

# Explicit Renegotiation in Repeated Games

Mikhail Safronov  
Northwestern University

Bruno Strulovici\*  
Northwestern University

November 14, 2014

## Abstract

Cooperative concepts of renegotiation in repeated games have typically assumed that Pareto-ranked equilibria could not coexist within the same renegotiation-proof set. With explicit renegotiation, however, a proposal to move to a Pareto-superior equilibrium can be deterred by a different continuation equilibrium which harms the proposer and rewards the refuser. This paper introduces a simple protocol of renegotiation for repeated games and defines the stability of social norms and renegotiation-proof outcomes in terms of a simple equilibrium refinement. We provide distinct necessary and sufficient conditions for renegotiation-proofness, which converge to each other as renegotiation frictions become negligible. Renegotiation-proof outcomes always exist and can be all included within a single, most permissive social norm that is straightforward to characterize graphically.

## 1 Introduction

The punishment equilibria used to sustain cooperation in repeated games are often Pareto inefficient. This puts into question their credibility and the implementability of cooperative outcomes based on such punishments when players are able to freely renegotiate the continuation of the game. In fact, incorporating renegotiation satisfactorily in repeated games has been a longstanding challenge.

To address this question, economists have introduced various concepts of renegotiation-proofness based on the following idea: roughly speaking, an equilibrium is *not* renegotiation-proof if it entails a continuation play that is Pareto dominated by some “credible” equilibrium (Pearce (1987), Bernheim and Ray (1989), Farrell and Maskin (1989), Abreu and Pearce (1991), and Asheim (1991)).<sup>1</sup>

---

\*We are grateful for comments from Ilya Segal, Larry Samuelson, and Joel Watson. Strulovici acknowledges financial support from the NSF (Grant No.1151410) and the Alfred P. Sloan Foundation.

<sup>1</sup>The first discussion along these lines is due to Farrell (1983), which is subsumed by Farrell and Maskin (1989). Other approaches to renegotiation include DeMarzo (1988), Benoît and Krishna (1993), and Bergin and MacLeod (1993). All these papers follow axiomatic approaches.

These concepts mainly differ regarding what “credible” means and yield contrasted results: while Pearce (1987) argued,<sup>2</sup> as in the first paragraph, that maximal cooperation may not be sustained due to the lack of a credible and severe enough punishment, Farrell and Maskin (1989) found that most renegotiation-proof outcomes, as players become arbitrarily patient, *had* to be on the Pareto frontier of the feasible set.<sup>3</sup>

Owing to their cooperative (i.e., non-strategic) nature, these concepts have left unexplored an aspect of renegotiation which becomes crucial, perhaps even obvious, when one considers *explicit* renegotiation: what happens when a player rejects another player’s proposal? Suppose that during the punishment phase of a two-player repeated game, the continuation payoffs are  $(X_1, X_2)$  and player 1 proposes a Pareto-improving equilibrium with payoffs  $(Y_1, Y_2)$ . Clearly, such a Pareto-improvement need not be accepted if, by rejecting 1’s proposal, player 2 gets rewarded by a higher continuation payoff  $Z_2 > Y_2$ . Moreover, if 1’s continuation payoff  $Z_1$  after 2 has rejected her offer is less than  $X_1$ , then it is suboptimal for 1 to propose the Pareto improvement in the first place. With explicit renegotiation, a bad equilibrium may perfectly withstand renegotiation as long as any (credible) deviating proposal can be deterred in this fashion. Punishing a player who deviates (here, in proposals) and rewarding the player who counters the deviation is standard thinking in the analysis of dynamic games. It is also realistic: for example, if a player tries to bribe another player to obtain some favor (a Pareto improving scheme for those players!), the player who rejects and exposes the bribe may be rewarded for doing so, while the corrupting player may be punished.<sup>4</sup>

This paper considers explicit renegotiation in repeated games by appending a simple stage at the end of each period: after actions and payoffs have been chosen and observed for period  $t$ , one of the players may be selected, with a fixed probability, to propose a continuation *plan*. A plan is more easily defined recursively,<sup>5</sup> as prescribing players’ actions and proposals in period  $t + 1$ , as

---

<sup>2</sup>See also Abreu and Pearce (1991) and Abreu, Pearce and Stacchetti (1993).

<sup>3</sup>Farrell and Maskin, like Bernheim and Ray, introduce weak and strong concepts of renegotiation-proofness. The strong notion is arguably the more satisfactory one as it allows external comparisons (for example, the repetition of any static Nash equilibrium forms a weakly renegotiation-proof of the repeated game, but that equilibrium can be challenged by other equilibria according to the strong concept). The strong concept is well-behaved (existence, uniqueness) when players are arbitrarily patient, although the set of strongly-renegotiation proof equilibria may be very small due to the lack of punishments outside of a line that goes through the Pareto frontier.

<sup>4</sup>It may be tempting for Player 1 to approach Player 2 nonetheless and beg him to ignore all equilibrium conventions and simply implement the Pareto-improving equilibrium. Note however that such a proposal is precisely the kind of deviation considered *within* our model, and that a player on the receiving end of the proposal can be rewarded by rejecting it, as long as a suitably rewarding continuation is available for that player within the social norm.

<sup>5</sup>The definition can also be stated in extensive form. By time homogeneity of the setting, the set of plans is also time invariant. As will be clear, however, the set of continuations in a given equilibrium, even allowing for renegotiation, may be history dependent.

well as the actual plan implemented in period  $t + 2$  as a function of the actions, proposals, and acceptance/rejection decisions made in period  $t + 1$ .<sup>6</sup>

A set  $\mathcal{N}$  of equilibria is called a *norm* if, roughly speaking, all continuations of these equilibria (except possibly deviating proposals that are accepted) belong to  $\mathcal{N}$ . This definition aims to capture the idea of a social norm, known to all players, describing all equilibria which players see as feasible and acceptable.

To assess the stability of a norm, we allow players to propose innovations which are equilibria of the game with renegotiation. A norm  $\mathcal{N}$  is *stable* if it can withstand innovations outside of the norm, by rewarding the player who rejects the proposal and punishing the player who made it.<sup>7</sup>

Our notion of stability is strategic:<sup>8</sup> it imposes an equilibrium refinement requiring that if a player makes a proposal and which the other player accepts, then that proposal is played. This definition rules out, in particular, cheap talk equilibria in which all proposals are ignored and any equilibrium of the underlying repeated game can be implemented. In fact, we find that as players become arbitrarily patient, the folk theorem need not hold any more when renegotiation is allowed.

We characterize the set of all *renegotiation-proof* equilibrium payoffs, which are those payoffs that belong to the payoff set of some stable norm, when players are patient and renegotiation frictions (modeled as the probability that no one gets to propose within any fixed time window) become negligible. This set is well behaved: it is always non empty and is very simple to characterize. In fact, its shape depends on only three points of the feasible payoff space: the minmax payoff vector  $V$ , and the payoff vectors  $P_1$  and  $P_2$  that delimit the Pareto frontier. The set of renegotiation-proof payoffs is the intersection of two positive orthants with the feasible set: the orthant with vertex  $V$  (as in the Folk Theorem) and the orthant whose boundaries go through  $P_1$  and  $P_2$ . Our characterization is based on distinct necessary and sufficient conditions, holding at all friction levels of renegotiation, which converge to each other as renegotiation frictions become negligible.

---

<sup>6</sup>Because these proposals can a priori be treated as cheap talk, any equilibrium of the underlying repeated game without renegotiation has an equivalent equilibrium in the game with renegotiation. In particular, there always exist equilibria in this game.

<sup>7</sup>Our results do not hinge on allowing players to propose arbitrary innovations. In fact, our characterization is rigorously identical if players are instead restricted to propose “credible” innovations, in the following sense. Given a norm  $\mathcal{N}$ , an equilibrium is  *$\mathcal{N}$ -credible* if any deviation from this equilibrium (in action or proposal) triggers reversal to the norm. For example, if a punishment equilibrium belongs to the norm  $\mathcal{N}$ , meaning that both players see such a punishment as feasible (despite the possibility of renegotiation), then players may reasonably entertain the idea of a cooperative equilibrium outside of the norm, with the understanding that if that equilibrium doesn’t work out (i.e., some player makes a deviation), they will revert to the norm and implement the punishment equilibrium.

<sup>8</sup>We also provide a simple set-theoretic notion which is payoff-equivalent to the strategic one.

Moreover, in our setting all renegotiation-proof points can be implemented within the same stable norm. Therefore, there is no issue of competition or indeterminacy between multiple norms, at least not for the most permissive norm.<sup>9</sup>

Our construction also implies *path dependence* for the set of proposals considered acceptable within the norm. For example, the cooperative proposal  $(Y_1, Y_2)$  above may have been acceptable at the beginning of the game, but not after a deviation. The relevance of this path dependence has been emphasized earlier (see, Abreu and Pearce (1991) and Asheim (1991)), and arises naturally in our explicit model of renegotiation. We also show that Farrell and Maskin's and Bernheim and Ray's common notion of *weak* renegotiation proofness (or internal consistency), often considered a minimum requirement for any renegotiation-proof concepts, *does* rule out sustainable equilibria in our game. This is intuitive: in these reduced-form concepts, renegotiation-proofness requires that no equilibrium Pareto dominates another one in the same norm. As we argued earlier, however, the pursuit of a Pareto dominated equilibrium may withstand renegotiation as long as the social norm can reward a player for rejecting any proposal outside of the norm.

