

ZERO-SUM STOCHASTIC STACKELBERG GAMES

Denizalp Goktas

Department of Computer Science
Brown University
Providence, RI 02906, USA
denizalp_goktas@brown.edu

Jiayi Zhao

Department of Computer Science
Pomona College
Pomona, CA, USA
jzae2019@mymail.pomona.edu

Amy Greenwald

Brown University
Providence, RI 02906, USA
amy_greenwald@brown.edu

ABSTRACT

Min-max optimization problems (i.e., zero-sum games) have been used to model problems in a variety of fields in recent years, from machine learning to economics. The literature to date has mostly focused on static zero-sum games, assuming independent strategy sets. In this paper, we study a form of dynamic zero-sum games, called stochastic games, with dependent strategy sets. Just as zero-sum games with dependent strategy sets can be interpreted as zero-sum Stackelberg games, stochastic zero-sum games with dependent strategy sets can be interpreted as zero-sum stochastic Stackelberg games. We prove the existence of an optimal solution in zero-sum stochastic Stackelberg games (i.e., a recursive Stackelberg equilibrium), provide necessary and sufficient conditions for a solution to be optimal, and show that a recursive Stackelberg equilibrium can be computed in polynomial time via value iteration. Finally, we show that stochastic Stackelberg games can model the problem of pricing and allocating goods across agents and time; more specifically, we propose a stochastic Stackelberg game whose solutions correspond to a recursive competitive equilibrium in a stochastic Fisher market. We close with a series of experiments which confirm our theoretical results and show how value iteration performs in practice.

Min-max optimization has paved the way for recent progress in a variety of fields, from machine learning to economics. These applications require computing solutions to **the min-max operator**, i.e., solving a **constrained min-max optimization problem**, also known as zero-sum games (with independent strategy sets): i.e., $\min_{x \in \mathcal{X}} \max_{y \in \mathcal{Y}} f(x, y)$, where the objective function $f : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$ is continuous, and the strategy sets $\mathcal{X} \subset \mathbb{R}^n$ and $\mathcal{Y} \subset \mathbb{R}^m$ are nonempty and compact. When f is convex-concave, the seminal minimax theorem (Neumann, 1928; Sion et al., 1958) holds, and such a problem can be interpreted as a simultaneous-move zero-sum game between an outer player x and an inner player y , with the solutions $(x^*, y^*) \in \mathcal{X} \times \mathcal{Y}$ of the min-max operator corresponding to a Nash equilibrium. More generally, one can consider **zero-sum stochastic games** (with independent strategy sets) $\mathcal{X} \subset \mathbb{R}^n$ and $\mathcal{Y} \subset \mathbb{R}^m$, nonempty and compact, and a state-dependent payoff function $r(s, x, y)$, for all $s \in \mathcal{S}$, where the players seek to optimize their cumulative (discounted) payoffs, in expectation. When $r(s, x, y)$ is bounded, continuous, and concave-convex in (x, y) , for all $s \in \mathcal{S}$, these games are guaranteed to have a unique solution (Shapley, 1953),¹ while the **optimal policies**, i.e., the per-state collection of solutions to the min-max operator can be interpreted as a **recursive Nash equilibrium** of a zero-sum stochastic game, and can be computed in polynomial time by iterative application of the min-max operator (Shapley, 1953). Zero-sum *stochastic* games generalize zero-sum games from a single state to multiple states, and have found even more applications in a variety of fields (Jaśkiewicz & Nowak, 2018).

¹Although Shapley’s original results concern payoffs which are bilinear in the outer and inner players’ actions, they extend directly to payoffs which are convex-concave in the players’ actions.

Recently, Goktas & Greenwald (2021) studied the computation of a **generalized min-max operator** i.e., solving a **constrained min-max game with dependent strategy sets**: i.e., $\min_{\mathbf{x} \in \mathcal{X}} \max_{\mathbf{y} \in \mathcal{Y}: h(\mathbf{x}, \mathbf{y}) \geq \mathbf{0}} r(\mathbf{x}, \mathbf{y})$ where, in addition to the aforementioned assumptions, $h : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$ is continuous. Goktas & Greenwald have shown that even more problems of interest can be captured by solutions to the generalized min-max operator. However, a minimax theorem is not guaranteed to hold in this setting; as a result, these problems are best interpreted as zero-sum sequential, i.e., Stackelberg (Von Stackelberg, 1934), games, and their solutions, as Stackelberg equilibria. One can likewise consider **zero-sum stochastic games with dependent feasible strategy sets** $\mathcal{X} \subset \mathbb{R}^n$ and $\mathcal{Y} \subset \mathbb{R}^m$, nonempty and compact and a state-dependent payoff function $r(\mathbf{s}, \mathbf{x}, \mathbf{y})$, as well as a state dependent constraint function $\mathbf{g}(\mathbf{s}, \mathbf{x}, \mathbf{y})$, for all $\mathbf{s} \in \mathcal{S}$, where the players seek to optimize their cumulative (discounted) payoffs, in expectation, while satisfying the constraint $\mathbf{g}(\mathbf{s}, \mathbf{x}, \mathbf{y}) \geq \mathbf{0}$ at each state $\mathbf{s} \in \mathcal{S}$. Such a problem is best interpreted as a zero-sum stochastic Stackelberg game, for which the appropriate solution concept is the **recursive Stackelberg equilibrium (recSE)**. Although very little is known about the properties of such problems, they generalize both min-max games with dependent strategy sets and zero-sum stochastic games (with independent strategy sets), and, as we show, have novel applications.

In this paper, we prove the existence of recSE in zero-sum stochastic Stackelberg games, provide necessary and sufficient conditions for a solution to be a recSE, and show that a recSE can be computed in (weakly) polynomial time via value iteration. We further show that stochastic Stackelberg games can be used to solve problems of pricing and allocating goods across agents and time. In particular, we introduce **stochastic Fisher markets**, a stochastic generalization of the Fisher market (Brainard et al., 2000), and a special case of Friesen’s (1979) financial market model, which itself is a stochastic generalization of the Arrow & Debreu model of a competitive economy (1954). We then prove the existence of recursive competitive equilibrium (Mehra & Prescott, 1977) in this model, under the assumption that consumers have continuous and homogeneous utility functions, by characterizing the competitive equilibria of any stochastic Fisher market as the Stackelberg equilibria of the corresponding stochastic Stackelberg game. Finally, we use value iteration to solve various stochastic Fisher markets, highlighting the issues that value iteration might face as a consequence of the smoothness properties of the utility functions.

Related Work Algorithms for min-max optimization problems (i.e., zero-sum games) with independent strategy sets have been extensively studied; for a summary see Goktas & Greenwald (2021), Section G. Goktas & Greenwald (2021) and (2022) studied zero-sum games with dependent strategy sets, proposing polynomial-time nested gradient descent ascent (GDA) and simultaneous GDA algorithms for such problems.

The computation of Stackelberg equilibrium in dynamic Stackelberg games has been studied in several interesting settings, but always with independent strategy sets. Bensoussan et al. (2015) study continuous-time stochastic Stackelberg games with continuous action spaces, and prove existence of a solution in their setting. Vasal (2020) and Vorobeychik & Singh (2012) study discrete-time stochastic Stackelberg games with discrete action and state spaces, and provide algorithms to solve such games. DeMiguel & Xu (2009) consider a stochastic Stackelberg game-like market model with n leaders and m followers; they prove the existence of a Stackelberg equilibrium in their model, and provide (without theoretical guarantees) algorithms that converge to such an equilibrium in experiments. Dynamic Stackelberg games (Li & Sethi, 2017) have been applied to a wide range of problems, including security (Vasal, 2020; Vorobeychik & Singh, 2012), insurance provision (Chen & Shen, 2018; Yuan et al., 2021), advertising (He et al., 2008), robust agent design (Rismiller et al., 2020), allocating goods across time intertemporal pricing (Oksendal et al., 2013).

The study of algorithms that compute competitive equilibria in Fisher markets was initiated by Devanur et al. (2002), who provided a polynomial-time method for solving these markets assuming linear utilities. More recently, there have been efforts to study markets in dynamic settings (Cheung et al., 2019; Gao et al., 2021; Goktas & Greenwald, 2021), in which the goal is to either track the changing equilibrium of a changing market, or minimize some regret-like quantity for the market. The models considered in these earlier works differ from ours as they do not have stochastic structure and do not invoke a dynamic solution concept.

1 PRELIMINARIES

Notation We use caligraphic uppercase letters to denote sets (e.g., \mathcal{X}), bold uppercase letters to denote matrices (e.g., \mathbf{X}), bold lowercase letters to denote vectors (e.g., \mathbf{p}), bold uppercase letters to denote vector-valued random variables (e.g., $\mathbf{\Gamma}$), lowercase letters to denote scalar quantities, (e.g., x), and uppercase letters to denote scalar-valued random variables (e.g., X). We denote the i th row vector of a matrix (e.g., \mathbf{X}) by the corresponding bold lowercase letter with subscript i (e.g., \mathbf{x}_i). Similarly, we denote the j th entry of a vector (e.g., \mathbf{p} or \mathbf{x}_i) by the corresponding Roman lowercase letter with subscript j (e.g., p_j or x_{ij}). We denote functions by a letter: e.g., f if the function is scalar valued, and \mathbf{f} if the function is vector valued. We denote the vector of ones of size n by $\mathbf{1}_n$. We denote the set of integers $\{1, \dots, n\}$ by $[n]$, the set of natural numbers by \mathbb{N} , the set of real numbers by \mathbb{R} . We denote the positive and strictly positive elements of a set by a $+$ and $++$ subscript respectively, e.g., \mathbb{R}_+ and \mathbb{R}_{++} . We denote the orthogonal projection operator onto a set C by Π_C , i.e., $\Pi_C(\mathbf{x}) = \arg \min_{\mathbf{y} \in C} \|\mathbf{x} - \mathbf{y}\|^2$. We denote by $\Delta_n = \{\mathbf{x} \in \mathbb{R}_+^n \mid \sum_{i=1}^n x_i = 1\}$, and by $\Delta(A)$ the set of probability measures on the set A .

A **stochastic Stackelberg game** $(\mathcal{S}, \mathcal{X}, \mathcal{Y}, \mu^{(0)}, r_x, r_y, \mathbf{g}, p, \gamma)$ is a dynamic two-player sequential, i.e., Stackelberg, game where one player we call the outer-player (resp. inner-player) moves through a set of states \mathcal{S} picking an action to play from their continuous set of actions $\mathcal{X} \subset \mathbb{R}^n$ (resp. $\mathcal{Y} \subset \mathbb{R}^m$). Players start the game at an initial state determined by an initial state distribution $\mu^{(0)} : \mathcal{S} \rightarrow [0, 1]$ s.t. for all states $s \in \mathcal{S}$, $\mu^{(0)}(s) \geq 0$ denotes the probability of the game being initialized at state s . At each state $s \in \mathcal{S}$ the action $\mathbf{x} \in \mathcal{X}$ chosen by the outer player determines the set of **feasible** actions $\{\mathbf{y} \in \mathcal{Y} \mid \mathbf{g}(s, \mathbf{x}, \mathbf{y}) \geq \mathbf{0}\}$ that the inner player can in turn play. Once the outer and inner players make their moves, they receive payoffs $r_x : \mathcal{S} \times \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$ and $r_y : \mathcal{S} \times \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$, respectively, and the game either ends with probability $1 - \gamma$, where $\gamma \in (0, 1)$ is called the **discount factor**, or transitions to a new state $s' \in \mathcal{S}$, according to a **transition** probability function $p : \mathcal{S} \times \mathcal{S} \times \mathcal{X} \times \mathcal{Y} \rightarrow [0, 1]$ s.t. $p(s' \mid s, \mathbf{x}, \mathbf{y}) \in [0, 1]$ denotes the probability of transitioning to state $s' \in \mathcal{S}$ from state $s \in \mathcal{S}$ when a strategy profile $(\mathbf{x}, \mathbf{y}) \in \mathcal{X} \times \mathcal{Y}$ is chosen by the players.

In this paper, we focus on stochastic Stackelberg **zero-sum** games $\mathcal{G}^{(0)} = (\mathcal{S}, \mathcal{X}, \mathcal{Y}, \mu^{(0)}, r, \mathbf{g}, p, \gamma)$, in which the outer player's loss is the inner player's gain, i.e., $r_x = -r_y$. We say that a zero-sum stochastic Stackelberg game is **convex-concave** if, for all $s \in \mathcal{S}$, $r(s, \mathbf{x}, \mathbf{y}), g_1(s, \mathbf{x}, \mathbf{y}), \dots, g_d(s, \mathbf{x}, \mathbf{y})$ are convex-concave in (\mathbf{x}, \mathbf{y}) . A zero-sum stochastic Stackelberg game reduces to zero-sum stochastic game (Shapley, 1953) if for all state-action tuples $(s, \mathbf{x}, \mathbf{y}) \in \mathcal{S} \times \mathcal{X} \times \mathcal{Y}$, $\mathbf{g}(s, \mathbf{x}, \mathbf{y}) \geq \mathbf{0}$. As we will show, it suffices to focus on deterministic policies in which a **policy** for the outer (resp. inner) player is a mapping from states to actions $\pi_x : \mathcal{S} \rightarrow \mathcal{X}$ (resp. $\pi_y : \mathcal{S} \rightarrow \mathcal{Y}$). The **outcome** of a zero-sum stochastic Stackelberg game $\mathcal{G}^{(0)}$ is a policy profile $(\pi_x, \pi_y) \in \mathcal{X}^{\mathcal{S}} \times \mathcal{Y}^{\mathcal{S}}$ consisting of policies for the outer and inner players, respectively. An outcome $(\pi_x, \pi_y) \in \mathcal{X}^{\mathcal{S}} \times \mathcal{Y}^{\mathcal{S}}$ is said to be **feasible** if, for all states $s \in \mathcal{S}$, $\mathbf{g}(s, \pi_x(s), \pi_y(s)) \geq \mathbf{0}$. For simplicity, we introduce a function $\mathbf{G} : \mathcal{X}^{\mathcal{S}} \times \mathcal{Y}^{\mathcal{S}} \rightarrow \mathbb{R}^{|\mathcal{S}| \times d}$ such that $\mathbf{G}(\pi_x, \pi_y) = (\mathbf{g}(s, \pi_x(s), \pi_y(s)))_{s \in \mathcal{S}}$, and define feasible outcomes as those $(\pi_x, \pi_y) \in \mathcal{X}^{\mathcal{S}} \times \mathcal{Y}^{\mathcal{S}}$ s.t. $\mathbf{G}(\pi_x, \pi_y) \geq \mathbf{0}$. For the rest of this paper, we assume:

Assumption 1.1. 1. For all states $s \in \mathcal{S}$, the functions $r(s, \cdot, \cdot), g_1(s, \cdot, \cdot), \dots, g_d(s, \cdot, \cdot)$ are continuous in $(\mathbf{x}, \mathbf{y}) \in \mathcal{X} \times \mathcal{Y}$ with payoffs r bounded, i.e., $\|r\|_{\infty} \leq \alpha < \infty$, for some $\alpha \in \mathbb{R}_+$, and 2 \mathcal{X}, \mathcal{Y} are non-empty and compact.

Given a zero-sum stochastic Stackelberg game $\mathcal{G}^{(0)}$, the **state-value function**, $v : \mathcal{S} \times \mathcal{X}^{\mathcal{S}} \times \mathcal{Y}^{\mathcal{S}} \rightarrow \mathbb{R}$, and the **action-value function**, $q : \mathcal{S} \times \mathcal{X} \times \mathcal{Y} \times \mathcal{X}^{\mathcal{S}} \times \mathcal{Y}^{\mathcal{S}} \rightarrow \mathbb{R}$ are respectively defined as: $v(s; \pi_x, \pi_y) = \mathbb{E}^{\pi_x, \pi_y}_{\mathbf{S}^{(t+1)} \sim p(\cdot \mid \mathbf{S}^{(t)}, \mathbf{X}^{(t)}, \mathbf{Y}^{(t)})} [\sum_{t=0}^{\infty} (1 - \gamma) \gamma^t r(\mathbf{S}^{(t)}, \mathbf{X}^{(t)}, \mathbf{Y}^{(t)}) \mid \mathbf{S}^{(0)} = s]$; $q(s, \mathbf{x}, \mathbf{y}; \pi_x, \pi_y) = \mathbb{E}^{\pi_x, \pi_y}_{\mathbf{S}^{(t+1)} \sim p(\cdot \mid \mathbf{S}^{(t)}, \mathbf{X}^{(t)}, \mathbf{Y}^{(t)})} [\sum_{t=0}^{\infty} (1 - \gamma) \gamma^t r(\mathbf{S}^{(t)}, \mathbf{X}^{(t)}, \mathbf{Y}^{(t)}) \mid \mathbf{S}^{(0)} = s, \mathbf{X}^{(0)} = \mathbf{x}, \mathbf{Y}^{(0)} = \mathbf{y}]$.

For clarity, we write expectations conditional on $\mathbf{X}^{(t)} = \pi_x(\mathbf{S}^{(t)})$ and $\mathbf{Y}^{(t)} = \pi_y(\mathbf{S}^{(t)})$ as $\mathbb{E}^{\pi_x, \pi_y}$, and denote the state- and action-value functions by $v^{\pi_x, \pi_y}(s)$, and $q^{\pi_x, \pi_y}(s, \mathbf{x}, \mathbf{y})$, respectively. Additionally, we let $\mathcal{V} = [-\alpha, \alpha]^{\mathcal{S}}$ be the space of all state-value functions of the form $v : \mathcal{S} \rightarrow [-\alpha, \alpha]$, and we let $\mathcal{Q} = [-\alpha, \alpha]^{\mathcal{S} \times \mathcal{X} \times \mathcal{Y}}$ be the space of all action-value functions of the form $q : \mathcal{S} \times \mathcal{X} \times \mathcal{Y} \rightarrow [-\alpha, \alpha]$. Note that by Assumption 1.1 the range of

the state- and action-value functions is $[-\alpha, \alpha]$. The cumulative payoff function of the game $u : \mathcal{X}^S \times \mathcal{Y}^S \rightarrow \mathbb{R}$ is the total expected loss (resp. gain) of the outer (resp. inner) player is then given by $u(\pi_x, \pi_y) = \mathbb{E}_{s \sim \mu^{(0)}(s)} [v^{\pi_x \pi_y}(s)]$.

Definition 1.2. A feasible policy profile $(\pi_x^*, \pi_y^*) \in \mathcal{X}^S \times \mathcal{Y}^S$ is said to be a **recursive Stackelberg equilibrium (recSE)** of a zero-sum stochastic Stackelberg game $\mathcal{G}^{(0)}$ iff $\max_{\pi_y \in \mathcal{Y}^S: \mathbf{G}(\pi_x^*, \pi_y) \geq \mathbf{0}} u(\pi_x^*, \pi_y) \leq u(\pi_x^*, \pi_y^*) \leq \min_{\pi_x \in \mathcal{X}^S} \max_{\pi_y \in \mathcal{Y}^S: \mathbf{G}(\pi_x, \pi_y) \geq \mathbf{0}} u(\pi_x, \pi_y)$.

Remark 1.3. This is a strong definition of a recSE, since we require the constraints $\mathbf{g}(s, \pi_x, \pi_y) \geq \mathbf{0}$ to be satisfied at all states $s \in \mathcal{S}$, not only states which are reached with strictly positive probability.

Mathematical Preliminaries A probability measure $q_1 \in \Delta(\mathcal{S})$ **convex stochastically dominates (CSD)** $q_2 \in \Delta(\mathcal{S})$ if $\int_{\mathcal{S}} v(s)q_1(s)ds \geq \int_{\mathcal{S}} v(s)q_2(s)ds$ for all continuous, bounded, and convex functions v on \mathcal{S} . A transition function p is termed **CSD convex** in \mathbf{x} if for all $\lambda \in (0, 1)$, $\mathbf{y} \in \mathcal{Y}$ and any (s', \mathbf{x}') , $(s^\dagger, \mathbf{x}^\dagger) \in \mathcal{S} \times \mathcal{X}$, with $(s, \mathbf{x}) = \lambda(s', \mathbf{x}') + (1 - \lambda)(s^\dagger, \mathbf{x}^\dagger)$, we have $\lambda p(\cdot | s', \mathbf{x}', \mathbf{y}) + (1 - \lambda)p(\cdot | s^\dagger, \mathbf{x}^\dagger, \mathbf{y})$ CSD $p(\cdot | s, \mathbf{x}, \mathbf{y})$. A transition function p is termed **CSD concave** in \mathbf{y} if for all $\lambda \in (0, 1)$ and any (s', \mathbf{y}') , $(s^\dagger, \mathbf{y}^\dagger) \in \mathcal{S} \times \mathcal{X} \times \mathcal{Y}$, with $(s, \mathbf{y}) = \lambda(s', \mathbf{y}') + (1 - \lambda)(s^\dagger, \mathbf{y}^\dagger)$, we have $p(\cdot | s, \mathbf{x}, \mathbf{y})$ CSD $\lambda p(\cdot | s', \mathbf{x}, \mathbf{y}') + (1 - \lambda)p(\cdot | s^\dagger, \mathbf{x}, \mathbf{y}^\dagger)$. A mapping $L : \mathcal{A} \rightarrow \mathcal{B}$ is said to be a **contraction mapping** (resp. **non-expansion**) w.r.t. norm $\|\cdot\|$ iff for all $\mathbf{x}, \mathbf{y} \in \mathcal{A}$, and for $k \in [0, 1)$ (resp. $k = 1$) such that $\|L(\mathbf{x}) - L(\mathbf{y})\| \leq k \|\mathbf{x} - \mathbf{y}\|$.

