

# Interaction and imitation with heterogeneous agents: a misleading evolutionary equilibrium\*

Francisco Cabo<sup>†</sup>

Ana García González<sup>‡</sup>

**Submitted version (not the final accepted version)**

## Abstract

In a two-population evolutionary game we analyze the interaction between individuals belonging to two populations with the same strategy set but different payoffs. Agents play a game against individuals in the two populations. They imitate agents belonging to the same and also the alternative population. When a revising agent is matched with an individual in the alternative population who plays differently, his expected payoff and the observed payoff of his partner diverge. Hence, he conjectures the payoff from switching to the other strategy by weighing what he expected and what he observes. The evolutionary dynamics has a unique asymptotically stable fixed point, which typically differs from the evolutionary stable equilibrium without inter-population imitation. For a collective action game we analyze to what extent the compliance rate and the social welfare differ from the Nash equilibrium, and how these gaps depend on the confidence that agents assign to what they see.

*Keywords:* Two-population evolutionary game, inter-population imitation, evolutionary stable equilibrium.

*JEL Classification:* C73, D91.

---

\*We would like to thank Fabio Lamantia, as well as Georges Zaccour and attendants to the GERAD seminar for their useful comments and suggestions. This study was funded by the Spanish Government (projects ECO2014-52343-P and ECO2017-82227-P), as well as financial aid from Junta de Castilla y León (projects VA024P17 and VA105G18), co-financed by FEDER funds.

<sup>†</sup>Corresponding author: IMUVa, Universidad de Valladolid. Mail address: Dept. Economía Aplicada (Matemáticas). Avda. Valle Esgueva, 6. 47011 Valladolid, Spain. Fax: +34 983423299. Tel.: +34 983186662. e-mail: pcabo@eco.uva.es

<sup>‡</sup>IMUVa, Universidad de Valladolid, Spain.

# 1 Introduction

In his work *Don Quixote de la Mancha*, Cervantes presents a delirious Don Quixote, whose “brain dries up” after reading too many books on chivalry and decides to become a knight and to fight injustice. The second character in this novel, Sancho, is the archetype of a perfectly rational individual or *homo oeconomicus*. And yet, when Don Quixote starts his travel in search of fame and glory Sancho joins him. For an outside observer, the dilemma between staying at home or following a madman and face deprivation seems trivial. And it is this decision of a rational individual to join Don Quixote, when his senses warn him of the negative consequences of this journey, what thrills the reader of the novel. Our explanation for such behavior relies on Sanchos’ inability to clearly anticipate the result of joining the knight-errant life of Don Quixote. Between his own expectation of an unfortunate adventure, and Quixote’s strong confidence in success, the latter is powerful enough to convince Sancho to imitate Don Quixote’s madness and act in a way which will not be in his best interest.

The apparently strange behavior of Sancho is not uncommon in real life situations. Consider, as an example, that an “ordinary” individual, S, meets a “sophisticated” person, Q, who is highly enthusiastic about a new fancy restaurant. Then S, who enjoys his regular pizza place, conjectures how he will enjoy the new restaurant by taking into account the low expectations he has on the new place together with the strong satisfaction he observes in Q. If this conjecture exceeds the satisfaction from his favorite pizza place, S could decide to visit the new restaurant. Because he is not a “sophisticated” person, he will not enjoy it as much as Q. Nonetheless, in a densely populated society, the following Saturday another “ordinary” person, following some other “sophisticated” person’s advice, will show up at the restaurant. At the equilibrium, the restaurant serves food to a given share of “sophisticated” individuals, plus a positive share of “ordinary” individuals who mistakenly imitate them. Similar situations are also common, for example, in marketing. We imitate the buying behavior of the models we observe in the commercials, refusing to acknowledge that the clothes that we buy will not fit as well on us as they fit on them. Another example is the use of internet sites like tripadvisor, that guide us to make buying decisions based on the experiences of other costumers who, typically, do not share the exact same preferences with us. In fact, in real life situations, anytime one individual follows someone else’s advice, he is uncertain on whether he will get the exact same satisfaction, and indeed, commonly does not.

To analyze the strategic interaction between heterogeneous individuals who imitate one another, we consider a two-population evolutionary game. When dealing with two distinct populations, the standard approach is to consider that players belonging to one group play a game with players from the other population, see, for example, Gong, Gao and Cao (2018) and references therein; Antoci *et al.* (2009) and (2012) in finance and environmental prob-

lems; or Antoci *et al.* (2011) in traffic congestion. This approach defines an asymmetric two-population game, in which agents in one population are paired with agents from the other, as in the incumbent-intruder, predator-prey or buyer-seller settings. In contrast, we assume that agents play against agents belonging to the same or the other population: intra and inter-population interaction. Thus, the payoff that one agent gets depends on what all other individuals are doing in both populations. Although uncommon, some authors have also considered heterogeneous agents who interact with each other. An example is the seminal work by Matsuyama (1991) and the subsequent literature on fashion cycles. They consider a population composed of conformists and non-conformists with different preferences: the former love to join the crowd and the later love individuality. Similarly, Giovinazzo and Naimzada (2015) or Naimzada and Pireddu (2018) also analyze the fashion cycle considering a population share updating mechanism between heterogeneous agents exhibiting both bandwagon and snob behaviors. Radi and Gardini (2018) analyze residential segregation for two groups of people (black and white) with different preferences for segregation. In biology, when analyzing mutualism, typically only between-species interaction is considered. Interestingly, Gokhale *et al.* (2019) find that the inclusion of within-species interaction can lead to stability.

It is important to highlight that in the standard setting in multi-population evolutionary games, when considering an imitative revision protocol, individuals are exclusively paired with and can only imitate the strategies played by individuals from their same population. This is clearly the case when referring to natural species, since genetic traits cannot jump between species. Likewise, this is also the standard hypothesis when referring to general evolutionary games applied in social sciences: the strategies which provide higher payoffs tend to be imitated by the individuals within the population (see, for example, Sandholm [2010]). In contrast, we allow individuals in one population to imitate the behavior of agents within and also in the other population. This hypothesis of inter-population imitation, proposed in Cabo and García-González (2019) and Cabo *et al.* (2019), helps to escape the well known result in asymmetric evolutionary games that strictly mixed-strategy Nash equilibria (NE) are not asymptotically stable.<sup>1</sup> Under this hypothesis it is shown that an equilibrium at which one of the two populations plays a mixed-strategy is an asymptotically stable fixed point of the evolutionary dynamics. This equilibrium typically differs from the Nash equilibrium.

Starting from the Sancho-Quixote metaphor, we study a collective action problem, as is compliance with social norms. There is a large literature on imitation in social dilemmas. The imitation revision protocol considered in the paper is quite common in the literature ever since Schlag (1998), which also considers a two-population game but without inter-population

---

<sup>1</sup>This was first highlighted by Selten (1980) who showed that every evolutionary stable equilibrium (ESS) must be a pure Nash equilibrium, and later proven by Hofbauer and Sigmund (1998) and Samuelson and Zhang (1992).

imitation. Other authors, starting from the work by Kandori *et al.* (1993), analyze the role of the matching mechanism when imitation is globally considered, as in Robson and Vega-Redondo (1997), or locally, as in Eshel *et al.* (1998) or Khan (2014). The standard assumption in these works is to assume homogeneous agents. By contrast, Apesteguia *et al.* (2010) highlights the lack of robustness when imitation involves agents with preferences slightly heterogeneous. Imitation involve not fully homogeneous agents also in networks, where agents differ in their position in the network and the opponents faced (see, for example, Khan 2014, or Cui and Wang 2016).

The idea of considering different individuals to analyze social dilemmas can be found in Ostrom (2000), who postulated the existence of individuals more inclined to comply than the rational egoist *homo oeconomicus*.<sup>2</sup> Thus, we associate selfish individuals to Sanchos, while norm-using Quixotes are more inclined to comply with the norms. The latter get a higher reward from compliance, either due to pure altruism, interest in gaining prestige, respect or friendship, or just an inner satisfaction or “warm-glow” from the mere fact of compliance (as stated by Andreoni [1990]). Likewise, defection harms Quixotes more strongly, either because they want to avoid scorn, or because doing bad involves a dis-utility or a “cold-prickle” (again in the terminology of Andreoni [1995]).

The two agents have different payoff matrices and face a different social dilemma. For Sanchos, defection is a dominant strategy, and the compliance with social norms is a prisoner’s dilemma game: everyone prefers to defect and the collective choice of the dominant strategy results in a bad equilibrium. On the other hand, we distinguish two type of Quixotes. The standard-Quixote would comply if his opponent defects, although he still has a free-riding incentive if the other complies, i.e., for him the social dilemma is represented as a snowdrift game. Alternatively, more extreme “mad-Quixotes” want to comply no matter what others are doing (these agents face no dilemma). In either case, from the assumption of intra and inter-population interaction, the welfare of an individual depends on the compliance decisions of all other individuals within his own and the other population. Moreover, intra and inter-population imitation allows an individual to imitate the decision made by any other agent he is paired with, regardless of whether he belongs to the same population or not.

