

REINFORCEMENT LEARNING AND COLLUSION

CLEMENS POSSNIG

UNIVERSITY OF WATERLOO

ABSTRACT. This paper presents an analytical characterization of the long run policies learned by algorithms that interact repeatedly. The algorithms observe a state variable and update policies to maximize long term discounted payoffs. I show that their long run policies correspond to equilibria that are stable points of a tractable differential equation. As an example, I consider a repeated Bertrand game, for which learning the stage game Nash equilibrium serves as non-collusive benchmark. I give necessary and sufficient conditions for this Nash equilibrium to be learned. I show how the interplay between monitoring technology (state variables) and market conditions such as price elasticities and markups determine whether stage game Nash, or collusive equilibria may be learned. I apply this framework to analyze two key regulatory policies: limiting data inputs for algorithms, and imposing competition in the software provider market. My results demonstrate that the former strategy holds more promise.

JEL classification. C62, C73, D43, D83.

Keywords. Multi-Agent Reinforcement Learning, Repeated Games, Collusion, Learning in Games.

Date: July 12, 2025.

⁰I thank my committee members Li Hao, Vitor Farinha Luz, and Michael Peters for years of guidance and conversations. I am grateful to Rohit Lamba, Alexander Frankel, Kevin Leyton-Brown, Wei Li, Vadim Marmer, Jesse Perla, Chris Ryan, and Kevin Song for many helpful discussions. I thank the participants at EC 22, GTA22, CORS/INFORMS 22, and CETC 22 for insightful comments. I also thank participants of the theory lunches at VSE for their extensive feedback. I gratefully acknowledge support through a University of Waterloo SSHRC Institutional Grant (SIG).

1. INTRODUCTION

More and more companies are using artificial intelligence-based tools in their pursuit of profit maximization. Such algorithms take market data to determine current price levels, updating in real-time. Algorithms can help firms adapt to rapidly changing market environments, and potentially better serve their markets. However, algorithms appear to have an inherent ability to collude. Studying the German gasoline retail market, Assad et al. (2024) observe that after a critical mass of firms deployed pricing algorithms, profit margins rose by 28%. This finding, alongside numerous simulation-based studies¹, raises an urgent question: Which algorithms, and which markets, are likely to support such outcomes?

I first introduce a common family of RL algorithms that repeatedly play a game. While the results are more general, the leading application considered is Bertrand price competition. The algorithms observe a common state variable without knowing their payoff function or state transition likelihoods, and adapt by repeatedly experimenting with price choices and estimating a value function. I show that to pin down the long-run behavior of the algorithmic learners, it is enough to find the stable rest points of a differential equation.

Next, I use this characterization to study whether the algorithms can learn to repeat the static Nash equilibrium, which we can think of as the non-collusive benchmark. It turns out that the answer depends on properties of the state variable and market conditions. Learning to repeat the static Nash equilibrium may be impossible, even within classical, well-studied² stage games.

Allowing for arbitrary, finite, state variables, the possibilities for collusive behavior are ample. I therefore focus on the payoff-optimal collusive equilibrium as a bound to the possibilities of cartelization in this setting. Through an approximation exercise, I show that optimal collusive equilibria are closely related to optimal imperfect monitoring equilibria of the bang-bang kind, as characterized in Abreu, Pearce, and Stacchetti (1986). Conditions for such equilibria to be learned also come down to properties of the state variable and the market competed in.

¹Klein (2021), Calvano, Calzolari, Denicoló, et al. (2021) show that algorithms may learn to play repeated game strategies akin to typical carrot-and-stick type strategies studied in the economic theory literature.

²Linear and logit demand, constant marginal cost.

The algorithms’ state variable serves as a model of the data input chosen by the algorithm designer.³ Even without the designer’s intent, this data acts as a monitoring technology through the correlation between data outcomes and price choices of the algorithms. Whether static Nash can be learned by algorithms depends both on the sensitivity of the monitoring technology, and market conditions such as own- and opponents-price elasticities.

One can think of monitoring sensitivity as a measure of the ability to identify deviations in the continuous space, related to Fudenberg, D. Levine, and Maskin (1994)’s pairwise identifiability condition. Fixing market conditions, higher monitoring sensitivity makes learning static Nash less likely. Fixing the monitoring technology, this equilibrium is more likely to be learned if there is weaker competition and higher own-price elasticity. Turning to collusion, conditions for such equilibria to be learned are, again, related to the sensitivity of the monitoring technology, and market conditions. A less sensitive monitoring technology may make it more likely for this optimal equilibrium to be learned, but at the same time may decrease the payoffs to that equilibrium. The latter insight is based on Kandori (1992).

The intuition for these insights comes down to learning, and the possibility of making mistakes. Mistakes induce perturbations in currently played policy profiles. For perturbations to have any impact down the line, they must be detected, i.e. monitoring must be sensitive enough. On the other side, given that perturbations are detected, they must matter; here market conditions come into play. Perturbations matter less if competitor’s prices don’t have much of an impact (weak competition), or if small adjustments are enough to correct for perturbations (large own-price elasticity).

I then consider two policy options: a restriction to the data algorithms may be allowed to employ as their inputs; and imposing competition in the algorithmic software provider market. I show that the former has potential, as it can be used to decrease the sensitivity of the monitoring technology, enabling the learning of the static Nash equilibrium. On the other hand, the latter comes only with ambiguous predictions.

I finish by extending the Bertrand application to one allowing for market dynamics. This setting is especially interesting as it directly applies as a model of Assad et al. (2024)’s study on algorithmic collusion in the German gasoline retail market. I conclude that also in

³E.g., aggregate market conditions, time of day, sales volume, consumer data, etc.

this market, a restriction on data inputs can facilitate the learning of the most competitive option, static Nash.

To the best of my knowledge, this is the first theoretical work on algorithmic collusion that uncovers the relationship between standard economic concepts such as price elasticity and monitoring technologies, and the learning of collusive behavior.

Related Literature. This project speaks to results in the fast-growing literature on algorithmic collusion, the theory of learning in games, as well as the study of asymptotic behavior of algorithms in the computer science literature. A more detailed discussion can be found in the online appendix.

Firstly, the literature on algorithmic collusion has received increasing attention in recent years. Assad et al. (2024) provide an empirical study supporting the hypothesis that algorithms may learn to play collusively, while there are many simulation studies suggesting the same, of which Calvano, Calzolari, Denicolo, et al. (2020), Calvano, Calzolari, Denicoló, et al. (2021), and Klein (2021) are important examples. A paper close in spirit to this study is Banchio and Mantegazza (2022). They consider a fluid approximation technique related to the stochastic approximation approach applied here, and recover interesting phenomena regarding the learning of cooperation for finite-action class of RL algorithms without memory, with a focus on Q -learning. The family of algorithms studied here concerns learning in continuous-action games with repeated game strategies, such as repeated Bertrand or Cournot competition. While a main example is an extension of Q -learning, it cannot accommodate Q -learning as a special case, nor is it a special case of Q -learning (ACQ, see discussion below). Cartea et al. (2022) show how stochastic approximation can be applied in the analysis of finite-action reinforcement learners such as Q -learners as well. Meylahn and V. den Boer (2022), Loots and V. den Boer (2023) use ODE methods related to the ones applied in this paper to prove that specific algorithms can learn to collude in a pricing game. Further important recent work in the area of algorithmic collusion includes Lamba and Zhuk (2022), Brown and MacKay (2021), Johnson, Rhodes, and Wildenbeest (2020), and Salcedo (2015). These papers feature stylized models of algorithmic competition, abstracting away from issues of learning and estimation, which are an important aspect of my analysis.

Secondly, this paper connects to a the theory of learning in games. Classically, this literature has been concerned with the ability of agents to learn a Nash equilibrium of the stage game when following a given learning rule (e.g. Milgrom and Roberts (1991), Fudenberg and Kreps (1993)). More recent results concern learning in stochastic games (e.g. Leslie, Perkins, and Xu (2020)), where the state variable is taken as an exogenous object. The class of algorithms studied here has the ability to learn *repeated game strategies*, i.e. strategies that condition on summaries of the history of the game, implemented as automaton strategies. The games that can be studied here therefore contain stochastic games as a special case, but also allow for the case where the state that agents observe represents a finite history of the repeated interaction.

My class contains algorithms that impose little informational assumptions as a special case, known as “model free”. Such algorithms do not carry a model of opponent behavior, and also no model of their environment and own payoffs. Thus, this class falls into the family of adaptive uncoupled learning rules as defined in Hart and Mas-Colell (2003). Further foundational papers in this literature include Milgrom and Roberts (1990) , Fudenberg and D. K. Levine (2009), Gaunersdorfer and Hofbauer (1995), and many more.

Thirdly, this paper makes use and slightly generalizes an extensive body of research related to stochastic approximation theory (see e.g. Borkar (2009)) and hyperbolic theory (Palis Jr, Melo, et al. (1982)). The generalizations are relegated to the online appendix. There is a growing strand of the computer science literature devoted to establishing convergence proofs in multi-agent algorithmic environments. The paper in that area closest to this one is Mazumdar, Ratliff, and Sastry (2020).

2. MULTI-AGENT LEARNING

This section introduces the updating rule (algorithm) and main assumptions used as running example in this paper. The algorithm is known as actor-critic Q -learning (ACQ). These algorithms keep track of an estimated performance criterion (the “critic”, or Q -function, essentially a value function) and a policy function (the “actor”) that is updated towards the maximizer of the performance criterion. The policy is a mapping from observables (states), such as past prices or other market data, to actions, in this case prices. A main advantage

of ACQ over the simpler and more commonly known Q -learning (Watkins (1989)) is that it directly applies to continuous-action problems, which are the focus of this paper.

In basic Q -learning, a Q -matrix is the algorithmic updating object. As these matrices can only take finitely many elements, Q -learning is ideally suited for finite-action problems, requiring discretization in the continuous-action world. This learning rule then uses an estimation-based Bellman iteration to estimate a state-action value function, known as the Q -function. Importantly, the policy found by such an algorithm can be directly ‘read-off’ of the Q -function. This implies that Q -policies move more erratically than actor-critic methods, as such methods dampen the relationship between policy and critic. Advantages of actor-critic methods therefore include variance reduction of updates, among others.

The results presented in this paper are not unique to the case of a Q -function used as the critic. A broader characterisation of algorithms for which the results stated here hold is found in the online appendix, which also extends results to allow for non-vanishing bias in the critic estimation. As critic estimation generally involves function approximation (especially in the continuous action case), non-vanishing estimation bias is an important concern.

Actor-critic reinforcement learning, which is a superset of the class studied here, has become popular in the reinforcement learning community, due to the variance reduction property mentioned before and higher flexibility than pure critic- or actor-based methods (Q -learning being a critic-based method). See e.g. the substantial popularity of PPO (Schulman et al. (2017)). Actor-critic learning has been studied in stage game settings before, see e.g. Leslie and Collins (2003).

2.1. Payoffs and Strategies. There are N algorithms indexed by i , each having as action space an interval $\mathbf{Y}_i \subseteq \mathbb{R}$, with profile space $\mathbf{Y} = \times_i \mathbf{Y}_i$. A state variable S taking values in space \mathbf{S} with $|\mathbf{S}| = K < \infty$ comes with a transition probability function, twice differentiable in \mathbf{Y} , $T : \mathbf{S}^2 \times \mathbf{Y} \rightarrow [0, 1]$. Each algorithm has a stage game payoff function $\pi^i : \mathbf{Y} \times \mathbf{S} \rightarrow \mathbb{R}$, \mathcal{C}^2 in \mathbf{Y} , and common discount factor $\delta \in (0, 1)$.⁴

Throughout, it is important to keep in mind that I define an environment competed on not by rational agents, but by algorithms constrained to play policies based on a fixed domain:

⁴Let $\mathcal{C}^i[\mathbf{Y}, \mathbf{C}]$ be the set of functions that are i times continuously differentiable, with domain \mathbf{Y} and range \mathbf{C} . When domain and range are clear, I write \mathcal{C}^i .

S. I will take S as an exogenous object chosen by whoever initialized the algorithm. I will assume throughout that the state variable and current state s is a common observable to all algorithms.

Algorithms update a policy function $\rho_{i,t} : \mathbf{S} \rightarrow \mathbf{Y}_i$. Since states are finite, policy profile $\rho_t \in \bar{\mathbf{Y}} = \mathbf{Y}^{NL}$ can be represented as a vector in \mathbb{R}^{NL} .

Assumption 1. *For all $\rho \in \bar{\mathbf{Y}}$, the Markov chain induced by $T_{ss'}[\rho(s)]$ is irreducible and aperiodic.*⁵

In fact, one can view any such policy as a stationary Markov strategy given state variable S . Next, define $\bar{\mathbf{Y}}_i = \mathbf{Y}_i^L$, and $\bar{\mathbf{Y}}_{-i} = \times_{j \neq i} \bar{\mathbf{Y}}_j$.

Expected future discounted payoffs $W^i(\rho_i, \rho_{-i}, s_0)$ are defined given stationary policy profiles $(\rho_i, \rho_{-i}) \in \bar{\mathbf{Y}}$:

$$W^i(\rho_i, \rho_{-i}, s_0) = \mathbb{E} \sum_{t=0}^{\infty} \delta^t \pi^i(\rho(s_t), s_t), \quad (1)$$

where the expectation is taken over the randomness in the stage game payoffs and state transitions. Define $B_S^i(\rho_{-i})$ as the optimal policy for i given a profile $\rho_{-i} \in \bar{\mathbf{Y}}_{-i}$, chosen from the constraint set of stationary, S -state policies:

$$B_S^i(\rho_{-i}) = \arg \max_{\rho \in \bar{\mathbf{Y}}_i} W^i(\rho, \rho_{-i}, s_0), \quad (2)$$

where due to our assumption on irreducibility of the state space the optimal policy does not depend on the initial state s_0 . The optimal policy is indeed optimal over all possible history-dependent policies since given a Markov stationary opponent profile ρ_{-i} there must be a Markov stationary best response. In what follows, write $\bar{B}_S(\rho)$ as the stacked best response correspondence over i .