Although our main objective is to characterize the set of sustainable equilibria and renegotiation-proof norms as renegotiation frictions become arbitrarily small, we also provide distinct necessary and sufficient conditions for arbitrary levels of frictions. We then show that these conditions converge to each other as frictions become negligible.

As mentioned above, our stable norms differ not only conceptually but also physically from previous definitions. Compared to weakly renegotiation-proof or internally consistent sets, they are both more stringent by allowing proposals that lie outside of the norms, but also more permissive by allowing Pareto-ranked equilibria to coexist within a given norm and the path-dependence of acceptable proposals and equilibria. Our concept is weaker than strongly renegotiation-proof sets when agents get arbitrarily patient, precisely because we allow punishments that are Pareto dominated within the norm. Conceptually, however, our work is related to both notions of internal and external consistency: stable norms are internally consistent in the sense that one has to be able to punish deviation by a continuation that lies within the norm. They are also related to external consistency because the norm is challenged by any equilibrium, and the set of equilibria is typically much larger than the norm itself.

In accordance to the accepted standard in the modern analysis of repeated games, we allow players to use a public randomization device and private mix strategies. This feature distinguishes our analysis from some of the earlier work on renegotiation. For example, Farrell and Maskin (1989)

---

<sup>9</sup>There are many more restrictive norms.

assume that players can observe each other's mixing strategies, rather than merely observing the action outcomes of the randomization.<sup>10</sup> Bernheim and Ray (1989) rule out mixing altogether, focusing the analysis on pure-strategy equilibria.<sup>11</sup>

In a recent paper, Miller and Watson (2013) study equilibrium selection in a repeated game with an explicit bargaining protocol and transfers. Their goals and analysis are quite different from this paper's. In particular, that paper allows transfers and proposes an axiomatic restriction for disagreements outcomes, which radically changes the analysis of punishments.

## 2 Setting

Consider a repeated game between two players indexed by  $i \in \{1, 2\}$ . Player  $i$ 's stage-game action,  $a_i$ , lies in a finite set  $\mathcal{A}_i$ . The vector  $\mathbf{a} = (a_1, a_2)$  of actions determines the period's payoffs  $\mathbf{u}(\mathbf{a}) = (u_1(\mathbf{a}), u_2(\mathbf{a}))$ . A distribution  $\alpha_i$  over  $\mathcal{A}_i$  will be called a *mixed action* for  $i$ , and  $\alpha = (\alpha_1, \alpha_2)$  will denote a vector of mixed actions for both players. Players have a common discount factor  $\delta \in (0, 1)$ , and we will often find it convenient to work with the current-period weight  $\varepsilon = 1 - \delta$ .

Each period is composed of the following stages:

- 1) Players observe the realization  $P$  of a public randomization device taking values in  $[0, 1]$ ;
- 2) They simultaneously choose mixed actions  $\alpha_i \in \Delta(\mathcal{A}_i)$ ,  $i \in \{1, 2\}$ . Mixing probabilities are not observable. Conditional on the realization  $P$  of the public randomization device, players choose their mixed actions independently from each other;
- 3) The vector  $\mathbf{a}$  of actions is observed and the period's payoffs are realized;
- 4) With probability  $p < 1$ , one of the players is chosen to propose a new *plan* describing the continuation of the game. Each player has the same probability of  $p/2$  being chosen. The chosen player may, however, conceal his proposal opportunity and remain silent instead, or mix between proposing or not. The object of the proposal is an infinite-horizon plan  $m$  from the set  $\mathcal{M}$  of all possible plans, and will be described shortly;

---

<sup>10</sup>That paper contains a claim that observing mixed strategies is without loss. However, it is possible to find counter-examples showing that this claim is erroneous. Intuitively, when players observe mixing, there is without loss only single continuation payoff vector, conditional on players' mixing strategies. When mixtures are unobservable, however, there must be a continuation vector for every possible outcome of the mixture, and all of these vectors must belong to the renegotiation-proof set. This is problematic because some of these continuations may have Pareto-ranked payoffs, violating weak renegotiation-proofness.

<sup>11</sup>Other papers have made different restrictions, such as focusing on symmetric equilibria or a finite horizon.

5) If  $i$  makes a proposal, player  $-i$  decides whether to accept it, possibly mixing between acceptance and rejection. The resulting decision  $D_{-i} \in \{0, 1\}$  is equal to 1 (0) if  $-i$  accepts (rejects) the proposal;

The *public* history for the stage consists of the realisation  $P$  of the randomisation device, the action vector  $\mathbf{a}$ , a proposal  $m_i$  (or absence thereof) and, if applicable, an acceptance decision  $D_{-i}$ . In addition, each player privately observes the mixing probability used for each of her decisions.

A *plan* at period  $t$  describes players' strategy for the infinite repetition of the stage-game described above, *from period  $t + 1$  onwards*. Those decisions (actions, proposals, and acceptance mixtures) are history-dependent. Because the setting is time invariant, the set  $\mathcal{M}$  of plans can be more conveniently defined recursively.

Specifically, a plan  $m \in \mathcal{M}$  at any period  $t$  is characterized as follows:

a) For each realization  $P$  of the public randomization device, a pair  $\alpha = \alpha[m](P)$  of mixed actions that players should play in period  $t + 1$ ;

b1) For each player  $i$ , a distribution  $\bar{\mu}_i = \bar{\mu}_i[m](P, \mathbf{a}) \in \Delta(\mathcal{M} \cup \emptyset)$  over proposals, where the outcome  $\emptyset$  means that  $i$  abstains from proposing (unbeknownst to player  $-i$ ). We assume for simplicity that distributions have finite support over plans.<sup>12</sup> The proposer's choice of a proposal distribution is conditional on the realization  $P$  of the public randomization device and on the pair  $\mathbf{a}$  of observed actions. Because  $p < 1$ , not observing any proposal from either player is always consistent with "on-path" behavior. The realized proposal is denoted  $\mu_i$ ;

b2) A probability  $q_{-i} = q_{-i}[m](P, \mathbf{a}, \mu_i)$  that  $-i$  accepts  $i$ 's proposal (whenever  $\mu_i \neq \emptyset$ ), conditional on  $P$ ,  $\mathbf{a}$ , and  $\mu_i$ ;

b3) If no one made a proposal, the acceptance stage is skipped. To economize on notation, we assume that some player  $i$  is, even in that case, conventionally selected (randomly or deterministically) as the proposer and let  $\mu_i = \emptyset$  and  $D_{-i} = 0$ . (So,  $-i$ 's conventional response is to systematically "reject"  $i$ 's non proposal.)

c) A continuation plan  $m_{+1} = m_{+1}[m](P, \mathbf{a}, i, \mu_i, D_{-i}) \in \mathcal{M}$  for period  $t + 2$  onwards, as a function of  $P$ ,  $\mathbf{a}$ ,  $i$ ,  $\mu_i$ ,  $D_{-i}$ , where  $i$  indicates the identity of the last proposer. Obviously, this plan must be independent of  $i$  whenever  $\mu_i = \emptyset$ , so that the convention chosen for the proposer in the absence of any actual proposal be indeed irrelevant. This restriction is applied throughout.

---

<sup>12</sup>This assumption sidesteps measurability issues over plans. It could easily be relaxed to, say, any subset  $\mathcal{P}$  of plans which is in bijection with a Borel subset  $\mathcal{B}$  of  $\mathbb{R}^n$ , say, in which case  $\bar{\mu}_i$  would be the pushback measure over  $\mathcal{P}$  which corresponds to any distribution over  $\mathcal{B}$ .

Notation:

- 1) Actual plans are denoted by  $m$  and proposals by  $\mu$ .
- 2) The subscript +1 in the plan notation  $m_{+1}[m](P, \mathbf{a}, i, \mu_i, D_{-i})$  indicates that this plan concerns the next period. This plan is, like the initial plan  $m$  itself, an element from the set  $\mathcal{M}$ .

### 3 Concepts

The previous section has introduced an infinite-horizon game which we call “negotiated game”. Any strategy profile of that game can be identified with a plan defined in that section. Indeed, a plan defines – explicitly or recursively – an arbitrary history-dependent mixture of actions at each node of the game.<sup>13</sup> Accordingly, the subgame perfect equilibria (SPEs) of the renegotiated repeated game can be identified as a subset  $\mathcal{S}$  of  $\mathcal{M}$ . Unless stated otherwise, in this paper “SPE” refers to an equilibrium of our game with renegotiation (not to be confused with the subgame perfect equilibria of the underlying repeated game without renegotiation).