For a revising agent paired with someone from the other population playing differently, the observed payoff diverges from his expected payoff. While the agent is certain of his current payoff, the payoff to the alternative strategy is uncertain given the discrepancy between observations and expectations. He conjectures what he would get by balancing the confidence he assigns to what he sees and to what he expected. By comparing this conjecture against the payoff to his current strategy, the agent decides whether or not to shift to the strategy

---

<sup>2</sup>A similar idea of a society divided between standard Nashian and Kantian individuals is presented in Grafton *et al.* (2017).

played by his partner. It is important to state clearly that, if he chooses to shift, his payoff is not that conjecture, but his expected payoff (defined in the payoff matrix). Thus, the more he trusts what he sees, the less accurate the conjecture and the more misled is the agent. For that reason, we call the equilibrium emerging from the assumption of inter-population imitation a misleading evolutionary equilibrium.

The equilibrium is unique and asymptotically stable, and typically different from the Nash equilibrium. Depending on parameters value, we can distinguish several possible situations: An equilibrium where all Quixotes comply and Sanchos imitate their compliant behavior partially or completely; an equilibrium where all Sanchos defect and Quixotes imitate their non-compliant behavior, again partially or completely; and finally, an equilibrium where all Sanchos defect and all Quixotes comply. The objective of the paper is twofold: First, we seek to analyze how the conjecture formation determines the type of equilibrium reached. And second, we study how the social welfare compares to the ESS equilibrium which would be reached without inter-population imitation and which coincides with the NE. In particular, we study how this comparison is affected by the confidence that agents assign to what they see.

As one would expect, Sanchos imitate compliant Quixotes if the former strongly trust what they see when paired with the latter, and Quixotes greatly differ from Sanchos. In particular, a solution where all Sanchos comply is only possible when playing against mad-Quixotes. The equilibrium where Quixotes imitate non-compliant Sanchos occurs if Quixotes highly trust what they see and are not very different from Sanchos. Finally, the less the agents trust what they see in others, the more likely is a solution in which Quixotes comply and Sanchos do not, just as in the NE or the ESS with no inter-population imitation.

Interestingly, social welfare in the misleading evolutionary equilibrium is not necessarily lower than in the NE. At the equilibrium where Quixotes imitate the non-compliant behavior of Sanchos, the latter are clearly worse off because the compliance rate in the global population is smaller. Because Quixotes value compliance and global compliance decreases, the social welfare for Quixote is also reduced. On the other hand, when Sanchos imitate compliant Quixotes, these latter are better off, since more people comply. The greater the confidence that Sanchos assign to what they see, the more they imitate compliance. A higher imitation typically enhances (resp. reduces) the welfare of Sanchos when their share in the global population is large (resp. small). However, when the two populations are of similar size, a higher confidence in what they see will initially increase the welfare of Sanchos, but as more and more Sanchos imitate compliance, the initial rise in welfare reverses, eventually leading to a worse off situation.

Section 2 explains what we understand by intra and inter-population interaction and imitation. The evolutionary dynamics and the different possible equilibria are presented in Section 3. Section 4 compares the misleading evolutionary equilibria to the NE and the average pay-

off obtained in each population. Section 5 concludes. Technical details are explained in the appendix.

## 2 Two-population evolutionary game

This section presents a two-population evolutionary game, in which individuals from the two populations share the same strategy set and differ in preferences. We denote these two populations Sanchos (Ss) and Quixotes (Qs), and analyze a collective action problem like the compliance with social norms. Individuals within either population have the same decision to make, to comply or to defect, but they differ in the payoffs from this decision. Individuals from one population play a game against individuals within their own as well as the other population. Moreover, when they revise their strategy, they can imitate individuals within or in the other population.

### 2.1 Intra-population and inter-population interaction

Sancho is the archetype of the *homo oeconomicus*. Thus, in a social dilemma, the strategic interaction between every two rational Ss can be represented by a prisoner's dilemma. The payoffs matrix for these players is given by

	C	D
C	0	$-d$
D	$b$	$-\phi$

Table 1: Payoffs matrix for Ss

We assume that

$$b, d, \phi > 0, \quad \text{and} \quad 0 < d - \phi < b. \quad (1)$$

Thus, if the opponent complies and S defects, the opponent bears the cost  $-d$ , and S gets  $b > 0$ , which defines the free-riding incentive. If the opponent defects, S still has an incentive to defect since we are assuming  $d > \phi$ , and hence defection is the dominant strategy. Moreover, since we further assume  $d - \phi < b$ , then defection, as opposed to compliance, is relatively more rewarding when the opponent complies. This defines a prisoner's dilemma in which mutual defection is the unique Nash equilibrium. It would be the evolutionary stable strategy if an evolutionary game with pairwise imitation were defined for Ss as an isolated population.

Qs are pro-social individuals who obtain an inner satisfaction or warm-glow from compliance and, furthermore, defection induces discomfort or cold-prickle. Thus, Qs attach a higher payoff to compliance and a stronger dissatisfaction to defection than Ss. The addition of the

warm-glow plus the cold prickle, denoted by  $\varepsilon$ , defines the absolute distance which separates the two populations. If we denote by  $\alpha \in (0, 1)$  the relative importance of the warm-glow in contrast with the cold-prickle, these can be written as  $\alpha\varepsilon$  and  $(1 - \alpha)\varepsilon$ , respectively. Therefore, the payoff matrix for Qs is

	C	D
C	$\alpha\varepsilon$	$\alpha\varepsilon - d$
D	$b - (1 - \alpha)\varepsilon$	$-\phi - (1 - \alpha)\varepsilon$

Table 2: Payoffs matrix for Qs

Depending on the strength of these two effects, we can distinguish two types of Qs. For standard-Qs, the warm-glow from compliance plus the cold-prickle from defection,  $\varepsilon$ , is not enough to counteract the free riding incentive, although together they make compliance attractive when the opponent defects. In consequence, Qs still have an incentive to free-ride on the compliance of others but prefer compliance when the opponent defects. This defines the social dilemma as a snowdrift game. For mad-Qs, the warm-glow from compliance plus the cold-prickle from defection are so strong that the dilemma disappears, as compliance becomes the dominating strategy. Conditions for one type or the other can be written as

$$d - \phi < \varepsilon < b \text{ (standard-Qs)}, \quad b \leq \varepsilon \text{ (mad-Qs)}. \quad (2)$$

The Nash equilibrium for Qs is given by

$$(C, D) = \begin{cases} (\Delta, 1 - \Delta) & \text{if } d - \phi < \varepsilon < b, \\ (1, 0) & \text{if } \varepsilon \geq b, \end{cases} \quad \text{with } \Delta = \frac{\varepsilon - (d - \phi)}{\sigma}. \quad (3)$$

The term  $\sigma = b - (d - \phi) > 0$  is defined as the sum of the off-diagonal payoffs minus the sum of the diagonal payoffs, and is positive under condition (1). For mad-Qs compliance is the dominant strategy and the NE is given by  $(1, 0)$ , while for standard-Qs facing a snowdrift game the NE in mixed strategies reads  $(\Delta, 1 - \Delta)$ . Under condition (2-left)  $\Delta \in (0, 1)$  and the NE would determine the ESS in an evolutionary game for an isolated population of Qs.

We do not consider isolated populations, but analyze a two-population evolutionary game in which agents play an intra-population as well as an inter-population game. The total mass of population is normalized to one, with a share  $s \in (0, 1)$  of Ss and hence a share  $1 - s \in (0, 1)$  of Qs. These shares do not vary because we assume unchanging preferences.<sup>3</sup> The ratio of Ss and Qs who comply is denoted by  $x$  and  $y$ , respectively, and correspondingly,  $1 - x$  and  $1 - y$  are the ratios of non-compliant Ss and Qs. In consequence, the set of social states can

<sup>3</sup>A interesting challenging extension would be to allow for evolution to also operate at the level of preferences as, for example, in Alger and Weibull (2013) and references therein.

be defined as  $X = \{\bar{x}' = (xs, (1-x)s, y(1-s), (1-y)(1-s)), | x, y \in [0, 1]\}$ . Let us denote  $\mathbb{S}, \mathbb{Q} \in \mathcal{M}_{2 \times 2}(\mathbb{R})$  the payoff matrices in Tables 1 and 2. Then, for a given state  $\bar{x}$ , since we allow for intra-population and inter-population interaction, the vector of payoffs can be computed as<sup>4</sup>

$$\bar{\pi}^t = (\pi_C^S, \pi_D^S, \pi_C^Q, \pi_D^Q) = \left( \begin{array}{c|c} \mathbb{S} & \mathbb{S} \\ \hline \mathbb{Q} & \mathbb{Q} \end{array} \right) \bar{x}.$$

Superscripts S and Q denote Sanchos and Quixotes, while subscripts C and D denote compliance and defection. For conciseness, these payoffs can be rewritten as a function of the share of compliance in the global population,  $q = xs + y(1-s)$ ,

$$\bar{\pi}^t(q) = (dq - d, (b + \phi)q - \phi, \alpha\varepsilon + dq - d, (b + \phi)q - \phi - (1 - \alpha)\varepsilon). \quad (4)$$

The payoff of compliance for Qs is equal to the payoff attained by Ss plus the warm-glow:  $\pi_C^Q(q) = \pi_C^S(q) + \alpha\varepsilon$ . Likewise, a Q endures a cold-prickle from defection which does not affect Ss:  $\pi_D^Q(q) = \pi_D^S(q) - (1 - \alpha)\varepsilon$ .

