# Micro-microfoundations: Strategic Preference Formation and Policy Design

(preliminary version prepared for the 66th AFSE meeting)

Guilhem Lecouteux[*]

March 6, 2017

## Abstract

This paper develops a model of strategic preference formation: I assume that players can choose in a first stage the weights they assign to the other players' material payoff, and then determine the optimal weights each player should choose so as to maximise her material payoff. I highlight a systematic relation between supermodularity (submodularity) and the formation of cooperative (competitive) preferences. I then investigate the implications of this framework for the design of public policies, and show in the case of climate change negotiations that international agreements relying on technology standards with trade sanctions rather than objectives of pollution abatement are more likely to succeed, since they create a coordination game and cut the strategic substituability of the initial game. *Journal of Economic Literature* classification numbers: C72, D01, Q20.

**Keywords:** preference formation, strategic commitment, interdependent preferences, supermodularity, climate change negotiation, microfoundations.

# 1 How public policies shape our preferences

The Lucas critique argues that econometric models should integrate individual optimisation behaviours when assessing the implications of economic policies, and constitutes a core argument of the microfoundations program in macroeconomics:

> 'This essay has been devoted to an exposition and elaboration of a single syllogism: given that the structure of all econometric model consists of optimal decision rules of economic agents, and that optimal decision rules vary systematically with changes in the structure of series relevant to the decision maker, it follows that any change in policy will systematically alter the structure of econometric models.' (Lucas, 1976, p.41)

The critique stresses that public policies may impact individual behaviours, since altering the strategic environment is likely to induce a new optimal behaviour regarding the satisfaction of one's preferences. We can however notice that this interpretation implicitly assumes that the underlying preferences remain identical across policy regimes: individual behaviour may evolve *only* because the new institutional setting implies a new optimal decision, and not because those underlying preferences may also change.

Consider for instance crowding-out effects and the impact of the introduction of pecuniary incentives on pro-social motives: numerous works highlighted that those kinds of policies can backfire, since the intrinsic motivation for a social objective disappears and is replaced by an extrinsic motivation. This suggests that individual preferences are likely to change, and the satisfaction of those new preferences can lead to a worse situation than before the implementation of the policy (see for instance Titmuss (1970) on blood donation, Frey & Oberholzer-Gee (1997) on the willingness to accept a NIMBY project, and also Ostman (1998) and Cardenas et al. (2000) on the management of common-pool resources).

So as to properly calibrate public policies, we should therefore not only anticipate the new optimal choice of the individuals, but also anticipate the new preferences induced by the policy. As an illustration of the approach developed in this paper, consider the following illustration:

**Symmetric Cournot game**: we are playing a symmetric Cournot game. Suppose that I decide to maximise my profit minus $\sigma\%$ of your profit (I want therefore to maximise my profit as well as the difference between our profits): if you know that I am an aggressive player, then you know that I am likely to produce more than my Nash output. You then reduce your production and we end up in a situation in which you play your best reply to my somewhat 'irrational' action (my production is indeed not a best reply to yours). I am then producing more than my Nash output and you less: if I choose the adequate level of $\sigma$, then the resulting equilibrium can actually correspond to the Stackelberg equilibrium in which I would be the leader and you the follower.

In this game, although I did not directly maximise my profit, we reached an outcome in which I obtained a higher profit than if I had directly maximised my profit. The idea that players can benefit from such *strategic commitments* — i.e. voluntary deviations from payoff maximisation — has been studied since at least von Stackelberg (1934), with the introduction of timing in oligopoly, and Schelling (1960) in the context of coordination games. It is also central in the literature on strategic delegation (e.g. Fershtman & Kalai (1997), Fershtman & Gneezy (2001), Sengul et al. (2012)), and is the core mechanism in the indirect evolutionary approach (e.g. Güth & Yaari (1992), Samuelson (2001), Heifetz et al. (2007a,b)). Furthermore, some experimental findings suggest that players progressively learn to make the optimal strategic commitment (Fischer et al., 2006, Poulsen & Roos, 2012).

A common feature of those approaches is the distinction between the function that determines the *gain* of the individual (her material payoff or fitness for instance) and the function that determines the *choice* of the individual (her preferences). While players (or the individuals who take the decision on their behalf) are *utility* maximisers, they are not necessarily *payoff* maximisers. In the case of the symmetric Cournot game, the strategy chosen by firm 1 does not maximise her payoff $\Pi_1$ (given the strategy of player 2), but her utility $U_1 = \Pi_1 - \sigma\Pi_2$. It is then assumed that there exists a 'preferences game', in which individual utility functions are chosen in a first stage (this game is not conscious in the indirect evolutionary approach, and is the result of evolutionary pressures), so as to maximise *in fine* the material payoff in the second stage.

The object of the paper is not to focus on a specific type of strategic commitment,

but rather to study the equilibrium of this preferences game, i.e. an optimal strategic commitment. I will therefore assume that the players are able to directly choose their preferences in the first stage game, without discussing *how* they keep their commitment: the optimal strategic commitments characterised in this paper will therefore be implemented if and only if the players have at their disposal a mechanism allowing them to keep their strategic commitments (contracts, evolutionary pressures for instance).

This paper is organised as follows. I firstly present a model of strategic preference formation in which players choose the weights they attribute to other players in their utility function so as to maximise their payoff (section 2). I then show that the players generally choose non-null weights and highlight that supermodular (respectively submodular) games tend to generate cooperative (aggressive) preferences (section 3). I then argue that an objective of public policies should be to alter the strategic environment of the game so as to facilitate the emergence of cooperative behaviours. I illustrate this point by studying climate change negotiations (section 4).

## 2   Model

In this section, I start by clarifying the notion of strategic preferences. I then introduce technical notations and define a notion of equilibrium characterising a strategy profile immune to individual strategic commitments. I illustrate those different notions by studying a public good game.

### 2.1   Bluff and commitment

Recall the Cournot competition discussed above: since being aggressive with firm 2 may be *in fine* beneficial to firm 1, it is in the interest of firm 1 to choose its level of production not only to maximise its profit, but also to maximise the difference between its profit and the profit of its opponent. From the perspective of firm 1, its *relative* success may therefore matter more than its *absolute* success, since by trying to outperform firm 2 rather than merely maximising its payoff, firm 1 is likely to obtain a higher profit than if it has adopted a profit-maximising strategy.

4

A question that may arise is then the status we should give to those different types of preferences: if firm 1 decided to adopt strategic preferences different from its material payoff, can we still say that the true objective of the firm is payoff maximisation? It would actually be more accurate to say that the true objective of firm 1 is to beat the competition, and that a fortunate by-product of this objective is the maximisation of firm 1's profit: it is only because firm 1 is committed to be aggressive that its profit is maximised (otherwise firm 2 would anticipate that firm 1 is bluffing, and therefore that firm 1 will play *in fine* its best reply).

Bargaining situations offer a salient illustration of this difference between commitment and bluff. Consider for instance the recent negotiations between Greece and the European Union concerning Greek national debt. In a Op-Ed article in the *New York Times*, former Greek finance minister (and game theorist) Yanis Varoufakis (2015a) claims that 'it would be pure folly to think of the current deliberations between Greece and our partners as a bargaining game to be won or lost via bluffs and tactical subterfuge', because, unlike within standard game theory in which the motives of the players are taken for granted (maximising one's material payoff), 'the whole point [of the current deliberations between Greece's European partners and the new government] is *to forge new motives*' (my emphasis). The main motive of the Greek government is to implement its social policy agenda, to 'do what is right not as a strategy but simply because it is ... right'. Varoufakis emphasises this commitment not to cross this 'red line', by stating that 'we are *determined* to clash with mighty vested interests in order to reboot Greece and gain our partners' trust. We are also *determined* not to be treated as a debt colony that should suffer what it must' (my emphasis). He concludes by claiming that:

> 'One may think that this retreat from game theory is motivated by some radical-left agenda. Not so. The major influence here is Immanuel Kant, the German philosopher who taught us that the rational and the free escape the empire of expediency by doing what is right.'

My point here is not to discuss whether the policy defended by the Greek government is the right one or not, but to offer an analysis of Varoufakis's argument in terms of the model of preferences developed in this paper. The negotiation can be roughly described as a game between Athens ($A$) and Brussels ($B$): $A$ faces a debt crisis and can be helped by $B$. $B$ has the choice between two strategies: to give a financial aid to $A$ (strategy $a$) or not (strategy $na$). $A$ can also implement

austerity measures (strategy $m$) or not (strategy $nm$). $\succ_i$ denotes the material payoff relation of player $i$. We have:

- $A$: $\{nm; a\} \succ_A \{m; a\} \succ_A \{m; na\} \succ_A \{nm; na\}$

- $B$: $\{m; na\} \succ_B \{m; a\} \succ_B \{nm; a\} \succ_B \{nm; na\}$

A possible representation of those material payoffs in a matrix is the following:

|      | $a$   | $na$  |
|------|-------|-------|
| $m$  | (2;2) | (1;3) |
| $nm$ | (3;1) | (0;0) |

The game in material payoff is therefore a chicken game: both players are ready to make an effort to avoid $A$'s default, but prefer that the others make the effort[1]. We have therefore two Nash equilibria in pure strategies, $\{nm, a\}$ and $\{m, na\}$. An additional difficulty of this game is that, although $A$ really needs the aid from $B$ to solve its crisis, $A$ has been elected on the promise that it would not accept any additional austerity measures. A possible strategy for $A$ would be to *pretend* that the respect of its political promise matters more than the payment of the debt, and therefore that — unlike its material payoff that only represents the financial interest of $A$ — $A$'s preferences are such that $\{nm; na\} \succ_A \{m; na\}$, i.e. that $A$ will respect its promise even if it implies no aid from $B$, and therefore default. This kind of threat is however not credible, since playing $m$ is still $A$'s best reply to $na$.

However, in his article, Varoufakis is trying to convince $B$ that $A$ is truthful when claiming that what matters for $A$ is not its 'material payoff' (i.e. the payment of $A$'s debt), but another motive[2] (the respect of its promise). Indeed, if $B$ believes

---

[1]We make the simplifying assumption here that austerity measures or a financial aid *alone* are sufficient to solve the crisis.