**DEFINITION 1** *A subset  $\mathcal{N}$  of  $\mathcal{S}$  is a **norm** if for any  $m \in \mathcal{N}$  such that  $\mu_i \in \bar{\mu}_i[m](P, \mathbf{a})$  or  $D_{-i} = 0$ ,  $m_{+1}[m](P, \mathbf{a}, i, \mu_i, D_{-i}) \in \mathcal{N}$ ;*

We interpret  $\mathcal{N}$  as a social norm: it describes the set of all continuation plays which players consider possible under “business as usual”. According to this norm, players may be punished if they deviate from the equilibrium path, but they are always punished *within* the norm, regardless of the history. However, players may in principle agree to play something outside of the norm. This happens if a player makes a deviation in proposal (hence creating an “innovation”), which the other player accepts.

Finally, we define the key concept of the paper.

**DEFINITION 2** *A norm  $\mathcal{N}$  is **stable** if for any SPE of  $\mathcal{N}$ , whenever  $i$  proposes an equilibrium  $\mu \in \mathcal{S}$  and  $-i$  accepts it,  $\mu$  is implemented.*

Stability thus amounts to a simple equilibrium refinement which rules out pure cheap talk, giving some bite to proposals. This requires that any SPE of the norm  $\mathcal{N}$  be able to withstand arbitrary proposals. Recall on-path continuations must all belong to the norm, stability implicitly requires

---

<sup>13</sup>In the actual game, the absence of a proposal triggers the next period. Therefore, a plan’s independence from the conventionally chosen proposer in case of a silence is not restriction on the set of strategy profiles being considered.

that any proposals that is *not* in the norm be rejected. We allow considerable leeway in proposals. As it turns out, both of our necessary and sufficient conditions, stated in the next section, are rigorously identical if one restricts proposals to a much smaller subset of “credible” proposals, which are roughly speaking equilibria such that any deviation triggers a reversal to the norm. Since such a restriction is not needed for the results, however, we start with the simplest concept and postpone credible innovations to Section 5.

Before defining renegotiation-proofness, we need to introduce notation to distinguish players’ payoffs at different stages of the game. Given a subset  $\mathcal{L}$  of SPEs, one may define the set of expected payoffs for both players at different times of the game. The set  $\mathcal{U}(\mathcal{L})$  (or just  $\mathcal{U}$ , when there is no confusion) denotes the set of expected payoffs for the players, across all possible SPEs in  $\mathcal{L}$ , computed *before* public randomization.  $\mathcal{V}$  is defined identically but computed *after* the realization of the randomization device  $P$ . In particular,  $\mathcal{U}$  is included in the convex hull of  $\mathcal{V}$ . Finally,  $\mathcal{W}$  consists of continuation payoffs *after* actions and payoffs for the current periods occurred, but *before* the proposal stage. Any payoff vector of  $\mathcal{W}$  is a mixture of three payoff vectors of  $\mathcal{U}$ , seen as continuation payoffs for the next period, according to whether player 1 or 2 gets to propose, or no one does. Thus, payoffs of  $\mathcal{W}$  are computed in terms of the next period.

Elements of  $\mathcal{U}$ ,  $\mathcal{V}$ , and  $\mathcal{W}$  are points of two-dimensional sets. For any point  $U$  of such a set, we let  $\pi_i(U)$  denote the  $i^{\text{th}}$  component of  $U$ , i.e.,  $i$ ’s continuation payoff at point  $U$ .

**DEFINITION 3** *A point  $A$  is  $q$ -renegotiation-proof if there exists  $\bar{\varepsilon} \in (0, 1/q)$  such that for all  $\varepsilon \leq \bar{\varepsilon}$  and  $p = q\varepsilon$ , there exists a stable norm  $\mathcal{N}$  such that  $A \in \mathcal{U}(\mathcal{N})$ . Moreover,  $A$  is renegotiation-proof if it is  $q$ -renegotiation-proof for all  $q$ ’s large enough.*

The coefficient  $q$  is inversely related to the amount of renegotiation frictions in the game: when  $q = 0$ , players never get a chance to renegotiate and the game reduces to a standard repeated game. As with the standard Folk theorem, any point of the feasible IR set (in the usual sense of repeated games, absent any renegotiation) is 0-sustainable. Our main objective is to characterize the set of sustainable payoffs. To do so, we first study the set of  $q$ -sustainable payoffs for any fixed  $q$ , and then let the renegotiation frictions go to zero.

In the definition above,  $A$  needs only belong to  $\mathcal{U}$  which, unlike  $\mathcal{V}$ , includes the initial use of a public randomization device. As it turns out however, this distinction is unimportant for the theorems below.



## 4 Main Result

### 4.1 Statement

We let  $v_i$  denote  $i$ 's minmax payoff in the stage game of the repeated game (absent any renegotiation).<sup>14</sup> The set of all feasible stage-game payoffs is a convex polygon. Similarly, let  $P^i$  denote the feasible payoff vector that provides  $i$ 's with his maximal payoff in the stage game.<sup>15</sup> The 'weak' Pareto frontier (consisting of all points which are not *strictly* Pareto dominated) lies between  $P_1$  and  $P_2$ .

Let  $v_1 = \max\{v_1; \pi_1(P_2)\}$  and  $v_2 = \max\{v_2; \pi_2(P_1)\}$ .

**THEOREM 1 (RENEGOTIATION-PROOF SET)** *Suppose that  $P_1 \neq P_2$ . Then, the following holds:*

**Sufficiency** *If*

$$\pi_i(A) > v_i \quad \text{for } i \in \{1, 2\}, \quad (1)$$

*then the point  $A$  is  $q$ -renegotiation-proof for all  $q \in \mathbb{R}_+$  and hence renegotiation-proof.*

**Necessity** *If  $A$  is  $q$ -renegotiation-proof, then*

$$\pi_i(A) \geq v_i + \max\left\{0; \frac{q/2}{1 + (q/2)} (\pi_i(P_{-i}) - v_i)\right\} \quad (2)$$

*for  $i \in \{1, 2\}$ . If  $A$  is renegotiation-proof, inequalities in (1) must hold for both players as weak inequalities.*

If  $P_1 = P_2$ , all stable norms are payoff equivalent and reduced to the singleton  $\{P_1\}$ . The only renegotiation-proof point is  $P_1$ , which is played forever. In this case, note players have perfectly aligned interests, as they both want to implement  $P_1$ . Our concept of renegotiation-proofness selects that point as the only possible outcome, as renegotiation frictions become negligible.

The statement of Theorem 1 can be visualized on Figure 1 for a fixed friction level of renegotiation. The orange domain represents the set of points which are known to be renegotiation-proof (i.e., part of a stable norm), while the red domain represents the additional points which *may* be renegotiation-proof. When  $q = 0$  (no renegotiation), the red domain extends all the way back to the minmax

<sup>14</sup>As usual, player  $-i$  is allowed to mix across actions to minmax  $i$ .

<sup>15</sup>If several such points exist, we choose the point among those with the lowest payoff for  $-i$ . In that case,  $P^i$  is not strictly Pareto dominated, but it will be Pareto dominated by points giving the same payoff to  $i$  and a strictly higher payoff to  $-i$ .