As shown in expression (4), intra-population and inter-population interaction imply that the payoff obtained by individuals in one population depends on what the agents in this and the other population are doing. The intra-population and inter-population dimensions are also present at the imitation process, when a revising individual must decide whether to maintain or to modify his current strategy.

## 2.2 Intra-population and inter-population imitation

Assuming a pairwise imitation revision protocol, a revising individual can be randomly paired with and imitate someone belonging to his own or the other population. To explain this process, consider as an example that a compliant Q receives a revising opportunity. The standard process in the literature assumes that he can only be paired with another Q. If his partner plays differently (defects), he will compare the payoff obtained by his partner against his own payoff:  $\pi_D^Q(q) - \pi_C^Q(q) = \sigma(q - \Delta)$ . If positive, the wider the payoff gap between the alternative and the current strategy, the more clearly the revising agent perceives that switching strategies is worthwhile. For this reason it is commonly assumed that the conditional imitation rate is proportional to the gap between payoffs when positive. Thus, the conditional imitation rate of an agent  $h \in \{S, Q\}$ , playing strategy  $i \in \{C, D\}$ , of switching to strategy  $-i$  when paired with someone within his own population is defined as  $r_{-i}^h(q) = [\pi_{-i}^h(q) - \pi_i^h(q)]_+$ . Notice that

---

<sup>4</sup>The symmetric and the asymmetric game would correspond to the block diagonal and the block off-diagonal matrices of the form:

$$\left( \begin{array}{c|c} \mathbb{S} & (0) \\ \hline (0) & \mathbb{Q} \end{array} \right), \quad \left( \begin{array}{c|c} (0) & \mathbb{S} \\ \hline \mathbb{Q} & (0) \end{array} \right).$$



by definition  $r_{-i}^h(q) > 0$  implies  $r_i^h(q) = 0$ . From (4), the conditional imitation rates read:

$$r_C^Q(q) = \sigma [\Delta - q]_+, \quad r_D^Q(q) = \sigma [q - \Delta]_+, \quad r_C^S(q) = 0, \quad r_D^S(q) = \sigma q + d - \phi > 0. \quad (5)$$

In our formulation inter-population imitation is also allowed. The revising compliant Q could be also paired with a S. Assume that this S defects. The revising agent knows what he is getting from compliance and learns the payoff of defection obtained by his partner. Now he detects a disparity between what he observes his partner is getting,  $\pi_D^S$ , and what he expected he would get by switching to defection,  $\pi_D^Q$ . Therefore, he is uncertain about the reward if he chooses to switch to defection. Uncertainty did not appear when paired with someone from his own population, and only becomes apparent now when paired with someone with different preferences. The revising agent could realize that his partner is an individual belonging to a different population with different preferences and ignore what he has learned. This is the standard assumption in two-population evolutionary games. At the other extreme, like St. Thomas the Apostle, he could base his decision only on what he sees, assuming that his expectations are wrong. Finally, we assume that he conjectures the payoff of switching to defection by weighing what he expected, given the population state, and what he observes:  $E_D^Q(q) = p^Q \pi_D^S(q) + (1 - p^Q) \pi_D^Q(q)$ . The weight  $p^Q$  represents how much Q trusts what he sees in contrast to  $1 - p^Q$ , which represents his confidence in his own expectations. This conjectured payoff is contrasted against his current payoff in order to decide whether to switch strategies.

Following identical reasoning, the conjectured payoff of an  $h$ -type<sup>5</sup> revising individual playing strategy  $i$ , who is paired with an individual from the other population playing the alternative strategy  $-i$  reads:

$$E_{-i}^h(q) = p^h \pi_{-i}^{-h}(q) + (1 - p^h) \pi_{-i}^h(q), \quad h \in \{S, Q\}, \quad i \in \{C, D\}, \quad (6)$$

where  $p^h$  represents the confidence of an  $h$ -type in what he sees, or to what extent he distrusts his own expectations - not necessarily the same for both agent types.

Under inter-population imitation, a revising  $h$  agent paired with an agent from the other population builds his conditional imitation rate by comparing his current payoff from strategy  $i$  against the conjectured payoff of the alternative strategy  $-i$ :<sup>6</sup>

$$R_{-i}^h(q) = \left[ E_{-i}^h(q) - \pi_i^h(q) \right]_+, \quad h \in \{S, Q\}, \quad i \in \{C, D\}. \quad (7)$$

Under pairwise imitation inter-population imitation makes possible three different situations: one strategy is imitated and the other is not (as in the standard multi-population

---

<sup>5</sup>Note that an agent type refers to the population he belongs to, and not to the strategy he plays.

<sup>6</sup>Small  $r$  in (5)/capital  $R$  in (7) refers to the conditional imitation rate when an individual in one population is paired with someone from his own/the other population. The expressions for  $R_{-i}^h(q)$  are given in (25)-(30) in the Appendix.

evolutionary game); the two strategies are imitated at the same time; or none of the strategies is imitated.

**Proposition 1** *Under inter-population imitation one can distinguish three scenarios depending on  $q$ .*

1.  $q \leq \min\{L^h, U^h\}$ : then  $R_C^h(q) > 0, R_D^h(q) = 0$ .

2.  $q(\min\{L^h, U^h\}, \max\{L^h, U^h\})$ :

*i) Ss (for  $\alpha > 1/2$ ) and Qs (for  $\alpha < 1/2$ ) imitate both compliance and defection if:*

$$L^h < U^h \quad \Rightarrow \quad R_C^h(q), R_D^h(q) > 0, \forall q \in (L^h, U^h), h \in \{S, Q\}.$$

*ii) Ss (for  $\alpha < 1/2$ ) and Qs (for  $\alpha > 1/2$ ) do neither imitate compliance nor defection if:*

$$U^h < L^h \quad \Rightarrow \quad R_C^h(q) = R_D^h(q) = 0, \forall q \in (U^h, L^h), h \in \{S, Q\}.$$

3.  $q \geq \max\{L^h, U^h\}$ : then  $R_C^h(q) = 0, R_D^h(q) > 0$ .

**Proof.** See Appendix.<sup>7</sup> ■

Figure 1 summarizes the different situations when a revising  $h$ -type agent is paired with a  $-h$ -type agent playing the alternative strategy. If the global compliance rate is quite small, a revising individual will imitate compliance and not defection. The conjectured payoff to compliance is always attractive, while the conjectured payoff to defection never is. The opposite is true if the global compliance rate is very large: defection is imitated and compliance is not. For intermediate values two situations are possible. One possibility is that, regardless of the strategy played by the revising agent, the alternative strategy never pays more and, hence, he does not switch his current strategy (Figure 1 left). This situation occurs for Ss if the warm-glow is of relatively little importance compared to the cold-prickle ( $\alpha < 1/2$ ), and for Qs when the the warm-glow is relatively more important ( $\alpha > 1/2$ ). The other possible situation involves a revising agent playing any strategy, who finds it attractive to switch to the alternative strategy (Figure 1 right). This occurs in the opposite case, for Ss when  $\alpha > 1/2$  and for Qs when  $\alpha < 1/2$ .

### 3 Evolutionary dynamics and MEE

This section characterizes the evolutionary dynamics and the different possible equilibria, depending on parameters values. We denote these equilibria as Misleading Evolutionary Equilibria

<sup>7</sup>The expressions for  $U^h, L^h$  are given in (29) and (30) in the Appendix.