[2]Note that $B$ was also probably trying to choose strategic preferences here: although $B$ wanted to avoid $A$'s default, several members of the Eurozone (e.g. Germany) made clear that any aid from $B$ should be conditioned on the implementation of austerity measures by $A$. In other words, $B$ states that it would play $a$ only if $A$ plays $m$, and therefore that $\{nm; na\} \succ_B \{nm; a\}$. With this threat, $B$ hopes that $A$ will eventually choose $m$, and then reach $B$'s preferred Nash equilibrium (although this argument is valid in our simplified model, the 'real' outcome would probably be $\{m; a\}$, since Brussels should still give a financial aid to save Athens from default).

that $A$ is simply bluffing, in the sense that $A$ pretends to want to implement its social policy under any circumstances (although $A$ would prefer to implement a austerity policy so as to get $B$'s aid), then $A$'s threat of implementing the social policy is not credible. On the other hand, if $A$ *truly* wants to implement the social policy (and even convince itself that this is the *right* thing to do), then $A$'s threat becomes credible, and $B$ is likely to offer its aid. Varoufakis indeed puts a strong emphasis on $A$'s *determination*, and even invokes Kant and his categorical imperative (i.e. that respecting this promise is an *unconditional* requirement). Varoufakis therefore tries to convince $B$ that the game should be represented as follows:

|     | a     | na    |
|-----|-------|-------|
| m   | (0;2) | (0;3) |
| nm  | (3;1) | (1;0) |

$nm$ (no austerity measure) is now a strictly dominant strategy for $A$, and $B$ best reply is to offer its aid. It seems here that Varoufakis is probably *not* bluffing: he genuinely decided to follow another objective than the initial one in terms of material payoff[3]. The fact that $A$ may *in fine* benefit from $B$'s aid without having to break its electoral promise is only a fortunate by-product of its new preferences. $A$'s commitment to respect its policy agenda is therefore not a bluff, it is a *rational* commitment. It is therefore possible that $A$, by showing its determination to do 'what is right', chose to 'forge new motives' as a means to satisfy its primary objective ($A$'s new preferences — described in the second matrix — are therefore well *strategic* preferences), but it is also possible that, in line with a more Kantian argument, $A$ deliberately chose its new preferences as an end in itself.

Throughout the rest of this paper, I will assume that strategic preferences are purely instrumental: *ex ante*, each individual has a well-identified objective (her material payoff), and chooses her optimal commitment so as to satisfy this objective. One's strategic preferences are therefore only valuable as a means to satisfy one's material payoff — in particular, I will not investigate whether one's strategic preferences may progressively become one's material payoff. So as to illustrate this point, think for instance of a philanthropist who helps the needy, because it is

---

[3]It should indeed be noticed that an austerity plan was eventually implemented, but only after Varoufakis resignation: while he was probably committed not to implement an austerity plan, his position became in minority within Tsipras government (see Varoufakis (2015b) interview in *New Statesman*).

her moral duty (and not merely by sympathy): she is then satisfying preferences that are different from her material payoff. She may then progressively find 'inner satisfaction in spreading joy' (Kant, 1785), implying that her strategic preferences progressively become her material payoff. She is then acting *in accordance with*, rather than *from*, her duty. Considering this kind of preference evolution would offer a basis to develop a more general model of preference formation, but this is however beyond the scope of the present work.

## 2.2 Preliminaries

Let $N = \{1, \ldots, n\}$ denote the set of players, with $n \geq 2$. $X = \prod_{i \in N} X_i$ denotes the set of pure strategy profiles where each set $X_i \subset \mathbb{R}$ denote the strategy space of player $i$. The material payoff of a player $i \in N$ is given by a function $\Pi_i : X \mapsto \mathbb{R}$, $\forall i \in N$. Assume that players may present interdependent preferences, i.e. that their utility function $U_i : X \mapsto \mathbb{R}$ — whose maximisation determines their choice — is a weighted sum of the material payoff functions $\Pi_j(x)$:

$$U_i(x|S) = \sum_{j \in N} \sigma_{ij} \Pi_j(x), \tag{1}$$

with $S = \{\sigma_{ij}\}_{i,j \in N} \in \mathbb{R}^{n \times n}$ a set of real parameters. $\sigma_{ij}$ therefore represents the weight player $i$ gives to player $j$ in her utility function, and its sign indicates whether player $i$ tries to cooperate or not with player $j$. $\Pi_i$ therefore measures the ultimate objective of player $i$ (her material payoff), while $U_i$ represents her strategic preferences, i.e. the optimal interdependent preferences player $i$ should choose so as to maximise *in fine* her material payoff $\Pi_i$.

For any game in normal form $\Gamma = \langle N, X, \Pi \rangle$, define a two-stage game $\Gamma^*$ as follows:

- in the second stage game $\Gamma_2(S) = \langle N, X, U(.|S) \rangle$, player $i \in N$ chooses a strategy $x_i \in X_i$ so as to maximise her utility function $U_i(x|S)$, $\forall i \in N$;

- in the first stage game $\Gamma_1 = \langle N, \mathbb{R}^{n \times n}, V \rangle$, player $i \in N$ chooses a vector of real parameters $S_i = \{\sigma_{i1}; \ldots; \sigma_{in}\}$ so as to maximise her indirect payoff function $V_i(S) = \Pi_i(\bar{x}(S))$, with $\bar{x}(S)$ a Nash equilibrium of $\Gamma_2(S)$, $\forall i \in N$.

For convenience, suppose that $\Pi_i$ is a $C^3$ function $\forall i \in N$. Furthermore:

**Assumption :** **$A0$**. $\forall S \in \mathbb{R}^{n \times n}$, $\Gamma(S)$ *has a unique Nash equilibrium in pure strategies* $\bar{x}(S)$, *i.e.* $\exists! \bar{x}(S) \in X$ *such that*, $\forall i \in N$:

$$\frac{\partial U_i}{\partial x_i}(\bar{x}(S)|S) = 0, \tag{2}$$

$$\frac{\partial^2 U_i}{\partial x_i^2}(\bar{x}(S)|S) < 0. \tag{3}$$

$A0$ is a very strong assumption, but it considerably alleviates the presentation of the main results, and appears to be not necessary. The two main results of this paper would indeed remain unchanged: the demonstration of proposition 1, according to which the players generally have an incentive in presenting strategic preferences different from their material payoff, could indeed easily be extended to a more general framework, while proposition 3, according to which the players choose cooperative (resp. aggressive) preferences in symmetric supermodular (submodular) games is proven under conditions that would ensure the existence of a unique Nash equilibrium in the second stage game (I indeed assume a strong form of diagonal dominance of the Jacobian matrix of marginal utilities).

I introduce the following notations:

- The partial derivatives of $\Pi_i : X \mapsto \mathbb{R}$ are denoted:

$$\Pi_i^{jk}(x) = \frac{\partial^2 \Pi_i}{\partial x_j \partial x_k}(x_1; \ldots; x_n). \tag{4}$$

- $J(S) \in \mathbb{R}^{n \times n}$ denotes the Jacobian matrix of the marginal utilities evaluated at the Nash equilibrium of $\Gamma_2(S)$:

$$J(S) = \{U_i^{ij}(\bar{x}(S))\}_{i,j \in N} \tag{5}$$

- For a $n \times n$ matrix $S \in \mathbb{R}^{n \times n}$, $S_{ij}$ denotes a $(n-1) \times (n-1)$ matrix that results from deleting row $i$ and column $j$ of $S$.

- For a $n \times n$ matrix $S \in \mathbb{R}^{n \times n}$, $C_{ij}^S = (-1)^{i+j} |S_{ij}|$ denotes the $(i; j)$ cofactor of $S$.

9

The notation for the derivatives also holds for the utility function $U_i$. The game $\Gamma(S)$ is supermodular (respectively submodular) if and only if, $\forall i \in N$:

$$U_i^{ij}(x) \geq (\leq)0 \qquad \forall x \in X, \quad \forall j \neq i. \tag{6}$$

I make the additional assumption that $J(S)$ and its minors $J_{ii}(S)$ are generically non singular $\forall S \in \mathbb{R}^{n \times n}$.

## 2.3   Subgame perfect equilibrium of commitment

Suppose that the players can choose their own weights $\sigma_{ij}$: they can therefore make strategic commitments, since their choice is determined by the maximisation of their utility function $U_i(x|S)$, while their payoff is determined by their material payoff $\Pi_i(x)$.

**Definition 2.1.** *Let* $\Gamma = \langle N, X, \Pi \rangle$ *denote a game in normal form. A strategy profile* $(\bar{x}; \bar{S}) \in X \times \mathbb{R}^{n \times n}$ *is a subgame perfect equilibrium of commitment (SPEC) of* $\Gamma$ *if and only if:*

- $\bar{x} \in X$ *is a Nash equilibrium of* $\Gamma_2(\bar{S})$,

- $\bar{S} \in \mathbb{R}^{n \times n}$ *is a Nash equilibrium of* $\Gamma_1$.

A SPEC is therefore a specific utility function (defined by the degree of interdependence with the other players) and a strategy profile of the initial game such that no player obtains a strictly higher material payoff by changing her strategic commitment $S_i$, i.e. there exists no game $\Gamma_2(S_i; \bar{S}_{-i})$ with $S_i \neq \bar{S}_i$ such that $i$ is better off at the Nash equilibrium of $\Gamma_2(S_i; \bar{S}_{-i})$ than at the Nash equilibrium of $\Gamma_2(\bar{S})$. A SPEC therefore characterises a strategy profile (and underlying strategic preferences) immune to individual deviations from the underlying strategic preferences.

## 2.4 Illustration

So as to illustrate this equilibrium notion, consider a game $\Gamma = \langle \{1, 2\}, \{\mathbb{R}^+\}^2, \Pi \rangle$, with:

$$\Pi_i(x_1, x_2) = ay + \frac{b}{2}y^2 - \frac{c}{2}x_i^2, \qquad a, c > 0, \; 4b < c, \tag{7}$$

with $y = (x_1 + x_2)$ if $b \geq 0$ and $y = \min\{(x_1 + x_2); |a/b|\}$ if $b < 0$ (this last condition ensures that the function $ay + \frac{b}{2}y^2$ is always increasing). $\Gamma$ is a public good game, in which each player chooses a level $x_i$ that generates a collective benefit and an individual cost.

