

Algorithms for Stochastic Games with Perfect Monitoring*

Dilip Abreu Benjamin Brooks Yuliy Sannikov

December 31, 2019

Abstract

We study the pure-strategy subgame-perfect Nash equilibria of stochastic games with perfect monitoring, geometric discounting, and public randomization. We develop novel algorithms for computing equilibrium payoffs, in which we combine policy iteration when incentive constraints are slack with value iteration when incentive constraints bind. We also provide software implementations of our algorithms. Preliminary simulations indicate that they are significantly more efficient than existing methods. The theoretical results that underlie the algorithms also imply bounds on the computational complexity of equilibrium payoffs when there are two players. When there are more than two players, we show by example that the number of extreme equilibrium payoffs may be countably infinite.

Keywords: Stochastic game, perfect monitoring, algorithm, computation.

JEL classification: C63, C72, C73, D90.

*Abreu: Department of Economics, New York University, dilip.abreu@nyu.edu; Brooks: Department of Economics, University of Chicago, babrooks@uchicago.edu; Sannikov: Graduate School of Business, Stanford University, sannikov@gmail.com. This work has benefited from the comments of numerous seminar audiences. We have also benefited from the superb research assistance of Mathieu Cloutier, Moshe Katzwer, and Kai Hao Yang. We are very grateful to Joel Sobel for his guidance and to several anonymous referees for their valuable input. Finally, we would like to acknowledge financial support from the National Science Foundation Grant #1530823.

1 Introduction

This paper develops new algorithms for computing equilibrium payoffs in stochastic games. Specifically, we study the payoffs that can be attained in pure-strategy subgame-perfect Nash equilibria of repeated games with perfect monitoring, public randomization, and a stochastically evolving state variable. The current state determines which actions are feasible for the players as well as the payoffs of those actions. The chosen actions in turn influence the future evolution of the state. This classical structure is used to model a wide variety of phenomena in economics and in other disciplines. The range of applications include: dynamic oligopoly with investment (in, e.g., capacity, research and development, advertising), risk sharing, and the dynamics of political bargaining and compromise (cf. Ericson and Pakes, 1995; Kocherlakota, 1996; Dixit, Grossman, and Gul, 2000).

The standard methodology for computing subgame-perfect equilibrium payoffs in repeated games comes from Abreu, Pearce, and Stacchetti (1986, 1990), hereafter APS. They showed that the set of equilibrium payoffs satisfies a recursive relationship that is analogous to the Bellman equation from dynamic programming. In particular, any equilibrium payoff can be decomposed into a flow payoff from the first period of play plus the expected discounted payoff from the next period onward, which, by subgame perfection, is also an equilibrium payoff. Just as the value function is the fixed point of the Bellman operator, so too the equilibrium payoff set is the largest fixed point of an operator that produces the set of payoffs which can be generated using continuation values chosen from a given set. Moreover, APS show that iterating this operator on any set that contains all equilibrium payoffs yields a sequence of sets that asymptotically converges to the set of equilibrium payoffs. Although APS wrote explicitly about games with imperfect monitoring and without a state variable, their results mechanically extend to the class of games studied here, where payoffs are generated in each state using continuation payoffs drawn from a received payoff correspondence.¹

Our main contribution is a refinement of the APS algorithm. In the analogy with dynamic programming, the APS algorithm is identified with value function iteration. We combine this approach with a form of policy iteration, which is used to partially solve out equilibrium payoffs when incentive constraints are slack. The resulting hybrid algorithm converges faster than existing methods, and the hybrid operator that is used in this algorithm is of bounded computational complexity when there are two players. The approach also leads to new structural insights about equilibria that generate extreme equilibrium payoffs, namely that

¹For early extensions involving a state variable see Atkeson (1991) and Phelan and Stacchetti (2001). A more recent application is Hörner et al. (2011). For a more complete description of the self-generation methodology for stochastic games, see Mailath and Samuelson (2006).

play must be stationary until the first period in which an incentive constraint binds.

The approach has several novel elements, some of which apply even to repeated games, while others are tailored to the stochastic setting. To motivate a defining element of our refinement, consider an infinitely repeated Prisoners' Dilemma. We normalize the Nash payoffs to zero and the payoff from mutual cooperation to 1. Suppose that the discount factor δ is such that static Nash is the only action profile that can be supported, even when all feasible payoffs can be promised as continuation values. A fortiori, we can conclude that the equilibrium payoff set is $\{(0,0)\}$. Nonetheless, this fact will only be discovered by the APS operator asymptotically: at iteration k , the payoff (δ^k, δ^k) will still be in the APS approximation. For even though mutual cooperation is not used to generate new payoffs, it is still implicitly being played k periods in the future through the received set of continuation payoffs. There is, in a sense, an internal inconsistency in the way that the APS operator is generating new payoffs in this example: It is claimed that static Nash generates an equilibrium payoff that maximizes the sum of the players' payoffs, and that this sum is strictly positive. But the sum of payoffs in the first period is zero, meaning that the sum of the continuation payoffs must be even higher than the best payoff we can generate! This is obviously not sustainable in equilibrium.

This example illustrates a more general principle in repeated games: Fix a set of welfare weights, i.e., a direction in payoff space, and consider the equilibrium with the highest payoffs in this direction. The action profile played in the first period of this equilibrium must have flow payoffs that are weakly above the highest equilibrium payoffs. Indeed, this principle extends to stochastic games: Fix a direction in welfare space. For each state, there is a highest equilibrium payoff in this direction, which is generated by playing some action profile in the first period, followed by continuation equilibrium payoffs in every continuation state. The continuation payoffs are bounded by the highest equilibrium payoff in their respective states, and sometimes further if required by incentive compatibility. As a result, the highest equilibrium payoff in a given state is below a *recursive level*, obtained by playing the optimal action profile for one period with the highest equilibrium payoffs as continuation payoffs.

This principle motivates our refinement. For a given direction, consider a policy of action profiles meant to attain the highest payoffs in each state. The highest payoffs in this direction are attained by recursively continuing this policy, with the highest equilibrium continuation values, until incentive compatibility requires some value burning. This happens when the best incentive compatible payoff, i.e., the level generated by the APS operator, is lower than that attained by following the policy. This leads to the following *max-min-max* problem to

bound the level of payoffs:

$$\text{max over actions of } \left(\text{min of (the “recursive” level and max level generated by APS)} \right).$$

The max-min-max operator maps a payoff correspondence into the correspondence of payoffs that are below the max-min-max level in all directions. Our first main result shows that iterative application of this operator can be used to compute equilibrium payoffs.

Because our operator uses levels that are weakly below those of APS, it generates smaller correspondences than APS. As a result, the sequence it generates converges weakly faster than the APS sequence. Indeed, there are even examples, such as the aforementioned Prisoners’ Dilemma, where the APS sequence only converges asymptotically, but the max-min-max sequence converges after finitely many rounds (although this is not generally the case). That our operator is theoretically novel and cuts more sharply is clear, but perhaps its main advantage is that *it is significantly easier to compute*. This derives from additional theoretical structure that the operator embodies, as we now explain.

First, as long as the max-min-max sequence decreases at the first iteration in the sense of set containment,² our operator has the following crucial property: it relies on the maximum level generated by APS only when an incentive constraint binds for some player. Thus, while our operator nominally requires us to know the APS levels, in fact it is sufficient to know the maximal level attained with payoffs in which an incentive constraint binds—a considerably easier task.

In addition, when computing the max-min-max levels, we maximize over a tuple of action profiles, and we minimize over a tuple of what we call *regimes*, which indicate for each state whether the minimum level is APS or recursive. We refer to the action profiles and regimes collectively as a *policy*. We show that for each direction, there exists a policy that optimizes payoffs simultaneously in all states. For any direction, the optimal policy can be found through a form of policy iteration. The policy is optimal if and only if for every state there is no alternative action profile which, if substituted in, would increase the level, and there is no regime substitution that would lower the level. Once we have found the optimal policy for one direction, it is easy to compute a range of directions for which it remains optimal, and also the improving substitution when the optimum changes.

When there are two players, these properties yield an especially powerful implementation. We find the optimal policy for a starting direction. After that, we move clockwise. By considering one-state substitutions, we endogenously identify directions where the optimal

²There are many initial correspondences that guarantee this will happen. Two examples are the feasible payoff correspondence and a correspondence that in every state is equal to a large hypercube that contains all of the flow payoffs for all states. The latter is what we use in our numerical simulations.

policy changes and update the optimal policy. After a full revolution, we have bounded payoffs in all directions. Moreover, for two players, each action profile can generate at most four extreme binding payoffs, generalizing the result of Abreu and Sannikov (2014) for repeated games. This leads to a bound on both the complexity of our operator and of the equilibrium payoff correspondence.

We have implemented this algorithm as a software package that is freely available through the online supplement and an author’s website.³ We report a number of numerical examples, including a risk-sharing game à la Kocherlakota (1996).

When there are more than two players, we show by example that the number of extreme equilibrium payoffs may be countably infinite. As a result, exact computation of equilibrium payoffs may be impossible. We therefore propose a more flexible operator that bounds payoffs with the max-min-max level for a subset of directions, which is dynamically updated between applications. We show that binding payoffs will remain sufficient to determine the APS level as long as any legacy directions we drop are not needed to determine the binding payoffs or the local frontier around them. There are no restrictions on how directions can be added. For any such sequence of direction sets, the algorithm is guaranteed to converge to a correspondence that contains all equilibrium payoffs. There are many ways to use this characterization, and we focus on one simple implementation: At every round, we drop some directions that are redundant for computing binding payoffs. If the number of directions is below a fixed bound, we randomly add new directions that correspond to “faces” of the exact max-min-max correspondence. These directions are computed via a generalization of the two-player direction rotation procedure.

In Online Appendix C, we present simulations where the equilibrium payoff correspondence has a finite number of faces that are successfully discovered by the algorithm, so that the sequence of approximations converges exactly to the equilibrium payoff correspondence. We also solve a three-player risk-sharing game to show a new result of independent economic interest, which is that formal insurance contracts between a subset of the players can lead to lower payoffs for all players.

All of the aforementioned algorithms converge to equilibrium payoffs from the outside, thus providing upper bounds. As a last topic, we show how our methodology can be adapted produce a lower bound. In particular, we construct an algorithm that necessarily converges, after finitely many steps, to a correspondence that *strictly self-generates in every direction*, thus robustly certifying that payoffs in this correspondence can be attained in equilibrium.

A notable antecedent of our work is Abreu and Sannikov (2014) who studied repeated games with two players, perfect monitoring, and public randomization. They proposed a

³www.benjaminbrooks.net/software.shtml

distinct refinement of the APS algorithm, and our procedure does not reduce to theirs when there is a single state. As mentioned above, they give a bound on the number of extreme equilibrium payoffs, which is tighter than our bound (due to the specialization to repeated games) but is based on the same geometry.

This paper supersedes our earlier work (Abreu, Brooks, and Sannikov, 2016), in which we studied the same class of games but restricted to two players. Based on similar ideas, we proposed a related but distinct algorithm, which also proceeds by iteratively modifying a payoff tuple, one state at a time, to obtain a sequence of payoffs and corresponding bounds. In contrast to the present work, that operator did not have bounded computational complexity in the two-player case and did not apply to many-player games.

Another key reference is Judd, Yeltekin, and Conklin (2003), hereafter JYC, who proposed the approximation of the APS operator by bounding payoffs in a fixed and finite set of directions. While they wrote about repeated games, their methodology readily generalizes to the class of stochastic games we consider.⁴ Key differences between the approaches are that we use the max-min-max operator rather than the APS operator, we endogenize the directions in order to identify faces, and we compute our operator exactly when there are two players. We report runtime comparisons between the implementations of our algorithm with our implementation of JYC for stochastic games. Preliminary simulations indicate that our algorithm is significantly faster than that of JYC. We note that our approach makes heavy use of perfect monitoring, whereas the JYC approach can be easily adapted to games with imperfect public monitoring.

There are other lines of work on computing equilibria without public randomization (Berg and Kitti, 2019) or with mixed strategies (Berg, 2019). There is also a large body of work computing Markov perfect equilibria (Pakes and McGuire, 1994). Our methodology can be used to bound payoffs in Markov equilibria, but it cannot be used to compute just the set of Markov equilibrium payoffs. Renner and Scheidegger (2018) propose a machine learning algorithm for approximating feasible payoffs in dynamic principal-agent problems, whose efficacy is suggested by simulations. In contrast, our algorithms pertain to stochastic games, and we show analytically that our procedure is guaranteed to converge to the equilibrium payoff correspondence.

Finally, our algorithms exploit the linear structure of equilibria. Many of the concepts we use are evocative of similar concepts in linear programming. The connection between linear programming and dynamic programming is well known, so this is not altogether unexpected. Online Appendix D studies the connection in detail. Our conclusion is that while there are

⁴Such an extension is done by Yeltekin, Cai, and Judd (2017). A related approach is taken by Sleet and Yeltekin (2016), using what they call block correspondences instead of bounding in fixed directions.

deep connections, but there are also fundamental differences because the minimization over regimes makes our optimization program non-convex, so it cannot simply be reduced to a linear program.

The rest of this paper is organized as follows. Section 2 describes the basic model and background material from APS. Section 3 presents our algorithm and its key properties. Section 4 studies implementation and examples when there are two players. Section 5 presents the infinite extreme point example and studies implementation for many players. Section 6 adapts our methodology to bound equilibrium payoffs from below. Section 7 is a conclusion.

2 Setting and background

Players $i = 1, \dots, N$ interact over infinitely many periods. There is a finite set of states S . If the current state is $s \in S$, player i takes an action a_i in a finite set $\mathbf{A}_i(s)$.⁵ The set of action profiles in state s is $\mathbf{A}(s) = \times_{i=1}^N \mathbf{A}_i(s)$. Player i 's flow utility from action profile $a \in \mathbf{A}(s)$ is $g_i(a|s)$. The resulting probability that the next state is s' is $\pi(s'|a, s)$. We will henceforth assume that each a is available in a single state and write $g_i(a)$ and $\pi(s'|a)$.⁶ Players discount future payoffs at the common rate $\delta \in (0, 1)$. Actions and the state are perfectly observable.

We will study the equilibrium payoff correspondence $\mathbf{V} : S \rightarrow 2^{\mathbb{R}^N}$, where $\mathbf{V}(s)$ is the set of expected discounted payoffs that can be achieved in some pure-strategy subgame-perfect Nash equilibrium with public randomization, when the initial state of the world is s . For a formal definition of an equilibrium in this setting, see Mailath and Samuelson (2006, Sections 5.5 and 5.7), in particular Corollary 5.7.1.⁷

Subgame-perfection implies that any equilibrium payoff can be decomposed into the discount-weighted sum of a flow payoff which is obtained in the first period and expected continuation equilibrium payoffs from the second period onwards. The technique of APS is to generalize this recursive relationship in a manner that is analogous to how the Bellman operator generalizes the recursive characterization of the value function in dynamic programming. Explicitly, fix a compact-valued payoff correspondence $\mathbf{W} : S \rightarrow 2^{\mathbb{R}^N}$. Note that the assumption of compactness of \mathbf{W} is maintained throughout. The associated *threat tuple* $\underline{\mathbf{w}}(\mathbf{W})$ is

$$\underline{\mathbf{w}}_i(\mathbf{W})(s) = \min \{w_i | (w_i, w_{-i}) \in \mathbf{W}(s) \text{ for some } w_{-i} \in \mathbb{R}^{N-1}\}.$$

⁵In general, we will use boldface to denote functions whose domain is the state space.

⁶This is without loss, since we could simply redefine an action to be the ordered pair (a_i, s) .

⁷Strictly speaking, the definition of an equilibrium in Mailath and Samuelson (2006) differs slightly from the one which we are implicitly using. They assume that there is a probability distribution over the initial state, while we define equilibrium payoffs *conditional* on the initial state.

For an action profile $a \in \mathbf{A}(s)$, let

$$\underline{u}_i(a, \mathbf{W}) = \max_{a'_i} \left[(1 - \delta)g_i(a'_i, a_{-i}) + \delta \sum_{s' \in S} \pi(s'|a'_i, a_{-i}) \underline{\mathbf{w}}_i(\mathbf{W})(s') \right].$$

We say that v is *generated in state s by the action profile $a \in \mathbf{A}(s)$ and the correspondence \mathbf{W}* if there exist $\mathbf{w} \in \mathbf{W}$ —meaning that $\mathbf{w}(\cdot)$ is a selection from $\mathbf{W}(\cdot)$ —such that

$$v = (1 - \delta)g(a) + \delta \sum_{s' \in S} \pi(s'|a) \mathbf{w}(s'); \quad (1)$$

$$(1 - \delta)g_i(a) + \delta \sum_{s' \in S} \pi(s'|a) \mathbf{w}_i(s') \geq \underline{u}_i(a, \mathbf{W}) \quad \forall i = 1, \dots, N. \quad (2)$$

This equation implicitly assumes that a deviation from a by player i will be punished by a transition to the worst continuation equilibrium for the deviator, which results in a payoff of $\underline{\mathbf{w}}_i(\mathbf{W})(s')$ if the next state is s' . This is without loss due to perfect monitoring.

Let $B(a, \mathbf{W})$ denote the set of payoffs that are generated by $a \in \mathbf{A}(s)$ and \mathbf{W} , and let $B(\mathbf{W})(s) = \text{co}(\cup_{a \in \mathbf{A}(s)} B(a, \mathbf{W}))$, where co denotes the convex hull. It is a fact that B is increasing in \mathbf{W} and maps compact-valued correspondences to compact-valued correspondences. We say that \mathbf{W} is *self-generating* if $\mathbf{W} \subseteq B(\mathbf{W})$, i.e., $\mathbf{W}(s) \subseteq B(\mathbf{W})(s)$ for all $s \in S$. APS's arguments, extended to stochastic games, show that if \mathbf{W} is bounded and self-generating, then $B(\mathbf{W}) \subseteq \mathbf{V}$. These properties imply that \mathbf{V} is the largest bounded self-generating payoff correspondence.⁸ Moreover, the following algorithm can be used to compute \mathbf{V} : Let \mathbf{W}^0 be any correspondence that contains \mathbf{V} , and generate the sequence $\mathbf{W}^k = B(\mathbf{W}^{k-1})$ for $k \geq 1$. Then $\cap_{k \geq 0} \mathbf{W}^k = \mathbf{V}$. Moreover, if $B(\mathbf{W}^0) \subseteq \mathbf{W}^0$, then the sequence is decreasing: $\mathbf{W}^k \subseteq \mathbf{W}^{k-1}$ for all $k > 0$.

3 A refinement of the algorithm of APS

We now describe our refinement of the APS algorithm. This algorithm will similarly generate a sequence of payoff correspondences via iterative application of a new operator \tilde{B} . This algorithm converges faster, as the operator \tilde{B} generates smaller correspondences and is significantly easier to compute than the APS operator B . Note that Online Appendix A contains pseudocode for the algorithms developed over the next three sections.

⁸Note that a pure-strategy subgame-perfect equilibrium need not exist. Subgame-perfection requires that the continuation equilibrium is an equilibrium after every history. Thus, an equilibrium exists in some state if and only if an equilibrium exists in every state. As a result, \mathbf{V} is either empty in all states or non-empty in all states.

3.1 Our operator

Preliminary to defining \tilde{B} , it is useful to reformulate the APS operator. Let $\Lambda = \{\lambda \in \mathbb{R}^N \mid \|\lambda\| = 1\}$ denote the set of N -dimensional *directions*, endowed with the subspace topology. For each $\lambda \in \Lambda$, $s \in S$, and $a \in \mathbf{A}(s)$, we define

$$x^{APS}(a, \lambda, \mathbf{W}) = \max \{\lambda \cdot v \mid v \in B(a, \mathbf{W})\},$$

where our convention is that the max of an empty set is $-\infty$. In addition, we say that an action profile a is *supportable (at \mathbf{W})* if $B(a, \mathbf{W}) \neq \emptyset$, so that $x^{APS}(a, \lambda, \mathbf{W})$ is finite for all λ . Let $\mathbf{A}(\mathbf{W})(s)$ be the set of supportable action profiles in state s . We further define $x^{APS}(s, \lambda, \mathbf{W}) = \max_{a \in \mathbf{A}(\mathbf{W})(s)} x^{APS}(a, \lambda, \mathbf{W})$ to be the maximum level in the direction λ that is attained by the APS operator, which again is $-\infty$ if there are no supportable actions in state s . Then

$$B(\mathbf{W})(s) = \{v \mid \lambda \cdot v \leq x^{APS}(s, \lambda, \mathbf{W}) \ \forall \lambda \in \Lambda\}.$$

Note that $B(\mathbf{W})(s)$ is empty if there are no supportable action profiles in state s . In addition, as per Footnote 8, if \mathbf{W} is empty in some state, then there are no continuation value profiles, no action can be supported, and $B(\mathbf{W})$ is empty in every state.

Our operator will be defined similarly in terms of bounding hyperplanes, but with tighter bounds than x^{APS} . To motivate the bounds, let us briefly consider which payoffs would be attainable in the absence of incentive constraints. In other words, what are the *feasible payoffs* that can be generated using some strategy profile? In a repeated game, the answer is simply the convex hull of the flow payoffs. In a stochastic game, things are more complicated because of how action profiles influence the evolution of the state. For any fixed welfare weights $\lambda \in \Lambda$, however, the problem of maximizing the λ -weighted sum of expected discounted payoffs is simply a Markov decision problem. Blackwell (1965) showed that there is a stationary solution, i.e., a selection of action profiles $\mathbf{a} \in \mathbf{A}$ such that an optimal strategy is to play $\mathbf{a}(s)$ whenever the state is s . Equivalently, the optimal strategy involves playing an optimal action for one period, followed by recursively starting the optimal strategy over in the next period. This strategy simultaneously attains the highest levels in all states. The associated optimal levels $x(s)$ are the unique solution to

$$x(s) = (1 - \delta)\lambda \cdot g(\mathbf{a}(s)) + \delta \sum_{s' \in S} \pi(s' \mid \mathbf{a}(s))x(s') \quad (3)$$

for all $s \in S$.

The situation is more complicated with incentive constraints, because it may not be possible to attain the levels in (3) without giving some player an incentive to deviate. In particular, it may be necessary to burn some surplus in some states in the direction λ in order to give sufficiently large continuation values to deter deviations. Exactly how much surplus needs to be burnt depends on the threat point and the shape of the frontier, which are things we do not know until we actually compute \mathbf{V} . We can, however, use an approximation \mathbf{W} that contains \mathbf{V} to bound the amount of value burning required to enforce any action profile $\mathbf{a}(s)$. At the same time, as we discussed in the introduction, there are cases where the APS bound is also too generous, because of spuriously large continuation values in \mathbf{W} .

These considerations motivate the hybrid approach that we now adopt, which is to use the recursive methodology of Markov decision problems in some states and APS-style bounds in other states. Holding fixed the actions played in the first period, we will select the configuration of recursive or APS for each state in order to minimize the bound on payoffs, thus ensuring that our approximation is not too generous.

To that end, let us define a *policy* to be a pair (\mathbf{a}, \mathbf{r}) , where $\mathbf{a} \in \mathbf{A}(\mathbf{W})$,⁹ and $\mathbf{r} : S \rightarrow \{R, APS\}$ is a *regime* (where the R stands for *recursive*). Let \mathbf{R} denote the set of regimes. For given λ and \mathbf{W} , consider the system

$$\mathbf{y}(s) = \begin{cases} (1 - \delta)\lambda \cdot g(\mathbf{a}(s)) + \delta \sum_{s' \in S} \pi(s' | \mathbf{a}(s)) \mathbf{y}(s') & \text{if } \mathbf{r}(s) = R; \\ x^{APS}(\mathbf{a}(s), \lambda, \mathbf{W}) & \text{if } \mathbf{r}(s) = APS \end{cases} \quad (4)$$

for all $s \in S$. Standard arguments can be used to show (4) has a unique solution: given $\mathbf{y} : S \rightarrow \mathbb{R}$, let $T(\mathbf{y}, \lambda, \mathbf{a}, \mathbf{r}, \mathbf{W})$ be the tuple defined by

$$T(\mathbf{y}, \lambda, \mathbf{a}, \mathbf{r}, \mathbf{W})(s) = \begin{cases} (1 - \delta)\lambda \cdot g(\mathbf{a}(s)) + \delta \sum_{s' \in S} \pi(s' | \mathbf{a}(s)) \mathbf{y}(s') & \text{if } \mathbf{r}(s) = R; \\ x^{APS}(\mathbf{a}(s), \lambda, \mathbf{W}) & \text{if } \mathbf{r}(s) = APS. \end{cases}$$

Clearly, any \mathbf{y} that satisfies (4) must be a fixed point of the operator $T(\cdot, \lambda, \mathbf{a}, \mathbf{r}, \mathbf{W})$. Moreover, this operator is a contraction of modulus δ in \mathbf{y} , so that a fixed point exists and is unique. We denote it by $x(s, \lambda, \mathbf{a}, \mathbf{r}, \mathbf{W})$. We record some properties of T for future reference.

Lemma 1. *Fix λ , \mathbf{W} , $\mathbf{a} \in \mathbf{A}(\mathbf{W})$, and \mathbf{r} . As a function of $\mathbf{y} : S \rightarrow \mathbb{R}$, T is*

(L1.i) *increasing;*

(L1.ii) *a contraction with modulus δ and hence has a unique fixed point \mathbf{y}^* ;*

⁹Our focus is primarily on computing which payoffs can be generated, taking as a given which action profiles are supportable. The computation of $\mathbf{A}(\mathbf{W})$ is a straightforward by-product of other necessary calculations. See a discussion in Section 4.2.3.

(L1.iii) if $T(\mathbf{y}) \leq (\geq) \mathbf{y}$ then $\mathbf{y}^* \leq (\geq) T(\mathbf{y})$.

Proof of Lemma 1.

(L1.i) This is immediate from positive linearity of T in \mathbf{y} when $\mathbf{r}(s) = R$ and the fact that it is independent of \mathbf{y} when $\mathbf{r}(s) = APS$.

(L1.ii) Let $\|\cdot\|$ denote the sup norm. Then

$$\begin{aligned} \|T(\mathbf{y}, \lambda, \mathbf{a}, \mathbf{r}, \mathbf{W}) - T(\mathbf{y}', \lambda, \mathbf{a}, \mathbf{r}, \mathbf{W})\| &= \delta \max_{\{s|\mathbf{r}(s)=R\}} \sum_{s' \in S} \pi(s'|\mathbf{a}(s)) |\mathbf{y}(s') - \mathbf{y}'(s')| \\ &\leq \delta \|\mathbf{y} - \mathbf{y}'\| \end{aligned}$$

as desired. The rest of the result follows from the Banach fixed point theorem.

(L1.iii) If $T(\mathbf{y}, \lambda, \mathbf{a}, \mathbf{r}, \mathbf{W}) \leq (\geq) \mathbf{y}$ then from (L1.i), we conclude that the sequence \mathbf{y}^k generated by iterative application of T starting with $\mathbf{y}^0 = \mathbf{y}$ is monotonically decreasing (increasing) and, by (L1.ii), must converge to the unique fixed point \mathbf{y}^* . The result then follows.

□

The next key definition mirrors that of x^{APS} :

$$x(s, \lambda, \mathbf{W}) = \max_{\mathbf{a} \in \mathbf{A}(\mathbf{W})} \min_{\mathbf{r} \in \mathbf{R}} x(s, \lambda, \mathbf{a}, \mathbf{r}, \mathbf{W}). \quad (5)$$

Finally, the operator \tilde{B} is defined according to

$$\tilde{B}(\mathbf{W})(s) = \{v | \lambda \cdot v \leq x(s, \lambda, \mathbf{W}) \forall \lambda \in \Lambda\}.$$

We refer to this as the *max-min-max operator*, since for each direction, we maximize over action tuples, minimize over regimes, and maximize over APS payoffs. Because of the minimization over regimes, \tilde{B} generates smaller correspondences than B . Furthermore, the simultaneous determination of levels in states for which the R regime is specified collapses multiple rounds of the APS operator into a single step.

We now verify that the operator \tilde{B} satisfies all of the critical properties of the APS operator, so that it can in fact be used to compute \mathbf{V} :

Theorem 1 (The max-min-max operator). *\tilde{B} has the following properties:*

(T1.i) \tilde{B} is increasing in \mathbf{W} , and if \mathbf{W} is compact, then $\tilde{B}(\mathbf{W})$ is compact;

(T1.ii) $\tilde{B}(\mathbf{W}) \subseteq B(\mathbf{W})$, and if $\mathbf{W} \subseteq \tilde{B}(\mathbf{W})$, then $\tilde{B}(\mathbf{W}) \subseteq \mathbf{V}$;

(T1.iii) $\mathbf{V} = \tilde{B}(\mathbf{V})$;

(T1.iv) Fix a correspondence $\tilde{\mathbf{W}}^0$ that contains \mathbf{V} . Define the sequence $\{\tilde{\mathbf{W}}^k\}_{k=0}^{\infty}$ by $\tilde{\mathbf{W}}^k = \tilde{B}(\tilde{\mathbf{W}}^{k-1})$. Then $\mathbf{V} = \bigcap_{k=0}^{\infty} \tilde{\mathbf{W}}^k$.

Remark 1. Recall that it is possible that no pure-strategy equilibria exist, and $\mathbf{V}(s) = \emptyset$ for all s . If this happens, then since the $\tilde{\mathbf{W}}^k$ correspondences are closed and decreasing, there must be some k at which $\tilde{\mathbf{W}}^k(s) = \emptyset$ for some s . In the next iteration, there will be no supportable action profiles in any state, and $\tilde{\mathbf{W}}^{k+1}$ will be empty in every state, at which point the algorithm converges. Theorem 1 and our subsequent results encompass the case where payoff correspondences are empty and no payoffs can be generated, and our proofs remain correct with the convention that the maximum of an empty set is $-\infty$. The explicit focus is, however, on the non-trivial case where \mathbf{W} is non-empty valued.