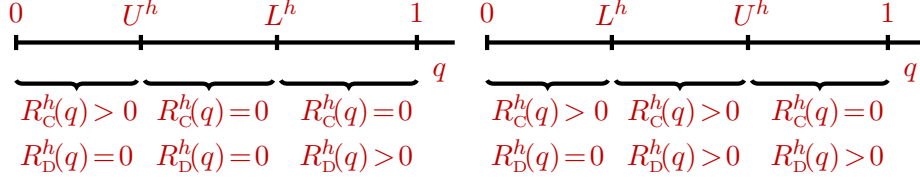


Figure 1: Conditional imitation rates across populations

(henceforth MEE). The evolutionary dynamics for population  $h$  is given by the share of defecting agents times the probability of their switching to compliance,  $\rho_C^h$ , minus the share of those who comply times the probability with which they switch to defection,  $\rho_D^h$ . Thus, the dynamics of the share of compliance for Ss and Qs reads:

$$\dot{x} = (1 - x)\rho_C^S - x\rho_D^S, \quad (8)$$

$$\dot{y} = (1 - y)\rho_C^Q - y\rho_D^Q, \quad (9)$$

From the assumption of inter-population imitation, the probability of switching strategies must take into account two possible encounters. A revising agent can be paired with someone from his own or from the other population:<sup>8</sup> the probability of being paired with an individual in population  $h \in \{S, Q\}$  playing strategy  $i \in \{C, D\}$  is given by the shares of  $h$ -type individuals playing strategy  $i$  (in the vector of social states,  $\bar{x}$ ). Thus, the probability of a revising individual to switching to the alternative strategy, can be written as<sup>9</sup>

$$\rho_C^S = xsr_C^S + y(1 - s)R_C^S, \quad \rho_D^S = (1 - x)sr_D^S + (1 - y)(1 - s)R_D^S, \quad (10)$$

$$\rho_C^Q = y(1 - s)r_C^Q + xsR_C^Q, \quad \rho_D^Q = (1 - y)(1 - s)r_D^Q + (1 - x)sR_D^Q. \quad (11)$$

Plugging these probabilities into (8)-(9) and taking into account the definition of  $r_i^h(q)$  and  $R_i^h(q)$  in (5), (7), the evolutionary dynamics is described by a system of two differential equations. The expressions which define this system of equations vary depending on the value taken by the global compliance rate.<sup>10</sup> For these dynamics, five different MEE are feasible depending on the parameter values, as described in the following proposition.

**Proposition 2** *The system of differential equations given in (8) and (9), with conditional imitation rates in (5), (7) and switching probabilities in (10)-(11), presents five different MEE depending on the parameter values.*

<sup>8</sup>We treat the two probabilities equally. One could alternatively assume that it is more likely to meet someone from your own population. While this complicates notation, we do not believe it would fundamentally change the results.

<sup>9</sup>The  $q$  argument has been removed for simplicity, however, these probabilities depend on the global and on each population's compliance rates.

<sup>10</sup>From Figure 1, it follows that the ordering of the upper and lower bounds  $U^h, L^h \forall h \in \{S, Q\}$ , (and whether they are actually positive) determines different regions for  $q$  within which the evolutionary dynamics behaves differently. These bounds and the dynamics for the different scenarios are described in the Appendix.

- $MEE_O = (0, 0)$ : Neither  $Ss$  nor  $Qs$  complies.
- $MEE_y = (0, y^*)$ : Only some  $Qs$  comply, with

$$y^* = \frac{\sigma + \varepsilon - (d - \phi) - \sqrt{(b - \varepsilon)^2 + 4\varepsilon\sigma p^Q(1 - \alpha)s}}{2\sigma(1 - s)} \in (0, 1). \quad (12)$$

- $MEE_{01} = (0, 1)$ : All  $Qs$  comply and all  $Ss$  defect.
- $MEE_x = (x^*, 1)$ : Some  $Ss$  comply together with all  $Qs$ , with

$$x^* = \frac{-(d - \phi) - 2\sigma(1 - s) + \sqrt{(d - \phi)^2 + 4\varepsilon\sigma p^S\alpha(1 - s)}}{2\sigma s} \in (0, 1). \quad (13)$$

- $MEE_I = (1, 1)$ : All  $Ss$  and all  $Qs$  comply.

Each of these  $MEE$  is asymptotically stable under the mutually excluding conditions:

$$MEE_O \Leftrightarrow \varepsilon(1 - (1 - \alpha)p^Q s) \leq b - \sigma, \quad (14)$$

$$MEE_y \Leftrightarrow b - \sigma < \varepsilon(1 - (1 - \alpha)p^Q s), \quad \varepsilon(1 - (1 - \alpha)p^Q) < b - \sigma s, \quad (15)$$

$$MEE_{01} \Leftrightarrow \alpha\varepsilon p^S \leq b - \sigma s \leq \varepsilon(1 - (1 - \alpha)p^Q), \quad (16)$$

$$MEE_x \Leftrightarrow b - \sigma s < \alpha\varepsilon p^S < \frac{b}{1 - s}. \quad (17)$$

$$MEE_I \Leftrightarrow \frac{b}{1 - s} < \alpha\varepsilon p^S. \quad (18)$$

**Proof.** See Appendix. ■

Conditions (14)-(18), divide the parameter space

$$\mathcal{P} = \{(b, d, \phi, s, p^S, p^Q, \alpha, \varepsilon), | b > d - \phi > 0, d, \phi > 0, s \in (0, 1), \alpha, p^S, p^Q \in [0, 1], \varepsilon > d - \phi\}$$

into five disjoint subsets  $\mathcal{P}_e$ , with  $\mathcal{P}_e \cap \mathcal{P}_j = \{\emptyset\}$   $e \neq j$ , and  $\cup_e \mathcal{P}_e = \mathcal{P}$ ,  $e, j \in \{O, y, 01, x, I\}$ . For a given vector of parameter values in  $\mathcal{P}$  there always exists an asymptotically stable equilibrium and this equilibrium is unique. We do not present the proof here, although the phase diagrams in Figure 2 show that, for each particular parameters constellation, the corresponding  $MEE_e$  is asymptotically stable.<sup>11</sup>

Note that  $(0, 0)$  and  $(1, 1)$  are always steady-state equilibria of the system of differential equations (8)-(11), although typically unstable, except under conditions (14) and (18), respectively. The equilibria in Proposition 2 are represented in Figure 2 as round solid black points when stable, and as round blank red points if unstable. Some of them are equilibria in pure strategies:  $MEE_O = (0, 0)$ ,  $MEE_{01} = (0, 1)$  and  $MEE_I = (1, 1)$ . In these equilibria all individuals in one population follow the same strategy. Likewise, under conditions (15) or (17)

<sup>11</sup>In Cabo and García-González (2019) it is proven that either the  $MEE_y$  or the  $MEE_x$  is the unique asymptotically stable fixed point of the evolutionary dynamics for specific parameters values. This analysis is carried out for the extreme case  $p^S = p^Q = 1$ .

agents in one population play mixed strategies.  $MEE_y$  is characterized by a population of Qs within which a share  $y^* \in (0, 1)$  complies ( $1 - y^* > 0$  defects). Similarly,  $MEE_x$  is compatible with a positive share,  $x^*$ , of compliant Ss. Among these equilibria, we can distinguish two main situations. Equilibria  $MEE_0$  and  $MEE_y$  correspond to situations in which some Qs imitate the non-compliant behavior of all Ss. Conversely, equilibria  $MEE_x$  and  $MEE_I$  correspond to situations in which some or all Ss imitate the compliant behavior of all Qs. In the intermediate case  $MEE_{01}$  individuals in one population never imitate individuals from the other population.