We associate to $\Gamma$ a two-stage game $\Gamma^*$. In the second stage game $\Gamma_2(S)$, the players maximise their utility functions $U_i$:

$$U_i(x_1, x_2 | S) = \sigma_{i1} \Pi_1(x_1, x_2) + \sigma_{i2} \Pi_2(x_1, x_2), \tag{8}$$

$$\Leftrightarrow U_i(x_1, x_2 | S) = (\sigma_{i1} + \sigma_{i2}) \left[ ay + \frac{b}{2}y^2 \right] - \frac{c}{2} \left[ \sigma_{i1} x_1^2 + \sigma_{i2} x_2^2 \right]. \tag{9}$$

We can easily check that $\sigma_{ii} = 0$ cannot be a first stage equilibrium (if $b > 0$, player $i$ chooses $x_i \to +\infty$ and gets her worst possible payoff; if $b < 0$, player $i$ chooses $x_i = |a/b|$ and supports all the costs). This means that each player necessarily cares about her own payoff at the SPEC. We can now normalise $\sigma_{ii}$ to $1$, $\forall i \in N$. The unique Nash equilibrium of $\Gamma_2(S)$ is then:

$$\bar{x}_i(S) = \frac{a(1 + \sigma_{ij})}{c - (2 + \sigma_{12} + \sigma_{21})b}, \qquad \forall i \in N. \tag{10}$$

(10) gives the optimal effort of each player given the weights they attribute to the other player within their utility function (we can verify that $(\bar{x}_1 + \bar{x}_2) < |a/b|$ when $b < 0$).

Suppose now that both players are able to directly choose those weights in a first stage game $\Gamma_1 = \langle N, \mathbb{R}^2, V \rangle$, with $V_i = \Pi_i(\bar{x}(S))$ the indirect payoff function of player $i$:

$$V_i(S) = \frac{a^2(2 + \sigma_{12} + \sigma_{21})((1 - 2\sigma_{12} - \sigma_{12}^2)c - (2 + \sigma_{12} + \sigma_{21})b)}{2(c - (2 + \sigma_{12} + \sigma_{21})b)}. \tag{11}$$

We can then compute the Nash equilibrium of the first stage game:

$$\bar{\sigma}_{ij} = \frac{c - 2b - \sqrt{c(c - 4b)}}{2b}, \qquad \forall i \in N. \tag{12}$$

We therefore obtain a unique[4] SPEC $(\bar{x}; \bar{S}) \in \{\mathbb{R}^+\}^2 \times \mathbb{R}^2$:

$$\begin{cases} \bar{x}_i = \dfrac{2ab - c + \sqrt{c(c - 4b)}}{2b\sqrt{c(c - 4b)}}, & \forall i \in N, \\[3mm] \dfrac{\bar{\sigma}_{ij}}{\bar{\sigma}_{ii}} = \dfrac{c - 2b - \sqrt{c(c - 4b)}}{2b}, & \forall i \in N, \ j \neq i. \end{cases} \tag{13}$$

The interpretation of this equilibrium is the following: *if the players can keep their commitments* (e.g. thanks to contracts, or if their utility function is the result of an unconscious evolutionary process), then they should form interdependent preferences and choose their strategy so as to maximise $U_i(x) = \Pi_i(x) + \frac{\bar{\sigma}_{ij}}{\bar{\sigma}_{ii}}\Pi_j(x)$. An interesting follow-up question would be to determine whether the players tend to form cooperative (i.e. to choose $\sigma_{ij} > 0$) or aggressive ($\sigma_{ij} < 0$) preferences. We can easily check that:

$$sign(\bar{\sigma}_{ij}) = sign(b). \tag{14}$$

This means that the sign of $b$ (whether the benefit function is concave or convex, and therefore whether the game is submodular or supermodular) determines the nature of the strategic preferences of both players. We can also check that the equilibrium output $\bar{x}$ will be higher than the Nash equilibrium if and only if $b > 0$, i.e. when the game is supermodular, and that players form cooperative preferences — the cooperation is however not full, since $\bar{\sigma}_{ij} < 1$.

With a convex benefit function ($b > 0$), the game is supermodular, both players partially cooperate and reach a higher payoff than the Nash payoff. Conversely, for a game with a concave benefit function ($b < 0$), the players will be more competitive at the first stage equilibrium and will therefore get a lower outcome. Indeed, in presence of strategic substitutes, each player has an incentive to 'blackmail' the other one — i.e. to unilaterally decrease her own output — in order to force the other player to increase her output. Since both players have

---

[4]There is a continuum of first stage equilibria since $\sigma_{ii}$ has been normalised to 1, but the SPEc is unique once the weight each player attributes to herself has been normalised.

the same reasoning, they enter in a vicious circle and end up with a deteriorated situation.

This illustration highlights the possible connection between supermodularity and the endogenous formation of cooperative preferences. We can indeed find a similar result within the literature on the indirect evolutionary approach: Bester & Güth (1998) for instance argue on the one hand that altruism is evolutionary stable in some games presenting strategic complementarities, while Bolle (2000) and Possajennikov (2000) notice on the other hand that relaxing this assumption will lead to the evolutionary stability of spite and anti-social motives.

# 3    Optimal preferences

I firstly introduce the notions of Stackelberg best reply and payoff functions, and then determine the optimal weights $\bar{\sigma}_{ij}$.

## 3.1    Stackelberg best reply and payoff functions

Before presenting the main results, I need to introduce the notion of *Stackelberg best reply function* and *Stackelberg payoff function*. The Stackelberg best reply function of player $j$ is her best reply to the strategy of player $i$, knowing that the players $k \neq i, j$ are maximising their utility: it is the reply function a Stackelberg leader would use so as to predict the behaviour of her followers. The Stackelberg payoff function is simply the material payoff of player $i$ that integrates the Stackelberg best reply functions of the other players, and whose maximisation determines the strategy chosen by a Stackelberg leader with $(n-1)$ followers.

**Definition 3.1.** *Let $\Gamma_2(S) \backslash \hat{x}_i = \langle N \backslash i, X_{-i}, U_{-i}(.|x_i = \hat{x}_i) \rangle$ denote the game $\Gamma_2(S)$ when $i$'s strategy is fixed to $\hat{x}_i$. The function $f_j : X_i \times S \mapsto X_j$ is the Stackelberg best reply function of player $j$ for $S \in \mathbb{R}^{n \times n}$ if and only if $\{f_1(\hat{x}_i|S); \dots ; f_n(\hat{x}_i|S)\} \in X_{-i}$ is a Nash equilibrium of $\Gamma_2(S) \backslash \hat{x}_i$.*

The set of Stackelberg best reply functions for players $j \neq i$ corresponds therefore to their optimal choice (i.e. a Nash equilibrium) for a given strategy of $i$. Note that the existence of a Nash equilibrium in $\Gamma_2(S)$ implies that the best reply functions $f_j(.|S)$ are defined on a non empty subset of $X_i$. Indeed, if it was not the case, then a second stage equilibrium could not exist, since $\bar{x} \in X$ is a Nash equilibrium of $\Gamma_2(S)$ if and only if:

$$f_j(\bar{x}_i|S) = \bar{x}_j, \qquad \forall i, j \in N. \tag{15}$$

For the same reasons motivating the assumption that each second stage game has a unique Nash equilibrium, I assume that there always exists a unique function $f_j : X_i \times S \mapsto X_j$, $\forall i, j \in N$. The reasoning supporting proposition 1 can indeed easily be extended to a more general framework with several functions $f_j$ (their existence being ensured by the existence of a Nash equilibrium in mixed strategies for each game $\Gamma_2(S)$), and the conditions establishing the relation between supermodularity and the formation of cooperative preferences would typically imply the uniqueness of the Stackelberg best reply function.

I can now define the Stackelberg function:

**Definition 3.2.** *Let $f_j(x_i|S)$ denote the Stackelberg best reply function of player $j$ for $S$. The function $\Psi_i : X_i \times S \mapsto \mathbb{R}$ is the Stackelberg function of player $i$ if and only if:*

$$\Psi_i(x_i|S) = \Pi_i(f_1(x_i|S); \ldots; f_n(x_i|S)). \tag{16}$$

$\Psi_i(x_i|S)$ corresponds to the material payoff of player $i$ when she anticipates the best reply of the other players (given their utility functions $U(x|S)$). It is the function that a player would maximise if she had a first mover advantage.

As a preparation for the propositions, I show the following lemmas (the proofs are provided in appendix):

**Lemma 1.** *Let $f_j(x_i|S)$ be the Stackelberg best reply function of player $j$ for $S$. We have:*

$$\frac{\partial f_j}{\partial x_i}(x_i|S) = \frac{C_{ij}^{J(S)}}{C_{ii}^{J(S)}}, \tag{17}$$

*with $C_{ij}^{J(S)}$ the $(i;j)$ cofactor of $J(S)$, the Jacobian matrix of the marginal utility functions $U_i^i(x|S)$, evaluated at the Nash equilibrium of $\Gamma_2(S)$.*

**Lemma 2.** *If $\forall j, k \neq i$:*

    *(i)* $|U_i^{ii}(\bar{x}(S)|S)| > (n-1)\left|U_i^{ij}(\bar{x}(S)|S)\right|$,

    *(ii)* $\left|U_i^{ik}(\bar{x}(S)|S)\right| < (n-1)\left|U_i^{ij}(\bar{x}(S)|S)\right|$,

    *then:*

$$sign\left(\frac{\partial f_j}{\partial x_i}(x_i|S)\right) = sign\left(U_j^{ji}(\bar{x}(S)|S)\right), \qquad \forall j \neq i. \tag{18}$$

**Lemma 3.** *If $\forall j, k \neq i$:*

    *(i)* $|U_i^{ii}(\bar{x}(S)|S)| > (n-1)\left|U_i^{ij}(\bar{x}(S)|S)\right|$,

    *(ii)* $\left|U_i^{ik}(\bar{x}(S)|S)\right| < (n-1)\left|U_i^{ij}(\bar{x}(S)|S)\right|$,

    *then:*

$$\sum_{j \neq i}\left|C_{ij}^{J(S)}\right| < \left|C_{ii}^{J(S)}\right|, \tag{19}$$

$$\Longleftrightarrow \sum_{j \neq i}\left|\frac{\partial f_j}{\partial x_i}\right| < 1. \tag{20}$$

Lemma 1 gives the expression of the first order derivative of the Stackelberg best reply function $f_j(x_i)$ as a function of the second order derivatives $U_i^{ij}$. We can then show that the sign of $\frac{\partial f_j}{\partial x_i}$ is the same than $U_j^{ji}(\bar{x}(S)|S)$, when conditions (i) and (ii) are verified (lemma 2). Condition (i) means that $j$'s impact on $i$'s marginal utility is relatively low compared to $i$'s impact on her own marginal utility (this is a strong form of row diagonal dominance — instead of asking ), and (ii) the cross derivatives $U_i^{ij}$ and $U_i^{ik}$ are relatively 'close' in absolute value $\forall j, k \neq i$, i.e. there is no player $j$ with a significantly higher importance from $i$'s perspective. Those conditions are typically verified for public good games and two-player games with $|U_i^{ii}| > |U_i^{ij}|$. Lemma 3 states that, under the same conditions, the sum of the $(i; j)$ cofactors, $\forall j \neq i$, is lower than the principal minor of $J(S)$.