2 PROPERTIES OF RECURSIVE STACKELBERG EQUILIBRIUM

We first state a necessary condition that the state-value function induced by a policy profile needs to satisfy in order for that policy profile to be a recSE. It is not directly clear if there exists a policy for any zero-sum stochastic Stackelberg game which satisfies this necessary condition since there is no way to guarantee that such an inequality will hold at all states.²

Lemma 2.1. Consider a zero-sum stochastic Stackelberg game $\mathcal{G}^{(0)}$. A policy profile $(\pi_x^*, \pi_y^*) \in \mathcal{X}^S \times \mathcal{Y}^S$ is a recSE if, for all $s \in \mathcal{S}$, we have that: $\max_{\mathbf{y} \in \mathcal{Y}: \mathbf{g}(s, \pi_x^*(s), \mathbf{y}) \geq \mathbf{0}} q^{\pi_x^* \pi_y^*}(s, \pi_x^*(s), \mathbf{y}) \leq q^{\pi_x^* \pi_y^*}(s, \pi_x^*(s), \pi_y^*(\mathbf{y})) \leq \min_{\mathbf{x} \in \mathcal{X}} \max_{\mathbf{y} \in \mathcal{Y}: \mathbf{g}(s, \mathbf{x}, \mathbf{y}) \geq \mathbf{0}} q^{\pi_x^* \pi_y^*}(s, \mathbf{x}, \mathbf{y})$. Equivalently, a policy profile (π_x^*, π_y^*) is a recSE if $(\pi_x^*(s), \pi_y^*(s))$ is a Stackelberg equilibrium: i.e., a solution to $\min_{\mathbf{x} \in \mathcal{X}} \max_{\mathbf{y} \in \mathcal{Y}: \mathbf{g}(s, \mathbf{x}, \mathbf{y}) \geq \mathbf{0}} q^{\pi_x^* \pi_y^*}(s, \mathbf{x}, \mathbf{y})$ at each $s \in \mathcal{S}$.

We define an operator $C : \mathcal{V} \rightarrow \mathcal{V}$ associated with a zero-sum stochastic Stackelberg game $\mathcal{G}^{(0)}$ whose fixed points satisfy the condition given in Lemma 2.1, and hence correspond to the value function associated with a recSE of $\mathcal{G}^{(0)}$. We then show that this operator is a contraction mapping, thereby establishing the existence of such a fixed point. This result generalizes a result first shown by Shapley (1953) for zero-sum stochastic games, i.e., zero-sum stochastic Stackelberg games in which $\mathbf{G}(\pi_x, \pi_y) \geq \mathbf{0}$, for all $(\pi_x, \pi_y) \in \mathcal{X}^S \times \mathcal{Y}^S$. Define $C : \mathcal{V} \rightarrow \mathcal{V}$ for a stochastic Stackelberg game $\mathcal{G}^{(0)}$ as the operator $(Cv)(s) = \min_{\mathbf{x} \in \mathcal{X}} \max_{\mathbf{y} \in \mathcal{Y}: \mathbf{g}(s, \mathbf{x}, \mathbf{y}) \geq \mathbf{0}} \mathbb{E}_{S' \sim p(\cdot | s, \mathbf{x}, \mathbf{y})} [r(s, \mathbf{x}, \mathbf{y}) + \gamma v(S')]$.

Theorem 2.2. (π_x^*, π_y^*) is a recSE of $\mathcal{G}^{(0)}$ of $v^{\pi_x^* \pi_y^*}$ iff it induces a value function which is a fixed point of C : i.e., (π_x^*, π_y^*) is a Stackelberg equilibrium iff, for all $s \in \mathcal{S}$, $(Cv^{\pi_x^* \pi_y^*})(s) = v^{\pi_x^* \pi_y^*}(s)$.

The following technical lemma is crucial to proving that C is a contraction mapping. It tells us that the generalized min-max operator is non-expansive; in other words, the generalized min-max operator is 1-Lipschitz w.r.t. the sup-norm.

Lemma 2.3. Suppose that $f, h : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$, $\mathbf{g} : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}^d$ are continuous functions, and \mathcal{X}, \mathcal{Y} are compact sets, we then have: $|\min_{\mathbf{x} \in \mathcal{X}} \max_{\mathbf{y} \in \mathcal{Y}: \mathbf{g}(\mathbf{x}, \mathbf{y}) \geq \mathbf{0}} f(\mathbf{x}, \mathbf{y}) - \min_{\mathbf{x} \in \mathcal{X}} \max_{\mathbf{y} \in \mathcal{Y}: \mathbf{g}(\mathbf{x}, \mathbf{y}) \geq \mathbf{0}} h(\mathbf{x}, \mathbf{y})| \leq \max_{(\mathbf{x}, \mathbf{y}) \in \mathcal{X} \times \mathcal{Y}} |f(\mathbf{x}, \mathbf{y}) - h(\mathbf{x}, \mathbf{y})|$.

With the above lemma in hand, we can now prove that C is a contraction mapping.

²All omitted results and proofs can be found in the appendix.

Theorem 2.4. Consider the operator C associated with a stochastic Stackelberg game $\mathcal{G}^{(0)}$. If Assumption 1.1 holds, then C is a contraction mapping w.r.t. to the sup norm $\|\cdot\|_\infty$ with constant γ .

Given some initial state-value function $v^{(0)} \in \mathcal{V}$, we define the **value iteration** process as $v^{(t+1)} = Cv^{(t)}$, for all $t \in \mathbb{N}_+$ (Algorithm 1). One way to interpret $v^{(t)}$ is as the function that returns the value $v^{(t)}(s)$ of each state $s \in \mathcal{S}$ in the t stage zero-sum Stackelberg game starting at the last stage t and continuing until stage 0, with terminal payoffs given by $v^{(0)}$. The following theorem, which is a consequence of Theorems 2.2 and 2.4, not only states the existence of a recSE but also provides us with a means of computing a recSE via value iteration.

Theorem 2.5. Consider a zero-sum stochastic Stackelberg game $\mathcal{G}^{(0)}$. If Assumption 1.1 holds, then $\mathcal{G}^{(0)}$ has at least one recSE (π_x^*, π_y^*) with unique value function $v^{\pi_x^* \pi_y^*}$. Further, $v^{\pi_x^* \pi_y^*}$ can be computed by iteratively applying C to any initial state-value function $v^{(0)} \in \mathcal{V}$: $\lim_{t \rightarrow \infty} v^{(t)} = v^{\pi_x^* \pi_y^*}$.

Remark 2.6. Unlike Shapley’s existence theorem for recursive Nash equilibria in zero-sum stochastic games, the above theorem does not require that the payoff function be convex-concave. The only conditions needed are continuity of the payoffs and constraints, and bounded payoffs. This makes the recSE a potentially useful solution concept, even for non-convex-non-concave stochastic games.

Note that the proof of Theorem 2.5 shows that there exists a policy profile which satisfies the conditions of Lemma 2.1. Since this recSE definition is independent of the initial state distribution, we can infer that the recSE of any zero-sum Stackelberg game $\mathcal{G}^{(0)} = (\mathcal{S}, \mathcal{X}, \mathcal{Y}, \mu^{(0)}, r, \mathbf{g}, p, \gamma)$ is independent of the initial state distribution $\mu^{(0)}$. Hence, from now on, we denote a zero-sum Stackelberg game by \mathcal{G} .

It seems that a stronger equilibrium concept for zero-sum stochastic Stackelberg games would be one that requires the strategy profile be a recSE of each subgame, aptly called **subgame perfect Stackelberg equilibrium**. A feasible policy profile $(\pi_x^*, \pi_y^*) \in \mathcal{S}^{\mathcal{X}} \times \mathcal{S}^{\mathcal{Y}}$ is a subgame perfect Stackelberg equilibrium of a zero-sum stochastic Stackelberg game \mathcal{G} if, for all $s \in \mathcal{S}, t \in \mathbb{N}, \mu^{(t+1)}(s) = \sum_{s' \in \mathcal{S}} p(s | s', \pi_x^*(s'), \pi_y^*(s')) \mu^{(t)}(s')$, (π_x^*, π_y^*) is a recSE of $(\mathcal{S}, \mathcal{X}, \mathcal{Y}, \mu^{(t)}, r, \mathbf{g}, p, \gamma)$. By definition, the set of subgame perfect Stackelberg equilibria of a zero-sum stochastic Stackelberg game is a subset of the set of recSE of the game. Since the set of recSE is independent of the initial state distribution $\mu^{(0)}$, the reverse inclusion also holds. We thus have the following corollary:

Corollary 2.7. The set of recSE and subgame perfect Stackelberg equilibria coincide in zero-sum stochastic Stackelberg games.

Theorem 2.5 tells us that value iteration converges to the value function associated with a recSE. Additionally, under Assumption 1.1, recSE is computable in (weakly) polynomial time.³

Theorem 2.8. [Convergence of Value Iteration] Suppose value iteration is run on input \mathcal{G} . If Assumption 1.1 holds, and if we initialize $v^{(0)}(s) = 0$, for all $s \in \mathcal{S}$, then for $k \geq \frac{1}{1-\gamma} \log \frac{2\alpha}{\epsilon(1-\gamma)}$, we have $v^{(k)}(s) - v^{\pi_x^* \pi_y^*}(s) \leq \epsilon$.

3 SUBDIFFERENTIAL ENVELOPE THEOREMS AND OPTIMALITY CONDITIONS FOR RECURSIVE STACKELBERG EQUILIBRIUM

In this section, we derive optimality conditions for recursive Stackelberg equilibria. In particular, we provide necessary conditions for a policy profile to be a recSE of any zero-sum stochastic Stackelberg game, and we provide sufficient conditions for a policy profile to be a recursive Stackelberg equilibrium of a zero-sum convex-concave stochastic Stackelberg game. Using these results, we prove in the next section that recursive market equilibrium (Mehra & Prescott, 1977) is an instance of recSE in a large class of stochastic markets.

³This convergence is only weakly polynomial time, since the computation of the generalized min-max operator applied to an arbitrary continuous function is an NP-hard problem; it is at least as hard as non-convex optimization. If, however, one restricts themselves to convex-concave stochastic Stackelberg games, Stackelberg equilibrium is computable in polynomial time, making the problem much more tractable.

The Benveniste-Scheinman theorem characterizes the derivative of the optimal value function associated with a recursive optimization problem w.r.t. its parameters, when it is differentiable (Benveniste & Scheinkman, 1979). Our proofs of the necessary and sufficient optimality conditions rely on a novel subdifferential generalization (Theorem C.3, Appendix C) of this theorem, which applies even when the optimal value function is not differentiable. A consequence of our subdifferential version of the Benveniste-Scheinman theorem is that we can easily derive the first-order necessary conditions for a policy profile to be a recSE of any stochastic Stackelberg game \mathcal{G} satisfying Assumption 1.1, under standard regularity conditions.

Theorem 3.1. *Consider a zero-sum stochastic Stackelberg game \mathcal{G} , where $\mathcal{X} = \{\mathbf{x} \in \mathbb{R}^n \mid q_1(\mathbf{x}) \leq 0, \dots, q_p(\mathbf{x}) \leq 0\}$ and $\mathcal{Y} = \{\mathbf{y} \in \mathbb{R}^m \mid r_1(\mathbf{y}) \geq 0, \dots, r_l(\mathbf{y}) \geq 0\}$. Let $\mathcal{L}_{s,\mathbf{x}}(\mathbf{y}, \boldsymbol{\lambda}) = r(\mathbf{s}, \mathbf{x}, \mathbf{y}) + \gamma \mathbb{E}_{\mathbf{S}' \sim p(\cdot | \mathbf{s}, \mathbf{x}, \mathbf{y})} [v(\mathbf{S}', \mathbf{x})] + \sum_{k=1}^d \lambda_k g_k(\mathbf{s}, \mathbf{x}, \mathbf{y})$. Suppose that Assumption 1.1 holds, and that 1. for all $\mathbf{s} \in \mathcal{S}, \mathbf{y} \in \mathcal{Y}$, $r(\mathbf{s}, \mathbf{x}, \mathbf{y})$, $g_1(\mathbf{s}, \mathbf{x}, \mathbf{y}), \dots, g_d(\mathbf{s}, \mathbf{x}, \mathbf{y})$ are concave in \mathbf{x} , 2. $\nabla_{\mathbf{x}} r(\mathbf{s}, \mathbf{x}, \mathbf{y}), \nabla_{\mathbf{x}} g_1(\mathbf{s}, \mathbf{x}, \mathbf{y}), \dots, \nabla_{\mathbf{x}} g_d(\mathbf{s}, \mathbf{x}, \mathbf{y})$, $\nabla_{\mathbf{y}} r(\mathbf{s}, \mathbf{x}, \mathbf{y}), \nabla_{\mathbf{y}} g_1(\mathbf{s}, \mathbf{x}, \mathbf{y}), \dots, \nabla_{\mathbf{y}} g_d(\mathbf{s}, \mathbf{x}, \mathbf{y})$ exist, for all $\mathbf{s} \in \mathcal{S}, \mathbf{x} \in \mathcal{X}, \mathbf{y} \in \mathcal{Y}$, 4. $p(\mathbf{s}' | \mathbf{s}, \mathbf{x}, \mathbf{y})$ is continuous CSD convex and differentiable in (\mathbf{x}, \mathbf{y}) , and 5. Slater's condition holds, i.e., $\forall \mathbf{s} \in \mathcal{S}, \mathbf{x} \in \mathcal{X}, \exists \hat{\mathbf{y}} \in \mathcal{Y}$ s.t. $g_k(\mathbf{s}, \mathbf{x}, \hat{\mathbf{y}}) > 0$, for all $k = 1, \dots, d$ and $r_j(\hat{\mathbf{y}}) > 0$, for all $j = 1, \dots, l$, and $\exists \mathbf{x} \in \mathbb{R}^n$ s.t. $q_k(\mathbf{x}) < 0$ for all $k = 1, \dots, p$. Then, there exists $\boldsymbol{\mu}^* : \mathcal{S} \rightarrow \mathbb{R}_+^p, \boldsymbol{\lambda}^* : \mathcal{S} \times \mathcal{X} \rightarrow \mathbb{R}_+^d$, and $\boldsymbol{\nu}^* : \mathcal{S} \times \mathcal{X} \rightarrow \mathbb{R}_+^l$ s.t. a policy profile $(\boldsymbol{\pi}_{\mathbf{x}}^*, \boldsymbol{\pi}_{\mathbf{y}}^*) \in \mathcal{X}^{\mathcal{S}} \times \mathcal{Y}^{\mathcal{S}}$ is a recSE of \mathcal{G} only if it satisfies the following conditions, for all $\mathbf{s} \in \mathcal{S}$:*

$$\nabla_{\mathbf{x}} \mathcal{L}_{\mathbf{s}, \boldsymbol{\pi}_{\mathbf{x}}^*(\mathbf{s})}(\boldsymbol{\pi}_{\mathbf{y}}^*(\mathbf{s}), \boldsymbol{\lambda}^*(\mathbf{s}, \boldsymbol{\pi}_{\mathbf{x}}^*(\mathbf{s}))) + \sum_{k=1}^p \mu_k^*(\mathbf{s}) \nabla_{\mathbf{x}} q_k(\boldsymbol{\pi}_{\mathbf{x}}^*(\mathbf{s})) = 0 \quad (1)$$

$$\nabla_{\mathbf{y}} \mathcal{L}_{\mathbf{s}, \boldsymbol{\pi}_{\mathbf{x}}^*(\mathbf{s})}(\boldsymbol{\pi}_{\mathbf{y}}^*(\mathbf{s}), \boldsymbol{\lambda}^*(\mathbf{s}, \boldsymbol{\pi}_{\mathbf{x}}^*(\mathbf{s}))) + \sum_{k=1}^l \nu_k^*(\mathbf{s}, \boldsymbol{\pi}_{\mathbf{x}}^*(\mathbf{s})) \nabla_{\mathbf{x}} r_k(\boldsymbol{\pi}_{\mathbf{y}}^*(\mathbf{s})) = 0 \quad (2)$$

$$\mu_k^*(\mathbf{s}) q_k(\boldsymbol{\pi}_{\mathbf{x}}^*(\mathbf{s})) = 0 \quad q_k(\boldsymbol{\pi}_{\mathbf{x}}^*(\mathbf{s})) \leq 0 \quad \forall k \in [p] \quad (3)$$

$$g_k(\mathbf{s}, \boldsymbol{\pi}_{\mathbf{x}}^*(\mathbf{s}), \boldsymbol{\pi}_{\mathbf{y}}^*(\mathbf{s})) \geq 0 \quad \lambda_k^*(\mathbf{s}, \boldsymbol{\pi}_{\mathbf{x}}^*(\mathbf{s})) g_k(\mathbf{s}, \boldsymbol{\pi}_{\mathbf{x}}^*(\mathbf{s}), \boldsymbol{\pi}_{\mathbf{y}}^*(\mathbf{s})) = 0 \quad \forall k \in [d] \quad (4)$$

$$\nu_k^*(\mathbf{s}, \boldsymbol{\pi}_{\mathbf{x}}^*(\mathbf{s})) \nabla_{\mathbf{x}} r_k(\boldsymbol{\pi}_{\mathbf{y}}^*(\mathbf{s})) = 0 \quad r_k(\boldsymbol{\pi}_{\mathbf{x}}^*(\mathbf{s})) \geq 0 \quad \forall k \in [l] \quad (5)$$

Note, that in Theorem 3.1, when one additionally assumes concavity of $\max_{\mathbf{y} \in \mathcal{Y}: g(\mathbf{s}, \mathbf{x}, \mathbf{y}) \geq 0} \{r(\mathbf{s}, \mathbf{x}, \mathbf{y}) + \gamma \mathbb{E}_{\mathbf{S}' \sim p(\cdot | \mathbf{s}, \mathbf{x}, \mathbf{y})} [v(\mathbf{S}', \mathbf{x})]\}$ in \mathbf{x} , Equations (74) to (78) become necessary and sufficient optimality conditions. For completeness, the reader can find the necessary and sufficient optimality conditions for convex-concave stochastic Stackelberg games under standard regularity conditions in Theorem C.4 (Appendix C). The proof follows exactly as that of Theorem 2.2.

4 RECURSIVE MARKET EQUILIBRIUM

We now introduce an application of zero-sum stochastic Stackelberg games, which generalizes a well known market model, the Fisher market (Brainard et al., 2000), to a dynamic setting in which buyers not only participate in markets across time, but their wealth persists. A **(static) Fisher market** consists of n buyers and m divisible goods (Brainard et al., 2000). Each buyer $i \in [n]$ is endowed with a budget $b_i \in \mathcal{B}_i \subset \mathbb{R}_+$ and a utility function $u_i : \mathbb{R}_+^m \times \mathcal{T}_i \rightarrow \mathbb{R}$, which is parameterized by a type $t_i \in \mathcal{T}_i$ that defines a preference relation over the consumption space \mathbb{R}_+^m . Each good is characterized by a supply $q_j \in \mathcal{Q}_j \subset \mathbb{R}_+$. A **stochastic Fisher market** is a dynamic market in which each state corresponds to a static Fisher market: i.e., each state $\mathbf{s} \in \mathcal{S}$ is characterized by a tuple $(\mathbf{t}, \mathbf{b}, \mathbf{q})$. In each state, the buyers choose their allocations $\mathbf{X} = (\mathbf{x}_1, \dots, \mathbf{x}_n)^T \in \mathbb{R}_+^{n \times m}$ and the market determines prices, after which the market comes to an end with probability $(1 - \gamma)$, or it moves into a new state with probability $p(\mathbf{s}' | \mathbf{s}, \mathbf{X})$.⁴ A **stochastic Fisher market with savings** is a stochastic Fisher market in which, at each state, each buyer i , in addition to choosing their allocation, can set aside some **savings** $\beta_i \in \mathbb{R}_+$ to spend at some future state, and transitions depend on their savings choices as well: i.e., $p(\mathbf{s}' | \mathbf{s}, \mathbf{X}, \boldsymbol{\beta})$.

⁴Note that, as is standard in the literature, we assume that prices do not determine the next state since market prices are set by a “fictional auctioneer.” There is no market participant who determined prices!