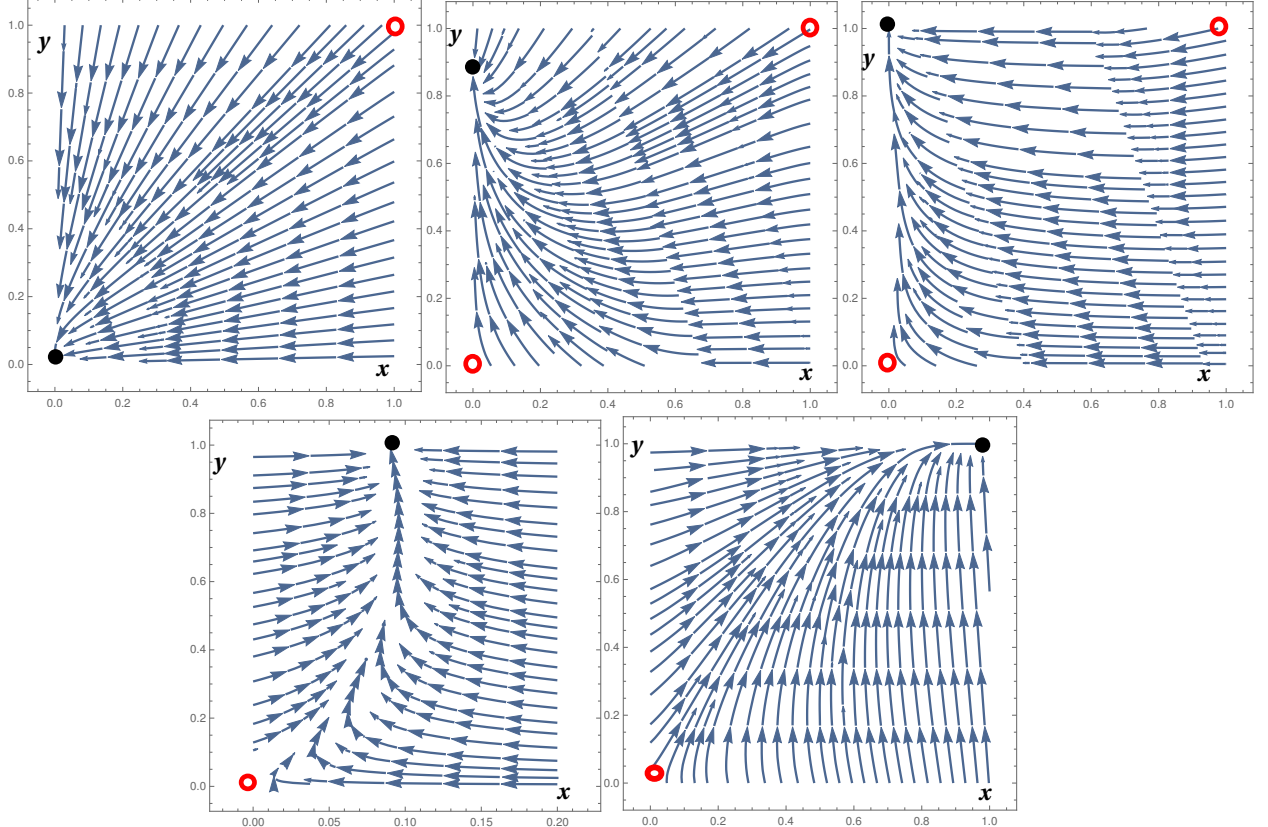


Figure 2: Stable equilibria  $MEE_0$  (up-left),  $MEE_y$  (up-center),  $MEE_{01}$  (up-right),  $MEE_x$  (down-left),  $MEE_I$  (down-right).

To clarify the meaning of conditions (14)-(17), note that they can be stated in terms of the absolute gap in preferences between Qs and Ss,  $\varepsilon$ , considering all other parameters fixed. Equilibrium  $MEE_0$ ,  $MEE_y$ ,  $MEE_{01}$ ,  $MEE_x$  or  $MEE_I$  is reached when  $\varepsilon$  takes values within the sub-interval:  $(d - \phi, \varepsilon_0)$ ,  $(\varepsilon_0, \varepsilon_y)$ ,  $(\varepsilon_y, \varepsilon_{01})$ ,  $(\varepsilon_{01}, \varepsilon_x)$ , or  $(\varepsilon_x, \infty)$ , respectively, with<sup>12</sup>

$$d - \phi \leq \varepsilon_0 \equiv \frac{b - \sigma}{1 - (1 - \alpha)p^{Q_S}} \leq \varepsilon_y \equiv \frac{b - \sigma s}{1 - (1 - \alpha)p^Q} \leq \varepsilon_{01} \equiv \frac{b - \sigma s}{\alpha p^S} \leq \varepsilon_x \equiv \frac{b}{\alpha p^S(1 - s)}. \quad (19)$$

As  $\varepsilon$  increases, we move from the equilibrium in which all Ss defect and all Qs imitate their non-compliant behavior, to successive equilibria in which only some Qs imitate defection; Ss

<sup>12</sup>With strict inequality for any  $s \in (0, 1)$ .

defect but Qs comply; all Qs comply and some Ss imitate this compliant behavior; all Qs comply and all Ss imitate compliance. The length of these intervals depends on  $p^Q$  and  $p^S$ . Therefore, not only does the absolute distance,  $\varepsilon$ , matter, but so does the confidence that Ss and Qs assign to what they see when paired with one another.

**Proposition 3** *Comparing the premium to defection for Ss at state  $(0, 1)$ ,  $b - \sigma s$ , with the warm-glow that Qs get from compliance,  $\alpha\varepsilon$ , we can distinguish two situations:*

1. *Ss do not imitate compliant Qs when  $\alpha\varepsilon \leq b - \sigma s$ . Then, three equilibria are possible:*

$$MEE_O \Leftrightarrow f_O(\varepsilon) \leq p^Q, \quad MEE_y \Leftrightarrow f_y(\varepsilon) < p^Q < f_O(\varepsilon), \quad MEE_{01} \Leftrightarrow p^Q \leq f_y(\varepsilon).$$

with

$$f_y(\varepsilon) \equiv \frac{\varepsilon - (b - \sigma s)}{(1 - \alpha)\varepsilon} < f_O(\varepsilon) \equiv \frac{\varepsilon - (d - \phi)}{\varepsilon(1 - \alpha)s}, \quad \forall s \in (0, 1), \quad (f_y(\varepsilon) = f_O(\varepsilon), \text{ if } s = 1).$$

2. *Qs do not imitate defecting Ss when  $\alpha\varepsilon > b - \sigma s$ . Then, three equilibria are possible:*

$$MEE_{01} \Leftrightarrow p^S \leq f_{01}(\varepsilon), \quad MEE_x \Leftrightarrow f_{01}(\varepsilon) < p^S < f_x(\varepsilon), \quad MEE_I \Leftrightarrow f_x(\varepsilon) \leq p^S.$$

with

$$f_{01}(\varepsilon) \equiv \frac{b - \sigma s}{\alpha\varepsilon} < f_x(\varepsilon) \equiv \frac{b}{\alpha\varepsilon(1 - s)}, \quad \forall s \in (0, 1), \quad (f_{01}(\varepsilon) = f_x(\varepsilon), \text{ if } s = 0).$$

**Proof.** See Appendix. ■

According to this proposition, if the warm-glow from compliance for Qs is small with respect to the premium to defection for Ss at state  $(0, 1)$ , then Ss will never imitate compliant Qs, regardless of the value of  $p^S$ . In that case, if Qs are reluctant to believe what they see,  $p^Q \in [0, f_y(\varepsilon)]$ , all Qs comply in  $MEE_{01}$ . If their level of confidence in what they see is moderate,  $p^Q \in [f_y(\varepsilon), f_O(\varepsilon)]$ , only some Qs comply in  $MEE_y$ . Finally, if they strongly believe in what they see,  $p^Q \geq f_O(\varepsilon)$ , all Qs imitate defecting Ss.

Conversely, if the warm-glow is strong enough, then Qs never defect and positive compliance rates within the population of Ss are possible, regardless of the value of  $p^Q$ . If Ss do not trust that they can get the same payoff from compliance observed in Qs,  $p^S \leq f_{01}(\varepsilon)$ , then even though all Qs comply, no S will imitate them. Conversely, if Ss are quite confident that they can get the payoff observed in Qs,  $f_{01}(\varepsilon) < p^S < f_x(\varepsilon)$ , some will imitate the compliant behavior of Qs. Finally, if their confidence in what they see is very high,  $f_x(\varepsilon) \leq p^S$ , all Ss will imitate compliance.

Proposition 3 is summarized in Figures 3 and 4. The first item in this Proposition refers to the case  $\alpha\varepsilon \leq b - \sigma s$ , depicted in the left panel of these Figures. The left panel in Figure

3 shows that the confidence that Qs assign to what they see,  $p^Q$ , facilitates a stable solution at which all or some Qs imitate defection, and makes more difficult a solution where all Qs comply, ignoring the non-compliant behavior of Ss. The left panel in Figure 4 shows that  $p^S$  has no influence on whether and to what extent Qs imitate Ss' defection. On the other hand, the right panel in Figures 3 and 4 displays the results collected in the second item of Proposition 3. Whenever the warm-glow is sufficiently strong to satisfy  $\alpha\varepsilon > b - \sigma s$ , Figure 4 shows that the more Ss trust what they see, the easier a long-run equilibrium at which they partially or completely imitate compliant Qs is. Logically, the confidence that Qs give to what they see has no effect on the compliance rate of Ss, as displayed in the right panel of Figure 3.

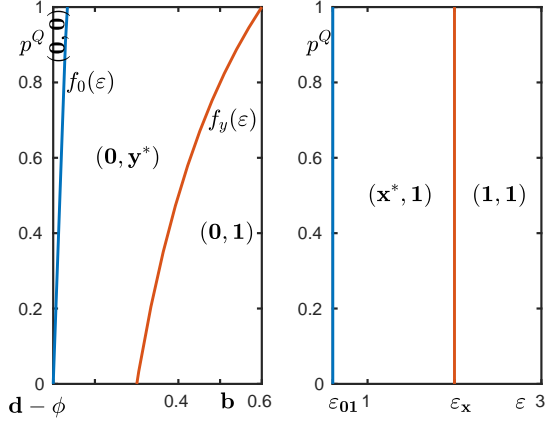


Figure 3: Type of equilibria in  $\varepsilon - p^Q$ .