We now determine the expression of the weights $\bar{\sigma}_{ij}$ at the Nash equilibrium of $\Gamma_1$, and determine their sign: we will then be able to define a class of games in which players endogenously adopt cooperative preferences, or conversely try to maximise the difference between their payoff and the payoff of their opponents.

## 3.2   Optimal weights

Let $\Gamma = \langle N, X, \Pi \rangle$ be a game in normal form, and $\Gamma^*$ its associated two-stage game. We firstly determine the conditions under which it is rational for all the players to choose null weights $\bar{\sigma}_{ij}$, $i \neq j$, i.e. all the players prefer to maximise their material payoff rather than adopting interdependent preferences:

**Proposition 1.** *Let $\bar{x} \in X$ be the Nash equilibrium of $\Gamma$, and $I_n$ a matrix in $\mathbb{R}^{n \times n}$ such that $\sigma_{ij} \neq 0$ iff $i = j$. $(\bar{x}; I_n)$ is a SPEC of $\Gamma$ if and only if:*

  *(i)  either $\Psi_i^i = 0$, $\forall i \in N$,*

  *(ii)  or $\Pi_i^j(\bar{x}) = 0$, $\forall j \in N$, and $\forall i$ such that $\Psi_i^i \neq 0$*

Proposition 1 states that, unless the interests of the players are perfectly aligned or opposed (in the sense that maximising one's payoff implies also maximising or

minimising the payoffs of all the other players — as implied by condition (i)), or that no-one can benefit from a first mover advantage (condition (ii)), there is at least one player who will be better off by choosing a non null weight $\sigma_{ij}$. The intuition behind this result is the following: in a strategic interaction with payoff maximisers, the highest payoff I can achieve is my Stackelberg payoff, i.e. the payoff I would get if I was able to play before the others (suppose here for the sake of argument that a player with a first mover advantage can always obtain the Nash payoff — this means for instance that, in a zero-sum game, a Stackelberg leader could play in mixed strategies). If I have the opportunity to choose strategic preferences different from my material payoff, then I can manipulate the Nash equilibrium of the game such that the strategy that satisfies my strategic preferences actually satisfies my Stackelberg payoff, i.e. such that the Nash equilibrium with strategic preferences corresponds to the Stackelberg equilibrium with my initial material payoff.

We can now provide the expression of the optimal weights:

**Proposition 2.** $(\bar{x}; \bar{S})$ *is a SPEC of* $\Gamma$ *if,* $\forall i, j \in N$:

$$\bar{\sigma}_{ij} = \frac{\Pi_i^j}{\Pi_j^i}(\bar{x}) \frac{\partial f_j}{\partial x_i}(\bar{x}_i | \bar{S}). \tag{21}$$

Proposition 2 gives the expression of the optimal weights a player $i$ should give to the other players so as to maximise her material payoff (we can check that $i$ necessarily maximises her own payoff, since $\bar{\sigma}_{ii} = 1$). This condition is not necessary, since — as shown in the proof — the vector $\bar{S}_i$ is determined by a single equation. Although other specifications were possible, I chose here to define $\bar{\sigma}_{ij}$ as a function of the Stackelberg best reply of player $j$ when $i$ is the leader, since it captures the idea that the attitude of $i$ towards $j$ fundamentally depends on the way $j$ reacts when $i$ changes her strategy.

The weights $\bar{\sigma}_{ij}$ are therefore chosen such that satisfying my strategic preferences is formally equivalent to maximising my Stackelberg function. It can then be interesting to determine under which conditions the choice of one's preferences implies a more cooperative or competitive behaviour with the other players, i.e.

to determine the sign of the optimal weights $\bar{\sigma}_{ij}$. Thanks to lemmas 1 and 2, we can see that the sign of $\bar{\sigma}_{ij}$ is determined by the sign of $U_j^{ji}$ (they have the same sign if and only if $sign(\Pi_i^j(\bar{x}(\bar{S}))) = sign(\Pi_j^i(\bar{x}(\bar{S})))$, which is for instance the case in public good games or Cournot oligopoly). It means that player $i$ will cooperate with player $j$ if and only if there is a strategic complementarity between $i$ and $j$ in the game $\Gamma_2(\bar{S})$. This implies in particular that, in supermodular games, players have an interest in presenting cooperative preferences, since this will generate a positive best reply from the other players: cooperating is therefore beneficial because it gives an incentive to other players to reciprocate. On the contrary, games with strategic substitutes will generate more competitive behaviours, the players having an incentive in maximising the difference between the payoffs rather than their sum.

Note however that proposition 2, lemma 1 and lemma 2 are not sufficient to ensure that players will necessarily cooperate if the initial game $\Gamma$ is supermodular: the condition holds only for the resulting game $\Gamma_2(\bar{S})$. It is in fact not impossible that there exists a SPEC in a supermodular game such that all players present negative $\sigma_{ij}$ (it can be consistent if the resulting game is submodular): there can therefore exist Nash equilibria in the first stage game that create an artificial competition between the players, although the initial game was supermodular. We can however notice that the only reason for $i$ to compete with $j$ is that $j$ competes at equilibrium with $i$: it seems quite unlikely that players will effectively converge to such an equilibrium.

A corollary of those results is that, if only one player $i$ is able to make strategic commitments (as in a game with a Stackelberg leader), then the strategy chosen by this Stackelberg leader would correspond to the satisfaction of cooperative (resp. competitive) preferences in supermodular (submodular) games: Stackelberg leadership therefore leads to a greater cooperation in supermodular games, and a greater competition in submodular games.

I now show that the connection between the supermodularity of the initial game $\Gamma$ and positive $\bar{\sigma}_{ij}$ holds for symmetric games.

## 3.3   Symmetric games

I focus here on symmetric games to establish a direct connection between supermodularity and the choice of cooperative preferences at the equilibrium of $\Gamma_1$.

**Definition 3.3.** *A game in normal form* $\Gamma = \langle N, X, \Pi \rangle$ *is symmetric if and only if:*

- $X_i = X_j$, $\forall i, j \in N$,

- *for any permutation* $s : N \mapsto N$:

$$\Pi_i(x_1; \ldots; x_i; \ldots; x_n) = \Pi_{s(i)}(x_{s(1)}; \ldots; x_{s(i)}; \ldots; x_{s(n)}). \tag{22}$$

We have the following proposition (proof in appendix):

**Proposition 3.** *Let* $\Gamma$ *be a symmetric game. If* $\forall j, k \neq i$:

(i) $\left| U_i^{ii} \right| > (n-1) \left| U_i^{ij} \right|$,

(ii) $\left| U_i^{ik} \right| < (n-1) \left| U_i^{ij} \right|$,

(iii) $\left| \Pi_j^{ji} \right| \geq \left| \Pi_k^{ji} \right|$,

*then a symmetric SPEC* $(\bar{x}; \bar{S})$ *verifies:*

$$sign\left( \bar{\sigma}_{ij} \right) = sign\left( \Pi_j^{ji}(\bar{x}(S)) \right), \qquad \forall j \neq i. \tag{23}$$

Proposition 3 states that the connection between supermodularity of the initial game and cooperation in second stage game is true for symmetric $n$-player games, under the assumptions (i) and (ii) introduced in the previous section, and under the additional assumption that the second order derivative $\Pi_k^{ji}$ (for different $i$, $j$ and $k$) is relatively low in absolute value. This result means that, in a symmetric game, player $i$ will choose to put a positive weight on the material payoff of player $j$ if and only if $\Pi_j^{ji}$ is positive in the initial game: supermodular games will then endogenously generate cooperative preferences. The cooperation is not full, since lemma 3 implies that:

$$\sum_{j \neq i} |\bar{\sigma}_{ij}| < 1, \qquad \forall i \in N. \tag{24}$$

19

This means that player $i$ will never assign a higher weight to the set of the other players compared to her own material payoff in her preferences. On the contrary, games with strategic substitutes will exacerbate the competition between the players and lead to more aggressive behaviours.

# 4    Application to climate change negotiations

Within the present framework, we can state — paraphrasing Lucas — that given that individual preferences consist of optimal decision rules in the first stage game, and that optimal decision rules vary systematically with changes in the structure of the strategic environment, then any change in policy will systematically alter individual preferences. The design of public policies should therefore integrate the possibility that the players will adapt their preferences. Propositions 2 and 3 imply that, in games characterised by strategic substituability, such as public good games with a concave benefit function, the players have an incentive to become more aggressive: it is therefore possible that the total contribution progressively decreases with the emergence of more competitive preferences, leading *in fine* to a deteriorated situation (worse than the initial Nash equilibrium). An interesting policy recommendation in this kind of situation would be to change the incentives of the initial game such that the players are not tempted any more to adopt such preferences. A solution to promote cooperation would then be to transform the game into one presenting strategic complementarities: this should indeed endogenously lead to more cooperative behaviours.

The aim of this section is to illustrate this point by studying climate change negotiations: we can see that the different solutions suggested up to now consist in designing economic incentives to reduce greenhouse gas emissions (with for instance the Kyoto Protocol or the European Union Emissions Trading Scheme). Those approaches however keep the strategic substituability of the initial game of pollution abatement, since there is a perfect substituability between the emissions of two countries $i$ and $j$. This may in turn give an incentive to the countries to adopt more aggressive positions in international negotiations: they can indeed threaten the other countries to lower their contribution, so as to force them to provide a greater effort. We can indeed reasonably assume that the countries are able to make strategic commitments, since international negotiations are not only

a matter of economic interest, but also of political influence. Efficient international agreements should therefore build a system of incentives that give a coordination structure to the game of pollution abatement, by relying for instance on the adoption of technological standards and trade sanctions to punish the countries not respecting the agreement (this argument is in line with the recommendations of Barrett (2003, 2007) concerning the design of international agreements).