A stochastic Fisher market with savings is denoted by $(\mathcal{S}, \mathcal{U}, \mathbf{b}^{(0)}, p, \gamma)$. Given such a market, a **recursive competitive equilibrium (recCE)** (Mehra & Prescott, 1977) is a tuple $(\mathbf{X}^*, \beta^*, \mathbf{p}^*) \in \mathbb{R}_+^{n \times m \times \mathcal{S}} \times \mathbb{R}_+^{n \times \mathcal{S}} \times \mathbb{R}_+^{m \times \mathcal{S}}$, which consists of an allocation, savings, and price system s.t. 1) the buyers are expected utility maximizing, constrained by their savings and spending constraints, i.e., for all buyers $i \in [n]$, $(\mathbf{x}_i^*, \beta_i^*)$ is the optimal policy that, for all states $(\mathbf{t}, \mathbf{b}, \mathbf{q}) \in \mathcal{S}$, solves the Bellman equation $\nu_i(\mathbf{t}, \mathbf{b}, \mathbf{q}) =$

$$\max_{(\mathbf{x}_i, \beta_i) \in \mathbb{R}_+^{m+1}} \max_{\mathbf{x}_i, \mathbf{p}^*(\mathbf{t}, \mathbf{b}, \mathbf{q}) + \beta_i \leq \mathbf{b}_i} \left\{ u_i(\mathbf{x}_i, t_i) + \gamma \mathbb{E}_{(\mathbf{t}', \mathbf{b}', \mathbf{q}') \sim p(\cdot | \mathbf{t}, \mathbf{b}, \mathbf{q}, (\mathbf{x}_i, \mathbf{X}_{-i}^*(\mathbf{s}), (\beta_i, \beta_{-i}^*(\mathbf{s})))} [\nu_i(\mathbf{t}', \mathbf{b}' + \beta_i, \mathbf{q}')] \right\},$$

where $\mathbf{X}_{-i}^*, \beta_{-i}^*$ denote the allocation and saving systems excluding buyer i , and 2) the market clears in each state so that unallocated goods in each state are priced at 0, i.e., for all $j \in [m]$ and $\mathbf{s} \in \mathcal{S}$, $p_j^*(\mathbf{t}, \mathbf{b}, \mathbf{q}) > 0 \implies \sum_{i \in [n]} x_{ij}^*(\mathbf{t}, \mathbf{b}, \mathbf{q}) = 1$ and $p_j^*(\mathbf{t}, \mathbf{b}, \mathbf{q}) \geq 0 \implies \sum_{i \in [n]} x_{ij}^*(\mathbf{t}, \mathbf{b}, \mathbf{q}) \leq 1$.

The following theorem shows that the recSE of a stochastic Fisher market with savings are in fact recursive competitive equilibria. We note that if one ignores the savings terms, then the recSE are recursive competitive equilibria of the same market without savings.

Theorem 4.1. *A stochastic Fisher market with savings $(\mathcal{S}, \mathcal{U}, \mathbf{b}^{(0)}, p, \gamma)$ in which \mathcal{U} is a vector of continuous and homogeneous utility functions has at least one recCE. Additionally, the recSE $(\mathbf{p}^*, \mathbf{X}^*, \beta^*)$ that solves the following Bellman equation corresponds to the recCE of $(\mathcal{S}, \mathcal{U}, \mathbf{b}^{(0)}, p, \gamma)$:*

$$v(\mathbf{t}, \mathbf{b}, \mathbf{q}) = \min_{\mathbf{p} \in \mathbb{R}_+^m} \max_{(\mathbf{X}, \beta) \in \mathbb{R}_+^{n \times (m+1)}: \mathbf{X}\mathbf{p} + \beta \leq \mathbf{b}} \sum_{j \in [m]} q_j p_j + \sum_{i \in [n]} (b_i - \beta_i) \log(u_i(\mathbf{x}_i, t_i)) + \gamma \mathbb{E}_{(\mathbf{t}', \mathbf{b}', \mathbf{q}') \sim p(\cdot | \mathbf{t}, \mathbf{b}, \mathbf{q}, \mathbf{X}, \beta)} [v(\mathbf{t}', \mathbf{b}' + \beta, \mathbf{q}')] \quad (6)$$

Remark 4.2. *This result could not be obtained by modifying the Lagrangian formulation, i.e., the simultaneous-move game form, of the Eisenberg-Gale program, because the saving problem is a convex-non-concave problem, and existence of recursive Nash equilibrium in deterministic policies requires convex-concavity of payoffs (Jaśkiewicz & Nowak, 2018).*

5 EXPERIMENTS

In order to better understand the iteration complexity of value iteration, and to understand how it performs in a continuous state space, we computed the recursive Stackelberg equilibria of various stochastic Fisher market with savings. We created markets with three different classes of utility functions, each of which endowed the state-value function with different smoothness properties.⁵ Let $\theta_i \in \mathbb{R}^m$ be a vector of parameters that describes the utility function of buyer $i \in [n]$. We considered the following (standard) utility function classes: 1. linear: $u_i(\mathbf{x}_i) = \sum_{j \in [m]} \theta_{ij} x_{ij}$; 2. Cobb-Douglas: $u_i(\mathbf{x}_i) = \prod_{j \in [m]} x_{ij}^{\theta_{ij}}$; and 3. Leontief: $u_i(\mathbf{x}_i) = \min_{j \in [m]} \left\{ \frac{x_{ij}}{\theta_{ij}} \right\}$.

Since the state space is continuous, the value function is also continuous in stochastic Fisher markets. As a result, we used fitted value iteration, finding a fit via linear regression (e.g., Boyan & Moore (1994)). To solve the generalized min-max operator at each step of value iteration, we used two methods: 1. **nested gradient descent ascent (GDA)** (Goktas & Greenwald (2021); Algorithm 2), which is not guaranteed to converge to a global optimum since the zero-sum Stackelberg game for stochastic Fisher markets is convex-non-concave; and 2. **max-oracle gradient descent** (Goktas & Greenwald (2021); Algorithm 3), where we used simulated annealing (Bertsimas & Tsitsiklis, 1993), a metaheuristic which aims to find a global optimum, as the max-oracle. Although simulated annealing is not guaranteed to converge to a global optimum, we observed that it outperformed nested GDA, more often finding a global maximum for the inner player.

To check whether the value function computed was optimal, we measured the exploitability of the market, meaning the distance between the recursive competitive equilibrium computed and the actual competitive equilibrium. To do so required that we check two conditions: 1) if each buyer's expected utility is maximized at the computed allocation and saving system at the price system output by the algorithm, and 2) if the market always clears. For both settings, given the value

⁵Our code can be found here, and details of our experimental setup can be found in Appendix E.

function computed by value iteration, we extracted the greedy policy and unrolled it across time to obtain the greedy actions $(\mathbf{X}^{(t)}, \beta^{(t)}, \mathbf{p}^{(t)})$ at each state $s^{(t)}$. We then computed the cumulative utility of the allocation and saving systems computed by the algorithms, i.e., for all $i \in [n]$, $\sum_{t=0}^T \gamma^t u_i(\mathbf{x}_i^{(t)})$, and compared this value to the expected maximum utility u_i^* , obtained by solving the consumption/saving problem for individual buyers given the price systems computed by our algorithms. We report the normalized distance between these two values: e.g., in the case of two buyers, we report $\frac{\|(u_1, u_2) - (u_1^*, u_2^*)\|}{\|(u_1^*, u_2^*)\|}$. Finally, we measured the excess demand, which we took as the distance to market clearance, i.e., $\frac{1}{T} \sum_{t=1}^T \|\sum_{i \in [n]} \mathbf{x}_i^{(t)} - \mathbf{q}^{(t)}\|$.

The exploitability of the recursive competitive equilibrium computed by both nested GDA and max-oracle gradient descent in all market types is shown in Figure 2, while Figure 1 depicts the average value of the value function across all states as it varies with time. We observe that value iteration converges for both nested GDA and max-oracle gradient descent in Cobb-Douglas markets. Max-oracle gradient descent achieves a higher value, because the oracle does a better job at finding a global maximum for the inner player. In linear markets, only the max-oracle gradient descent algorithm seems to converge. This is most likely due to the non-differentiability of the payoffs in linear Fisher markets, making it difficult for nested GDA to find a global optimum. In Leontief markets, neither algorithm seems to affect the value function very much, which might lead one to conclude that value iteration converges. In terms of exploitability, both algorithms seem to perform more or less as they do in the linear markets, which suggests that the initial value function is close to the optimal value function. Surprisingly, even though value iteration converges in both settings for Cobb-Douglas markets, the allocation and savings computed seem to be far from optimal. This is most likely due to the fact that we used linear regression to approximate the value function at each state, even though the value function associated with each individual buyer’s consumption/saving problem, and likewise the market’s value function, is strictly concave in Cobb-Douglas markets. We note that the buyers’ value functions are linear in linear and Leontief markets, and hence compatible with our fitted linear regression, which might explain the greater exploitability of Cobb-Douglas markets even though value iteration seems to converge.

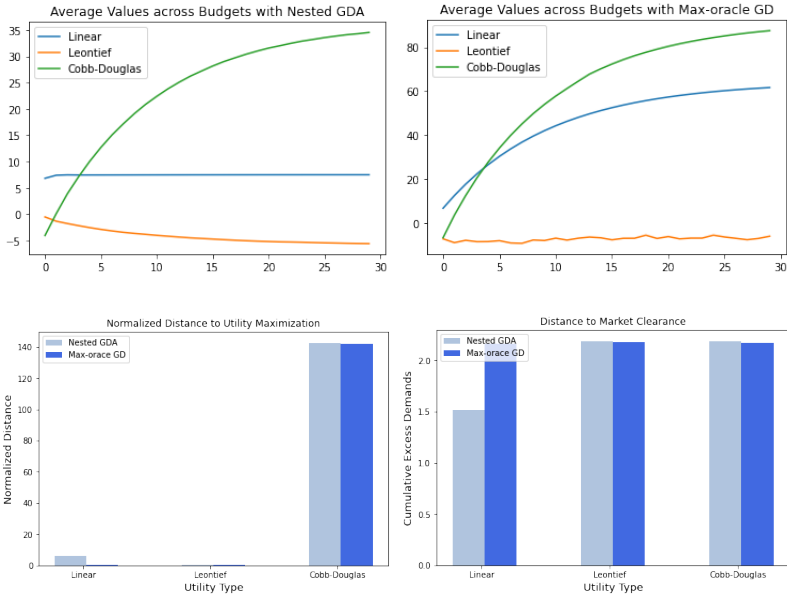


Figure 1: The average value of the value function under Nested GDA (left) and Max-oracle gradient descent with simulated annealing (right).

Figure 2: Distance to utility maximization (left) and market clearance (right) of the computed recCE.

6 CONCLUSION

In this paper, we proved the existence of recursive Stackelberg equilibria in zero-sum stochastic Stackelberg games, provided necessary and sufficient conditions for a policy profile to be a recursive Stackelberg equilibrium, and showed that a Stackelberg equilibrium can be computed

in (weakly) polynomial time via value iteration. Further, we provided a characterization of recursive competitive equilibria in stochastic Fisher markets with savings via zero-sum stochastic Stackelberg games. Finally, we used value iteration to solve for recursive competitive equilibria in stochastic Fisher markets. Future work seeking to improve our solution methods for stochastic Fisher markets should aim to replace the linear regression step in our fitted value iteration method with concave regression, since the value function is guaranteed to be concave in budgets for stochastic Fisher markets. Additionally, it is conceivable that policy-network-based deep reinforcement learning methods may be able to resolve the difficulties that our method has in solving for global optima in stochastic Fisher markets.

REFERENCES

- Mohammad Alkousa, Darina Dvinskikh, Fedor Stonyakin, Alexander Gasnikov, and Dmitry Kovalev. Accelerated methods for composite non-bilinear saddle point problem, 2020.
- Kenneth Arrow and Gerard Debreu. Existence of an equilibrium for a competitive economy. *Econometrica: Journal of the Econometric Society*, pp. 265–290, 1954.
- Alp E. Atakan. Stochastic convexity in dynamic programming. *Economic Theory*, 22(2):447–455, 2003. ISSN 09382259, 14320479. URL <http://www.jstor.org/stable/25055693>.
- Stefan Banach. Sur les opérations dans les ensembles abstraits et leur application aux équations intégrales. *Fund. math*, 3(1):133–181, 1922.
- Richard Bellman. On the theory of dynamic programming. *Proceedings of the National Academy of Sciences of the United States of America*, 38(8):716, 1952.
- Alain Bensoussan, Shaokuan Chen, and Suresh P Sethi. The maximum principle for global solutions of stochastic stackelberg differential games. *SIAM Journal on Control and Optimization*, 53(4): 1956–1981, 2015.
- L. M. Benveniste and J. A. Scheinkman. On the differentiability of the value function in dynamic models of economics. *Econometrica*, 47(3):727–732, 1979. ISSN 00129682, 14680262. URL <http://www.jstor.org/stable/1910417>.
- Dimitris Bertsimas and John Tsitsiklis. Simulated annealing. *Statistical science*, 8(1):10–15, 1993.
- Justin Boyan and Andrew Moore. Generalization in reinforcement learning: Safely approximating the value function. *Advances in neural information processing systems*, 7, 1994.
- William C Brainard, Herbert E Scarf, et al. *How to compute equilibrium prices in 1891*. Citeseer, 2000.
- Lv Chen and Yang Shen. On a new paradigm of optimal reinsurance: a stochastic stackelberg differential game between an insurer and a reinsurer. *ASTIN Bulletin: The Journal of the IAA*, 48 (2):905–960, 2018.
- Yun Kuen Cheung, Martin Hoefer, and Paresh Nakhe. Tracing equilibrium in dynamic markets via distributed adaptation. In *Proceedings of the 18th International Conference on Autonomous Agents and MultiAgent Systems*, pp. 1225–1233, 2019.
- Victor DeMiguel and Huifu Xu. A stochastic multiple-leader stackelberg model: analysis, computation, and application. *Operations Research*, 57(5):1220–1235, 2009.
- N. R. Devanur, C. H. Papadimitriou, A. Saberi, and V. V. Vazirani. Market equilibrium via a primal-dual-type algorithm. In *The 43rd Annual IEEE Symposium on Foundations of Computer Science, 2002. Proceedings.*, pp. 389–395, 2002. doi: 10.1109/SFCS.2002.1181963.
- Harley Flanders. Differentiation under the integral sign. *The American Mathematical Monthly*, 80 (6):615–627, 1973. ISSN 00029890, 19300972. URL <http://www.jstor.org/stable/2319163>.
- Peter H. Friesen. The arrow-debreu model extended to financial markets. *Econometrica*, 47(3): 689–707, 1979. ISSN 00129682, 14680262. URL <http://www.jstor.org/stable/1910415>.

- Yuan Gao, Christian Kroer, and Alex Peysakhovich. Online market equilibrium with application to fair division. *Advances in Neural Information Processing Systems*, 34, 2021.
- Gauthier Gidel, Hugo Berard, Gaëtan Vignoud, Pascal Vincent, and Simon Lacoste-Julien. A variational inequality perspective on generative adversarial networks, 2020.
- Denizalp Goktas and Amy Greenwald. Convex-concave min-max stackelberg games. *Advances in Neural Information Processing Systems*, 34, 2021.
- Denizalp Goktas and Amy Greenwald. Robust no-regret learning in min-max Stackelberg games, 2022.
- Erfan Yazdandoost Hamedani and Necdet Serhat Aybat. A primal-dual algorithm for general convex-concave saddle point problems. *arXiv preprint arXiv:1803.01401*, 2, 2018.
- Xiuli He, Ashutosh Prasad, and Suresh P Sethi. Cooperative advertising and pricing in a dynamic stochastic supply chain: Feedback stackelberg strategies. In *PICMET'08-2008 Portland International Conference on Management of Engineering & Technology*, pp. 1634–1649. IEEE, 2008.
- Mingyi Hong, Junyu Zhang, , Shuzhong Zhang, et al. On lower iteration complexity bounds for the saddle point problems, 2020.
- Adam Ibrahim, Waiss Azizian, Gauthier Gidel, and Ioannis Mitliagkas. Lower bounds and conditioning of differentiable games. *arXiv preprint arXiv:1906.07300*, 2019.
- Anna Jaśkiewicz and Andrzej S Nowak. Zero-sum stochastic games. *Handbook of dynamic game theory*, pp. 1–64, 2018.
- Chi Jin, Praneeth Netrapalli, and Michael I. Jordan. What is local optimality in nonconvex-nonconcave minimax optimization?, 2020.
- Anatoli Juditsky, Arkadi Nemirovski, et al. First order methods for nonsmooth convex large-scale optimization, ii: utilizing problems structure. *Optimization for Machine Learning*, 30(9):149–183, 2011.
- HW Kuhn and AW Tucker. Proceedings of 2nd berkeley symposium, 1951.
- Tao Li and Suresh P Sethi. A review of dynamic stackelberg game models. *Discrete & Continuous Dynamical Systems-B*, 22(1):125, 2017.
- Tianyi Lin, Chi Jin, and Michael Jordan. On gradient descent ascent for nonconvex-concave minimax problems. In *International Conference on Machine Learning*, pp. 6083–6093. PMLR, 2020a.
- Tianyi Lin, Chi Jin, and Michael I Jordan. Near-optimal algorithms for minimax optimization. In *Conference on Learning Theory*, pp. 2738–2779. PMLR, 2020b.
- Songtao Lu, Ioannis Tsaknakis, and Mingyi Hong. Block alternating optimization for non-convex min-max problems: algorithms and applications in signal processing and communications. In *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 4754–4758. IEEE, 2019.
- R Mehra and EC Prescott. Recursive competitive equilibria and capital asset pricing. *Essays in Financial Economics*, 1977.
- Aryan Mokhtari, Asuman Ozdaglar, and Sarath Pattathil. Convergence rate of $\mathcal{O}(1/k)$ for optimistic gradient and extra-gradient methods in smooth convex-concave saddle point problems, 2020.
- Arkadi Nemirovski. Prox-method with rate of convergence $o(1/t)$ for variational inequalities with lipschitz continuous monotone operators and smooth convex-concave saddle point problems. *SIAM Journal on Optimization*, 15(1):229–251, 2004.
- Yurii Nesterov. Dual extrapolation and its applications to solving variational inequalities and related problems. *Mathematical Programming*, 109(2):319–344, 2007.

- Yurii Nesterov and Laura Scrimali. Solving strongly monotone variational and quasi-variational inequalities. *Discrete & Continuous Dynamical Systems*, 31(4):1383–1396, 2011.
- J v Neumann. Zur theorie der gesellschaftsspiele. *Mathematische annalen*, 100(1):295–320, 1928.
- Maher Nouiehed, Maziar Sanjabi, Tianjian Huang, Jason D Lee, and Meisam Razaviyayn. Solving a class of non-convex min-max games using iterative first order methods. *arXiv preprint arXiv:1902.08297*, 2019.
- Bernt Oksendal, Leif Sandal, and Jan Uboe. Stochastic stackelberg equilibria with applications to time-dependent newsvendor models. *Journal of Economic Dynamics and Control*, 37(7):1284–1299, 2013.
- Dmitrii M Ostrovskii, Andrew Lowy, and Meisam Razaviyayn. Efficient search of first-order nash equilibria in nonconvex-concave smooth min-max problems. *arXiv preprint arXiv:2002.07919*, 2020.
- Yuyuan Ouyang and Yangyang Xu. Lower complexity bounds of first-order methods for convex-concave bilinear saddle-point problems, 2018.
- Hassan Rafique, Mingrui Liu, Qihang Lin, and Tianbao Yang. Non-convex min-max optimization: Provable algorithms and applications in machine learning, 2019.
- Zhao Renbo. Optimal algorithms for stochastic three-composite convex-concave saddle point problems, 2019.
- Sean C Rismiller, Jonathan Cagan, and Christopher McComb. Stochastic stackelberg games for agent-driven robust design. In *International Design Engineering Technical Conferences and Computers and Information in Engineering Conference*, volume 84003, pp. V11AT11A039. American Society of Mechanical Engineers, 2020.
- Maziar Sanjabi, Jimmy Ba, Meisam Razaviyayn, and Jason D. Lee. On the convergence and robustness of training gans with regularized optimal transport. In *Proceedings of the 32nd International Conference on Neural Information Processing Systems, NIPS’18*, pp. 7091–7101, Red Hook, NY, USA, 2018. Curran Associates Inc.
- Lloyd S Shapley. Stochastic games. *Proceedings of the national academy of sciences*, 39(10):1095–1100, 1953.
- Maurice Sion et al. On general minimax theorems. *Pacific Journal of mathematics*, 8(1):171–176, 1958.
- Kiran Koshy Thekumparampil, Prateek Jain, Praneeth Netrapalli, and Sewoong Oh. Efficient algorithms for smooth minimax optimization, 2019.
- Paul Tseng. On linear convergence of iterative methods for the variational inequality problem. *Journal of Computational and Applied Mathematics*, 60(1):237–252, 1995. ISSN 0377-0427. doi: [https://doi.org/10.1016/0377-0427\(94\)00094-H](https://doi.org/10.1016/0377-0427(94)00094-H). URL <https://www.sciencedirect.com/science/article/pii/037704279400094H>. Proceedings of the International Meeting on Linear/Nonlinear Iterative Methods and Verification of Solution.
- Paul Tseng. On accelerated proximal gradient methods for convex-concave optimization. *submitted to SIAM Journal on Optimization*, 1, 2008.
- Deepanshu Vasal. Stochastic stackelberg games, 2020.
- Heinrich Von Stackelberg. *Marktform und gleichgewicht*. J. springer, 1934.
- Yevgeniy Vorobeychik and Satinder Singh. Computing stackelberg equilibria in discounted stochastic games. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 26, pp. 1478–1484, 2012.
- Yu Yuan, Zhibin Liang, and Xia Han. Robust reinsurance contract with asymmetric information in a stochastic stackelberg differential game. *Scandinavian Actuarial Journal*, pp. 1–28, 2021.
- Renbo Zhao. A primal dual smoothing framework for max-structured nonconvex optimization, 2020.