Proof of Theorem 1.

(T1.i) It is clear that $x^{APS}(a, \lambda, \mathbf{W})$ is increasing in \mathbf{W} for every $a \in \mathbf{A}(\mathbf{W})(s)$ and λ . Hence, if $\mathbf{W}' \subseteq \mathbf{W}$, then for all $(\mathbf{y}, \lambda, \mathbf{a}, \mathbf{r})$, where $\mathbf{a} \in \mathbf{A}(\mathbf{W})$, $T(\mathbf{y}, \lambda, \mathbf{a}, \mathbf{r}, \mathbf{W}) \geq T(\mathbf{y}, \lambda, \mathbf{a}, \mathbf{r}, \mathbf{W}')$, which immediately implies that the fixed point of $T(\cdot, \lambda, \mathbf{a}, \mathbf{r}, \mathbf{W})$ is greater than the fixed point of $T(\cdot, \lambda, \mathbf{a}, \mathbf{r}, \mathbf{W}')$. Thus, $x(s, \lambda, \mathbf{a}, \mathbf{r}, \mathbf{W})$ is increasing in \mathbf{W} . As a result, $\min_{\mathbf{r} \in \mathbf{R}} x(s, \lambda, \mathbf{a}, \mathbf{r}, \mathbf{W})$, and $x(s, \lambda, \mathbf{W})$ are also increasing in \mathbf{W} . If \mathbf{W} is compact, then $x^{APS}(a, \lambda, \mathbf{W})$ is bounded above for every λ . Thus, $\tilde{B}(\mathbf{W})(s)$ is bounded and closed, being the intersection of closed half-spaces.

(T1.ii) Clearly, $x(s, \lambda, \mathbf{W}) \leq x^{APS}(s, \lambda, \mathbf{W})$, which implies that \tilde{B} is always contained in B . Thus, if $\mathbf{W} \subseteq \tilde{B}(\mathbf{W})$, then $\mathbf{W} \subseteq B(\mathbf{W})$ and hence, by APS, $B(\mathbf{W}) \subseteq \mathbf{V}$. Consequently, $\tilde{B}(\mathbf{W}) \subseteq \mathbf{V}$.

(T1.iii) From (T1.ii), it suffices to show that $\mathbf{V} \subseteq \tilde{B}(\mathbf{V})$, i.e., for all λ , $x(s, \lambda, \mathbf{V}) \geq x^{APS}(s, \lambda, \mathbf{V})$. To that end, fix λ , and for all s , let $\mathbf{a}(s)$ be an action that maximizes $x^{APS}(a, \lambda, \mathbf{V})$ and let $\mathbf{w}(s')$ be the associated continuation values as a function of the next-period state s' . We will show that $\min_{\mathbf{r} \in \mathbf{R}} x(s, \lambda, \mathbf{a}, \mathbf{r}, \mathbf{V}) \geq x^{APS}(s, \lambda, \mathbf{V})$, so that $x(s, \lambda, \mathbf{V}) \geq x^{APS}(s, \lambda, \mathbf{V})$, which implies the result. Since $\mathbf{V} = B(\mathbf{V})$, $x^{APS}(s, \lambda, \mathbf{V}) \geq \lambda \cdot u$ for all $u \in \mathbf{V}(s')$ for all s' . Since $\mathbf{w}(s') \in \mathbf{V}(s')$ for all s' ,

$$\begin{aligned} x^{APS}(s, \lambda, \mathbf{V}) &= (1 - \delta)\lambda \cdot g(\mathbf{a}(s)) + \delta \sum_{s' \in S} \pi(s' | \mathbf{a}(s)) \lambda \cdot \mathbf{w}(s') \\ &\leq (1 - \delta)\lambda \cdot g(\mathbf{a}(s)) + \delta \sum_{s' \in S} \pi(s' | \mathbf{a}(s)) x^{APS}(s', \lambda, \mathbf{V}). \end{aligned}$$

Thus, if we let $\mathbf{y}(s) = x^{APS}(s, \lambda, \mathbf{V})$ for all s , then for *any* regimes \mathbf{r} , $T(\mathbf{y}, \lambda, \mathbf{a}, \mathbf{r}, \mathbf{V}) \geq \mathbf{y}$ (with equality if $\mathbf{r}(s) = APS$ and weak inequality if $\mathbf{r}(s) = R$). By (L1.iii), we conclude that $\mathbf{y}(s) = x^{APS}(s, \lambda, \mathbf{V}) \leq x(s, \lambda, \mathbf{a}, \mathbf{r}, \mathbf{V}) = \mathbf{y}^*(s)$, as required.

(T1.iv) (T1.ii) implies that $\widetilde{\mathbf{W}}^k \subseteq \mathbf{W}^k$, where the latter is the k th element of the APS sequence starting from $\widetilde{\mathbf{W}}^0$. Also, the fact that $\widetilde{\mathbf{W}}^0$ contains \mathbf{V} , (T1.i), and (T1.iii) imply that $\mathbf{V} \subseteq \widetilde{\mathbf{W}}^k$. Thus, $\mathbf{V} \subseteq \cap_k \widetilde{\mathbf{W}}^k \subseteq \cap_k \mathbf{W}^k = \mathbf{V}$.

□

Remark 2. Roughly, our definition of a policy (\mathbf{a}, \mathbf{r}) allows us to treat separately states in which incentive constraints are slack from those in which they bind. When incentive constraints are slack, optimal behavior is stationary until a constraint binds, payoffs are defined recursively, and R is the optimal regime. When incentive constraints bind, value burning is required to provide incentives and the minimal regime is necessarily APS . At the fixed point (where $\mathbf{W} = \mathbf{V}$) and at the corresponding optimal policy in the direction λ , this description is exactly correct: the recursive regimes are the most permissive and yield upper bounds on attainable levels. This fundamental observation underlies the proof of (T1.iii). However, the algorithm approaches the fixed point from the “outside,” i.e., $\mathbf{V} \subseteq \widetilde{\mathbf{W}}^k$, and it is possible that recursion along the path of the algorithm yields lower levels than the corresponding APS levels, since the latter may rely on spuriously generous continuation payoffs. This motivates our requirement that the regimes are chosen to minimize levels. In the next section, we show that this minimization reduces to simply setting the regime to R in states for which this yields a level below that of APS. As we are dealing with a simultaneous equation system, this test is modestly more complicated than it might appear at first.

Remark 3. $\widetilde{B}(\mathbf{W})(s)$ is the intersection of hyperplanes where the levels are given by $x(s, \cdot, \mathbf{W})$. While the former is necessarily a convex set, the latter is generally *not* a convex function, so that $x(s, \cdot, \mathbf{W})$ need not be the support function of $\widetilde{B}(\mathbf{W})(s)$. An example in which this happens is presented in Section 4.3.1. Thus, one must be careful in applying intuition from convex geometry to \widetilde{B} and the sets $\widetilde{\mathbf{W}}^k$.

Remark 4. It is a straightforward consequence of (T1.iv) that the sequence $\widetilde{\mathbf{W}}^k$ converges to \mathbf{V} in the Hausdorff metric. In practice, we terminate the algorithm when the distance between successive iterates falls below some threshold (which it must eventually, since the sequence is Cauchy). The distance between iterates has no simple relationship with the distance to \mathbf{V} , and there is no guarantee that the final correspondence is close to the fixed point. In Section 6, we adapt our algorithm to produce a sub-correspondence of \mathbf{V} , which can be used to bound the error in the approximation.

The remainder of this section develops further properties of \tilde{B} that make it especially tractable for computation, namely, the *state independence of the optimal policy* and the *sufficiency of binding payoffs* when the APS level is minimal.

3.2 State-independence of the optimal policy

The definition of $x(s, \lambda, \mathbf{W})$ in (5) leaves open the possibility that the optimal policy in the direction λ depends on s . We now show that there exists a policy that is simultaneously optimal for all states, and we lay the foundations for a simple algorithm to compute it.

3.2.1 Minimal regimes

For notational economy, we shall temporarily suppress the dependence of x and other objects on \mathbf{W} , and simply write $x(s, \lambda)$, etc. For $a \in \mathbf{A}(s)$ define

$$x^R(a, \lambda, \mathbf{a}, \mathbf{r}) = (1 - \delta)\lambda \cdot g(a) + \delta \sum_{s'} \pi(s'|a)x(s', \lambda, \mathbf{a}, \mathbf{r}).$$

Given λ and \mathbf{a} , we say that the regime \mathbf{r} is *minimal* if for all $s \in S$, $x(s, \lambda, \mathbf{a}, \mathbf{r}) = \min_{\mathbf{r}' \in \mathbf{R}} x(s, \lambda, \mathbf{a}, \mathbf{r}')$. In other words, they minimize the level in all states simultaneously. In addition, given \mathbf{r} , let $\mathbf{r} \setminus s$ be the regime that is the same as \mathbf{r} in every state except s , i.e., we flip the regime in state s .

We now show that there exist minimal regimes. Moreover, there is a simple set of inequalities that characterize when regimes are minimal. These inequalities will be central to our implementations in Sections 4 and 5.

Lemma 2 (Minimal regimes). *For all $\mathbf{a} \in \mathbf{A}(\mathbf{W})$ and λ ,*

(L2.i) *there exist minimal regimes;*

(L2.ii) *\mathbf{r} is minimal if and only if for all $s \in S$,*

$$x(s, \lambda, \mathbf{a}, \mathbf{r}) = \min \{x^{APS}(\mathbf{a}(s), \lambda), x^R(\mathbf{a}(s), \lambda, \mathbf{a}, \mathbf{r})\}; \quad (6)$$

(L2.iii) *if (6) is violated for some s , then \mathbf{r} is not minimal. Moreover, for all $s' \in S$, $x(s', \lambda, \mathbf{a}, \mathbf{r} \setminus s) \leq x(s', \lambda, \mathbf{a}, \mathbf{r})$, with strict inequality in state s .*

Proof of Lemma 2. These results follow directly from Lemma 1:

(L2.iii) If (6) is violated at s , then letting $\mathbf{y} = x(\lambda, \mathbf{a}, \mathbf{r})$, we conclude that $T(\mathbf{y}, \lambda, \mathbf{a}, \mathbf{r} \setminus s) \leq \mathbf{y}$.

By (L1.iii), $x(\lambda, \mathbf{a}, \mathbf{r} \setminus s) = \mathbf{y}^* \leq T(\mathbf{y}, \lambda, \mathbf{a}, \mathbf{r} \setminus s) \leq \mathbf{y} = x(\lambda, \mathbf{a}, \mathbf{r})$, where the penultimate inequality is strict by assumption in state s .

(L2.ii) Only if follows from (L2.iii). For the if direction, suppose that \mathbf{r} satisfies (6) for all s . Then for any $\mathbf{r}' \in \mathbf{R}$, it follows that $\mathbf{y} \leq T(\mathbf{y}, \lambda, \mathbf{a}, \mathbf{r}')$, where $\mathbf{y} = x(\lambda, \mathbf{a}, \mathbf{r})$. Then by (L1.iii), $\mathbf{y}^* = x(\lambda, \mathbf{a}, \mathbf{r}') \geq x(\lambda, \mathbf{a}, \mathbf{r})$, so that \mathbf{r} is indeed minimal.

(L2.i) Let \mathbf{r} solve $\min_{\mathbf{r}' \in \mathbf{R}} \sum_{s \in S} x(s, \lambda, \mathbf{a}, \mathbf{r}')$. We argue that \mathbf{r} is minimal. Suppose not. Then by (L2.ii), (6) is violated at some $s \in S$, and by (L2.iii), $\sum_{s' \in S} x(s', \lambda, \mathbf{a}, \mathbf{r} \setminus s) < \sum_{s' \in S} x(s', \lambda, \mathbf{a}, \mathbf{r})$, a contradiction.

□

3.2.2 Maximal actions

We now extend these results to actions: as long as there are supportable action profiles in every state, there exists a selection of action profiles that maximizes the level for all states simultaneously, and maximal actions are characterized by a simple set of inequalities.

We will prove this result using an operator that is analogous to T but directly imposes minimality. For $\mathbf{y} : S \rightarrow \mathbb{R}$ let

$$T^{\min}(\mathbf{y}, \lambda, \mathbf{a})(s) = \min \left\{ x^{APS}(\mathbf{a}(s), \lambda), (1 - \delta)\lambda \cdot g(\mathbf{a}(s)) + \delta \sum_{s' \in S} \pi(s' | \mathbf{a}(s)) \mathbf{y}(s') \right\}.$$

Lemma 3. Fix λ and $\mathbf{a} \in \mathbf{A}(\mathbf{W})$. As a function of $\mathbf{y} : S \rightarrow \mathbb{R}$, T^{\min} is

(L3.i) increasing;

(L3.ii) a contraction with modulus δ , and hence has a unique fixed point \mathbf{y}^* ;

(L3.iii) if $T^{\min}(\mathbf{y}) \leq (\geq) \mathbf{y}$ then $\mathbf{y}^* \leq (\geq) T^{\min}(\mathbf{y})$;

Proof of Lemma 3. The proof is identical to that of Lemma 1, replacing T with T^{\min} . □

Now, let us define $x(s, \lambda, \mathbf{a}) = \min_{\mathbf{r} \in \mathbf{R}} x(s, \lambda, \mathbf{a}, \mathbf{r})$ to be the minimal levels associated with the action tuple $\mathbf{a} \in \mathbf{A}(\mathbf{W})$. The definition of $x(s, \lambda)$ in (5) implies that $x(s, \lambda) = \max_{\mathbf{a} \in \mathbf{A}(\mathbf{W})} x(s, \lambda, \mathbf{a})$. For a given λ , we say that \mathbf{a} is *maximal* if for all $s \in S$, $x(s, \lambda) = x(s, \lambda, \mathbf{a})$, i.e., \mathbf{a} attains the max-min-max level in all states simultaneously. Also define

$$x^R(a, \lambda, \mathbf{a}) = (1 - \delta)\lambda \cdot g(a) + \delta \sum_{s' \in S} \pi(s' | a) x(s', \lambda, \mathbf{a}).$$

Finally, for $\mathbf{a} \in \mathbf{A}$, $s \in S$, and $a \in \mathbf{A}(s)$, let us define the substituted action tuple $\mathbf{a} \setminus (s, a) \in \mathbf{A}$ to be the action a in state s and $\mathbf{a}(s')$ in each state $s' \neq s$.

Lemma 4 (Maximal actions). Suppose that $\mathbf{A}(\mathbf{W})$ is non-empty valued. Then for all λ ,

(L4.i) *there exist maximal actions;*

(L4.ii) $\mathbf{a} \in \mathbf{A}(\mathbf{W})$ *is maximal if and only if for all* $s \in S$ *and* $a \in \mathbf{A}(\mathbf{W})(s)$,

$$x(s, \lambda, \mathbf{a}) \geq \min \{x^{APS}(a, \lambda), x^R(a, \lambda, \mathbf{a})\}, \quad (7)$$

(L4.iii) *if (7) is violated for some* $s \in S$ *and* $a \in \mathbf{A}(\mathbf{W})(s)$, *then* \mathbf{a} *is not maximal. Moreover, for all* $s' \in S$, $x(s', \lambda, \mathbf{a} \setminus (s, a)) \geq x(s', \lambda, \mathbf{a})$, *with strict inequality in state* s .

Proof of Lemma 4. The proof mirrors that of Lemma 2, and follows directly from Lemma 3:

(L4.iii) If (7) is violated at (s, a) , then letting $\mathbf{y} = x(\lambda, \mathbf{a})$, we conclude that $T^{min}(\mathbf{y}, \lambda, \mathbf{a} \setminus (s, a)) \geq \mathbf{y}$. By (L3.iii), $x(\lambda, \mathbf{a} \setminus (s, a)) = \mathbf{y}^* \geq T^{min}(\mathbf{y}, \lambda, \mathbf{a} \setminus (s, a)) \geq \mathbf{y} = x(\lambda, \mathbf{a})$, where the penultimate inequality is strict by assumption, in state s .

(L4.ii) Only if follows from (L4.iii). For the if direction, suppose that \mathbf{a} satisfies (7) for all s and $a \in \mathbf{A}(\mathbf{W})(s)$. Then for any $\mathbf{a}' \in \mathbf{A}(\mathbf{W})$ it follows that $\mathbf{y} \geq T^{min}(\mathbf{y}, \lambda, \mathbf{a}')$ where $\mathbf{y} = x(\lambda, \mathbf{a})$. Then by (L3.iii), $\mathbf{y}^* = x(\lambda, \mathbf{a}') \leq x(\lambda, \mathbf{a})$, so that \mathbf{a} is indeed maximal.

(L4.i) Let \mathbf{a} solve $\max_{\mathbf{a}' \in \mathbf{A}(\mathbf{W})} \sum_{s \in S} x(s, \lambda, \mathbf{a}')$. We argue that \mathbf{a} is maximal. Suppose not. Then by (L4.ii), (7) is violated at some $s \in S$ and $a \in \mathbf{A}(\mathbf{W})(s)$ and by (L4.iii) $\sum_{s' \in S} x(s', \lambda, \mathbf{a} \setminus (s, a)) > \sum_{s' \in S} x(s', \lambda, \mathbf{a})$, a contradiction.

□

To summarize, there exists a state-independent optimal policy. This is extremely useful for computation, since it means we can solve for optimal levels in all states simultaneously.

3.3 The sufficiency of binding payoffs

In the discussion around Theorem 1, we identified the *APS* regime with a situation in which incentive constraints bind, so that value burning is required to deter deviations. There is nothing in the definition of x^{APS} , however, that requires that incentive constraints bind, and we have left open the possibility that the *APS* regime is minimal even though incentive constraints are slack. While this may happen, it does not occur in equilibrium or along the sequence $\widetilde{\mathbf{W}}^k$, as long as $B(\widetilde{\mathbf{W}}^0) \subseteq \widetilde{\mathbf{W}}^0$. This fact is immensely useful for computation, since it means that we do not need to compute optimal APS levels in all directions (which would amount to computing the APS operator itself). Instead, we can just compute the optimal APS level when incentive constraints bind.

We now develop this result formally. Let us reintroduce the payoff correspondence \mathbf{W} as an argument in optimal levels. For any $s \in S$ and $a \in \mathbf{A}(\mathbf{W})(s)$, define

$$\widehat{x}^{APS}(a, \lambda, \mathbf{W}) = \max \{ \lambda \cdot w \mid w \in B(a, \mathbf{W}) \text{ and } \exists i \text{ s.t. (2) holds as equality} \}. \quad (8)$$

We refer to the difference

$$\gamma(a, \lambda, \mathbf{W}) = x^{APS}(a, \lambda, \mathbf{W}) - \widehat{x}^{APS}(a, \lambda, \mathbf{W})$$

as the *APS gap*.¹⁰ We shall see that for any λ and $\mathbf{a} \in \mathbf{A}(\mathbf{W})$, if $\gamma(\mathbf{a}(s), \lambda, \mathbf{W})$ is strictly positive in state s , then R is a minimal regime in state s . On the other hand, if $\gamma(\mathbf{a}(s), \lambda, \mathbf{W})$ is zero, then by definition we can restrict attention to APS payoffs for which at least one player's incentive constraint binds.

To prove this result, we need two intermediate lemmas. We say that the APS operator B *sub-generates at \mathbf{W} in the direction λ* if for all $s \in S$, $x^{APS}(s, \lambda, \mathbf{W}) \leq \max \{ \lambda \cdot w \mid w \in \mathbf{W}(s) \}$, and B *sub-generates at \mathbf{W}* if $B(\mathbf{W}) \subseteq \mathbf{W}$. These notions extend to \widetilde{B} in the obvious way.

Lemma 5. *Suppose that B sub-generates at \mathbf{W} in the direction λ . Then for any $\mathbf{a} \in \mathbf{A}(\mathbf{W})$ and s , if $\gamma(\mathbf{a}(s), \lambda, \mathbf{W}) > 0$, then*

$$x^{APS}(\mathbf{a}(s), \lambda, \mathbf{W}) \geq x^R(\mathbf{a}(s), \lambda, \mathbf{a}, \mathbf{W}) = x(s, \lambda, \mathbf{a}, \mathbf{W}). \quad (9)$$

Moreover, there exist minimal regimes such that $\mathbf{r}(s) = R$ for all s with $\gamma(\mathbf{a}(s), \lambda, \mathbf{W}) > 0$.

Proof of Lemma 5. Suppose that $\gamma(\mathbf{a}(s), \lambda, \mathbf{W}) > 0$. Then any maximal continuation values in \mathbf{W} in the direction λ , denoted \mathbf{w} , must be incentive compatible for $\mathbf{a}(s)$, and

$$x^{APS}(\mathbf{a}(s), \lambda, \mathbf{W}) = (1 - \delta) \lambda \cdot g(\mathbf{a}(s)) + \delta \sum_{s' \in S} \pi(s' | \mathbf{a}(s)) \lambda \cdot \mathbf{w}(s').$$

Sub-generation and the definition of x imply that for all s' , $\lambda \cdot \mathbf{w}(s') \geq x^{APS}(\mathbf{a}(s'), \lambda, \mathbf{W}) \geq x(s', \lambda, \mathbf{W})$. Hence,

$$\begin{aligned} x^{APS}(\mathbf{a}(s), \lambda, \mathbf{W}) &\geq (1 - \delta) \lambda \cdot g(\mathbf{a}(s)) + \delta \sum_{s' \in S} \pi(s' | \mathbf{a}(s)) x(s', \lambda, \mathbf{W}) \\ &\geq (1 - \delta) \lambda \cdot g(\mathbf{a}(s)) + \delta \sum_{s' \in S} \pi(s' | \mathbf{a}(s)) x(s', \lambda, \mathbf{a}, \mathbf{W}) = x^R(\mathbf{a}(s), \lambda, \mathbf{a}, \mathbf{W}) \end{aligned}$$

¹⁰Note that if a is supportable, then there must exist a payoff where (2) holds as an equality, so that \widehat{x}^{APS} is finite. This follows from the definition of $\underline{u}_i(a, \mathbf{W})$, which is at least the payoff obtained by playing a with the worst continuation values in \mathbf{W} . Thus, it cannot be that $\underline{u}_i(a, \mathbf{W})$ is strictly below every payoff that satisfies (1) for some $\mathbf{w} \in \mathbf{W}$.

as desired.

Finally, suppose \mathbf{r} is minimal and $\gamma(\mathbf{a}(s), \lambda, \mathbf{W}) > 0$. If $x^{APS}(\mathbf{a}(s), \lambda, \mathbf{W}) > x^R(\mathbf{a}(s), \lambda, \mathbf{a}, \mathbf{W})$, then $\mathbf{r}(s) = R$. Otherwise, (9) implies that $x^{APS}(\mathbf{a}(s), \lambda, \mathbf{W}) = x^R(\mathbf{a}(s), \lambda, \mathbf{a}, \mathbf{W})$. Thus, if we set $\mathbf{r}'(s) = R$ for all states with $\gamma(\mathbf{a}(s), \lambda, \mathbf{W}) > 0$ and $\mathbf{r}'(s') = \mathbf{r}(s')$ otherwise, then $x(\cdot, \lambda, \mathbf{a}, \mathbf{r}', \mathbf{W})$ is a fixed point of $T(\cdot, \lambda, \mathbf{a}, \mathbf{r}', \mathbf{W})$, so that \mathbf{r}' satisfies (6) and is minimal. \square

Remark 5. For the remainder of our analysis, we work with payoff correspondences at which B sub-generates. We shall therefore without loss restrict attention to minimal regimes that satisfy the selection of Lemma 5, i.e., ones for which $\mathbf{r}(s) = R$ whenever the APS gap is positive. This allows us to avoid the computation of non-binding APS payoffs, which would basically entail computing all of $B(\mathbf{W})$, whereas the sign of the APS gap is automatically computed in the process of finding the optimal binding APS payoffs. This is discussed further in Section 4.2.3.

Note that monotonicity of B implies that if B sub-generates at \mathbf{W} , then it will sub-generate at $B(\mathbf{W})$ as well. It does not follow that B will sub-generate at other sub-correspondences of \mathbf{W} . However:

Lemma 6. *If \tilde{B} sub-generates at \mathbf{W} , then B sub-generates at $\tilde{B}(\mathbf{W})$.*

Proof of Lemma 6. Towards a contradiction, suppose that some action profile $a \in \mathbf{A}(\mathbf{W})(s)$, with continuation values $\mathbf{w} \in \tilde{B}(\mathbf{W})$, generates a payoff outside of $\tilde{B}(\mathbf{W})$. Then for some λ ,

$$\begin{aligned} x(s, \lambda, \mathbf{W}) < x^{APS}(a, \lambda, \tilde{B}(\mathbf{W})) &= \lambda \cdot \left((1 - \delta)g(a) + \delta \sum_{s' \in S} \pi(s'|a)\mathbf{w}(s') \right) \\ &\leq (1 - \delta)\lambda \cdot g(a) + \delta \sum_{s' \in S} \pi(s'|a)x(s', \lambda, \mathbf{W}), \end{aligned}$$

where the last inequality holds because $\lambda \cdot \mathbf{w}(s') \leq x(s', \lambda, \mathbf{W})$, since $\mathbf{w}(s') \in \tilde{B}(\mathbf{W})(s')$. The right-hand side of this inequality equals $x^R(a, \lambda, \mathbf{a}, \mathbf{W})$ for any $\mathbf{a} \in \mathbf{A}(\mathbf{W})$ that is maximal in the direction λ (given \mathbf{W}). Since $\tilde{B}(\mathbf{W}) \subseteq \mathbf{W}$, we know that $x^{APS}(s, \lambda, \mathbf{W}) > x(s, \lambda, \mathbf{W})$ as well. That is, $x(s, \lambda, \mathbf{a}, \mathbf{W}) < \min\{x^{APS}(a, \lambda, \mathbf{W}), x^R(a, \lambda, \mathbf{a}, \mathbf{W})\}$, contradicting (L4.ii). \square

This leads to the following result about the sequence generated by \tilde{B} :

Proposition 1 (Sufficiency of binding payoffs). *Let $\tilde{\mathbf{W}}^k$ be the sequence from (T1.iv). Suppose B sub-generates at $\tilde{\mathbf{W}}^0$. Then for any $k \geq 0$, B sub-generates at $\tilde{\mathbf{W}}^k$. Hence, for any λ and $\mathbf{a} \in \mathbf{A}(\tilde{\mathbf{W}}^k)$, there exist minimal regimes \mathbf{r} such that if $\gamma(\mathbf{a}(s), \lambda, \tilde{\mathbf{W}}^k) > 0$, then $\mathbf{r}(s) = R$.*

Proof of Proposition 1. By assumption, $B(\widetilde{\mathbf{W}}^0) \subseteq \widetilde{\mathbf{W}}^0$. Hence, $\widetilde{B}(\widetilde{\mathbf{W}}^0) \subseteq \widetilde{\mathbf{W}}^0$. Since \widetilde{B} is increasing (T1.i), it follows that $\widetilde{B}(\widetilde{\mathbf{W}}^k) \subseteq \widetilde{\mathbf{W}}^k$ for all $k \geq 0$. By Lemma 6, $B(\widetilde{\mathbf{W}}^k) \subseteq \widetilde{\mathbf{W}}^k$ for all $k \geq 0$. The second part of the proposition then follows from Lemma 5. \square

3.4 Computing $x(s, \lambda, \mathbf{W})$

In the next sections, we use the characterization of optimal policies and the sufficiency of binding payoffs to construct simple algorithms for computing \widetilde{B} . These algorithms depend on a subroutine that computes $x(\lambda, \mathbf{W})$. We have already referenced pieces of this routine, but we now synthesize these results into a unified algorithm, together with various simplifications that are possible when B sub-generates at \mathbf{W} .

The computation of $x(\lambda, \mathbf{W})$ consists of an outer routine, in which we maximize over \mathbf{a} , and an inner routine, where we minimize over \mathbf{r} . For the inner routine (Algorithm 1 in Online Appendix A), we use (L2.ii), which says that the regimes \mathbf{r} are not minimal if and only if (6) is violated in some state s , and (L2.iii), which says that $\mathbf{r} \setminus s$ has lower levels in all states. This suggests a simple iterative procedure: Starting from any \mathbf{r} , check for violations of (6). If there are none, then \mathbf{r} is minimal. Otherwise, replace \mathbf{r} with $\mathbf{r} \setminus s$, where s is associated with a violation of (6), and continue. By (L2.iii), The levels decrease at each substitution, so the regimes must converge after at most $|\mathbf{R}|$ substitutions.¹¹

This routine nominally requires us to compute $x^{APS}(a, \lambda, \mathbf{W})$ to check (6). But under the hypothesis that B sub-generates at \mathbf{W} , we can use Lemma 5 to simplify this process: In any state where $\gamma(\mathbf{a}(s), \lambda, \mathbf{W}) > 0$, the minimal regime can be taken to be recursive. Otherwise, by definition $x^{APS} = \widehat{x}^{APS}$, so we can use binding payoffs in (6).

The outer routine (Algorithm 2) is analogous, with actions instead of regimes: Starting from some \mathbf{a} , check for violations of (7). If there are none, then by (L4.ii), \mathbf{a} is maximal. Otherwise, replace \mathbf{a} with $\mathbf{a} \setminus (s, a)$ where (s, a) is associated with a violation, compute new minimal regimes, and continue. By (L4.iii), the levels increase at every substitution, so the actions converge in at most $|\mathbf{A}|$ steps.