Definition 1. *Define*

(i) $E_S \subset \bar{\mathbf{Y}}$ *to be the set of Nash equilibria in policy profiles based on payoff functions W^i .*

In other words, E_S is the set of profiles ρ^ s.t. $\rho^* \in \bar{B}_S(\rho^*)$.*

⁵For definitions, see e.g. Appendix A in Puterman (2014)

(ii) $\rho^* \in E_S$ as 'differential Nash equilibrium' if ρ^* is interior, first order conditions hold for each agent at ρ^* , and the Hessian of each agent's optimization problem at ρ^* is negative definite.

If $\rho^* \in E_S$ is a differential Nash equilibrium, there is an open neighborhood U_{ρ^*} of ρ^* such that best responses must be single valued for all $\rho \in U_{\rho^*}$. Let $\mathcal{U}_S = \bigcup_{\rho^* \in E_S} U_{\rho^*}$. Given these definitions of the underlying payoff environment, assume:

Assumption 2 (Equilibrium existence and differentiability).

- (i) Given state variable S , stationary equilibrium profiles $\rho^* \in \overline{\mathbf{Y}}$ exist.
- (ii) There exist $\rho^* \in E_S$ that are differential Nash equilibria.

A sufficient condition for both points in Assumption 2 to hold is the existence of a differential static Nash equilibrium, given $\pi(a, s)$ for all $s \in \mathbf{S}$. As our analysis of limiting strategies will depend on a smoothness condition of an underlying differential equation at some equilibria, the second point will prove crucial.

2.2. Algorithmic Learning in Reduced Form. The purpose of the following analysis is to understand long-run outcomes of machine learning algorithms that satisfy consistency in their critic estimation step, in a manner introduced below. This can be seen as a benchmark for ideal individual behavior of a family of algorithms. For each algorithm i , let $F_S^i : \overline{\mathbf{Y}} \rightarrow \mathbb{R}^K$ refer to the performance criterion, i.e. critic term representing the function approximation target the learning algorithm intends to estimate.

For ACQ learning, one would have $F_S^i(\rho) = B_S^i(\rho_{-i}) - \rho_i$ ⁶, which is an extension of best response dynamics (Gilboa and Matsui (1991)) to a setting of repeated-games payoffs. A common alternative are gradient updating schemes, in which case one would take $F_S^i(\rho) = \frac{\partial}{\partial \rho_i} W^i(\rho_i, \rho_{-i}, s_0)$, the vector of payoff gradients of the learner. General results in the next section apply to both⁷. Some results in Section 4 apply only partly to both, which will be noted. When specificity is required, I will refer to the former as $F_{S,B}^i$ and to the latter as $F_{S,G}^i$.

⁶The results presented in the following can be generalized to a setting where $B_S^i(\cdot)$ is not single-valued on a positive-measure set of opponent strategies, as long as focus remains on differential Nash equilibria.

⁷And their combination, as discussed in the policy subsection.

This paper remains agnostic about the specificities of the critic estimation part of the algorithms. The goal is to gain insights about what can be learned as long as this function approximation step is reasonably well behaved, a property to be defined below. Note that well-behavedness in the critic estimation does not imply convergence of the algorithms to a Nash equilibrium, or convergence to some specific ideal strategy.

For each i , policies $\rho_{i,t}$ update according to

$$\rho_{i,t+1} = \rho_{i,t} + \alpha_t^i [F_S^i(\rho_t) + d_{t+1}^i + M_{t+1}^i], \quad (3)$$

where $\alpha_t^i > 0$ is a sequence of stepsizes converging to $\underline{\alpha} \geq 0$, d_{t+1}^i is a bias term converging to zero, and M_{t+1}^i is an error term of bounded variance. Bias and error term represent the estimation error involved in the estimation of F_S^i . For some estimator $\hat{F}_{S,t}^i$, one can write $d_{t+1}^i + M_{t+1}^i = \hat{F}_{S,t}^i(\rho_t) - F_S^i(\rho_t)$. I assume throughout that stepsizes of all agents lie within an order of magnitude of each other. Detailed sufficient conditions on stepsizes, bias, and error terms are relegated to Appendix A. The conditions ensure that (3) can be interpreted as a Robbins-Monro scheme (Robbins and Monro 1951), to which an extensive machinery for asymptotic results has been developed (c.f. Borkar 2009).

3. LONG RUN BEHAVIOR: MAIN RESULTS

This section presents the main results regarding characterisation of long run behavior of the algorithms. For a set A , let $cl(A)$ be its closure.

Definition 2. *Take the algorithm defined in (3). The limit set is defined as*

$$L_S = \bigcap_{t \geq 0} cl(\{\rho_\ell \mid \ell \geq t\}),$$

the set of limits of convergent subsequences ρ_{t_k} .

I write S as subscript to underline the dependence of the limiting set on the state variable S . As the characterizations introduced here will require properties of a differential equation, I present next some useful definitions:

Definition 3. Given some ODE $\dot{\rho} = f(\rho)$, let ρ^* be a rest point of $f(\rho)$. Let $\Lambda = \text{eig}[Df(\rho^*)]$ the set of eigenvalues of the linearization of f at ρ^* . For a complex number z , let $\mathbf{Re}[z] \in \mathbb{R}$ be the real part. ρ^* is

- Asymptotically stable if $\mathbf{Re}[\lambda] < 0$ holds for all $\lambda \in \Lambda$.
- Linearly unstable if $\mathbf{Re}[\lambda] > 0$ holds for at least one $\lambda \in \Lambda$.

One can think of $\mathbf{Re}[\lambda] < 0$ as a contraction property of the dynamical system around the rest point. Asymptotically stable rest points are *attractors* of the ODE. In other words, if the dynamical system were to start close to such a rest point, it will converge to it. On the other side, linearly unstable rest points don't come with a contraction property. There is at least some repelling direction of the ODE around the rest point.

To save notation, write $F_S(\rho)$ as stacked version of critic terms $F_S^i(\rho)$ for $\rho \in \bar{\mathbf{Y}}$.

Theorem 1. Let $\rho^* \in E_S$ be asymptotically stable for F_S . Then

$$\mathbb{P}[L_S = \{\rho^*\}] > 0.$$

Proof Sketch of Theorem 1

The full proof for this and the following Theorems can be found in the online appendix.

Firstly, I make a connection between the recursion in (3) and the differential equation induced by F_S . One can relate a time-interpolated version of the recursion ρ_t to solutions to the ordinary differential equation

$$\dot{\rho} = AF_S(\rho(t)),$$

where A is a diagonal matrix of strictly positive weights, representing the limiting relative stepsizes of different algorithms i . When considering that the updating rates α_t^i converge to zero, one may convince oneself that the recursion looks similar to a discrete approximation to a time-derivative. The idea is to show that the time-interpolated version of ρ_t must stay close, almost surely, to solutions of $AF_S(\rho)$. Attracting points of the differential system are then natural candidates to also attract ρ_t .

On the other hand, learning to play unstable rest points is an issue:

Theorem 2. *Let $\rho^* \in E_S$ be linearly unstable for AF_S . Then there exists an open neighborhood U of ρ^* such that*

$$\mathbb{P}[L_S \in U] = 0.$$

Proof Sketch of Theorem 2

ρ^* being unstable implies that there exists an unstable manifold that ρ^* lies on, which acts as a repeller to the differential equation based on F_S . I go on to show that due to the instability of ρ^* and nonvanishing variance of noise term M_{t+1} , no matter how close the algorithmic process gets to ρ^* , and no matter how large t is, there is always a nonzero probability that ρ_t lands on the unstable manifold and therefore must move away from ρ^* .

Hence, asymptotically stable equilibria are equilibria that can be limiting points of the RL learning procedure, while unstable equilibria are not. The intuition is related to how RL learn to play: since such agents make errors due to estimation and also to explore their action space, opponent’s strategy profiles are constantly perturbed. In other words, out of the view of a fixed agent i , the other agents are frequently deviating to policies nearby in the policy space. Now suppose the current profile ρ_t is close to an equilibrium ρ^* . Since i ’s updating rule tracks F_S , their policy will only stay close to ρ^* if the dynamics of F_S are somehow robust to deviations. This robustness is implied by asymptotic stability, and broken by unstable equilibria.

There is a caveat here, however: Theorem 1 does not state that all limiting points in L_S will be equilibria of the underlying repeated game as played by rational players. Depending on details of the stage game and state variable, one may or may not be able to rule out the case where algorithm updates get trapped in a cycle, or other more complex behavior not involving rest points (see Papadimitriou and Piliouras (2018)). I do not include cycles in the above definition, however it is straightforward to extend Theorem 1 to the case of attracting cycles as in Faure and Roth (2010), and there exist results considering linearly unstable cycles (Michel Benaïm and Faure (2012)) that suggest one may extend Theorem 2 to such linearly unstable cycles also.⁸ Notice that this observation implies that the Folk

⁸The inclusion of an analysis of limit cycles is an interesting avenue of further research, but would be beyond the scope of this paper.

theorem is neither necessary nor sufficient in describing the possible payoffs achievable by learning algorithms.

4. BERTRAND APPLICATION

This section introduces a differentiated goods Bertrand model (Bertrand model for short) to apply the more general results introduced previously. Throughout, I will say a given profile ρ is *learnable*, or *may be learned*, if there is positive probability that algorithms will converge to it; and *will not be learned* if algorithms will converge to it with zero probability.

There are two firms, $i \in \{1, 2\}$. Firms choose prices $p_i \in \mathbf{Y} \subseteq \mathbb{R}_+$. After prices are posted, consumers choose to consume, and at which firm with some idiosyncratic noise; this leads to a random aggregate demand shock observable by all. Let $\tilde{A} \in \Omega \subset \mathbb{R}_+$ be the random variable representing the aggregate demand shock. Assume Ω is compact. Conditional on the price vector p , this random variable has an absolutely continuous distribution with a density $g(a; p)$. Demand for each individual firm is therefore a random variable, denoted \tilde{X}_i . \tilde{X}_i depends on other firm's actions only through \tilde{A} . The expectation of \tilde{X}_i given p is then written as $X_i(p) = \mathbb{E}[\tilde{X}_i \mid p]$.

Expected profits given p are then given by $\pi_i(p) = (p_i - c)X_i(p)$. For ease of exposition, we assume that $\pi_i(p)$ are symmetric.

Assumption 3. For all $p \in \mathbb{R}_+^2$, $p_{-i} \geq 0$:

- (1) $\frac{\partial}{\partial p_i} X_i(p) \leq 0$, $\frac{\partial}{\partial p_{-i}} X_i(p) \geq 0$.
- (2) $\frac{\partial}{\partial p_i} X_i(0, p_{-i}) < 0$.
- (3) $\pi_i(p)$ is quasiconcave in p_i .
- (4) $\left| \frac{\partial}{\partial p_i} X_i(p) \right| \geq \left| \frac{\partial}{\partial p_{-i}} X_i(p) \right|$
- (5) $\left| \frac{\partial^2}{(\partial p_i)^2} X_i(p) \right| \geq \left| \frac{\partial^2}{\partial p_i \partial p_{-i}} X_i(p) \right|$.
- (6) For all $a \in \Omega$, $g(a; p)$ is twice differentiable in p .

(1) is a natural assumption in the Bertrand game. (2) ensures that positive prices will optimally be played. (3) implies that first order conditions are sufficient for best responses in the stage game. (4) and (5) ensure that effects of own price decisions dominate opponent's

decision's impact on own demand. (6) is a regularity condition that will become useful for the results on stability of equilibria.

Example 1. Consider the linear demand differentiated Bertrand model, with $X_i(p) = A - bp_i + \gamma p_{-i}$, $b > \gamma > 0$. If one takes \tilde{A} to feature linearly in the individual demand \tilde{X}_i , and has $\mathbb{E}[\tilde{A} \mid p] = Z + p_1 + p_2$ for some $Z \geq 0$, this model accomodates our motivation above, and satisfies Assumption 3.

Example 2. A slight extension to the logit-demand model also accomodates this setting. Suppose a mass of consumers $Z > 0$, upon observing a price index $P = p_1 + p_2$, stochastically decide to participate in the market. The expected mass of consumers given P is $M(P) = Z - P \geq 0$. Conditional on participation, the standard logit-demand model is carried out. Then we have $X_i(p) = M(P) \exp(\mu_i - \beta_i p_i) / \left(\sum_{j=1}^2 \exp(\mu_j - \beta_j p_j) \right)$, where $\mu_i, \beta_i > 0$.⁹

Under the trivial state variable S_0 which takes only one value (i.e. only stage game learning is possible), refer to $F_{S_0,B}, F_{S_0,G}$ as static best response dynamics (where payoffs are stage game payoffs, not repeated game payoffs as in payoff function W), and static gradient dynamics, respectively.

Lemma 1. Under Assumption 3, there is a unique interior Nash equilibrium p_N , which is symmetric. Given state S_0 , this Nash equilibrium will be learned with probability 1 by ACQ learners, and with positive probability by gradient learners.