A PSEUDO-CODE FOR ALGORITHMS

Algorithm 1 Value Iteration (with Min-Max Oracle)**Inputs:** $\mathcal{S}, \mathcal{X}, \mathcal{Y}, r, g, p, \gamma, T$ **Outputs:** $v^{(T)}$

- 1: Initialize $v^{(0)}$ arbitrarily, e.g. $v^{(0)} = 0$
- 2: **for** $t = 1, \dots, T$ **do**
- 3: **for** $s \in \mathcal{S}$ **do**
- 4: $v^{(t+1)}(s) = \min_{\mathbf{x} \in \mathcal{X}} \max_{\mathbf{y} \in \mathcal{Y}: g(\mathbf{s}, \mathbf{x}, \mathbf{y}) \geq 0} \mathbb{E}_{S' \sim p(\cdot | \mathbf{s}, \mathbf{x}, \mathbf{y})} \left[r(\mathbf{s}, \mathbf{x}, \mathbf{y}) + \gamma v^{(t)}(S') \right]$
- 5: **end for**
- 6: **end for**
- 7: **return** $v^{(T)}$

Algorithm 2 Nested GDA on Action Value Function**Inputs:** $v, \mathbf{b}, \eta_{\mathbf{x}}, \eta_{\mathbf{y}}, T_p, T_X$ **Output:** $(\mathbf{p}^{(t)}, \mathbf{X}^{(t)})_{t=1}^T$

- 1: **for** $t = 1, \dots, T_p$ **do**
- 2:
- 3: **for** $s = 1, \dots, T_X$ **do**
- 4: For all $i \in [n]$,
- 5: $\mathbf{x}_i^{(s)} = \Pi_{\{(\mathbf{x}_i, \beta_i) \in \mathbb{R}_+^m \times \mathbb{R}_+ : \mathbf{x}_i \mathbf{p}^{(t-1)} + \beta_i \leq b_i\}} \left(\mathbf{x}_i^{(s-1)} + \eta_{\mathbf{X}} \left(\frac{b_i - \beta_i^{(s-1)}}{u_i^{(t)}(\mathbf{x}_i^{(t)})} \nabla_{\mathbf{x}_i} u_i^{(t)}(\mathbf{x}_i^{(t)}) \right) \right)$
- 6: $\beta_i^{(s)} = \Pi_{\{(\mathbf{x}_i, \beta_i) \in \mathbb{R}_+^m \times \mathbb{R}_+ : \mathbf{x}_i \mathbf{p}^{(t-1)} + \beta_i \leq b_i\}} \left(\beta_i^{(s-1)} + \eta_{\mathbf{X}} \left(-\log(u_i(\mathbf{x}_i^{(s-1)})) + \gamma \frac{\partial v(\mathbf{b})}{\partial b_i} \right) \right)$
- 7: **end for**
- 8: $\mathbf{p}^{(t)} = \Pi_{\mathbb{R}_+^m} \left(\mathbf{p}^{(t-1)} - \eta_{\mathbf{p}} (\mathbf{1} - \sum_{i \in [n]} \mathbf{x}_i) \right)$
- 9: **end for**
- 10: **return** $(\mathbf{p}^{(t)}, \mathbf{X}^{(t)})_{t=1}^T$

Algorithm 3 Max-Oracle GDA with Simulated Annealing on Action Value Function**Inputs:** $v, \mathbf{b}, \eta_{\mathbf{x}}, \eta_{\mathbf{y}}, T_p, T_X$ **Output:** $(\mathbf{p}^{(t)}, \mathbf{X}^{(t)})_{t=1}^T$

- 1: **for** $t = 1, \dots, T_p$ **do**
- 2: $q(\mathbf{p}, \mathbf{X}, \beta) = \sum_{j \in [m]} q_j p_j + \sum_{i \in [n]} (b_i - \beta_i) \log(u_i(\mathbf{x}_i, t_i)) + \gamma \mathbb{E}_{(t', \mathbf{b}', \mathbf{q}') \sim p(\cdot | t, \mathbf{b}, \mathbf{q}, \mathbf{X}, \beta)} [v(t', \mathbf{b}' + \beta, \mathbf{q}')]$
- 3: For all $i \in [n]$, $\mathbf{x}_i^{(t)}, \beta_i^{(t)} = \text{Simulated_Annealing}(q)$
- 4: $\mathbf{x}_i^{(t)}, \beta_i^{(t)} = \Pi_{\{(\mathbf{x}_i, \beta_i) \in \mathbb{R}_+^m \times \mathbb{R}_+ : \mathbf{x}_i \mathbf{p}^{(t-1)} + \beta_i \leq b_i\}} \left[\left(\mathbf{x}_i^{(t)}, \beta_i^{(t)} \right) \right]$
- 5: **end for**
- 6: $\mathbf{p}^{(t)} = \Pi_{\mathbb{R}_+^m} \left(\mathbf{p}^{(t-1)} + \eta_{\mathbf{p}} (\mathbf{1} - \sum_{i \in [n]} \mathbf{x}_i) \right)$
- 7: **return** $(\mathbf{p}^{(t)}, \mathbf{X}^{(t)})_{t=1}^T$

Algorithm 4 Value Iteration on Stochastic Fisher Market

```

1: Initialize  $v^{(0)}$  arbitrarily, e.g.  $v^{(0)} = 0$ 
2: for  $k = 1, \dots, T_v$  do
3:   for  $s \in \mathcal{S}$  do
4:      $v^{(k+1)}(s) = (1 - \gamma) \min_{\mathbf{p} \in \mathbb{R}_+^m} \max_{(\mathbf{X}, \boldsymbol{\beta}) \in \mathbb{R}_+^n \times \mathbb{R}_+^m : \mathbf{X}\mathbf{p} + \boldsymbol{\beta} \leq \mathbf{b}} \sum_{j \in [m]} q_j p_j +$ 
        $\sum_{i \in [n]} (b_i - \beta_i) \log(u_i(\mathbf{x}_i)) + \gamma v^*(\boldsymbol{\beta})$ 
5:   end for
6: end for

```

B OMITTED RESULTS AND PROOFS SECTION 2

We first note the following fundamental relationship between the state-value and action-value functions which is an analog of Bellman's Theorem (Bellman, 1952) and which follows from their definitions:

Theorem B.1. *Given a stochastic min-max Stackelberg game $(\mathcal{S}, \mathcal{X}, \mathcal{Y}, \mu^{(0)}, r, \mathbf{g}, p, \gamma)$, for all $v \in \mathcal{V}$, $q \in \mathcal{Q}$, $\pi_x \in \mathcal{X}^{\mathcal{S}}$, and $\pi_y \in \mathcal{Y}^{\mathcal{S}}$, $v = v^{\pi_x \pi_y}$ and $q = q^{\pi_x \pi_y}$ iff:*

$$v(s) = q(s, \pi_x(s), \pi_y(s)) \quad (7)$$

$$q(s, \mathbf{x}, \mathbf{y}) = \mathbb{E}_{S' \sim p(\cdot | s, \mathbf{x}, \mathbf{y})} [r(s, \mathbf{x}, \mathbf{y}) + \gamma v(S')] \quad (8)$$

Proof of Theorem B.1. By the definition of the state value function we have $v_i^{\pi_x \pi_y} = q_i(s, \pi_x(s), \pi_y(s))$, hence by Equation (7) we must have that $v_i = v_i^{\pi_x \pi_y}$. Additionally, by Equation (8) and the definition of the action-value functions this also implies that $q_i(s, \mathbf{x}, \mathbf{y}) = q_i^{\pi_x \pi_y}(s, \mathbf{x}, \mathbf{y})$ \square

Lemma B.2. *Suppose that there exists $(\pi_x^*, \pi_y^*) \in \mathcal{S}^{\mathcal{X}} \times \mathcal{S}^{\mathcal{Y}}$ such that for all states $s \in \mathcal{S}$:*

$$\max_{\pi_y \in \mathcal{S}^{\mathcal{Y}}: G(\pi_x^*, \pi_y) \geq 0} v^{\pi_x^* \pi_y}(s) \leq v^{\pi_x^* \pi_y^*}(s) \leq \min_{\pi_x \in \mathcal{S}^{\mathcal{X}}} \max_{\pi_y \in \mathcal{S}^{\mathcal{Y}}: G(\pi_x, \pi_y) \geq 0} v^{\pi_x \pi_y}(s) . \quad (9)$$

Then, (π_x^, π_y^*) is a recSE.*

Lemma B.2. *Suppose that there exists $(\pi_x^*, \pi_y^*) \in \mathcal{S}^{\mathcal{X}} \times \mathcal{S}^{\mathcal{Y}}$ such that for all states $s \in \mathcal{S}$:*

$$\max_{\pi_y \in \mathcal{S}^{\mathcal{Y}}: G(\pi_x^*, \pi_y) \geq 0} v^{\pi_x^* \pi_y}(s) \leq v^{\pi_x^* \pi_y^*}(s) \leq \min_{\pi_x \in \mathcal{S}^{\mathcal{X}}} \max_{\pi_y \in \mathcal{S}^{\mathcal{Y}}: G(\pi_x, \pi_y) \geq 0} v^{\pi_x \pi_y}(s) . \quad (9)$$

Then, (π_x^, π_y^*) is a recSE.*

Proof of Lemma B.2. We prove the right hand side inequality first:

$$v^{\pi_x^* \pi_y^*}(s) \leq \min_{\pi_x \in \mathcal{S}^{\mathcal{X}}} \max_{\pi_y \in \mathcal{S}^{\mathcal{Y}}: G(\pi_x, \pi_y) \geq 0} v^{\pi_x \pi_y}(s) \quad (10)$$

$$\sum_{s \in \mathcal{S}} \mu^{(0)}(s) v^{\pi_x^* \pi_y^*}(s) \leq \sum_{s \in \mathcal{S}} \mu^{(0)}(s) \min_{\pi_x \in \mathcal{S}^{\mathcal{X}}} \max_{\pi_y \in \mathcal{S}^{\mathcal{Y}}: G(\pi_x, \pi_y) \geq 0} v^{\pi_x \pi_y}(s) \quad (11)$$

$$u(\pi_x^*, \pi_y^*) \leq \sum_{s \in \mathcal{S}} \mu^{(0)}(s) \min_{\pi_x \in \mathcal{S}^{\mathcal{X}}} \max_{\pi_y \in \mathcal{S}^{\mathcal{Y}}: G(\pi_x, \pi_y) \geq 0} v^{\pi_x \pi_y}(s) \quad (12)$$

Since (π_x^*, π_y^*) is well-defined, it must mean that the same policy minimizes (resp. maximizes) the value of each state for the x player (y -player), hence the min-max operator can be pulled out of the sum, giving us:

$$u(\pi_x^*, \pi_y^*) \leq \min_{\pi_x \in \mathcal{S}^{\mathcal{X}}} \max_{\pi_y \in \mathcal{S}^{\mathcal{Y}}: G(\pi_x, \pi_y) \geq 0} \sum_{s \in \mathcal{S}} \mu^{(0)}(s) v^{\pi_x \pi_y}(s) \quad (13)$$

$$u(\pi_x^*, \pi_y^*) \leq \min_{\pi_x \in \mathcal{S}^{\mathcal{X}}} \max_{\pi_y \in \mathcal{S}^{\mathcal{Y}}: G(\pi_x, \pi_y) \geq 0} u(\pi_x, \pi_y) \quad (14)$$

Similarly, manipulating the the left hand side inequality , we obtain:

$$\min_{\pi_x \in \mathcal{S}^{\mathcal{X}}} \max_{\pi_y \in \mathcal{S}^{\mathcal{Y}}: G(\pi_x, \pi_y) \geq 0} u(\pi_x, \pi_y) \leq u(\pi_x^*, \pi_y^*) \quad (15)$$

□

We will show, there always exists a policy profile which satisfies the conditions of Lemma B.2. To show the existence of such a policy profile, we first re-state the condition given in Lemma B.2, in terms of action value functions.

Lemma 2.1. *Consider a zero-sum stochastic Stackelberg game $\mathcal{G}^{(0)}$. A policy profile $(\pi_x^*, \pi_y^*) \in \mathcal{S}^{\mathcal{X}} \times \mathcal{S}^{\mathcal{Y}}$ is a recSE if, for all $s \in \mathcal{S}$, we have that: $\max_{y \in \mathcal{Y}: g(s, \pi_x^*(s), y) \geq 0} q^{\pi_x^* \pi_y^*}(s, \pi_x^*(s), y) \leq q^{\pi_x^* \pi_y^*}(s, \pi_x^*(x), \pi_y^*(y)) \leq \min_{x \in \mathcal{X}} \max_{y \in \mathcal{Y}: g(s, x, y) \geq 0} q^{\pi_x^* \pi_y^*}(s, x, y)$. Equivalently, a policy profile (π_x^*, π_y^*) is a recSE if $(\pi_x^*(s), \pi_y^*(s))$ is a Stackelberg equilibrium: i.e., a solution to $\min_{x \in \mathcal{X}} \max_{y \in \mathcal{Y}: g(s, x, y) \geq 0} q^{\pi_x^* \pi_y^*}(s, x, y)$ at each $s \in \mathcal{S}$.*

Proof of Lemma 2.1.

$$\max_{y \in \mathcal{Y}: g(s, \pi_x^*(s), y) \geq 0} q^{\pi_x^* \pi_y^*}(s, \pi_x^*(s), y) \leq q^{\pi_x^* \pi_y^*}(s, \pi_x^*(s), \pi_y^*(s)) \leq \min_{x \in \mathcal{X}} \max_{y \in \mathcal{Y}: g(s, x, y) \geq 0} q^{\pi_x^* \pi_y^*}(s, x, y) \quad (16)$$

$$\max_{\pi_y \in \mathcal{S}^{\mathcal{Y}}: G(\pi_x^*, \pi_y) \geq 0} q^{\pi_x^* \pi_y^*}(s, \pi_x^*(s), \pi_y(s)) \leq q^{\pi_x^* \pi_y^*}(s, \pi_x^*(s), \pi_y^*(s)) \leq \min_{\pi_x \in \mathcal{S}^{\mathcal{X}}} \max_{\pi_y \in \mathcal{S}^{\mathcal{Y}}: G(\pi_x, \pi_y) \geq 0} q^{\pi_x^* \pi_y^*}(s, \pi_x(s), \pi_y(s)) \quad (17)$$

$$\max_{\pi_y \in \mathcal{S}^{\mathcal{Y}}: G(\pi_x^*, \pi_y) \geq 0} q^{\pi_x^* \pi_y^*}(s, \pi_x^*(s), \pi_y(s)) \leq q^{\pi_x^* \pi_y^*}(s, \pi_x^*(s), \pi_y^*(s)) \leq \min_{\pi_x \in \mathcal{S}^{\mathcal{X}}} \max_{\pi_y \in \mathcal{S}^{\mathcal{Y}}: G(\pi_x, \pi_y) \geq 0} q^{\pi_x \pi_y}(s, \pi_x(s), \pi_y(s)) \quad (18)$$

$$\max_{\pi_y \in \mathcal{S}^{\mathcal{Y}}: G(\pi_x^*, \pi_y) \geq 0} v^{\pi_x^* \pi_y^*}(s) \leq v^{\pi_x^* \pi_y^*}(s) \leq \min_{\pi_x \in \mathcal{S}^{\mathcal{X}}} \max_{\pi_y \in \mathcal{S}^{\mathcal{Y}}: G(\pi_x, \pi_y) \geq 0} v^{\pi_x \pi_y}(s) \quad (19)$$

By Lemma B.2, we must then have that (π_x^*, π_y^*) is a recursive Stackelberg equilibrium of $(\mathcal{S}, \mathcal{X}, \mathcal{Y}, \mu^{(0)}, r, g, p, \gamma)$.

□

Theorem 2.2. (π_x^*, π_y^*) is a recSE of $\mathcal{G}^{(0)}$ of $v^{\pi_x^* \pi_y^*}$ iff it induces a value function which is a fixed point of C : i.e., (π_x^*, π_y^*) is a Stackelberg equilibrium iff, for all $s \in \mathcal{S}$, $(Cv^{\pi_x^* \pi_y^*})(s) = v^{\pi_x^* \pi_y^*}(s)$.

Proof of Theorem 2.2. (recursive Stackelberg equilibrium \implies Fixed Point):

Suppose that (π_x^*, π_y^*) is a recursive Stackelberg equilibrium of $(\mathcal{S}, \mathcal{X}, \mathcal{Y}, \mu^{(0)}, r, g, p, \gamma)$, then:

$$(Cv^{\pi_x^* \pi_y^*})(s) = \min_{x \in \mathcal{X}} \max_{y \in \mathcal{Y}: g(s, x, y) \geq 0} \mathbb{E}_{S' \sim p(\cdot | s, x, y)} \left[r(s, x, y) + \gamma v^{\pi_x^* \pi_y^*}(S') \right] \quad (20)$$

$$= \min_{x \in \mathcal{X}} \max_{y \in \mathcal{Y}: g(s, x, y) \geq 0} q^{\pi_x^* \pi_y^*}(s, x, y) \quad (21)$$

$$= v^{\pi_x^* \pi_y^*}(s) \quad (\text{Lemma 2.1}) \quad (22)$$

(Fixed Point \implies recursive Stackelberg equilibrium) Suppose that a value function $v^{\pi_x^* \pi_y^*}$ which is induced by a policy profile (π_x^*, π_y^*) is a fixed point of C , we then have for all states $s \in \mathcal{S}$:

$$v^{\pi_x^* \pi_y^*}(s) = (Cv^{\pi_x^* \pi_y^*})(s) \quad (23)$$

$$= \min_{x \in \mathcal{X}} \max_{y \in \mathcal{Y}: g(s, x, y) \geq 0} \mathbb{E}_{S' \sim p(\cdot | s, x, y)} \left[r(s, x, y) + \gamma v^{\pi_x^* \pi_y^*}(S') \right] \quad (24)$$

$$= \min_{x \in \mathcal{X}} \max_{y \in \mathcal{Y}: g(s, x, y) \geq 0} q^{\pi_x^* \pi_y^*}(s, x, y) \quad (25)$$

Hence, by Lemma 2.1, (π_x^*, π_y^*) is recursive Stackelberg equilibrium. \square

Lemma 2.3. *Suppose that $f, h : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$, $g : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}^d$ are continuous functions, and \mathcal{X}, \mathcal{Y} are compact sets, we then have:*

$$\left| \min_{x \in \mathcal{X}} \max_{y \in \mathcal{Y}: g(x, y) \geq 0} f(x, y) - \min_{x \in \mathcal{X}} \max_{y \in \mathcal{Y}: g(x, y) \geq 0} h(x, y) \right| \leq \max_{(x, y) \in \mathcal{X} \times \mathcal{Y}} |f(x, y) - h(x, y)|.$$

Proof of Lemma 2.3. Let (x^*, y^*) be a Stackelberg equilibrium of $\min_{x \in \mathcal{X}} \max_{y \in \mathcal{Y}: g(x, y) \geq 0} f(x, y)$, and (x', y') be a Stackelberg equilibrium of $\min_{x \in \mathcal{X}} \max_{y \in \mathcal{Y}: g(x, y) \geq 0} h(x, y)$. Additionally, let $\bar{y} \in \arg \max_{y \in \mathcal{Y}: g(x', y) \geq 0} f(x', y)$, then by the definition of a Stackelberg equilibrium, we have $f(x^*, y^*) = \min_{y \in \mathcal{Y}} \max_{x \in \mathcal{X}: g(x, y) \geq 0} f(x, y) \leq \max_{y \in \mathcal{Y}: g(x', y) \geq 0} f(x', y) = f(x', \bar{y})$, and $h(x', y') = \max_{y \in \mathcal{Y}: g(x', y) \geq 0} h(x', y) \geq h(x', \bar{y})$.