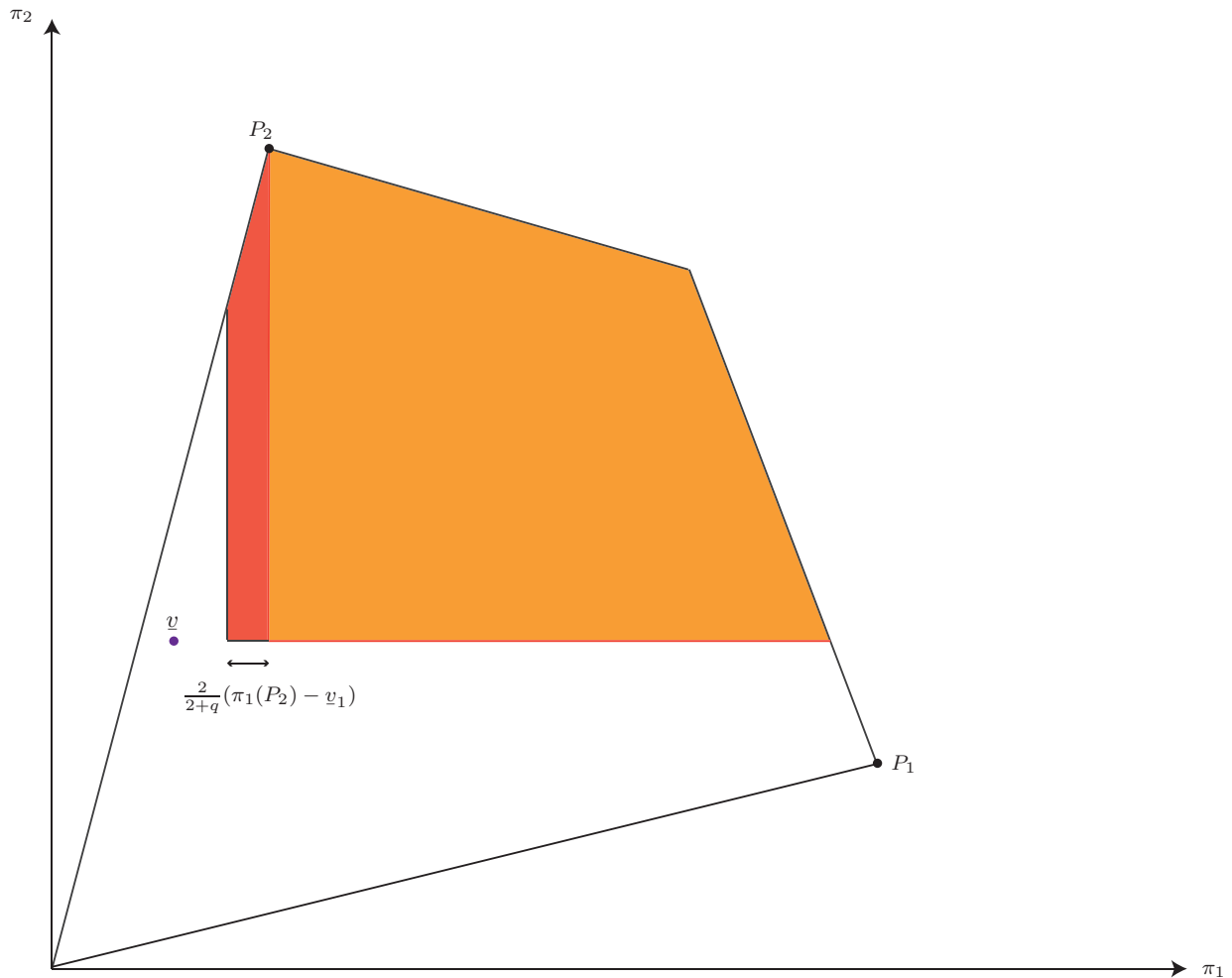


Figure 1: Necessary and sufficient conditions for renegotiation factor  $q$

point  $v$ , and we obtain the Folk Theorem. As the renegotiation frictions vanish ( $q \rightarrow +\infty$ ), the red domain disappears: necessary and sufficient conditions become identical (up to the boundary).

It is straightforward to characterize renegotiation-proof points when frictions vanish. This may be done graphically, and Figure 2 represents the corresponding sets for various configuration. In configuration (a), renegotiation constrains the set of implementable payoffs because the deterrence points  $P_1$  and  $P_2$  are too close to each other (all this is relative to the vector of minmax payoffs). Configuration (b) represents a fully cooperative game. The only renegotiation-proof outcome is the Pareto efficient point. In configuration (c), the punishment/reward points are sufficiently far apart, and the Folk Theorem is restored.

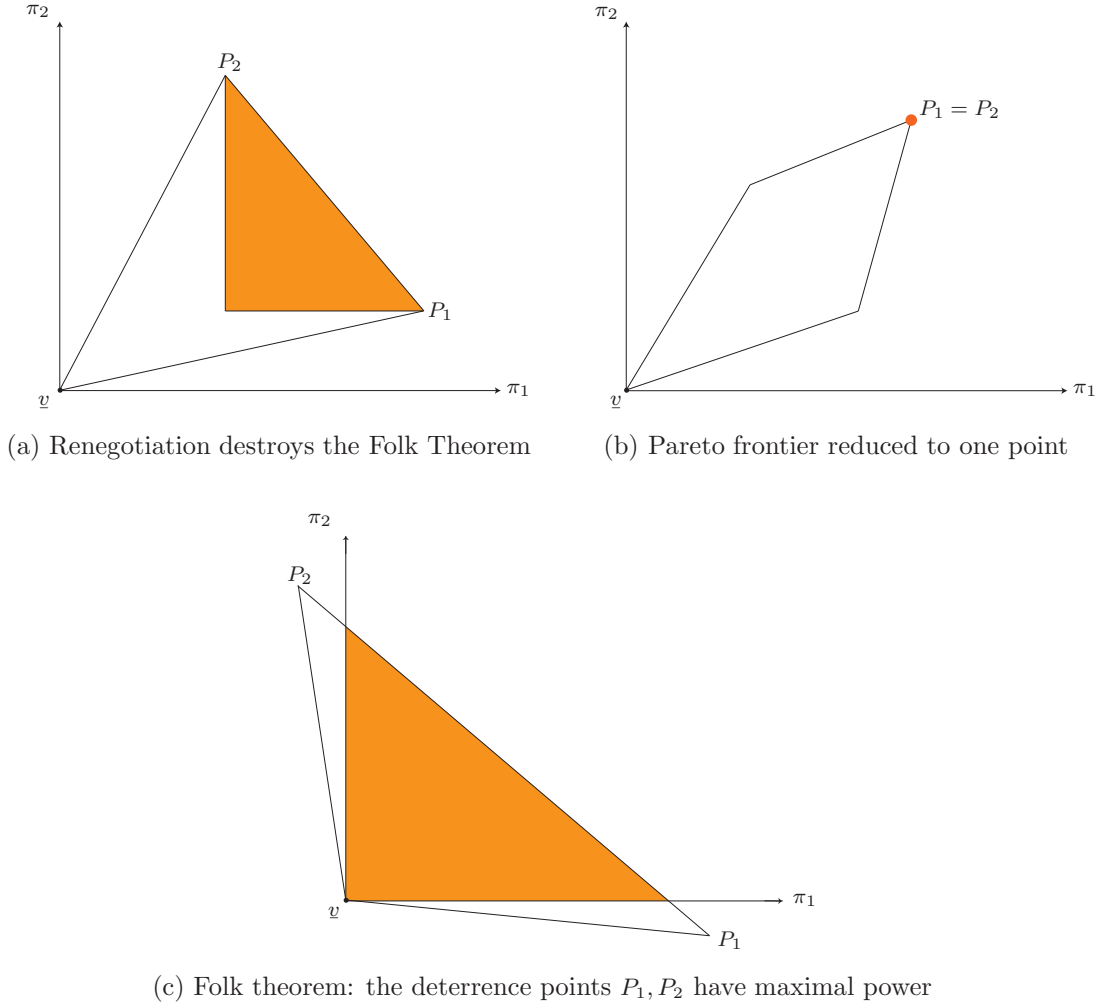


Figure 2: Renegotiation-proof sets for various configurations

## 4.2 Proof: Sufficient Conditions

*Outline* We construct, for any point  $A$  satisfying (1), a norm containing  $A$  and which is stable for any  $q \geq 0$ . The construction starts by choosing two points  $A_1, A_2$  such that  $A_i$  defines  $i$ 's worst possible payoff according to the norm.<sup>16</sup> Given any continuation payoff far away from  $A_i$ , it is always in  $i$ 's interest to follow the prescribed play in action, since any deviation provides a gain of order  $\varepsilon$  and can be punished by moving to  $A_i$ . The key is to choose  $A_i$  so that  $i$  is adequately incentivized near  $A_i$  and to complete the norm with enough equilibria to guarantee that the norm is stable. That last part is achieved by including some Pareto-optimal points  $Q_1, Q_2$  in the norm so

<sup>16</sup>Unless stated otherwise, payoffs are elements of the set  $\mathcal{U}$ , i.e., at the beginning of a period.

that for any credible proposal that  $-i$  may deviate to,  $i$  can always be rewarded, and  $-i$  punished, by rejecting  $-i$ 's proposal and have  $Q_i$  implemented.

For each player  $i$ , there are two cases to consider, depending on whether  $i$ 's minmax payoff  $v_i$  lies above or below  $\pi_i(P_{-i})$ . We treat the former case first.

Suppose, thus, that  $v_1 > \pi_1(P_2)$  and  $v_2 > \pi_2(P_1)$ , and consider any point  $A$  satisfying (1).

For  $\varepsilon$  small enough, the points  $A_1$  and  $A_2$  with coordinates

$$\pi_1(A_1) = v_1 + \varepsilon^{1/2}; \quad \pi_2(A_1) = \pi_2(A)$$

and

$$\pi_1(A_2) = \pi_1(A); \quad \pi_2(A_2) = v_2 + \varepsilon^{1/2}$$

are feasible and such that  $\pi_1(A_1) < \pi_1(A)$  and  $\pi_2(A_2) < \pi_2(A)$ .

The point  $A_1$  is implemented as follows, with a similar implementation for  $A_2$ .

1) Action stage: player 2 minmaxes player 1, possibly mixing between several actions  $a_{2j}$  and 1 best responds by a pure action  $a_{1,minmax}$  achieving his minmax payoff.

1a) If no deviation in action is observed – i.e., 1 chooses  $a_{1,minmax}$  and 2's realized action  $a_{2j}$  is in the support of the mixture minmaxing 1 – the continuation payoff  $B_{1j} \in \mathcal{W}$  is a function of  $a_{2j}$ , chosen so that i) 2 is indifferent between mixing actions  $a_{2j}$ , and ii) all  $B_{1j}$ 's give 1 the same continuation payoff. This latter condition implies that

$$\pi_1(A_1) = \varepsilon v_1 + (1 - \varepsilon)\pi_1(B_{1j}) \tag{3}$$

Note that all  $B_{1j}$ 's are within an  $\varepsilon$ -proportional distance of  $A_1$ .