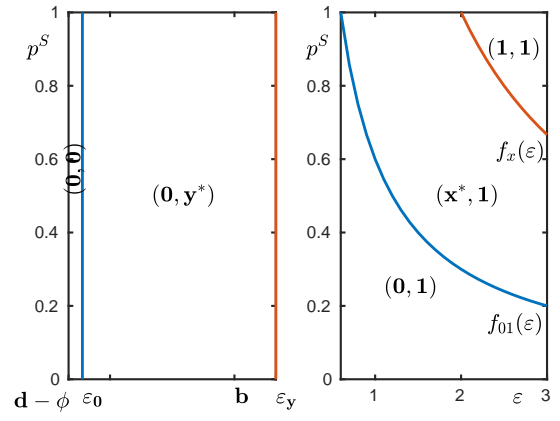


Figure 4: Type of equilibria in  $\varepsilon - p^S$ .

### 3.1 Compliance at each MEE

In what follows, we analyze how compliance reacts to the main parameters of the model. In particular, we focus on the confidence that Ss and Qs give to what they see,  $p^S$  and  $p^Q$ ; on how their preferences differ, either in absolute terms,  $\varepsilon$ , or in the bias that Qs have towards a large warm-glow from compliance ( $\alpha$  large) or a strong cold-prickle from defection ( $\alpha$  small); and finally, on the size of each of the two populations,  $s$ . With that aim, we plot the level curves of the global compliance rate,  $q^*$ , for different values of  $\alpha$  and  $s$  in the  $\varepsilon - p^Q$  space in Figure 5 and in the  $\varepsilon - p^S$  space in Figures 6 and 7. The expression of  $q^*$  for each of the five equilibria  $MEE_O$ ,  $MEE_y$ ,  $MEE_{01}$ ,  $MEE_x$ , and  $MEE_I$ , is given by  $0, y^*(1 - s), 1 - s, x^*s + 1 - s$  and  $1$ , respectively. The numerical illustration considers the following parameters values:

$$b = d = 0.5, \phi = 0.4, \quad (20)$$

with  $p^S = 0.5$  in Figure 5, and  $p^Q = 0.5$  in Figures 6 and 7.

The discrepancy in preferences between Ss and Qs (in absolute terms,  $\varepsilon$ , and in relative terms,  $\alpha$ ) determines the feasible regions for each equilibrium. The ordering of the different

equilibria, from lowest compliance,  $MEE_O$ , to highest compliance,  $MEE_I$ , as a function of  $\varepsilon$  in (19) is also clear in Figures 3-7. The more socially oriented Qs are,  $\varepsilon$ , the less feasible a stable equilibrium at which Qs imitate the non-compliant behavior of Ss is, and the more feasible an equilibrium at which Ss imitate the compliant behavior of Qs is.<sup>13</sup> On the other hand, the bias towards a higher warm-glow from compliance rather than a strong cold-prickle from defection,  $\alpha$ , has a twofold effect. As shown in Figure 5, a bias towards the cold-prickle (i.e. small  $\alpha$ ) makes easier the equilibria where Qs imitate defecting Ss, even if Qs are very different from Ss in absolute terms ( $MEE_O$  and  $MEE_y$ ). Conversely, Figure 6 shows that a higher bias towards the warm-glow (i.e. large  $\alpha$ ) widens the region where Ss partially or completely imitate compliant Qs ( $MEE_x$  and  $MEE_I$ ). The discrepancy between Ss and Qs also determines the actual share of compliance within the population of Qs in  $MEE_y$ , or within the population of Ss in  $MEE_x$ . As shown in Figures 5 and 6, compliance increases with the absolute distance between the two populations,  $\varepsilon$ , and the bias towards a large warm-glow,  $\alpha$ . The analytical proof is straightforward by computing the corresponding partial derivatives of expressions (12) and (13).

The type of equilibrium and the compliance rate also depend crucially on how much the agents trust what they see,  $p^S$ ,  $p^Q$ . Compliance in  $MEE_x$  increases with the confidence that Ss give to what they see when paired with Qs. Conversely, Qs will reduce compliance the more they trust in what they see when paired with Ss,  $p^Q$ , in  $MEE_y$ . The proof is straightforward by computing the corresponding partial derivatives of expressions (12) and (13). We do not present them here; instead we illustrate these interrelations with the help of Figures 5-7.

A bottom-up movement in Figures 3 and 5 corresponds to an increment in  $p^Q$ . If  $\varepsilon$  is sufficiently small (left panel in Figure 3 and Figure 5 left), as the confidence that Qs give to what they see runs from 0 to 1, it leads Qs from complete to partial and even to zero compliance. Moreover, within the  $MEE_y$  region, the greater  $p^Q$  is, the less the Qs comply. Conversely, the compliance decisions of Qs are independent of the confidence that the agents from the other population give to what they see,  $p^S$ , as can be observed in the left panel in Figure 4 and in Figures 6 and 7 for small  $\varepsilon$ . On the other hand, if  $\varepsilon$  is large, the right panel of Figure 4, as well as Figures 6 and 7, illustrate that starting in region  $MEE_{O1}$ , a higher confidence in what they see can trigger a compliant behavior for some or for all Ss. Moreover, within region  $MEE_x$ , the more strongly Ss trust that they can attain the same warm-glow they observe in Qs, the more will they imitate compliant Qs (moving to higher level curves).

Finally, the effect of a larger share of Ss on the global population (i.e. lower share of Qs) on compliance decisions is illustrated in Figure 7, moving from the left to the right graph. A higher  $s$  makes it more likely for revising Qs to be paired with defecting Ss, but at the same time, a greater share of Ss implies a lower global compliance rate, and hence a stronger payoff

---

<sup>13</sup>By more or less feasible we mean a wider or a narrower range of values for parameters  $p^S$ ,  $p^Q$ ,  $\alpha$ ,  $s$ .



to compliance relative to defection for Qs. All in all, the region where Qs imitate the defecting behavior of Ss shrinks. This is analytically obvious as  $\varepsilon_y$  decreases with  $s$ . Likewise, a higher share  $s$  makes it less likely that revising Ss are paired with compliant Qs, although it also makes compliance more attractive. The region where Ss imitate compliance widens. However, it is the region with partial compliance which widens greatly, while the region where all Ss comply narrows. Comparing corresponding points within the  $\varepsilon - p^S$  plane in Figures 7 left and right, one can conclude that the global compliance rate decreases with the share Ss in the global population.

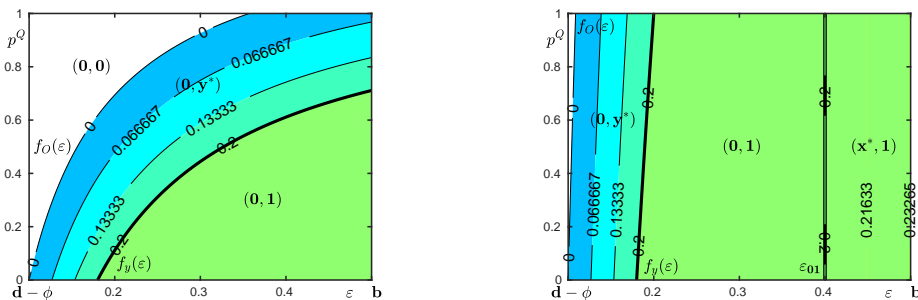


Figure 5:  $q^*$  level curves for  $s = 0.8$  and  $\alpha = 0.1$  (left);  $\alpha = 0.9$  (right)

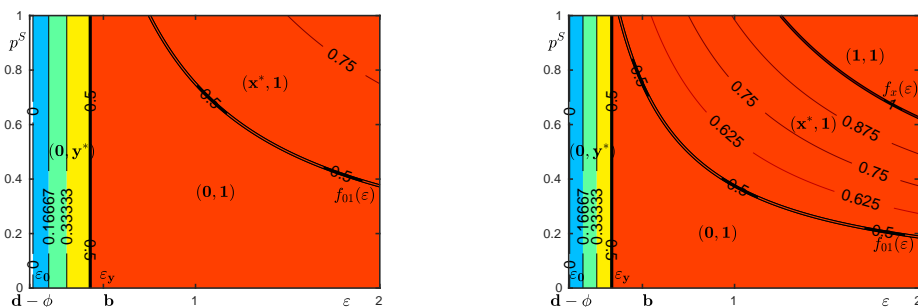


Figure 6:  $q^*$  level curves for  $s = 0.5$  and  $\alpha = 0.4$  (left);  $\alpha = 0.8$  (right)