Suppose that $\min_{x \in \mathcal{X}} \max_{y \in \mathcal{Y}: g(x, y) \geq 0} f(x, y) \geq \min_{x \in \mathcal{X}} \max_{y \in \mathcal{Y}: g(x, y) \geq 0} h(x, y)$ this gives us:

$$\left| \min_{x \in \mathcal{X}} \max_{y \in \mathcal{Y}: g(x, y) \geq 0} f(x, y) - \min_{x \in \mathcal{X}} \max_{y \in \mathcal{Y}: g(x, y) \geq 0} h(x, y) \right| \quad (26)$$

$$= |f(x^*, y^*) - h(x', y')| \quad (27)$$

$$\leq |f(x', \bar{y}) - h(x', y')| \quad (28)$$

$$\leq |f(x', \bar{y}) - h(x', \bar{y})| \quad (29)$$

$$\leq \max_{(x, y) \in \mathcal{X} \times \mathcal{Y}} |f(x, y) - h(x, y)| \quad (30)$$

The opposite case follows similarly by symmetry. \square

Theorem 2.4. *Consider the operator C associated with a stochastic Stackelberg game $\mathcal{G}^{(0)}$. If Assumption 1.1 holds, then C is a contraction mapping w.r.t. to the sup norm $\|\cdot\|_\infty$ with constant γ .*

Proof of Theorem 2.4. We will show that C is a contraction mapping, which then by Banach fixed point theorem establish the result. Let $v, v' \in \mathcal{V}$ be any two state value functions and $q, q' \in \mathcal{Q}$ be the respective associated action-value functions. We then have:

$$\|Cv - Cv'\|_\infty \quad (31)$$

$$\leq \max_{s \in \mathcal{S}} \left| \min_{x \in \mathcal{X}} \max_{y \in \mathcal{Y}: g(s, x, y) \geq 0} q(s, x, y) - \min_{x \in \mathcal{X}} \max_{y \in \mathcal{Y}: g(s, x, y) \geq 0} q'(s, x, y) \right| \quad (32)$$

$$\leq \max_{s \in \mathcal{S}} \max_{(x, y) \in \mathcal{X} \times \mathcal{Y}} |q(s, x, y) - q'(s, x, y)| \quad (\text{Lemma 2.3})$$

$$\quad (33)$$

$$\leq \max_{s \in \mathcal{S}} \max_{(x, y) \in \mathcal{X} \times \mathcal{Y}} \left| \mathbb{E}_{S' \sim p(\cdot | s, x, y)} [r(s, x, y) + \gamma v(S')] - \mathbb{E}_{S' \sim p(\cdot | s, x, y)} [r(s, x, y) + \gamma v'(S')] \right| \quad (34)$$

$$\leq \max_{s \in \mathcal{S}} \max_{(x, y) \in \mathcal{X} \times \mathcal{Y}} \left| \mathbb{E}_{S' \sim p(\cdot | s, x, y)} [\gamma v(S') - \gamma v'(S')] \right| \quad (35)$$

$$\leq \gamma \max_{s \in \mathcal{S}} \max_{(x, y) \in \mathcal{X} \times \mathcal{Y}} \left| \mathbb{E}_{S' \sim p(\cdot | s, x, y)} [v(S') - v'(S')] \right| \quad (36)$$

$$\leq \gamma \max_{s \in \mathcal{S}} \max_{(x, y) \in \mathcal{X} \times \mathcal{Y}} |v(s) - v'(s)| \quad (37)$$

$$= \gamma \|v - v'\|_\infty \quad (38)$$

Since $\gamma \in (0, 1)$, C is a contraction mapping. \square

Theorem 2.5. *Consider a zero-sum stochastic Stackelberg game $\mathcal{G}^{(0)}$. If Assumption 1.1 holds, then $\mathcal{G}^{(0)}$ has at least one recSE (π_x^*, π_y^*) with unique value function $v^{\pi_x^* \pi_y^*}$. Further, $v^{\pi_x^* \pi_y^*}$ can be computed by iteratively applying C to any initial state-value function $v^{(0)} \in \mathcal{V}$: $\lim_{t \rightarrow \infty} v^{(t)} = v^{\pi_x^* \pi_y^*}$.*

Proof of Theorem 2.5. By combining Theorem 2.4 and the Banach fixed point theorem Banach (1922), we obtain that a fixed point of C exists. Hence, by Theorem 2.2, a recursive Stackelberg equilibrium of $(\mathcal{S}, \mathcal{X}, \mathcal{Y}, \mu^{(0)}, r, \mathbf{g}, p, \gamma)$ exists and the value function induced by all recursive Stackelberg equilibria is the same, i.e., the optimal value function is unique. Additionally, by the second part of the Banach fixed point theorem, we must then also have $\lim_{t \rightarrow \infty} v^{(t)} \rightarrow v^{\pi_x^* \pi_y^*}$. \square

For any $q \in \mathcal{Q}$, we define a **greedy policy profile** with respect to q as a policy profile (π_x^q, π_y^q) such that $\pi_x^q \in \arg \min_{x \in \mathcal{X}} \max_{y \in \mathcal{Y}: g(s, x, y) \geq 0} q(s, x, y)$ and $\pi_y^q \in \arg \max_{y \in \mathcal{Y}: g(s, \pi_x^q(x), y) \geq 0} q(s, \pi_x^q(x), y)$. The following lemma provides a progress bound for each iteration of value iteration which is expressed in terms of the value function associated with the greedy policy profile.

Lemma B.3. *Let (π_x^*, π_y^*) be the recSE of a zero-sum Stochastic Stackelberg game $(\mathcal{S}, \mathcal{X}, \mathcal{Y}, \mu^{(0)}, r, \mathbf{g}, p, \gamma)$. For any $q \in \mathcal{Q}$, let π_x^q, π_y^q denote the greedy policy with respect to q of x -player and y -player respectively. Then, the following bound holds:*

$$v^{\pi_x^q \pi_y^q}(s) - v^{\pi_x^* \pi_y^*}(s) \leq \frac{2\|q - q^*\|_\infty}{(1 - \gamma)} \quad (39)$$

Using the above progress bound, we can then obtain a convergence rate for value iteration under Assumption 1.1.

Proof of Lemma B.3. For any state s , we have:

$$v^{\pi_x^* \pi_y^*}(s) - v^{\pi_x^q \pi_y^q}(s) \quad (40)$$

$$= q^{\pi_x^* \pi_y^*}(s, \pi_x^*(s), \pi_y^*(s)) - q^{\pi_x^q \pi_y^q}(s, \pi_x^q(s), \pi_y^q(s)) \quad (41)$$

$$= q^{\pi_x^* \pi_y^*}(s, \pi_x^*(s), \pi_y^*(s)) - q^{\pi_x^* \pi_y^*}(s, \pi_x^q(s), \pi_y^q(s)) + q^{\pi_x^* \pi_y^*}(s, \pi_x^q(s), \pi_y^q(s)) - q^{\pi_x^q \pi_y^q}(s, \pi_x^q(s), \pi_y^q(s)) \quad (42)$$

$$= q^{\pi_x^* \pi_y^*}(s, \pi_x^*(s), \pi_y^*(s)) - q^{\pi_x^* \pi_y^*}(s, \pi_x^q(s), \pi_y^q(s)) + \gamma \mathbb{E}_{s' \sim p(\cdot | s, \pi_x^q(s), \pi_y^q(s))} [v^{\pi_x^* \pi_y^*}(s') - v^{\pi_x^q \pi_y^q}(s')] \quad (43)$$

$$\geq q^{\pi_x^* \pi_y^*}(s, \pi_x^*(s), \pi_y^*(s)) + \left[q(s, \pi_x^q(s), \pi_y^q(s)) - \max_{y \in \mathcal{Y}: g(s, \pi_x^*(s), y) \geq 0} q(s, \pi_x^*(s), y) \right] - q^{\pi_x^* \pi_y^*}(s, \pi_x^q(s), \pi_y^q(s))$$

$$+ \gamma \mathbb{E}_{s' \sim p(s, \pi_x^q(s), \pi_y^q(s))} [v^{\pi_x^* \pi_y^*}(s') - v^{\pi_x^q \pi_y^q}(s')] \quad (44)$$

$$(45)$$

where the last line was obtained by the definition of a recursive Stackelberg equilibrium which ensures that $q(s, \pi_x^q(s), \pi_y^q(s)) \leq \max_{y \in \mathcal{Y}: g(s, \pi_x^*(s), y) \geq 0} q(s, \pi_x^*(s), y)$.

Re-organizing terms, we have:

$$v^{\pi_x^* \pi_y^*}(s) - v^{\pi_x^q \pi_y^q}(s) \quad (46)$$

$$\geq \left[q^{\pi_x^* \pi_y^*}(s, \pi_x^*(s), \pi_y^*(s)) - \max_{y \in \mathcal{Y}: g(s, \pi_x^*(s), y) \geq 0} q(s, \pi_x^*(s), y) \right] + \left[q(s, \pi_x^q(s), \pi_y^q(s)) - q^{\pi_x^* \pi_y^*}(s, \pi_x^q(s), \pi_y^q(s)) \right]$$

$$+ \gamma \mathbb{E}_{s' \sim p(s, \pi_x^q(s), \pi_y^q(s))} [v^{\pi_x^* \pi_y^*}(s') - v^{\pi_x^q \pi_y^q}(s')] \quad (47)$$

$$= \left[\max_{y \in \mathcal{Y}: g(s, \pi_x^*(s), y) \geq 0} q^{\pi_x^* \pi_y^*}(s, \pi_x^*(s), y) - \max_{y \in \mathcal{Y}: g(s, \pi_x^*(s), y) \geq 0} q(s, \pi_x^*(s), y) \right] + \left[q(s, \pi_x^q(s), \pi_y^q(s)) - q^{\pi_x^* \pi_y^*}(s, \pi_x^q(s), \pi_y^q(s)) \right] + \gamma \mathbb{E}_{s' \sim p(s, \pi_x^q(s), \pi_y^q(s))} [v^{\pi_x^* \pi_y^*}(s') - v^{\pi_x^q \pi_y^q}(s')] \quad (48)$$

Taking the minimum of both sides over states, and re-collecting terms, we obtain:

$$\min_{\mathbf{s} \in \mathcal{S}} \left(v^{\pi_{\mathbf{x}}^* \pi_{\mathbf{y}}^*}(\mathbf{s}) - v^{\pi_{\mathbf{x}}^q \pi_{\mathbf{y}}^q}(\mathbf{s}) \right) \geq -2\|q - q^*\|_{\infty} + \gamma \min_{\mathbf{s} \in \mathcal{S}} \left[v^{\pi_{\mathbf{x}}^* \pi_{\mathbf{y}}^*}(\mathbf{s}) - v^{\pi_{\mathbf{x}}^q \pi_{\mathbf{y}}^q}(\mathbf{s}) \right] \quad (49)$$

$$(1 - \gamma) \min_{\mathbf{s} \in \mathcal{S}} \left(v^{\pi_{\mathbf{x}}^* \pi_{\mathbf{y}}^*}(\mathbf{s}) - v^{\pi_{\mathbf{x}}^q \pi_{\mathbf{y}}^q}(\mathbf{s}) \right) \geq -2\|q - q^*\|_{\infty} \quad (50)$$

$$\min_{\mathbf{s} \in \mathcal{S}} \left(v^{\pi_{\mathbf{x}}^* \pi_{\mathbf{y}}^*}(\mathbf{s}) - v^{\pi_{\mathbf{x}}^q \pi_{\mathbf{y}}^q}(\mathbf{s}) \right) \geq -\frac{2\|q - q^*\|_{\infty}}{(1 - \gamma)} \quad (51)$$

$$v^{\pi_{\mathbf{x}}^* \pi_{\mathbf{y}}^*}(\mathbf{s}) - v^{\pi_{\mathbf{x}}^q \pi_{\mathbf{y}}^q}(\mathbf{s}) \geq -\frac{2\|q - q^*\|_{\infty}}{(1 - \gamma)} \quad \forall \mathbf{s} \in \mathcal{S} \quad (52)$$

□

Theorem 2.8. [Convergence of Value Iteration] Suppose value iteration is run on input \mathcal{G} . If Assumption 1.1 holds, and if we initialize $v^{(0)}(\mathbf{s}) = 0$, for all $\mathbf{s} \in \mathcal{S}$, then for $k \geq \frac{1}{1-\gamma} \log \frac{2\alpha}{\epsilon(1-\gamma)}$, we have $v^{(k)}(\mathbf{s}) - v^{\pi_{\mathbf{x}}^* \pi_{\mathbf{y}}^*}(\mathbf{s}) \leq \epsilon$.

Proof of Theorem 2.8. First note that by Assumption 1.1, we have that $\|v^{\pi_{\mathbf{x}}^* \pi_{\mathbf{y}}^*}\|_{\infty} \leq \alpha$. Applying the operator C repeatedly and using the fact that $v^{\pi_{\mathbf{x}}^* \pi_{\mathbf{y}}^*} = C v^{\pi_{\mathbf{x}}^* \pi_{\mathbf{y}}^*}$ from Theorem 2.2, we obtain

$$\begin{aligned} & \|v^{(k)} - v^{\pi_{\mathbf{x}}^* \pi_{\mathbf{y}}^*}\|_{\infty} \\ &= \|(C)^k v^{(0)} - (C)^k v^*\|_{\infty} \\ &\leq \gamma^k \|v^{(0)} - v^{\pi_{\mathbf{x}}^* \pi_{\mathbf{y}}^*}\|_{\infty} \\ &= \|(C)^k q^{(0)} - (C)^k q^*\|_{\infty} \\ &\leq \gamma^k \|v^{(0)} - v^{\pi_{\mathbf{x}}^* \pi_{\mathbf{y}}^*}\|_{\infty} \\ &= \gamma^k \max_{\mathbf{s} \in \mathcal{S}} \left| \min_{\mathbf{x} \in \mathcal{X}} \max_{\mathbf{y} \in \mathcal{Y}: g(\mathbf{s}, \mathbf{x}, \mathbf{y}) \geq 0} q^{(0)}(\mathbf{s}, \mathbf{x}, \mathbf{y}) - \min_{\mathbf{x} \in \mathcal{X}} \max_{\mathbf{y} \in \mathcal{Y}: g(\mathbf{s}, \mathbf{x}, \mathbf{y}) \geq 0} q^{\pi_{\mathbf{x}}^* \pi_{\mathbf{y}}^*}(\mathbf{s}, \mathbf{x}, \mathbf{y}) \right| \\ &= \gamma^k \max_{\mathbf{s} \in \mathcal{S}} \max_{(\mathbf{x}, \mathbf{y}) \in \mathcal{X} \times \mathcal{Y}} \left| q^{(0)}(\mathbf{s}, \mathbf{x}, \mathbf{y}) - q^{\pi_{\mathbf{x}}^* \pi_{\mathbf{y}}^*}(\mathbf{s}, \mathbf{x}, \mathbf{y}) \right| \quad (\text{Lemma 2.3}) \\ &= \gamma^k \|q^*\|_{\infty} \\ &\leq \gamma^k \alpha \end{aligned}$$

Since $1 - x \leq e^{-x}$ for any $x \in \mathbb{R}$, we have

$$\gamma^k = (1 - (1 - \gamma))^k \leq (e^{-(1-\gamma)})^k \leq e^{-(1-\gamma)k}$$

Combining this bound with Equation (39) for all $q \in \mathcal{Q}$:

$$\begin{aligned} v^{(k)}(\mathbf{s}) - v^{\pi_{\mathbf{x}}^* \pi_{\mathbf{y}}^*}(\mathbf{s}) &\leq \frac{2\|q - q^*\|_{\infty}}{1 - \gamma} \\ &\leq \frac{2\gamma^k \alpha}{(1 - \gamma)} \\ &\leq \frac{2\alpha}{(1 - \gamma)} e^{-(1-\gamma)k} \end{aligned}$$

Thus it suffices to solve for k such that

$$\frac{2\alpha}{(1 - \gamma)} e^{-(1-\gamma)k} \leq \epsilon .$$

which concludes the proof. □

C OMITTED RESULTS AND PROOFS SECTION 3

Our characterization of the subdifferential of the value function associated with a Stackelberg equilibrium w.r.t. its parameters relies on a slightly generalized version of the subdifferential envelope theorem (Theorem C.1, Appendix C) of Goktas & Greenwald (2021), which characterizes the set of subdifferentials of parametrized constrained optimization problems, i.e., the set of subgradients w.r.t. \mathbf{x} of $f^*(\mathbf{x}) = \max_{\mathbf{y} \in \mathcal{Y}: h(\mathbf{x}, \mathbf{y}) \geq 0} f(\mathbf{x}, \mathbf{y})$. In particular, we note that Goktas & Greenwald’s proof goes through even without assuming the concavity of $f(\mathbf{x}, \mathbf{y}), h_1(\mathbf{x}, \mathbf{y}), \dots, h_d(\mathbf{x}, \mathbf{y})$ in \mathbf{y} , for all $\mathbf{x} \in \mathcal{X}$.

Theorem C.1 (Subdifferential Envelope Theorem). *Consider the function $f^*(\mathbf{x}) = \max_{\mathbf{y} \in \mathcal{Y}: h(\mathbf{x}, \mathbf{y}) \geq 0} f(\mathbf{x}, \mathbf{y})$ where $f : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$, and $\mathbf{h} : \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}^d$. Let $\mathcal{Y}^*(\hat{\mathbf{x}}) = \arg \max_{\mathbf{y} \in \mathcal{Y}: h(\hat{\mathbf{x}}, \mathbf{y}) \geq 0} f(\hat{\mathbf{x}}, \mathbf{y})$. Suppose that 1. $f(\mathbf{x}, \mathbf{y}), h_1(\mathbf{x}, \mathbf{y}), \dots, h_d(\mathbf{x}, \mathbf{y})$ are continuous in (\mathbf{x}, \mathbf{y}) and convex in \mathbf{x} ; 2. $\nabla_{\mathbf{x}} f, \nabla_{\mathbf{x}} h_1, \dots, \nabla_{\mathbf{x}} h_d$ exist for all $\mathbf{x} \in \mathcal{X}, \mathbf{y} \in \mathcal{Y}$; 3. \mathcal{Y} is non-empty and compact, and 4. (Slater’s condition) $\forall \mathbf{x} \in \mathcal{X}, \exists \hat{\mathbf{y}} \in \mathcal{Y}$ s.t. $g_k(\mathbf{x}, \hat{\mathbf{y}}) > 0$, for all $k = 1, \dots, d$. Then, f^* is subdifferentiable and at any point $\hat{\mathbf{x}} \in \mathcal{X}, \partial_{\mathbf{x}} f^*(\hat{\mathbf{x}}) =$*

$$\text{conv} \left(\bigcup_{\mathbf{y}^*(\hat{\mathbf{x}}) \in \mathcal{Y}^*(\hat{\mathbf{x}})} \bigcup_{\lambda_k^*(\hat{\mathbf{x}}, \mathbf{y}^*(\hat{\mathbf{x}})) \in \Lambda^*(\hat{\mathbf{x}}, \mathbf{y}^*(\hat{\mathbf{x}}))} \left\{ \nabla_{\mathbf{x}} f(\hat{\mathbf{x}}, \mathbf{y}^*(\hat{\mathbf{x}})) + \sum_{k=1}^d \lambda_k^*(\hat{\mathbf{x}}, \mathbf{y}^*(\hat{\mathbf{x}})) \nabla_{\mathbf{x}} g_k(\hat{\mathbf{x}}, \mathbf{y}^*(\hat{\mathbf{x}})) \right\} \right), \quad (53)$$

where conv is the convex hull operator and $\lambda^*(\hat{\mathbf{x}}, \mathbf{y}^*(\hat{\mathbf{x}})) = (\lambda_1^*(\hat{\mathbf{x}}, \mathbf{y}^*(\hat{\mathbf{x}})), \dots, \lambda_d^*(\hat{\mathbf{x}}, \mathbf{y}^*(\hat{\mathbf{x}})))^T \in \Lambda^*(\hat{\mathbf{x}}, \mathbf{y}^*(\hat{\mathbf{x}}))$ are the optimal KKT multipliers associated with $\mathbf{y}^*(\hat{\mathbf{x}}) \in \mathcal{Y}^*(\hat{\mathbf{x}})$.

We note the following known lemma which provides the necessary conditions for the Bellman equation associated with a recursive optimization problem to be continuous and convex in its parameters.⁶

Lemma C.2. *Consider the Bellman equation associated with a parametric recursive stochastic optimization problem, where $r : \mathcal{S} \times \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$, with state space \mathcal{S} and parameter set \mathcal{X} , and $\gamma \in (0, 1)$:*

$$v(\mathbf{s}, \mathbf{x}) = \max_{\mathbf{y} \in \mathcal{Y}: g(\mathbf{s}, \mathbf{x}, \mathbf{y}) \geq 0} \left\{ r(\mathbf{s}, \mathbf{x}, \mathbf{y}) + \gamma \mathbb{E}_{\mathbf{S}' \sim p(\cdot | \mathbf{s}, \mathbf{x}, \mathbf{y})} [v(\mathbf{S}', \mathbf{x})] \right\} \quad (54)$$

Suppose that Assumption 1.1 holds, and that 1. for all $\mathbf{s} \in \mathcal{S}, \mathbf{y} \in \mathcal{Y}$, $r(\mathbf{s}, \mathbf{x}, \mathbf{y}), g_1(\mathbf{s}, \mathbf{x}, \mathbf{y}), \dots, g_d(\mathbf{s}, \mathbf{x}, \mathbf{y})$ are concave in \mathbf{x} , 2. $p(\mathbf{s}' | \mathbf{s}, \mathbf{x}, \mathbf{y})$ is continuous (Feller property) and CSD convex in \mathbf{x} . Then, v is continuous and convex in \mathbf{x} .

Theorem C.3. [Subdifferential Benveniste-Scheinkman Theorem] *Consider the Bellman equation associated with a recursive stochastic optimization problem where $r : \mathcal{S} \times \mathcal{X} \times \mathcal{Y} \rightarrow \mathbb{R}$, with state space \mathcal{S} and parameter set \mathcal{X} , and $\gamma \in (0, 1)$:*

$$v(\mathbf{s}, \mathbf{x}) = \max_{\mathbf{y} \in \mathcal{Y}: g(\mathbf{s}, \mathbf{x}, \mathbf{y}) \geq 0} \left\{ r(\mathbf{s}, \mathbf{x}, \mathbf{y}) + \gamma \mathbb{E}_{\mathbf{S}' \sim p(\cdot | \mathbf{s}, \mathbf{x}, \mathbf{y})} [v(\mathbf{S}', \mathbf{x})] \right\} \quad (55)$$

Suppose that Assumption 1.1 holds, and that 1. for all $\mathbf{s} \in \mathcal{S}, \mathbf{y} \in \mathcal{Y}$, $r(\mathbf{s}, \mathbf{x}, \mathbf{y}), g_1(\mathbf{s}, \mathbf{x}, \mathbf{y}), \dots, g_d(\mathbf{s}, \mathbf{x}, \mathbf{y})$ are concave in \mathbf{x} , 2. $\nabla_{\mathbf{x}} r(\mathbf{s}, \mathbf{x}, \mathbf{y}), \nabla_{\mathbf{x}} g_1(\mathbf{s}, \mathbf{x}, \mathbf{y}), \dots, \nabla_{\mathbf{x}} g_d(\mathbf{s}, \mathbf{x}, \mathbf{y}), \nabla_{\mathbf{x}} p(\mathbf{s}' | \mathbf{s}, \mathbf{x}, \mathbf{y})$ exist for all $\mathbf{s}, \mathbf{s}' \in \mathcal{S}, \mathbf{x} \in \mathcal{X}, \mathbf{y} \in \mathcal{Y}$, 3. $p(\mathbf{s}' | \mathbf{s}, \mathbf{x}, \mathbf{y})$ is continuous CSD convex, and differentiable in \mathbf{x} , and 4. Slater’s condition holds for the optimization problem, i.e., $\forall \mathbf{x} \in \mathcal{X}, \mathbf{s} \in \mathcal{S}, \exists \hat{\mathbf{y}} \in \mathcal{Y}$ s.t. $g_k(\mathbf{s}, \mathbf{x}, \hat{\mathbf{y}}) > 0$, for all $k = 1, \dots, d$.