1b) If 2 deviates in action (i.e., chooses an action outside of the mixture used to minmax 1), the continuation payoffs jump to the point  $A_2$ , mentioned above, which gives him the lowest possible payoff in the norm.<sup>17</sup> This punishment is clearly enough to incentivize 2, because any gain is of order  $\varepsilon$ , whereas  $A_2$  is arbitrarily close to 2's minmax payoff and thus at an  $\varepsilon$ -independent distance from  $\pi_2(A_1)$  (and, hence,  $\pi_2(B_{1j})$ 's)

1c) If 1 deviates in action, disregard this. Such a deviation is obviously suboptimal, since 1 was prescribed to best respond to being minmaxed by 2.

2) Proposal stage: the point  $B_{1j}$  is implemented as follows: if either 2 gets a chance to propose, or no player does, the continuations payoffs return to  $A_1 \in \mathcal{U}$ . (2 is prescribed to remain silent.) If

---

<sup>17</sup>More precisely, it jumps to the point  $B_{21}$ , which is the analogue of the point  $B_{11}$ , following the implementation of  $A_2$ .

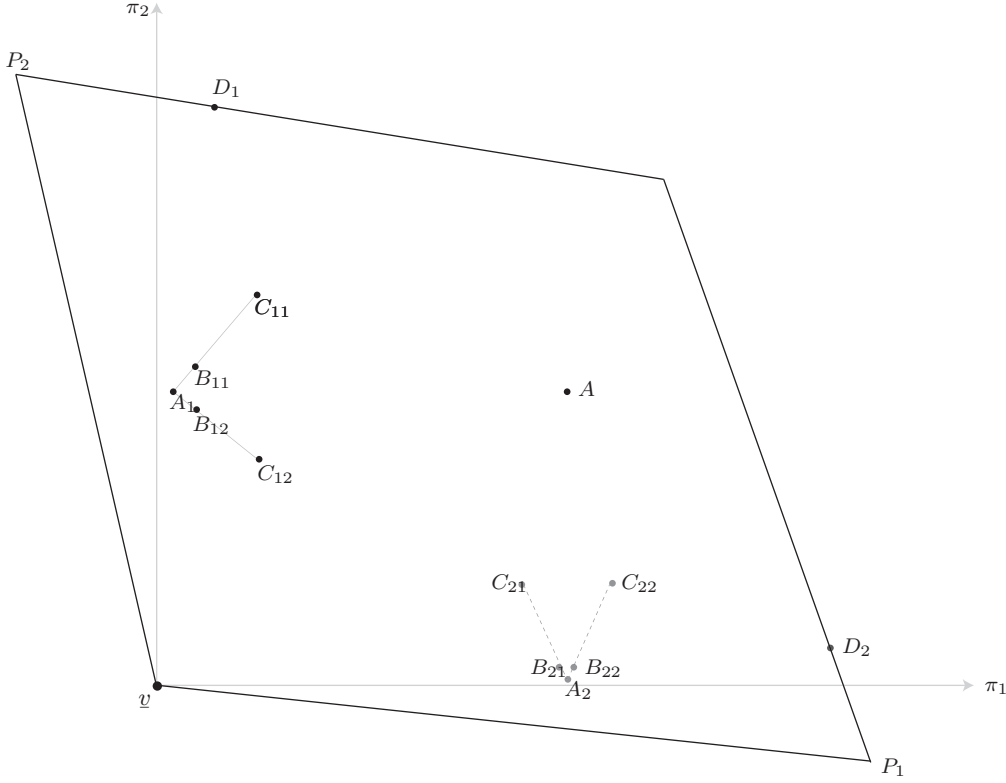


Figure 3: Construction of a stable norm

1 gets a chance to propose, he proposes a point  $C_{1j}$ , which lies on the line going through  $A_1$  and  $B_{1j}$ , and is chosen so as to satisfy the promise-keeping condition holds:

$$\pi_1(B_{1j}) = (1 - [p/2])\pi_1(A_1) + [p/2]\pi_1(C_{1j}) \quad (4)$$

Player 2 is prescribed to accept proposal  $C_{1j}$ . The points  $\{C_{1j}\}_j$  give the same payoff to 1, independently of  $j$ . Their implementation is described in 3) below.

2a) If 1 proposes any plan other than the one implementing  $C_{1j}$ , he is punished by the implementation of a point  $D_1$  chosen so that i)  $\pi_1(D_1) < \pi_1(C_{1j})$  and ii) 2 prefers  $D_1$  to 1's proposal. Precisely,  $D_1$  is defined as the point of the Pareto frontier that gives 1 a payoff of:

$$\frac{\pi_1(A_1) + \pi_1(C_{1j})}{2} \quad (5)$$

2b) If 2 deviates by making a proposal or rejecting 1's offer to move to  $C_{1j}$ , he is punished by the player-2 analogue of point  $D_1$ .

This construction is represented on Figure 3.

3) Next periods: the point  $C_{1j}$  is easily implemented, because it gives 1 a payoff of the order of  $\sqrt{\varepsilon}$  above what  $A_1$  and, hence,  $B_{1j}$ 's give him, as explained shortly. A deviation in action by 1 brings a gain of order  $\varepsilon$  and is punished by a drop of order  $\sqrt{\varepsilon}$  in 1's continuation payoff, and is thus suboptimal, for  $\varepsilon$  small enough. More precisely, the point  $C_{1j}$  can be implemented by a deterministic sequence of actions keeping players' continuation payoffs within a distance  $K\varepsilon$  from  $C_{1j}$ . The rules implementing that sequence are simple: play a deterministic action profile keeping continuation payoffs  $\varepsilon$ -close to  $C_{1j}$  and do not allow any proposal. If 1 deviates in actions, move to one of the points  $B_{1j}$ ; if he deviates in proposals, move to  $D_1$ . A similar rule is applied for player 2, who has even more to lose from deviating.

4) Finally, the point  $D_1$  lies at a distance of order  $\sqrt{\varepsilon}$  from  $A_1$ , and can therefore be implemented similarly to  $C_{1j}$  by a deterministic sequence of actions which keep the continuation payoff within a distance  $K\varepsilon$  from  $D_1$ . Any proposal is ignored.

We verify the claim that all  $C_{1j}$ 's lie at a  $\sqrt{\varepsilon}$ -proportional distance to the right of  $A_1$ . From 3 and 4, we get

$$\pi_1(A_1) = \varepsilon v_1 + (1 - \varepsilon)\pi_1(B_j) = \varepsilon v_1 + (1 - \varepsilon)[(1 - [q\varepsilon/2])\pi_1(A_1) + [q\varepsilon/2]\pi_1(C_{1j})]$$

Ignoring terms of order  $\varepsilon^2$  and higher, this implies that

$$\pi_1(A_1) = \varepsilon v_1 + (1 - [1 + \frac{q}{2}]\varepsilon)\pi_1(A_1) + [q\varepsilon/2]\pi_1(C_{1j}).$$

Subtracting  $\pi_1(A_1)$  from both sides and dividing by  $\varepsilon$  yields

$$\varepsilon^{1/2} = \pi_1(A_1) - v_1 = (q/2)[\pi_1(C_{1j}) - \pi_1(A_1)], \quad (6)$$

which shows the claim.

The direction of each vector  $\overrightarrow{A_1 C_{1j}}$ , which is also the direction of the vector  $\overrightarrow{A_1 B_{1j}}$  depends only on 2's action  $a_{2j}$ ; it does not change when  $\varepsilon$  goes to 0. This shows that for  $\varepsilon$  small enough

1.  $C_{1j}$  is a feasible payoff;
2.  $\pi_2(C_{1j})$  exceeds  $\pi_2(A_2)$  by an  $\varepsilon$ -independent value.

As explained above, the system of actions and proposals implementing  $A_i$ 's,  $B_{ij}$ 's and  $C_{ij}$ 's and  $D_i$ 's is incentive compatible in actions and proposals.

To conclude the construction of the norm, note that  $A$  gives each player  $i$  a payoff higher than  $A_i$ , by an amount that is bounded below away from zero and thus independent of  $\varepsilon$  as  $\varepsilon$  goes to

zero. One may therefore implement  $A$  by a deterministic sequence of actions, chosen so that the continuation payoffs stay within a distance  $K\varepsilon$  of  $A$ . Deviations in actions are punished by moving to  $B_{11}$  or  $B_{21}$ , depending on which of the players, 1 or 2, has deviated. Deviations in proposals are similarly punished by moving to  $D_1$  or  $D_2$ .

We now verify that the norm is stable. Within the norm, notice that whenever 1 gets a chance to make a proposal his payoff is at least  $\pi_1(D_1)$ . Since  $D_1$  is on the Pareto frontier, any proposal giving 1 strictly more than  $\pi_1(D_1)$  must give 2 less than  $\pi_2(D_1)$ . This means that  $D_1$  can serve as a punishment in case 1 makes such a proposal. The norm thus constructed is stable.