Hence, if algorithms were not able to learn based on state variables correlated to past actions, learning to play Nash is a likely outcome. The introduction of a state variable can be seen as the introduction of a monitoring technology, as policies can now condition on random variables correlated to each agent's past actions. This interpretation is useful also in that the results in this section directly relate to the effectiveness of the monitoring technology available to the algorithms.

⁹Assumption 3 needs to be weakened for this example. The remaining claims in this section carry over here as well for bounded p_i , and large enough β_i , which ensures that signs in the assumption aren't flipped for a large enough set of prices.

The results in this section are stated and proved under the assumption that $\frac{\alpha_t^1}{\alpha_t^2} \rightarrow 1$ as $t \rightarrow \infty$, i.e. stepsizes of the two agents are asymptotically equal. The more general setting is discussed later on.

4.1. Learning Nash. Define for any $s, s' \in \mathbf{S}$, and $p \geq 0$:

$$T_{ss'}(p_1, p_2) = \mathbb{P}[s' | s; p_1, p_2],$$

the transition probability of moving to s' given current state s and quantity choices p_i in state s . Also assume that $T_{ss'}(p_1, p_2) = \mathbb{P}[s' | s; p_1 + p_2]$ for all s, p_i , i.e. transition probabilities only depend on aggregate prices, what one can think of as a price index level. I will therefore commonly write $T_{ss'}(p_1, p_2) = T_{ss'}(P)$ with $P = p_1 + p_2$.

Lemma 1 implies that there exist trivial state variables (taking a constant value) under which static Nash can be expected to be learned by ACQ learners and gradient learners. I show next that even though that is true, a larger family of state variables exist so that when they are used, ACQ learners will not converge to this Nash equilibrium.

For the remainder of this section, if not further specified, results are stated regarding ACQ learners ($F_{S,B}$) only. To ease intuitions, consider a state variable S the transitions of which depend on realizations of the aggregate shock \tilde{A} . Since \tilde{A} is the only payoff-relevant variable correlated to all agent's choices, we call such state variables PR-states, or just PR (payoff-relevant). For such state variables, transitions are pinned down by a deterministic mapping $f_S(s, \tilde{A}) \in \mathbf{S}$ for all $s \in \mathbf{S}$. Define own-and-opponent's price elasticity of demand as

$$\xi_o(p) = \left(\frac{\partial}{\partial p_1} X_1(p) \right) \frac{p_1}{X_1(p)}, \quad \xi_c(p) = \left(\frac{\partial}{\partial p_2} X_1(p) \right) \frac{p_2}{X_1(p)},$$

and let $L(p) = \frac{p-c}{c}$, with $L_N = L(p_N)$ being the Lerner index at Nash. Finally, define the growth rate of the Lerner index as $G_N = \left(\frac{\partial}{\partial p_1} L(p) \right) \Big|_{p=p_N} / L_N$. Write

$$d\xi_o(p) = \frac{\partial}{\partial p_1} \xi_o(p) + \frac{\partial}{\partial p_2} \xi_o(p).$$

$d\xi_o(p) < 0$ implies that price hikes by both agents lead to an overall increase in magnitude of elasticity for agent 1, as $\xi_o(p) \leq 0$ by Assumption 3. Call such markets *balanced*.¹⁰ Finally, note that π_i depends on conditional density $g(a; p)$ only through expected demand $X_i(p)$. Say two densities g, g' induce π if both densities lead to the same expected demand function.

Theorem 3. *Consider a PR-state variable S_{PR} with $|\mathbf{S}_{PR}| = K \geq 1$ states, and let $\rho_N : \mathbf{S}_{PR} \rightarrow \mathbb{R}_+$ be the policy that plays p_N in every state. Given π , there exists $0 < C < \infty$ such that for any g inducing π ,*

- (i) ρ_N will be learned with positive probability if $\sup_{a \in \Omega} \left| \frac{\partial}{\partial P} g(a; P) \right|_{P=P_N} < C$.
- (ii) ρ_N will not be learned if $\inf_{a \in \Omega} \left| \frac{\partial}{\partial P} g(a; P) \right|_{P=P_N} > C$.
- (iii) C is proportional to

$$\frac{1}{\xi_c(p_N)} (|\xi_o(p_N)| G_N - d\xi_o(p_N)).$$

Whether ρ_N can be approached by learning algorithms comes down to sensitivity of transition probabilities to deviations in prices from p_N , in conjunction with market properties. Crucially, note that derivatives of $g(a; P)$ that are small in magnitude indicate that deviations in price would be harder to detect, i.e. result in an ineffective monitoring technology. Point (iii) indicates how such market properties may lead to an increase in the set of transition probabilities that may allow for the learning of ρ_N . Overall, this point indicates an interesting dichotomy: on the one hand, weaker competition (low ξ_c , balanced market) facilitates learning, while on the other hand, higher own-price elasticity $|\xi_o|$ and more sensitive Lerner indices (G_N large) also facilitate learning. Note that this results holds for *any* PR-state variable.¹¹

The conclusion is concerning, as achieving favorable market conditions for the learning of ρ_N may be difficult, weak competition being commonly associated also with low own-price elasticity. An example of markets favorable to the learning of ρ_N according to point (iii) would be luxury goods markets with high brand recognition. Strong branding leads to weak

¹⁰Examples 1, 2 satisfy this when increases in own price have large enough effects on own demand relative to increases in opponent price.

¹¹The bound readily generalizes to any state variable with differentiable transition function by replacing the bound in (i), (ii) by one that depends on transition probabilities directly.

competition, while luxury goods may be avoided by some consumers if prices are too high, which constitutes high ξ_o .

In terms of intuitions, recall that attractiveness of ρ_N with respect to F_S determines whether this equilibrium may be learned by algorithms. Attractiveness is a property of robustness to perturbations in the policy space; being able to detect likely deviations (i.e. having an effective monitoring technology) is necessary for perturbations to have any bite at all. Furthermore, the properties in (iii) all come down to robustness should deviations be detected, by firstly ensuring that opponent's deviations don't have too large of an effect (see 'small enough' ξ_c , 'large enough' $d\xi_o$), while secondly own deviations must readily correct for perturbations, requiring a 'large enough' $|\xi_o|$.

It is clear that state variables of arbitrary $K > 1$ have the potential to support complex collusive schemes featuring a variety of different prices over time. The bounds to producer and consumer surplus due to such equilibria are characterised by the profit-maximizing collusive scheme. This, in turn, can be pinned down by simple binary strategies, as is known from Abreu, Pearce, and Stacchetti (1990), henceforth APS. Under additional conditions, binary schemes will be the only way to support optimal collusive schemes. The stability of such optimal collusive schemes then serves as a bound to the extend of possible collusion among algorithms.

4.2. Relationship to the best Equilibrium. APS provide a result stating that the best strongly symmetric sequential equilibrium (SSE) of the repeated game can be supported by a bang-bang solution, under their setting.

Such a bang-bang strategy, by definition, is constructed using subsets of Ω , intended as punishment and reward regions. Translated into this paper's setup, there exists a state variable S^* with $\mathbf{S}^* = \{A, B\}$ and sets Ω_A, Ω_B such that $s = A$ and $\tilde{A} \in \Omega_A$ implies the next period's state is A , and $s = B$ and $\tilde{A} \in \Omega_B$ implies next period's state is B . The reverse holds for $\Omega \setminus \Omega_s$.

Notice that any binary partition of Ω affects payoffs of players only by pinning down their transition probability functions $T_{ss'}(P)$. As $g(a; p)$ is twice differentiable in p by Assumption 3, we can restrict attention to twice continuously differentiable $T_{ss'}$. Call the space of such

transition probability functions \mathcal{T} . For any $T_{ss'} \in \mathcal{T}$, let $E^*(T_{ss'})$ be the set of symmetric Nash equilibria given expected discounted payoffs $W(\sigma)$ when $T_{ss'}$ governs state transitions. This set is nonempty for all $T_{ss'} \in \mathcal{T}$ due to the repetition of the static Nash equilibrium p_N .

As APS' result was shown only for finite strategy sets, I introduce an approximation result to the continuous action case. Let \mathbf{Y}_L be a discretization of cardinality $0 < L < \infty$ of price-set \mathbf{Y} , where for any discretization I impose that $p_N \in \mathbf{Y}_L$.

Define $W(\rho, T)$ for symmetric profiles $\rho \in \bar{\mathbf{Y}}$, transition probabilities $T \in \mathcal{T}$ as long run expected payoffs. Define $E_L(T)$ as the set of symmetric equilibria given discretization \mathbf{Y}_L . Here, APS's bang-bang result directly applies. By the observation above, we can alternatively characterise the maximal SSE as

$$V_L = \sup_{\substack{\rho \in E_L(T) \\ T \in \mathcal{T}}} W(\rho, T).$$

Analogously, define

$$V = \sup_{\substack{\rho \in E^*(T) \\ T \in \mathcal{T}}} W(\rho, T). \quad (4)$$

Define V^* to be the best SSE payoff among all SSE of Γ^∞ .

Proposition 1. *Given additional regularity conditions on W , \mathcal{T} ¹², there exists an SSE ρ of Γ^∞ supported by a binary-state policy, under some $T^* \in \mathcal{T}$ such that $V = W(\sigma, T^*)$. It holds that*

$$(1) \ V \leq V^*.$$

$$(2) \ \text{For any } \varepsilon \text{ there exists } \bar{L} \text{ such that for all } L \geq \bar{L}, |V - V_L| < \varepsilon.$$

Proposition 1 tells us that there exist binary state variables such that if used by algorithms, they may learn to play strategies that achieve the best SSE payoff for any discretization of their game. Whether such strategies may be learned comes down to sensitivity of transition probabilities in a similar fashion as in Theorem 3. To state this formally, a few more variables require defining.

¹²These stronger conditions ensure continuity of the equilibrium correspondence, see Assumption 6

Let $T_{ss'}^*(P)$ be the transition probability function supporting an SSE that achieves V . Call the payoff-maximizing equilibrium policy based on this state variable ρ^* , then $\rho^*(A) = p_A \geq p_N \geq p_B = \rho^*(B)$. Call a state variable *generic* if its associated transition matrix $T_{ss'}(\rho^*(s))$ is generic in the space of matrices.

Given PR-state variable S with $K > 1$ states, say policy $\bar{\rho} : \mathbf{S} \rightarrow \mathbf{Y}$ supports ρ^* if all the prices played under $\bar{\rho}$ equal either p_A or p_B . To save notation, we also say $\bar{\rho}$ is generic if there is a generic state variable generating $\bar{\rho}$ supporting ρ^* .

Theorem 4. *There exist $0 < C_{\pi,1}, C_{\pi,2} < \infty$, such that for any generic $\bar{\rho}$ supporting ρ^* ,*

(i) *$\bar{\rho}$ will be learned with positive probability if*

$$\max \left\{ \max_{s \in \{A,B\}} \sup_{a \in \Omega} \left| \frac{\partial}{\partial P} g(a; P) \Big|_{P=P_s} \right|, \max_{s \in \{A,B\}} \sup_{a \in \Omega} \left| \frac{\partial^2}{(\partial P)^2} g(a; P) \Big|_{P=P_s} \right| \right\} < C_{\pi,1}.$$

(ii) *$\bar{\rho}$ will not be learned if*

$$\min_{s \in \{A,B\}} \inf_{a \in \Omega} \left| \frac{\partial}{\partial P} g(a; P) \Big|_{P=P_s} \right| > C_{\pi,2}.$$

Hence, a similar insight to Theorem 3 applies here as well: If state transitions (i.e. derivatives of g in the case of PR-states) are insensitive enough to deviations in prices, then collusive equilibria such as ρ^* may be learned by algorithms. Note however that less sensitive transitions also affect the equilibrium set E_S , associated equilibrium payoffs, and consumer welfare down the line. This is discussed in the next section.

4.3. Policy Options.

4.3.1. *Covariate Restriction.* It follows from Theorem 3 that if one were able to affect the distribution over observed states, one may be able to ensure learnable ρ_N . This channel represents a feasible, and realistic policy instrument. Indeed, restrictions to the inputs (i.e. state variables) of algorithms have been implemented in the United States after successful lawsuits (see e.g. the Supreme Court decision in *Students for Fair Admissions, Inc. v. President and Fellows of Harvard College*).

It is important to note, however, that according to Theorem 4, de-sensitizing the density g of aggregate shock \tilde{A} may also introduce the possibility of collusive outcomes being

learned. Nevertheless, less sensitive g can be interpreted as less accurate monitoring technology available to players; this in turn can only lead to (weakly) less concerning collusive possibilities.

Formally, suppose initially that algorithms follow some PR-state variable given commonly observable \tilde{A} . I model restrictions to state variables of the algorithm as a garbling of observable \tilde{A} . Introduce $\beta \in (0, 1]$, and let $\tilde{U} \sim U([0, 1])$ be a uniformly distributed random variable, independent of \tilde{A} . Algorithms can only condition actions on $\tilde{C} \sim (1 - \beta) \circ \tilde{A} + \beta \circ \tilde{U}$, i.e. \tilde{C} is distributed as a convex combination of \tilde{A} and \tilde{U} . The support of \tilde{C} is still compact and positive, just as for \tilde{A} , so that all previous results go through for PR-state variables S given commonly observed \tilde{C} . For $\beta = 0$, we recover the unrestricted case, and $\beta = 1$ leads to E_S consisting of only stage game Nash equilibria as the state would carry no information about the past. With some abuse of notation, let $V(\beta)$ be the sup bang-bang payoff as defined in (4), when \tilde{C} given β is used as common observable.