As with regimes, we can simplify the computation using Lemma 5: if $\gamma(a, \lambda, \mathbf{W}) > 0$, then Lemma 5 implies that the minimal regime can taken to be recursive. Lemma 3 then implies that $x(s, \lambda, \mathbf{a} \setminus (s, a), \mathbf{W}) > x(s, \lambda, \mathbf{a}, \mathbf{W})$ if and only if $x^R(a, \lambda, \mathbf{a}, \mathbf{W}) > x(s, \lambda, \mathbf{a}, \mathbf{W})$. Otherwise, binding APS payoffs are maximal, and we can replace x^{APS} with \widehat{x}^{APS} in (7).

Thus, we are able to compute an optimal policy and the optimal levels in each direction using only binding payoffs and the sign of the APS gap. We summarize this discussion in

¹¹In fact, since the levels are decreasing, it is not hard to see that if $x^R(\mathbf{a}(s), \lambda, \mathbf{a}, \mathbf{r}, \mathbf{W}) \leq x^{APS}(\mathbf{a}(s), \lambda, \mathbf{W})$, the recursive regime will be minimal at all subsequent rounds. Thus, a switch from APS to R is irreversible, and the maximum number of regime substitutions is $2|S|$.

the following result:

Proposition 2. *Suppose that $\mathbf{A}(\mathbf{W})$ is non-empty valued and B sub-generates at \mathbf{W} . Then the preceding algorithm converges to an optimal policy in at most $|\mathbf{A}||\mathbf{R}|$ steps.*

Remark 6. Given a policy (\mathbf{a}, \mathbf{r}) , there may be multiple action substitutions that lead to higher levels or regime substitutions which lead to lower levels. The arguments above shows that the procedure is order independent, and it will converge to an optimal policy as long as at least one improving substitution is implemented at each round. The pseudocode in Online Appendix A makes a substitution in each state where there is an improvement, although it does not specify how to select from multiple improving substitutions in a given state.

By the discussion preceding Proposition 2, we obtain a refinement of (L2.ii) and (L4.ii), which we record for future reference:

Lemma 7. *Suppose that B sub-generates at \mathbf{W} in the direction λ , and fix $\mathbf{a} \in \mathbf{A}(\mathbf{W})$. The regimes \mathbf{r} are minimal for (λ, \mathbf{a}) if and only if*

$$x(s, \lambda, \mathbf{a}, \mathbf{r}, \mathbf{W}) = \begin{cases} \min \{ \hat{x}^{APS}(\mathbf{a}(s), \lambda, \mathbf{W}), x^R(\mathbf{a}(s), \lambda, \mathbf{a}, \mathbf{r}, \mathbf{W}) \} & \text{if } \gamma(a, \lambda, \mathbf{W}) = 0; \\ x^R(\mathbf{a}(s), \lambda, \mathbf{a}, \mathbf{r}, \mathbf{W}) & \text{if } \gamma(\mathbf{a}(s), \lambda, \mathbf{W}) > 0. \end{cases} \quad (10)$$

Also, \mathbf{a} is optimal if and only if for all $s \in S$ and $a \in \mathbf{A}(\mathbf{W})(s)$,

$$x(s, \lambda, \mathbf{a}, \mathbf{W}) \geq \begin{cases} \min \{ \hat{x}^{APS}(a, \lambda, \mathbf{W}), x^R(a, \lambda, \mathbf{a}, \mathbf{W}) \} & \text{if } \gamma(a, \lambda, \mathbf{W}) = 0; \\ x^R(a, \lambda, \mathbf{a}, \mathbf{W}) & \text{if } \gamma(a, \lambda, \mathbf{W}) > 0. \end{cases} \quad (11)$$

3.5 Optimal payoffs

Our analysis thus far has simplified matters by focusing on optimal levels and suppressing the payoffs, either APS or recursive, that attain these levels. This perspective is useful when computing the bound in a single direction, but it is insufficient for understanding how $x(s, \lambda, \mathbf{W})$ changes as we vary λ . As noted in Remark 3, $x(s, \lambda, \mathbf{W})$ may be non-convex in λ , so that the payoffs that attain the optimal level need not be elements of $\tilde{B}(\mathbf{W})$. Nonetheless, they turn out to be an essential part of our analysis. We re-introduce payoffs and establish basic results that will be used in the following sections. To do so, we will enrich the regimes to include information on which payoff is used when the regime is APS (in which case we presume that incentive constraints bind, per the hypothesis that B sub-generates at \mathbf{W} and Lemma 5).

Given \mathbf{W} and $a \in \mathbf{A}(\mathbf{W})(s)$, we define

$$C(a, \mathbf{W}) = \text{ext} \{v | v \in B(a, \mathbf{W}) \text{ and } (2) \text{ binds for some } i\},$$

where $\text{ext } X$ denotes the set of extreme points of the set X . Thus, for every λ , $\widehat{x}^{APS}(a, \lambda, \mathbf{W}) = \max\{\lambda \cdot v | v \in C(a, \mathbf{W})\}$. For any $\mathbf{a} \in \mathbf{A}(\mathbf{W})$, denote by $\mathbf{P}(\mathbf{a}, \mathbf{W})$ the set of selections $\mathbf{p}(s) \in \{R\} \cup C(\mathbf{a}(s), \mathbf{W})$. Thus, \mathbf{p} encodes whether the regime is recursive or APS, and a choice of binding payoff if the regime is APS. Given a tuple of payoffs \mathbf{u} and $a \in \mathbf{A}(s)$, let

$$u^R(a, \mathbf{u}) = (1 - \delta)g(a) + \delta \sum_{s' \in S} \pi(s'|a)\mathbf{u}(s').$$

Next, we refer to a triple (s, a, p) with $a \in \mathbf{A}(\mathbf{W})(s)$ and $p \in \{R\} \cup C(a, \mathbf{W})$ as a *substitution*. Given a substitution (s, a, p) , let $u(s, a, p, \mathbf{u}) = u^R(a, \mathbf{u})$ if $p = R$, and $u(s, a, p, \mathbf{u}) = p$ if $p \in C(a, \mathbf{W})$. A pair (\mathbf{a}, \mathbf{p}) with $\mathbf{p} \in \mathbf{P}(\mathbf{a}, \mathbf{W})$ induces a tuple of payoffs \mathbf{u} that solve the recursive system

$$\mathbf{u}(s) = u(s, \mathbf{a}(s), \mathbf{p}(s), \mathbf{u}) \tag{12}$$

for all $s \in S$. Note that (12) has a unique solution, by analogous arguments as for (L1.ii).

We now explain how to translate an optimal policy into a corresponding pair and payoffs. Fix λ and $\mathbf{a} \in \mathbf{A}(\mathbf{W})$. We maintain that B sub-generates at \mathbf{W} , so that there exist minimal \mathbf{r} satisfying the selection of Lemma 5, i.e., $\mathbf{r}(s) = APS$ only if $x^{APS}(\mathbf{a}(s), \lambda, \mathbf{W}) = \widehat{x}^{APS}(\mathbf{a}(s), \lambda, \mathbf{W})$. We say that $\mathbf{p} \in \mathbf{P}(\mathbf{a}, \mathbf{W})$ is *min-max for* (λ, \mathbf{a}) if for such minimal regimes \mathbf{r} , $\mathbf{p}(s) = R$ if $\mathbf{r}(s) = R$, and $\lambda \cdot \mathbf{p}(s) = x^{APS}(\mathbf{a}(s), \lambda, \mathbf{W}) = \widehat{x}^{APS}(\mathbf{a}(s), \lambda, \mathbf{W})$ otherwise. In addition, the pair (\mathbf{a}, \mathbf{p}) is *optimal for* λ if \mathbf{a} is maximal for λ and \mathbf{p} is min-max for (λ, \mathbf{a}) . Finally, we say that the payoffs \mathbf{u} are *optimal for* λ if they are induced by a pair that is optimal for λ .

We next establish that optimal pairs exist and attain the optimal level:

Lemma 8. *Fix λ and \mathbf{W} , and suppose that B sub-generates at \mathbf{W} . Then for all $\mathbf{a} \in \mathbf{A}(\mathbf{W})$, there exists a $\mathbf{p} \in \mathbf{P}(\mathbf{a}, \mathbf{W})$ that is min-max for (λ, \mathbf{a}) , and $x(\lambda, \mathbf{a}, \mathbf{W}) = \lambda \cdot \mathbf{u}$ where \mathbf{u} is induced by (\mathbf{a}, \mathbf{p}) . As a result, if $\mathbf{A}(\mathbf{W})$ is non-empty valued, then there exists a pair (\mathbf{a}, \mathbf{p}) that is optimal for λ , and for any payoffs \mathbf{u} that are optimal for λ , $x(s, \lambda, \mathbf{W}) = \lambda \cdot \mathbf{u}$.*

Proof of Lemma 8. The existence of a min-max \mathbf{p} is immediate from the preceding discussion. Let \mathbf{u} denote the payoffs induced by (\mathbf{a}, \mathbf{p}) . From the definition of a min-max \mathbf{p} , it follows directly that if $\mathbf{p}(s) \neq R$, then $\lambda \cdot \mathbf{u}(s) = \widehat{x}^{APS}(\mathbf{a}(s), \lambda, \mathbf{W}) = x(s, \lambda, \mathbf{a}, \mathbf{W})$. And if $\mathbf{r}(s) = \mathbf{p}(s) = R$, then $\mathbf{u}(s) = u^R(\mathbf{a}(s), \mathbf{u})$. A routine calculation shows that $\lambda \cdot \mathbf{u}$ is the fixed

point of $T(\cdot, \lambda, \mathbf{a}, \mathbf{r}, \mathbf{W})$, so that the minimal levels are attained. Finally, when \mathbf{a} are optimal actions, these levels must coincide with $x(s, \lambda, \mathbf{W})$. \square

We next state optimality conditions for a pair (\mathbf{a}, \mathbf{p}) that are analogous to (10) and (11), which will allow us to work directly with pairs rather than policies in computing $\tilde{B}(\mathbf{W})$.

Lemma 9. *Suppose that B sub-generates at \mathbf{W} , and fix λ , $\mathbf{a} \in \mathbf{A}(\mathbf{W})$, and $\mathbf{p} \in \mathbf{P}(\mathbf{a}, \mathbf{W})$, and let \mathbf{u} be induced by (\mathbf{a}, \mathbf{p}) . Then \mathbf{p} is min-max for (λ, \mathbf{a}) if and only if for all s :*

$$\lambda \cdot \mathbf{u}(s) = \min \left\{ \max \{ \lambda \cdot v \mid v \in C(\mathbf{a}(s), \mathbf{W}) \}, \lambda \cdot u^R(\mathbf{a}(s), \mathbf{u}) \right\}$$

if $\gamma(\mathbf{a}(s), \lambda, \mathbf{W}) = 0$, and $\mathbf{p}(s) = R$ if $\gamma(\mathbf{a}(s), \lambda, \mathbf{W}) > 0$ (so that $\mathbf{u}(s) = u^R(\mathbf{a}(s), \mathbf{u})$). Moreover, (\mathbf{a}, \mathbf{p}) is optimal for λ if and only if: \mathbf{p} is min-max and for all s and $a \in \mathbf{A}(\mathbf{W})(s)$,

$$\lambda \cdot \mathbf{u}(s) \geq \begin{cases} \min \{ \max \{ \lambda \cdot v \mid v \in C(a, \mathbf{W}) \}, \lambda \cdot u^R(a, \mathbf{u}) \} & \text{if } \gamma(a, \lambda, \mathbf{W}) = 0; \\ \lambda \cdot u^R(a, \mathbf{u}) & \text{if } \gamma(a, \lambda, \mathbf{W}) > 0. \end{cases}$$

Proof of Lemma 9. This follows from Lemmas 7 and 8, with the observation that $x(s, \lambda, \mathbf{a}) = \lambda \cdot \mathbf{u}(s)$ for the payoffs \mathbf{u} induced by (\mathbf{a}, \mathbf{p}) , and that $x^R(a, \lambda, \mathbf{a}) = \lambda \cdot u^R(a, \mathbf{u})$. \square

Remark 7. The algorithm of Proposition 2 is easily adapted to directly compute an optimal pair. Each policy (\mathbf{a}, \mathbf{r}) satisfying the selection of Lemma 5 can be identified with an equivalent pair (\mathbf{a}, \mathbf{p}) , where $\mathbf{p}(s) = R$ if $\mathbf{r}(s) = R$ and otherwise $\mathbf{p}(s)$ is a highest binding payoff for $\mathbf{a}(s)$. It is immediate that the induced payoffs \mathbf{u} satisfy $\lambda \cdot \mathbf{u}(s) = x(s, \lambda, \mathbf{a}, \mathbf{r}, \mathbf{W})$ for all s . In iterating over pairs, just as with policies, there is an inner min-maximization over \mathbf{p} and an outer maximization over \mathbf{a} . The computation of the APS gap remains the same (and we will comment below on how to do this efficiently). But when comparing the recursive and binding APS levels to test for an improvement, instead of using $x^R(a, \lambda, \mathbf{a}, \mathbf{r}, \mathbf{W})$ and $\hat{x}^{APS}(a, \lambda, \mathbf{W})$, we use $\lambda \cdot u^R(a, \mathbf{u})$ and $\lambda \cdot v$, where v is a highest binding payoff. Starting from (\mathbf{a}, \mathbf{p}) where $\mathbf{p}(s)$ is recursive or maximal in $C(a, \mathbf{W})$, the computed sequence of pairs corresponds to a sequence of policies generated by the algorithm of Proposition 2, as long as ties are broken the same way (cf. Remark 6).

4 Two players:

Further implications and implementation

We now specialize to two-player games which, as we shall see, have an especially simple structure: the number of extreme equilibrium payoffs is bounded, \tilde{B} has bounded computational

complexity, and it can be computed via a simple procedure.

4.1 The complexity of \mathbf{V}

We first establish a bound on the number of optimal pairs.

Lemma 10. *If $N = 2$ and \mathbf{W} is convex valued, then $|C(a, \mathbf{W})| \leq 4$, and the number of pairs is at most*

$$\bar{L} = 5^{|S|} \times_{s \in S} |\mathbf{A}(s)|. \quad (13)$$

Proof of Lemma 10. Recall the definition of $C(a, \mathbf{W})$ and condition (2). If w is a feasible and incentive compatible expected continuation value for which player i 's incentive constraint binds, then $w_i = \underline{u}_i(a, \mathbf{W})$. The set of such points is either empty, a singleton, or it has at most two extreme points. Thus, $|C(a, \mathbf{W})| \leq 4$. Finally, each optimal payoff tuple is associated with some $\mathbf{a} \in \mathbf{A}(\mathbf{W})$ and $\mathbf{p} \in \mathbf{P}(\mathbf{a}, \mathbf{W})$. But by the preceding argument, $|C(\mathbf{a}(s), \mathbf{W})| \leq 4$, so $|\mathbf{P}(\mathbf{a}, \mathbf{W})| \leq 5^{|S|}$. The bound immediately follows. \square

We immediately obtain a bound on the complexity of \mathbf{V} , which generalizes an analogous result of Abreu and Sannikov (2014) for repeated games:

Corollary 1. *If $N = 2$, then \mathbf{V} has at most \bar{L} extreme points.*

Proof of Corollary 1. The proof of (T1.iii) shows that $x(s, \lambda, \mathbf{V})$ is the support function of \mathbf{V} . Since there are at most \bar{L} pairs, $x(s, \lambda, \mathbf{V})$ can have at most \bar{L} linear segments. As a result, \mathbf{V} has at most \bar{L} extreme points. \square

Note that Lemma 10 does not immediately yield a bound on the complexity of $\widetilde{\mathbf{W}}^k$. The fact that these correspondences have bounded complexity will, however, be proven using the computational procedure we describe next.

4.2 Implementation with two players

Under the hypothesis that B sub-generates at \mathbf{W} , there is a simple procedure to compute $\widetilde{B}(\mathbf{W})$, which consists of iterative application of two subroutines. Starting from an optimal pair for some initial direction, we compute a set of clockwise rotations of the direction for which that pair remains optimal. We then rotate the direction of optimization clockwise as far as we can, subject to the incumbent pair remaining optimal, and compute the pair that would become optimal if the direction were to rotate further by a small amount. We then continue iteratively: compute a new range of directions, rotate, and re-optimize. The main

result for this section, Theorem 2, shows that this procedure maps out the entire frontier of $\tilde{B}(\mathbf{W})$ in a bounded number of steps.

The key step in the algorithm is to compute a range of directions for which an optimal pair remains optimal. We will show that a suitably chosen optimal pair can only cease to be optimal at what we call a *test direction*. These directions are identified with the substitutions of either actions or regimes that are used in the optimization routines of Proposition 2 and Remark 7. Specifically, the test directions corresponding to a substitution are the critical directions at which that substitution would leave the level unchanged. Proposition 3 below establishes two fundamental results: First, except when $\tilde{B}(\mathbf{W})$ is degenerate, a test direction always exists. Second, the incumbent optimal pair remains optimal for all directions between the initial direction and the “shallowest” test direction, i.e., the test direction with the smallest clockwise angle of rotation from the initial direction.

This high-level summary neglects two complications that are dealt with in the following pages. First, we will impose a refinement on the test directions considered by our algorithm that we call *legitimacy*. This condition is used to weed out some spurious test directions at which the optimal pair would not change. The legitimacy condition can, however, be dropped without affecting our results on convergence and computational complexity.

A more substantive issue is what is meant by a “suitably chosen optimal pair.” The test direction methodology requires us to start from an optimal pair that is *robust*, in the sense that it remains optimal for small clockwise perturbations of the direction. Proposition 4 characterizes a procedure called *lexicographic optimization*, which is guaranteed to produce a robustly optimal pair for any direction.

Thus, Propositions 3 and 4 establish three results that are used to prove Theorem 2: Given a robustly optimal pair, there exists a legitimate test direction, and the incumbent pair remains optimal to the next shallowest legitimate test direction (Proposition 3). Moreover, for any direction, there exists a robustly optimal pair and a method to compute it (Proposition 4). As a result, iterative application of direction rotation and optimization will necessarily identify an optimal pair for every direction.

4.2.1 Rotating the direction clockwise

We now proceed formally. For rotating the direction, we will work with pairs that satisfy a mild refinement: A pair (\mathbf{a}, \mathbf{p}) with induced payoffs \mathbf{u} is *canonical* if $\mathbf{p}(s) = R$ whenever $u^R(\mathbf{a}(s), \mathbf{u}) \in C(\mathbf{a}, \mathbf{W})$. In other words, if a binding payoff is exactly equal to the recursive payoff, then we break ties in favor of the recursive regime. In this situation, recursive is “canonically” minimal in the sense that it generates the lowest level in all directions, as the recursive payoff is also an APS payoff. (Note, however, that \mathbf{p} in a canonical pair

need not be min-max, nor does the definition depend on the direction of optimization). It is immediate that for any optimal pair, there is another optimal pair with the same payoffs that is canonical, which is obtained by switching the regime to \mathbf{r} whenever $\mathbf{p}(s) = u^R(\mathbf{a}(s), \mathbf{u})$:

Lemma 11. *If (\mathbf{a}, \mathbf{p}) is optimal for λ and induces \mathbf{u} , then there exists a $\mathbf{p}' \in \mathbf{P}(\mathbf{a}, \mathbf{W})$ such that $(\mathbf{a}, \mathbf{p}')$ is optimal for λ , induces \mathbf{u} , and is canonical.*

Now, fix a payoff tuple \mathbf{u} and substitution (s, a, p) , where $a \in \mathbf{A}(\mathbf{W})(s)$ and $p \in \{R\} \cup C(a, \mathbf{W})$. Recall that $u(s, a, p, \mathbf{u})$ is equal to p if $p \in C(a, \mathbf{W})$ and is equal to $u^R(a, \mathbf{u})$ if $p = R$. We say that λ' is *test direction* for (s, a, p) at \mathbf{u} if $u(s, a, p, \mathbf{u}) \neq \mathbf{u}(s)$ and

$$\lambda' \cdot (u(s, a, p, \mathbf{u}) - \mathbf{u}(s)) = 0. \quad (14)$$

In other words, λ' is normal to the direction in which the substitution moves payoffs. Furthermore, we say that (s, a, p) and a corresponding test direction λ' are *legitimate* if

$$\lambda' \cdot u(s, a, p, \mathbf{u}) \leq \min \{ \lambda' \cdot u^R(s, a, p, \mathbf{u}), x^{APS}(a, \lambda, \mathbf{W}) \}.$$

We note for future reference that the number of test directions is bounded. For there are at most \bar{L} pairs, and the number of substitutions is at most

$$\bar{M} = 5 \sum_{s \in S} |\mathbf{A}(s)|.$$

Moreover, each pair and substitution has at most two associated test directions that solve (14). As a result, there are at most $2\bar{L}\bar{M}$ test directions.

Next, let us define $[\lambda, \lambda']$ to be the *closed arc* of directions obtained by moving clockwise from λ to λ' . We extend this convention to open and half-closed arcs in the obvious way. We say that a pair (\mathbf{a}, \mathbf{p}) is *robustly optimal* at λ if there exists $\lambda' \neq \lambda$ such that (\mathbf{a}, \mathbf{p}) is optimal for all directions in $[\lambda, \lambda']$. A payoff \mathbf{u} is *robustly optimal* if it is induced by a robustly optimal pair.

The following proposition characterizes the search for shallowest legitimate test directions and their relationship with robustly optimal pairs.

Proposition 3 (Test directions). *Suppose that $N = 2$ and that B sub-generates at \mathbf{W} . Suppose further that (\mathbf{a}, \mathbf{p}) is robustly optimal at λ and induces \mathbf{u} . Then either \mathbf{u} is optimal for all directions, or there exists a legitimate test direction at \mathbf{u} that is not equal to λ . Moreover, if $\lambda' \neq \lambda$ is the shallowest test direction and (\mathbf{a}, \mathbf{p}) is canonical, then (\mathbf{a}, \mathbf{p}) is optimal for all directions in $[\lambda, \lambda']$.*

We now prove Proposition 3. Note that the remaining results of this section all maintain the implicit hypotheses that B sub-generates at \mathbf{W} , so that binding payoffs are sufficient, and that $N = 2$. We first record a simple technical result.

Lemma 12. *For each (\mathbf{a}, \mathbf{p}) , the set of directions in which \mathbf{p} is min-max for \mathbf{a} is closed, and the set of directions in which (\mathbf{a}, \mathbf{p}) is optimal is closed. For each \mathbf{u} , the set of directions in which \mathbf{u} is optimal is closed.*

Proof of Lemma 12. This is a consequence of Lemma 9 and the fact that γ , \widehat{x}^{APS} , and x^R are all continuous in the direction. As a result, if (\mathbf{a}, \mathbf{p}) satisfies either the min-max or optimality conditions of Lemma 9 along a convergent sequence of directions, then it also satisfies them in the limit. Finally, \mathbf{u} is optimal on the union of the finitely many closed sets for which pairs that induce \mathbf{u} are optimal. \square

The next two lemmas establish the first part of Proposition 3.

Lemma 13. *If \mathbf{u} and $\mathbf{u}' \neq \mathbf{u}$ are both optimal at some direction λ' , then λ' is a legitimate test direction at \mathbf{u} .*

Proof of Lemma 13. Let $(\mathbf{a}', \mathbf{p}')$ be an optimal pair in the direction λ' that induces \mathbf{u}' . It must be that $u(s, \mathbf{a}'(s), \mathbf{p}'(s), \mathbf{u}) \neq \mathbf{u}(s)$ for some s . Otherwise, \mathbf{u} is a solution to (12) for $(\mathbf{a}', \mathbf{p}')$, and uniqueness of the solution would imply that $\mathbf{u} = \mathbf{u}'$, a contradiction. As a result, λ' is a test direction for the substitution $(s, \mathbf{a}'(s), \mathbf{p}'(s))$. Moreover, $\lambda' \cdot u(s, \mathbf{a}'(s), \mathbf{p}'(s), \mathbf{u}) = \lambda' \cdot \mathbf{u}'(s)$ for all s , so that legitimacy follows from the optimality conditions for $(\mathbf{a}', \mathbf{p}')$. \square

Lemma 14. *Suppose that \mathbf{u} is robustly optimal at λ and is not optimal at $\widehat{\lambda} \neq \lambda$. Then there is a legitimate test direction in $(\lambda, \widehat{\lambda})$.*

Proof of Lemma 14. By Lemma 12, the set of directions at which \mathbf{u} is optimal is closed, so that there is a largest closed arc $[\lambda, \lambda']$ on which \mathbf{u} is optimal. We will show that there exists $\mathbf{u}' \neq \mathbf{u}$ that is optimal at λ' , so the result follows from Lemma 13.

By assumption, \mathbf{u} is not optimal at $\widehat{\lambda}$, and since \mathbf{u} is robustly optimal at λ , we conclude that $\lambda' \in (\lambda, \widehat{\lambda})$. The definition of λ' then implies that there is a sequence of directions in $(\lambda', \widehat{\lambda}]$ converging to λ' at which \mathbf{u} is not optimal. Since there is an optimal pair for every direction (Lemma 8) and only finitely many pairs (Lemma 10), there must be a pair $(\mathbf{a}', \mathbf{p}')$ which induces payoffs $\mathbf{u}' \neq \mathbf{u}$, such that $(\mathbf{a}', \mathbf{p}')$ is optimal for directions arbitrarily close to λ' . Lemma 12 then implies that $(\mathbf{a}', \mathbf{p}')$ and \mathbf{u}' are also optimal at λ' . \square

The next lemma is used to prove the second part of Proposition 3.

Lemma 15. *Suppose that (\mathbf{a}, \mathbf{p}) is robustly optimal at λ and induces \mathbf{u} , and that $\lambda' \neq \lambda$ is the shallowest legitimate test direction at \mathbf{u} . If (\mathbf{a}, \mathbf{p}) is canonical, then it is optimal for all directions in $[\lambda, \lambda']$.*

Proof of Lemma 15. The proof consists of three steps.

Step 1: \mathbf{u} is optimal for all directions in (λ, λ') . If not, then there is a $\hat{\lambda} \in (\lambda, \lambda')$ at which it is not optimal, and Lemma 14 then implies that there is a legitimate test direction in $(\lambda, \hat{\lambda})$, which contradicts the hypothesis that λ' is shallowest.

Step 2: If (\mathbf{a}, \mathbf{p}) is canonical and is optimal at $\hat{\lambda} \in (\lambda, \lambda')$, then there is a closed neighborhood of $\hat{\lambda}$ on which (\mathbf{a}, \mathbf{p}) is optimal.

By Step 1 and the fact that (\mathbf{a}, \mathbf{p}) induces \mathbf{u} , (\mathbf{a}, \mathbf{p}) is optimal in a neighborhood of $\hat{\lambda}$ if and only if \mathbf{p} is min-max for \mathbf{a} in this neighborhood. The following two cases establish that for every s , there is a neighborhood of $\hat{\lambda}$ on which the min-max condition in Lemma 9 is satisfied. Thus, \mathbf{p} is min-max on the intersection of these finitely many neighborhoods, which, together with Lemma 12, implies the result.

Case 1: $\mathbf{p}(s) \neq R$. Then for all $v \in C(\mathbf{a}(s), \mathbf{W}) \setminus \{\mathbf{p}(s)\}$, it must be that $\hat{\lambda} \cdot v < \hat{\lambda} \cdot \mathbf{p}(s)$. Otherwise $\mathbf{p}(s)$ would not be maximal or, if there is a tie for maximal payoff, then Lemma 13 implies that $\hat{\lambda}$ is a legitimate test direction, contradicting the hypothesis that λ' is shallowest. Thus, there is a neighborhood of $\hat{\lambda}$ for which $\mathbf{p}(s)$ is the highest binding APS payoff. Next, since (\mathbf{a}, \mathbf{p}) is canonical, it must be that $u^R(\mathbf{a}(s), \mathbf{u}) \neq \mathbf{p}(s)$.¹² Moreover, it must be that $\hat{\lambda} \cdot u^R(\mathbf{a}(s), \mathbf{u}) > \hat{\lambda} \cdot \mathbf{p}(s)$, for otherwise $\hat{\lambda}$ would again be a legitimate test direction. Finally, if there are directions arbitrarily close to $\hat{\lambda}$ for which the APS gap is positive, then Lemma 5 implies that $\hat{\lambda} \cdot u^R(\mathbf{a}(s), \mathbf{u}) \leq x^{APS}(\hat{\lambda}, \mathbf{W})$ for $\hat{\lambda}$ arbitrarily close to $\hat{\lambda}$. Continuity of x^{APS} then implies that $\hat{\lambda} \cdot u^R(\mathbf{a}(s), \mathbf{u}) \leq x^{APS}(\hat{\lambda}, \mathbf{W}) = \hat{\lambda} \cdot \mathbf{p}(s)$, again a contradiction. We conclude that there is a neighborhood of $\hat{\lambda}$ for which $u^R(\mathbf{a}(s), \mathbf{u})$ is strictly above $\mathbf{p}(s)$.

Case 2: $\mathbf{p}(s) = R$. If $\gamma(\mathbf{a}(s), \hat{\lambda}, \mathbf{W}) > 0$, then continuity of the APS gap implies that $\mathbf{p}(s) = R$ is minimal for a neighborhood of $\hat{\lambda}$. Otherwise, the APS gap is zero and there is a maximal APS payoff at $\hat{\lambda}$ that is binding. If $\hat{x}^{APS}(\hat{\lambda}, \mathbf{W}) > \hat{\lambda} \cdot \mathbf{u}(s)$, then again, continuity implies that the recursive regime is minimal in a neighborhood of $\hat{\lambda}$. If $\hat{x}^{APS}(\hat{\lambda}, \mathbf{W}) = \hat{\lambda} \cdot \mathbf{u}(s)$, then $\hat{\lambda} \cdot v = \hat{\lambda} \cdot \mathbf{u}(s)$ for some maximal $v \in C(\mathbf{a}(s), \mathbf{W})$. If $\mathbf{u}(s) \neq v$, then $\hat{\lambda}$ is a legitimate test direction corresponding to the substitution $(s, \mathbf{a}(s), v)$, which contradicts λ' being shallowest. Otherwise, it must be that $\mathbf{u}(s) = v$. As a result, the recursive payoff is also an APS payoff,

¹²This is the only step in the argument that uses the hypothesis (\mathbf{a}, \mathbf{p}) is canonical. Without this hypothesis, it could be that $\mathbf{p}(s) \neq R$ and $\mathbf{p}(s) = \mathbf{u}(s) = u^R(\mathbf{a}(s), \mathbf{u})$. Hence, the recursive regime is minimal for all directions, but the recursive and APS regimes happen to tie at $\hat{\lambda}$. As the direction rotates, the recursive regime may become *uniquely* minimal, because the APS gap becomes positive. This need not happen at a test direction. Instead of selecting canonical pairs, we could have included additional test directions when the APS gap switches sign. These are the normals to the non-binding APS frontier at a binding payoff.

so that $\mathbf{p}(s) = R$ is min-max for \mathbf{a} in all directions.