Consider now the second case, in which  $\underline{v}_1 \leq \pi_1(P_2)$  and/or  $\underline{v}_2 \leq \pi_2(P_1)$ . The construction of the norm is almost identical to the previous case. The only difficulty is that the difference  $\pi_1(A_1) - \underline{v}_1$  is now bounded below away from zero, whereas it was previously of order  $\varepsilon$ . This may lead to situations in which the points  $C_{1j}$  constructed above are no longer feasible and/or may give 2 a payoff lower than  $\pi_2(A_2)$ . This difficulty is easily addressed by adding, for each  $j$ , a new point  $E_{1j}$  lying on the segment  $[A_1B_{1j}]$  – and thus also on the line  $(A_1C_{1j})$  – such that if player 2 gets a chance to propose, or if nobody does, continuation payoffs jump to the point  $E_{1j}$ . The promise keeping condition (4) becomes

$$\pi_1(B_{1j}) = (1 - [p/2])\pi_1(E_{1j}) + [p/2]\pi_1(C_{1j}) \quad (7)$$

By choosing  $E_{1j}$  close enough to  $B_{1j}$ , one can make the point  $C_{1j}$  within a distance  $\sqrt{\varepsilon}$  of  $B_{1j}$  and, hence, of  $A_1$ . This guarantees that  $C_{1j}$  is feasible and does not drop below  $\pi_2(A_2)$ , so that the rest of the argument for the first case can be applied. Finally, whenever the point  $E_{1j}$  must be implemented in the next period, we use the public randomization device to implement it as a probabilistic mixture between  $A_1$  and  $C_{1j}$ .

### 4.3 Proof: Necessary Conditions

When player  $i$ 's minmax  $\underline{v}_i$  is higher than  $\pi_i(P_{-i})$ , the necessary condition simply states that  $A$ 's payoff must belong to the feasible IR set (i.e., be above the minmax). The only interesting case, therefore, is when  $\underline{v}_i < \pi_i(P_{-i})$ . We will establish the necessary condition corresponding to player 1. The proof for the other case is identical and independent.

Let us thus assume that  $\pi_1(P_2) > \underline{v}_1$  and suppose, by contradiction, that there is a point  $A$  such that  $\pi_1(A) < \underline{v}_1 = \underline{v}_1 + \frac{q/2}{1+(q/2)}(\pi_1(P_2) - \underline{v}_1)$ , which is  $q$ -renegotiation-proof. This means that one can construct, for any  $\varepsilon$  small enough, a stable norm  $\mathcal{N}$  that contains  $A$  as one of its continuation payoffs.

In this case, we first build a proposal which yields the payoff vector  $P_2$ . The point  $P_2$  is easily shown to be the payoff vector of an equilibrium of the game, and hence has to be considered in the definition 2 to check that  $A$  belongs to a stable norm.<sup>18</sup> However, we argue later in the paper (Section 5) that our necessary and sufficient conditions are unchanged if one restricts attention, in Definition 2, to *credible* proposals instead of arbitrary ones. A proposal is credible with respect to some norm  $\mathcal{N}$  if it is an SPE such that any deviation (in action or proposal) leads to reversal to an equilibrium of the norm  $\mathcal{N}$ .<sup>19</sup> To this end, we now show that  $P_2$  is the payoff of a credible proposal. The SPE implementing  $P_2$  is constructed as follows: players are prescribed to play, in all periods, the pure action profile with payoff  $P_2$ , and to abstain from making any proposal. Any deviation, whether in action or in proposal, triggers the implementation of point  $A$ . Clearly, player 2 cannot benefit from deviating as he is getting his highest possible payoff in the game. Moreover, the difference  $\pi_1(P_2) - \pi_1(A)$  is by assumption bounded below by  $\frac{q/2}{1+(q/2)}(\pi_1(P_2) - v_1)$ , which is  $\varepsilon$ -independent. Therefore, 1 cannot benefit from deviating either: a deviation in action may create an immediate gain of order  $\varepsilon$ , but triggers a drop in continuation payoffs that is  $\varepsilon$ -independent and dominates the gain. A deviation in proposal triggers  $A$ , which again is detrimental to 1. We thus have a constructed an equilibrium of the game with payoff  $P_2$ .

Let  $C_1$  denote 1's infimum payoff in  $\mathcal{N}$  when it is his turn to propose. Since  $P_2$  is a possible proposal payoff, and since it lies on the Pareto frontier,  $A$  is a payoff of  $\mathcal{N}$  only if

$$\pi_1(P_2) \leq C_1$$

We will show that it is impossible.

Let us denote  $A_1 = \inf_{V \in \mathcal{V}(\mathcal{N})} \pi_1(v)$ ,  $B_1 = \inf_{W \in \mathcal{W}(\mathcal{N})} \pi_1(w)$ ,  $D_1 = \inf_{U \in \mathcal{U}(\mathcal{N})} \pi_1(u)$ .

Consider a sequence  $\{V_m\} \in \mathcal{V}(\mathcal{N})$  such that  $\pi_1(V_m) \rightarrow A_1$ . For any  $V_m$  there is an action that implements it the first period of the corresponding SPE. However, if player 1 deviates, he can guarantee himself an immediate payoff of at least  $v_1$ , and the worst punishment for him after deviation gives him at least  $B_1$ . Therefore,  $\pi_1(V_m) \geq \varepsilon v_1 + (1 - \varepsilon)B_1$ . Since this inequality holds for all  $V_m$  we obtain, taking the limit:

$$A_1 \geq \varepsilon v_1 + (1 - \varepsilon)B_1 \tag{8}$$

Because any element of  $\mathcal{U}(\mathcal{N})$  lies in the convex hull of  $\mathcal{V}(\mathcal{N})$ , and  $C_1$  is a mixture of points in

---

<sup>18</sup>By the Folk Theorem,  $P_2$  can be implemented by an SPE of the repeated game without renegotiation. By treating all proposals as cheap talk,  $P_2$  can they also be implemented as SPE of the game with renegotiation.

<sup>19</sup>See Definition 4 for a precise definition.



$\mathcal{U}(\mathcal{N})$ ,<sup>20</sup> we have

$$C_1 \geq D_1 \geq A_1$$

Consider now a sequence  $\{W_m\} \in \mathcal{W}(\mathcal{N})$  such that  $\pi_1(W_m) \rightarrow B_1$ . Any element  $W_m$  is a weighted average of an expected payoff vector  $EU_m^1$  whenever 1 gets a chance to propose, an expected payoff vector  $EU_m^2$  when it is 2's turn to propose, and a payoff vector  $U_m^0$  in case no one gets to propose:

$$W_m = (p/2)(EU_m^1) + (p/2)(EU_m^2) + (1-p)(U_m^0) \quad (9)$$

We note that  $EU_m^1$  is a mixture of elements of  $\mathcal{U}$  resulting from 1's mixture over proposals and 2's mixture over his acceptance decision. Similarly,  $EU_m^2$  is a mixture of elements of  $\mathcal{U}$ .

Since all elements  $U_m$ 's belong to  $\mathcal{U}(\mathcal{N})$ , we have

$$\pi_1(EU_m^2) \geq A_1,$$

$$\pi_1(U_m^0) \geq A_1.$$

Equation (9) then implies that

$$\pi_1(W_m) \geq (1 - [p/2])A_1 + [p/2]\pi_1(EU_m^1)$$

Recalling that  $C_1$  denotes 1's infimum payoff when he gets to propose, we get

$$\pi_1(W_m) \geq (1 - [p/2])A_1 + (p/2)C_1.$$

Taking limits,

$$B_1 \geq (1 - [p/2])A_1 + (p/2)C_1$$

or

$$B_1 \geq (1 - [q\varepsilon/2])A_1 + (q\varepsilon/2)C_1. \quad (10)$$

Combining (8) and (10), we conclude that

$$A_1 \geq \varepsilon v_1 + (1 - \varepsilon)[(1 - [q\varepsilon/2])A_1 + (q\varepsilon/2)C_1]$$

or, ignoring terms of order  $\varepsilon^2$  in right-hand side,

$$A_1 \geq \varepsilon v_1 + (1 - [1 + (q/2)]\varepsilon)A_1 + (q\varepsilon/2)C_1.$$

Subtracting  $A_1$  on both sides of the last equation and dividing by  $\varepsilon$ , we obtain

$$0 \geq v_1 - [1 + (q/2)]A_1 + (q/2)C_1 \quad (11)$$

---

<sup>20</sup>It is a mixture over the payoffs obtained following  $-i$ 's acceptance/rejection decision to  $i$ 's proposal.

From  $A_1 \leq \pi_1(A)$ ,  $C_1 \geq \pi_1(P_2)$ , and  $\pi_1(A) < v_1 = v_1 + \frac{q/2}{1+(q/2)}(\pi_1(P_2) - v_1)$ , we get

$$0 < v_1 - [1 + (q/2)]A_1 + (q/2)C_1$$

which contradicts (11). This shows the necessary condition for player 1. An identical reasoning for player 2 shows the second necessary condition. This proves the result. When  $P_1 = P_2$ , a similar reasoning implies the result.