Proposition 2. *Take any stage game and \tilde{A} satisfying Assumption 3. When PR-state variables given \tilde{C} are used by algorithms,*

- (i) *there exists $\beta \in [0, 1)$ s.t. ρ_N will be learned with positive probability,*
- (ii) *$V(\beta)$ weakly decreases in β .*

4.3.2. Upstream Competition. Here I consider an extension where the two agents may use algorithms that differ either via their stepsizes α_t^i , or their critic function F_S . This analysis can be interpreted as one of competition among algorithmic software providers selling learning algorithms to firms; in such a case it is more likely to find asymmetric learning rates, or critic functions.

Recall that I assume throughout that stepsizes of all agents lie within an order of magnitude from each other (see Assumption 5 in the appendix). It turns out that, in the case of symmetric equilibria, all insights stated for equal stepsizes carry over to the more general case. Suppose that $\frac{\alpha_t^1}{\alpha_t^2} \rightarrow \bar{\alpha}$ as $t \rightarrow \infty$. So far in this section, results have been stated under

$\bar{\alpha} = 1$. Also define the system asymmetric in critic functions

$$\dot{\rho} = F_{A,S} \equiv \begin{bmatrix} B_S^1(\rho_2) - \rho_1 \\ W_{1,S}^2(\rho_2, \rho_1) \end{bmatrix},$$

where $W_{1,S}^2$ refers to player 2's gradient vector in ρ_2 , for long-run payoffs given the same state variable S as used by player 1.

Proposition 3. *Take any state variable S , and any symmetric $\rho^* \in E_S$. Consider best response dynamics $F_{S,B}$. Under a regularity condition on $W(\rho^*)$,¹³*

- (i) ρ^* will be approached with positive probability given $\bar{\alpha} \in (0, 1)$ if and only if this is true given $\bar{\alpha} = 1$.*
- (ii) When $\rho^* = \rho_N$, (i) holds also when replacing $F_{B,S}$ with gradient dynamics $F_{G,S}$.*
- (iii) When $\rho^* = \rho_N$, ρ^* will be approached with positive probability given $F_{A,S}$ if and only if ρ^* will be approached with positive probability given $\bar{\alpha} = 1$ under $F_{B,S}$.*

This result indicates that upstream competition among software providers, selling algorithmic pricing solutions to firms, may not affect likely learning outcomes of those firms, as long as stepsizes remain within an order of magnitude of each other and firms' incentives are sufficiently symmetric. When restricting attention to the benchmark ρ_N , the result extends to gradient dynamics, and also to the case of asymmetric critics $F_{A,S}$.

If it were true that, e.g., $\bar{\alpha} = 0$, the limiting behavior of algorithms could be quite different. The machinery applied in this paper extends to this case, with the important change that the interaction among algorithms can be interpreted as sequential: if 2's updates are an order of magnitude faster than 1's updates to their policy, in the limit, the behavior observed will be as if 1 commits to a (Markov-) policy, 2 observes this, and best responds (in the case of ACQ learning). Clearly, the equilibrium set of such a game would be inherently different. Whether such a setting would lead to less collusive outcomes is beyond the scope of this paper.

4.4. Extension to Market Dynamics. Practical applications of algorithmic pricing commonly involve the need to adjust prices to fluctuating market conditions. Examples include

¹³All eigenvalues of the jacobian of the best response function at ρ^* are real.

gasoline retail pricing, which involves changing demand over day and night, and weekdays to weekend as discussed in the introduction (Assad et al. (2024)). To accomodate this setting, consider the following extension:

There is a random variable S_D representing aggregate market conditions $s_D \in \mathbf{S}_D$, which takes finitely many values $K_D \geq 1$. S_D evolves as an irreducible Markov chain with transition matrix T_D . At every period t , before price choices are made, $s_D \in \mathbf{S}_D$ is revealed to all firms. As before, $\tilde{A} \in \Omega$ is a random variable representing market conditions that are affected by current price choices, and is revealed at the end of each period. Assume that \tilde{A} is independent of S_D .¹⁴ Stage game payoffs are then written as conditional on S_D , in expectation over \tilde{A} . Define $X_i(p, s_D)$ as the expected demand given s_D and price vector p , where expectation is taken over \tilde{A} .

$$\pi_i(p, s_D) = X_i(p, s_D)(p_i - c).$$

We can extend Assumption 3 to this setting in a straightforward manner:

Assumption 4. For all $s_D \in \mathbf{S}_D$, all $p \in \mathbb{R}_+^2$, $p_{-i} \geq 0$:

- (1) $\frac{\partial}{\partial p_i} X_i(p, s_D) \leq 0$, $\frac{\partial}{\partial p_{-i}} X_i(p, s_D) \geq 0$.
- (2) $\frac{\partial}{\partial p_i} X_i(0, p_{-i}, s_D) < 0$.
- (3) $\pi_i(p, s_D)$ is quasiconcave in p_i .
- (4) $\left| \frac{\partial}{\partial p_i} X_i(p, s_D) \right| \geq \left| \frac{\partial}{\partial p_{-i}} X_i(p, s_D) \right|$
- (5) $\left| \frac{\partial^2}{(\partial p_i)^2} X_i(p, s_D) \right| \geq \left| \frac{\partial^2}{\partial p_i \partial p_{-i}} X_i(p, s_D) \right|$.
- (6) For all $a \in \Omega$, $g(a; p)$ is twice differentiable in p .

An argument analogous to that of Lemma 1 gives that for all $s_D \in \mathbf{S}_D$, there is a unique symmetric (static) Nash equilibrium $p_N(s_D)$. Given some K -sized PR-state variable \mathbf{S}_{PR} , policies of algorithms would map as $\rho : \mathbf{S}_{PR} \times \mathbf{S}_D \rightarrow \mathbb{R}_+$. As before, let $\rho_N \in \mathbb{R}_+^{K_D K}$ be the repetition of the static Nash equilibria, so that $\rho_N(s, s_D) = p_N(s_D)$ for all $s \in \mathbf{S}_{PR}$, $s_D \in \mathbf{S}_D$.

¹⁴This can be weakened to \tilde{A} having a distribution that depends on the realization s_D , while not affecting transitions of S_D in the future.

Analogously to Theorem 3 we define here

$$\begin{aligned}\gamma_1^m &= X^m \left(\xi_{o,1}^m L_N^m + \xi_o^m \frac{\partial}{\partial p_1} L_N^m \right) \\ \gamma_2^m &= X^m \xi_{o,2}^m L_N^m \\ \gamma_3^m &= X^m \xi_c^m L_N^m,\end{aligned}$$

where superscript m denotes payoff terms given state s_D^m , all evaluated at $p_N(s_D^m)$, which is dropped for ease of notation. Then let

$$\lambda_1 = \max_m \left| \frac{\gamma_2^m}{\gamma_1^m} \right|; \quad \lambda_2 = \max_m \left| \frac{1}{\gamma_1^m} \right|; \quad \lambda_3 = \max_m |\gamma_3^m|,$$

where notably $-\frac{\gamma_2^m}{\gamma_1^m}$ is the slope of the static best response function under s_d^m , evaluated at $p_N(s_D^M)$. Hence, under Assumption 4 we have $\lambda_1 < 1$.

Theorem 5. *For any S_D state variable satisfying this setting, given π , there exists $0 < C < \infty$ such that for any g inducing π , ρ_N will be learned with positive probability if*

$$\sup_{a \in \Omega} \left| \frac{\partial}{\partial P} g(a; P) \Big|_{P=P_N} \right| < C.$$

Furthermore, given any S_{PR} state variable,

- (i) *there exist λ_1 and $\lambda_2 \lambda_3$ small enough so that ρ_N will be learned with positive probability.*
- (ii) *there exist m_1, m_2, m_3 and $|\gamma_1^{m_1}|, |\gamma_2^{m_2}|, |\gamma_3^{m_3}|$ large enough so that ρ_N will not be learned.*

Hence, intuitions extend from the $K_D = 1$ setting discussed in Theorem 3 to this more general setting: whether competitive outcomes can be learned depends on the sensitivity of the monitoring technology and the magnitudes of elasticities, lerner indices and their growth rates.

5. CONCLUSION

This paper considers the long-run behavior of a class of RL algorithms and shows how it can be interpreted via the stability of repeated game equilibria according to an underlying differential equation. The application of collusion in repeated games is employed to show the usefulness of this framework.

An important insight from my analysis is the dependence of the attractability of a given equilibrium of the repeated game on state variables observed by algorithms, i.e. their implied monitoring technology. This insight, as discussed, may serve as a tool to curb algorithmic collusion.

Interesting future research directions include more detailed considerations of asymmetric learning settings, as touched upon in the discussion on upstream competition.

Furthermore, the characterization of long-run behaviors serves as a methodology that can allow for a variety of interesting economic application. The method enables researches to pick a given interaction of interest, e.g. an auction, a stock market, or multilateral platform, then pick a class of algorithms, and evaluate long-run outcomes in the chosen setting.

APPENDIX A. REDUCED FORM ALGORITHM ASSUMPTIONS

The following assumptions are sufficient for the results stated in Section 3 to go through, upon minor extensions to known results from stochastic approximation theory, to be found in Benaïm (1999), Borkar (2009), Michel Benaïm and Faure (2012). A thorough argument generalizing results further to the non-vanishing bias case can be found in the online appendix.

For notational ease, write $F(\rho) = F_S(\rho)$, as the stacking over i of $F_S^i(\rho)$. For all results to follow, state variables will be fixed. The algorithm (3) can be written as

$$\rho_{n+1} = \rho_n + \alpha_n [F(\rho_n) + \delta_n + M_{n+1}], \quad (5)$$

where δ_n is a vector of δ_n^i stacked over i . We switch to an identification of time periods by n in order to distinguish the continuous timescale t used in the associated continuous time systems.

Assumption 5. *Let \mathcal{F}_n be the σ -field generated by $\{\rho_n, \delta_n, M_n, \rho_{n-1}, \delta_{n-1}, M_{n-1} \dots, \rho_0, \delta_0, M_0\}$, i.e. all the information available to the updating rule at a given period n .*

- (i) *Stepsizes α_n^i satisfy, for all i , to be square-summable, but not summable.*
- (ii) *For all i, j , $\lim_{n \rightarrow \infty} \frac{\alpha_n^i}{\alpha_n^j}$ exists and lies in (c, ∞) , for some $c > 0$.*

(iii) F is Lipschitz continuous and grows sublinearly, i.e.

$$\limsup_{\|\rho\| \rightarrow \infty} \frac{\|F(\rho)\|}{\|\rho\|} < \infty.$$

(iv) M_{n+1} is a Martingale-difference noise. There is $0 < \bar{M} < \infty$ and $q \geq 2$ such that for all n

$$\mathbb{E}[M_{n+1} | \mathcal{F}_n] = 0; \quad \mathbb{E}[\|M_{n+1}\|^q | \mathcal{F}_n] < \bar{M} \quad \mathcal{F}_0 - \text{almost surely.}$$

(v) There exists a continuous function

$$\Omega : \bar{\mathbf{Y}} \mapsto \mathcal{J}(\bar{\mathbf{Y}}),$$

where $\mathcal{J}(\bar{\mathbf{Y}})$ is the space of positive definite matrices given vectors in $\bar{\mathbf{Y}}$, such that for all n

$$\mathbb{E}[M_{n+1} M'_{n+1} | \mathcal{F}_n] = \Omega(\rho_n),$$

whenever $\rho_n \in \mathcal{U}$.

(vi)

$$E[\|\delta_n\|] = o(b_n),$$

where $b_n \rightarrow 0$ satisfies $\max_i \lim_{n \rightarrow \infty} \frac{\alpha_n^i}{b_n} = 0$, α_n^i being i 's stepsize.

(vii)

$$\sup_{n \geq 0} \mathbb{E}[\|\delta_n\|^2] < \infty,$$

(viii) For all $n' < n''$, $\delta_{n'}, \delta_{n''}$ are uncorrelated conditional on $\mathcal{F}_{n'}$.

(ix) Iterates stay bounded almost surely:

$$\sup_n \|\rho_n\| < \infty, \text{ a.s..}$$

Point (i) is known as the Robbins-Monro condition (Robbins and [Monro 1951](#)) on stepsizes. It ensures that stepsizes converge slowly enough so that the whole real line can be mapped (as a continuous-time interval), while converging not too slowly in order for error terms to be averaged out. (ii) ensures that all stepsizes lie within the same order of magnitude. Point (iii) ensures global integrability and uniqueness of solutions to $\dot{\rho} = F(\rho)$. In the example

of ACQ, it is an assumption on payoffs W^i , and that best responses can't grow too quickly. Point (iv) implies that given current information in period t , new errors due to $t + 1$'s estimator of F are well-behaved. It is a common assumption in stochastic approximation theory. Point (v) ensures that some variance in error terms remains for all n ; this is satisfied e.g. if the estimation of F involves exploratory noise, or stochasticity during the estimation as is true under randomized Bellman-iteration schemes. This assumption will be the main driver that pushes iterations away from unstable equilibria. Point (vi) ensures that the bias term vanishes faster than stepsizes. Points (vii), (viii) are further regularity conditions on the bias term. Even though commonly made, point (ix) is often difficult to verify. It is common for results to be stated conditioning on the event that (ix) holds, see for example Michel Benaïm and Faure (2012). For a more general discussion of sufficient conditions for bounded iterates, see Borkar (2009), Chapter 2.

APPENDIX B. PROOFS

First, note the following fact about block symmetric matrices.