## 4.1   Model

Consider two identical countries $i \in N$, in which a firm produces and sells a consumption good in quantity $q_i$. There is no international trade, and the national firm takes the national price $p_i$ as given. The production of the final good generates a pollution $D(q_1 + q_2)$ that negatively affects both countries. Each country can tax the production of its firm (tax rate of $\tau_i$ per unit of production), and is therefore able to indirectly set its level of production. The countries are facing a public good game: they indeed choose their level of production (*via* their national regulation) so as to maximise their national payoff $\Pi_i$, knowing that this generates a national gain (in terms of surplus) but a global loss (pollution). This leads in turn to an over-production and a Pareto-dominated Nash equilibrium.

Suppose now that the countries want to implement an international agreement in order to maximise the global payoff $\Pi = \Pi_1 + \Pi_2$, knowing that — once the agreement is signed — both countries will choose their level of production so as to maximise $\Pi_i$. In this model, firms are in perfect competition, the two countries play a game, and they try to reach an agreement from the social planner's perspective. The objective of this illustration is to compare several alternatives of international systems and to argue in favour of systems creating a game of coordination between the countries rather than keeping the strategic substituability of the initial game. For convenience, I assume that the countries can implement an international tax system such that no fraud is possible, and the funds are collected by an international fund (we can assume for instance that those funds are then used to indemnify the victims of the pollution). This model offers a very simple picture of the current negotiations on climate change: the national productions generate the emission of greenhouse gases, and the different countries try to establish an international system so as to reduce the environmental damage of their production. The international fund in the model can be assimilated to the Green Climate Fund, which is funded by the developed countries emitting a higher

quantity of greenhouse gases. I suggest now comparing two main scenarios:

- International carbon tax (ICT): each country must pay a constant tax $t_{ICT}$ per unit of pollution emitted

- Trade sanction (TS): the country with the less demanding regulation must pay a tax $t_{TS}(\tau_1; \tau_2)$ to the other country per unit of pollution emitted; this tax depends on the difference between national taxes

ICT seems to correspond to the ideal solution from an economic perspective: the tax internalises the negative externality of the production, and can be defined so as to reach a Pareto optimal outcome. In the second scenario, although there is no international trade, the situation can be related to a mechanism of trade sanction: if a country is in a situation of environmental dumping, then its partners may impose additional taxes on the goods exported by this country (such that no firm can be eventually advantaged by a less restricting regulation). The situation here is relatively similar, since the country directly pays to the other an additional tax in case of environmental dumping. Formally, the material payoff of country $i$ is the following:

$$\Pi_i(q|ICT) = CS_i(q_i) + \pi_i(q_i) + \tau_i q_i - D(q_1 + q_2) - t_{ICT} q_i, \qquad (25)$$

$$\Pi_i(q|TS) = CS_i(q_i) + \pi_i(q_i) + \tau_i q_i - D(q_1 + q_2) - t_{TS}(\tau_j - \tau_i) q_i, \qquad (26)$$

with $CS_i$ the consumer surplus and $\pi_i$ the profit of the firm from country $i$. Assume a linear demand, convex costs for the firms, and a convex damage function:

$$p_i = a - b q_i, \qquad (27)$$

$$\pi_i = (p_i - \tau_i) q_i - \frac{c}{2} q_i^2, \qquad (28)$$

$$D(q) = \frac{\delta}{2} (q_1 + q_2)^2. \qquad (29)$$

Since the firms are in perfect competition on each national market, we can easily compute the consumer surplus $CS_i$, as well as the production $q_i$ as a function of $\tau_i$:

$$CS_i(q_i) = \frac{b}{2}q_i^2, \tag{30}$$

$$CS_i(q_i) + \pi_i(q_i) + \tau_i q_i = aq_i - \frac{b+c}{2}q_i^2, \tag{31}$$

$$\text{with} \qquad q_i = \frac{a - \tau_i}{b + c}. \tag{32}$$

Without international agreement, the material payoff of each country is therefore:

$$\Pi_i(q) = aq_i - \frac{b+c}{2}q_i^2 - \frac{\delta}{2}(q_1 + q_2)^2, \tag{33}$$

and the Nash equilibrium, $\forall i \in N$:

$$\begin{cases} \bar{q}_i = \dfrac{a}{b + c + 2\delta}, \\ \bar{\tau}_i = \dfrac{2a\delta}{b + c + 2\delta}. \end{cases} \tag{34}$$

Both countries are therefore producing too much, since the social optimum (maximising the sum of material payoff) is, $\forall i \in N$:

$$\begin{cases} \tilde{q}_i = \dfrac{a}{b + c + 4\delta}, \\ \tilde{\tau}_i = \dfrac{4a\delta}{b + c + 4\delta}. \end{cases} \tag{35}$$

We now compare the two alternatives in a 'naive' scenario, i.e. if the countries do not anticipate (or simply cannot make) strategic commitments. We will then study the case of a 'sophisticated' scenario in which countries are able to make strategic commitments (the appropriate equilibrium solution concept would then be a SPEC and not the Nash equilibrium characterised above)

## 4.2   Naive scenario

Consider firstly that both countries implement ICT: they therefore pay to a third party a tax $t_{ICT}$ per unit of pollution. Their material payoff is now:

$$\Pi_i(q) = aq_i - \frac{b+c}{2}q_i^2 - \frac{\delta}{2}(q_1 + q_2)^2 - t_{ICT}q_i. \qquad (36)$$

Suppose that the countries agree on a naive tax, i.e. a tax that does not take into account the possibility that players may make strategic commitments, once the agreement is signed. In the absence of strategic commitment, the level of the tax that allows the countries to reach the social optimum (35) is:

$$t_{ICT,n} = \frac{2a\delta}{b+c+4\delta}. \qquad (37)$$

If the countries agree to implement this international tax, then they have the adequate incentives to reach the optimal production $\tilde{q}$, given their initial preferences.

Consider now the scenario TS. Assume that $\tau_i < \tau_j$: the country $i$ (with the lowest national tax) must pay a tax $t_{TS}(\tau_j - \tau_i)$ per unit of pollution, and country $j$ collects this tax within the limits of its own production. The residual is collected by the international organisation[5]. Within this scenario, the material payoff of country $i$ is, $\forall i \in N$:

$$\Pi_i(q) = aq_i - \frac{b+c}{2}q_i^2 - \frac{\delta}{2}(q_1 + q_2)^2 - t_{TS}(\tau_j - \tau_i)q_i, \qquad (38)$$

$$\Longleftrightarrow \quad \Pi_i(q) = aq_i - \frac{b+c}{2}q_i^2 - \frac{\delta}{2}(q_1 + q_2)^2 - t_{TS}(b+c)(q_i - q_j)q_i. \qquad (39)$$

As previously, consider that the countries agree on a naive tax rate, that does not take into account the possibility for the countries to make strategic commitments *ex post*. The expected Nash equilibrium if scenario TS is chosen is then:

$$q_i = \frac{a}{(b+c)(1+t_{TS}) + 2\delta}, \qquad \forall i \in N, \qquad (40)$$

$$\tau_i = \frac{a(2\delta + (b+c)t_{TS})}{2\delta + (b+c)(1+t_{TS})}, \qquad \forall i \in N. \qquad (41)$$

---

[5]I adopt this specific framework to be consistent with a scenario of trade sanctions, in particular if one country is bigger than the other: the funds collected by $j$ must indeed be in a scale consistent with its own production.

So as to reach the social optimum (35), we should therefore implement the following tax rate:

$$t_{TS} = \frac{2\delta}{b + c}.$$ (42)

In a naive scenario, both taxation systems are strictly equivalent, since the adequate definition of the taxation levels allow the countries to reach the social optimum.

## 4.3 Sophisticated scenario

Suppose now that both countries can make strategic commitments (similarly to the Greek case discussed previously, a possible way to implement this kind of commitment would be to form electoral promises, knowing that not respecting them may induce a sanction from the voters during the forthcoming elections). The game faced by the countries is a submodular game (public good game with a concave benefit function). This means that the players are likely to form aggressive preferences and choose negative weights $\sigma_{ij}$. The unique SPEC of the game is

$$\begin{cases} \bar{\sigma}_{ij} = \dfrac{\sqrt{(b + c)(b + c + 4\delta)} - b - c - 2\delta}{2\delta}, \\ \bar{q}_i = \dfrac{a}{\sqrt{(b + c)(b + c + 4\delta)}}. \end{cases}$$ (43)

Both countries try to get the upper hand on the other, and end up *in fine* with a deteriorated situation, in which they both produce more than at Nash equilibrium.

Consider now that both countries agree on ICT. We can now notice that this policy keeps the submodularity of the initial game: a player can therefore benefit from a strategic commitment such that she eventually produces $\hat{q} > \tilde{q}$, since the best reply of the other country will be to increase her effort (by reducing her production). In particular, since the implementation of a tax per unit of production does not affect the cross derivatives $\Pi_i^{ij}$, the players will choose the exact same weights as previously, and the tax (37) will not give the adequate incentives to reach the social optimum.

Suppose therefore that the countries anticipate that they will make strategic commitments once the agreement is signed. The unique SPEC is then:

$$
\begin{cases}
\bar{\sigma}_{ij} = \dfrac{\sqrt{(b+c)(b+c+4\delta)} - b - c - 2\delta}{2\delta}, \\
\bar{q}_i = \dfrac{(a - t_{ICT})(b+c+\delta)}{(b+c)^2 + 2\delta(b+c)}.
\end{cases}
\tag{44}
$$

The optimal tax rate is therefore:

$$
t^*_{ICT} = \frac{2a\delta(1 - \bar{\sigma})}{b + c + 4\delta},
\tag{45}
$$

$$
\Leftrightarrow \quad t^*_{ICT} = a\left[1 - \sqrt{\frac{b+c}{b+c+4\delta}}\right],
\tag{46}
$$

which is strictly higher than the naive tax $t_{ICT,n}$. An international agreement must therefore take into account the possibility *ex post* for the countries to benefit from strategic commitments. The question that arises then is to know whether the countries are likely to make the optimal strategic commitment: a crucial condition for choosing an optimal commitment is indeed that we anticipate that the others know that we will respect our commitment. If we do not believe that the other is sufficiently rational to play the first stage game, or alternatively that preferences are not directly chosen, but are the result of evolutionary pressures, then it is not certain that $t^*_{ICT}$ will be well-suited. In addition of preventing the players from adopting competitive strategic commitments, we should also implement a tax system such that the optimal policy does not depend on the level of $\sigma_{ij}$, on which the social planner has no direct control.