Let $\mathcal{Y}^*(\mathbf{s}, \mathbf{x}) = \max_{\mathbf{y} \in \mathcal{Y}: g(\mathbf{s}, \mathbf{x}, \mathbf{y}) \geq 0} \{ r(\mathbf{s}, \mathbf{x}, \mathbf{y}) + \gamma \mathbb{E}_{\mathbf{S}' \sim p(\cdot | \mathbf{s}, \mathbf{x}, \mathbf{y})} [v(\mathbf{S}', \mathbf{x})] \}$, then v is subdifferentiable and $\partial_{\mathbf{x}} v(\mathbf{s}, \hat{\mathbf{x}}) =$

$$\text{conv} \left(\bigcup_{\mathbf{y}^*(\mathbf{s}, \hat{\mathbf{x}}) \in \mathcal{Y}^*(\mathbf{s}, \hat{\mathbf{x}})} \bigcup_{\lambda_k^*(\mathbf{s}, \hat{\mathbf{x}}, \mathbf{y}^*(\mathbf{s}, \hat{\mathbf{x}})) \in \Lambda^*(\mathbf{s}, \hat{\mathbf{x}}, \mathbf{y}^*(\mathbf{s}, \hat{\mathbf{x}}))} \left\{ \nabla_{\mathbf{x}} r(\mathbf{s}, \hat{\mathbf{x}}, \mathbf{y}^*(\mathbf{s}, \hat{\mathbf{x}})) + \gamma \nabla_{\mathbf{x}} \mathbb{E}_{\mathbf{S}' \sim p(\cdot | \mathbf{s}, \hat{\mathbf{x}}, \mathbf{y}^*(\mathbf{s}, \hat{\mathbf{x}}))} [v(\mathbf{S}', \hat{\mathbf{x}})] + \sum_{k=1}^d \lambda_k^*(\mathbf{s}, \hat{\mathbf{x}}, \mathbf{y}^*(\mathbf{s}, \hat{\mathbf{x}})) \nabla_{\mathbf{x}} g_k(\mathbf{s}, \hat{\mathbf{x}}, \mathbf{y}^*(\mathbf{s}, \hat{\mathbf{x}})) \right\} \right). \quad (56)$$

⁶The interested reader can infer a proof from the proof of Theorem 2 of Atakan (2003)

Suppose additionally, that for all $s, s' \in \mathcal{S}, \mathbf{x} \in \mathcal{X}, \mathbf{y}^*(s, \mathbf{x}) \in \mathcal{Y}^*(s, \mathbf{x}) \nabla_{\mathbf{x}} p(s' | s, \mathbf{x}, \mathbf{y}^*(s, \mathbf{x})) > 0$, then $\partial_{\mathbf{x}} v(s, \hat{\mathbf{x}}) =$

$$\text{conv} \left(\bigcup_{\mathbf{y}^*(s, \hat{\mathbf{x}}) \in \mathcal{Y}^*(s, \hat{\mathbf{x}})} \bigcup_{\lambda_k^*(s, \hat{\mathbf{x}}, \mathbf{y}^*(s, \hat{\mathbf{x}})) \in \Lambda^*(s, \hat{\mathbf{x}}, \mathbf{y}^*(s, \hat{\mathbf{x}}))} \left\{ \nabla_{\mathbf{x}} r(s, \hat{\mathbf{x}}, \mathbf{y}^*(s, \hat{\mathbf{x}})) + \gamma \mathbb{E}_{S' \sim p(\cdot | s, \hat{\mathbf{x}}, \mathbf{y}^*(s, \hat{\mathbf{x}}))} [\nabla_{\mathbf{x}} v(S', \hat{\mathbf{x}})] \right. \right. \\ \left. \left. + \gamma \mathbb{E}_{S' \sim p(\cdot | s, \hat{\mathbf{x}}, \mathbf{y}^*(s, \hat{\mathbf{x}}))} [v(S', \hat{\mathbf{x}}) \nabla_{\mathbf{x}} \log(p(S' | s, \hat{\mathbf{x}}, \mathbf{y}^*(s, \hat{\mathbf{x}})))] + \sum_{k=1}^d \lambda_k^*(s, \hat{\mathbf{x}}, \mathbf{y}^*(s, \hat{\mathbf{x}})) \nabla_{\mathbf{x}} g_k(s, \hat{\mathbf{x}}, \mathbf{y}^*(s, \hat{\mathbf{x}})) \right\} \right). \quad (57)$$

where conv is the convex hull operator and $\lambda^*(s, \hat{\mathbf{x}}, \mathbf{y}^*(s, \hat{\mathbf{x}})) = (\lambda_1^*(s, \hat{\mathbf{x}}, \mathbf{y}^*(s, \hat{\mathbf{x}})), \dots, \lambda_d^*(s, \hat{\mathbf{x}}, \mathbf{y}^*(s, \hat{\mathbf{x}})))^T \in \Lambda^*(s, \hat{\mathbf{x}}, \mathbf{y}^*(s, \hat{\mathbf{x}}))$ are the optimal KKT multipliers associated with $\mathbf{y}^*(s, \hat{\mathbf{x}}) \in \mathcal{Y}^*(s, \hat{\mathbf{x}})$.

Proof of Theorem C.3. First, note that by Lemma C.2, v is continuous and convex in \mathbf{x} , hence it is subdifferentiable. From Theorem C.1, we then obtain the first part of the theorem:

$$\partial_{\mathbf{x}} v(s, \hat{\mathbf{x}}) = \text{conv} \left(\bigcup_{\mathbf{y}^*(s, \hat{\mathbf{x}}) \in \mathcal{Y}^*(s, \hat{\mathbf{x}})} \bigcup_{\lambda_k^*(s, \hat{\mathbf{x}}, \mathbf{y}^*(s, \hat{\mathbf{x}})) \in \Lambda^*(s, \hat{\mathbf{x}}, \mathbf{y}^*(s, \hat{\mathbf{x}}))} \left\{ \nabla_{\mathbf{x}} r(s, \hat{\mathbf{x}}, \mathbf{y}^*(s, \hat{\mathbf{x}})) + \gamma \nabla_{\mathbf{x}} \mathbb{E}_{S' \sim p(\cdot | s, \hat{\mathbf{x}}, \mathbf{y}^*(s, \hat{\mathbf{x}}))} [v(S', \mathbf{x})] \right. \right. \\ \left. \left. + \sum_{k=1}^d \lambda_k^*(s, \hat{\mathbf{x}}, \mathbf{y}^*(s, \hat{\mathbf{x}})) \nabla_{\mathbf{x}} g_k(s, \hat{\mathbf{x}}, \mathbf{y}^*(s, \hat{\mathbf{x}})) \right\} \right). \quad (58)$$

By the Leibniz integral rule Flanders (1973), the gradient of the expectation can instead be expressed as an expectation of the gradient under continuity of the function whose expectation is taken, in this case v . In particular, if for all $s, s' \in \mathcal{S}, \mathbf{x} \in \mathcal{X}, \mathbf{y}^*(s, \mathbf{x}) \in \mathcal{Y}^*(s, \mathbf{x}) \nabla_{\mathbf{x}} p(s' | s, \mathbf{x}, \mathbf{y}^*(s, \mathbf{x})) > 0$ we have:

$$\nabla_{\mathbf{x}} \mathbb{E}_{S' \sim p(\cdot | s, \mathbf{x}, \mathbf{y})} [v(S', \mathbf{x})] \quad (59)$$

$$= \nabla_{\mathbf{x}} \int_{z \in \mathcal{S}} p(z | s, \mathbf{x}, \mathbf{y}) v(z, \mathbf{x}) dz \quad (60)$$

$$= \int_{z \in \mathcal{S}} \nabla_{\mathbf{x}} [p(z | s, \mathbf{x}, \mathbf{y}) v(z, \mathbf{x})] dz \quad (\text{Leibniz Integral Rule}) \quad (61)$$

$$= \int_{z \in \mathcal{S}} [p(z | s, \mathbf{x}, \mathbf{y}) \nabla_{\mathbf{x}} v(z, \mathbf{x}) + v(z, \mathbf{x}) \nabla_{\mathbf{x}} p(z | s, \mathbf{x}, \mathbf{y})] dz \quad (\text{Product Rule}) \quad (62)$$

$$= \int_{z \in \mathcal{S}} \left[p(z | s, \mathbf{x}, \mathbf{y}) \nabla_{\mathbf{x}} v(z, \mathbf{x}) + v(z, \mathbf{x}) p(z | s, \mathbf{x}, \mathbf{y}) \frac{\nabla_{\mathbf{x}} p(z | s, \mathbf{x}, \mathbf{y})}{p(z | s, \mathbf{x}, \mathbf{y})} \right] dz \quad (p(z | s, \mathbf{x}, \mathbf{y}) > 0) \quad (63)$$

$$= \int_{z \in \mathcal{S}} [p(z | s, \mathbf{x}, \mathbf{y}) \nabla_{\mathbf{x}} v(z, \mathbf{x}) + v(z, \mathbf{x}) p(z | s, \mathbf{x}, \mathbf{y}) \nabla_{\mathbf{x}} \log p(z | s, \mathbf{x}, \mathbf{y})] dz \quad (64)$$

$$= \int_{z \in \mathcal{S}} [p(z | s, \mathbf{x}, \mathbf{y}) \nabla_{\mathbf{x}} v(z, \mathbf{x})] dz + \int_{z \in \mathcal{S}} [v(z, \mathbf{x}) p(z | s, \mathbf{x}, \mathbf{y}) \nabla_{\mathbf{x}} \log p(z | s, \mathbf{x}, \mathbf{y})] dz \quad (65)$$

$$= \mathbb{E}_{S' \sim p(\cdot | s, \mathbf{x}, \mathbf{y})} [\nabla_{\mathbf{x}} v(S', \mathbf{x})] + \mathbb{E}_{S' \sim p(\cdot | s, \mathbf{x}, \mathbf{y})} [v(S', \mathbf{x}) \nabla_{\mathbf{x}} \log p(S' | s, \mathbf{x}, \mathbf{y})] \quad (66)$$

This gives us $\partial_{\mathbf{x}}v(\mathbf{s}, \hat{\mathbf{x}}) =$

$$\text{conv} \left(\bigcup_{\mathbf{y}^*(\mathbf{s}, \hat{\mathbf{x}}) \in \mathcal{Y}^*(\mathbf{s}, \hat{\mathbf{x}})} \bigcup_{\lambda_k^*(\mathbf{s}, \hat{\mathbf{x}}, \mathbf{y}^*(\mathbf{s}, \hat{\mathbf{x}})) \in \Lambda^*(\mathbf{s}, \hat{\mathbf{x}}, \mathbf{y}^*(\hat{\mathbf{x}}))} \left\{ \nabla_{\mathbf{x}} r(\mathbf{s}, \hat{\mathbf{x}}, \mathbf{y}^*(\mathbf{s}, \hat{\mathbf{x}})) + \gamma \mathbb{E}_{\mathbf{S}' \sim p(\cdot | \mathbf{s}, \hat{\mathbf{x}}, \mathbf{y}^*(\mathbf{s}, \hat{\mathbf{x}}))} [\nabla_{\mathbf{x}} v(\mathbf{S}', \hat{\mathbf{x}})] \right. \right. \\ \left. \left. + \gamma \mathbb{E}_{\mathbf{S}' \sim p(\cdot | \mathbf{s}, \hat{\mathbf{x}}, \mathbf{y}^*(\mathbf{s}, \hat{\mathbf{x}}))} [v(\mathbf{S}', \hat{\mathbf{x}}) \nabla_{\mathbf{x}} \log(p(\mathbf{S}' | \mathbf{s}, \hat{\mathbf{x}}, \mathbf{y}^*(\mathbf{s}, \hat{\mathbf{x}})))] + \sum_{k=1}^d \lambda_k^*(\mathbf{s}, \hat{\mathbf{x}}, \mathbf{y}^*(\mathbf{s}, \hat{\mathbf{x}})) \nabla_{\mathbf{x}} g_k(\mathbf{s}, \hat{\mathbf{x}}, \mathbf{y}^*(\mathbf{s}, \hat{\mathbf{x}})) \right\} \right). \quad (67)$$

□

Note that in the special case that the probability transition function is representing a deterministic recursive parametrized optimization problem, $v(\mathbf{s}, \mathbf{x}) = \max_{\mathbf{y} \in \mathcal{Y}: g(\mathbf{s}, \mathbf{x}, \mathbf{y}) \geq 0} \{r(\mathbf{s}, \mathbf{x}, \mathbf{y}) + \gamma [v(\tau(\mathbf{s}, \mathbf{x}, \mathbf{y}), \mathbf{x})]\}$ i.e., $p(\mathbf{s}' | \mathbf{s}, \mathbf{x}, \mathbf{y}) \in \{0, 1\}$ for all $\mathbf{s}, \mathbf{s}' \in \mathcal{S}, \mathbf{x} \in \mathcal{X}, \mathbf{y} \in \mathcal{Y}$, and $\tau: \mathcal{S} \times \mathcal{X} \times \mathcal{Y} \rightarrow \mathcal{S}$ is such that $\tau(\mathbf{s}, \mathbf{x}, \mathbf{y}) = \mathbf{s}'$ iff $p(\mathbf{s}' | \mathbf{s}, \mathbf{x}, \mathbf{y}) = 1$, the CSD convexity assumption reduces to the linearity of the deterministic state transition function τ (Proposition 1 of Atakan (2003)). In this case, the subdifferential of the Bellman equation reduces to

$$\partial_{\mathbf{x}}v(\mathbf{s}, \hat{\mathbf{x}}) = \text{conv} \left(\bigcup_{\mathbf{y}^*(\mathbf{s}, \hat{\mathbf{x}}) \in \mathcal{Y}^*(\mathbf{s}, \hat{\mathbf{x}})} \bigcup_{\lambda_k^*(\mathbf{s}, \hat{\mathbf{x}}, \mathbf{y}^*(\mathbf{s}, \hat{\mathbf{x}})) \in \Lambda^*(\mathbf{s}, \hat{\mathbf{x}}, \mathbf{y}^*(\mathbf{s}, \hat{\mathbf{x}}))} \left\{ \nabla_{\mathbf{x}} r(\mathbf{s}, \hat{\mathbf{x}}, \mathbf{y}^*(\mathbf{s}, \hat{\mathbf{x}})) + \gamma \nabla_{\mathbf{x}} \tau(\mathbf{s}, \hat{\mathbf{x}}, \mathbf{y}) \nabla_{\mathbf{s}} v(\tau(\mathbf{s}, \hat{\mathbf{x}}, \mathbf{y}), \hat{\mathbf{x}}) \right. \right. \\ \left. \left. + \nabla_{\mathbf{x}} v(\tau(\mathbf{s}, \hat{\mathbf{x}}), \hat{\mathbf{x}}) + \sum_{k=1}^d \lambda_k^*(\mathbf{s}, \hat{\mathbf{x}}, \mathbf{y}^*(\mathbf{s}, \hat{\mathbf{x}})) \nabla_{\mathbf{x}} g_k(\mathbf{s}, \hat{\mathbf{x}}, \mathbf{y}^*(\mathbf{s}, \hat{\mathbf{x}})) \right\} \right) \quad (68)$$

(69)

Theorem 3.1. Consider a zero-sum stochastic Stackelberg game \mathcal{G} , where $\mathcal{X} = \{\mathbf{x} \in \mathbb{R}^n \mid q_1(\mathbf{x}) \leq 0, \dots, q_p(\mathbf{x}) \leq 0\}$ and $\mathcal{Y} = \{\mathbf{y} \in \mathbb{R}^m \mid r_1(\mathbf{y}) \geq 0, \dots, r_l(\mathbf{y}) \geq 0\}$. Let $\mathcal{L}_{\mathbf{s}, \mathbf{x}}(\mathbf{y}, \boldsymbol{\lambda}) = r(\mathbf{s}, \mathbf{x}, \mathbf{y}) + \gamma \mathbb{E}_{\mathbf{S}' \sim p(\cdot | \mathbf{s}, \mathbf{x}, \mathbf{y})} [v(\mathbf{S}', \mathbf{x})] + \sum_{k=1}^d \lambda_k g_k(\mathbf{s}, \mathbf{x}, \mathbf{y})$. Suppose that Assumption 1.1 holds, and that 1. for all $\mathbf{s} \in \mathcal{S}, \mathbf{y} \in \mathcal{Y}$, $r(\mathbf{s}, \mathbf{x}, \mathbf{y}), g_1(\mathbf{s}, \mathbf{x}, \mathbf{y}), \dots, g_d(\mathbf{s}, \mathbf{x}, \mathbf{y})$ are concave in \mathbf{x} , 2. $\nabla_{\mathbf{x}} r(\mathbf{s}, \mathbf{x}, \mathbf{y}), \nabla_{\mathbf{x}} g_1(\mathbf{s}, \mathbf{x}, \mathbf{y}), \dots, \nabla_{\mathbf{x}} g_d(\mathbf{s}, \mathbf{x}, \mathbf{y}), \nabla_{\mathbf{y}} r(\mathbf{s}, \mathbf{x}, \mathbf{y}), \nabla_{\mathbf{y}} g_1(\mathbf{s}, \mathbf{x}, \mathbf{y}), \dots, \nabla_{\mathbf{y}} g_d(\mathbf{s}, \mathbf{x}, \mathbf{y})$ exist, for all $\mathbf{s} \in \mathcal{S}, \mathbf{x} \in \mathcal{X}, \mathbf{y} \in \mathcal{Y}$, 4. $p(\mathbf{s}' | \mathbf{s}, \mathbf{x}, \mathbf{y})$ is continuous CSD convex and differentiable in (\mathbf{x}, \mathbf{y}) , and 5. Slater's condition holds, i.e., $\forall \mathbf{s} \in \mathcal{S}, \mathbf{x} \in \mathcal{X}, \exists \hat{\mathbf{y}} \in \mathcal{Y}$ s.t. $g_k(\mathbf{s}, \mathbf{x}, \hat{\mathbf{y}}) > 0$, for all $k = 1, \dots, d$ and $r_j(\hat{\mathbf{y}}) > 0$, for all $j = 1, \dots, l$, and $\exists \mathbf{x} \in \mathbb{R}^n$ s.t. $q_k(\mathbf{x}) < 0$ for all $k = 1 \dots, p$. Then, there exists $\boldsymbol{\mu}^*: \mathcal{S} \rightarrow \mathbb{R}_+^p, \boldsymbol{\lambda}^*: \mathcal{S} \times \mathcal{X} \rightarrow \mathbb{R}_+^d$, and $\boldsymbol{\nu}^*: \mathcal{S} \times \mathcal{X} \rightarrow \mathbb{R}_+^l$ s.t. a policy profile $(\boldsymbol{\pi}_{\mathbf{x}}^*, \boldsymbol{\pi}_{\mathbf{y}}^*) \in \mathcal{X}^{\mathcal{S}} \times \mathcal{Y}^{\mathcal{S}}$ is a recSE of \mathcal{G} only if it satisfies the following conditions, for all $\mathbf{s} \in \mathcal{S}$:

$$\nabla_{\mathbf{x}} \mathcal{L}_{\mathbf{s}, \boldsymbol{\pi}_{\mathbf{x}}^*(\mathbf{s})}(\boldsymbol{\pi}_{\mathbf{y}}^*(\mathbf{s}), \boldsymbol{\lambda}^*(\mathbf{s}, \boldsymbol{\pi}_{\mathbf{x}}^*(\mathbf{s}))) + \sum_{k=1}^p \mu_k^*(\mathbf{s}) \nabla_{\mathbf{x}} q_k(\boldsymbol{\pi}_{\mathbf{x}}^*(\mathbf{s})) = 0 \quad (1)$$

$$\nabla_{\mathbf{y}} \mathcal{L}_{\mathbf{s}, \boldsymbol{\pi}_{\mathbf{x}}^*(\mathbf{s})}(\boldsymbol{\pi}_{\mathbf{y}}^*(\mathbf{s}), \boldsymbol{\lambda}^*(\mathbf{s}, \boldsymbol{\pi}_{\mathbf{x}}^*(\mathbf{s}))) + \sum_{k=1}^l \nu_k^*(\mathbf{s}, \boldsymbol{\pi}_{\mathbf{x}}^*(\mathbf{s})) \nabla_{\mathbf{x}} r_k(\boldsymbol{\pi}_{\mathbf{y}}^*(\mathbf{s})) = 0 \quad (2)$$

$$\mu_k^*(\mathbf{s}) q_k(\boldsymbol{\pi}_{\mathbf{x}}^*(\mathbf{s})) = 0 \quad q_k(\boldsymbol{\pi}_{\mathbf{x}}^*(\mathbf{s})) \leq 0 \quad \forall k \in [p] \quad (3)$$

$$g_k(\mathbf{s}, \boldsymbol{\pi}_{\mathbf{x}}^*(\mathbf{s}), \boldsymbol{\pi}_{\mathbf{y}}^*(\mathbf{s})) \geq 0 \quad \lambda_k^*(\mathbf{s}, \boldsymbol{\pi}_{\mathbf{x}}^*(\mathbf{s})) g_k(\mathbf{s}, \boldsymbol{\pi}_{\mathbf{x}}^*(\mathbf{s}), \boldsymbol{\pi}_{\mathbf{y}}^*(\mathbf{s})) = 0 \quad \forall k \in [d] \quad (4)$$

$$\nu_k^*(\mathbf{s}, \boldsymbol{\pi}_{\mathbf{x}}^*(\mathbf{s})) \nabla_{\mathbf{x}} r_k(\boldsymbol{\pi}_{\mathbf{y}}^*(\mathbf{s})) = 0 \quad r_k(\boldsymbol{\pi}_{\mathbf{y}}^*(\mathbf{s})) \geq 0 \quad \forall k \in [l] \quad (5)$$

Proof of Theorem 3.1. By Theorem 2.2 and Theorem 2.4 we know that $(\boldsymbol{\pi}_{\mathbf{x}}^*, \boldsymbol{\pi}_{\mathbf{y}}^*)$ is a recursive Stackelberg equilibrium iff

$$v^{\boldsymbol{\pi}_{\mathbf{y}}^* \boldsymbol{\pi}_{\mathbf{x}}^*}(\mathbf{s}) = \left(C v^{\boldsymbol{\pi}_{\mathbf{y}}^* \boldsymbol{\pi}_{\mathbf{x}}^*} \right) (\mathbf{s}). \quad (70)$$

Note that for any policy profile (π_x^*, π_y^*) that satisfies $v^{\pi_y^* \pi_x^*}(s) = (Cv^{\pi_y^* \pi_x^*})(s)$ by definition of C we have that $(\pi_x^*(s), \pi_y^*(s))$ is a Stackelberg equilibrium of

$$\min_{x \in \mathcal{X}} \max_{y \in \mathcal{Y}: g(s, x, y) \geq 0} \left\{ r(s, x, y) + \gamma \mathbb{E}_{S' \sim p(\cdot | s, x, y)} [v(S')] \right\}$$

for all $s \in \mathcal{S}$.