Remark 1. *Suppose A, B are square matrices of the same dimension. Let*

$$T = \begin{bmatrix} A & B \\ B & A \end{bmatrix}.$$

Then one can show

$$\det(T) = \det(A - B)\det(A + B).$$

Given a square matrix A , define Λ as the set of eigenvalues of the A . Then define

$$\kappa(A) = \max\{|\lambda| : \lambda \in \Lambda\},$$

as the spectral radius of A . Define $J^i(\sigma^*)$ as the Jacobian of $B_S^i(\sigma^*)$, which is the matrix of best response derivatives of a given player. For symmetric σ^* , we drop the i superscript to save notation.

Lemma 2. *Suppose $\alpha^* = \beta^* = \sigma^* \in E_S$ is a differential, symmetric Nash equilibrium. Let $\bar{\kappa}$ be the real part of the spectral radius of $J(\sigma^*)$. Then σ^* is asymptotically stable under $F_{S,B}$ if $\bar{\kappa} < 1$, and unstable if $\bar{\kappa} > 1$.*

Proof. Using Remark 1, we get that

$$ch(\lambda) = \det(J(\sigma^*) - (1 + \lambda)I_2)\det(J(\sigma^*) + (1 + \lambda)I_2).$$

Thus, if μ is an eigenvalue of $J(\sigma^*)$, then $\pm|\mu-1|$ is an eigenvalue of $X(\sigma^*)$, and the conclusion follows, since asymptotic stability requires that all eigenvalues of $X(\sigma^*)$ have negative real parts. \square

Hence, it is enough to characterize the eigenvalues $\lambda_{1,2}$ of the matrix of best-response derivatives J of player 1 at symmetric Nash policies, which are the main focus of this paper.

The following Lemma will help with bounding eigenvalues of some Jacobians J involved in later proofs:

Lemma 3. *Consider a block-matrix M with*

$$M = \begin{bmatrix} M_1 & M_2 \\ M_3 & M_4 \end{bmatrix}$$

is a matrix consisting of blocks M_i . For m, k s.t. $m + k = n$, suppose we can write those blocks as

$$M_1 = A; \quad M_2 = B$$

$$M_3 = C; \quad M_4 = D,$$

A is a $m \times m$ matrix of strictly positive entries s.t, D is a $k \times k$ matrix, and C, B are $m \times k, k \times m$ matrices respectively, so that A, D are invertible. Let $E_1 = \max(\|A\|_\infty, \|D\|_\infty)$, $E_2 = \|B\|_\infty \|C\|_\infty$. Then, for all $\lambda \in \text{eig}(M)$:

$$|\lambda| \leq \frac{E_1}{2} + \sqrt{\left(\frac{E_1}{2}\right)^2 + E_2}.$$

Proof. Using the Schur complement and genericity of the blocks M_i , we can write the characteristic equation of M as

$$\text{char}(\lambda) = \det(D - \lambda I_k) \det(A - \lambda I_m - B((D - \lambda I_k)^{-1} C)).$$

Now, any eigenvalue $\lambda \in \text{eig}(M)$ with $\lambda \notin \text{eig}(D)$ must satisfy

$$\lambda \in \text{eig}(A - B(D - \lambda I_k)^{-1} C) \equiv \text{eig}(\bar{M}(\lambda)).$$

Letting $\rho(A)$ be the spectral radius of some matrix A , recall that for any (sub-multiplicative) matrix norm $\|\cdot\|$, we have $\rho(A) \leq \|A\|$. Consider $\|A\|_\infty = \max_i \sum_j |A_{ij}|$. Then,

$$\begin{aligned} |\lambda| &\leq \|A - B(D - \lambda I_k)^{-1} C\|_\infty \\ &\leq \|A\|_\infty + \|B\|_\infty \|C\|_\infty \|(D - \lambda I_k)^{-1}\|_\infty. \end{aligned}$$

Similarly, we have, as long as $\lambda \notin \text{eig}(A)$,

$$\begin{aligned} |\lambda| &\leq \|D - B(A - \lambda I_m)^{-1} C\|_\infty \\ &\leq \|D\|_\infty + \|B\|_\infty \|C\|_\infty \|(A - \lambda I_m)^{-1}\|_\infty. \end{aligned}$$

We need to accurately bound the inverse terms. As we seek an approach that relates small $|\lambda|$ with small enough $|\partial_{p_1} g(a; p)|$, write $D = hD$, for some $h \neq 0$. $\lambda \neq 0$ by genericity of M_i . Write $\lambda(I_k - \frac{h}{\lambda} D) = (-\lambda)I_k + D$. For $|h| < |\lambda|$ small enough, we will have $\|\frac{h}{\lambda} D\|_\infty < 1$. Then, an application of the von Neumann series approximation gives

$$\frac{1}{|\lambda|} \left\| \left(I_k + \frac{h}{\lambda} D \right)^{-1} \right\|_\infty \leq \sum_{k=0}^{\infty} \left\| \frac{h}{\lambda} D \right\|_\infty^k = \frac{1}{|\lambda|} \frac{1}{1 - \frac{h\|D\|_\infty}{|\lambda|}} = \frac{1}{|\lambda| - h\|D\|_\infty}.$$

An analogous bound can be constructed for $A - (\lambda)I_m$. By genericity, $\text{eig}(A) \cap \text{eig}(D) = \emptyset$. Then, for all $\lambda \in \text{eig}(M)$, when $|h| > 0$ small enough,

$$|\lambda| \leq \max(\|A\|_\infty, \|D\|_\infty) + \|B\|_\infty \|C\|_\infty \frac{1}{|\lambda|}.$$

Writing $E_1 = \max(\|A\|_\infty, \|D\|_\infty)$ and $E_2 = \|B\|_\infty\|C\|_\infty$, we can solve the above as an equality. This leads to a quadratic equation in $|\lambda|$, to larger root of which equals

$$\frac{E_1}{2} + \sqrt{\left(\frac{E_1}{2}\right)^2 + E_2}.$$

The result follows. □

Now consider the general situation in which $\mathbf{S} = \{s_1, \dots, s_K\}$. Let $\rho : \mathbf{S} \rightarrow \mathbb{R}_+$, $\gamma : \mathbf{S} \rightarrow \mathbb{R}_+$ be two policies using such states. For long term payoffs define recursively, for all $s_i \in \mathbf{S}$:

$$W(\rho, \gamma, s_i) = (1 - \delta)\pi(\rho(s_i), \gamma(s_i)) + \delta \sum_{k=1}^K T_{kk'}(\rho(s_i) + \gamma(s_i)) W(\rho, \gamma, s_k),$$

the discounted expected payoff from playing ρ if the opponent plays γ . Fixing the profile ρ, γ , we can also use the vector form

$$\begin{aligned} \tilde{W} &= [W(\rho, \gamma, s_1), \dots, W(\rho, \gamma, s_K)]^\top, \\ U &= [\pi(\rho(s_1), \gamma(s_1)), \dots, \pi(\rho(s_K), \gamma(s_K))]^\top, \end{aligned}$$

to write

$$W = (1 - \delta)U + \delta TW \iff W = [I_K - \delta T]^{-1} (1 - \delta)U,$$

where $T = (T_{kk'})_{k, k' \in \{1, \dots, K\}}$ is the Markov transition matrix given a fixed profile ρ, γ . Note that for all $\delta < 1$, $I_K - \delta T$ is an M -matrix. The inverse of $I_K - \delta T$ exists and has all elements non-negative. As a result, we then have that all rows of $[I_K - \delta T]^{-1}$ sum to $\frac{1}{1-\delta}$. In the remainder of this section, to save notation write $\rho_k = \rho(s_k)$, and similarly for γ . Counting arguments of W as the first K arguments referring to own strategy ρ , the next K arguments referring to γ , indicate derivatives and cross-derivatives of W with respect to a specific argument $1 \leq j \leq 2K$ using a subscript j . Then we have:

Corollary 1. *The derivatives of vector W can be written as, for $i \leq K < j$:*

$$\begin{aligned}
\tilde{W}_i &= [I_K - \delta T]^{-1} \delta \frac{\partial T}{\partial \rho_i} \tilde{W} + [I_K - \delta T]^{-1} (1 - \delta) \frac{\partial U}{\partial \rho_i} \\
\tilde{W}_j &= [I_K - \delta T]^{-1} \delta \frac{\partial T}{\partial \gamma_{j-K}} \tilde{W} + [I_K - \delta T]^{-1} (1 - \delta) \frac{\partial U}{\partial \gamma_{j-K}} \\
\tilde{W}_{ii} &= [I_K - \delta T]^{-1} \delta \frac{\partial^2 T}{(\partial \rho_i)^2} \tilde{W} + [I_K - \delta T]^{-1} (1 - \delta) \frac{\partial^2 U}{(\partial \rho_i)^2} - 2 [I_K - \delta T]^{-1} \delta \frac{\partial T}{\partial \rho_i} \tilde{W}_i \\
\tilde{W}_{ij} &= [I_K - \delta T]^{-1} \delta \frac{\partial^2 T}{\partial \rho_i \partial \gamma_{j-K}} \tilde{W} + [I_K - \delta T]^{-1} (1 - \delta) \frac{\partial^2 U}{\partial \rho_i \partial \gamma_{j-K}} \\
&\quad + [I_K - \delta T]^{-1} \delta \frac{\partial T}{\partial \gamma_{j-K}} \tilde{W}_i + [I_K - \delta T]^{-1} \delta \frac{\partial T}{\partial \rho_i} \tilde{W}_j.
\end{aligned}$$

Proof. This follows from some matrix algebra, importantly using the following fact:

For a matrix function X of variable y , let ∂X be the partial derivative of X with respect to y . Then $\partial(X^{-1}) = -(X^{-1})(\partial X)(X^{-1})$. \square

If $\tilde{W}_i = \mathbf{0}_K$, we can further simplify these:

$$\begin{aligned}
\tilde{W}_j &= [I_K - \delta T]^{-1} (1 - \delta) \left[\frac{\partial U}{\partial \gamma_{j-K}} - \frac{\partial U}{\partial \rho_{j-K}} \right] \\
\tilde{W}_{ii} &= [I_K - \delta T]^{-1} \delta \frac{\partial^2 T}{(\partial \rho_i)^2} \tilde{W} + [I_K - \delta T]^{-1} (1 - \delta) \frac{\partial^2 U}{(\partial \rho_i)^2} \\
\tilde{W}_{ij} &= [I_K - \delta T]^{-1} \delta \frac{\partial^2 T}{\partial \rho_i \partial \gamma_{j-K}} \tilde{W} + [I_K - \delta T]^{-1} (1 - \delta) \frac{\partial^2 U}{\partial \rho_i \partial \gamma_{j-K}} \\
&\quad + [I_K - \delta T]^{-1} \delta \frac{\partial T}{\partial \rho_i} \tilde{W}_j.
\end{aligned} \tag{6}$$

B.1. Proof of Theorem 3. To determine the stability of ρ_N , we need to compute the eigenvalues of the linearized best-response dynamics at ρ_N . Since by assumption, states are irreducible, we can fix an arbitrary initial state s_1 when computing first order conditions to pin down best responses. The implicit function theorem and (6) can then be used to find the gradient of the best response.

First, write $W_{ij} = \tilde{W}_{ij}(s_1)$ for all $i \leq K, j \leq 2K$. Then we can write the Hessian as $H = \text{diag}(W_{ii})$, the diagonal matrix having the second derivatives W_{ii} for $i = 1, \dots, K$ on

the diagonal. H is diagonal by the derivation of W_i : write

$$W_i = [I_K - \delta T]_1^{-1} \left[\delta \frac{\partial T}{\partial \rho_i} \tilde{W} + (1 - \delta) \frac{\partial U}{\partial \rho_i} \right],$$

where for any matrix A we write A_i as the i 'th row of A . Then taking another derivative with respect to a variable $j \leq K$:

$$W_{ij} = (\partial [I_K - \delta T]^{-1}) [I_K - \delta T] W_i + [I_K - \delta T]^{-1} \left[\delta \frac{\partial^2 T}{\partial \rho_i \partial \rho_j} \tilde{W} + (1 - \delta) \frac{\partial^2 U}{\partial \rho_i \partial \rho_j} \right]. \quad (7)$$

Notice that $\frac{\partial^2 T}{\partial \rho_i \partial \rho_j}$, $\frac{\partial^2 U}{\partial \rho_i \partial \rho_j}$ are matrices of all zeros if $i \neq j$ and $i, j \leq K$. Then plugging in that $W_i = 0$, we get that indeed $W_{ij} = 0$ whenever $i \neq j$ and $i, j \leq K$. So H must be diagonal.

Now define $R = [W_{ij}]_{i \leq K < j}$ as the matrix of cross derivatives between an agent's own strategy $\rho(s_i)$ and an opponent's strategy $\gamma(s_{j-K})$. Then we can define, using the implicit function theorem, the best response derivative matrix as

$$M = -H^{-1}R.$$

Since we evaluate this at ρ_N , we can make multiple observations that will greatly simplify the structure of M :

Firstly, long term payoffs $\tilde{W} = \mathbf{U}^N$, since ρ_N prescribes the same action in each state.

Secondly, by the nature of p_N , $\frac{\partial U}{\partial \rho_i} = 0$ for all $i \leq K$.

