zhang, Hong, and Zhang [70]. Subsequently, Lin, Jin, and Jordan [38] and Alkousa et
 779 al. [3] analyzed algorithms that converge in $\tilde{O}(\sqrt{\mu_y \mu_x})$ and $\tilde{O}(\min \{\mu_x \sqrt{\mu_y}, \mu_y \sqrt{\mu_x}\})$ iterations,
 780 respectively. For the special case where f is μ_x -strongly-convex-linear, Juditsky, Nemirovski,
 781 et al. [35], Hamedani and Aybat [27], and Zhao [72] all present methods that converge to an ε -
 782 approximate solution in $O(\sqrt{\mu_x/\varepsilon})$. When assumptions on $f(\mathbf{x}, \cdot)$ are dropped and it is assumed to
 783 be μ_x -strongly-convex-concave, Thekumparampil et al. [61] provide an algorithm that converges to
 784 an approximate solution in $\tilde{O}(\mu_x/\varepsilon)$, and Ouyang and Xu [49] provide a lower bound of $\tilde{\Omega}(\sqrt{\mu_x/\varepsilon})$.
 785 Lin, Jin, and Jordan then went on to develop a faster algorithm, with iteration complexity of
 786 $\tilde{O}(\sqrt{\mu_x/\varepsilon})$. When f is simply assumed to be convex-concave, Nemirovski [43], Nesterov [44], and
 787 Tseng [62] describe an algorithm with $\tilde{O}(\varepsilon^{-1})$ and Ouyang and Xu [49] prove a lower bound of
 788 $\Omega(\varepsilon^{-1})$. We include a detailed summary table of these results in Table 4.

789 When f is assumed to be non-convex- μ_y -strongly-concave, and the goal is to compute a first-order
 790 Nash or “local” Stackelberg equilibrium, Sanjabi et al. [54] provide an algorithm that converges
 791 to ε -an approximate solution in $O(\varepsilon^{-2})$ iterations. Jin, Netrapalli, and Jordan [34], Rafique et al.
 792 [52], Lin, Jin, and Jordan [37], and Lu, Tsaknakis, and Hong [39] provide algorithms that converge
 793 in $\tilde{O}(\mu_y^2 \varepsilon^{-2})$, while Lin, Jin, and Jordan [38] provide an even faster algorithm, with an iteration
 794 complexity of $\tilde{O}(\sqrt{\mu_y} \varepsilon^{-2})$. When f is non-convex-non-concave and the goal to compute is an
 795 approximate first-order Nash equilibrium, Lu, Tsaknakis, and Hong [39] provide an algorithm
 796 with iteration complexity $\tilde{O}(\varepsilon^{-4})$, while Nouiehed et al. [47] provide an algorithm with iteration

complexity $\tilde{O}(\varepsilon^{-3.5})$. More recently, Ostrovskii, Lowy, and Razaviyayn [48] and Lin, Jin, and Jordan [38] proposed an algorithm with iteration complexity $\tilde{O}(\varepsilon^{-2.5})$. When f is non-convex-non-concave and the desired solution concept is a “local” Stackelberg equilibrium, Jin, Netrapalli, and Jordan [34], Rafique et al. [52], and Lin, Jin, and Jordan [37] provide algorithms with a $\tilde{O}(\varepsilon^{-6})$ complexity. More recently, Thekumparampil et al. [61], Zhao [71], and Lin, Jin, and Jordan [38] have proposed algorithms that converge to an ε -approximate solution in $\tilde{O}(\varepsilon^{-3})$ iterations. We include a detailed summary table of these results in Table 5

Table 4: Iteration complexities for min-max games with independent strategy sets in convex-concave settings. Note that these results assume that the objective function is Lipschitz-smooth.

Setting	Reference	Iteration Complexity
μ_x -Strongly-Convex- μ_y -Strongly-Concave	[63]	$\tilde{O}(\mu_x + \mu_y)$
	[69]	
	[24]	
	[42]	
	[3]	$\tilde{O}(\min\{\mu_x\sqrt{\mu_y}, \mu_y\sqrt{\mu_x}\})$
	[38]	$\tilde{O}(\sqrt{\mu_x\mu_y})$
	[32]	$\tilde{\Omega}(\sqrt{\mu_x\mu_y})$
	[70]	
μ_x -Strongly-Convex-Linear	[35]	$O(\sqrt{\mu_x/\varepsilon})$
	[27]	
	[72]	
μ_x -Strongly-Convex-Concave	[61]	$\tilde{O}(\mu_x/\sqrt{\varepsilon})$
	[38]	$\tilde{O}(\sqrt{\mu_x/\varepsilon})$
	[49]	$\tilde{\Omega}(\sqrt{\mu_x/\varepsilon})$
Convex-Concave	[43]	$O(\varepsilon^{-1})$
	[44]	
	[62]	
	[38]	$\tilde{O}(\varepsilon^{-1})$
	[49]	$\Omega(\varepsilon^{-1})$

Table 5: Iteration complexities for min-max games with independent strategy sets in non-convex-concave settings. Note that although all these results assume that the objective function is Lipschitz smooth, some authors make more assumptions, e.g., [47] prove their result for objective functions that satisfy the Lojasiwicz condition.

Setting	Reference	Iteration Complexity
Nonconvex- $\mu_{\mathbf{y}}$ -Strongly-Concave, First Order Nash Equilibrium or Local Stackelberg Equilibrium	[34]	$\tilde{O}(\mu_{\mathbf{y}}^2 \varepsilon^{-2})$
	[52]	
	[37]	
	[39]	
	[38]	$\tilde{O}(\sqrt{\mu_{\mathbf{y}}} \varepsilon^{-2})$
Nonconvex-Concave, First Order Nash Equilibrium	[39]	$\tilde{O}(\varepsilon^{-4})$
	[47]	$\tilde{O}(\varepsilon^{-3.5})$
	[48]	$\tilde{O}(\varepsilon^{-2.5})$
	[38]	
Nonconvex-Concave Local Stackelberg Equilibrium	[34]	$\tilde{O}(\varepsilon^{-6})$
	[47]	
	[38]	
	[61]	$\tilde{O}(\varepsilon^{-3})$
	[71]	
	[38]	