Although an international tax may lead *in fine* to the social optimum, we can notice that the players necessarily adopt aggressive preferences at equilibrium — which may be an undesirable property (in a non-welfarist perspective) of international relations. Furthermore, such a system is highly sensitive to the propensity of the players to respect their optimal commitment: if one player does not keep her optimal commitment, either because she is not rational enough, or because she does not think the other will keep her commitment, or because preferences evolve over time *via* an evolutionary dynamics, then the tax will probably not be well adapted.

Consider finally the adoption of TS when players are able able to keep their commitments. We can notice that the game with the payoff functions described in (39) is supermodular if and only if $t_{TS} > \frac{\delta}{b+c}$. This is for instance the case for the naive tax rate (42): the game with trade sanctions in the naive scenario is therefore supermodular, and should lead to the formation of cooperative preferences. The weights at the SPEC are indeed, $\forall i \in N$:

$$\bar{\sigma}_{ij} = \frac{\sqrt{(b+c)^2(t_{TS}+1)^2 + 4\delta(b+c)(t+1)} - (b+c)(t_{TS}+1) - 2\delta}{2\delta}, \qquad (47)$$

which are well positive if and only if $t_{TS} > \frac{\delta}{b+c}$. We saw that the naive tax in the case of scenario ICT was not adapted when players were making a strategic commitment (the optimal level of tax indeed directly depended on the level of $\sigma_{ij}$). A crucial difference between scenarios ICT and TS is that the naive tax implemented with TS still remains optimal at the SPEC. We can indeed show that the level of tax $t_{TS}$ that maximises social welfare does not depend on $\sigma_{ij}$ when $\sigma_{12} = \sigma_{21} = \sigma$ (this condition is verified in the present game, since the game is symmetric and has a unique SPEC). We have indeed in this situation the optimal production for $i$:

$$q_i = \frac{a}{(b+c)(1+t_{TS}(1-\sigma)) + 2\delta(1+\sigma)}, \qquad (48)$$

which is equal to the social optimum $\tilde{q}$ if and only if:

$$t_{TS} = \frac{2\delta}{b+c}. \qquad (49)$$

The symmetry of the game can help the countries to tackle the two issues faced in scenario ICT, (i) that countries were likely to adopt aggressive preferences, and (ii) that, from a more practical perspective, the level of the tax depended on the likelihood for both countries to keep their optimal commitment. Firstly, adopting a mechanism of trade sanctions is likely to generate cooperative preferences between countries, since their interests are now aligned: if a country increases its level of effort (by increasing $\tau_i$, and therefore diminishing $q_i$) then the cost is shared between the countries. The country with the lower effort must indeed now pay a compensation to the other country, and has now a new

27

incentive to increase its own effort: in addition to diminishing the environmental damage, the country will also stop paying the other country. Secondly, as long as the preferences of the country evolve in a similar way (and therefore $\sigma_{12} = \sigma_{12}$), the tax rate $t_{TS}$ given by (42) remains optimal over time. It is therefore possible to implement a naive tax rate, since this level of taxation will still be optimal *ex post*, once the countries have made symmetric commitments.

International agreements based on a mechanism of trade sanctions rather than an international tax are more likely to succeed, since they align the interests of the different countries: since the game is supermodular (both countries indeed know that increasing one's effort will increase the incentives for the other to increase its own effort), cooperative behaviours are likely to emerge, unlike with scenario ICT, in which the initial submodularity of the game is preserved, leading to the emergence of aggressive behaviours. We can indeed notice that the weights $\bar{\sigma}_{ij}$ chosen at the SPEC are increasing with $t_{TS}$, and that $\lim_{t_{TS} \to +\infty} \bar{\sigma}_{ij} = 1$: this means that by increasing the level of sanction, the players will naturally converge to cooperative preferences and directly maximise the global welfare.
A last interesting property of scenario TS is that, at a symmetric equilibrium, there is no transfer between the countries or with the international fund (both countries have indeed the same level of production). While the collect of the tax with scenario ICT is likely to generate additional costs, it is possible to achieve the same results in terms of individual incentives without having to make any transfer between countries.

# 5    Conclusion

It is generally implicitly assumed that, so as to get the highest level of payoff, players should choose the strategy that maximises their payoff: however, payoff-maximising behaviours are generally indirectly self-defeating, suggesting that rational players should be able to form *strategic* preferences, such that the satisfaction of those preferences leads *in fine* to an equilibrium with a higher material payoff (see (Parfit, 1984, 17-19) for a detailed argument supporting the idea of rational irrationality). The analysis developed in this paper provides a formal framework to study the choice of such strategic preferences. I identified the optimal weights each player should assign to the others in her utility function, and then the conditions

under which this process could lead to the formation of cooperative or competitive preferences. I highlighted a strong connection between supermodularity and the emergence of cooperative preferences, and then discussed the implications in terms of policy design. I argued that the efficient design of public policies should take into account the possible change in individual preferences induced by the policy, and suggested that public policies should privilege incentives that create a coordination game between the players and cut the possible submodularity of the initial game (as in climate change negotiations): this type of approach may indeed facilitate the formation of cooperative preferences, and hence facilitate the achievement of the policy objective.

# References

Barrett, S. (2003). "Environment and Statecraft: The Strategy of Environmental Treaty-Making". *American Economic Review*, 90, 166–193.

Barrett, S. (2007). *Why Cooperate? The Incentive to Supply Global Public Goods.* Oxford University Press.

Bester, H. & Güth, W. (1998). "Is Altruism Evolutionary Stable?". *Journal of Economic Behavior and Organisation*, 34, 211–221.

Bolle, F. (2000). "Is Altruism Evolutionary Stable? And Envy and Malevolence? - Remarks on Bester and Güth". *Journal of Economic Behavior and Organisation*, 42, 131–133.

Cardenas, J., Stranlund, J., & Willis, C. (2000). "Local Environmental Control and Institutional Crowding-out". *World Development*, 28(10), 1719–1733.

Fershtman, C. & Gneezy, U. (2001). "Strategic Delegation: an Experiment". *RAND Journal of Economics*, 32(2), 352–368.

Fershtman, C. & Kalai, E. (1997). "Unobserved Delegation". *International Economic Review*, 38(4), 763–774.

Fischer, S., Güth, W., Müller, W., & Stiehler, A. (2006). From ultimatum to nash bargaining: Theory and experimental evidence. *Experimental Economics*, 9(1), 17–33.

Frey, B. & Oberholzer-Gee, F. (1997). "The Cost of Price Incentives : An Empirical Analysis of Motivation Crowding-Out". *American Economic Review*, 87(4), 746–755.

Güth, W. & Yaari, M. (1992). "Explaining Reciprocal Behavior in Simple Strategic Games: an Evolutionary Approach". In U. Witt (Ed.), *Explaining Forces and Changes: Approaches to Evolutionary Economics*. University of Michigan Press.

Heifetz, A., Shannon, C., & Spiegel, Y. (2007a). "The Dynamic Evolution of Preferences". *Economic Theory*, 32, 251–286.

Heifetz, A., Shannon, C., & Spiegel, Y. (2007b). "What to Maximize if You Must". *Journal of Economic Theory*, 133, 31–57.

Kant, I. (1785). *Groundwork of the Metaphysics of Morals*. Cambridge: Cambridge Univsity Press. Translated by Mary Gregor (1997).

Lucas, R. (1976). Econometric policy evaluation: A critique. *Carnegie-Rochester conference series on public policy*.

Ostman, A. (1998). "External Control May Destroy the Commons". *Rationality and Society*, 10(1), 103–122.

Parfit, D. (1984). *Reasons and Persons*. Oxford: Oxford University Press.

Possajennikov, A. (2000). "On the Evolutionary Stability of Altruistic and Spiteful Preferences". *Journal of Economic Behavior and Organization*, 42, 125–129.

Poulsen, A. U. & Roos, M. W. (2012). "Do People Make Strategic Commitments? Experimental Evidence on Strategic Information Avoidance". *Experimental Economics*, 13, 206–225.

Samuelson, L. (2001). "Introduction to the Evolution of Preferences". *Journal of Economic Theory*, 97(2), 225–230.

Schelling, T. (1960). *The Strategy of Conflict*. Cambridge, Mass.: Harvard University Press.

Sengul, M., Gimeno, J., & Dial, J. (2012). "Strategic Delegation: A Review, Theoretical Integration, and Research Agenda". *Journal of Management*, 38(1), 375–414.

Titmuss, R. (1970). *The Gift Relationship: From Human Blood to Social Policy.* London: Allen and Unwin.

Varoufakis, Y. (2015a). "No Time for Games in Europe". *The New York Times.* `http://www.nytimes.com/2015/02/17/opinion/yanis-varoufakis-no-time-for-games-in-europe.html?_r=0`, accessed 08/04/2015.

Varoufakis, Y. (2015b). "Our Battle to Save Greece". *NewStatesman.* http://www.newstatesman.com/world-affairs/2015/07/yanis-varoufakis-full-transcript-our-battle-save-greece , accessed 29/09/15.

von Stackelberg, H. (1934). *Marktform und Gleichgewicht.* Vienna, Berlin: Springer.

# A   Lemma 1

We show that:

$$\frac{\partial f_j}{\partial x_i}(x_i|S) = \frac{C_{ij}^{J(S)}}{C_{ii}^{J(S)}}, \tag{50}$$

with $f_j(x_i|S)$ the Stackelberg best reply function of player $j$ for $S$, $C_{ij}^{J(S)}$ the $(i;j)$ cofactor of $J(S)$, the Jacobian matrix of the marginal utility functions $U_i^i(x|S)$, evaluated at the Nash equilibrium of $\Gamma_2(S)$.