Now note that since by definition each row of T sums to one, and therefore each row of $\frac{\partial T}{\partial \rho_i}$ and $\frac{\partial^2 T}{\partial \rho_i \partial \rho_j}$ must sum to zero. Therefore, at ρ_N , we can simplify the elements of H, R to

$$\begin{aligned} W_{ii} &= [I_K - \delta T]_1^{-1} (1 - \delta) \frac{\partial^2 U}{(\partial \rho_i)^2}, \\ W_{ij} &= [I_K - \delta T]_1^{-1} (1 - \delta) \left[\frac{\partial^2 U}{\partial \rho_i \partial \gamma_{j-K}} + \delta \frac{\partial T}{\partial \rho_i} [I_K - \delta T]^{-1} \frac{\partial U}{\partial \gamma_{j-K}} \right], \end{aligned}$$

where for any matrix A we write A_i as the i 'th row of A . Let e_i be the K -vector that is one in entry i , and zero in all others. Using the fact that ρ_N is constant for all states, we

can write this down in the more simple form

$$\begin{aligned} W_{ii} &= [I_K - \delta T]_1^{-1} (1 - \delta) \pi_{11}^N e_i, \\ W_{ij} &= [I_K - \delta T]_1^{-1} (1 - \delta) \left[\pi_{12}^N e_i + \delta \frac{\partial T}{\partial \rho_i} [I_K - \delta T]^{-1} \pi_2^N e_i \right], \end{aligned}$$

if $i = j - K$, and

$$W_{ij} = [I_K - \delta T]_1^{-1} (1 - \delta) \left[\delta \frac{\partial T}{\partial \rho_i} [I_K - \delta T]^{-1} \pi_2^N e_{j-K} \right],$$

otherwise, following the notation used in section 4. To save notation, write $Z = [I_K - \delta T]^{-1}$. Then,

$$\frac{W_{ij}}{W_{ii}} = \begin{cases} \frac{1}{\pi_{11}^N} \left[\pi_{12}^N + \pi_2^N \delta \sum_{k=1}^K T'_{s_i s_k} Z_{k,i} \right] & \text{if } i = j - K \\ \frac{1}{\pi_{11}^N} \left[\pi_2^N \delta \sum_{k=1}^K T'_{s_i s_k} Z_{k,j-K} \right] & \text{o.w.} \end{cases}$$

For the proof of point (1), we will upper bound eigenvalues of this system. Note from the above that we can write

$$M = -\frac{\pi_{12}^N}{\pi_{11}^N} I_K + \frac{\pi_2^N}{\pi_{11}^N} B,$$

where terms in B depend on T, Z . If $\pi_2^N = 0$, the bound is trivial. Assume $\pi_2^N > 0$ (positivity follows from Assumption 3). M 's simple form implies that $v \in \mathbb{R}^K$ is an eigenvector of M if and only if it is an eigenvector of B . It follows that eigenvalues $\lambda \in \text{eig}(M)$ and $\mu \in \text{eig}(B)$ are related through the equation

$$\lambda = -\frac{\pi_{12}^N}{\pi_{11}^N} + \frac{\pi_2^N}{\pi_{11}^N} \mu.$$

As $-\frac{\pi_{12}^N}{\pi_{11}^N} \in (0, 1)$ by Assumption 3, it is sufficient to bound $\frac{\pi_2^N}{\pi_{11}^N} \mu$.

From this, we can derive that $|\lambda| < 1$ is equivalent to

$$\mu \in \left(\frac{\pi_{12}^N}{\pi_2^N} + \frac{\pi_{11}^N}{\pi_2^N}, \frac{\pi_{12}^N}{\pi_2^N} - \frac{\pi_{11}^N}{\pi_2^N} \right), \quad (8)$$

an interval that contains 0. Note that $B = \delta T' Z$, where T' is the transition-derivative matrix, where each row i corresponds to the derivative of row i of T with respect to ρ_i , all

evaluated at ρ_N . Increasing $|T'|$ in absolute value by a scalar c will increase $|\mu|$ by that scalar c . Generically in the space of transition probability matrices, there exists at least one eigenvalue $\mu \in \text{eig}(B)$ s.t. $\mu \neq 0$. Without loss, let $\mu > 0$. Then for any $c > 0$, $c\mu \in \text{eig}(cB)$. For any T' finite, cT' corresponds to another transition probability matrix, as summing to zero over rows is still satisfied, and the perturbation only needs to be carried out at the point ρ_N . One can directly construct a transition probability matrix D such that $D = T$ at ρ_N , and $D' = cT'$, by writing for ρ close enough to ρ_N : $D(\rho) = T + (\rho - \rho_N)cT'$. Hence, one can find a transition probability matrix with $|T'|$ small enough such that all eigenvalues μ will ensure $|\lambda| < 1$. Conversely, taking c large enough will lead to some eigenvalue of B large enough s.t. $|\lambda| > 1$ for some associated eigenvalue λ .

For point (1), interpret c as scaling down $\sup_{a \in \Omega} \left| \frac{\partial}{\partial P} g(a; P) \right|_{P=P_N}$, which in turn scales down $|T'|$. For point (2), let c scale up $\inf_{a \in \Omega} \left| \frac{\partial}{\partial P} g(a; P) \right|_{P=P_N}$. Finally, recall that we have $\pi_{12} \geq 0 > \pi_{11}$ under Assumption 3. Thus, $|\pi_{12}^N + \pi_{11}^N| < \pi_{12}^N - \pi_{11}^N$, so that the negative lower bound of the interval in (8) is closer to zero than the positive upper bound. This allows to conclude that for points (1) and (2), the cutoff C stated in the Theorem will be the same: for (1), all eigenvalues must satisfy $|\lambda| < 1$ - having all $|\mu|$ small enough to lie above the absolute value of the lower bound is sufficient. On the other hand, for instability, at least one λ must satisfy $|\lambda| > 1$. Finding one negative μ small enough is then enough, which allows to check the same lower bound. Note that the sign of μ can always be imposed by a perturbation using scalar c as indicated above, since ρ_N remains an equilibrium no matter the sign of transition probabilities.

For point (3), note first that we can write $\pi_1 = X(p) \left(1 + \xi_o(p) \frac{p_1 - c}{p_1} \right)$, so that at p_N , we have $-\xi_o(p_N) = \frac{1}{L_N}$, where $L_N = \frac{p_N - c}{p_N}$ is the Lerner index at p_N . Define

$$\xi_{o,1}(p) = \frac{\partial}{\partial p_1} \xi_o(p), \quad \xi_{o,2}(p) = \frac{\partial}{\partial p_2} \xi_o(p).$$

Then,

$$\begin{aligned}\pi_2 &= \xi_c(p)X(p)\frac{p_1 - c}{p_2} \\ \pi_{11} &= \xi_o(p)\frac{X(p)}{p_1}\pi_1 + \xi_{o,1}(p)\frac{p_1 - c}{p_1}X(p) + \xi_o(p)\frac{c}{p_1^2}X(p) \\ \pi_{12} &= \xi_c(p)\frac{X(p)}{p_2}\pi_1 + \xi_{o,2}(p)\frac{p_1 - c}{p_1}X(p).\end{aligned}$$

At p_N , these simplify to

$$\begin{aligned}\pi_2 &= \xi_c(p_N)X(p_N)L_N \\ \pi_{11} &= X(p_N)\left(\xi_{o,1}(p_N)L_N + \xi_o(p_N)\frac{\partial}{\partial p_1}L_N\right) \\ \pi_{12} &= X(p_N)\xi_{o,2}(p_N)L_N.\end{aligned}$$

Now consider the lower bound for (8) above. It is clear that the constant C identified for points (1),(2) is proportional to this lower bound. Using the above, and as $\pi_{12} + \pi_{11} \leq 0$ by Assumption 3, we can write the magnitude of the lower bound. Define $G_N = \frac{\partial}{\partial p_1}L_N/L_N$ as growth rate of the Lerner index at p_N . Then,

$$\begin{aligned}\frac{|\pi_{11}^N + \pi_{12}^N|}{\pi_2^N} &= \frac{-1}{\xi_c(p_N)}(\xi_o(p_N)G_N + \xi_{o,1}(p_N) + \xi_{o,2}(p_N)) \\ &= \frac{1}{\xi_c(p_N)}(|\xi_o(p_N)|G_N - d\xi_o(p_N)).\end{aligned}$$

The conclusion follows: the bound grows as ξ_c falls. G_N grows as c grows, $p_N - c$ falls, p_N falls. The bound grows as $|\xi_o(p_N)|$ grows. In a balanced market, $d\xi_o \leq 0$. Here the bound grows as $|d\xi_o|$ grows. In unbalanced markets, the bound grows as $|d\xi_o|$ falls. In general, as the market becomes ‘more balanced’, the bound grows. ■

B.2. Proof of Theorem 4. For point (i), recall that J refers to the $K \times K$ linearization of $F_S(\rho)$ at $\bar{\rho}$. Let $\mathbf{S}_A \subset \mathbf{S}$ refer to those states s for which $\bar{\rho}(s) = p_A$, and let $\mathbf{S}_B = \mathbf{S} \setminus \mathbf{S}_A$. Note that as $\bar{\rho}$ supports ρ^* , we may write J just as matrix M in the statement of Lemma 3, where k corresponds to $k = |\mathbf{S}_A|$, and $m = |\mathbf{S}_B|$. The eigenvalue bound of Lemma 3 then applies to J .

We can then follow a similar argument as in the proof of Theorem 3. As derived in (6), entries of J grow in magnitude firstly with transition probability terms, which are pinned down by the density derivatives $\left| \frac{\partial}{\partial P} g(a; P) \right|_{P=P_s}$, $\left| \frac{\partial^2}{(\partial P)^2} g(a; P) \right|_{P=P_s}$, $s \in \{A, B\}$. Secondly, terms of J grow with payoff terms $|\pi_i^s|, |\pi_{ij}^s|$ for $1 \leq i, j \leq 2$, $s \in \{A, B\}$. For any such payoff terms, we can push down density derivatives to a small enough value so that stability is accomplished. This follows, as in the extreme case where $g(a; p)$ does not vary with p , the best SSE satisfies $\rho^*(s) = p_N$ for all s trivially. Then, the bound collapses to the bound of Theorem 3.

As for point (ii), recall that for any matrix, the sum of eigenvalues equals the trace of the matrix. In this case, using (6) we can write

$$\text{tr}(J) = -K + \sum_{k=1}^K \frac{\pi_{11}^k - \pi_{12}^k + [I_K - \delta T]_1^{-1} \delta \frac{\partial T}{\partial \rho_k} e_k (\pi_1^k - \pi_2^k)}{\delta \frac{\partial^2 T}{(\partial \rho_k)^2} e_k W + \pi_{11}^k},$$

where π_i^k, π_{ij}^k refer to the derivatives of π evaluated at $\bar{\rho}_k$, and e_k is the K -dimensional column vector equal to 1 at element k , and zeros everywhere else. From this it is clear that $|\text{tr}(J)|$ can be made to grow by growing $|\frac{\partial T}{\partial \rho_k}|$. The result follows. \blacksquare

B.3. Proof of Lemma 1. Given Assumption 3, the Nash equilibrium must be interior. (3) ensures that best responses are unique, while (4) and (5) imply that the slope of the best response is less than one for all $p_{-i} \geq 0$. This implies that the Nash equilibrium p_N is unique, and symmetry of the payoffs implies that this equilibrium is symmetric. For static best response dynamics, it is well known that these conditions imply global attraction to the Nash equilibrium (Milgrom and Roberts (1990)). As for gradient dynamics, it is straightforward to show that p_N must be asymptotically stable: let M_G, M_B be the linearization of gradient dynamics and best response dynamics, respectively at p_N . Then by symmetry of payoffs, we can write $M_G = -\pi_{11}^N M_B$. Any eigenvalue $\lambda \in \text{eig}(M_G)$ is such that $-(\lambda/\pi_{11}^N) \in \text{eig}(M_B)$. As $-\pi_{11}^N > 0$, all eigenvalues of M_G are negative if and only if all of M_B are, so that stability carries over. \blacksquare

Proof of Proposition 1. For any interior $\sigma \in E^*(T_{ss'})$, let $J(\sigma)$ be the 2×2 matrix of best response derivatives, i.e. the Jacobian of the best response function at the equilibrium σ , $B_S^1(\sigma_2)$. We require the following additional assumption:

Assumption 6. *For all $T_{ss'} \in \mathcal{T}$, all $\sigma \in E^*(T_{ss'})$ are interior, with negative definite Hessian, and all eigenvalues of $J(\sigma)$ are different from 1.*

For any fixed $T_{ss'}$, the above assumption is a generic property over the space of regular payoff functions. The assumption has additional strength as it imposes that given the regular stage game payoff function π , there exists no $T_{ss'} \in \mathcal{T}$ that could lead to a singular Hessian at some equilibrium, or a $J(\sigma)$ with eigenvalue equal to 1. For any discretization \mathbf{Y}_L , define $W^L(\rho, T) : \mathbf{Y}^2 \times \mathcal{T} \rightarrow \mathbb{R}$ as restriction of the payoff function to \mathbf{Y}_L , given some T .

$$W^L(\rho, T) = W(f^L(\rho), T),$$

where

$$f^L(\rho) = \arg \min_{\rho' \in \mathbf{Y}_L^2} \|\rho - \rho'\|,$$

for any norm on \mathbf{Y}^2 , the projection of ρ onto discrete space \mathbf{Y}_L . For every sequence \mathbf{Y}_L there is an associated sequence $\alpha_L(T)$ with

$$\alpha_L(T) = \max_{\rho \in \mathbf{Y}^2} \|W^L(\rho, T) - W(\rho, T)\|.$$

Continuity of W and the Lipschitz property of density $g(a; p)$ implies that $\alpha_L(T) \rightarrow 0$ for all $T \in \mathcal{T}$. Write $\alpha_L(\mathbf{Y}_L, T)$ for a sequence of α_L given a fixed sequence of discretizations and transition function T . Say that a discretization sequence \mathbf{Y}_L is *covering* if $\alpha_L(\mathbf{Y}_L, T) \rightarrow 0$ (and $p_N \in \mathbf{Y}_L$). From now on fix a covering sequence of discretizations \mathbf{Y}_L and transition probability T .