## 5 Equivalent notions of stability

### 5.1 Credible innovations

Our definition of stability allowed players to implement any continuation equilibrium that is proposed and accepted at some point of the game. When players are used to a given norm  $\mathcal{N}$ , however, one may question why they would take such a proposal seriously. As it turns, both the necessary and the sufficient conditions of Theorem 1 remain identical if one restricts proposals to a much smaller subset.

**DEFINITION 4** *Given a norm  $\mathcal{N}$ , an  $\mathcal{N}$ -credible (or just “credible”, when there is no confusion) proposal is an SPE such that any off-equilibrium play (action, proposal, or acceptance decision) is followed by a continuation in  $\mathcal{N}$  at the next period;*

A credible proposal is thus an SPE which can be supported under the assumption that any deviation will be followed by a reversal to the norm. For example, if a norm includes a harsh punishment equilibrium for both players. Then, it supports many credible equilibria, any deviation of which triggers a reversal to the norm and, more precisely, to the punishment equilibrium.

**DEFINITION 5** *A norm  $\mathcal{N}$  is stable with respect to credible innovations if it satisfies the refinement of Definition 2 for all  $\mathcal{N}$ -credible proposals.*

Clearly Definition 5 is more permissive than Definition 4, because it imposes the refinement over a smaller set of proposals. However, we get the following result.

**THEOREM 2** *All the conclusions of Theorem 1 continue to hold if the norms sustaining renegotiation-proof points are only required to be stable with respect to credible innovations.*

The proof is straightforward. Because this second definition of stability is more permissive, our construction for the sufficiency condition also works in this case. Moreover, it is easy to check that the proposals considered to derive the necessary conditions are credible.

## 5.2 Set-theoretic definition

The norms that we defined earlier were **open** in the sense that they allowed the possibility that players might depart from the norm in case an equilibrium outside of the norm is proposed and accepted. This openness is necessary if we want to treat such proposals seriously. It is possible however to bring our work closer to the set-theoretic approach that was studied in the late eighties and early nineties. In fact, we show in this section that our earlier analysis can be entirely reexpressed in terms of set-theoretic definitions, yielding exactly the same characterization.

In this section, players do not actually take equilibria outside of the norm seriously. They consider their norm as the only possible outcomes. We thus begin by “closing” our definition of a norm:

**DEFINITION 6** *A subset  $\mathcal{N}$  of  $\mathcal{S}$  is a **closed norm** if for any  $m \in \mathcal{N}$ ,  $m_{+1}[m](P, \mathbf{a}, i, \mu_i, D_{-i}) \in \mathcal{N}$ .*

The only difference with the earlier definition of a norm is that continuations belong to the norm, including when a proposal outside of the norm is made and accepted. Next, our earlier definition of stability is translated into set-theoretic terms. To keep in line with the previous section, we state the definition for credible proposals. It should be noted however that the same definition dropping “credible” yields the same set.

**DEFINITION 7** *A closed norm  $\mathcal{N}$  is **stable** if the following holds: consider any SPE of  $\mathcal{N}$  and history at which  $i$  gets a chance to propose and let  $\hat{U}_i$  denote  $i$ 's continuation payoff. Then, for any credible proposal with payoff vector  $U$  which gives  $i$  a payoff  $U_i > \hat{U}_i$ , there exists a payoff vector  $U'$  in the norm such that  $\pi_{-i}(U') \geq \pi_{-i}(U)$  and  $\pi_i(U') \leq \hat{U}_i$ .*

**THEOREM 3**

1. *For any closed norm  $\mathcal{N}^c$ , there exists an open norm  $\mathcal{N}^o$  which has the same payoff set, and vice versa.*
2. *For any stable closed norm  $\mathcal{N}^c$ , there exists a stable open norm  $\mathcal{N}^o$  which has the same payoff set, and vice versa.*

*Proof.*

1. Any closed norm  $\mathcal{N}^c$  is an open norm as well, so the first statement is trivially true. Now consider any open norm  $\mathcal{N}^o$ . To construct a payoff-equivalent closed norm  $\mathcal{N}^c$ , we modify each plan/equilibrium  $m$  of  $\mathcal{N}^o$  as follows:  $m$ 's rules on and off the equilibrium path are kept unchanged except when a player, say  $i$ , makes a proposal  $\mu_i$  which is off the equilibrium path. In this case, because  $\mathcal{N}^o$  is only an open norm, the continuation equilibrium if  $-i$  accepts the proposal need not lie in  $\mathcal{N}^o$ . Following such a proposal, players are instead prescribed to behave as if  $i$  had remained silent. The new rules define an equilibrium: when playing the original equilibrium  $m$ ,  $i$  was not making the proposal  $\mu_i$  anyway, so removing this option does not affect *equilibrium* behavior and payoffs. By construction, the set of modified equilibria form a closed norm  $\mathcal{N}^c$ , and because each equilibrium of  $\mathcal{N}^o$  has been modified into a single payoff-equivalent equilibrium of  $\mathcal{N}^c$ , the norms are payoff equivalent.

2. We start with the observation that if two norms  $\mathcal{N}^c$  and  $\mathcal{N}^o$  have the same payoff sets, then any proposal that is credible according to either norm is credible according to the other norm.

We now consider any renegotiation-proof open norm  $\mathcal{N}^o$  and construct the corresponding closed norm  $\mathcal{N}^c$  as in Part 1. To show that  $\mathcal{N}^c$  is renegotiation-proof, consider any SPE  $m$  of  $\mathcal{N}^c$ , history, and credible proposal  $U$  such that  $\pi_i(U)$  is strictly greater than  $i$ 's continuation payoff  $\hat{U}_i$ . From the above observation,  $U$  is also credible for  $\mathcal{N}^o$ . For the equilibrium  $\tilde{m}$  of  $\mathcal{N}^o$  corresponding to  $m$ , and the same history, if  $i$  proposes  $U$ ,  $-i$  must reject it with positive probability (for otherwise  $\pi_i(U)$  would coincide with  $\hat{U}_i$ ). Let  $U'$  denote the continuation payoff if  $-i$  rejects  $U$ . By renegotiation-proofness of  $\mathcal{N}^o$ ,  $-i$  knows that if he accepts  $U$  it will be implemented. Since it is weakly optimal for  $-i$  to reject  $U$ , it must therefore be the case that  $\pi_{-i}(U') \geq \pi_{-i}(U)$ . Moreover, it must also be the case that  $\pi_i(U') \leq \hat{U}_i$ , for otherwise it would be strictly optimal for  $i$  to deviate by proposing  $U$ , and  $\tilde{m}$  would not be an equilibrium. Using this  $U'$  in Definition 7, this implies that  $\mathcal{N}^c$  is renegotiation-proof.

Next, consider any renegotiation-proof closed norm  $\mathcal{N}^c$ . To construct a payoff-equivalent renegotiation-proof open norm  $\mathcal{N}^o$ , we simultaneously modify all SPE's of  $\mathcal{N}^c$ . The modification proceeds in two steps, and is based on the recursive definition a plan. Recall that a plan at time is a prescription of actions, proposals and acceptance decisions for the next period (each depending on what happened in earlier stages), along with a continuation plan resulting from these stages to applied in the period after next. In Step 1, we modify the prescriptions for time  $t + 1$ , and still use plans of  $\mathcal{N}^c$  as continuation plans. The purpose of this step is to a prescription compatible with the requirement that if a credible proposal is made and accepted, then it has to be played. In Step 2, we replace these continuation plans of  $\mathcal{N}^c$  by those built in Step 1, to guarantee that the rule applies at all periods, guaranteeing that off-path proposals which are accepted are implemented,

so that Definition 5 holds at all periods.

Consider any SPE  $m$  of  $\mathcal{N}^c$ . We modify  $m$  as follows. For the modified SPE  $\tilde{m}$ , the action stage and on-path proposals are prescribed exactly as in  $m$ .<sup>21</sup> Now consider a history at which  $i$ 's gets a chance to propose and makes any proposal  $U$  which is not prescribed by  $m$  but which is  $\mathcal{N}^c$ -credible. If  $-i$  accepts the proposal, we construct  $\tilde{m}$  by prescribing that players implement this proposal.<sup>22</sup> If the proposal gives  $i$  a strictly higher payoff than his equilibrium continuation payoff  $\hat{U}_i$ , then by stability of  $\mathcal{N}^c$ , there must exist a payoff vector  $U'$  corresponding to some equilibrium  $m'$  of  $\mathcal{N}^c$ , which gives player  $-i$  at least as much as  $U$ , and which gives player  $i$  at most  $\hat{U}_i$ . We prescribe playing the equilibrium corresponding to  $U'$  in case player  $-i$  rejects the proposal. If  $U$  does not improve upon  $i$ 's equilibrium continuation payoff, we prescribe playing the continuation equilibrium corresponding to any of  $i$ 's equilibrium proposal in case  $-i$  rejects  $U$ . Finally, if  $i$  makes a non-credible proposal, the proposal is ignored as if  $i$  had stayed silent.