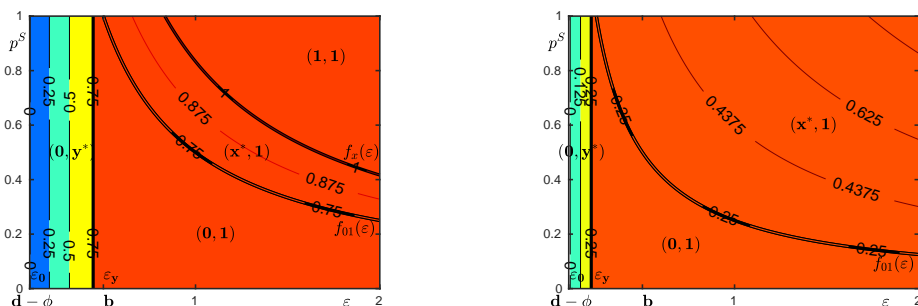


Figure 7:  $q^*$  level curves for  $\alpha = 0.8$  and  $s = 0.25$  (left);  $s = 0.75$  (right)

## 4 Comparison with the Nash Equilibrium

In this section, we first compute the Nash equilibrium of the two-population game. Since defection is the dominant strategy for Ss, their best reply function at any state  $\bar{x}$  is zero compliance:  $B^S(\bar{x}) = (x, 1 - x)^b = (0, 1)$ . Moreover, from the payoff function in (4), the best reply for Qs at state  $\bar{x}$  is

$$B^Q(\bar{x}) = (y, 1 - y)^b = \begin{cases} \left( \frac{\Delta - xs}{1 - s}, 1 - \frac{\Delta - xs}{1 - s} \right) & \text{if } \Delta < 1 - s + xs, \\ (1, 0) & \text{if } \Delta \geq 1 - s + xs. \end{cases}$$

with  $\Delta$  given in (3). From these reply functions the Nash equilibrium follows:

$$NE = (x^{NE}, y^{NE}) = \begin{cases} \left( 0, \frac{\Delta}{1 - s} \right) & \text{if } \varepsilon < b - \sigma s, \\ (0, 1) & \text{if } \varepsilon \geq b - \sigma s. \end{cases} \quad (21)$$

This equilibrium highlights a “snowdrift effect”: the greater the ratio of Ss in the overall population,  $s$ , whose dominant strategy is defection, the stronger the incentive to comply for Qs. Therefore, Qs comply more in the two-population game than in the case of a simple population of Qs, as given in (3).

The global compliance rate for the NE is given by

$$q^{NE} = x^{NE}s + y^{NE}(1 - s) = \begin{cases} \Delta & \text{if } \varepsilon < b - \sigma s, \\ 1 - s & \text{if } \varepsilon \geq b - \sigma s. \end{cases} \quad (22)$$

**Proposition 4** *The NE in (21) coincides with the ESS of the evolutionary game without inter-population imitation,  $p^S = p^Q = 0$ .*

**Proof.** See Appendix. ■

According to proposition 4, the NE in (21) corresponds to the ESS of the two-population evolutionary game with intra-population and inter-population interaction, but where imitation only occurs within but not between populations. This is the standard view in the literature, and it also arises in our formulation in the extreme case that agents fully ignore what they see when paired with someone different,  $p^S = p^Q = 0$ . In contrast, in the general case with  $p^S, p^Q > 0$ , the asymptotically stable equilibria  $MEE_e$  with  $e \in \{O, y, 01, x, I\}$  typically depart from the NE. The following proposition collects the comparison between the  $MEE$  in the general case when agents trust, at least partially, in what they see, and the NE (or ESS without inter-population imitation).

**Proposition 5** *The compliance rates in each of the  $MEE$  compare against the NE as:*

- *Condition (14):  $MEE_O = (0, 0)$  versus  $NE = (0, y^{NE})$ , with  $0 < y^{NE}$ .*

- Condition (15):  $MEE_y = (0, y^*)$  versus  $NE = (0, y^{NE})$ , with  $y^* < y^{NE} \leq 1$ .
- Condition (16):  $MEE_{01} = (0, 1) = NE$ .
- Condition (17):  $MEE_x = (x^*, 1)$ , with  $x^* > 0$  versus  $NE = (0, 1)$ .
- Condition (18):  $MEE_I = (1, 1)$  versus  $NE = (0, 1)$ .

Under conditions (14) or (15), the warm-glow from compliance for Qs is weak in comparison with the premium to defection for Ss. Ss do not comply and Qs comply below Nash, since some of them are tempted to imitate the highly rewarding defecting behavior observed in Ss,  $x^* = x^{NE} = 0$  while  $y^* < y^{NE}$ . Under the NE, even full compliance for Qs is possible. Condition (16) guarantees  $MEE_{01}$  with full compliance among Qs, although can not induce compliance among Ss, just as in the NE. Finally, under conditions (17) or (18) the warm-glow that defecting Ss conjecture they will get if they switch to compliance is strong enough to induce some of them (in  $MEE_x$ ) or all of them (in the  $MEE_I$ ) to imitate compliance, contrary to the zero compliance under the NE.

Depending on whether Qs imitate the defecting behavior of Ss, or Ss imitate the compliant behavior of Qs, Proposition 5 states that, in the  $MEE$ , Qs can comply below or Ss can comply above their NE. The next question is: how do these differences in compliance rates translate into differences in social welfare (or the average payoff for each population)?

The social welfare for the social state  $\bar{x}$  is computed as

$$\pi_e^S = \pi_C^S(q)x + \pi_D^S(q)(1-x), \quad \pi_e^Q = \pi_C^Q(q)y + \pi_D^Q(q)(1-y), \quad e \in \{O, y, 01, x, I, NE\}. \quad (23)$$

To compare the social welfares we distinguish three situations: The equilibria where Qs partially or completely imitate defecting Ss are analyzed in Proposition 6. At equilibrium  $MEE_{01}$  Ss do not imitate Qs and neither do Qs imitate Ss. Hence, this equilibrium coincides with the NE and so do populations' welfares. Finally, Proposition 7 analyzes the two equilibria where Ss partially or completely imitate compliant Qs.

**Proposition 6** *Ss do not comply and Qs imitate defecting Ss in  $MEE_O$  and in  $MEE_y$ . Then, for each population, the social welfare compares to the NE as*

$$\pi_O^h < \pi_{NE}^h, \quad \pi_y^h < \pi_{NE}^h, \quad \text{and} \quad \frac{d(\pi_y^h - \pi_{NE}^h)}{dp^Q} < 0, \quad \forall h \in \{S, Q\}.$$

**Proof.** See the Appendix. ■

Under conditions (14) and (15) defection is so rewarding that Qs imitate the defecting behavior of Ss and comply below the NE in  $MEE_y$ , or do not comply at all in  $MEE_O$ . In either equilibria, Ss are worse off with a lower global compliance rate and less free-riding. For Qs, a lower compliance has a two-fold effect: lower costs from compliance, and also lower

benefits from a smaller global compliance rate,  $q$ . Proposition 7 proves that the second effect is stronger and Qs are also worse than under Nash.

Under  $MEE_O$  the payoffs in either population are unaffected by  $p^Q$ . In contrast, when  $MEE_y$  is the stable equilibrium, the rate of compliance within the population of Qs,  $y^*$ , negatively depends on  $p^Q$ . The greater this value, the more Qs imitate defecting Ss, and hence the less they comply. A lower  $y^*$  below the NE,  $y^{NE}$ , widens the negative gap  $\pi_x^h - \pi_{NE}^h$  both for Ss and for Qs. Thus, the more inclined Qs are to imitate others, the less the average payoff obtained by Ss and Qs.

**Proposition 7** *All Qs comply and Ss imitate compliance in  $MEE_x$  and in  $MEE_I$ . The social welfare in these equilibria compare to the NE as*

1. For Qs,

$$\pi_1^Q > \pi_x^Q > \pi_{NE}^Q \quad \text{and} \quad \frac{d(\pi_x^Q - \pi_{NE}^Q)}{dp^S} > 0.$$

2. For Ss,

$$\pi_x^S \geq \pi_{NE}^S \Leftrightarrow s \geq \hat{s}; \quad \pi_1^S > \pi_{NE}^S \Leftrightarrow s \geq \frac{b}{b + \phi}. \quad (24)$$

Moreover,

$$\frac{d(\pi_x^S - \pi_{NE}^S)}{dp^S} \geq 0 \Leftrightarrow p^S \leq \hat{p}^S.$$

with  $\hat{s} = \hat{s}(\alpha, \varepsilon, p^S)$  and  $\hat{p}^S = \hat{p}^S(\alpha, \varepsilon, s)$  given in the proof, in the Appendix.

**Proof.** See the Appendix. ■

When the warm-glow from compliance is highly attractive to guarantee full compliance among Qs, and Ss trust sufficiently in what they see then, some or even all Ss are induced to comply in the  $MEE_x$  or the  $MEE_I$ . Without inter-population imitation, this incentive disappears and Ss never comply, while the strong warm-glow still induces full compliance for Qs,  $(x^{NE}, y^{NE}) = (0, 1)$ . Thus, Qs bear the same costs from compliance with or without inter-population imitation. In contrast, they benefit from a higher global compliance rate in the  $MEE_x$  and even higher in the  $MEE_I$ . In consequence,  $\pi_1^Q > \pi_x^Q > \pi_{NE}^Q$ .