Fix a state $s \in \mathcal{S}$, under the assumptions of the theorem, the conditions of Theorem C.3 are satisfied and $u^*(s, x) = \max_{y \in \mathcal{Y}: g(s, x, y) \geq 0} \{r(s, x, y) + \gamma \mathbb{E}_{S' \sim p(\cdot | s, x, y)} [v(S')]\}$ is subdifferentiable in x . Since $u^*(s, x)$ is convex in x , and Slater's condition are satisfied by the assumptions of the theorem, the necessary and sufficient conditions for $\pi_x^*(s)$ to be an optimal solution to $\min_{x \in \mathcal{X}} u^*(s, x)$ are given by the KKT conditions Kuhn & Tucker (1951) for $\min_{x \in \mathcal{X}} u^*(s, x)$. Note that we can state the first order KKT conditions explicitly thanks to the subdifferential Benveniste-Scheinkman theorem (Theorem C.3). That is, $\pi_x^*(s)$ is an optimal solution to $\min_{x \in \mathcal{X}} u^*(s, x)$ if there exists $\mu^*(s) \in \mathbb{R}_+^p$ such that:

$$\nabla_x \mathcal{L}_{s, \pi_x^*(s)}(y^*(s, \pi_x^*(s)), \lambda^*(s, \pi_x^*(s), y^*(s, \pi_x^*(s)))) + \sum_{k=1}^p \mu_k^*(s) \nabla_x q_k(\pi_x^*(s)) = 0 \quad (71)$$

$$\mu_k^*(s) q_k(\pi_x^*(s)) = 0 \quad \forall k \in [p] \quad (72)$$

$$q_k(\pi_x^*(s)) \leq 0 \quad \forall k \in [p] \quad (73)$$

where $y^*(s, \pi_x^*(s)) \in \arg \max_{y \in \mathcal{Y}: g(s, \pi_x^*(s), y) \geq 0} \{r(s, \pi_x^*(s), y) + \gamma \mathbb{E}_{S' \sim p(\cdot | s, \pi_x^*(s), y)} [v(S')]\}$ and

$\lambda^*(s, \pi_x^*(s), y^*(s, \pi_x^*(s))) = (\lambda_1^*(s, \pi_x^*(s), y^*(s, \pi_x^*(s))), \dots, \lambda_d^*(s, \pi_x^*(s), y^*(s, \pi_x^*(s))))^T \in \Lambda^*(s, \pi_x^*(s), y^*(s, \pi_x^*(s)))$ are the optimal KKT multipliers associated with $y^*(s, \pi_x^*(s)) \in \mathcal{Y}^*(s, \pi_x^*(s))$ which are guaranteed to exist since Slater's condition is satisfied for $\max_{y \in \mathcal{Y}: g(s, x, y) \geq 0} \{r(s, x, y) + \gamma \mathbb{E}_{S' \sim p(\cdot | s, x, y)} [v(S')]\}$.

Similarly, fix a state $s \in \mathcal{S}$ and action for the outer player $x \in \mathcal{X}$, since Slater's condition is satisfied for $\max_{y \in \mathcal{Y}: g(s, x, y) \geq 0} \{r(s, x, y) + \gamma \mathbb{E}_{S' \sim p(\cdot | s, x, y)} [v(S')]\}$, the necessary conditions for $\pi_y^*(s)$ to be a Stackelberg equilibrium strategy for the inner player at state s are given by the KKT conditions for $\max_{y \in \mathcal{Y}: g(s, x, y) \geq 0} \{r(s, x, y) + \gamma \mathbb{E}_{S' \sim p(\cdot | s, x, y)} [v(S')]\}$. That is, there exists $\lambda^*(s, x) \in \mathbb{R}_+^d$ and $\nu^*(s, x) \in \mathbb{R}_+^l$ such that:

$$\nabla_y \mathcal{L}_{s, x}(\pi_y^*(s), \lambda^*(s, x)) + \sum_{k=1}^l \nu_k^*(s) \nabla_x r_k(\pi_y^*(s)) = 0 \quad (74)$$

$$g_k(s, x, \pi_y^*(s)) \geq 0 \quad \forall k \in [d] \quad (75)$$

$$\lambda_k^*(s, x) g_k(s, x, \pi_y^*(s)) = 0 \quad \forall k \in [d] \quad (76)$$

$$\nu_k^*(s, x) \nabla_x r_k(\pi_y^*(s)) = 0 \quad \forall k \in [l] \quad (77)$$

$$r_k(x) \geq 0 \quad \forall k \in [l] \quad (78)$$

Combining the necessary and sufficient conditions in Equations (71) to (73) with the necessary conditions in Equations (74) to (78), we obtain the necessary conditions for (π_x^*, π_y^*) to be a recursive Stackelberg equilibrium. \square

Theorem C.4 (Recursive Stackelberg Equilibrium Necessary and Sufficient Optimality Conditions). *Consider a zero-sum Stochastic Stackelberg game $(\mathcal{S}, \mathcal{X}, \mathcal{Y}, r, g, p, \gamma)$, where $\mathcal{X} = \{x \in \mathbb{R}^n \mid q_1(x) \leq 0, \dots, q_p(x) \leq 0\}$ and $\mathcal{Y} = \{y \in \mathbb{R}^m \mid r_1(y) \geq 0, \dots, r_l(y) \geq 0\}$. Let $\mathcal{L}_{s, x}(y, \lambda) = r(s, x, y) + \gamma \mathbb{E}_{S' \sim p(\cdot | s, x, y)} [v(S', x)] + \sum_{k=1}^d \lambda_k g_k(s, x, y)$. Suppose that Assumption 1.1 holds, and that 1. for all $s \in \mathcal{S}, y \in \mathcal{Y}, r(s, x, y), g_1(s, x, y), \dots, g_d(s, x, y)$*

are convex-concave in (\mathbf{x}, \mathbf{y}) , 2. $\nabla_{\mathbf{x}} r(\mathbf{s}, \mathbf{x}, \mathbf{y}), \nabla_{\mathbf{x}} g_1(\mathbf{s}, \mathbf{x}, \mathbf{y}), \dots, \nabla_{\mathbf{x}} g_d(\mathbf{s}, \mathbf{x}, \mathbf{y}), \nabla_{\mathbf{y}} r(\mathbf{s}, \mathbf{x}, \mathbf{y}), \nabla_{\mathbf{y}} g_1(\mathbf{s}, \mathbf{x}, \mathbf{y}), \dots, \nabla_{\mathbf{y}} g_d(\mathbf{s}, \mathbf{x}, \mathbf{y})$ exist for all $\mathbf{s} \in \mathcal{S}, \mathbf{x} \in \mathcal{X}, \mathbf{y} \in \mathcal{Y}$, 4. $p(\mathbf{s}' \mid \mathbf{s}, \mathbf{x}, \mathbf{y})$ is continuous, CSD convex in \mathbf{x} , CSD concave in \mathbf{y} , and differentiable in (\mathbf{x}, \mathbf{y}) , and 5. Slater's condition holds, i.e., $\forall \mathbf{s} \in \mathcal{S}, \mathbf{x} \in \mathcal{X}, \exists \hat{\mathbf{y}} \in \mathcal{Y}$ s.t. $g_k(\mathbf{s}, \mathbf{x}, \hat{\mathbf{y}}) > 0$, for all $k = 1, \dots, d$ and $r_j(\hat{\mathbf{y}}) > 0$ for all $j = 1, \dots, l$, and $\exists \mathbf{x} \in \mathbb{R}^n$ s.t. $q_k(\mathbf{x}) < 0$ for all $k = 1 \dots, p$. Then, there exists $\boldsymbol{\mu}^* : \mathcal{S} \rightarrow \mathbb{R}_+^p, \boldsymbol{\lambda}^* : \mathcal{S} \times \mathcal{X} \rightarrow \mathbb{R}_+^d$, and $\boldsymbol{\nu}^* : \mathcal{S} \times \mathcal{X} \rightarrow \mathbb{R}_+^l$ such that a policy profile $(\boldsymbol{\pi}_{\mathbf{x}}^*, \boldsymbol{\pi}_{\mathbf{y}}^*) \in \mathcal{X}^{\mathcal{S}} \times \mathcal{Y}^{\mathcal{S}}$ is a recursive Stackelberg equilibrium of $(\mathcal{S}, \mathcal{X}, \mathcal{Y}, r, \mathbf{g}, p, \gamma)$ only if it satisfies the following conditions for all $\mathbf{s} \in \mathcal{S}$:

$$v^{\boldsymbol{\pi}_{\mathbf{y}}^*} \boldsymbol{\pi}_{\mathbf{y}}^*(\mathbf{s}) = \left(C v^{\boldsymbol{\pi}_{\mathbf{y}}^*} \boldsymbol{\pi}_{\mathbf{y}}^* \right) (\mathbf{s}) \quad (79)$$

$$\nabla_{\mathbf{x}} \mathcal{L}_{\mathbf{s}, \boldsymbol{\pi}_{\mathbf{x}}^*(\mathbf{s})}(\boldsymbol{\pi}_{\mathbf{y}}^*(\mathbf{s}), \boldsymbol{\lambda}^*(\mathbf{s}, \boldsymbol{\pi}_{\mathbf{x}}^*(\mathbf{s}))) + \sum_{k=1}^p \mu_k^*(\mathbf{s}) \nabla_{\mathbf{x}} q_k(\boldsymbol{\pi}_{\mathbf{x}}^*(\mathbf{s})) = 0 \quad (80)$$

$$\nabla_{\mathbf{y}} \mathcal{L}_{\mathbf{s}, \boldsymbol{\pi}_{\mathbf{x}}^*(\mathbf{s})}(\boldsymbol{\pi}_{\mathbf{y}}^*(\mathbf{s}), \boldsymbol{\lambda}^*(\mathbf{s}, \boldsymbol{\pi}_{\mathbf{x}}^*(\mathbf{s}))) + \sum_{k=1}^l \nu_k^*(\mathbf{s}, \boldsymbol{\pi}_{\mathbf{x}}^*(\mathbf{s})) \nabla_{\mathbf{y}} r_k(\boldsymbol{\pi}_{\mathbf{y}}^*(\mathbf{s})) = 0 \quad (81)$$

$$\mu_k^*(\mathbf{s}) q_k(\boldsymbol{\pi}_{\mathbf{x}}^*(\mathbf{s})) = 0 \quad \forall k \in [p] \quad (82)$$

$$q_k(\boldsymbol{\pi}_{\mathbf{x}}^*(\mathbf{s})) \leq 0 \quad \forall k \in [p] \quad (83)$$

$$g_k(\mathbf{s}, \boldsymbol{\pi}_{\mathbf{x}}^*(\mathbf{s}), \boldsymbol{\pi}_{\mathbf{y}}^*(\mathbf{s})) \geq 0 \quad \forall k \in [d] \quad (84)$$

$$\lambda_k^*(\mathbf{s}, \boldsymbol{\pi}_{\mathbf{x}}^*(\mathbf{s})) g_k(\mathbf{s}, \boldsymbol{\pi}_{\mathbf{x}}^*(\mathbf{s}), \boldsymbol{\pi}_{\mathbf{y}}^*(\mathbf{s})) = 0 \quad \forall k \in [d] \quad (85)$$

$$\nu_k^*(\mathbf{s}, \boldsymbol{\pi}_{\mathbf{x}}^*(\mathbf{s})) \nabla_{\mathbf{x}} r_k(\boldsymbol{\pi}_{\mathbf{y}}^*(\mathbf{s})) = 0 \quad \forall k \in [l] \quad (86)$$

$$r_k(\boldsymbol{\pi}_{\mathbf{x}}^*(\mathbf{s})) \geq 0 \quad \forall k \in [l] \quad (87)$$

D OMITTED RESULTS AND PROOFS SECTION 4

Before, we introduce the stochastic Stackelberg game whose recursive Stackelberg equilibria correspond to recursive competitive equilibria of an associate stochastic Fisher market, we introduce the following technical lemma, which provides the necessary and sufficient conditions for an allocation and saving system of a buyer to be expected utility maximizing.

Lemma D.1. For any price system $\mathbf{p} \in \mathbb{R}_+^{\mathcal{S} \times m}$, a tuple $(\mathbf{x}_i^*, \beta_i^*) \in \mathbb{R}_+^{\mathcal{S} \times n \times m} \times \mathbb{R}_+^{n \times m}$ consisting of an allocation system and saving system for a buyer i , given by a continuous, and homogeneous utility function $u_i : \mathbb{R}_+^m \times \mathcal{T} \rightarrow \mathbb{R}$ representing a locally non-satiated preference is expected utility maximizing constrained by the saving and spending constrained, i.e., the first condition of a competitive equilibrium is satisfied only we have for all states $\mathbf{s} \in \mathcal{S}, \mathbf{x}_i^*(\mathbf{t}, \mathbf{b}, \mathbf{q}) \cdot \mathbf{p}(\mathbf{t}, \mathbf{b}, \mathbf{q}) + \beta_i^*(\mathbf{t}, \mathbf{b}, \mathbf{q}) \leq b_i$, and,

$$x_{ij}^*(\mathbf{s}) > 0 \implies \frac{\partial u_i(\mathbf{x}_i^*(\mathbf{s}); t_i)}{\partial x_{ij}} = \frac{u_i(\mathbf{x}_i^*(\mathbf{s}); t_i)}{p_j} = \frac{u_i(\mathbf{x}_i^*(\mathbf{s}); t_i)}{b_i - \beta_i^*(\mathbf{s})} \quad \forall j \in [m] \quad (88)$$

$$\beta_i^*(\mathbf{s}) > 0 \implies \frac{\partial \nu_i}{\partial b_i}(\mathbf{s}) = \gamma \frac{\partial}{\partial \beta_i} \mathbb{E}_{(\mathbf{t}', \mathbf{b}', \mathbf{q}')} [\nu_i(\mathbf{t}', \mathbf{b}' + \beta_i^*(\mathbf{s}), \mathbf{q}')] \quad (89)$$

If additionally, u_i is concave, then the above condition become also sufficient.

Proof of Lemma D.1. Let $\mathcal{L}(\mathbf{x}_i, \beta_i, \lambda, \boldsymbol{\mu}; \mathbf{p}) = u_i(\mathbf{x}_i; t_i) + \gamma \mathbb{E}_{(\mathbf{t}', \mathbf{b}', \mathbf{q}') \sim p(\cdot \mid \mathbf{t}, \mathbf{b}, \mathbf{q}, (\mathbf{x}_i, \mathbf{X}_{-i}^*(\mathbf{s}), (\beta_i, \beta_{-i}^*(\mathbf{s})))} [\nu_i(\mathbf{t}', \mathbf{b}' + \beta_i, \mathbf{q}')] + \lambda(b_i - \mathbf{x}_i \mathbf{p}) + \sum_{j \in [m]} \mu_j x_{ij} + \mu_{m+1} \beta_i$ be the Lagrangian associated with

$$\nu_i(\mathbf{t}, \mathbf{b}, \mathbf{q}) = \max_{(\mathbf{x}_i, \beta_i) \in \mathbb{R}_+^{m+1} : \mathbf{x}_i \cdot \mathbf{p}(\mathbf{t}) + \beta_i \leq b_i} u_i(\mathbf{x}_i, t_i) + \gamma \mathbb{E}_{(\mathbf{t}', \mathbf{b}', \mathbf{q}')} [\nu_i(\mathbf{t}', \mathbf{b}' + \beta_i, \mathbf{q}')] .$$

Then, for any $b_i > 0$, Slater's condition holds and as the objective is concave (by the Weierstrass M-test and the uniform limit theorem), the KKT first order necessary and sufficient optimality conditions for an allocation $\mathbf{x}_i^* \in \mathbb{R}_+^m$, saving $\beta_i^* \in \mathbb{R}_+$ and associated Lagrangian multipliers $\lambda^* \in \mathbb{R}_+, \boldsymbol{\mu}^* \in \mathbb{R}^{m+1}$ to be optimal for any prices $\mathbf{p} \in \mathbb{R}_+^m$ are given by the following pair of

equations:

$$\frac{\partial u_i}{\partial x_{ij}}(\mathbf{x}_i^*(\mathbf{s}); t_i) - \lambda^*(\mathbf{s})p_j(\mathbf{s}) + \mu_j^* \doteq 0 \quad (90)$$

$$\gamma \frac{\partial}{\partial \beta_i} \mathbb{E}_{(\mathbf{t}', \mathbf{b}', \mathbf{q}')} [\nu_i(\mathbf{t}', \mathbf{b}' + \beta_i^*, \mathbf{q}')] - \lambda_i^*(\mathbf{s}) + \mu_{m+1}^* \doteq 0 \quad (91)$$

Additionally, by the KKT complementarity conditions, we have for all $j \in [m]$, $\mu_j^* x_{ij}^* = 0$ and $\mu_{m+1}^* \beta_i^* = 0$, which gives us:

$$x_{ij}^*(\mathbf{s}) > 0 \implies \lambda^* = \frac{\frac{\partial u_i}{\partial x_{ij}}(\mathbf{x}_i^*(\mathbf{s}); t_i)}{p_j} \quad \forall i \in [n], j \in [m] \quad (92)$$

$$\beta_i^*(\mathbf{s}) > 0 \implies \gamma \frac{\partial}{\partial \beta_i} \mathbb{E}_{(\mathbf{t}', \mathbf{b}', \mathbf{q}')} [\nu_i(\mathbf{t}', \mathbf{b}' + \beta_i^*, \mathbf{q}')] - \lambda^*(\mathbf{s}) = 0 \quad \forall i \in [n], j \in [m] \quad (93)$$

Re-organizing expressions, yields:

$$x_{ij}^*(\mathbf{s}) > 0 \implies \lambda^* = \frac{\frac{\partial u_i}{\partial x_{ij}}(\mathbf{x}_i^*(\mathbf{s}); t_i)}{p_j} \quad \forall i \in [n], j \in [m] \quad (94)$$

$$\beta_i^*(\mathbf{s}) > 0 \implies \lambda^*(\mathbf{s}) = \gamma \frac{\partial}{\partial \beta_i} \mathbb{E}_{(\mathbf{t}', \mathbf{b}', \mathbf{q}')} [\nu_i(\mathbf{t}', \mathbf{b}' + \beta_i^*, \mathbf{q}')] \quad \forall i \in [n] \quad (95)$$

Using the envelope theorem, we can also compute $\frac{\partial \nu_i}{\partial b_i}(\mathbf{s})$ as follows:

$$\frac{\partial \nu_i}{\partial b_i}(\mathbf{s}) = \lambda^*(b_i) \quad (96)$$

Hence, combining the above with Equation (94) and Equation (95), we get:

$$x_{ij}^*(\mathbf{s}) > 0 \implies \lambda^* = \frac{\frac{\partial u_i}{\partial x_{ij}}(\mathbf{x}_i^*(\mathbf{s}); t_i)}{p_j} \quad \forall i \in [n], j \in [m] \quad (97)$$

$$\beta_i^*(\mathbf{s}) > 0 \implies \frac{\partial \nu_i}{\partial b_i}(\mathbf{s}) = \gamma \frac{\partial}{\partial \beta_i} \mathbb{E}_{(\mathbf{t}', \mathbf{b}', \mathbf{q}')} [\nu_i(\mathbf{t}', \mathbf{b}' + \beta_i^*, \mathbf{q}')] \quad \forall i \in [n] \quad (98)$$