Notice that $E_L(T)$ is closed-valued, trivially by finiteness of \mathbf{Y}_L . Furthermore, $E^*(T)$ is closed-valued: W is continuous, \mathbf{Y} compact, and thus Berge gives us that the best-response correspondence is closed and compact-valued. Then, applying the closed-graph theorem gives us that the equilibrium set $E^*(Y)$, as a set of fixed points of a closed and compact correspondence, must be closed. To get to claim (1), I will show that any converging sequence

$\rho_L \in E_L(T)$ has its limit in $E^*(T)$, and any $\rho \in E^*(T)$ has a converging sequence in $E_L(T)$ approaching it. In other words, upper and lower hemicontinuity properties hold for the equilibrium correspondence along sequences of covering discretizations.

Lemma 4. *For all covering sequences \mathbf{Y}_L ,*

$$\lim_{K \rightarrow \infty} H(E_L(T), E^*(T)) = 0,$$

where $H(\cdot, \cdot)$ is the Hausdorff-distance.

Proof. We first show upper hemicontinuity in K . Suppose u.h.c. is not satisfied. Then there exists a subsequence $\rho_{L_t} \in E_{L_t}(T)$ with $\rho_{L_t} \rightarrow_t \bar{\rho} \notin E^*(T)$. The converging subsequence exists since \mathbf{Y}^2 is compact. To ease notation, re-define $L = L_t$ for the rest of the proof. Not being an equilibrium, we have that there exists $\rho_T \neq \bar{\rho}$ that maximizes the deviation payoff

$$\Delta_T = W(\rho_T, \bar{\rho}, T) - W(\bar{\rho}, \bar{\rho}, T) > 0.$$

Pick $\varepsilon \in (0, \Delta_T)$. By convergence of ρ_L , and by continuity of W , we have that there exists $L_{1,T}$ such that for all $L \geq L_{1,T}$,

$$\left| W(\rho_T, \rho_L, T) - W(\rho_T, \bar{\rho}, T) \right| \leq \frac{\varepsilon}{3}. \quad (9)$$

By the same argument, there is a $L_{2,T}$ s.t. for all $L \geq L_{2,T}$,

$$\left| W(\rho_L, \rho_L, T) - W(\bar{\rho}, \bar{\rho}, T) \right| \leq \frac{\varepsilon}{3}. \quad (10)$$

Furthermore, we can always choose $\bar{L}_T \geq \max\{L_{1,T}, L_{2,T}\}$ large enough so that $\alpha_L \leq \frac{\varepsilon}{3}$, implying

$$\left| W(f^L(\rho_T), \rho_L, T) - W(\rho_z, \rho_L, T) \right| \leq \frac{\varepsilon}{3}. \quad (11)$$

Take $L \geq \bar{L}_T$. Define the best deviation under the discrete game as

$$\hat{\rho}_L = \arg \max_{\rho \in \mathbf{Y}_L^2 \setminus \rho_L} W(\rho, \rho_L, T).$$

Now we have

$$\begin{aligned}
W(\hat{\rho}_L, \rho_L, T) - W(\rho_L, \rho_L, T) &\geq W(f^L(\rho_T), \rho_L, T) - W(\rho_L, \rho_L, T) \\
&= W(\rho_T, \rho_L, T) - W(\rho_L, \rho_L, T) + \beta_{1,L} \\
&= W(\rho_T, \bar{\rho}, T) - W(\bar{\rho}, \bar{\rho}, T) + \beta_{1,L} + \beta_{2,L} + \beta_{3,L} \\
&\geq \Delta_T + \beta_{1,L} + \beta_{2,L} + \beta_{3,L},
\end{aligned}$$

where $\beta_{1,L}$ corresponds to the projection error (11), and $\beta_{2,L}, \beta_{3,L}$ correspond to (9), (10) respectively. We have that $|\beta_{i,L}| \leq \frac{\varepsilon}{3}$, and thus

$$W(\hat{\rho}_L, \rho_L, T) - W(\rho_L, \rho_L, T) \geq \Delta_T - \varepsilon > 0,$$

implying that $\rho_L \notin E_L(T)$, a contradiction.

For lower hemicontinuity, Assumption 6 imposes that for all equilibria in $E^*(T)$ for all players, Hessians at the equilibrium are negative definite. Thus, small deviations must lead to strict payoff loss. Fix T , then the proof is via contradiction: there exists some equilibrium $\rho \in E^*(T)$ that is not approximated by any sequence of equilibria in $E_L(T)$. The proof can be done analogously to the one above; defining $\Delta_T > 0$ as the best deviation payoff:

$$\Delta_T = W(\rho, \rho, T) - \max_{\mathbf{Y}^2 \setminus \rho} W(\rho_T, \rho, T) > 0.$$

Since $\Delta_T > 0$, we can find a fine enough discretizations s.t. ρ can be approximated arbitrarily closely, in which case incentives must also align, by continuity of W . The contradiction follows as before. \square

Continuity of the equilibrium correspondence gives us that for all $\varepsilon > 0$ there is $K > 0$ large enough so that

$$\|V_L(T) - V^*(T)\| < \varepsilon,$$

with $V_L(T), V^*(T)$ being the maximal payoff over the equilibrium sets $E_L(T), E^*(T)$.

To make judgements about $\sup_{T \in \mathcal{T}} V_L(T)$, a uniform continuity property of $V_L(T)$ will be useful. By Assumption 6, all equilibria in E_L and E^* are hyperbolic, for K large enough.

Hyperbolicity implies that an implicit function theorem holds: For any $\rho \in E_L$, there exists neighborhoods $\mathcal{N}_\rho, \mathcal{N}_T$ of ρ, T and a continuous map $h : \mathcal{N}_T \rightarrow \mathcal{N}_\rho$ such that $h(T) \in E_L(T)$ for all $T \in \mathcal{N}_T$. Thus, the equilibrium correspondences $E_L(T), E^*(T)$ are continuous in T for all K large enough.

As $W(\rho, T)$ is continuous both in ρ and T ¹⁵, and equilibrium correspondences are continuous in T , Berge's Maximum Theorem applies. We have that $V_L(T)$ is continuous in T . Moreover, as the payoff functions $W(\rho, T)$ are bounded, twice differentiable in ρ and transition probabilities T (evaluated at ρ), these payoff functions are Lipschitz both in ρ and T . Then, $V_L(T)$ are bounded, Lipschitz as well for K large enough. This follows first from continuity of $E^*(T)$, and then by the Lipschitz property of $W(\rho, T)$. Finally, boundedness and the Lipschitz property imply that $V_L(T)$ are equicontinuous in T .

By the Arzelá-Ascoli Theorem, boundedness and equicontinuity of $V_L(T)$ implies the existence of a subsequence K_m so that V_{K_m} converges uniformly to some V . Pointwise convergence of $V_L(T)$ implies that this limit satisfies $V = V^*$. A simple contradiction argument with another application of Arzelá-Ascoli shows that indeed $V_L(T)$ converges uniformly to V^* .

Twice differentiability of T and boundedness implies, by the Arzelá-Ascoli Theorem, the relative compactness of \mathcal{T} . Hence, we have that $V_L(T)$ converges uniformly to V^* over a relatively compact set \mathcal{T} . It follows that

$$\lim_{K \rightarrow \infty} \sup_{T \in \mathcal{T}} V_L(T) = \sup_{T \in \mathcal{T}} V^*(T).$$

■

B.4. Proof of Proposition 2. Point (i) follows directly from Theorem 3. Point (ii) is an application of the analysis in Kandori (1992).

B.5. Proof of Proposition 3. We are proving the following statement, restated in more technical terms:

Proposition 4. *Take any state variable S , and any symmetric $\rho^* \in E_S$. Consider best response dynamics $F_{B,S}$. When all eigenvalues of $J(\rho^*)$ are real,*

¹⁵Regarding functions T , consider the sup-norm as metric on \mathcal{T} .

- (i) ρ^* is asymptotically stable given $\bar{\alpha} \in (0, 1)$ if and only if it is asymptotically stable given $\bar{\alpha} = 1$.
- (ii) When $\rho^* = \rho_N$, ρ^* is asymptotically stable given $\bar{\alpha} \in (0, 1)$ if and only if it is asymptotically stable given $\bar{\alpha} = 1$ also under gradient dynamics $F_{G,S}$.
- (iii) When $\rho^* = \rho_N$, ρ^* is asymptotically stable when one player uses best response, and the other uses gradient dynamics if and only if ρ^* is asymptotically stable given $\bar{\alpha} = 1$ under $F_{B,S}$.

When there are $K \geq 1$ states, for any $\alpha \in (0, 1]$ we can write the linearization $M_B(\alpha)$ of F_S at symmetric $\rho^* \in E_S$ as follows

$$M_B(\alpha) = \begin{bmatrix} -\alpha I_K & \alpha J \\ J & -I_K \end{bmatrix},$$

where J is the best-response derivative matrix at ρ^* of players 1 and 2, by symmetry. As $-I_K$ and J commute, the characteristic equation of $M_B(\alpha)$ can be written as

$$\begin{aligned} \text{char}(\lambda) &= \det(\alpha J J - (\alpha + \lambda)(1 + \lambda)I_K) \\ &= \det\left(\alpha^{\frac{1}{2}}J - ((\alpha + \lambda)(1 + \lambda))^{\frac{1}{2}}I_K\right) \det\left(\alpha^{\frac{1}{2}}J + ((\alpha + \lambda)(1 + \lambda))^{\frac{1}{2}}I_K\right). \end{aligned}$$

Thus, for any μ such that $|\mu| \in \text{eig}(J)$, $\lambda_{1,2} \in \text{eig}(M_B(\alpha))$ where $\lambda_{1,2}$ are the solutions to

$$\lambda^2 + (1 + \alpha)\lambda + \alpha(1 - \mu^2) = 0,$$

i.e.

$$\lambda_{1,2} = -\frac{1 + \alpha}{2} \pm \sqrt{\left(\frac{1 + \alpha}{2}\right)^2 - \alpha(1 - \mu^2)}.$$

$M_B(1)$ has all eigenvalues negative if and only if $|\mu| < 1$. $\lambda < 0$ if and only if

$$\alpha(1 - \mu^2) > 0,$$

which is equivalent to $|\mu| < 1$. Thus, stability under $\alpha = 1$ carries over to all $\alpha \in (0, 1)$.

As for point (ii), note that the Hessian of $W(\rho)$ at ρ_N equals $\pi_{11}^N I_K$. An analogous argument to the proof of Lemma 1 can be used to show that ρ_N is asymptotically stable

under gradient learning if and only if it is under F_S . Thus, the above conclusion remains under gradient learning for ρ_N . Similarly, note that for point (iii), at ρ_N we can consider the difference between $F_{S,B}$ and $F_{S,G}$ as being down to scaling every value in $F_{S,B}$ by the constant $-\pi_{11}^N I_K$, which one can interpret as some $\alpha > 0$, and reach the required conclusion.

Regarding other symmetric equilibria $\rho^* \in E_S$, the connection between gradient learning and best response dynamics is more tenuous. For any symmetric $\rho^* \neq \rho_N$, we may not write

$$M_G = \nu M_B,$$

for some scalar $\nu \neq 0$. In general, we'd have

$$M_G = D M_B,$$

where $D \geq 0$ is a diagonal matrix of positive, but varying, entries. In this case, it is easy to construct examples where M_B has negative eigenvalues only, but M_G does not. \blacksquare

B.6. Proof of Theorem 5. We start by a derivation analogous to the proof of Theorem 3, using (6). First, as transitions of \mathbf{S}_{PR} and \mathbf{S}_D are independent, it is useful to look into properties of the joint transition matrix $M \in [0, 1]^{KL}$, which now represents an irreducible Markov chain over states $r \in \mathbf{R} = \mathbf{S}_D \times \mathbf{S}_{PR}$. Note that M is a function of policies ρ through transition function T over S_{PR} . With some abuse of notation, let $W(\rho, r)$ be the expected discounted payoffs of player 1 in this setting, given initial state $r \in \mathbf{R}$. Given this notation of states, we denote F_R as the best-response dynamic solutions of which ACQ learners would converge to in this setting. As in the proof of Theorem 3, we can derive, at an equilibrium ρ^* ,

$$W_{ij} = [I_K - \delta M]^{-1} \left[\delta \frac{\partial^2 M}{\partial \rho_i \partial \rho_j} \tilde{W} + (1 - \delta) \frac{\partial^2 U}{\partial \rho_i \partial \rho_j} \right],$$

for $1 \leq i, j \leq KL$, using the notation as before, noting that now W, U are KL -dimensional vectors. Order states r so that $r_1, \dots, r_K = (s_D^1, s_1), \dots, (s_D^1, s_K)$, i.e. s_D is kept fixed while $s_k \in \mathbf{S}_{PR}$ advances. Then, at ρ_N it follows that $U_N = U(\rho_N)$ is such that the first K elements equal $\pi_N^1 = \pi(p_N^1, s_D^1)$, etc, until the last K elements equal $\pi_N^L = \pi(p_N^L, s_D^L)$.