Step 3: Let $[\underline{\lambda}, \bar{\lambda}]$ be a largest closed arc in $[\lambda, \lambda']$ on which (\mathbf{a}, \mathbf{p}) is optimal, which by hypothesis is non-empty. If $\bar{\lambda} \neq \lambda'$, then Step 2 with $\hat{\lambda} = \bar{\lambda}$ implies that there is a closed neighborhood of $\bar{\lambda}$, denoted U , on which (\mathbf{a}, \mathbf{p}) is optimal. Then (\mathbf{a}, \mathbf{p}) is optimal on $[\underline{\lambda}, \bar{\lambda}] \cup U$, which is a strict superset of $[\underline{\lambda}, \bar{\lambda}]$, a contradiction. We similarly conclude that $\underline{\lambda} = \lambda$, so that (\mathbf{a}, \mathbf{p}) is optimal for all directions in $[\lambda, \lambda']$. \square

Proof of Proposition 3. The proposition follows directly from Lemmas 14 and 15. \square

4.2.2 Finding a robustly optimal pair

The last task is to find a robustly optimal pair and payoffs in a direction λ . This is accomplished with a procedure that we refer to as *lexicographic optimization*: Starting from (\mathbf{a}, \mathbf{p}) , we optimize according to the procedure of Proposition 2, using payoffs as described in Remark 7, except that in ranking any pair of payoffs v and v' , we use the *lexicographic ordering*, whereby a payoff v is greater than v' if $\lambda \cdot v > \lambda \cdot v'$ or if $\lambda \cdot v = \lambda \cdot v'$ and $\tilde{\lambda} \cdot v > \tilde{\lambda} \cdot v'$, where $\tilde{\lambda}$ is equal to λ rotated 90 degrees clockwise. In this case, we write $v >_{\lambda} v'$. Note that the use of the lexicographic order only affects how ties are broken in the algorithm (cf. Remark 6). In particular, instead of using any highest payoff in $C(a, \mathbf{W})$, we use the lexicographically highest, and when the APS and recursive payoffs are tied, we select whichever is lexicographically minimal. In addition, we break ties in favor of the recursive regime if the APS gap is zero but would become strictly positive for small clockwise rotations: In particular, wherever the condition $\gamma(a, \lambda, \mathbf{W}) > 0$ was used previously, we now use the condition that there exists a $\lambda' \neq \lambda$ such that $\gamma(a, \lambda', \mathbf{W}) > 0$ for all $\lambda'' \in (\lambda, \lambda']$. In this case, we say that the APS gap for a is *lexicographically positive* at λ' . (Note that continuity of γ in λ implies that the APS gap is lexicographically positive whenever it is positive.) This procedure is fully described in Algorithms 4 and 5 in Online Appendix A.

The following proposition characterizes lexicographic optimization and, as a corollary, shows that a robustly optimal pair exists in every direction:

Proposition 4. *Suppose that $N = 2$, $\mathbf{A}(\mathbf{W})$ is non-empty valued, and B sub-generates at \mathbf{W} . Then lexicographic optimization in a direction λ terminates in finitely many steps at a pair (\mathbf{a}, \mathbf{p}) that is robustly optimal at λ . Moreover, once the algorithm has reached an optimal pair for λ , every subsequent pair is also optimal for λ .*

Proof of Proposition 4. Since lexicographic optimization only affects how ties are broken, convergence of the algorithm described in Remark 7 implies that lexicographic optimization will reach a pair that is optimal for λ . Once such a pair is reached, there are no substitutions

that strictly improve in the direction λ . Thus, any substitution considered by lexicographic optimization will keep the λ level the same and move payoffs in the direction $\tilde{\lambda}$, which is λ rotated 90 degrees clockwise. A straightforward adaptation of the argument for Proposition 2 shows that lexicographic regime minimization produces a sequence of payoffs that monotonically decrease in the direction $\tilde{\lambda}$, and lexicographic action maximization produces payoffs that monotonically increase in the direction $\tilde{\lambda}$, so that no pair is repeated and the sequence converges after finitely many steps to a limit (\mathbf{a}, \mathbf{p}) .

We claim that (\mathbf{a}, \mathbf{p}) is robustly optimal at λ . An important fact used below is that $v >_{\lambda} v'$ if and only if there exists $\lambda' \neq \lambda$ such that $\lambda'' \cdot v > \lambda'' \cdot v'$ for all $\lambda'' \in (\lambda, \lambda']$. Now, we will argue that (\mathbf{a}, \mathbf{p}) satisfies the optimality conditions of Lemma 9 for small clockwise rotations from λ . Since lexicographic optimization has converged, we know that the analogue of these conditions is satisfied at λ , where we select the lexicographically lowest of the recursive payoff and the lexicographically highest binding payoff, and we select the recursive payoff if the APS gap is lexicographically positive at λ . We shall argue that this implies that the conditions in Lemma 9 are satisfied for small clockwise rotations.

Let us first consider the min-max conditions. If the APS gap is lexicographically positive for $\mathbf{a}(s)$, then $\mathbf{u}(s) = u^R(\mathbf{a}(s), \mathbf{u})$, and, moreover, the APS gap is positive for some arc $(\lambda, \lambda']$ with $\lambda' \neq \lambda$, so that $\mathbf{p}(s) = R$ is min-max on $(\lambda, \lambda']$. If the APS gap is not lexicographically positive, then there is an arc $(\lambda, \lambda'']$ over which (i) the APS gap is zero, (ii) the payoff v that is lexicographically highest at λ remains highest, and (iii) the ranking between v and $u^R(\mathbf{a}(s), \mathbf{u})$ is strict if and only if there is the same strict lexicographic ranking at λ . The fact that we have selected the lexicographically minimal of v and $u^R(\mathbf{a}(s), \mathbf{u})$ then implies that the Lemma 9 min-max conditions are satisfied on $(\lambda, \lambda']$.

The analysis for optimality is entirely analogous. We conclude that for every regime or action substitution, there is a non-trivial clockwise arc over which the Lemma 9 conditions are satisfied. Since there are finitely many substitutions, there exists a non-trivial clockwise arc over which (\mathbf{a}, \mathbf{p}) is optimal, so that (\mathbf{a}, \mathbf{p}) is robustly optimal. \square

4.2.3 Computing $\tilde{B}(\mathbf{W})$

We now summarize the computation of $\tilde{B}(\mathbf{W})$ with two players, assuming that B subgenerates at \mathbf{W} and there is a supportable action profile in each state (otherwise $\tilde{B}(\mathbf{W})(s)$ is empty for some state, so that there are no pure-strategy subgame-perfect equilibria). A preliminary step is to compute, for each action profile, whether it is supportable, the sets $C(a, \mathbf{W})$, and the APS gap. This is done as follows. To compute $C(a, \mathbf{W})$, we simply intersect each of the binding incentive rays with the half spaces that define \mathbf{W} , and a is supportable

if and only if $C(a, \mathbf{W})$ is non-empty (cf. Footnote 10).¹³ In addition, for each $v \in C(a, \mathbf{W})$, we record the direction on the frontier of $B(a, \mathbf{W})$ that points into the incentive compatible region, denoted $d(v)$. The APS gap is positive at λ if there is a $v \in \arg \max_{v \in C(a, \mathbf{W})} \lambda \cdot v'$ such that $\lambda \cdot d(v) > 0$, and it is lexicographically positive if $\lambda \cdot d(v) > 0$ or $d(v) = \alpha \tilde{\lambda}$ for some $\alpha > 0$, where $\tilde{\lambda}$ is equal to λ rotated 90 degrees clockwise.

After these preliminary calculations, we pick an initial direction $\lambda^0 \in \Lambda$, at which we compute a robustly optimal pair $(\mathbf{a}^0, \mathbf{p}^0)$ and its payoffs \mathbf{u}^0 . If there are no legitimate test directions at \mathbf{u}^0 , then we stop. Otherwise, proceeding inductively from \mathbf{u}^{k-1} for $k \geq 1$, we set λ^k equal to the shallowest legitimate test direction from \mathbf{u}^{k-1} and λ^{k-1} . We then compute a robustly optimal pair $(\mathbf{a}^k, \mathbf{p}^k)$ and corresponding robustly optimal payoffs \mathbf{u}^k for λ^k , via lexicographic optimization starting from $(\mathbf{a}^{k-1}, \mathbf{p}^{k-1})$. We stop when the direction of optimization passes λ^0 , at step K . This procedure is described in Algorithm 6.

Theorem 2. *Suppose that $N = 2$, $\mathbf{A}(\mathbf{W})$ is non-empty valued, and B sub-generates at \mathbf{W} . Then the previously described procedure terminates in at most $2\bar{L}\bar{M}$ substitutions and runtime $O(\bar{L}\bar{M}^2)$. If there are no legitimate test directions at \mathbf{u}^0 , then $\tilde{B}(\mathbf{W})(s) = \{\mathbf{u}^0(s)\}$ for all s . Otherwise,*

$$\tilde{B}(\mathbf{W})(s) = \{v | \lambda^k \cdot v \leq \lambda^k \cdot \mathbf{u}^k(s) \ \forall k = 1, \dots, K\}. \quad (15)$$

Proof of Theorem 2. As there are supportable actions in every state, convergence to the initial robustly optimal pair is guaranteed by Proposition 4. If there are no legitimate test directions, Proposition 3 implies \mathbf{u}^0 is optimal in every direction and is equal to $\tilde{B}(\mathbf{W})$. Otherwise, there are legitimate test directions at every step of the algorithm. Proposition 3 implies that \mathbf{u}^{k-1} is optimal on $[\lambda^{k-1}, \lambda^k]$, so that $\tilde{B}(\mathbf{W})$ is equal to right-hand side of (32).

We next argue the complexity bound. The computation of $(\mathbf{a}^0, \mathbf{p}^0)$ takes at most \bar{L} substitutions. In addition, there are at most $2\bar{L}\bar{M}$ test directions in which we optimize, and to compute the shallowest legitimate substitution requires the consideration of \bar{M} substitutions. It remains to bound the number of substitutions that are actually made. If (s, a, p) is substituted into $(\mathbf{a}^k, \mathbf{p}^k)$, then both $(\mathbf{a}^k, \mathbf{p}^k)$ and $(\mathbf{a}^k, \mathbf{p}^k) \setminus (s, a, p)$ are optimal in a direction λ^k . As a result, λ^k is one of the two test directions satisfying (14). Since the direction of optimization rotates monotonically clockwise, it follows that (s, a, p) can be substituted into $(\mathbf{a}^k, \mathbf{p}^k)$ at most twice over the course of the algorithm. Thus, the total number of substitutions (and hence K) is at most $2\bar{L}\bar{M}$. Hence, the total runtime is $O(\bar{L}\bar{M}^2)$. \square

¹³It is easy to see that $\mathbf{A}(\mathbf{W})$ is decreasing in \mathbf{W} , so that when \tilde{B} is applied iteratively to generate the sequence $\tilde{\mathbf{W}}^k$ from (T1.iv), $C(a, \mathbf{W})$ need only be computed for action profiles in $\mathbf{A}(\tilde{\mathbf{W}}^{k-1})$.

Remark 8. We have not ruled out the possibility that there are more than \bar{L} bounding hyperplanes. This is because $x(s, \lambda, \mathbf{W})$ need not be convex, and the algorithm may come back to a pair that was previously found to be optimal. We note, however, that the run-time of the algorithm is sensitive to the actual number of bounding hyperplanes, which in practice we have found to be much smaller than \bar{L} .

4.2.4 Further improvements and single-state substitutions

The algorithm characterized by Section 4.2.3 nominally considers all substitutions when it searches for the next robust optimum. Proposition 4 shows that we can without loss restrict attention to substitutions which leave the pair optimal in the shallowest test direction, speeding up computation. It is possible to go even further, by restricting attention to test directions for which the associated substitutions which would become strict improvements for small clockwise rotations, using a lexicographic version of the legitimacy test. This is done in our software implementation.

In addition, there are many cases where we can find the next robustly optimal pair directly in a single step, as we now explain. If a “shallowest substitution” (that is, one that generates a shallowest test direction) entails a new action or a switch to or from a recursive regime, then frequently it will be unique, as ties across states or actions are exceptional. It is then trivial to check if the substitution generates the next robustly optimal pair, or whether the incumbent pair continues to be optimal at λ' .

On the other hand, if all shallowest substitutions entail a change from one binding APS payoff to another, then the next robustly optimal pair is obtained by switching all of the binding payoffs to the ones that are lexicographically optimal. For example, if the shallowest direction is $\lambda' = (-1, 0)$, which corresponds to punishing player 1, then payoffs may jump from minimizing to maximizing player 2’s payoff, subject to minimizing player 1’s payoff. This and the analogous situation for $\lambda' = (0, -1)$ occur frequently in our simulations.

Our code does not currently implement all of these simplifications. We view the software as a living entity which will continue to be improved by ourselves and the community. Such improvements can only lead to better performance than that reported in the next section.

4.3 Examples

We have implemented this algorithm in a software package called SGSolve, which is available through an author’s website and Github, under the terms of the GPLv3 license. The code consists of routines that implement the algorithm and graphical interface for specifying games and visualizing solutions. The following two examples were solved using this package. An

a_1/a_2	D	C
C	$(-1, 5)$	$(4, 2)$
D	$(0, 0)$	$(5, -1)$

Figure 1: An asymmetric Prisoners' Dilemma.

additional two-player example is reported in Online Appendix C.

4.3.1 Asymmetric Prisoners' Dilemma

Our first example is the asymmetric repeated Prisoners' Dilemma of Figure 1, with $\delta = 1/2$. As there is a single state, we temporarily drop the argument s and use regular font weight.

The computation is depicted in Figure 2. The flow payoffs are four circles. We take \widetilde{W}^0 to be the feasible and individually rational payoffs, which are in gray.¹⁴ In the center panel, the shaded polygons are the sets $B(a, \widetilde{W}^0)$ generated by the APS operator. For each a and λ , $x^{APS}(a, \lambda, \widetilde{W}^0)$ is the level of the highest payoff in $B(a, \widetilde{W}^0)$, whereas $x(a, \lambda, \widetilde{W}^0)$ is the minimum of $x^{APS}(a, \lambda, \widetilde{W}^0)$ and the recursive level, which in a repeated game is just the level of the flow payoff.

The initial direction of optimization is $\lambda^0 = (0, 1)$. For APS, the optimal level is $x^{APS}((C, D), \lambda^0, \widetilde{W}^0)$, which is attained at the top left corner of the red set. This is also the max-min-max level, since $g(C, D) = (-1, 5)$ is strictly higher in the direction λ^0 . Note that player 1's incentive constraint binds at the optimum, consistent with Proposition 1.

Rotating clockwise from λ^0 , the algorithm computes a sequence of test directions and robustly optimal pairs. Between λ^0 and λ^1 , $((C, D), APS)$ remains optimal. At λ^1 , both $((C, D), R)$ is tied with $((C, C), APS)$, but the latter optimum is robust. At λ^2 , $((D, C), R)$ becomes robustly optimal, as $g(D, C)$ is still below the best APS payoff for (D, C) . Between λ^3 and λ^4 , the robust optimum is $((D, C), APS)$. At λ^4 , the optimal switches to $((D, D), R)$, which remains optimal until λ^5 , at which point $((C, D), APS)$ is again optimal.

The corresponding hyperplanes for these directions and levels are then intersected to form \widetilde{W}^1 . It is not hard to see that at the second round, the exact same set will be generated, i.e., $\widetilde{B}(\widetilde{W}^1) = \widetilde{W}^1$. For even though the binding payoff that was used between λ^2 and λ^3 is no longer available, the half spaces in this direction were redundant anyway. Thus, the operator converges after exactly one round.

In contrast, the APS, JYC, and Abreu and Sannikov (2014) operators would all cut less in the directions between λ^2 and λ^3 . For APS, this is because the best APS payoffs for (D, C) are higher than the flow payoff. The operator of Abreu and Sannikov (2014) would

¹⁴If we started with \widetilde{W}^0 equal to the feasible set, then our algorithm would only compute the equilibrium threats asymptotically, and the algorithm would only converge asymptotically.

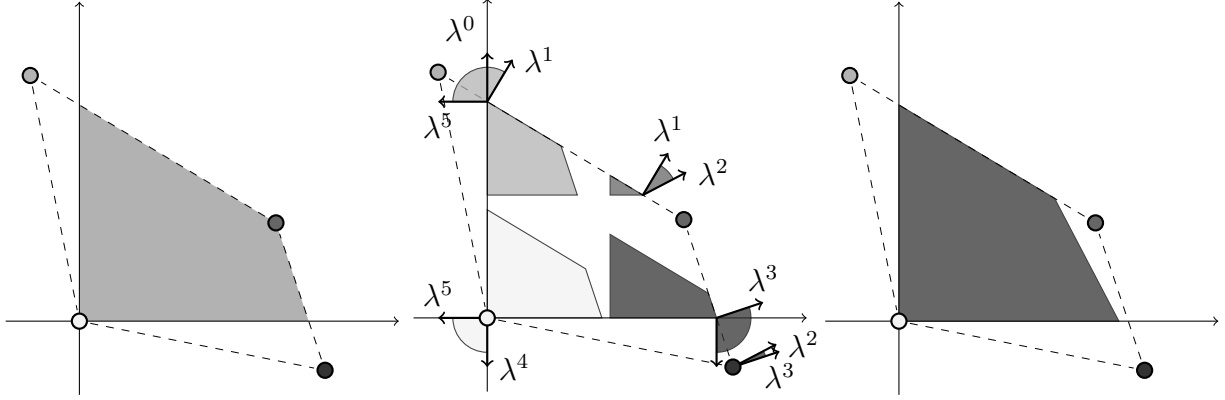


Figure 2: The asymmetric Prisoners' Dilemma. Left: Feasible and individually rational payoffs. Center: APS payoffs and optimal max-min-max levels. Right: Equilibrium payoffs.

use the best APS binding payoff for (D, C) between λ^2 and λ^3 , thus leading to a higher level. Hence, all of these operators would generate a strictly larger \widetilde{W}^1 . Moreover, in the second iteration, they would still be able to generate payoffs that are higher than $(5, -1)$ in directions between λ^2 and λ^3 , and in fact, the optimal level converges only asymptotically.

4.3.2 Risk sharing

Our second example is a risk sharing game in the style of Kocherlakota (1996). Two agents receive time-varying endowments of a consumption good. The total endowment is always equal to 1, and the state s represents player 1's share of the endowment.¹⁵ The state is i.i.d. uniform on an evenly spaced grid between 0 and 1. We let $e_i(s)$ denote player i 's endowment, i.e., $e_1(s) = s$, $e_2(s) = 1 - s$. The agents can make transfers to one another in increments of $1/(M(|S| - 1))$, up to their own endowment. Player i 's consumption is $c_i(a, s) = e_i(s) + a_j - a_i$. Flow utility is $g_i(a, s) = \sqrt{c_i(a, s)}$.

As is well-known, minimum equilibrium payoffs are attained in autarky, where players make no transfers and consume their endowments. The resulting payoffs are

$$\underline{v}_i(s) = (1 - \delta)\sqrt{e_i(s)} + \delta \frac{1}{|S|} \sum_{s' \in S} \sqrt{e_i(s')}.$$

There are more efficient equilibria, wherein the players use transfers to smooth consumption over time. Specifically, the set of feasible payoffs is the convex hull of the vectors $(\sqrt{c}, \sqrt{1 - c})$ for all feasible c .¹⁶ For δ sufficiently large, the folk theorem says that it is possible to attain

¹⁵This version of the model is studied by Ljungqvist and Sargent (2004, Chapter 20).

¹⁶It is a special feature of this game that the set of possible flow payoffs is the same across all states, so that the set of feasible payoffs is just the convex hull of the flow payoffs.

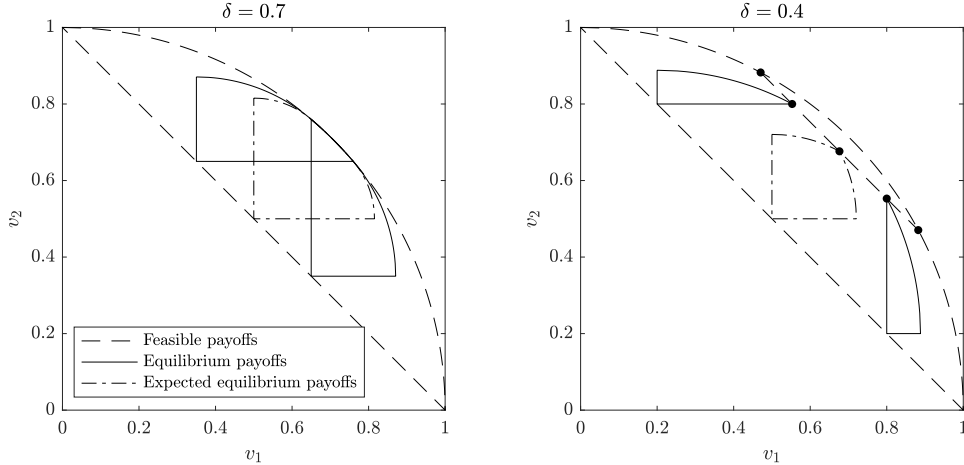


Figure 3: The two-state risk sharing game.

$ S $	M	# Faces	MMM	JYC-100	JYC-200
2	20	27	1.5s	24.6s	78.3s
2	40	47	4.3s	1m 54.9s	5m 30.0s
5	20	102	1m 19.7s	9m 55.4s	30m 55.0s
9	15	145	6m 31.2s	44m 18.7s	2h 25m 46.4s

Table 1: Run times on the risk-sharing game for max-min-max and JYC with 100 and 200 evenly spaced directions.

payoffs on the frontier, provided that both players' insured payoff is at least that of autarky.

Our first example has two states, $M = 200$, and $\delta \in \{0.4, 0.7\}$. Equilibrium payoffs are depicted in Figure 3. We stopped iterating when the Hausdorff distance between $\widetilde{\mathbf{W}}^k$ and $\widetilde{\mathbf{W}}^{k-1}$ was less than 10^{-8} . At $\delta = 0.4$, the algorithm converged in 33 iterations and 16 seconds.¹⁷ At $\delta = 0.7$, it converged in 36 iterations and 1 minute and 18 seconds.

At $\delta = 0.7$, full insurance can be supported at a range of consumption levels. This can be seen because the equilibrium payoff sets overlap with the frontier of feasible payoffs. The only way to generate these payoffs is to repeat the action profile that generates the given flow payoff forever. At the low discount factor ($\delta = 0.4$), equilibrium payoffs are bounded away from the feasible frontier, and efficient risk sharing cannot be supported.

We also computed equilibrium payoffs for a variety of $|S|$ and M , with $\delta = 0.7$ and a convergence threshold of 10^{-6} . The number of faces of \mathbf{V} and run times are reported in Table 1. For comparison, we have also solved this game using the JYC algorithm, which approximates the APS operator in a fixed grid of directions using linear programming. This methodology is readily adapted to stochastic games. Our implementation uses the commercial linear programming software Gurobi, but otherwise it is integrated into the rest of our

¹⁷All benchmarks were measured on a mid-2014 MacBook Pro.

program and uses the same data structures as our implementation of the max-min-max algorithm.¹⁸ Columns JYC-100 and JYC-200 in Table 1 report runtimes for our implementation of the JYC algorithm with 100 and 200 fixed directions, respectively.

We find that JYC takes between one and two orders of magnitude longer than the max-min-max operator. For example, when there are 5 states and $M = 20$ (which corresponds to 81 consumption levels), max-min-max takes 1 minute and 20 seconds, while JYC with 100 directions takes 10 minutes, and with 200 directions takes 31 minutes. Note that these algorithms produce different outputs: the max-min-max limit has faces which approximate those of \mathbf{V} , whereas the JYC limit bounds payoffs in exogenous directions that are unrelated to faces of the equilibrium payoff correspondence.

We should be cautious in drawing general conclusions from these simulations: Run times will vary from game to game, implementation to implementation, and computer to computer. In many applications, it may not be mission critical whether the algorithm terminates in three minutes or in three days. Nonetheless, these simulations strongly suggest that our algorithm will provide faster and more accurate solutions than other known methods.

5 Implementation with many players

5.1 The complexity of \mathbf{V}

We now return to the general setting with many players. A critical complication is that the number of extreme equilibrium payoffs is no longer bounded. Indeed, we will show that the following three-player repeated game has countably infinitely many extreme equilibrium payoffs.

The flow payoffs are depicted in Figure 4, where x is a very low negative payoff, and $\delta = 1/2$. This game has twenty-seven action profiles, but only four can be played in equilibrium. Specifically, (A, A, A) is a static Nash equilibrium, and the permutations of (C, B, B) can be sustained when the discount factor is sufficiently large (in particular when $\delta = 1/2$). Since each player can guarantee themselves a payoff of 0 by always playing A , no other action profile can be played in equilibrium as long as x is sufficiently low.

The equilibrium payoff set V is depicted in the left-hand panel of Figure 5. Online Appendix B contains detailed analysis of the game. We now present an informal overview. Since

¹⁸Our methodology uses the optimal solution in one direction as a starting point for computing solutions in adjacent directions, thereby accelerating convergence to a solution. JYC emphasize the separability of the linear programs for each action profile and direction, and their original Fortran implementation starts each computation from scratch. Our implementation of JYC starts each optimization from the adjacent solution, which considerably speeds up the computation.

a_1/a_2	$a_3 = A$			$a_3 = B$			$a_3 = C$		
	A	B	C	A	B	C	A	B	C
A	$(4, 4, 4)$	$(0, x, 0)$	$(0, x, 0)$	$(0, 0, x)$	$(3, x, x)$	$(0, x, x)$	$(0, 0, x)$	$(0, 0, x)$	$(0, 0, x)$
B	$(x, 0, 0)$	$(x, x, 3)$	$(x, x, 0)$	$(x, 3, x)$	(x, x, x)	$(8, 0, 8)$	$(x, 0, x)$	$(8, 8, 0)$	(x, x, x)
C	$(x, 0, 0)$	$(x, x, 0)$	$(x, x, 0)$	$(x, 0, x)$	$(0, 8, 8)$	(x, x, x)	$(x, 0, x)$	(x, x, x)	(x, x, x)

Figure 4: A three player game.

only four action profiles can be played in equilibrium, we know that V must be contained in the triangular pyramid with corners at $(4, 4, 4)$ and the permutations of $(0, 8, 8)$. It includes $(4, 4, 4)$ and a large flat on the efficient frontier where $v_1 + v_2 + v_3 = 16$. It is easy to show that the minimum equilibrium payoff is 3. For each $i = 1, 2, 3$, there is a face D_i where player i 's payoff is minimized. Besides the Nash payoff, all of the extreme points of V lie in one of these flats. The remaining faces are triangles with the Nash payoff as one vertex.

Let us focus on D_3 , which is depicted in the right panel of Figure 5. The set of extreme points of D_3 has two accumulation points, each of which is approached by two sequences. One sequence starts on the efficient frontier, at points denoted u and u' , and moves down, and the other starts at inefficient payoffs, denoted v and v' , and moves up. All of these payoffs are generated by playing (B, B, C) for one period, followed by a continuation payoff in $C(B, B, C)$.

The set $C(B, B, C)$ is generated as follows. The height of the blue plane in the left-hand panel of Figure 5 is the continuation value at which player 3's incentive constraint binds, i.e., $v_3 = 6$. For each D_k with $k = 1, 2$, exactly one of the accumulation points and its corresponding sequences lies above the plane. For each element \hat{v} of these sequences, there is a corresponding extreme point of $C(B, B, C)$, where the line between \hat{v} and $(4, 4, 4)$ crosses $v_3 = 6$. In addition, there are four more payoffs in $C(B, B, C)$, two of which are generated by randomizing between u and u' , and two of which are generated by randomizing between v and v' . Thus, half of the sequences in D_3 are generated from sequences in each D_k , for $k = 1, 2$.

We note that the same basic structure obtains if we perturb the players' payoffs, and it seems to be a generic possibility that the number of extreme equilibrium payoffs is infinite. As an aside, the development of this example illustrates the potential value of computational methods in repeated games. Only after being inspired by numerical results were we able to construct the example analytically.