Finally, going back to Equation (90), multiplying by $x_{ij}^*(\mathbf{s})$ and summing up across all $j \in [m]$, we obtain:

$$\sum_{j \in [m]} x_{ij}^*(\mathbf{s}) \frac{\partial u_i}{\partial x_{ij}}(\mathbf{x}_i^*(\mathbf{s}); t_i) - \lambda^*(\mathbf{s}) p_j(\mathbf{s}) x_{ij}^*(\mathbf{s}) + \mu_j^* x_{ij}^*(\mathbf{s}) = 0 \quad (99)$$

$$u_i(\mathbf{x}_i^*(\mathbf{s}); t_i) - \lambda^*(\mathbf{s}) \sum_{j \in [m]} p_j(\mathbf{s}) x_{ij}^*(\mathbf{s}) + \mu_j^* x_{ij}^*(\mathbf{s}) = 0 \quad (\text{Euler's Theorem for Homogeneous Functions}) \quad (100)$$

$$u_i(\mathbf{x}_i^*(\mathbf{s}); t_i) - \lambda^*(\mathbf{s})(b_i - \beta_i^*(\mathbf{s})) = 0 \quad (\text{Slack Complementarity}) \quad (101)$$

$$\lambda^*(\mathbf{s}) = \frac{u_i(\mathbf{x}_i^*(\mathbf{s}); t_i)}{b_i - \beta_i^*(\mathbf{s})} \quad (102)$$

Combining the above conditions, with Equation (102), and adding to it Equation (98), and ensuring that the KKT primal feasibility conditions hold as well, we obtain the following necessary and sufficient conditions that need to hold for all states $\mathbf{s} \in \mathcal{S}$:

$$\mathbf{x}_i^*(\mathbf{t}, \mathbf{b}, \mathbf{q}) \cdot \mathbf{p}(\mathbf{t}, \mathbf{b}, \mathbf{q}) + \beta_i^*(\mathbf{t}, \mathbf{b}, \mathbf{q}) \leq b_i \quad (103)$$

$$x_{ij}^*(\mathbf{s}) > 0 \implies \frac{\frac{\partial u_i}{\partial x_{ij}}(\mathbf{x}_i^*(\mathbf{s}); t_i)}{p_j} = \frac{u_i(\mathbf{x}_i^*(\mathbf{s}); t_i)}{b_i - \beta_i^*(\mathbf{s})} \quad \forall j \in [m] \quad (104)$$

$$\beta_i^*(\mathbf{s}) > 0 \implies \frac{\partial \nu_i}{\partial b_i}(\mathbf{s}) = \gamma \frac{\partial}{\partial \beta_i} \mathbb{E}_{(\mathbf{t}', \mathbf{b}', \mathbf{q}')} [\nu_i(\mathbf{t}', \mathbf{b}' + \beta_i^*(\mathbf{s}), \mathbf{q}')] \quad (105)$$

□

Theorem 4.1. A stochastic Fisher market with savings $(\mathcal{S}, \mathcal{U}, \mathbf{b}^{(0)}, p, \gamma)$ in which \mathcal{U} is a vector of continuous and homogeneous utility functions has at least one recCE. Additionally, the recSE $(\mathbf{p}^*, \mathbf{X}^*, \boldsymbol{\beta}^*)$ that solves the following Bellman equation corresponds to the recCE of $(\mathcal{S}, \mathcal{U}, \mathbf{b}^{(0)}, p, \gamma)$:

$$v(\mathbf{t}, \mathbf{b}, \mathbf{q}) = \min_{\mathbf{p} \in \mathbb{R}_+^m} \max_{(\mathbf{X}, \boldsymbol{\beta}) \in \mathbb{R}_+^{n \times (m+1)}: \mathbf{X}\mathbf{p} + \boldsymbol{\beta} \leq \mathbf{b}} \sum_{j \in [m]} q_j p_j + \sum_{i \in [n]} (b_i - \beta_i) \log(u_i(\mathbf{x}_i, t_i)) + \gamma \mathbb{E}_{(\mathbf{t}', \mathbf{b}', \mathbf{q}') \sim p(\cdot | \mathbf{t}, \mathbf{b}, \mathbf{q}, \mathbf{X}, \boldsymbol{\beta})} [v(\mathbf{t}', \mathbf{b}' + \boldsymbol{\beta}, \mathbf{q}')] \quad (6)$$

Proof of Theorem 4.1. Define $\mathcal{L}(\mathbf{p}, \mathbf{X}, \boldsymbol{\beta}, \boldsymbol{\lambda}) = \sum_{j \in [m]} q_j p_j + \sum_{i \in [n]} (b_i - \beta_i) \log(u_i(\mathbf{x}_i, t_i)) + \gamma \mathbb{E}_{(\mathbf{t}', \mathbf{b}', \mathbf{q}') \sim p(\cdot | \mathbf{t}, \mathbf{b}, \mathbf{q}, \mathbf{X}, \boldsymbol{\beta})} [v(\mathbf{t}', \mathbf{b}' + \boldsymbol{\beta}, \mathbf{q}')] + \sum_{j \in [m]} \lambda_j (b_j - \mathbf{x}_j \cdot \mathbf{p} + \beta_j)$. By Theorem 3.1, the necessary optimality conditions for the stochastic Stackelberg game

$$\max_{(\mathbf{X}, \boldsymbol{\beta}) \in \mathbb{R}_+^{n \times m} \times \mathbb{R}_+^n: \mathbf{X}\mathbf{p} + \boldsymbol{\beta} \leq \mathbf{b}} \sum_{j \in [m]} q_j p_j + \sum_{i \in [n]} (b_i - \beta_i) \log(u_i(\mathbf{x}_i, t_i)) + \gamma \mathbb{E}_{(\mathbf{t}', \mathbf{b}', \mathbf{q}') \sim p(\cdot | \mathbf{t}, \mathbf{b}, \mathbf{q}, \mathbf{X}, \boldsymbol{\beta})} [v(\mathbf{t}', \mathbf{b}' + \boldsymbol{\beta}, \mathbf{q}')] \quad (7)$$

are that for all states $\mathbf{s} \in \mathcal{S}$ there exists $\boldsymbol{\mu}^*(\mathbf{s}) \in \mathbb{R}^{n \times (m+1)}$, and $\boldsymbol{\nu}^*(\mathbf{s}) \in \mathbb{R}_+^{n \times (m+1)}$ associated with the non-negativity constraints for $(\mathbf{X}, \boldsymbol{\beta})$, and \mathbf{p} respectively, and $\boldsymbol{\lambda}^*(\mathbf{s}) \in \mathbb{R}_+^m$ such that:

$$q_j - \sum_{i \in [n]} \lambda_i^*(\mathbf{s}) x_{ij}^*(\mathbf{s}) - \nu_j^*(\mathbf{s}) \doteq 0 \quad (106)$$

$$\frac{b_i - \beta_i^*(\mathbf{s})}{u_i(\mathbf{x}_i^*(\mathbf{s}))} \frac{\partial u_i}{\partial x_{ij}}(\mathbf{x}_i^*(\mathbf{s}); t_i) - \lambda_i^*(\mathbf{s}) p_j(\mathbf{s}) + \mu_{ij}^*(\mathbf{s}) \doteq 0 \quad \forall i \in [n], j \in [m] \quad (107)$$

$$-\log(u_i(\mathbf{x}_i^*(\mathbf{s}))) + \gamma \frac{\partial}{\partial \beta_i} \mathbb{E}_{(\mathbf{t}', \mathbf{b}', \mathbf{q}')} [v(\mathbf{t}', \mathbf{b}' + \boldsymbol{\beta}^*(\mathbf{s}), \mathbf{q}')] - \lambda_i^*(\mathbf{s}) + \mu_{i(m+1)}^*(\mathbf{s}) \doteq 0 \quad \forall i \in [n] \quad (108)$$

Note that by Theorem 3.1, we also have $\mu_{i(m+1)}^*(\mathbf{s}) \beta_i^*(\mathbf{s}) = \mu_{i+1}^*(\mathbf{s}) x_{ij}^*(\mathbf{s}) = 0$ which gives us:

$$p_j(\mathbf{s}) > 0 \implies q_j - \sum_{i \in [n]} \lambda_i^*(\mathbf{s}) x_{ij}^*(\mathbf{s}) = 0 \quad \forall [m] \in [m] \quad (109)$$

$$x_{ij}^*(\mathbf{s}) > 0 \implies \frac{b_i - \beta_i^*(\mathbf{s})}{u_i(\mathbf{x}_i^*(\mathbf{s}))} \frac{\partial u_i}{\partial x_{ij}}(\mathbf{x}_i^*(\mathbf{s})) - \lambda_i^*(\mathbf{s}) p_j(\mathbf{s}) = 0 \quad \forall i \in [n], j \in [m] \quad (110)$$

$$\beta_i^*(\mathbf{s}) > 0 \implies -\log(u_i(\mathbf{x}_i^*(\mathbf{s}); t_i)) + \gamma \frac{\partial}{\partial \beta_i} \mathbb{E}_{(\mathbf{t}', \mathbf{b}', \mathbf{q}')} [v(\mathbf{t}', \mathbf{b}' + \boldsymbol{\beta}^*(\mathbf{s}), \mathbf{q}')] - \lambda_i^*(\mathbf{s}) + \mu_{i(m+1)}^*(\mathbf{s}) = 0 \quad \forall i \in [n] \quad (111)$$

Re-organizing expressions, we obtain:

$$p_j(\mathbf{s}) > 0 \implies q_j = \sum_{i \in [n]} \lambda_i^*(\mathbf{s}) x_{ij}^*(\mathbf{s}) \quad \forall j \in [m] \quad (112)$$

$$x_{ij}^*(\mathbf{s}) > 0 \implies \frac{u_i(\mathbf{x}_i^*(\mathbf{s}))}{b_i - \beta_i^*(\mathbf{s})} \lambda_i^*(\mathbf{s}) = \frac{\frac{\partial u_i}{\partial x_{ij}}(\mathbf{x}_i^*(\mathbf{s}); t_i)}{p_j(\mathbf{s})} \quad \forall i \in [n], j \in [m] \quad (113)$$

$$\beta_i^*(\mathbf{s}) > 0 \implies -\log(u_i(\mathbf{x}_i^*(\mathbf{s}))) + \gamma \frac{\partial}{\partial \beta_i} \mathbb{E}_{(\mathbf{t}', \mathbf{b}', \mathbf{q}')} [v(\mathbf{t}', \mathbf{b}' + \boldsymbol{\beta}^*(\mathbf{s}), \mathbf{q}')] - \lambda_i^*(\mathbf{s}) = 0 \quad \forall i \in [n] \quad (114)$$

Using the envelope theorem, we can compute $\frac{\partial v}{\partial b_i}$ as follows:

$$\frac{\partial v^*}{\partial b_i}(\mathbf{s}) = \log(u_i(\mathbf{x}_i^*(\mathbf{s}); t_i)) + \lambda_i^*(\mathbf{s}) \quad (115)$$

Re-organizing expressions, we get:

$$\lambda_i^*(\mathbf{s}) = \frac{\partial v}{\partial b_i}(\mathbf{s}) - \log(u_i(\mathbf{x}_i^*(\mathbf{s}); t_i)) \quad (116)$$

Combining Equation (116) and Equation (114), we obtain:

$$\beta_i^*(\mathbf{s}) > 0 \implies -\log(u_i(\mathbf{x}_i^*(\mathbf{s}); t_i)) + \gamma \frac{\partial}{\partial \beta_i} \mathbb{E}_{(\mathbf{t}', \mathbf{b}', \mathbf{q}')} [v(\mathbf{t}', \mathbf{b}' + \beta^*(\mathbf{s}), \mathbf{q}')] - \frac{\partial v}{\partial b_i}(\mathbf{b}) + \log(u_i(\mathbf{x}_i^*(\mathbf{b}); t_i)) = 0 \quad \forall i \in [n] \quad (117)$$

$$\beta_i^*(\mathbf{s}) > 0 \implies \gamma \frac{\partial}{\partial \beta_i} \mathbb{E}_{(\mathbf{t}', \mathbf{b}', \mathbf{q}')} [v(\mathbf{t}', \mathbf{b}' + \beta^*(\mathbf{s}), \mathbf{q}')] - \frac{\partial v}{\partial b_i}(\mathbf{b}) = 0 \quad \forall i \in [n] \quad (118)$$

$$\beta_i^*(\mathbf{s}) > 0 \implies \frac{\partial v}{\partial b_i}(\mathbf{b}) = \gamma \frac{\partial}{\partial \beta_i} \mathbb{E}_{(\mathbf{t}', \mathbf{b}', \mathbf{q}')} [v(\mathbf{t}', \mathbf{b}' + \beta^*(\mathbf{s}), \mathbf{q}')] \quad \forall i \in [n] \quad (119)$$

Going back to Equation (107), multiplying both sides by $x_{ij}^*(\mathbf{s})$ and summing up across all $j \in [m]$, we get:

$$\sum_{j \in [m]} \frac{b_i - \beta_i^*(\mathbf{s})}{u_i(\mathbf{x}_i^*(\mathbf{s}); t_i)} \sum_{j \in [m]} x_{ij}(\mathbf{s})^* \frac{\partial u_i}{\partial x_{ij}}(\mathbf{x}_i^*(\mathbf{s})) - \lambda_i^*(\mathbf{s}) \sum_{j \in [m]} p_j(\mathbf{s}) x_{ij}^*(\mathbf{s}) = 0 \quad (120)$$

$$\frac{b_i - \beta_i^*(\mathbf{s})}{u_i(\mathbf{x}_i^*(\mathbf{s}); t_i)} u_i(\mathbf{x}_i^*(\mathbf{s}); t_i) - \lambda_i^*(\mathbf{s}) \sum_{j \in [m]} p_j(\mathbf{s}) x_{ij}^*(\mathbf{s}) = 0 \quad (\text{Euler's Theorem}) \quad (121)$$

$$b_i - \beta_i^*(\mathbf{s}) - \lambda_i^*(\mathbf{s}) \sum_{j \in [m]} p_j(\mathbf{s}) x_{ij}^*(\mathbf{s}) = 0 \quad (122)$$

By Theorem 3.1, we have that $\lambda_i^*(\mathbf{s}) \left(b_i - \sum_{j \in [m]} p_j(\mathbf{s}) x_{ij}^*(\mathbf{s}) - \beta_i^*(\mathbf{s}) \right) = 0$, which gives us:

$$b_i - \beta_i^*(\mathbf{s}) - \lambda_i^*(\mathbf{s}) (b_i - \beta_i^*(\mathbf{s})) = 0 \quad (123)$$

$$\lambda_i^*(\mathbf{s}) = 1 \quad (124)$$

Combining the above with Equations (112) to (114) we obtain:

$$p_j(\mathbf{s}) > 0 \implies q_j = \sum_{i \in [n]} x_{ij}^*(\mathbf{s}) \quad \forall j \in [m] \quad (125)$$

$$x_{ij}^*(\mathbf{s}) > 0 \implies \frac{u_i(\mathbf{x}_i^*(\mathbf{s}))}{b_i - \beta_i^*(\mathbf{s})} = \frac{\frac{\partial u_i}{\partial x_{ij}}(\mathbf{x}_i^*(\mathbf{s}); t_i)}{p_j(\mathbf{s})} \quad \forall i \in [n], j \in [m] \quad (126)$$

$$\beta_i^*(\mathbf{s}) > 0 \implies \frac{\partial v}{\partial b_i}(\mathbf{b}) = \gamma \frac{\partial}{\partial \beta_i} \mathbb{E}_{(\mathbf{t}', \mathbf{b}', \mathbf{q}')} [v(\mathbf{t}', \mathbf{b}' + \beta^*(\mathbf{s}), \mathbf{q}')] \quad \forall i \in [n] \quad (127)$$

Since the utility functions are non-satiated, and by the second equation, the buyers are utility maximizing at state \mathbf{s} over all allocations, we must also have that Walras' law holds, i.e., $\mathbf{p} \cdot \left(\mathbf{q} - \sum_{i \in [n]} \mathbf{x}_i \right) - \sum_{i \in [n]} \beta_i$. Walras' law combined with the first equation above then imply the second condition of a recursive competitive equilibrium. Finally, by Lemma D.1, the last two equations imply the first condition of recursive competitive equilibrium.

□

E EXPERIMENT DETAILS

We initialized a stochastic Fisher market with $n = 2$ buyers and $m = 5$ goods. To simplify the analysis, we assumed deterministic transitions such that the buyers do not get new budgets at each time period, and their types/valuations as well as the supply of goods does not change at each state, i.e., the type/valuation space and supply space has cardinality 1. This reduced the market to a deterministic repeated market setting in which the amount of budget saved by the buyers differentiates different states of the market. To initialize the state space of the market, we first fixed a range of $[10, 50]^m$ for the buyers' valuations and drew for all buyers $i \in [n]$ valuations θ_i from that range uniformly at random at the beginning of the experiment. We have assumed the supply of goods is $\mathbf{1}_m$ and that the budget space was $[0, 1]^n$. This means that our state space for our experiments was $\mathcal{S} = \{(\theta_1, \theta_2)\} \times \{\mathbf{1}_m\} \times [0, 1]^n$. We note that although the assumption that buyers valuations/type space has cardinality one does simplify the problem, the supply of the goods being $\mathbf{1}$ at each state is wlog because goods are divisible and the allocation of goods to buyers at each state can then be interpreted as the percentage of a particular good allocated to a buyer. We assumed initial budgets of $\mathbf{b}^{(0)} = \mathbf{1}_n$ for both buyers.

Since the state space is continuous, the value function is also continuous in the stochastic Fisher market setting. As a result, we had to use fitted variant of value iteration. In particular, we assumed that the value function had a linear form at each state such that $v(\mathbf{t}, \mathbf{b}, \mathbf{q}; \mathbf{a}, c) = \mathbf{a}^T \mathbf{b} + c$ for some parameters $\mathbf{a} \in \mathbb{R}^n$, $c \in \mathbb{R}$, and we tried to approximate the value function at the next step of value iteration by using linear regression. That is, at each value iteration step, we uniformly sampled 25 budget vectors from the range $[0, 1]^n$. Next, for each sampled budget \mathbf{b} , we solved the min-max step given that budget as a state. This process gave us (budget, value) pairs on which we ran linear regression to approximate the value function at the next iterate.

To solve the generalized min-max operator at each step of value iteration, we used two different methods for comparison: **nested gradient descent ascent (GDA)** Goktas & Greenwald (2021) (Algorithm 2), which is not guaranteed to converge to a global optimum since the min-max Stackelberg game for stochastic Fisher markets is convex-non-concave and **max-oracle gradient descent** Goktas & Greenwald (2021) where the max-oracle is the simulated annealing algorithm Bertsimas & Tsitsiklis (1993), a metaheuristic which probabilistically aims to find a global maximum. Although, simulated annealing annealing is not guaranteed to converge to a global maximum, we observed that it often performed better than nested GDA at finding a global optimum for the inner player's strategy. For both methods, we have run value iteration for 30 iterations. We ran nested GDA with learning rates $\eta_{\mathbf{X}} = 0.5$, $\eta_{\mathbf{p}} = 0.02$ for linear, $\eta_{\mathbf{X}} = 0.5$, $\eta_{\mathbf{p}} = 0.01$ for leontief, and $\eta_{\mathbf{X}} = 0.5$, $\eta_{\mathbf{p}} = 0.0005$ for Cobb-Douglas. We also ran max-oracle gradient descent with a learning rate of $\eta_{\mathbf{p}} = 0.5$. The outer loop of nested GDA was run for $T_{\mathbf{p}} = 30$ iterations, while its inner loop was run for $T_{\mathbf{X}} = 100$ iterations, and the max-oracle gradient descent algorithm was run for $T_{\mathbf{p}} = 30$ iterations. We depict the trajectory of the average value of the value function at each iteration of value iteration, under nested GDA in Section 5, and under max-oracle gradient descent in Section 5.

Utility Type	Cumulative Utility based on VI	Maximized cumulative Utility	Distance to Utility Maximization	Normalized Distance to Utility Maximization	Distance to Market Clearance (Average Excess Demands)
Linear (Nested GDA)	[96.11864217, 872.44457586]	[84.79659719559233, 90.3457392997638]	782.1807839960541	6.312670581380343	1.51360722680293
Linear (Max-oracle GD)	[58.83670238, 65.37844446]	[96.96171820101995, 106.36975094410113]	55.98038976844064	0.388939329909976	2.1650851972578646
Leontief (Nested GDA)	[0.01241015, 0.01197241]	[0.009253314180356959, 0.008070244600082384]	0.005019216377258324	0.408792752480498	2.1885091150520655
Leontief (Max-oracle GD)	[0.0044322, 0.00609951]	[0.012505356366874554, 0.00961266013048025]	0.008804435107563718	0.5581970083770711	2.178489762142351
Cobb-Douglas (Nested GDA)	[-759048.31311689, -844533.23523943]	[-5410.446761817995, -5757.897160340484]	1127614.3407095883	142.71734067443717	2.1892506104891303
Cobb-Douglas (Max-oracle GD)	[-763010.21531314, -847841.51430998]	[-5482.763563167677, -5790.554467234857]	1132650.7224766547	142.0356741620785	2.1697254289647336

Figure 3: Exploitability of the computed recursive competitive equilibrium by both methods.