Consider that all players but $i$ are maximizing their utility functions, i.e. that they play their best reply strategy according to $x_i$; if player $i$ changes her strategy such that $dx_i \neq 0$, then, we must verify, $\forall j \neq i$ (the different functions are evaluated in $(f_1(x_i); \ldots; f_n(x_i))$, i.e. when all players but $i$ maximize their utility functions):

$$dU_j^j(x) = 0, \tag{51}$$

$$U_j^{ji} \ dx_i + \sum_{k \neq i} U_j^{jk} \ dx_k = 0. \tag{52}$$

We can rewrite this system of linear equations with $\mathrm{d}x_{-i} = {}^t\{\mathrm{d}x_k\}_{k \neq i}$, and $B_i = {}^t\{u_k^{ki}\,\mathrm{d}x_i\}_{k \neq i}$:

$$J_{ii}\,\mathrm{d}x_{-i} + B_i = 0. \tag{53}$$

Since we assumed that $J_{ii}$ is non singular, the system (53) has a unique solution, with $J_{ii}^j$ a $(n-1) \times (n-1)$ matrix identical to $J$ except for the column made of $U_k^{kj}$, $\forall k \neq i$ which is replaced by $-B_i$, and without row $i$ and column $i$:

$$\mathrm{d}x_j = \frac{\left|J_{ii}^j\right|}{|J_{ii}|}. \tag{54}$$

We can develop the determinant of $J_{ii}^j$ (we suppose that $i < j$ without loss of generality) and add a row and a column at the $i$th place as follows:

$$\left|J_{ii}^j\right| = \begin{vmatrix} U_1^{11} & \ldots & U_1^{1,i-1} & 0 & U_1^{1,i+1} & \ldots & U_1^{1,j-1} & -U_1^{1,i}\,\mathrm{d}x_i & U_k^{k,j+1} & \ldots & U_1^{1n} \\ \ldots & \ldots & \ldots & 0 & \ldots & \ldots & \ldots & \ldots & \ldots & \ldots & \ldots \\ U_{i-1}^{i-1,1} & \ldots & U_{i-1}^{i-1,i-1} & 0 & U_{i-1}^{i-1,i+1} & \ldots & U_{i-1}^{i-1,j-1} & -U_{i-1}^{i-1,i}\,\mathrm{d}x_i & U_{i-1}^{i-1,j+1} & \ldots & U_{i-1}^{i-1,n} \\ 0 & \ldots & 0 & 1 & 0 & \ldots & 0 & 0 & 0 & \ldots & 0 \\ U_1^{i+1,1} & \ldots & U_{i+1}^{i+1,i-1} & 0 & U_{i+1}^{i+1,i+1} & \ldots & U_{i+1}^{i+1,j-1} & -U_{i+1}^{i+1,i}\,\mathrm{d}x_i & U_{i+1}^{i+1,j+1} & \ldots & U_{i+1}^{i+1,n} \\ \ldots & \ldots & \ldots & 0 & \ldots & \ldots & \ldots & \ldots & \ldots & \ldots & \ldots \\ U_n^{n1} & \ldots & U_n^{n,i-1} & 0 & U_n^{n,i+1} & \ldots & U_n^{n,j-1} & -U_n^{ni}\,\mathrm{d}x_i & U_n^{n,j+1} & \ldots & U_n^{nn} \end{vmatrix} \tag{55}$$

We can then invert the $i^{th}$ with the $j^{th}$ column, and we obtain:

$$\left|J_{ii}^j\right| = (-1)^{i+j}\,|J_{ij}|\,\mathrm{d}x^i. \tag{56}$$

We can now rewrite the relation (54):

$$\mathrm{d}x^j = \frac{C_{ij}^J}{C_{ii}^J}(f^1(x^i); \ldots; f^n(x^i))\,\mathrm{d}x^i. \tag{57}$$

This last relation gives us the best reply of player $j$ to a given variation of strategy of player $i$ in order to maximize her utility function when all the other players but $i$ are maximizing their utility functions. We can notice that the primitive of

the best reply in terms of variation $dx_j$ is the Stackelberg best reply function of player $j$, i.e. the strategy $x_j$ which maximizes the utility function $U_j$ for a given strategy of player $i$, knowing the best reply of the other players $k \neq i, j$. We have therefore:

$$\frac{\partial f_j}{\partial x_i}(x_i|S) = \frac{C_{ij}^J}{C_{ii}^J}(f_1(x_i); \ldots; f_n(x_i)). \tag{58}$$

# B   Lemma 2

We now prove that $\frac{\partial f_j}{\partial x_i}$ has the same sign than $U_j^{ji}$ when the two following conditions are verified:

$$(n-1)\left|U_i^{ij}\right| < \left|U_i^{ii}\right|, \qquad \forall j \neq i, \tag{59}$$

$$(n-1)\left|U_i^{ij}\right| > \left|U_i^{ik}\right|, \qquad \forall j, k \neq i. \tag{60}$$

Condition (59) implies a strong form of diagonal dominance for $J$, since it means that the diagonal terms are all significantly greater than all the off-diagonal terms (this condition is identical to diagonal dominance if the off-diagonal terms are identical). Condition (60) implies a similar condition, i.e. that the off-diagonal terms are relatively close. In both situations, those conditions mean that there is no player $j$ who has a significantly higher importance for $i$ compared to the other players $k$. We can notice that the condition (60) has no sense when $n = 2$, since there does not exist a $k$ different from $i$ and $j$. The condition (59) is then sufficient.

We now determine the signs of $C_{ii}^J$ and $C_{ij}^J$. Since $J$ is diagonal dominant, we know that $J_{ii}$ is also diagonal dominant. We also know that, at the Nash equilibrium of $\Gamma_2(S)$, we have $U_i^{ii}(\bar{x}(S)) < 0 \; \forall i \in N$. We have therefore:

$$sign(C_{ii}^J) = (-1)^{n-1}. \tag{61}$$

We need now to determine the sign of $C_{ij}^J$. For clarity, we will illustrate our demonstration by focusing on the case $n = 4$. We have the following matrix $J$:

$$J = \begin{pmatrix} U_1^{11} & \dots & U_1^{14} \\ \dots & & \dots \\ U_4^{41} & \dots & U_4^{44} \end{pmatrix} \tag{62}$$

Without loss of generality, suppose that $i < j$. We have:

$$C_{ij}^J = (-1)^{(i+j)} |J_{ij}| \tag{63}$$

We now invert lines and columns in $|J_{ij}|$ such that the term $U_j^{ji}$ stand in the first row and first column. This required $(i + j - 3)$ operations: we indeed need $(j - 1)$ operations to reach the first column and $(i - 2)$ operations to reach the first line (we indeed deleted the $i$th row to obtain $J_{ij}$, and we assumed that $i < j$). $J'_{ij}$ denote the matrix that results from those operations. We obtain:

$$C_{ij}^J = (-1)^{(i+j)} (-1)^{(i+j-3)} \left| J'_{ij} \right|, \tag{64}$$
$$C_{ij}^J = - \left| J'_{ij} \right|. \tag{65}$$

An interesting property of $J'_{ij}$ is that its first principal minor is necessarily composed by the second order derivatives $U_k^{kk}$, $k \neq i, j$. In the case of $n = 4$, we have for instance:

$$J_{24} = \begin{pmatrix} U_1^{11} & U_1^{12} & U_1^{13} \\ U_3^{31} & U_3^{32} & U_3^{33} \\ U_4^{41} & U_4^{42} & U_4^{43} \end{pmatrix} \tag{66}$$

$$J'_{24} = \begin{pmatrix} U_4^{42} & U_4^{41} & U_4^{43} \\ U_1^{12} & U_1^{11} & U_1^{13} \\ U_3^{32} & U_3^{31} & U_3^{33} \end{pmatrix} \tag{67}$$

It implies that the first principal minor is row-diagonal dominant. We now operate on the columns of $J'_{ij}$ in order to have $(U_j^{ji}; 0; \dots; 0)$ as a first row. We obtain:

$$|J'_{24}| = \begin{vmatrix} U_4^{42} & 0 & 0 \\ U_1^{12} & U_1^{11} - \frac{U_4^{41}}{U_4^{42}}U_1^{12} & U_1^{13} - \frac{U_4^{43}}{U_4^{42}}U_1^{12} \\ U_3^{32} & U_3^{31} - \frac{U_4^{41}}{U_4^{42}}U_3^{32} & U_3^{33} - \frac{U_4^{43}}{U_4^{42}}U_3^{32} \end{vmatrix} \tag{68}$$

We can now check that, under the assumptions (59) and (60), the first principal minor is still diagonal dominant, and the diagonal terms are still negative. We have then:

$$sign(|J'_{ij}|) = sign\left(U_j^{ji}(-1)^{n-2}\right). \tag{69}$$

We obtain:

$$sign(C_{ij}^J) = -sign(|J'_{ij}|), \tag{70}$$
$$sign(C_{ij}^J) = (-1)^{n-1}sign(U_j^{ji}). \tag{71}$$

We can now complete our proof and determine the sign of $\frac{\partial f_j}{\partial x_i} = \frac{C_{ij}^J}{C_{ii}^J}$:

$$sign\left(\frac{\partial f_j}{\partial x_i}\right) = \frac{sign\left(C_{ij}^J\right)}{sign\left(C_{ii}^J\right)}, \tag{72}$$
$$sign\left(\frac{\partial f_j}{\partial x_i}\right) = sign(U_j^{ji}). \tag{73}$$

# C   Lemma 3

We now show that, under the same assumptions than lemma 2, we have

$$\sum_{j\neq i}\left|C_{ij}^{J(S)}\right| < \left|C_{ii}^{J(S)}\right|, \tag{74}$$
$$\Longleftrightarrow \sum_{j\neq i}\left|\frac{\partial f_j}{\partial x_i}\right| < 1. \tag{75}$$

We know that:

$$C_{ii}^J = \sum_{k \neq i} U_k^{kj} C_{kj}^{J_{ii}}, \qquad \forall j \neq i, \tag{76}$$

$$C_{ii}^J = \sum_{j \neq i} \left[ \frac{1}{n-1} \sum_{k \neq i} U_k^{kj} C_{kj}^{J_{ii}} \right], \tag{77}$$

$$C_{ij}^J = -\sum_{k \neq i} U_k^{ki} C_{kj}^{J_{ii}} \tag{78}$$

Without loss of generality, suppose that $C_{ii}^J$ is positive (which is true if $n$ is odd). We therefore have:

$$\left| C_{ii}^J \right| = C_{ii}^J, \tag{79}$$

$$\left| C_{ij}^J \right| = \sum_{k \neq i} U_k^{ki} C_{kj}^{J_{ii}}. \tag{80}$$

We therefore obtain:

$$\left| C_{ii}^J \right| - \sum_{j \neq i} \left| C_{ij}^J \right| = \sum_{j \neq i} \sum_{k \neq i} \left( \frac{U_k^{kj}}{n-1} - U_k^{ki} \right) C_{kj}^{J_{ii}}. \tag{81}$$

(74) is true if and only if:

$$\sum_{j \neq i} \sum_{k \neq i} \left( \frac{U_k^{kj}}{n-1} - U_k^{ki} \right) C_{kj}^{J_{ii}} > 0. \tag{82}$$

If we multiply by $(n-1)$ and divide on both sides by $C_{jj}^{J_{ii}}$ (negative by construction, since we assumed $C_{jj}^J > 0$), we obtain:

$$\sum_{j \neq i} \left[ U_j^{jj} - (n-1)U_j^{ji} + \sum_{k \neq i,j} (U_k^{kj} - (n-1)U_k^{ki}) \frac{C_{kj}^{J_{ii}}}{C_{jj}^{J_{ii}}} \right] < 0. \tag{83}$$

We can now check that under conditions (i) and (ii), this condition is verified (the second term is well negative, since $\frac{C_{kj}^{J_{ii}}}{C_{jj}^{J_{ii}}}$ has the same sign than $U_j^{jk}$ by lemma 2). (74) is therefore true.