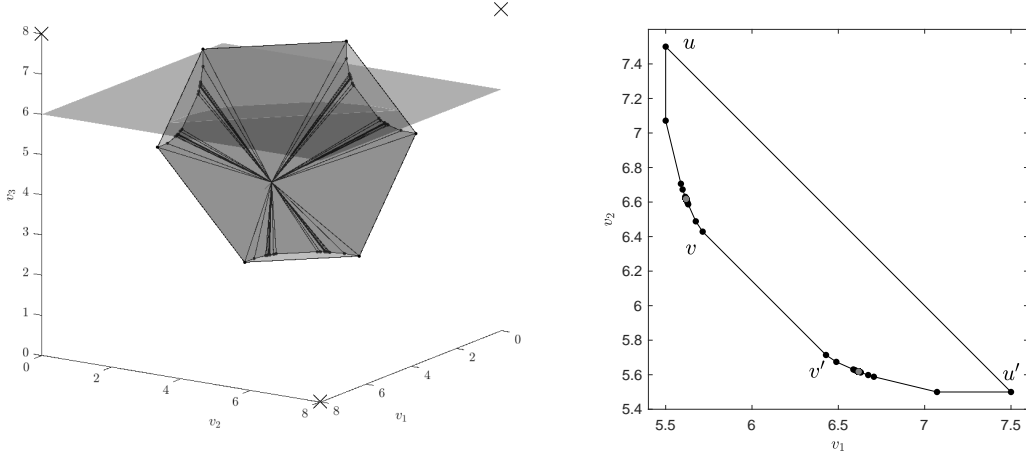


Figure 5: Left: Equilibrium payoffs set for the game in Figure 4 with $\delta = 1/2$, looking down on the Pareto frontier. Right: equilibrium payoffs in which $v_3 = 3$.

5.2 Implementation and approximation with $N \geq 3$

Given that \mathbf{V} may have infinitely many extreme points, exact computation of the sequence from (T1.iv) may be impossible. We now describe a procedure that can compute \tilde{B} exactly when the number of extreme points is small and approximates it when the number of extreme points blows up.

At every iteration, there will be a set of directions $\hat{\Lambda}^k$, with $\hat{\Lambda}^0$ being arbitrary, and we compute the new correspondence $\widehat{\mathbf{W}}^{k+1} = \tilde{B}(\widehat{\mathbf{W}}^k, \hat{\Lambda}^k)$, where

$$\tilde{B}(\mathbf{W}, \hat{\Lambda})(s) = \{v | \lambda \cdot v \leq x(s, \lambda, \mathbf{W}) \forall \lambda \in \hat{\Lambda}\}.$$

Thus, $\tilde{B}(\cdot, \hat{\Lambda})$ is analogous to \tilde{B} , but where we only bound payoffs for directions in $\hat{\Lambda}$. For future reference, we similarly define

$$B(\mathbf{W}, \hat{\Lambda})(s) = \{\lambda \cdot v \leq x^{APS}(s, \lambda, \mathbf{W}) \text{ for all } \lambda \in \hat{\Lambda}\}.$$

We will fully specify the procedure for updating $\hat{\Lambda}^k$ shortly. At a high level, we will drop directions that are redundant and add new directions that correspond to faces of $\tilde{B}(\widehat{\mathbf{W}}^k)$, in a manner analogous to what we did for two players. We will, however, cap the size of $\hat{\Lambda}^k$ so that the complexity is bounded. Note that for any sequence $\{\hat{\Lambda}^k\}$, each correspondence $\widehat{\mathbf{W}}^k$ must contain $\tilde{\mathbf{W}}^k$, and thus $\mathbf{V} \subseteq \bigcap_{k \geq 0} \widehat{\mathbf{W}}^k$. By updating the directions endogenously, the sequence $\widehat{\mathbf{W}}^k$ will converge to \mathbf{V} if memory permits, and otherwise we coarsen the approximation to satisfy the complexity bound.

It is critical, however, to choose $\widehat{\Lambda}^k$ so that we retain a key computational advantage of the max-min-max operator, which is the sufficiency of binding payoffs. When $\widehat{\Lambda}^k = \Lambda$, this was established by Lemmas 5 and 6. (Recall that Λ is the set of all directions.) These results still apply, but the following weaker formulation of Lemma 6 will be useful:

Lemma 16. *For any $\widehat{\Lambda} \subseteq \Lambda$, if $\widetilde{B}(\mathbf{W}, \widehat{\Lambda}) \subseteq \mathbf{W}$, then $B(\widehat{\mathbf{W}}) \subseteq \widehat{\mathbf{W}}$, where $\widehat{\mathbf{W}} = \widetilde{B}(\mathbf{W}, \widehat{\Lambda})$.*

Proof of Lemma 16. The proof of Lemma 6 directly implies that $B(\widehat{\mathbf{W}}, \widehat{\Lambda}) \subseteq \widehat{\mathbf{W}}$. Since B is decreasing in its second argument, the weaker conclusion of the lemma follows. \square

By Lemma 5, if $B(\mathbf{W}, \widehat{\Lambda}) \subseteq \mathbf{W}$, then a positive APS gap in some direction implies that the minimal regime can be taken to be recursive. As a consequence of Lemma 16, as long as we choose $\widehat{\Lambda}^k$ such that $\widetilde{B}(\widehat{\mathbf{W}}^k, \widehat{\Lambda}^k) \subseteq \widetilde{B}(\widehat{\mathbf{W}}^{k-1}, \widehat{\Lambda}^{k-1}) = \widehat{\mathbf{W}}^k$, the inductive hypothesis that B sub-generates will be satisfied, so that binding payoffs are sufficient and we need only compute binding APS payoffs and the local frontier around those payoffs.

In fact, we can go a step further. Recall that the computation of $\widetilde{B}(\mathbf{W}, \widehat{\Lambda})$ only depends on the binding APS level $\widehat{x}^{APS}(s, \lambda, \mathbf{W})$ and the sign of $\gamma(a, \lambda, \mathbf{W})$. If these two functions are the same for correspondences \mathbf{W} and \mathbf{W}' , then we say that they have the same *local binding frontier*. Now, suppose that $\widehat{\Lambda}^k$ is such that the correspondence $\widetilde{B}(\widehat{\mathbf{W}}^{k-1}, \widehat{\Lambda}^k)$ has the same local binding frontier as $\widehat{\mathbf{W}}^k$. Then $\widetilde{B}(\widehat{\mathbf{W}}^{k-1}, \widehat{\Lambda}^k)$ need not be a subset of $\widetilde{B}(\widehat{\mathbf{W}}^{k-1}, \widehat{\Lambda}^{k-1})$. But as long as we continue to use \widetilde{B} in subsequent iterations, the resulting computations will be *exactly the same*:

Lemma 17. *Suppose $\mathbf{W} \subseteq \mathbf{W}'$, both correspondences have the same local binding frontier, and B sub-generates at \mathbf{W} . Then for all s, λ , and $\mathbf{a} \in \mathbf{A}(\mathbf{W})$, $x(s, \lambda, \mathbf{a}, \mathbf{W}) = x(s, \lambda, \mathbf{a}, \mathbf{W}')$.*

Proof of Lemma 17. Since B sub-generates at \mathbf{W} , the hypothesis of Lemma 5 is satisfied. Let us then take \mathbf{r} to be minimal regimes for \mathbf{a} at \mathbf{W} in the direction λ that satisfy the refinement of Lemma 5, i.e., the regime is recursive whenever the APS gap is strictly positive. Since \mathbf{W} and \mathbf{W}' have the same sign APS gap and the same binding payoffs, we conclude that

$$x(s, \lambda, \mathbf{a}, \mathbf{r}, \mathbf{W}) = x(s, \lambda, \mathbf{a}, \mathbf{r}, \mathbf{W}'). \quad (16)$$

To verify that \mathbf{r} is minimal for \mathbf{a} at \mathbf{W}' , we need only check the minimality conditions (6). Equation (16) implies that

$$x^R(\mathbf{a}(s), \lambda, \mathbf{a}, \mathbf{r}, \mathbf{W}) = x^R(\mathbf{a}(s), \lambda, \mathbf{a}, \mathbf{r}, \mathbf{W}'). \quad (17)$$

When $\mathbf{r}(s) = APS$, the optimal APS level is binding, and by hypothesis this level is the same for \mathbf{W} and \mathbf{W}' , so that both sides of (6) are unchanged. When $\mathbf{r}(s) = R$, (6) implies

that $x^R(\mathbf{a}(s), \lambda, \mathbf{a}, \mathbf{r}, \mathbf{W}) \leq x^{APS}(\mathbf{a}(s), \lambda, \mathbf{W})$. As x^{APS} is increasing in its third argument, this implies that $x^R(\mathbf{a}(s), \lambda, \mathbf{a}, \mathbf{r}, \mathbf{W}) \leq x^{APS}(\mathbf{a}(s), \lambda, \mathbf{W}')$. Combining this last inequality with (17), we conclude that \mathbf{r} satisfies (6) at \mathbf{W}' as well. \square

In addition, if we *add* a new direction to $\widehat{\Lambda}^k$ that was not present in $\widehat{\Lambda}^{k-1}$, then this will only cause $\widehat{\mathbf{W}}^k$ to shrink further, so that the hypothesis for binding payoffs to be sufficient will still be satisfied:

Theorem 3. *Suppose the sequence $\{\widehat{\Lambda}^k\}_{k \geq 0}$ is such that for every $k \geq 1$, $\overline{\mathbf{W}}^k = \widetilde{B}(\widehat{\mathbf{W}}^{k-1}, \widehat{\Lambda}^k \cap \widehat{\Lambda}^{k-1})$ has the same local binding frontier as $\widehat{\mathbf{W}}^k = \widetilde{B}(\widehat{\mathbf{W}}^{k-1}, \widehat{\Lambda}^{k-1})$, and assume $B(\widehat{\mathbf{W}}^0, \widehat{\Lambda}^0) \subseteq \widehat{\mathbf{W}}^0$. Then for every $k \geq 0$, $B(\widehat{\mathbf{W}}^k) \subseteq \widehat{\mathbf{W}}^k$. As a result, for any $\mathbf{a} \in \mathbf{A}(\widehat{\mathbf{W}}^k)$, there exist minimal regimes such that $\mathbf{r}(s) = R$ whenever $\gamma(\mathbf{a}(s), \lambda, \widehat{\mathbf{W}}^k) > 0$.*

Proof of Theorem 3. We will argue that for every k , $B(\widehat{\mathbf{W}}^k) \subseteq \widehat{\mathbf{W}}^k$. The second part of the theorem then follows directly from Lemma 5.

Take as inductive hypotheses that (i) $\widetilde{B}(\overline{\mathbf{W}}^{k-1}, \widehat{\Lambda}^{k-1}) \subseteq \overline{\mathbf{W}}^{k-1}$ and (ii) $B(\widehat{\mathbf{W}}^{k-1}) \subseteq \widehat{\mathbf{W}}^{k-1}$. We set $\overline{\mathbf{W}}^0 = \widehat{\mathbf{W}}^0$ so that these hypotheses are true for the base case $k = 1$. Also note for current and later use that $\widetilde{B}(\cdot, \cdot)$ is increasing in its first argument and decreasing in its second, as is $B(\cdot, \cdot)$. As a result, $\widehat{\mathbf{W}}^{k-1} \subseteq \overline{\mathbf{W}}^{k-1}$. The assumption that $\overline{\mathbf{W}}^j$ and $\widehat{\mathbf{W}}^j$ have the same local binding frontier for every $j \geq 1$ is used freely below. Together with the inductive hypothesis (ii), Lemma 17 then implies that $\widehat{\mathbf{W}}^{k-1}$ and $\overline{\mathbf{W}}^{k-1}$ are interchangeable in the \widetilde{B} operator. Consequently,

$$\widetilde{B}(\overline{\mathbf{W}}^{k-1}, \widehat{\Lambda}^{k-1}) = \widetilde{B}(\widehat{\mathbf{W}}^{k-1}, \widehat{\Lambda}^{k-1}) = \widehat{\mathbf{W}}^k. \quad (18)$$

By the inductive hypothesis (i) and Lemma 16 (where $\widehat{\Lambda} = \widehat{\Lambda}^{k-1}$, $\mathbf{W} = \overline{\mathbf{W}}^{k-1}$, and $\mathbf{W}' = \widehat{\mathbf{W}}^k$), it follows that $B(\widehat{\mathbf{W}}^k) \subseteq \widehat{\mathbf{W}}^k$, thus extending the inductive hypothesis (ii) to k .

Lemma 17 then implies that $\widehat{\mathbf{W}}^k$ and $\overline{\mathbf{W}}^k$ are also interchangeable in \widetilde{B} . Also, (i) and interchangeability of $\widehat{\mathbf{W}}^{k-1}$ and $\overline{\mathbf{W}}^{k-1}$ imply that $\widehat{\mathbf{W}}^k = \widetilde{B}(\widehat{\mathbf{W}}^{k-1}, \widehat{\Lambda}^{k-1}) = \widetilde{B}(\overline{\mathbf{W}}^{k-1}, \widehat{\Lambda}^{k-1}) \subseteq \overline{\mathbf{W}}^{k-1}$. Hence,

$$\begin{aligned} \widetilde{B}(\overline{\mathbf{W}}^k, \widehat{\Lambda}^k) &\subseteq \widetilde{B}(\overline{\mathbf{W}}^k, \widehat{\Lambda}^k \cap \widehat{\Lambda}^{k-1}) \\ &= \widetilde{B}(\widehat{\mathbf{W}}^k, \widehat{\Lambda}^k \cap \widehat{\Lambda}^{k-1}) \\ &\subseteq \widetilde{B}(\overline{\mathbf{W}}^{k-1}, \widehat{\Lambda}^k \cap \widehat{\Lambda}^{k-1}) \\ &= \widetilde{B}(\widehat{\mathbf{W}}^{k-1}, \widehat{\Lambda}^k \cap \widehat{\Lambda}^{k-1}) = \overline{\mathbf{W}}^k, \end{aligned}$$

where we have used, in order, monotonicity of \widetilde{B} in its second argument, interchangeability of $\overline{\mathbf{W}}^k$ and $\widehat{\mathbf{W}}^k$, monotonicity of \widetilde{B} in its first argument, interchangeability of $\overline{\mathbf{W}}^{k-1}$ and

$\widehat{\mathbf{W}}^{k-1}$, and the definition of $\overline{\mathbf{W}}^k$. This extends the inductive hypothesis (i) to k . \square

Having established conditions on $\{\widehat{\Lambda}^k\}$ for binding payoffs to be sufficient, we now fill in the remaining details of the algorithm. Fix a positive integer $L > 0$, which is the maximum number of directions. At every iteration, we construct $\widehat{\Lambda}^k$ by first dropping directions in a set Λ' such that $\widetilde{B}(\widehat{\mathbf{W}}^{k-1}, \widehat{\Lambda}^{k-1})$ and $\widetilde{B}(\widehat{\mathbf{W}}^{k-1}, \widehat{\Lambda}^{k-1} \setminus \Lambda')$ have the same local binding frontier. This is in fact easy to do: in the process of computing $C(a, \widetilde{\mathbf{W}}^{k-1})$, we intersect the binding payoff sets with the half-spaces that define the expected continuation value set $\sum_{s \in S} \pi(s|a) \widehat{\mathbf{W}}^k(s)$. There is some order in which the directions are intersected. If we find that the half space $(\lambda, x(s, \lambda, \widehat{\mathbf{W}}^{k-1}))$ does not result in a change to $C(a, \mathbf{W})$ or the local frontier around those payoffs for any a , then the direction is redundant and can be dropped.¹⁹

After dropping directions, we can add new directions as long as the total number is less than L . The directions we add are faces of $\widetilde{B}(\widehat{\mathbf{W}}^{k-1})$, i.e., directions for which there are N different optimal payoffs which are not contained in an affine subspace of dimension less than $N - 1$. This is the natural generalization of test directions from the two-player case.

To compute a face direction, we first pick an initial direction λ^0 uniformly in Λ and compute optimal payoffs \mathbf{u} , which with probability one are uniquely optimal (and hence remain optimal for perturbations in all directions). We then draw a direction of rotation $\widetilde{\lambda}^1$ uniformly on Λ and compute the legitimate test direction with the smallest rotation from λ^0 in the direction $\widetilde{\lambda}^1$, which is denoted λ^1 . This step is analogous to that described in Section 4. The corresponding substitution is (s^1, a^1, p^1) , and let $d^1 = u(s^1, a^1, p^1, \mathbf{u}) - \mathbf{u}(s^1)$ be the direction in which the substitution moves payoffs. We then pick a new direction $\widetilde{\lambda}^2$ uniformly on the set of directions which are orthogonal to d^1 , and compute the shallowest legitimate test direction from λ^0 in the direction $\widetilde{\lambda}^2$, denoted λ^2 , with substitution (s^2, a^2, r^2) . Continuing inductively, after n steps, we will have a direction λ^n and n substitutions which move payoffs in linearly-independent directions d^1, \dots, d^n that are orthogonal to λ^n . After $N - 1$ steps, the process converges to a face direction λ^{N-1} . Pseudocode for the face-finding procedure is in Algorithms 8 and 9.

We run this procedure $L - |\widehat{\Lambda}^{k-1} \setminus \Lambda'|$ times and add the new face directions (skipping duplicates) to $\widehat{\Lambda}^{k-1}$ to obtain $\widehat{\Lambda}^k$. (Note that L may be greater than the number of face directions, in which case this procedure necessarily encounters duplicates.) Algorithm 10 gives pseudocode for this step. In practice, we have found it better to generate new directions once every five iterations or so, while redundant directions are dropped every iteration. This completes the many-player algorithm.

¹⁹Our code randomizes the order in which the half spaces are intersected, in a crude effort to identify a non-redundant set of directions. This process could in principle be systematized using standard convex hull algorithms.

We may compare this procedure with that of JYC. Both algorithms bound payoffs in a finite number of directions. While JYC use the APS level for each direction, we use the max-min-max level. The difference is significant, since there are many fewer payoffs that can be optimal for max-min-max, and we solve out the optimal payoffs when the regime is recursive. In addition, JYC hold the directions fixed, whereas our procedure adjusts the directions dynamically to achieve a sharper approximation.

We have implemented this procedure as part of our software package. Online Appendix C describes two numerical examples. The first is a simple three-player binary-action contribution game, in which the equilibrium payoff set has a small number of faces, to which the algorithm converges exactly. The second example is a three-player version of the risk-sharing game considered in Section 4. We use our algorithm to illustrate how partial contracts, in which two of the three players commit to perfectly insure one another, can lead to lower welfare for all parties, as the partial contract limits the players' ability to punish deviations.

The procedure we implemented is just one of many possible ways to generate a sequence $(\widetilde{\mathbf{W}}^k, \widehat{\Lambda}^k)$ that satisfies the hypotheses of Theorem 3. For example, we could take $\widetilde{\mathbf{W}}^0$ to be the feasible payoff correspondence, or to be large hypercubes that contain all flow payoffs (as we have done in our simulations). In either case, the initial correspondence has finitely many extreme points, so that $\widetilde{B}(\mathbf{W}^0)$ is guaranteed to have finitely many faces. As long as L is sufficiently large, we can set $\widehat{\Lambda}^0 = \Lambda$, i.e., all directions, and compute $\widetilde{B}(\mathbf{W}^0)$ exactly. This can be done via a generalization of the face-finding procedure, whereby we recursively map faces that are adjacent to one another via successive one-dimensional searches. This precision could be maintained as long as the number of face directions is less than L . One could even use the above stochastic algorithm temporarily, while the number of face directions is large, but reverts to exact computation when the payoff sets simplify. We have not attempted to explore all of these possibilities, and they are promising directions for future research.

6 Lower bounds on \mathbf{V}

Like that of APS, our methodology generates a sequence that converges to \mathbf{V} from the “outside,” meaning that every element is a superset of \mathbf{V} . This sequence may only converge asymptotically, so that if we stop iterating after finitely many rounds, we only obtain an upper bound on \mathbf{V} . As a final topic, we adapt our methodology to compute a corresponding lower bound.

For a fixed $\epsilon > 0$, consider the following perturbed APS operator: $B^\epsilon(\mathbf{W})(s) = \{v | \lambda \cdot v \leq x^{APS}(s, \lambda, \mathbf{W}) - \epsilon \forall \lambda \in \Lambda\}$. Just like the APS operator, B^ϵ takes compact correspondences to compact correspondences and is increasing in \mathbf{W} . As a result, there is a largest bounded

fixed point, denoted \mathbf{V}^ϵ , which can be computed by iterative application of B^ϵ on any correspondence that contains \mathbf{V}^ϵ . Moreover, B^ϵ is decreasing in ϵ , which implies that \mathbf{V}^ϵ is decreasing in ϵ , so that $\mathbf{V}^\epsilon \subseteq \mathbf{V}^0 = \mathbf{V}$.

We propose to compute a lower bound on \mathbf{V} by generating a sequence that converges to \mathbf{V}^ϵ . At first glance, this plan seems to have the same problem: How do we know when we have converged to \mathbf{V}^ϵ ? The difference is that we do not want to compute \mathbf{V}^ϵ ; we just want to find a set that self generates. In particular, since the iterates of B^ϵ converge to \mathbf{V}^ϵ , eventually we will obtain a correspondence \mathbf{W} such that the Hausdorff distance between \mathbf{W} and $B^\epsilon(\mathbf{W})$ is less than $\epsilon/2$. As a result, $B(\mathbf{W})$ will be strictly larger than \mathbf{W} by at least $\epsilon/2$ in every direction, so that we can robustly certify that $B(\mathbf{W}) \subseteq \mathbf{V}$. Note that we have no guarantee that \mathbf{V}^ϵ is close to \mathbf{V} , or even non-empty valued. In simulations discussed below, however, the lower bound we obtain seems to converge to \mathbf{V} as ϵ goes to zero.²⁰

At the same time, a premise of this paper is that B is hard to compute, which is why we developed the operator \tilde{B} , and the same computational difficulties arise when computing B^ϵ . Fortunately, our methodology can be adapted to compute \mathbf{V}^ϵ . Given a policy (\mathbf{a}, \mathbf{r}) , we define $x^\epsilon(s, \lambda, \mathbf{a}, \mathbf{r}, \mathbf{W})$ to be the unique solution to

$$x^\epsilon(s, \lambda, \mathbf{a}, \mathbf{r}, \mathbf{W}) = -\epsilon + \begin{cases} (1 - \delta)\lambda \cdot g(\mathbf{a}(s)) + \delta \sum_{s' \in S} \pi(s' | \mathbf{a}(s)) x^\epsilon(s', \lambda, \mathbf{a}, \mathbf{r}, \mathbf{W}) & \text{if } \mathbf{r}(s) = R; \\ x^{APS}(\mathbf{a}(s), \lambda, \mathbf{W}) & \text{if } \mathbf{r}(s) = APS. \end{cases}$$

We define $x^\epsilon(s, \lambda, \mathbf{W}) = \max_{\mathbf{a} \in \mathbf{A}(\mathbf{W})} \min_{\mathbf{r} \in \mathbf{R}} x^\epsilon(s, \lambda, \mathbf{a}, \mathbf{r}, \mathbf{W})$ and $\tilde{B}^\epsilon(\mathbf{W}) = \{v | \lambda \cdot v \leq x^\epsilon(s, \lambda, \mathbf{W}) \forall \lambda \in \Lambda\}$. The operator \tilde{B}^ϵ satisfies all the desirable properties of \tilde{B} . In particular, \tilde{B}^ϵ is increasing, maps compact correspondences to compact correspondences, and has \mathbf{V}^ϵ as a fixed point. As a result, if we fix a correspondence $\tilde{\mathbf{W}}^0$ that contains \mathbf{V}^ϵ and generate the sequence $\tilde{\mathbf{W}}^k = \tilde{B}^\epsilon(\mathbf{W}^{k-1})$, then $\mathbf{V}^\epsilon = \bigcap_{k=0}^\infty \tilde{\mathbf{W}}^k$. In fact, we can even take $\tilde{\mathbf{W}}^0$ to be the upper bound on \mathbf{V} obtained by iterative application of \tilde{B} .

This result is reported as Theorem 4 in Online Appendix E, where we rederive all of our key results, adding in ϵ 's where appropriate. All of our other results extend as well: For every direction, there exists a state-independent optimal policy. Binding payoffs are sufficient as long as B^ϵ sub-generates, and B^ϵ will sub-generate along the sequence we compute as long as it sub-generates at the first round.²¹ There is a simple set of conditions that characterize

²⁰We have good reason to think that \mathbf{V}^ϵ will not collapse in general. For example, if we were to drop incentive constraints, then the analogue of \mathbf{V}^ϵ is simply the contraction of the feasible payoff correspondence in every direction by $\epsilon/(1 - \delta)$.

²¹For this result, it is essential that the ϵ penalty is recursively compounded in the definition of x^ϵ . If we simply added an ϵ penalty in both regimes, the resulting operator *would* result in a sequence that converges to \mathbf{V}^ϵ , but the penalty attached to recursive payoffs would be too small, so that eventually, the minimal regime would always be *APS*.

optimal policies and optimal pairs. When there are two players, the correspondence \mathbf{V}^ϵ has at most \bar{L} extreme payoffs, and \tilde{B}^ϵ can be computed in runtime $O(\bar{L}\bar{M}^2)$.²²

We have implemented the operator \tilde{B}^ϵ for two players as part of our software package. In Online Appendix C, we report an application to the risk-sharing example of Section 4, for which the lower bound and upper bound are virtually indistinguishable. Beyond two players, further approximation of \tilde{B}^ϵ may be needed, as discussed in Section 5.

7 Conclusion

It has been our purpose to study the subgame perfect equilibria of stochastic games. We have developed a new fundamental structural property of extremal equilibria, namely that equilibrium play is stationary until incentive constraints bind. We developed a new “max-min-max” algorithm that exploits this structure, using policy iteration when incentive constraints are slack to obtain tighter bounds than the APS operator on which payoffs can be generated. The bounds can also be computed using only knowledge of binding payoffs and the slope of the frontier around those payoffs. Moreover, the optimal equilibrium structure changes in only one state at a time as the direction of optimization moves, which greatly simplifies the computation of the set of payoffs that can be generated. When there are two players, the resulting algorithm, and by extension the equilibrium payoff correspondence, are of bounded complexity. We have shown by example that the number of extreme equilibrium payoffs may be infinite with more than two players, but we have provided a flexible routine that can approximate equilibrium payoffs when computing power is limited and compute them exactly when the equilibrium correspondence is not too complicated.

The insights that we have developed are obviously particular to the special class of games we have considered. We have made heavy use of perfect monitoring, public randomization, and the restriction to pure-strategy equilibria. These assumptions are widely used both in theory and application and, in our view, are eminently worthy of study. While the basic results on the max-min-max operator can be extended, these particular insights will presumably be more or less useful for computation depending on the class of game being studied. At a broader level, our approach is to develop methods that are tailored to the special structure that arises in extremal equilibria. It is our hope that similar efforts will bear fruit for other classes of games and solution concepts, for example, those involving imperfect monitoring or mixed strategies.

²²A subtlety here is that the function x^ϵ is no longer linear in λ . Nonetheless, as we argue in Online Appendix E, B^ϵ can be computed by intersecting bounds at test directions.

References

- ABREU, D., B. BROOKS, AND Y. SANNIKOV (2016): “A “pencil-sharpening” algorithm for two-player stochastic games with perfect monitoring,” Tech. rep., New York University and University of Chicago and Stanford University.
- ABREU, D., D. PEARCE, AND E. STACCHETTI (1986): “Optimal cartel equilibria with imperfect monitoring,” *Journal of Economic Theory*, 39, 251–269.
- (1990): “Toward a theory of discounted repeated games with imperfect monitoring,” *Econometrica*, 58, 1041–1063.
- ABREU, D. AND Y. SANNIKOV (2014): “An algorithm for two-player repeated games with perfect monitoring,” *Theoretical Economics*, 9, 313–338.
- ATKESON, A. (1991): “International lending with moral hazard and risk of repudiation,” *Econometrica: Journal of the Econometric Society*, 1069–1089.
- BERG, K. (2019): “Set-valued games and mixed-strategy equilibria in discounted supergames,” *Discrete Applied Mathematics*, 255, 1–14.
- BERG, K. AND M. KITTI (2019): “Equilibrium paths in discounted supergames,” *Discrete Applied Mathematics*.
- BLACKWELL, D. (1962): “Discrete dynamic programming,” *The Annals of Mathematical Statistics*, 719–726.
- (1965): “Discounted dynamic programming,” *The Annals of Mathematical Statistics*, 226–235.
- DANTZIG, G. B. AND M. N. THAPA (2006): *Linear programming 1: introduction*, Springer Science & Business Media.
- DIXIT, A., G. M. GROSSMAN, AND F. GUL (2000): “The dynamics of political compromise,” *Journal of political economy*, 108, 531–568.
- ERICSON, R. AND A. PAKES (1995): “Markov-perfect industry dynamics: A framework for empirical work,” *The Review of Economic Studies*, 62, 53–82.
- HÖRNER, J., T. SUGAYA, S. TAKAHASHI, AND N. VIEILLE (2011): “Recursive methods in discounted stochastic games: An algorithm for $\delta \rightarrow 1$ and a folk theorem,” *Econometrica*, 79, 1277–1318.

- JUDD, K. L., S. YELTEKIN, AND J. CONKLIN (2003): “Computing supergame equilibria,” *Econometrica*, 71, 1239–1254.
- KOCHERLAKOTA, N. R. (1996): “Implications of efficient risk sharing without commitment,” *The Review of Economic Studies*, 63, 595–609.
- LJUNGQVIST, L. AND T. J. SARGENT (2004): *Recursive macroeconomic theory*, MIT press.
- MAILATH, G. J. AND L. SAMUELSON (2006): “Repeated games and reputations: long-run relationships,” *OUP Catalogue*.
- PAKES, A. AND P. MCGUIRE (1994): “Computing Markov-Perfect Nash Equilibria: Numerical Implications of a Dynamic Differentiated Product Model,” *RAND Journal of Economics*, 25, 555–589.
- PHELAN, C. AND E. STACCHETTI (2001): “Sequential equilibria in a Ramsey tax model,” *Econometrica*, 69, 1491–1518.
- RENNER, P. AND S. SCHEIDEGGER (2018): “Machine learning for dynamic incentive problems,” Tech. rep., Lancaster University and University of Lausanne.
- SLEET, C. AND Ş. YELTEKIN (2016): “On the computation of value correspondences for dynamic games,” *Dynamic Games and Applications*, 6, 174–186.
- YELTEKIN, Ş., Y. CAI, AND K. L. JUDD (2017): “Computing equilibria of dynamic games,” *Operations Research*, 65, 337–356.