We now verify that  $\tilde{m}$  is an SPE that yields the same payoff as  $m$ . Since  $\tilde{m}$  prescribes the same actions as  $m$ , players are incentivized to follow the same as actions as those prescribed by  $m$ . If  $i$  gets a chance to make a proposal, any proposal prescribed by  $m$  (and hence  $\tilde{m}$ ) yields the same continuation payoff as in  $m$ . If player  $i$  makes a credible, off-equilibrium proposal that improves upon his equilibrium payoff, then player  $-i$  is incentivized to reject it, and  $i$ 's continuation payoff is weakly lower than his continuation payoff. It is never optimal for  $i$  to make a credible proposal that lower's his equilibrium payoff, regardless of  $-i$ 's acceptance decision. Finally, we replace all continuation plans by their modified versions.

There remains to verify that the set consisting of all modified equilibria forms a stable open norm, denoted  $\mathcal{N}^o$ , which is payoff equivalent to  $\mathcal{N}^c$ . First, we notice that continuation equilibria outside of  $\mathcal{N}^o$  may only arise when a player make an off equilibrium proposal (which, by construction, also has to be credible) which is accepted by the other player. Thus,  $\mathcal{N}^o$  is an open norm. By construction, each element of  $\mathcal{N}^o$  corresponds to exactly one element of  $\mathcal{N}^c$ , which yields the same expected payoffs. Therefore, the norms are payoff equivalent. As observed earlier, this implies that they have the same set of credible proposals. This, in turn, implies that any credible proposal of

---

<sup>21</sup>We need to make another small modification to  $m$  issue whenever  $i$  is proposing a continuation  $\mu$  outside of the norm  $\mathcal{N}^c$ , which  $-i$  is supposed to accept, and which is followed by a continuation  $\mu'$  in the norm  $\mathcal{N}^c$  (as it should, since the norm is closed). In that case, we can replace this play by  $i$  proposing instead  $\mu'$  and have it accepted by  $-i$ . That this modification can be done while preserving the equilibrium is straightforward to check. In fact, any SPE of the game can be turned into a payoff equivalent "truthful" SPE of the game, i.e., one in which any proposal that is made and accepted *on the equilibrium path* is implemented. See Section 6.

<sup>22</sup>At this point, we do not know yet that the proposal is  $\mathcal{N}^o$ -credible. We only know that it is  $\mathcal{N}^c$ -credible. However, the norm  $\mathcal{N}^o$  that we are constructing will be payoff equivalent to  $\mathcal{N}^c$  and hence have the same credible proposals.

$\mathcal{N}^o$  that is accepted is played and, hence, that  $\mathcal{N}^o$  is stable. ■

## 6 Discussion

Understanding and tractably modeling renegotiation in repeated games has been a longstanding challenge. Nevertheless, the protocol and concepts studied in the paper lead to a particularly simple characterization of stable norms and renegotiation-proof equilibria. We hope our characterization can serve as a useful benchmark for applied economists who need to incorporate renegotiation in their models. This section discusses several extensions or variations which were omitted until now to simplify the exposition of our main ideas.

### Truthful equilibria

To be interesting any analysis of renegotiation based on an explicit protocol must give some bite to accepted proposals. We have kept our analysis as simple as possible by imposing our refinement on stable norms only. In fact, it is easy to see that we could have focused our analysis on “truthful equilibria,” in which any (equilibrium) proposal that is accepted is implemented. Indeed, it is easy to see to prove that any equilibrium of the game has a payoff-equivalent truthful equilibrium. For example, consider an equilibrium and history at which  $i$  proposes  $\mu$ ,  $-i$  accepts it, and another plan  $\mu'$  is implemented. It is easy to modify the initial equilibrium by having  $i$  propose  $\mu'$  instead and have it accepted. There are other minor issues to address, but it is easy to develop the argument to build a truthful equilibrium.

By a similar reasoning, any equilibrium is equivalent to another equilibrium in which player proposals are always accepted in equilibrium. For example, if a proposal is rejected with probability 1 after  $\mu$  is proposed, leading to a continuation  $\mu'$ , then one can replace this by an equilibrium in which  $i$  proposes  $\mu'$  instead and have it accepted. Even if  $-i$  was mixing between acceptance and rejection, leading to two continuations  $\mu'$  and  $\mu''$ , one can have  $i$  instead proposed the mixture of  $\mu'$  and  $\mu''$ , corresponding to  $-i$ 's acceptance probabilities, and have it accepted with probability 1 by  $-i$ . Such mixture is always achievable by using the public randomization device at the beginning of the next period.

While such a focus is natural, it is not necessary for the results and we chose not to discuss it until to avoid cluttering the analysis.

## Asymmetric proposing probabilities and bargaining power

It is easy to extend the analysis to a protocol in which one of the players has a higher probability factor  $q_i$  of proposal than the other player. The sufficient conditions are unchanged in this setting, but the necessary conditions become tighter for the player whose proposal probability is higher, which translates a higher minimal guaranteed payoff for that player, across all renegotiation-proof equilibria. To see this starkly suppose that  $v_1 < \pi_1(P_2)$  and  $v_2 < \pi_2(P_1)$  (configuration (a) in Figure 2), so that renegotiation potentially benefits both players, compared to the minmax payoffs, and consider the case in which 1 can make frequent proposals while 2 never gets a chance to make a proposal (i.e.,  $q_1$  is arbitrarily large while  $q_2 = 0$ ). Then, 2's minimal guaranteed renegotiation-proof payoff collapses to his minmax payoff, while 1 is guaranteed to get a payoff of at least  $\pi_1(P_2)$ .

Asymmetric proposing probabilities and bargaining power

## Fixed discount rate

Our objective is to understand how renegotiation affects the strategic behavior of patient players and in particular how it affects the players' ability to cooperate and to implement punishments. It would also be interesting to study how renegotiation affects more impatient players. Our protocol and analysis could be used to this end.

Impatience can have a dual effect on repeated games with renegotiation. As with standard repeated games, it reduces the force of future punishments and hence the incentives to cooperate in the short term in the presence of profitable deviations. However, impatience also weakens the impact of future renegotiation. It is therefore possible that impatience weakens the impact of renegotiation on repeated games. This conjecture can be explained most starkly when players are perfectly impatient ( $\delta = 0$ ). In this case, renegotiation has of course no impact and the set of equilibria collapses to the repetitions of static Nash equilibria.

As with the analysis of standard repeated games without renegotiation, however, we expect the analysis of this case to be much more difficult.

### Three or more players

The analysis has focused on the case of two players, which is a common restriction in studies of renegotiation given the complexity of the problem (e.g., Farrell and Maskin (1989), Benoît and Krishna (1993), and Miller and Watson (2013).<sup>23</sup>). A natural follow-up of this work is to extend the analysis to three or more players. This raises new conceptual issues: Can proposals be targeted toward a subset of individuals? What happens if only a subset of the agents agrees to renegotiate? Because the protocol and concepts of the present paper lead to a particularly tractable analysis, we are hopeful to successfully explore this extension in future research.

---

<sup>23</sup>Abreu et al. (1993) focus instead on symmetric equilibria.



## References

- ABREU, D., PEARCE, D. (1991) “A perspective on renegotiation in repeated games,” in R. Selten (ed.), *Game Equilibrium Models*, Springer Berlin Heidelberg.
- ABREU, D., PEARCE, D., AND E. STACCHETTI (1993) “Renegotiation and Symmetry in Repeated Games,” *Journal of Economic Theory*, Vol. 60, pp. 217–240.
- ASHEIM, G. (1991) “Extending Renegotiation-Proofness to Infinite Horizon Games,” *Games and Economic Behavior*, Vol. 3, pp. 278–294.
- BENOÎT, J.-P., KRISHNA, V. (1993) “Renegotiation in Finitely Repeated Games,” *Econometrica*, Vol. 61, pp. 303–323.
- BERGIN, J., MACLEOD, B. (1993) “Efficiency and Renegotiation in Repeated Games,” *Journal of Economic Theory*, Vol. 61, pp. 42–73.
- BERNHEIM, B.D., RAY, D. (1989) “Collective Dynamic Consistency in Repeated Games,” *Games and Economic Behavior*, Vol. 1, pp. 295–326.
- DEMARZO, P. (1988) “Coalitions and Sustainable Social Norms in Repeated Games,” *mimeo*, Stanford University.
- FARRELL, J. (1983) “Credible Repeated Game Equilibria,” *Unpublished Manuscript*.
- FARRELL, J., MASKIN, E. (1989) “Renegotiation in Repeated Games,” *Games and Economic Behavior*, Vol. 1, pp. 327–360.
- MILLER, D., WATSON, J. (2013) “A Theory of Disagreement in Repeated Games with Bargaining,” *Econometrica*, Vol. 81, pp. 2303–2350.
- PEARCE, D. (1987) “Renegotiation-Proof Equilibria: Collective Rationality and Intertemporal Cooperation,” *Cowles Foundation Discussion Paper*, No. 855.