# D Proposition 1

We show that $(\bar{x}; I_n)$ is a SPEC of $\Gamma$ if and only if:

(i) either $\forall i, j \in N$, $\Pi_i^j(\bar{x})\Pi_j^i(\bar{x}) = 0$,

(ii) or $\forall i, j \in N$, $i \neq j$, $\Psi_i^i(\bar{x}) = 0$,

with $\bar{x} \in X$ the Nash equilibrium of $\Gamma$, and $I_n$ a matrix in $\mathbb{R}^{n \times n}$ such that $\sigma_{ij} \neq 0$ if and only if $i = j$.
$\forall S \in \mathbb{R}^{n \times n}$, there exists a unique Nash equilibrium for $\Gamma_2(S)$ $\bar{x} \in X$ that verifies, $\forall i \in N$:

$$U_i^i(\bar{x}|S) = 0, \tag{84}$$

$$\sum_{j \in N} \sigma_{ij}\Pi_j^i(\bar{x}) = 0. \tag{85}$$

By definition of the Stackelberg best reply function, the indirect payoff function $V_i : S \mapsto \mathbb{R}$ can be rewritten as follows:

$$V_i(S) = \Pi_i(\bar{x}(S)), \tag{86}$$

$$V_i(S) = \Pi_i(f_1(\bar{x}_i(S)); \ldots; \bar{x}_i(S); \ldots; f_n(\bar{x}_i(S))), \tag{87}$$

$$V_i(S) = \Psi_i(\bar{x}_i|S). \tag{88}$$

The relation (88) implies that, at the Nash equilibrium of $\Gamma_1$, player $i$ maximises her Stackelberg payoff function when she maximises her indirect payoff $V_i$. We must therefore verify, at the Nash equilibrium of $\Gamma_1$:

$$\frac{\partial V_i}{\partial \sigma_{ij}}(\bar{S}) = \Psi_i^i(\bar{x}_i|\bar{S})\frac{\partial \bar{x}_i}{\partial \sigma_{ij}} = 0, \qquad \forall j \in N, \tag{89}$$

$$\text{i.e. either} \quad \frac{\partial \bar{x}_i}{\partial \sigma_{ij}}(\bar{S}) = 0, \qquad \forall j \in N, \tag{90}$$

$$\text{or} \quad \Psi_i^i(\bar{x}_i|\bar{S}) = 0.. \tag{91}$$

If (90) is not true, then $I_n$ is a first stage game equilibrium if and only if (91) holds for $\bar{x}$, the Nash equilibrium of $\Gamma$. We have therefore proven the condition (ii) of proposition 1.

We now determine the conditions under which the conditions (90) holds. We must therefore characterise $\bar{x}(S)$, the Nash equilibrium of $\Gamma_2(S)$. We identify the best reply of player $i$, $\forall i \neq j$, when a player $j$ unilaterally changes her strategy $S_j$. We consider here the differential of $U_i^i$, and look for the reactions $\mathrm{d}x_i$ that verify $\mathrm{d}U_i^i(\bar{x}) = 0$, $\forall i \in N$. We have the following relations:

$$\mathrm{d}U_i^i(\bar{x}) = 0, \qquad \forall i \in N, \tag{92}$$

$$\sum_{j \in N} \left[ U_i^{ij}(\bar{x}) \ \mathrm{d}x_j + \Pi_j^i(\bar{x}) \ \mathrm{d}\sigma_{ij} \right] = 0, \qquad \forall i \in N. \tag{93}$$

We solve this system of linear equations in $\mathrm{d}x_i$:

$$\begin{pmatrix} U_1^{11}(\bar{x}) & \ldots & U_1^{1n}(\bar{x}) \\ \ldots & & \ldots \\ U_n^{n1}(\bar{x}) & \ldots & U_n^{nn}(\bar{x}) \end{pmatrix} \begin{pmatrix} \mathrm{d}x_1 \\ \ldots \\ \mathrm{d}x_n \end{pmatrix} + \begin{pmatrix} \sum_{j \in N} \Pi_j^1(\bar{x}) \, \mathrm{d}\sigma_{1j} \\ \ldots \\ \sum_{j \in N} \Pi_j^n(\bar{x}) \, \mathrm{d}\sigma_{nj} \end{pmatrix} = 0, \tag{94}$$

$$J(S) \, \mathrm{d}x + \mathrm{d}A = 0, \tag{95}$$

with $\mathrm{d}x = {}^t\{\mathrm{d}x_i\}_{i \in N}$ the column vector of strategies' variations; $\mathrm{d}A = {}^t\{\mathrm{d}A_i\}_{i \in N}$. We make the additional assumption that $J(S)$ and its minors $J_{ii}(S)$ are generically non singular $\forall S \in \mathbb{R}^{n \times n}$. The system (95) has therefore a unique solution (for notational convenience, we do not mention on which set of parameters $S$ $J$ is defined, unless a confusion is possible):

$$\mathrm{d}x_i = \frac{|J^i|}{|J|} \qquad \forall i \in N, \tag{96}$$

with $J^i$ a $n \times n$ matrix identical to $J$, except for the $i^{th}$ column which is replaced by $- \mathrm{d}A$. We deduce the following relations:

$$\mathrm{d}x_i = -\frac{\sum_{k \in N} C_{ki}^J \, \mathrm{d}A^k}{|J|}, \tag{97}$$

$$\implies \quad \frac{\partial \bar{x}_i}{\partial \sigma_{ik}}(S) = -\Pi_k^i \frac{C_{ii}^J}{|J|}(\bar{x}(S)) \qquad \forall S \in \mathbb{R}^{n \times n}. \tag{98}$$

We can therefore see the condition (90) implies:

$$\Pi_k^i(\bar{x}) = 0 \qquad \forall k \in N. \tag{99}$$

This last condition means that the strategy profile that maximizes the utility function $U_i$ of player $i$ also maximizes her own payoff $\Pi_i$ as well as the payoff of all the other players $j \neq i$ (or minimizes it). It means therefore that, if $\forall i, k \in N$, $\Pi_k^i(\bar{x}) = 0$ at Nash equilibrium, then $(\bar{x}; I_n)$ is a SPEC.

# E   Proposition 2

We prove here that $\bar{S}$ is a Nash equilibrium of the first stage game $\Gamma$ when:

$$\bar{\sigma}_{ij} = \frac{\Pi_i^j}{\Pi_j^i}(\bar{x}) \frac{\partial f_j}{\partial x_i}(\bar{x}_i | \bar{S}). \tag{100}$$

We look for conditions under which (91) is verified at the first stage game equilibrium. We can therefore rewrite the first order condition of the first stage game equilibrium (89):

$$\sum_{j \in N} \Pi_i^j \frac{\partial f_j}{\partial x_i}(\bar{x}_i(S)) = 0. \tag{101}$$

Combining equations (85) and (101), we can obtain an expression of the parameters $\sigma_{ij}$ at equilibrium:

$$\sum_{j \in N} \Pi_i^j \frac{\partial f_j}{\partial x_i} = \sum_{j \in N} \sigma_{ij} \Pi_j^i. \tag{102}$$

We can then suggest the following specification for the first stage game equilibrium:

$$\sigma_{ij} = \frac{\Pi_i^j}{\Pi_j^i} \frac{\partial f_j}{\partial x_i}. \tag{103}$$

Note that, since we are maximising the Stackelberg function in the first stage game $\Gamma_1$, we only have $n$ equations to determine the $n^2$ parameters $\sigma_{ij}$. Although other specifications were possible, we chose here to define $\sigma_{ij}$ as a function of the Stackelberg best reply of player $j$ when $i$ is the leader, since it captures the idea that the behaviour of $i$ towards $j$ fundamentally depends on the way $j$ reacts when $i$ changes her strategy.

# F   Proposition 3

We prove that, under the following assumptions:

(i) $|U_i^{ii}| > (n-1)|U_i^{ij}|$,

(ii) $|U_i^{ik}| < (n-1)|U_i^{ij}|$,

(iii) $|\Pi_j^{ji}| \geq |\Pi_k^{ji}|$,

a symmetric SPEC $(\bar{x}; \bar{S})$ verifies:

$$sign\left(\bar{\sigma}_{ij}\right) = sign\left(\Pi_j^{ji}(\bar{x}(S))\right), \qquad \forall j \neq i. \tag{104}$$

A symmetric SPEC implies that $\Pi_i^j(\bar{x}(\bar{S})) = \Pi_j^i(\bar{x}(\bar{S}))$. The optimal weights $\bar{\sigma}_{ij}$ are therefore:

$$\bar{\sigma}_{ij} = \frac{\partial f_j}{\partial x_i}(\bar{x}(\bar{S})). \tag{105}$$

Lemma 3 ensures that, under conditions (i) and (ii), we have:

$$\sum_{j \neq i} |\bar{\sigma}_{ij}| < 1. \tag{106}$$

We can then easily deduce the following relation, when condition (iii) is verified:

$$\left|\Pi_i^{ij}\right| > \left|\sum_{k \neq i} \sigma_{ij} \Pi_k^{ij}\right|, \tag{107}$$

$$sign(U_i^{ij}) = sign(\Pi_i^{ij}). \tag{108}$$

By lemma 2, we know that $\bar{\sigma}_{ij}$ has the same sign than $U_j^{ji}$. We have therefore, at a symmetric SPEC:

$$sign\left(\bar{\sigma}_{ij}\right) = sign\left(\Pi_j^{ji}(\bar{x}(S))\right), \qquad \forall j \neq i. \tag{109}$$