A Online appendix:

Pseudocode for Sections 3, 4, and 5

Algorithm 1 Minimize regimes.

```

1: procedure MINIMIZEREGIMES( $\lambda, \mathbf{a}, \mathbf{r}, \mathbf{W}$ )
2:   define  $\tilde{S}$  to be the states with  $\gamma(\mathbf{a}(s), \lambda, \mathbf{W}) > 0$ 
3:   define  $\mathbf{r}' := \mathbf{r}$ 
4:   for all  $s \in \tilde{S}$  do
5:      $\mathbf{r}'(s) := R$  ▷ For these states, recursive can be taken to be minimal
6:   loop
7:      $\mathbf{r}'' := \mathbf{r}'$ 
8:     for all  $s \notin \tilde{S}$  do
9:       if  $\hat{x}^{APS}(\mathbf{a}(s), \lambda, \mathbf{W}) < x(s, \lambda, \mathbf{a}, \mathbf{r}', \mathbf{W})$  then
10:         $\mathbf{r}''(s) := APS$  ▷ The best APS payoff is lower
11:       else if  $x^R(\mathbf{a}(s), \lambda, \mathbf{a}, \mathbf{r}', \mathbf{W}) < x(s, \lambda, \mathbf{a}, \mathbf{r}', \mathbf{W})$  then
12:         $\mathbf{r}''(s) := R$  ▷ The recursive payoff is lower
13:       if  $\mathbf{r}'' \neq \mathbf{r}'$  then
14:          $\mathbf{r}' := \mathbf{r}''$  ▷ Continue updating
15:       else
16:         return  $\mathbf{r}'$  ▷ These regimes are minimal

```

Algorithm 2 Optimize the policy.

```
1: procedure OPTIMIZEPOLICY( $\lambda, \mathbf{W}$ )
2:   define  $\mathbf{a} \in \mathbf{A}(\mathbf{W})$ 
3:   define  $\mathbf{r} \in \mathbf{R}$ 
4:   loop
5:     define  $\mathbf{a}' := \mathbf{a}$ 
6:      $\mathbf{r} := \text{MINIMIZEREGIMES}(\lambda, \mathbf{a}, \mathbf{r}, \mathbf{W})$ 
7:     for all  $s \in S, a \in \mathbf{A}(\mathbf{W})(s)$  do
8:       if  $\text{and}(x^R(a, \lambda, \mathbf{a}, \mathbf{r}, \mathbf{W}) > x(s, \lambda, \mathbf{a}, \mathbf{r}, \mathbf{W}),$   

          $\text{or}(\gamma(a, \lambda, \mathbf{W}) > 0, \hat{x}^{APS}(a, \lambda, \mathbf{W}) > x(s, \lambda, \mathbf{a}, \mathbf{r}, \mathbf{W}))$  then
9:          $\mathbf{a}'(s) := a$ 
10:    if  $\mathbf{a} \neq \mathbf{a}'$  then
11:       $\mathbf{a} := \mathbf{a}'$  ▷ Continue updating
12:    else
13:      return  $(\mathbf{a}, \mathbf{r})$  ▷ The policy is optimal
```

Algorithm 3 Compute the shallowest legitimate test direction with $N = 2$.

Require: \mathbf{u} is robustly optimal for direction λ

```
1: procedure FINDNEXTDIRECTION( $\lambda, \mathbf{u}, \mathbf{W}$ )
2:   define  $\lambda' := \lambda$ 
3:   for all  $s \in S, a \in \mathbf{A}(\mathbf{W})(s), p \in \{R\} \cup C(a, \mathbf{W})$  do ▷ Iterate over all substitutions
4:     for all  $\lambda''$  that is a test direction for  $(s, a, p)$  given  $\mathbf{u}$  do
5:       if  $\text{and}(\lambda''$  is legitimate,  $\lambda''$  is shallower than  $\lambda')$  then
6:          $\lambda' := \lambda''$ 
7:   return  $\lambda'$  ▷ The optimal payoffs may change at  $\lambda'$ 
```

Given $N = 2$, let $\hat{u}^{APS}(a, \lambda, \mathbf{W})$ be the highest binding APS payoff in the direction λ where comparisons are made lexicographically using $>_\lambda$. In a slight abuse of notation, we write $\gamma(a, \lambda^+, \mathbf{W}) > 0$ if the APS gap for a is lexicographically positive at λ . The following procedure, analogous to Algorithm 1, uses lexicographic comparisons to choose the regimes which become minimal for the action tuple \mathbf{a} after direction λ

Algorithm 4 Lexicographically minimize regimes for $N = 2$.

```

1: procedure LEXMINIMIZEREGIMES( $\lambda, \mathbf{a}, \mathbf{p}, \mathbf{W}$ )
2:   define  $\tilde{S}$  to be the states where  $\gamma(a, \lambda^+, \mathbf{W}) > 0$ 
3:   define  $\mathbf{p}' := \mathbf{p}$ 
4:   for all  $s \in \tilde{S}$  do
5:      $\mathbf{p}'(s) := R$  ▷ For these states, recursive must be minimal
6:   loop
7:     define  $\mathbf{p}'' := \mathbf{p}'$ 
8:     define  $\mathbf{u} :=$  the payoffs induced by  $(\mathbf{a}, \mathbf{p}')$ 
9:     for all  $s \notin \tilde{S}$  do
10:      if  $\mathbf{u}(s) >_{\lambda} \hat{u}^{APS}(\mathbf{a}(s), \lambda, \mathbf{W})$  then
11:         $\mathbf{r}''(s) := APS$  ▷ The best APS payoff is lexicographically lower
12:      else if  $\mathbf{u}(s) >_{\lambda} u^R(\mathbf{a}(s), \lambda, \mathbf{u})$  then
13:         $\mathbf{p}''(s) := R$  ▷ The recursive payoff is lexicographically lower
14:      if  $\mathbf{p}'' \neq \mathbf{p}'$  then
15:         $\mathbf{p}' := \mathbf{p}''$  ▷ Continue updating
16:      else
17:        return  $\mathbf{p}'$  ▷ This  $\mathbf{p}$  is minimal

```

The following procedure, analogous to Algorithm 2, uses lexicographic comparisons to find the pair which is optimal after direction λ (i.e. the robustly optimal pair).

Algorithm 5 Lexicographically optimize the policy.

```

1: procedure LEXOPTIMIZEPOLICY( $\lambda, \mathbf{a}, \mathbf{p}, \mathbf{W}$ )
2:   define  $\mathbf{a}' := \mathbf{a}$ 
3:   define  $\mathbf{p}' := \mathbf{p}$ 
4:   loop
5:     define  $\mathbf{a}'' := \mathbf{a}'$ 
6:      $\mathbf{p}' := \text{LEXMINIMIZEREGIMES}(\lambda, \mathbf{a}', \mathbf{p}', \mathbf{W})$ 
7:     define  $\mathbf{u} :=$  the payoffs induced by  $(\mathbf{a}', \mathbf{p}')$ 
8:     for all  $s \in S, a \in \mathbf{A}(\mathbf{W})(s)$  do
9:       if and  $(u^R(a, \lambda, \mathbf{u}) >_\lambda \mathbf{u}(s),$ 
10:         or  $(\gamma(a, \lambda^+, \mathbf{W}) > 0, \hat{u}^{APS}(a, \lambda, \mathbf{W}) >_\lambda \mathbf{u}(s)))$  then
11:          $\mathbf{a}''(s) := a$ 
12:       if  $\mathbf{a}'' \neq \mathbf{a}'$  then
13:          $\mathbf{a}' := \mathbf{a}''$  ▷ Continue updating the actions
14:       else
15:         for all  $s \in S$  do
16:           if  $\mathbf{p}'(s) = u^R(\mathbf{a}'(s), \mathbf{u})$  then
17:              $\mathbf{p}'(s) := R$  ▷ Make the pair canonical
18:           return  $(\mathbf{a}', \mathbf{p}')$  ▷ Return the optimal pair

```

Algorithm 6 Compute \tilde{B} for $N = 2$.

Require: $B(\mathbf{W}) \subseteq \mathbf{W}$

```

1: procedure  $\tilde{B}(\mathbf{W})$ 
2:   for all  $s \in S$  do
3:     define  $\mathbf{A}(\mathbf{W})(s) = \emptyset$ 
4:     for all  $a \in \mathbf{A}(s)$  do
5:       Compute  $C(a, \mathbf{W})$ 
6:       if  $C(a, \mathbf{W}) \neq \emptyset$  then
7:          $\mathbf{A}(\mathbf{W})(s) := \mathbf{A}(\mathbf{W})(s) \cup \{a\}$ 
8:   if  $\mathbf{A}(\mathbf{W})(s) = \emptyset$  for some  $s$  then
9:     return an empty correspondence
10:  define  $\mathbf{W}' := (\mathbb{R}^N)^S$                                 ▷ There are supportable actions
11:  define  $\lambda := (1, 0)$                                     ▷ Begin pointing due east
12:  define  $(\mathbf{a}, \mathbf{p})$  to be an arbitrary pair
13:  loop
14:    define  $(\mathbf{a}', \mathbf{p}') := \text{LEXOPTIMIZEPOLICY}(\lambda, \mathbf{a}, \mathbf{p}, \mathbf{W})$ 
15:    define  $\mathbf{u} :=$  the payoffs induced by  $(\mathbf{a}, \mathbf{p})$ 
16:     $\lambda' := \text{FINDNEXTDIRECTION}(\lambda, \mathbf{u}, \mathbf{W})$ 
17:     $\mathbf{W}' := \mathbf{W}' \cap \{\lambda' \cdot \mathbf{v} \leq \lambda' \cdot \mathbf{u}\}$           ▷ Intersect  $\mathbf{W}'$  with the new half space
18:    if  $\lambda$  points strictly north and  $\lambda'$  points weakly south then
19:      return  $\mathbf{W}'$                                           ▷ Completed a full revolution
20:    else
21:       $\lambda := \lambda', \mathbf{a} := \mathbf{a}'$                             ▷ Continue with the new direction

```

Algorithm 7 Compute \mathbf{V} to a tolerance ϵ in the metric d . Returns the approximation.

Require: $B(\tilde{\mathbf{W}}^0) \subseteq \tilde{\mathbf{W}}^0$ and $\mathbf{V} \subseteq \tilde{\mathbf{W}}^0$

```

1: procedure SOLVE( $\tilde{\mathbf{W}}^0, \epsilon$ )
2:   define  $k := 0$ 
3:   do
4:      $k := k + 1$ 
5:      $\tilde{\mathbf{W}}^k := \tilde{B}(\tilde{\mathbf{W}}^{k-1})$ 
6:   while  $d(\tilde{\mathbf{W}}^k, \tilde{\mathbf{W}}^{k-1}) > \epsilon$                             ▷ Stop when the movement is small
7:   return  $\tilde{\mathbf{W}}^k$ 

```

Next, given directions λ and $\tilde{\lambda}$, we define the $(\lambda, \tilde{\lambda})$ -line to be the subset of directions in Λ of the form $\cos(\theta)\lambda + \sin(\theta)\tilde{\lambda}$, where $\theta \in (0, 2\pi]$. We order $(\lambda, \tilde{\lambda})$ -line according to θ in this parameterization. We also extend the notion of test directions for the substitution (s, a, p) given the payoffs \mathbf{u} to be any direction satisfying (14). Legitimacy also extends to this setting. Finally, we redefine robust optimality in the many player setting by saying that \mathbf{u} is robustly optimal if it remains optimal in a neighborhood of λ . (Note that this definition is more restrictive than what is used in Section 4, where robustly optimal payoffs only had to remain optimal for perturbations in one direction.)

Algorithm 8 Update the direction by rotating towards $\tilde{\lambda}$. Returns the new direction of optimization and the direction in which payoffs move.

Require: \mathbf{u} is robustly optimal at λ

```

1: procedure ROTATEDIRECTION( $\lambda, \tilde{\lambda}, \mathbf{u}, \mathbf{W}$ )
2:   for all  $s \in S, a \in \mathbf{A}(\mathbf{W})(s), p \in \{R\} \cup C(a, \mathbf{W})$  do
3:     for all test directions  $\lambda''$  for  $(s, a, p)$  and  $\mathbf{u}$  in the  $(\lambda, \tilde{\lambda})$ -line do
4:       if  $\lambda''$  is legitimate and a smaller rotation than  $\lambda'$  then
5:          $\lambda' := \lambda''$ 
6:          $d := u(s, a, p, \mathbf{u}) - \mathbf{u}(s)$   $\triangleright$  the direction in which  $(s, a, p)$  moves payoffs
7:   return  $(\lambda', d)$ 

```

Algorithm 9 Compute a randomly chosen face of $\tilde{B}(\mathbf{W})$. Return the direction and the corresponding half space.

Require: $B(\mathbf{W}, \hat{\Lambda}) \subseteq \mathbf{W}$

```

1: procedure FINDFACE( $\mathbf{W}$ )
2:   define  $\lambda^0$  randomly
3:   define  $(\mathbf{a}, \mathbf{r}) := \text{OPTIMIZEPOLICY}(\lambda^0, \mathbf{W})$ 
4:   define  $\mathbf{p} \in \mathbf{P}(\mathbf{a}, \mathbf{W})$  to be min-max for  $(\mathbf{a}, \mathbf{r}, \mathbf{W})$ 
5:   define  $\mathbf{u} :=$  payoffs induced by  $(\mathbf{a}, \mathbf{p})$ 
6:   for  $n = 1, \dots, N - 1$  do
7:     define  $\tilde{\lambda}^n$  randomly to be orthogonal to  $\{\lambda^0\} \cup \{d^l | l = 1, \dots, n - 1\}$ 
8:     define  $(\lambda^n, d^n) := \text{ROTATEDIRECTION}(\lambda^{n-1}, \tilde{\lambda}^n, \mathbf{u}, \mathbf{W})$ 
9:   define  $H := \{\mathbf{v} | \lambda^{N-1} \cdot \mathbf{v} \leq \lambda^{N-1} \cdot \mathbf{u}\}$ 
10:  return  $(\lambda^{N-1}, H)$ 

```

Algorithm 10 Approximate $\tilde{B}(\mathbf{W})$, given an incumbent set of directions $\hat{\Lambda}$. Returns a new approximation and a new set of directions.

Require: $B(\mathbf{W}, \hat{\Lambda}) \subseteq \mathbf{W}$

```

1: procedure  $\tilde{B}(\mathbf{W}, \hat{\Lambda}, L)$ 
2:   define  $\mathbf{W}' := (\mathbb{R}^N)^S$ 
3:   define  $\hat{\Lambda}' := \emptyset$ 
4:   for all  $\lambda \in \hat{\Lambda}$  do
5:      $(\mathbf{a}, \mathbf{r}) := \text{OPTIMIZEPOLICY}(\lambda, \mathbf{W})$ 
6:     define  $\mathbf{p} \in \mathbf{P}(\mathbf{a}, \mathbf{W})$  to be min-max for  $(\mathbf{a}, \mathbf{r}, \mathbf{W})$ 
7:     define  $H := \{\mathbf{v} \mid \mathbf{v} \cdot \lambda' \leq x(\lambda', \mathbf{p}, \mathbf{W})\}$ 
8:     if  $\mathbf{W}'$  and  $\mathbf{W}' \cap H$  do not have the same local binding frontier then
9:        $\hat{\Lambda}' := \hat{\Lambda}' \cup \{\lambda\}$ 
10:       $\mathbf{W}' := \mathbf{W}' \cap H$ 
11:   define  $K := |\hat{\Lambda}'|$ 
12:   for  $k = 1, \dots, L - K$  do
13:      $(\lambda, H) := \text{FINDFACE}(\mathbf{W})$ 
14:     if  $\mathbf{W}'$  and  $\mathbf{W}' \cap H$  do not have the same local binding frontier then
15:        $\mathbf{W}' := \mathbf{W}' \cap H$ 
16:        $\hat{\Lambda}' := \hat{\Lambda}' \cup \{\lambda\}$ 
17:   return  $(\mathbf{W}', \hat{\Lambda}')$ 

```

Algorithm 10 can be combined with an analogue of Algorithm 7 to approximate \mathbf{V} when $N > 2$.

B Online appendix:

A repeated game with infinitely many extreme equilibrium payoffs

In this Online Appendix, we give additional details on the three-player example from Section 5.1, depicted in Figure 4, that has infinitely many extreme equilibrium payoffs. First we construct a self-generating set that turns out to be V . We will then argue that this is in fact the equilibrium payoff set.

B.1 The equilibrium payoff set

Recall that only four action profiles can be played in equilibrium, which induce payoffs $(4, 4, 4)$ and permutations of $(8, 8, 0)$. Note that $(4, 4, 4)$ is one of the equilibrium payoffs.

We will generate two sequences of payoffs $\{u^l\}_{l=0}^\infty$ and $\{v^l\}_{l=0}^\infty$. The payoff u^0 corresponds to u in the right panel of Figure 5, and the subsequent sequence is the sequence of extreme payoffs that move counter-clockwise around the frontier. The payoff v^0 corresponds to v in the right panel of Figure 5, and the sequence of extreme points moves clockwise around the frontier. Aside from $(4, 4, 4)$, the extreme equilibrium payoffs are permutations of points in these sequences.

Every v^l is generated the same way, by randomizing between u^l and $(4, 4, 4)$, to make the incentive constraint for player 1 bind, i.e.,

$$v^l = \left(6 - \frac{1}{u_2^l - 4}, 6 + \frac{u_1^l - 4}{u_2^l - 4}, 3\right) = \frac{1}{2}(8, 8, 0) + \frac{1}{2} [\beta^l(3, u_1^l, u_2^l) + (1 - \beta^l)(4, 4, 4)], \quad (19)$$

where

$$\beta^l = \frac{2}{u_2^l - 4}. \quad (20)$$

The payoffs u^l are generated in three different ways. First, the permutations of u^0 , i.e., the extreme points on the efficient frontier comprise a self-generating set and are generated according to

$$u^0 = \left(\frac{11}{2}, \frac{15}{2}, 3\right) = \frac{1}{2}(8, 8, 0) + \frac{1}{2} \left[\frac{1}{4} \left(3, \frac{11}{2}, \frac{15}{2}\right) + \frac{3}{4} \left(3, \frac{15}{2}, \frac{11}{2}\right) \right],$$

i.e., by playing (B, B, C) for one period, followed by randomizing over two other efficient extreme payoffs to make the incentive constraint (player 1's in this case) bind.

Given u^0 , we can generate v^0 according to (19) and (20), which turns out to be $v^0 = (40/7, 45/7, 3)$, with $\beta^0 = 4/7$. The payoff u^1 is then generated by playing (B, B, C) for one period, followed by randomization between two permutations of v^0 :

$$u^1 = \left(\frac{11}{2}, \frac{99}{14}, 3 \right) = \frac{1}{2}(8, 8, 0) + \frac{1}{2} [\alpha^1(3, v_2^0, v_1^0) + (1 - \alpha^1)(3, v_1^0, v_2^0)].$$

where $\alpha^1 = 3/5$ is again chosen to make player 1's incentive constraint bind. Finally, the rest of the u^l sequence for $l \geq 2$ is generated according to

$$u^l = \left(6 - \frac{1}{v_2^{l-2} - 4}, 6 + \frac{v_1^{l-2} - 4}{v_2^{l-2} - 4}, 3 \right) = \frac{1}{2}(8, 8, 0) + \frac{1}{2} [\alpha^l(3, v_1^{l-2}, v_2^{l-2}) + (1 - \alpha^l)(4, 4, 4)],$$

where

$$\alpha^l = \frac{2}{v_2^{l-2} - 4}$$

is again chosen to make player 1's constraint bind.

Finally, these sequences converge to the accumulation points in Figure 5, which are permutations of $((9 + \sqrt{5})/2, (11 + \sqrt{5})/2, 3)$. These payoffs, together with $(4, 4, 4)$, comprise another self-generating set, where

$$\left(\frac{9 + \sqrt{5}}{2}, \frac{11 + \sqrt{5}}{2}, 3 \right) = \frac{1}{2}(8, 8, 0) + \frac{1}{2} \left[\alpha^* \left(3, \frac{9 + \sqrt{5}}{2}, \frac{11 + \sqrt{5}}{2} \right) + (1 - \alpha^*)(4, 4, 4) \right],$$

where $\alpha^* = 3 - \sqrt{5}$.

B.2 Feasible set

We next argue that the equilibrium payoff set is the convex hull of the points constructed heretofore. The analysis consists of several steps. First, since only these three action profiles can possibly be played in equilibrium, we know that the equilibrium payoff set must be contained in the triangular pyramid with peak at $(4, 4, 4)$ and base corners which are permutations of $(0, 8, 8)$. Thus, the pyramid "points" in the direction $(-1, -1, -1)$. In the sequel, we refer to this as the "feasible set."

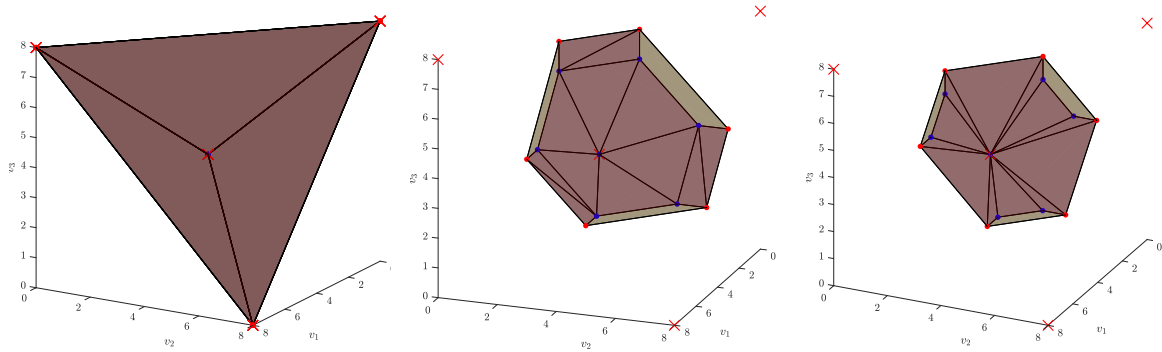


Figure 6: Three different bounds on the equilibrium payoff set. Left: The convex hull of the flow payoffs. Center: The left set less the payoffs that are below the threat point. Right: Additional payoffs removed to create the set \widetilde{W}_0 . Faces that coincide with incentive constraints are colored tan.

B.3 Equilibrium threats

Clearly, the equilibrium threat point \underline{v} must be less than 4 (since the Nash equilibrium is certainly an equilibrium payoff). Thus, from the definition of \widetilde{B} , the only way that player 3 can obtain a lower payoff is if (B, B, C) is played in the first period, with a flow payoff of $(8, 8, 0)$. Moreover, any payoff we generate with this action must be weakly above $(8, 8, 0)$ in the direction $(0, 0, -1)$, and therefore it must be generated with a binding incentive constraint. But players 1 and 2 are playing myopic best responses at (B, B, C) , so the only relevant incentive constraint is player 3's. Plugging in the specified payoffs and discount factor, we conclude that

$$\begin{aligned} \underline{v} &= \frac{1}{2}3 + \frac{1}{2}\underline{v} \\ \iff \underline{v} &= 3. \end{aligned}$$

B.4 The efficient frontier

In addition, we claim that no equilibrium payoff can lie above the plane that contains $(4, 4, 4)$, $(11/2, 15/2, 3)$, and $(3, 15/2, 11/2)$, i.e., with level x and direction λ such that

$$\begin{aligned} x &= \lambda \cdot (4, 4, 4) = \lambda \cdot \left(\frac{11}{2}, \frac{15}{2}, 3 \right) = \lambda \cdot \left(3, \frac{15}{2}, \frac{11}{2} \right) \\ \iff x &= -52, \quad \lambda = (-7, 1, -7). \end{aligned} \tag{21}$$

(The permutations of this statement also apply when we give the low payoff of 3 to player 1 or player 2). The reason is as follows. Consider maximizing payoffs in this direction. The optimal level must be at least -52 , which is that of $(4, 4, 4)$, the Nash equilibrium. But the flow payoff $(8, 0, 8)$ has level -112 , which is strictly below the Nash level, and hence cannot generate the optimal payoff. So, we may ask, what is the highest level that can be generated by $(0, 8, 8)$ or $(8, 8, 0)$? We will consider the former, and the case for the latter is symmetric. In this direction, the flow payoffs $(0, 8, 8)$ are maximal among all payoffs in the feasible pyramid, so that the minimal regime must be APS. To satisfy incentive compatibility, the continuation value of player 1 must be at least 6. Player 3's continuation value must be at least 3 from incentive compatibility. Finally, the sum of the payoffs is at most 16 (from feasibility). It follows that the highest level that can be attained in this direction is

$$\lambda \cdot \left(\frac{1}{2}(0, 8, 8) + \frac{1}{2}(6, 7, 3) \right) = -52.$$

Moreover, the permutations of $u^0 = (11/2, 15/2, 3)$ are a self-generating set. In particular,

$$\left(\frac{11}{2}, \frac{15}{2}, 3 \right) = \frac{1}{2}(8, 8, 0) + \frac{1}{2} \left[\frac{1}{4} \left(3, \frac{11}{2}, \frac{15}{2} \right) + \frac{3}{4} \left(3, \frac{15}{2}, \frac{11}{2} \right) \right].$$

We conclude that these are all extreme equilibrium payoffs (being at the corners of the hyperplanes in (21), the minimum payoff constraints, and the efficient frontier. Moreover, the convex hull of these points is the set of Pareto efficient payoffs.

The equilibrium payoff set must lie inside the polyhedron defined by the hyperplanes in (21), the constraints $v_i \geq 3$ for all i , and the constraint $\sum_i v_i \leq 16$. We denote this set by \widehat{W} .

B.5 Structure of minimal regimes

Note that since (A, A, A) is a Nash equilibrium, no matter what feasible set W we consider, as long as $V \subseteq W$, the recursive regime will be minimal for (A, A, A) , i.e., $x(\lambda, (A, A, A), W) = \lambda \cdot (4, 4, 4)$.

In addition, we claim that whenever (B, B, C) is maximal, the minimal regime must be APS. For we already know that the payoffs $(4, 4, 4)$, and permutations of u^0 can be generated. This pins down the optimal level exactly in all directions except those which are in the interior of $\widehat{\Lambda}_1 = \text{co}\{(-7, -7, 1), (-7, 1, -7), (-1, 0, 0)\}$, or permutations thereof. (Outside of these sets of directions, an optimal payoff must be one of the aforementioned extreme points). $\widehat{\Lambda}_i$ denotes the permutations of these directions, where we give the weight

-1 to a different player. For directions in $\widehat{\Lambda}_1$, it is easy to argue that $(0, 8, 8)$ is higher than all other payoffs in the feasible triangle, so that necessarily the minimal regime for (C, B, B) is *APS*. In addition, either $(8, 0, 8)$ or $(8, 8, 0)$ is minimal among all feasible payoffs, so that the corresponding minimal regimes are all recursive, and hence these action profiles cannot be maximal in directions in $\widehat{\Lambda}_1$.

This means that for directions in $\widehat{\Lambda}_1$, the optimal level is simply given by $\widehat{x}^{APS}((C, B, B), \lambda)$, and we can reduce the computation of \widetilde{B} to simply computing the sets $C(a)$ (where we drop the argument \mathbf{W} for notational simplicity) for each $a \neq (A, A, A)$. Specifically, for all W contained within \widehat{W} ,

$$\widetilde{B}(W) = \text{co} \left(\{(4, 4, 4)\} \cup_{a \in \{(C, B, B), (B, C, B), (B, B, C)\}} C(a) \right).$$

We also note for future reference that if $v \in C(a)$ and $v' \in C(a')$, then there is no direction in which both v and v' are both maximal. This comes from the fact that the sets of directions $\widehat{\lambda}_i$ are disjoint.

B.6 Two more bounds on binding payoffs

The focus of the analysis now shifts to the sets $C(a)$ where $a_i = C$ and $a_{-i} = (B, B)$ for some $i \in \{1, 2, 3\}$. Ultimately we will construct a sequence of iterates using the \widetilde{B} operator that converge to V and demonstrate that the limit set has infinitely many extreme points. Before doing so, we will slightly refine the approximation so that the sequence converges in an orderly manner.

It is straightforward that any $v \in C(B, B, C)$ must satisfy $v_j \geq 11/2$ for $j = 1, 2$. This follows from the fact that the flow payoff is 8 and the minimal equilibrium payoff is 3.

In addition, consider the direction $(-7, -7, -29)$. We claim that no payoff in $C(a)$ can be above the level -172 . This level is attained by the Nash payoff $(4, 4, 4)$ and also by $(8, 8, 0)$ with maximal continuation payoffs w such that $w_3 \geq 6$ and $w \cdot (1, -7, -7) \leq -52$. In particular, the solution is attained by payoffs

$$v^0 = \left(\frac{45}{7}, \frac{40}{7}, 3 \right) = \frac{1}{2}(8, 8, 0) + \frac{1}{2} \left[\frac{4}{7} \left(\frac{11}{2}, 3, \frac{15}{2} \right) + \frac{3}{7}(4, 4, 4) \right].$$

Note that since the permutations of $u^0 = (15/2, 11/2, 3)$ are already known to be part of a self-generating set, we know that v^0 are also equilibrium payoffs, and hence the plane with level -172 in direction $(-7, -7, -29)$ is a supporting hyperplane of V . In fact, it intersects V in a face that contains $(45/7, 40/7, 3)$, $(40/7, 45/7, 3)$, and $(4, 4, 4)$.