To simplify these derivatives, the following will be helpful: Letting T_D, T be the transition matrices of $\mathbf{S}_D, \mathbf{S}_{PR}$ respectively, due to the ordering over states we can write $M = T_D \otimes T$, where \otimes is the Kronecker product. As the spectral radius of δM is $\delta < 1$ (M is a row-stochastic matrix), we can apply the geometric series expansion and properties of the Kronecker product to get

$$[I - \delta M]^{-1} = \sum_{q=1}^{\infty} \delta^q M^q = \sum_{q=1}^{\infty} \delta^q (T_D^q \otimes T^q).$$

Let $1 \leq \alpha, \beta \leq L$ be indices for $K \times K$ blocks of $[I - \delta M]^{-1}$, which we denote $[I - \delta M]_{\alpha, \beta}^{-1}$. Note that for each such block, rows sum to the same constant. Let $\mathbf{1}_K$ be a K dimensional column of ones. Row sums can be written as

$$[I - \delta M]_{\alpha, \beta}^{-1} \mathbf{1}_K = \sum_{q=1}^{\infty} \delta^q (T_{D, \alpha, \beta}^q \otimes T^q) \mathbf{1}_K = \sum_{q=1}^{\infty} \delta^q T_{D, \alpha, \beta}^q \mathbf{1}_K = C_{\alpha, \beta} \mathbf{1}_K,$$

as rows of T^q sum to one for each $0 \leq q$, and where $C_{\alpha, \beta} = \frac{1}{1 - \delta T_{D, \alpha, \beta}}$. As also U_N is constant across each K -element block, it follows that $W = [I - \delta M]^{-1} U_N$ is constant across each K -element block. This significantly simplifies the cross-derivatives required for the Jacobian of F_R to be analysed here. Derivatives of M such as $\frac{\partial^2 M}{\partial \rho_i \partial \rho_j}$ sum to zero across each K -element block, which implies that $\frac{\partial^2 M}{\partial \rho_i \partial \rho_j} W = \mathbf{0}_{KL}$. We can thus write the derivatives of W evaluated at ρ_N in manner analogous to Theorem 3, taking r_1 as initial state without loss:

$$W_{ii} = [I_K - \delta M]_1^{-1} (1 - \delta) \frac{\partial^2 U_N}{(\partial \rho_i)^2},$$

$$W_{ij} = [I_K - \delta M]_1^{-1} (1 - \delta) \left[\frac{\partial^2 U_N}{\partial \rho_i \partial \gamma_j} + \delta \frac{\partial M}{\partial \rho_i} [I_K - \delta M]^{-1} \frac{\partial U_N}{\partial \gamma_j} \right],$$

where γ_j is some opponent strategy with $1 \leq i, j \leq KL$. For each such i, j , let $s_D(i), s_D(j)$ be the associated state $s_D \in \mathbf{S}_D$. With some abuse of notation, let $\pi_{11}^N(i) = \frac{\partial^2 U_N}{(\partial \rho_i)^2}$, $\pi_{12}^N(i) = \frac{\partial^2 U_N}{\partial \rho_i \partial \gamma_j}$, and $\pi_2^N(j) = \frac{\partial U_N}{\partial \gamma_j}$ similarly to Theorem 3. Then, letting e_i be KL -dimensional columns

which are zero everywhere but in element i ,

$$\begin{aligned} W_{ii} &= [I - \delta M]_1^{-1} (1 - \delta) \pi_{11}^N(i) e_i, \\ W_{ij} &= [I - \delta M]_1^{-1} (1 - \delta) \left[\pi_{12}^N(i) e_i + \delta \frac{\partial M}{\partial \rho_i} [I - \delta M]^{-1} \pi_2^N(i) e_i \right], \end{aligned}$$

if $i = j$, and

$$W_{ij} = [I - \delta M]_1^{-1} (1 - \delta) \left[\delta \frac{\partial M}{\partial \rho_i} [I - \delta M]^{-1} \pi_2^N(j) e_j \right],$$

otherwise, following the notation used in section 4. To save notation, write $Z = [I - \delta M]^{-1}$. Then,

$$\frac{W_{ij}}{W_{ii}} = \begin{cases} \frac{1}{\pi_{11}^N(i)} \left[\pi_{12}^N(i) + \pi_2^N(i) \delta \sum_{\ell=1}^{KL} \zeta_{r_i r_\ell} Z_{\ell, i} \right] & \text{if } i = j \\ \frac{1}{\pi_{11}^N(i)} \left[\pi_2^N(j) \delta \sum_{\ell=1}^{KL} \zeta_{r_i r_\ell} Z_{\ell, j} \right] & \text{o.w.} \end{cases}$$

where $\zeta_{rr'} = \frac{\partial}{\partial \rho(r)} M_{rr'}$. For the proof of point (1), we will upper bound eigenvalues of this system. Note from the above that we can write J , the Jacobian of F_R , as

$$J = (E_1 \otimes I_K) + (E_2 \otimes I_K) \delta \zeta Z (E_3 \otimes I_K),$$

where E_1, E_2, E_3 are L dimensional diagonal matrices with values equal to $-\frac{\pi_{12}^N(s_D^m)}{\pi_{11}^N(s_D^m)}, \frac{-1}{\pi_{11}^N(s_D^m)}, \pi_2^N(s_D^m)$, respectively over $1 \leq m \leq L$.

Now to the proof of point (1). As in the proof of Theorem 4, recall that the spectral radius κ of J can be upper bounded by any matrix norm of J . Hence,

$$\begin{aligned} \kappa &\leq \|E_1\| \|I_K\| + \delta \|E_2 \otimes I_K\| \|\zeta\| \|Z\| \|E_3 \otimes I_K\| = \|E_1\| + \frac{\delta}{1 - \delta} \|E_2\| \|E_3\| \|\zeta\| \\ &= \lambda_1 + \lambda_2 \lambda_3 \frac{\delta}{1 - \delta} \|\zeta\|, \end{aligned}$$

where $\lambda_1, \lambda_2, \lambda_3$ are the largest absolute entries of the diagonal matrices E_1, E_2, E_3 . Recall that we define, analogously to Theorem 3,

$$\begin{aligned}\pi_2^N(s_D^m) &= X^m \xi_c^m L_N^m \\ \pi_{11}^N(s_D^m) &= X^m \left(\xi_{o,1}^m L_N^m + \xi_o^m \frac{\partial}{\partial p_1} L_N^m \right) \\ \pi_{12}^N(s_D^m) &= X^m \xi_{o,2}^m L_N^m,\end{aligned}$$

where superscript m denotes payoff terms given state s_D^m , all evaluated at $p_N(s_D^m)$, which is dropped for ease of notation. Then,

$$\begin{aligned}\lambda_1 &= \max_m \left| \frac{\xi_{o,2}^m}{\xi_{o,1}^m + \xi_o^m G_N^m} \right| \\ \lambda_2 &= \max_m \left(X^m \left| \xi_{o,1}^m L_N^m + \xi_o^m \frac{\partial}{\partial p_1} L_N^m \right| \right)^{-1} \\ \lambda_3 &= \max_m X^m |\xi_c^m L_N^m|.\end{aligned}$$

For the first insight, given the upper bound on κ above, note that $\|\zeta\|$ shrinks with $\sup_{a \in \Omega} \left| \frac{\partial}{\partial P} g(a; P) \right|_{P=P_N}$, and the result follows. Point (1) then follows directly also from this upper bound.

For point (2), recall that for any matrix, the trace equals the sum of its eigenvalues. In this case,

$$\text{tr}(J) = K \text{tr}(E_1) + \delta \sum_{m=1}^L E_{2,m} E_{3,m} (\zeta Z)_m,$$

where $E_{2,m}, E_{3,m}$ are m -th elements of E_2, E_3 , and $(\zeta Z)_m$ is the m -th diagonal $K \times K$ block of ζZ . Thus, for all ζZ there are elements of E_1, E_2, E_3 large enough in absolute value so that $|\text{tr}(J)|$ is large, which can only be true if at least one eigenvalue of J is large. ■

REFERENCES

- Abreu, Dilip, David Pearce, and Ennio Stacchetti (1986). “Optimal cartel equilibria with imperfect monitoring”. In: *Journal of Economic Theory* 39.1, pp. 251–269.
- (1990). “Toward a theory of discounted repeated games with imperfect monitoring”. In: *Econometrica: Journal of the Econometric Society*, pp. 1041–1063.
- Assad, Stephanie et al. (2024). “Algorithmic pricing and competition: empirical evidence from the German retail gasoline market”. In: *Journal of Political Economy* 132.3, pp. 723–771.

- Banchio, Martino and Giacomo Mantegazza (2022). “Games of Artificial Intelligence: A Continuous-Time Approach”. In: *arXiv preprint arXiv:2202.05946*.
- Benaïm, M (1999). “Dynamics of Stochastic Approximation, Le Seminaire de Probabilite’, Springer Lecture Notes in Mathematics”. In.
- Benaïm, Michel and Mathieu Faure (2012). “Stochastic approximation, cooperative dynamics and supermodular games”. In: *The Annals of Applied Probability* 22.5, pp. 2133–2164.
- Borkar, Vivek S (2009). *Stochastic approximation: a dynamical systems viewpoint*. Vol. 48. Springer.
- Brown, Zach Y and Alexander MacKay (2021). *Competition in pricing algorithms*. Tech. rep. National Bureau of Economic Research.
- Calvano, Emilio, Giacomo Calzolari, Vincenzo Denicolo, et al. (2020). “Artificial intelligence, algorithmic pricing, and collusion”. In: *American Economic Review* 110.10, pp. 3267–97.
- Calvano, Emilio, Giacomo Calzolari, Vincenzo Denicoló, et al. (2021). “Algorithmic collusion with imperfect monitoring”. In: *International journal of industrial organization* 79, p. 102712.
- Cartea, Álvaro et al. (2022). “Algorithms can Learn to Collude: A Folk Theorem from Learning with Bounded Rationality”. In: *Available at SSRN 4293831*.
- Faure, Mathieu and Gregory Roth (2010). “Stochastic approximations of set-valued dynamical systems: Convergence with positive probability to an attractor”. In: *Mathematics of Operations Research* 35.3, pp. 624–640.
- Fudenberg, Drew and David M Kreps (1993). “Learning mixed equilibria”. In: *Games and economic behavior* 5.3, pp. 320–367.
- Fudenberg, Drew, David Levine, and Eric Maskin (1994). “The Folk Theorem with Imperfect Public Information”. In: *Econometrica* 62.5, pp. 997–1039.
- Fudenberg, Drew and David K Levine (2009). “Learning and equilibrium”. In: *Annu. Rev. Econ.* 1.1, pp. 385–420.
- Gaunersdorfer, Andrea and Josef Hofbauer (1995). “Fictitious play, Shapley polygons, and the replicator equation”. In: *Games and Economic Behavior* 11.2, pp. 279–303.
- Gilboa, Itzhak and Akihiko Matsui (1991). “Social stability and equilibrium”. In: *Econometrica: Journal of the Econometric Society*, pp. 859–867.
- Hart, Sergiu and Andreu Mas-Colell (2003). “Uncoupled dynamics do not lead to Nash equilibrium”. In: *American Economic Review* 93.5, pp. 1830–1836.
- Johnson, Justin, Andrew Rhodes, and Matthijs R Wildenbeest (2020). “Platform design when sellers use pricing algorithms”. In: *Available at SSRN 3753903*.
- Kandori, Michihiro (1992). “The use of information in repeated games with imperfect monitoring”. In: *The Review of Economic Studies* 59.3, pp. 581–593.
- Klein, Timo (2021). “Autonomous algorithmic collusion: Q-learning under sequential pricing”. In: *The RAND Journal of Economics* 52.3, pp. 538–558.
- Lamba, Rohit and Sergey Zhuk (2022). “Pricing with algorithms”. In: *arXiv preprint arXiv:2205.04661*.
- Leslie, David S and Edmund J Collins (2003). “Convergent multiple-timescales reinforcement learning algorithms in normal form games”. In: *The Annals of Applied Probability* 13.4, pp. 1231–1251.
- Leslie, David S, Steven Perkins, and Zibo Xu (2020). “Best-response dynamics in zero-sum stochastic games”. In: *Journal of Economic Theory* 189, p. 105095.

- Loots, Thomas and Arnoud V. den Boer (2023). “Data-driven collusion and competition in a pricing duopoly with multinomial logit demand”. In: *Production and Operations Management* 32.4, pp. 1169–1186.
- Mazumdar, Eric, Lillian J Ratliff, and S Shankar Sastry (2020). “On gradient-based learning in continuous games”. In: *SIAM Journal on Mathematics of Data Science* 2.1, pp. 103–131.
- Meylahn, Janusz M and Arnoud V. den Boer (2022). “Learning to collude in a pricing duopoly”. In: *Manufacturing & Service Operations Management* 24.5, pp. 2577–2594.
- Milgrom, Paul and John Roberts (1990). “Rationalizability, learning, and equilibrium in games with strategic complementarities”. In: *Econometrica: Journal of the Econometric Society*, pp. 1255–1277.
- (1991). “Adaptive and sophisticated learning in normal form games”. In: *Games and economic Behavior* 3.1, pp. 82–100.
- Palis Jr, J, W de Melo, et al. (1982). “Geometric Theory of Dynamical Systems”. In.
- Papadimitriou, Christos and Georgios Piliouras (2018). “From nash equilibria to chain recurrent sets: An algorithmic solution concept for game theory”. In: *Entropy* 20.10, p. 782.
- Puterman, Martin L (2014). *Markov decision processes: discrete stochastic dynamic programming*. John Wiley & Sons.
- Robbins, Herbert and Sutton Monro (1951). “A stochastic approximation method”. In: *The annals of mathematical statistics*, pp. 400–407.
- Salcedo, Bruno (2015). “Pricing algorithms and tacit collusion”. In: *Manuscript, Pennsylvania State University*.
- Schulman, John et al. (2017). “Proximal policy optimization algorithms”. In: *arXiv preprint arXiv:1707.06347*.
- Watkins, Christopher John Cornish Hellaby (1989). “Learning from delayed rewards”. In.