We thus conclude that $C(B, B, C)$ is contained within the trapezoid of payoffs v defined

by $v_3 = 3$, $\sum_i v_i \leq 16$, $v_1 \geq 11/2$, $v_2 \geq 11/2$, and $-29v_3 - 7(v_1 + v_2) \leq -172$. This trapezoid is denoted by \tilde{C}_3^0 (\tilde{C}_i^k will later denote a sequence of minimal payoff sets for player i). We note for future reference that \tilde{C}_3^0 is the convex hull of the payoffs $(15/2, 11/2, 3)$ and

$$w^0 = (11/2, 93/14, 3)$$

and the permutations obtained by interchanging the payoffs of players 1 and 2. Note that the payoff w^0 is at the intersection of the bounds $v_1 \geq 11/2$, $v_3 = 3$, and $-29v_3 - 7(v_1 + v_2) \leq -172$. We correspondingly define the sets \tilde{C}_1^0 and \tilde{C}_2^0 by permuting players' payoffs. We let

$$\tilde{W}^0 = \text{co} \left(\{(4, 4, 4)\} \cup_{i=1,2,3} \tilde{C}_i^0 \right),$$

which will serve as the initial set for the sequence we generate in the next and final subsection.

B.7 The sequence $\{\tilde{W}^k\}$

We now analyze the sequence of sets produced by iterative application of \tilde{B} to \tilde{W}^0 . The critical issue is to determine the shape of the sets \tilde{C}_i^{k+1} . In the following discussion, we take the perspective of minimum payoffs for player $i = 3$, but the case is symmetric for the other players.

We will argue that at iteration $k \geq 0$, the set \tilde{C}_3^k is the convex hull of the points $\{u^l\}_{l=0}^k$, $\{v^l\}_{l=0}^{k-1}$, the payoff w^k defined as above for $k = 0$ and by

$$w^k = \frac{1}{2} (8, 8, 0) + \frac{1}{2} \left[\frac{2}{w_2^{k-1} - 4} (3, w_1^{k-1}, w_2^{k-1}) + \left(1 - \frac{2}{w_2^{k-1}}\right) (4, 4, 4) \right]$$

for $k \geq 1$, and the permutations thereof obtained by interchanging the payoffs of players 1 and 2. The base case has already been given for $k = 0$ in the previous subsection.

Let us take as an inductive hypothesis that the set \tilde{W}^k is comprised of the following edges: First, there are edges between the payoffs in \tilde{C}_i^k . Second, there are edges between permutations of w^0 that are in different \tilde{C}_i^k sets. Finally, there are edges that connect all of the payoffs in \tilde{C}_i^0 with the Nash payoff $(4, 4, 4)$.

Given this inductive hypothesis for $k - 1$, we can easily compute the set \tilde{C}_3^k . First, we compute the intersection of \tilde{W}^{k-1} with the plane $w_3 = 6$ to find the extreme binding continuation values for player 3. We then average these payoff with the flow payoff $(8, 8, 0)$ to obtain \hat{C}_3^k . The intersections with the $w_3 = 6$ plane must lie on edges of \tilde{W}^{k-1} that have one point with v_3 higher than 6 and another point with v_3 less than 6. There are three kinds of such edges that have intersections with the $w_3 = 6$ plane: The edges between

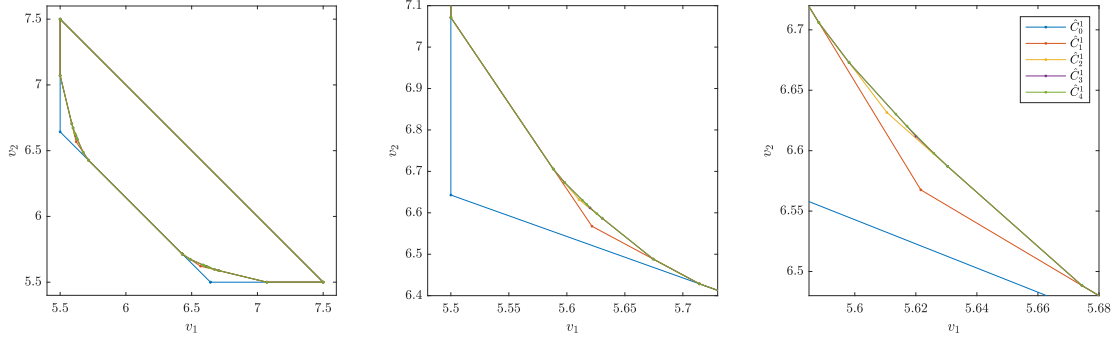


Figure 7: The sets \tilde{C}_i^k for $k \in \{0, 1, 2, 3\}$. At every iteration, four new faces are added to each \tilde{C}_i^k . The right two panels show the left-hand corner of the set at different levels of magnification.

permutations of u^0 , e.g., $(3, 11/2, 15/2)$ and $(3, 15/2, 11/2)$, which will generate the point u^0 ; The edges between permutations of w^0 (when $k = 1$) or between permutations of v^0 , which generate u^1 (when $k > 1$); and the edges between one of the payoffs whose permutation is in $\{u^l\}_{l=0}^{k-1} \cup \{v^l\}_{l=0}^{k-2} \cup \{w^{k-1}\}$, and the Nash payoff $(4, 4, 4)$, which generate a payoff v^l , u^{l+2} , or w^k , respectively. From the inductive hypothesis, all of these intersections must result in new extreme payoffs of \tilde{C}_3^k .

As an example, when $k = 1$, the payoffs generated will be u^0 , u^1 , v^0 and w^1 , as well as their permutations when we swap the payoffs of players 1 and 2. The first five elements of the \tilde{C}_3^k sequence are depicted in Figure 7.

Finally, it remains to argue that the inductive hypothesis will be true for k . The new payoffs generated in \tilde{C}_3^k can be divided into those where player 1's payoff is at least 6, and those where player 1's payoff is less than 6. Focus for now on the former. These payoffs are maximal for directions that are convex combinations of $(-7, -7, -29)$, $(1, -7, -7)$, and $(0, 0, -1)$. Note that we have already characterized supporting hyperplanes of V in these three directions, which are also necessarily supporting hyperplanes of \tilde{W}^k . For directions other than $(0, 0, -1)$ and $(1, -7, -7)$, the only other optimal payoff is $(4, 4, 4)$, so that edges on supporting hyperplanes in these directions will be composed of either two payoffs in \hat{C}_3^k with $v_1 \geq 3$, or one of the payoffs in \tilde{C}_3^k with $v_1 \geq 6$ and $(4, 4, 4)$. For the direction $(-7, -7, -29)$, when $k = 0$, the permutations of w^0 and the Nash equilibrium are all optimal, so there is one additional edge, between the permutations of w^0 . For $k \geq 1$, it is the permutations of v^0 and $(4, 4, 4)$ that are optimal in this direction. Finally, for the direction $(0, 0, -1)$, all of the payoffs in \tilde{C}_3^k are optimal, so edges here will be between points in \hat{C}_3^k . A similar analysis applies to other extreme points, so that the inductive hypothesis is true

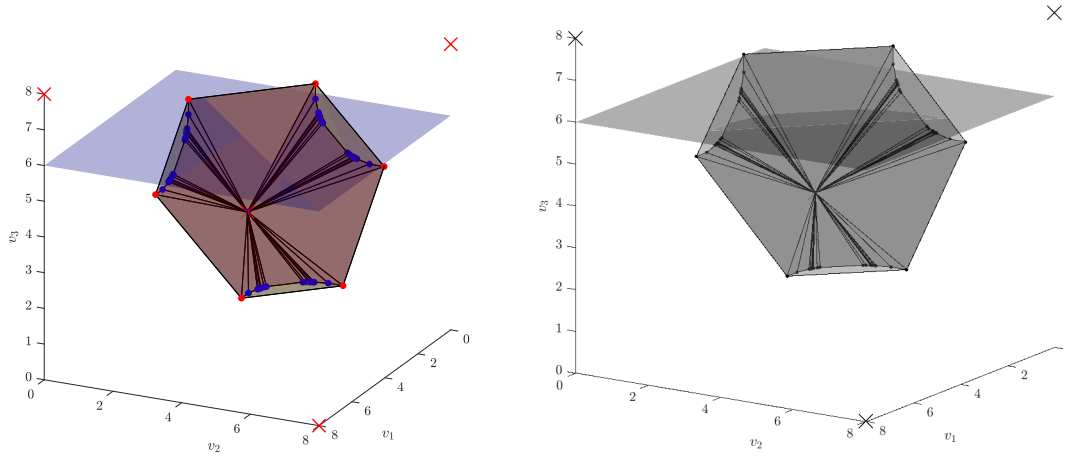


Figure 8: Two views of \widetilde{W}^5 . Flow payoffs are marked with red crosses. Efficient extreme payoffs are red dots, inefficient extreme payoffs are blue dots. The minimum incentive compatible continuation value for player 3 (whose payoff is on the z axis) is a blue plane. The intersection of this set with the payoff set, contracted towards the payoff $(8, 8, 0)$, generates the bottom flat of V .

for k .

Note that at the k th round, we drop the permutations of w^{k-1} , but add the permutations of w^k , v^{k-1} , and u^k . Thus, the number of extreme points increases by 12 on every iteration. Moreover, the points u^k and v^{k-1} , once added, are never dropped, so the set of extreme points increases without bound. In the limit, the sequence w^k converges to the accumulation point w^* , which is generated according to

$$w^* = \frac{1}{2}(8, 8, 0) + \frac{1}{2}[\alpha^*(3, w_1^*, w_2^*) + (1 - \alpha^*)(4, 4, 4)].$$

The weight α^* must solve

$$0 = -\frac{1}{8}\alpha^3 + \frac{1}{2}\alpha^2 + \alpha - 1.$$

This equation has three real roots, only one of which is between 0 and 1, which is $\alpha^* = 3 - \sqrt{5}$.

The resulting payoff is

$$w^* = \left(\frac{9 + \sqrt{5}}{2}, \frac{11 + \sqrt{5}}{2}, 3 \right).$$

The set \widetilde{W}^5 is depicted in Figure 8. At this resolution, this set is indistinguishable from V .

As a final note, while the analysis of this game is involved, in many ways it is the simplest example possible. Four is the minimum number of equilibrium action profiles such that the equilibrium payoff set is full dimension, which is necessary for the number of extreme points to be unbounded. The incentive constraints are also quite simple: One action profile is a Nash equilibrium, and for each other equilibrium action profile, only a single incentive constraint binds, that of the player whose payoff is being minimized.

C Online appendix:

Additional examples

C.1 Two-player two-state Prisoners' Dilemma

This example illustrates the utility of the test directions in iteratively computing optimal levels. There are two states, L and R , and the stage game in each state is a Prisoners' Dilemma with the payoffs in Figure 9. The probability of staying in the same state is $1/3$ if the players take the same action, and it is $1/2$ if the players take different actions. The discount factor is $\delta = 2/3$.

We computed the sequence $\widetilde{\mathbf{W}}^k$ until the Hausdorff distance between successive iterations was less than 10^{-8} . The computation took 0.37 seconds. The sequence of payoff correspondences is depicted in Figure 10. The final payoff set for the left state has six extreme points, and the right state has four.

It turns out that the equilibrium threat point is generated by a policy that plays (D, D) in both states in the recursive regime. The resulting threat point is

$$(\underline{v}_i(L), \underline{v}_i(R)) = \left(\frac{8}{11}, \frac{14}{11} \right).$$

The utilitarian efficient payoffs that are optimal in the direction $(1, 1)$ are generated by playing (C, C) in both states in the recursive regime. The resulting symmetric payoffs are $19/11$ in the left state and $25/11$ in the right state.

We may ask, how will the optimal policy change as the direction rotates clockwise from $(1, 1)$? A natural conjecture, which turns out to be correct, is that the optimal policy will change by switching from (C, C) to (D, C) in some state. But should this switch occur first in the left state or the right state? Both switches would move flow payoffs in the same direction of $(1, -3)$. But switching from (C, C) to (D, C) in state L would increase the probability of staying at $s = L$ where payoffs are lower, whereas switching in $s = R$ leads to a higher probability $s = R$, where payoffs are higher. Thus, less surplus is burnt by switching when $s = R$, and indeed this is the correct substitution.

	$s = L$			$s = R$	
a_1/a_2	D	C		D	C
C	$(-1, 2)$	$(1, 1)$		$(1, 4)$	$(3, 3)$
D	$(0, 0)$	$(2, -1)$		$(2, 2)$	$(4, 1)$

Figure 9: A two-state Prisoners' Dilemma. The probability of remaining in the same state is $1/3$ if $a_i = a_j$, and otherwise it is $1/2$.

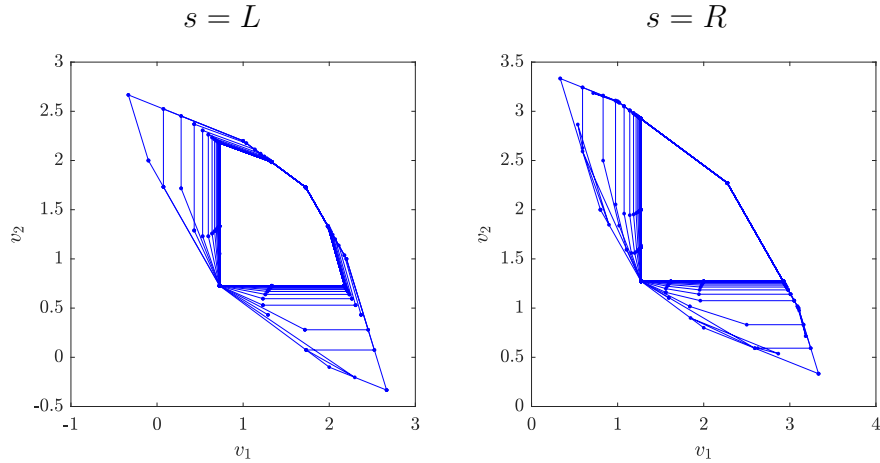


Figure 10: The sequence of correspondences generated by the max-min-max operator.

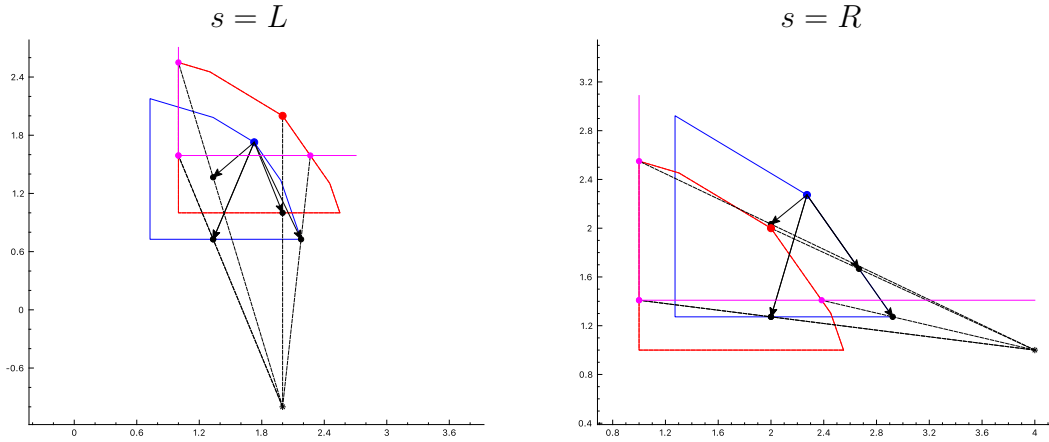


Figure 11: Test directions for (D, C) , relative to the symmetric efficient payoffs.

Our algorithm resolves this question mechanically using the test directions, which are depicted in Figure 11. The flow payoffs from (D, C) are depicted with black stars, payoff sets in blue, expected continuation payoff sets in red, binding incentive constraints in magenta. The test directions are black arrows. The shallowest legitimate test directions point along the frontier, and are generated by (D, C) in the right state. Note that there is a tie between the recursive and APS substitutions: both move payoffs along the frontier,²³ although only the recursive substitution is “lexicographically legitimate” in the sense described in Section 4.2.4.

²³In fact, there is a three-way tie, since there is a binding substitution for (C, C) in the right state that moves payoffs in the same direction. This test direction is not depicted in Figure 11.

C.2 A three-player contribution game

We implemented the stochastic algorithm for three players as part of the aforementioned SGSolve package.²⁴ Let us illustrate the algorithm with two examples. The first example is a simple contribution game: $N = 3$, $S = \{1, 2\}$, $\mathbf{A}_i(s) = \{0, 1\}$, and

$$u_i(a, s) = 2 \sum_{j=1}^N a_j - 3a_i + 20s.$$

The transition probabilities are $\pi(s|a) = 1/2$ for every s and a , and $\delta = 2/3$. The stage game in each state is effectively a three-player Prisoners' Dilemma.

This example illustrates how our algorithm can solve for the equilibrium payoff exactly. We initialized the algorithm with 214 directions that are distributed around the unit sphere. We used the convergence criterion that no directions were added or dropped between iterations, and the Hausdorff distance between consecutive iterations was less than 10^{-8} . Due to the stochastic nature of the algorithm, its performance varies on each run. On one series of five runs, the algorithm finished with 9 directions three times, and 10 directions the other two. Over the course of one of the runs that terminated with 9 face directions, the algorithm added 72 endogenous directions and dropped 277. In all cases, the algorithm converged in 45 iterations and took between 2.85 and 3.11 seconds.

One can analytically verify that the equilibrium payoff correspondence for this game has exactly 9 face directions. Thus, in the runs where the algorithm terminated with 9 directions, it correctly identified the structure of equilibrium payoffs, which are depicted in Figure 12. All sixteen action profiles can be sustained. The efficient points are generated by always playing $a = (1, 1, 1)$ in both states. There is also an inefficient point which corresponds to the Nash equilibrium $a = (0, 0, 0)$. The remaining points are generated by playing permutations of $a = (1, 1, 0)$ and $(1, 0, 0)$. The face in which a player's payoff is minimized is attained by $a_i = 1$ and $a_{-i} = (0, 0)$.

For comparison, we solved this game using our implementation of the JYC algorithm with the same set of 214 initial directions. The same tolerance was achieved in 49 iterations and 3 minutes and 38.45 seconds. So, the JYC code is between one and two orders of magnitude slower. All of our previous caveats still apply, but we find this suggestive that the stochastic max-min-max algorithm is significantly more efficient.

A natural question is, which features of the max-min-max operator explain the difference in performance? We also ran a version of our algorithm with the same 214 initial directions,

²⁴We note that the graphical interface currently only works for two-player games, but the three-player routines are part of the callable library.

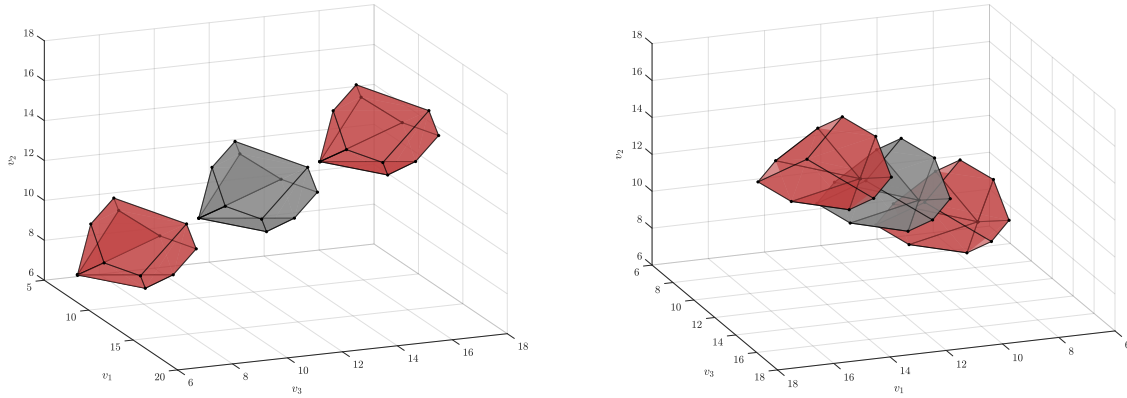


Figure 12: Equilibrium payoffs for the contribution game. The equilibrium payoff correspondence is in red, and the expected equilibrium payoffs are in gray.

but where we set $\widehat{\Lambda}^k = \widehat{\Lambda}^0$ for all k , i.e., the set of directions is held fixed. In this case, the algorithm converged in 44 iterations and 3.53 seconds. This suggests that most of the efficiency gain comes from using the max-min-max level rather than APS. The endogeneity of directions, however, leads to a tight limit set.

C.3 Three-player risk sharing and partial formal insurance

We solved a three-player risk-sharing game, as in Section 4. Each player now has an endowment e_i and their actions are vectors that specify how much they transfer to each other player. For this particular simulation, we used $u(c) = \sqrt{c}$, the endowment grid is $E = \{0, 0.5, 1\}$, and endowment distribution is independent across periods and uniform over endowment profiles that sum to 1. The discount factor is $\delta = 0.6$. For this simulation, we capped the algorithm at 300 directions and iterated until a convergence threshold of 10^{-8} . The algorithm converged in 68 seconds and 33 iterations. Over the course of the computation, 492 endogenous directions were added and 551 redundant directions were dropped. The computed expected equilibrium payoff set is depicted in the left-hand panel of Figure 13.

As a simple application, we used our algorithms to investigate the following question: What happens to equilibrium risk sharing and social welfare if the players can write formal insurance contracts? If all of the players can write a formal full insurance contract, so that they equally share their collective resources, then the welfare implications are obvious: The sum of the agents' surpluses must weakly increase. If only two of the three players can write such a contract, however, the implications are ambiguous. Suppose that players 2 and 3 write such a contract, so that $c_2 = c_3$ and each consumes half of their total endowment net of transfers to player 1, and their transfers to player 1 are chosen to maximize the joint

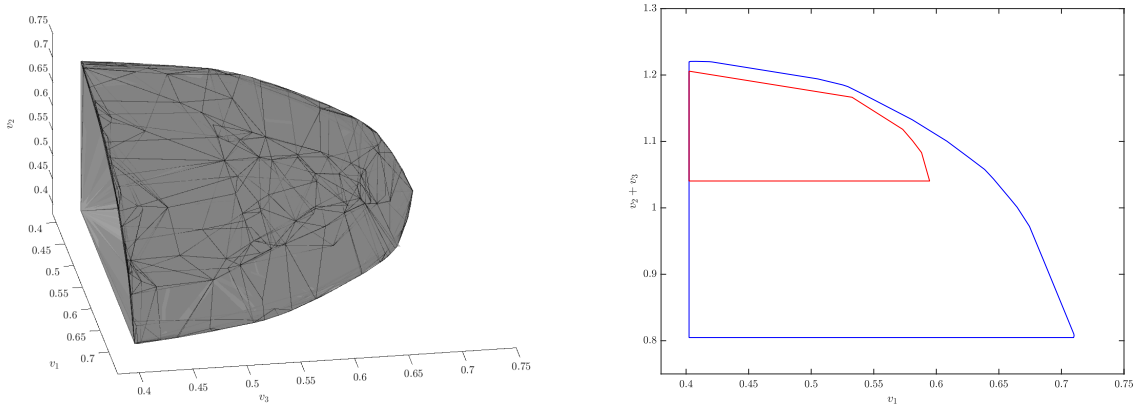


Figure 13: Risk sharing with three players. Left: Expected equilibrium payoffs. Right: Achievable $(v_1, v_2 + v_3)$ pairs. Non-cooperative play is in blue, and cooperation between players 2 and 3 is in red.

surplus $v_2 + v_3$. On the one hand, players 2 and 3 should be better off, because they are always guaranteed a minimal level of insurance, so that their autarkic payoffs are higher with such a contract than without. On the other hand, the higher autarkic payoffs tighten incentive constraints and may reduce risk sharing with player 1.

We used the two-player algorithm of Section 4 to investigate what would happen if players 2 and 3 behave cooperatively to maximize their joint surplus. Expected equilibrium payoffs are plotted in red in the right-hand panel of Figure 13. For comparison, the blue curve represents the possible $(v_1, v_2 + v_3)$ pairs in the game where players 2 and 3 behave non-cooperatively. The threat payoff for players 2 and 3 is clearly higher with the contract: their minimum joint surplus is approximately 0.805 in the non-cooperative case, and approximately 1.04 when they cooperate. A striking result is that the tightening of incentive constraints appears to overwhelm the benefits of additional risk sharing, and the Pareto frontier when players 2 and 3 cooperate is strictly below the Pareto frontier when they behave non-cooperatively. Thus, the example illustrates how partial insurance contracts may lead to lower social welfare.

C.4 Lower bounds on payoffs

Recall the two-state risk sharing example of Section 4 with $\delta = 0.7$. Figure 14 compares the upper and lower bounds on equilibrium payoffs, which are blue and red respectively. The lower bound was computed with $\epsilon = 0.005$. The red dotted line corresponds to the expansion of the lower bound by ϵ in every direction, which is contained in the correspondence that would be produced by the APS operator. The payoffs that induce the upper and lower

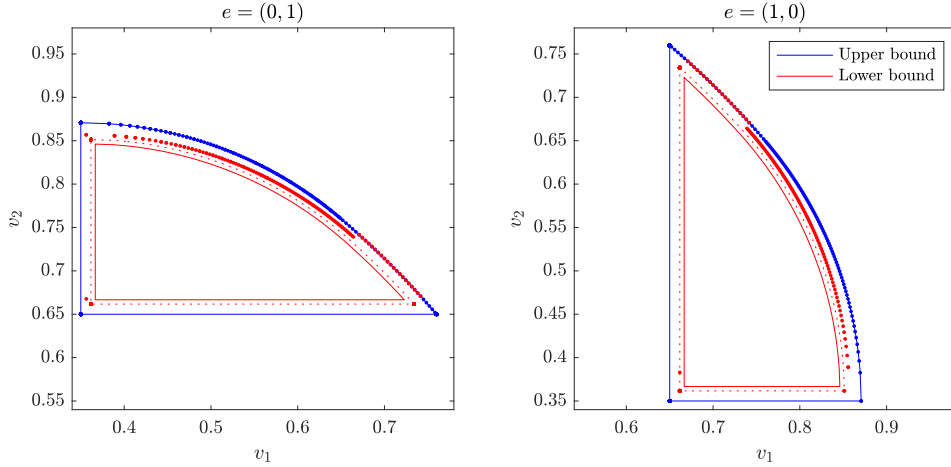


Figure 14: Lower bounds on \mathbf{V} for $\epsilon = 0.005$. Dots represent the actual payoffs used to generate the bounds.

bounds are represented as dots. Note that the distance from the payoffs to the lower bound set varies depending on the direction of the bound. This distance is greater when more states have a minimal regime that is recursive. When both states are binding, such as the payoffs that approximate the threat point, the penalty is ϵ in both states. When one state is binding, such as when we maximize one player's payoff, the penalty in the binding state is still ϵ , but the penalty in the recursive state is $\epsilon/(1 - \delta/2)$. When both states are recursive, which is when the direction is close to maximizing the sum of payoffs, the penalties in both states are $\epsilon/(1 - \delta)$. At directions where the minimal regimes change, the penalties (and hence the level of the optimal payoffs) change discontinuously.

The computation depicted in Figure 14 used a relatively large value for ϵ for visual effect. When ϵ is small, the distance between the outer and inner bounds shrinks as well and appears to go to zero. For example, we computed upper bounds on \mathbf{V} and \mathbf{V}^ϵ to a tolerance of 10^{-7} when $\epsilon = 10^{-6}$. The extreme points of the upper bound on \mathbf{V} are all within 10^{-6} of the bounds for \mathbf{V}^ϵ , so that the lower and upper bounds are indistinguishable (up to the tolerance for computing the extreme points of the upper bound).

D Online appendix:

Connections to linear programming and dynamic programming

To the student of linear programming, our procedure may evoke the simplex algorithm and sensitivity analysis. The choice of (\mathbf{a}, \mathbf{r}) bears a resemblance to the choice of a basis, and our use of test directions and optimization is similarly reminiscent of parametric programming in the theory of linear programming (see Dantzig and Thapa, 2006, for a comprehensive treatment). In this section, we attempt to elucidate the connection.

Suppose we were not concerned with incentives at all and simply wanted to compute the feasible payoff correspondence \mathbf{F} , i.e., payoffs that can be obtained with some pure-strategy profile starting in state s (still allowing public randomization). For a fixed direction λ , the problem of computing the optimal levels

$$x(s, \lambda) = \max\{\lambda \cdot v \mid v \in \mathbf{F}(s)\}$$

is a Markov decision problem. It is shown by Blackwell (1962) that there is an optimal strategy profile which is stationary and given by some $\mathbf{a} \in \mathbf{A}$. There are many ways to compute the solution, including value function iteration, policy function iteration, and linear programming. In particular, the levels $\{x(s, \lambda)\}_{s \in S}$ are the solution to the linear program

$$\min_{y^{R(\cdot)}} \sum_{s \in S} y^R(s) \tag{22a}$$

$$\text{s.t. } y^R(s) \geq (1 - \delta)\lambda \cdot g(a) + \delta \sum_{s' \in S} \pi(s'|a)y^R(s') \text{ for all } s \in S, a \in \mathbf{A}(\mathbf{W})(s). \tag{22b}$$

A solution can be computed via the simplex algorithm, which will select exactly $|S|$ of the constraints to bind, so that their intersection uniquely pins down the value of y^R . At an optimum, there must be a binding constraint in each state, since otherwise we could decrease $y^R(s)$, and simultaneously decrease the right-hand side of every constraint. The choice of binding constraints is therefore a choice of exactly one action profile for each state, i.e., an $\mathbf{a} \in \mathbf{A}(\mathbf{W})$, which is an optimal policy. The simplex algorithm would identify such an optimal policy as a basic solution to the LP (22).

It has long been understood that the output of the simplex algorithm can be used to conduct “sensitivity analysis”: how much can we perturb the original problem without changing the optimal basis? In our case, we are concerned with sensitivity to λ , and for what range of directions would the optimal solution remain the same. As we rotate λ ,

we change the constants in the constraints. Eventually the optimal basis will change, and generically a single constraint will leave the basis and be replaced by a new one. This corresponds to changing the optimal policy in a single state. The next action to enter can be determined using well-known techniques, as in Dantzig and Thapa (2006). Mapping out the set of solutions for all λ is known as *parametric programming*, which is also a well established concept in mathematical optimization. In fact, this is precisely how our algorithm would behave if we restricted ourselves to using $\mathbf{r}(s) = R$, in which case the algorithm would converge in exactly one iteration (provided we start with any compact and convex valued correspondence that contains the feasible correspondence, e.g., large boxes whose bounds are given by the maximum and minimum flow payoffs across all states and actions.)

This is not our program, since we do have incentive constraints. It is in that sense closer to the problem of APS, which can also be formulated as an LP thusly: $x^{APS}(s, \lambda)$ is the solution to

$$\min_{y^{APS}(\cdot)} \sum_{s \in S} y^{APS}(s) \quad (23a)$$

$$\text{s.t. } y^{APS}(s) \geq \max\{\lambda \cdot v \mid v \in B(a, \mathbf{W})\} \text{ for all } s \in S, a \in \mathbf{A}(\mathbf{W})(s). \quad (23b)$$

This is not an LP in standard form, because of the inner maximization which is also an LP. But that problem can be replaced with its dual, in which case we have a single minimization program. Suppose that \mathbf{W} has finitely many faces with normals $\{\lambda^l\}_{l=0}^L$ and corresponding levels $\{z_l(s)\}_{l=0}^L$. Let $\mu_l(a, s)$ denote the multiplier on feasibility of the continuation value for action a in state s in the direction λ , and let $\alpha_i(a)$ denote the multiplier on the incentive constraint for player i . Applying the strong duality theorem of linear programming, we conclude that the best APS payoff is equal to the minimum of

$$y^{APS}(a) = \sum_{s' \in S} \sum_{l=1}^L \mu_l(a, s') ((1 - \delta)\lambda^l \cdot g(a) + \delta z_l(s')) - \sum_{i=1}^N \alpha_i(a) \underline{u}_i(a)$$

across all μ_l and α_l that are non-negative. Thus, we can expand (23) to

$$\begin{aligned} \min_{y^{APS}(\cdot), \mu_l(\cdot), \alpha_i(\cdot)} \sum_{s \in S} y^{APS}(s) \\ \text{s.t. } y^{APS}(s) \geq y^{APS}(a) \text{ for all } s \in S, a \in \mathbf{A}(\mathbf{W})(s) \end{aligned} \quad (24)$$

Again, this LP could be solved using the simplex algorithm, and one can map out the set of all basic solutions for all λ using sensitivity analysis and parametric programming.

Again, this is not our program. Ours is in fact a hybrid of the two:

$$\min_{y(\cdot), y^R(\cdot), y^{APS}(\cdot), \mu_l(\cdot) \geq 0, \alpha_i(\cdot) \geq 0} \sum_{s \in S} y(s) \quad (25a)$$

s.t. (22b) and (23b) and (24)

$$y(s) \geq \min\{y^R(a), y^{APS}(a)\} \quad \forall s \in S, a \in \mathbf{A}(\mathbf{W})(s). \quad (25b)$$

This is *not* an LP, because of the min operator in (25b). However, we can modify this program to make it into a larger LP, so that one could again use sensitivity analysis and parametric programming to map out solutions.

Specifically, we can add parameters $r(a) \in \{R, APS\}$ (which are not variables in the LP) and replace (25a) and (25b) with

$$\min_{y(\cdot), y^R(\cdot), y^{APS}(\cdot), \mu_l(\cdot) \geq 0, \alpha_i(\cdot) \geq 0} \sum_{s \in S} \left[y(s) + \sum_{a \in \mathbf{A}(\mathbf{W})(s)} (y^R(a) + y^{APS}(a)) \right] \quad (26)$$

s.t. (22b) and (23b) and (24)

$$y(s) \geq y^{r(a)}(a) \quad \forall s \in S, a \in \mathbf{A}(\mathbf{W})(s)$$

This is now an LP, and the $y(s)$ in the solution corresponds to the optimal levels under a particular conjecture as to which are the minimizing regimes, action profile by action profile. We could compute the level $x(s, \lambda)$ for a fixed direction by solving a sequence of such LPs, where at each step, we replace $r(a)$ with $\arg \min_r y^r(a)$, where $y^r(a)$ is taken from the previous solution. This will necessarily produce a decreasing sequence of solutions, whose limit is $x(s, \lambda)$.

Now, once we reach the optimal solution regimes $r(a)$, if we add one more constraint:

$$y^{r(a)}(a) \leq y^{r'}(a) \quad \text{for all } s \in S, a \in \mathbf{A}(\mathbf{W})(s), r \in \{R, APS\}, \quad (27)$$

the optimal solution will not change. Moreover, if we do sensitivity analysis on this expanded program, we will exactly find the range of directions λ under which the optimal actions and level-minimizing regimes do not change, action profile by action profile. So, in principle, one way to map out $x(s, \lambda)$ is to do sensitivity analysis on the expanded program of (26) and (27) to find adjacent directions where the solution to that program would change, and for those adjacent directions, resolve (26), re-optimizing the regimes $r(a)$ as needed.

Overall, this is quite a bit more work than what we have done in our more direct implementation. Effectively, the LP-based approach involves computing optimal regimes for every

action profile, even those which are not optimal, whereas our main procedure only computes minimal regimes for maximal action profiles. We have even implemented the LP based algorithm for two players using Gurobi, a high-performance commercial linear programming package. We found that this program took an order of magnitude longer to solve than the more direct implementation described in Section 4.2.

Nonetheless, this discussion may help to explain where the linear structure comes from, and why we end up using similar objects as those which arise in linear programming. It may also explain why we cannot simply use off-the-shelf techniques from linear programming in determining the function $x(s, \lambda)$.

E Online appendix: Redux for \tilde{B}^ϵ

This appendix extends the key results from Sections 3 and 4 to the operator \tilde{B}^ϵ .

E.1 Convergence results

Define the operator

$$T^\epsilon(\mathbf{y}, \lambda, \mathbf{a}, \mathbf{r}, \mathbf{W})(s) = -\epsilon + \begin{cases} (1 - \delta)\lambda \cdot g(\mathbf{a}(s)) + \delta \sum_{s' \in S} \mathbf{y}(s') \pi(s' | \mathbf{a}(s)) & \text{if } \mathbf{r}(s) = R; \\ x^{APS}(\mathbf{a}(s), \lambda, \mathbf{W}) & \text{if } \mathbf{r}(s) = APS. \end{cases}$$

Lemma 18 (Operator T^ϵ). *Fix λ , \mathbf{a} , \mathbf{r} , and \mathbf{W} . As a function of $\mathbf{y} : S \rightarrow \mathbb{R}$, T^ϵ is*

(L18.i) *increasing;*

(L18.ii) *a contraction with modulus δ and hence has a unique fixed point \mathbf{y}^* ;*

(L18.iii) *if $T^\epsilon(\mathbf{y}) \leq (\geq) \mathbf{y}$ then $\mathbf{y}^* \leq (\geq) T^\epsilon(\mathbf{y})$.*

Proof. The proof coincides verbatim with that of Lemma 1, changing T to T^ϵ . \square

Theorem 4 (The perturbed max-min-max algorithm). *For every $\epsilon > 0$, as a function of $\mathbf{W} : S \rightarrow 2^{\mathbb{R}^N}$, the operator \tilde{B}^ϵ has the following properties:*

(T4.i) *\tilde{B}^ϵ is increasing in \mathbf{W} , and if \mathbf{W} is compact, then $\tilde{B}^\epsilon(\mathbf{W})$ is compact;*

(T4.ii) *$\tilde{B}^\epsilon(\mathbf{W}) \subseteq B^\epsilon(\mathbf{W})$. Thus, if $\mathbf{W} \subseteq \tilde{B}^\epsilon(\mathbf{W})$, then \mathbf{W} is self-generating and $\mathbf{W} \subseteq \mathbf{V}^\epsilon$;*

(T4.iii) *$\mathbf{V}^\epsilon = \tilde{B}^\epsilon(\mathbf{V}^\epsilon)$;*

(T4.iv) *Fix a correspondence $\tilde{\mathbf{W}}^0$ that contains \mathbf{V}^ϵ . Define the sequence $\{\tilde{\mathbf{W}}^k\}_{k=0}^\infty$ by $\tilde{\mathbf{W}}^k = \tilde{B}^\epsilon(\tilde{\mathbf{W}}^{k-1})$. Then $\mathbf{V}^\epsilon = \bigcap_{k=0}^\infty \tilde{\mathbf{W}}^k$.*

Proof of Theorem 1.

(T4.i) For every λ and (\mathbf{a}, \mathbf{r}) , we can write

$$\eta(s, \mathbf{a}, \mathbf{r}) = x(s, \lambda, \mathbf{a}, \mathbf{r}, \mathbf{W}) - x^\epsilon(s, \lambda, \mathbf{a}, \mathbf{r}, \mathbf{W}).$$

Then η uniquely solves the system of equations

$$\eta(s, \mathbf{a}, \mathbf{r}) = \epsilon + \begin{cases} 0 & \text{if } \mathbf{r}(s) = APS; \\ \delta \sum_{s' \in S} \pi(s' | \mathbf{a}(s)) \eta(s', \mathbf{a}, \mathbf{r}) & \text{otherwise.} \end{cases}$$

Note for future reference that η is independent of both λ and \mathbf{W} . Thus, since x is monotonic in \mathbf{W} , so is x^ϵ . This implies monotonicity of \tilde{B}^ϵ . $\tilde{B}^\epsilon(\mathbf{W})$ is also closed, being the intersection of closed half-spaces, and bounded because \hat{x}^{APS} is bounded, so that x^ϵ is bounded as well.

(T4.ii) Clearly, $x^\epsilon(s, \lambda, \mathbf{W}) \leq x^{APS}(s, \lambda, \mathbf{W}) - \epsilon$, which implies that \tilde{B}^ϵ is always contained in B^ϵ . Thus, if $\mathbf{W} \subseteq \tilde{B}^\epsilon(\mathbf{W})$, then $\mathbf{W} \subseteq B^\epsilon(\mathbf{W})$ and hence, by APS, $B^\epsilon(\mathbf{W}) \subseteq \mathbf{V}^\epsilon$. Consequently, $\tilde{B}^\epsilon(\mathbf{W}) \subseteq \mathbf{V}$.

(T4.iii) From (T4.ii), it suffices to show that $\mathbf{V}^\epsilon \subseteq \tilde{B}^\epsilon(\mathbf{V}^\epsilon)$, i.e., for all λ , $x^\epsilon(s, \lambda, \mathbf{V}^\epsilon) \geq x^{APS}(s, \lambda, \mathbf{V}^\epsilon) - \epsilon$. To that end, fix λ , and for all s , let $\mathbf{a}(s)$ be an action that maximizes $x^{APS}(a, \lambda, \mathbf{V}^\epsilon)$ and let $\mathbf{w}(\cdot)$ be the associated continuation values as a function of the next-period state s' . We will show that $\min_{\mathbf{r}} x^\epsilon(s, \lambda, \mathbf{a}, \mathbf{r}, \mathbf{V}^\epsilon) \geq x^{APS}(s, \lambda, \mathbf{V}^\epsilon) - \epsilon$, so that $x^\epsilon(s, \lambda, \mathbf{V}^\epsilon) \geq x^{APS}(s, \lambda, \mathbf{V}^\epsilon) - \epsilon$, which implies the result. Since $\mathbf{V}^\epsilon = B^\epsilon(\mathbf{V}^\epsilon)$, $x^{APS}(s, \lambda, \mathbf{V}^\epsilon) - \epsilon \geq \lambda \cdot u$ for all $u \in \mathbf{V}^\epsilon(s')$ for all s' . Since $\mathbf{w}(s') \in \mathbf{V}^\epsilon(s')$ for all s' ,

$$\begin{aligned} x^{APS}(s, \lambda, \mathbf{V}^\epsilon) &= (1 - \delta)\lambda \cdot g(\mathbf{a}(s)) + \delta \sum_{s' \in \mathcal{S}} \pi(s' | \mathbf{a}(s)) \lambda \cdot \mathbf{w}(s') \\ &\leq (1 - \delta)\lambda \cdot g(\mathbf{a}(s)) + \delta \sum_{s' \in \mathcal{S}} \pi(s' | \mathbf{a}(s)) (x^{APS}(s', \lambda, \mathbf{V}^\epsilon) - \epsilon). \end{aligned}$$

Thus, if we let $\mathbf{y}(s) = x^{APS}(s, \lambda, \mathbf{V}^\epsilon) - \epsilon$ for all s , then for *any* regimes \mathbf{r} , $T^\epsilon(\mathbf{y}, \lambda, \mathbf{a}, \mathbf{r}, \mathbf{V}^\epsilon) \geq \mathbf{y}$ (with equality if $\mathbf{r}(s) = APS$ and weak inequality if $\mathbf{r}(s) = R$). By (L18.iii), we conclude that $\mathbf{y}(s) = x^{APS}(s, \lambda, \mathbf{V}^\epsilon) - \epsilon \leq x^\epsilon(s, \lambda, \mathbf{a}, \mathbf{r}, \mathbf{V}^\epsilon) = \mathbf{y}^*(s)$, as required.

(T4.iv) (T4.ii) implies that $\tilde{\mathbf{W}}^k \subseteq \mathbf{W}^k$, where the latter is the k th element of the APS sequence for B^ϵ starting from $\tilde{\mathbf{W}}^0$. Also, the fact that $\tilde{\mathbf{W}}^0$ contains \mathbf{V} , (T4.i), and (T4.iii) imply that $\mathbf{V}^\epsilon \subseteq \tilde{\mathbf{W}}^k$. Thus, $\mathbf{V}^\epsilon \subseteq \cap_k \tilde{\mathbf{W}}^k \subseteq \cap_k \mathbf{W}^k = \mathbf{V}^\epsilon$.

□

E.2 State independence of the optimal policy

We now restate the results for minimal regimes. Let us define

$$x^{R, \epsilon}(a, \lambda, \mathbf{a}, \mathbf{r}) = (1 - \delta)\lambda \cdot g(a) + \delta \sum_{s' \in \mathcal{S}} \pi(s' | a) x^\epsilon(s', \lambda, \mathbf{a}, \mathbf{r}).$$

For given λ , \mathbf{W} , and $\mathbf{a} \in \mathbf{A}(\mathbf{W})$, we say that the regimes \mathbf{r} are *minimal* if and only if for all $s \in S$,

$$x^\epsilon(s, \lambda, \mathbf{a}, \mathbf{r}) = \min_{\mathbf{r}' \in \mathbf{R}} x^\epsilon(s, \lambda, \mathbf{a}, \mathbf{r}').$$

Lemma 19 (Minimal regimes). *For all $\mathbf{a} \in \mathbf{A}(\mathbf{W})$, λ , and $\epsilon > 0$,*

(L19.i) *there exists minimal regimes;*

(L19.ii) *\mathbf{r} is minimal if and only if for all $s \in S$,*

$$x^\epsilon(s, \lambda, \mathbf{a}, \mathbf{r}) = \{x^{APS}(\mathbf{a}(s), \lambda), x^{R,\epsilon}(\mathbf{a}(s), \lambda, \mathbf{a}, \mathbf{r})\} - \epsilon; \quad (28)$$

(L19.iii) *if (28) is violated for some s , then \mathbf{r} is not minimal. Moreover, for all $s' \in S$, $x^\epsilon(s', \lambda, \mathbf{a}, \mathbf{r} \setminus s) \leq x^\epsilon(s', \lambda, \mathbf{a}, \mathbf{r})$, with strict inequality in state s .*

Proof of Lemma 18. The proof follows verbatim that of Lemma 1, replacing T with T^ϵ . \square

Proof of Lemma 19. The proof follows verbatim that of Lemma 2, replacing T with T^ϵ , x with x^ϵ , references to equation (6) with (28), and references to Lemma 1 with references to Lemma 18. \square

We next extend the results for maximal actions. Define $x^\epsilon(s, \lambda, \mathbf{a})$ to be $x^\epsilon(s, \lambda, \mathbf{a}, \mathbf{r})$ for some minimal regimes \mathbf{r} . Also, define

$$T^{min,\epsilon}(\mathbf{y}, \lambda, \mathbf{a})(s) = \min \left\{ x^{APS}(\mathbf{a}(s), \lambda), (1 - \delta)\lambda \cdot g(\mathbf{a}(s)) + \delta \sum_{s' \in S} \mathbf{y}(s') \pi(s' | \mathbf{a}(s)) \right\} - \epsilon.$$

Lemma 20 (Operator $T^{min,\epsilon}$). *Fix $\epsilon > 0$, λ , and $\mathbf{a} \in \mathbf{A}(\mathbf{W})$. As a function of $\mathbf{y} : S \rightarrow \mathbb{R}$, $T^{min,\epsilon}$ is*

(L20.i) *increasing;*

(L20.ii) *a contraction with modulus δ , and hence has a unique fixed point \mathbf{y}^* ;*

(L20.iii) *if $T^{min,\epsilon}(\mathbf{y}) \leq (\geq) \mathbf{y}$ then $\mathbf{y}^* \leq (\geq) T^{min,\epsilon}(\mathbf{y})$;*

Proof of Lemma 20. The proof follows verbatim that of Lemma 3, replacing T^{min} with $T^{min,\epsilon}$. \square

We further define

$$x^{R,\epsilon}(a, \lambda, \mathbf{a}) = (1 - \delta)\lambda \cdot g(a) + \delta \sum_{s' \in S} \pi(s' | a) x^\epsilon(s', \lambda, \mathbf{a}, \mathbf{r}),$$

where \mathbf{r} is minimal for \mathbf{a} and λ .

Lemma 21 (Maximal actions). *Suppose that $\mathbf{A}(\mathbf{W})$ is non-empty valued. For all $\epsilon > 0$ and λ ,*

(L21.i) *there exist maximal actions;*

(L21.ii) *$\mathbf{a} \in \mathbf{A}(\mathbf{W})$ is maximal if and only if for all $s \in S$ and $a \in \mathbf{A}(\mathbf{W})(s)$,*

$$x^\epsilon(s, \lambda, \mathbf{a}) \geq \min \{x^{APS}(a, \lambda), x^{R, \epsilon}(a, \lambda, \mathbf{a})\} - \epsilon, \quad (29)$$

with equality when $a = \mathbf{a}(s)$;

(L21.iii) *if (29) is violated for some $s \in S$ and $a \in \mathbf{A}(\mathbf{W})(s)$, then \mathbf{a} is not maximal. Indeed, for all $s' \in S$, $x^\epsilon(s', \lambda, \mathbf{a} \setminus (s, a)) \geq x^\epsilon(s', \lambda, \mathbf{a})$, with strict inequality in state s .*

Proof of Lemma 21. Once again, this follows verbatim the proof of Lemma 4, replacing x with x^ϵ , T^{\min} with $T^{\min, \epsilon}$, references to (7) with references to (29), and references to Lemma 3 with references to Lemma 20. \square

E.3 Sufficiency of binding payoffs

Lemma 22. *For any direction λ , if B^ϵ sub-generates at \mathbf{W} in the direction λ , then for any $\mathbf{a} \in \mathbf{A}(\mathbf{W})$, if $\gamma(\mathbf{a}(s), \lambda, \mathbf{W}) > 0$, then*

$$x^{APS}(\mathbf{a}(s), \lambda, \mathbf{W}) - \epsilon \geq x^{R, \epsilon}(\mathbf{a}(s), \lambda, \mathbf{a}, \mathbf{W}) - \epsilon = x^\epsilon(s, \lambda, \mathbf{a}, \mathbf{W})$$

Moreover, there exist minimal regimes such that $\mathbf{r}(s) = R$ for s with $\gamma(\mathbf{a}(s), \lambda, \mathbf{W}) > 0$.

Proof of Lemma 22. Suppose that $\gamma(\mathbf{a}(s), \lambda, \mathbf{W}) > 0$. Then the best continuation values from \mathbf{W} in the direction λ , denoted \mathbf{w} , must be incentive compatible for $\mathbf{a}(s)$, and

$$x^{APS}(\mathbf{a}(s), \lambda, \mathbf{W}) = (1 - \delta) \lambda \cdot g(\mathbf{a}(s)) + \delta \sum_{s' \in S} \pi(s' | \mathbf{a}(s)) \lambda \cdot \mathbf{w}(s').$$

Sub-generation and the definition of x^ϵ imply that $\lambda \cdot \mathbf{w}(s') \geq x^{APS}(\mathbf{a}(s'), \lambda, \mathbf{W}) - \epsilon \geq x^\epsilon(s', \lambda, \mathbf{W})$. Hence,

$$\begin{aligned} x^{APS}(\mathbf{a}(s), \lambda, \mathbf{W}) - \epsilon &\geq (1 - \delta) \lambda \cdot g(\mathbf{a}(s)) + \delta \sum_{s' \in S} \pi(s' | \mathbf{a}(s)) x^\epsilon(s', \lambda, \mathbf{W}) \\ &\geq x^{R, \epsilon}(\mathbf{a}(s), \lambda, \mathbf{a}, \mathbf{W}) \end{aligned}$$

as desired.

Finally, suppose \mathbf{r} is minimal and $\gamma(\mathbf{a}(s), \lambda, \mathbf{W}) > 0$. If $x^{APS}(\mathbf{a}(s), \lambda, \mathbf{W}) > x^{R,\epsilon}(\mathbf{a}(s), \lambda, \mathbf{a}, \mathbf{W})$, then $\mathbf{r}(s)$ is necessarily R . Otherwise, (9) implies that $x^{APS}(\mathbf{a}(s), \lambda, \mathbf{W}) = x^{R,\epsilon}(\mathbf{a}(s), \lambda, \mathbf{a}, \mathbf{W})$. Thus, if we set $\mathbf{r}'(s) = R$ for all states with $\gamma(\mathbf{a}(s), \lambda, \mathbf{W}) > 0$ and $\mathbf{r}'(s') = \mathbf{r}(s')$ otherwise, then $x^\epsilon(\cdot, \lambda, \mathbf{a}, \mathbf{r}, \mathbf{W})$ is clearly a fixed point of $T^\epsilon(\cdot, \lambda, \mathbf{a}, \mathbf{r}', \mathbf{W})$, so that \mathbf{r}' also satisfies (6) and is minimal. \square

Lemma 23. *If \tilde{B}^ϵ sub-generates at \mathbf{W} , then B^ϵ sub-generates at $\tilde{B}^\epsilon(\mathbf{W})$.*

Proof of Lemma 23. Towards a contradiction, suppose that some action profile $a \in \mathbf{A}(\mathbf{W})(s)$, with continuation values $\mathbf{w} \in \tilde{B}^\epsilon(\mathbf{W})$, generates a payoff outside the convex set $\tilde{B}^\epsilon(\mathbf{W})$. Then for some direction λ , $x^{APS}(a, \lambda, \tilde{B}^\epsilon(\mathbf{W})) - \epsilon > x^\epsilon(s, \lambda, \mathbf{W})$, so

$$\begin{aligned} x^\epsilon(s, \lambda, \mathbf{W}) + \epsilon &< x^{APS}(a, \lambda, \tilde{B}^\epsilon(\mathbf{W})) = \lambda \cdot \left((1 - \delta)g(a) + \delta \sum_{s' \in S} \pi(s'|a)\mathbf{w}(s') \right) \\ &\leq (1 - \delta)\lambda \cdot g(a) + \delta \sum_{s' \in S} \pi(s'|a)x^\epsilon(s', \lambda, \mathbf{W}), \end{aligned}$$

where the last inequality holds because $\lambda \cdot \mathbf{w}(s') \leq x^\epsilon(s', \lambda, \mathbf{W})$, since $\mathbf{w}(s') \in \tilde{B}^\epsilon(\mathbf{W})(s')$. The right-hand side of this inequality equals $x^{R,\epsilon}(a, \lambda, \mathbf{a}, \mathbf{W})$ for any $a \in \mathbf{A}(\mathbf{W})(s)$ that is maximal in the direction λ (given \mathbf{W}). Since $\tilde{B}^\epsilon(\mathbf{W}) \subseteq \mathbf{W}$, we know that $x^{APS}(s, \lambda, \mathbf{W})$ is greater than $x^\epsilon(s, \lambda, \mathbf{W})$ as well. That is, $x^\epsilon(s, \lambda, \mathbf{a}, \mathbf{W}) < \min\{x^{APS}(a, \lambda, \mathbf{W}), x^{R,\epsilon}(a, \lambda, \mathbf{a}, \mathbf{W})\} - \epsilon$, contradicting (L21.ii). \square

Proposition 5 (Sufficiency of binding payoffs). *As long as B^ϵ sub-generates at $\tilde{\mathbf{W}}^0$, then for any $k \geq 0$, B^ϵ sub-generates at $\tilde{\mathbf{W}}^k$. As a result, for any λ and $\mathbf{a} \in \mathbf{A}(\mathbf{W})$, if $\gamma(a, \lambda, \tilde{\mathbf{W}}^k) > 0$ is strictly positive, then $\mathbf{r}^*(s) = R$.*

Proof of Proposition 5. Follows verbatim the proof of Proposition 1, replacing B and \tilde{B} with B^ϵ and \tilde{B}^ϵ , respectively. \square

Finally, we extend the characterizations of optimal policies and optimal pairs.

Lemma 24. *If B^ϵ sub-generates at \mathbf{W} in the direction λ , the actions $\mathbf{a} \in \mathbf{A}(\mathbf{W})$ are maximal if and only if for all (s, a) ,*

$$x^\epsilon(s, \lambda, \mathbf{a}) \geq \begin{cases} \min\{\hat{x}^{APS}(a, \lambda), x^{R,\epsilon}(a, \lambda, \mathbf{a})\} - \epsilon & \text{if } \gamma(a, \lambda, \mathbf{W}) = 0; \\ x^{R,\epsilon}(a, \lambda, \mathbf{a}) - \epsilon & \text{if } \gamma(a, \lambda, \mathbf{W}) > 0. \end{cases}$$

Note that Lemma 10 implies, via the same argument in Corollary 1, that \mathbf{V}^ϵ has at most \bar{L} extreme points. Moreover, we can adapt the algorithm in Section 4 to compute $\tilde{B}^\epsilon(\mathbf{W})$. It is still the case that direction where robustly optimal payoffs \mathbf{u} cease to be optimal corresponds to a substitution (s, a, p) . When $p = R$, the change must occur at a direction λ' such that

$$\lambda' \cdot (u^R(a, \mathbf{u}) - \mathbf{u}(s)) = \sum_{s' \in S} \pi(s'|a) \eta(s', \mathbf{a}, \mathbf{r}) + \epsilon - \eta(s, \mathbf{a}, \mathbf{r}), \quad (30)$$

so that the change in level is exactly offset by a change in penalty, and otherwise

$$\lambda' \cdot (p - \mathbf{u}(s)) = \epsilon - \eta(s, \mathbf{a}, \mathbf{r}), \quad (31)$$

where \mathbf{r} are the regimes associated with the incumbent optimal pair that induces \mathbf{u} . There are at most $2\bar{L}\bar{M}$ directions that satisfy (30) or (31). Such a direction is again called legitimate if (\mathbf{a}, \mathbf{p}) is optimal in that direction. We can therefore compute \tilde{B}^ϵ by finding the optimal pair in one direction, then iteratively computing the legitimate substitution direction with the smallest clockwise rotation, and then lexicographically optimizing the pair in the new direction. This produces a sequence of directions and optimal payoffs $\{(\lambda^k, \mathbf{u}^k)\}_{k=0}^K$.

Note that a subtle issue is that the new optimal level $x^\epsilon(s, \lambda, \mathbf{W})$ is no longer piecewise linear, but is piecewise affine of the form $\lambda \cdot u - \eta$, where $\eta > 0$. Because of the sign of the constant, it turns out that directions at which the optimal pair does not change are still redundant. In particular, if we have a clockwise sequence of directions λ , λ' , and λ'' at which u are the optimal payoffs in state s and η is the optimal penalty, then

$$\begin{aligned} \lambda' \cdot v &= \frac{\alpha\lambda + (1-\alpha)\lambda''}{\|\alpha\lambda + (1-\alpha)\lambda''\|} \cdot v \\ &\leq \frac{1}{\|\alpha\lambda + (1-\alpha)\lambda''\|} [\alpha(\lambda \cdot u - \eta) + (1-\alpha)(\lambda'' \cdot u - \eta)] \\ &= \lambda' \cdot u - \frac{1}{\|\alpha\lambda + (1-\alpha)\lambda''\|} \eta \\ &\leq \lambda' \cdot u - \eta, \end{aligned}$$

since $\|\alpha\lambda + (1-\alpha)\lambda''\| \leq 1$. As a result, we can simply intersect the half-spaces at legitimate test directions to compute \tilde{B}^ϵ . We therefore have:

Theorem 5. *Suppose that $N = 2$, $\mathbf{A}(\mathbf{W})$ is non-empty valued, and B^ϵ sub-generates at \mathbf{W} . Then the previously described procedure terminates in at most $2\bar{L}\bar{M}$ substitutions and runtime $O(\bar{L}\bar{M}^2)$. If there are no legitimate test directions at \mathbf{u}^0 , then $\tilde{B}(\mathbf{W})(s) = \{\mathbf{u}^0(s)\}$*

for all s . Otherwise,

$$\tilde{B}^\epsilon(\mathbf{W})(s) = \{v | \lambda^k \cdot v \leq \lambda^k \cdot \mathbf{u}^k(s) \ \forall k = 1, \dots, K\}. \quad (32)$$