The comparison of the social welfare for Ss when they imitate compliant Qs is more cumbersome. When moving from the NE to the  $MEE_x$  or  $MEE_I$ ,  $x$  rises from 0 to  $x^* > 0$  or 1. This has a positive effect from a higher global compliance rate, implying greater payoffs for both compliance and defection. However, it also has a negative effect associated with those  $x^*$  Ss who “mistakenly” imitate compliance, and suffer a reduction in their payoffs as they move from (highly rewarding) defection to (less rewarding) compliance. Which of these two effects prevails the other depends on the parameter values and, crucially, on each population’s size, defined by  $s$ . Note that condition (17) for the  $MEE_x$  can be written as  $s > \underline{s}$  ( $\underline{s}$  given in the Appendix). Above this bound, inter-population imitation reduces Ss welfare, for  $s \in (\underline{s}, \hat{s})$ ,

but increases Ss welfare for  $s \in (\hat{s}, 1)$ . To understand the role played by the number of Ss, note that if this is small, a positive  $x$  has little impact on global compliance,  $q$ , and its positive effect on the average payoff for Ss is moderate. Moreover, because Qs are relatively abundant ( $s$  small), and since all of them comply, the global compliance rate,  $q$ , is large, which is associated with a large premium to defection for Ss,  $\pi_D^S - \pi_C^S$ , as shown in (4). Thus, for a relatively small  $s$ , the negative effect surpasses the positive effect and Ss are worse off in the  $MEE_x$ . Conversely, when Ss are relatively abundant, the positive effect overcomes the negative one. Similar reasoning applies when comparing  $MEE_I$  to NE (where the bound  $\hat{s}$  is now replaced by the constant  $b/(b + \phi)$ ).

The comparison of the social welfare of Ss in the  $MEE_x$  versus the NE is presented in Figures 8-11. Figure 8 corresponds to  $\varepsilon = 0.3$ , where Qs have a small warm-glow  $d - \phi < \varepsilon < b$  ( $0.1 < 0.3 < 0.5$ ); Figure 9 considers a higher  $\varepsilon$ ,  $d - \phi < \varepsilon < b$  ( $0.1 < 0.49 < 0.5$ ), with a small free-riding incentive, close to 0; Figure 10 considers a very large  $\varepsilon$ ,  $d - \phi < b < \varepsilon$  ( $0.1 < 0.5 < 2$ ), and compliance is a dominant strategy for Qs. In the left panel these figures depict the level curves of  $\pi_x^S - \pi_{NE}^S$  for  $(p^S, s) \in (0, 1) \times (0, 1)$ . The shaded area in the left panel of Figures 8-10 is the region where  $MEE_x$  is the stable solution. Two different regions are separated by the white  $s = \hat{s}$  line, where  $\pi_x^S - \pi_{NE}^S = 0$ . Below this line Ss are relatively scarce in the global population and the negative level curves represents a worse average welfare for Ss,  $\pi_x^S < \pi_{NE}^S$ . Above this line the situation reverses and  $\pi_x^S > \pi_{NE}^S$ . Similar insights can be gained from Figure 11 where the  $MEE_I$  is compared to the NE (the  $s = \hat{s}$  line is replaced by  $s = b/(b + \phi)$ ).

Interestingly, on comparing Figures 8 to 10 we observe that, as the gap in preferences,  $\varepsilon$ , narrows, so the region for stable solution  $MEE_x$  narrows. And as this region shrinks, the lower part where Ss experience welfare losses narrows more rapidly. In fact, for the parameters considered in Figure 8, the lower region has disappeared. Thus, if the gap in preferences is low, is it less obvious that Ss imitate compliance:  $s$  and  $p^S$  need to be larger. However, when the  $MEE_x$  applies, it is more likely that it provides greater social welfare to Ss than the NE.<sup>14</sup>

A second important finding displayed by Figure 9 relates to the level of confidence that Ss give to what they see,  $p^S$ . From condition (17) the greater this level of confidence, the wider the interval  $(\underline{s}, 1)$  where equilibria  $MEE_x$  or  $MEE_I$  occur.<sup>15</sup> More importantly, if a non-empty interval  $(\underline{s}, \hat{s})$  exists, then as  $p^S$  increases, this interval where Ss are worse off in the  $MEE_x$  increases, while the interval  $(\hat{s}, 1)$ , where they are better off, decreases.<sup>16</sup> This indicates that the more Ss trust in what they see, the more likely it is that some of them in  $MEE_x$

<sup>14</sup>This does not contradict Proposition 7.2. The interval  $(\underline{s}, \hat{s})$  with lower welfare for Ss disappears when  $\hat{s}$  lies below  $\underline{s}$ .

<sup>15</sup>This is also displayed if the two left panels of Figures 10 and 11 are overlaid.

<sup>16</sup>From the definition of  $\underline{s}$  in the proof of Proposition 7 and  $\hat{s}$  in (36) it is not difficult to prove that the increment in subinterval  $(\underline{s}, \hat{s})$  is more pronounced than the increment in subinterval  $(\hat{s}, 1)$ , which could even decrease.

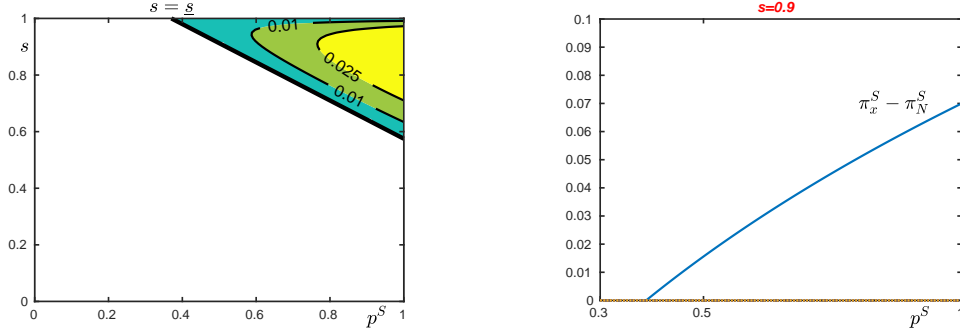


Figure 8: Level curves of  $\pi_x^S - \pi_{NE}^S$  (left) and  $\pi_x^S - \pi_{NE}^S$  for  $s = .9$  (right).

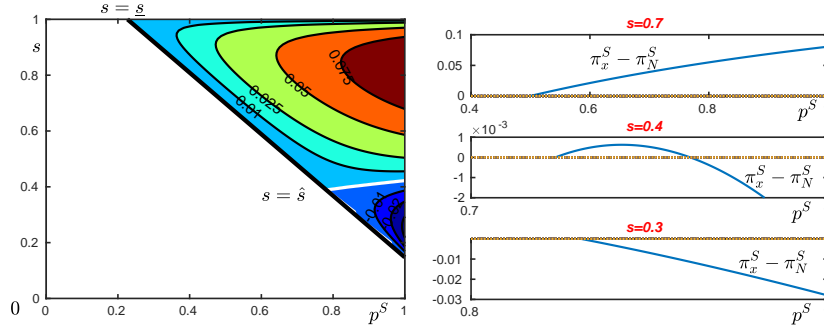


Figure 9: Level curves of  $\pi_x^S - \pi_{NE}^S$  (left) and  $\pi_x^S - \pi_{NE}^S$  for  $s = .3, .4, .7$  (right).

or all of them in  $MEE_I$  comply, and the greater the compliance rate within this population is. However, they are not necessarily better off. In fact, it becomes more likely that they are worse off than under the NE with zero compliance.

The effect of  $p^S$  crucially depends on the share of selfish agents in the global population. If  $s$  is large and the stable equilibrium is  $MEE_x$ , Ss are better off when some of them comply. Therefore, the more Ss believe in what they see, the more they will imitate compliance and the higher  $\pi_x^S$  grows above  $\pi_{NE}^S$ . See Figure 8 (right) and the upper graph in Figures 9 and 10 (right). If, conversely, the ratio of Ss is very small, Ss are worse off when some of them comply. In consequence, a higher  $p_s$ , which increases the number of compliant Ss, worsens the situation, widening the negative gap with the NE (see the lower graph in Figures 9 and 10 right). Finally and interestingly, if  $s$  takes a moderate value, as  $p_s$  grows the gap  $\pi_x^S - \pi_{NE}^S$  can move from positive to negative. This is shown in the middle graph in Figures 9 and 10 (right): as more Ss imitate compliance, the social welfare initially increases above the NE to reach a maximum, and decreases thenceforth and turns negative at some point.

Note finally that the confidence that Ss give to what they see,  $p^S$ , can determine whether  $MEE_x$  or  $MEE_I$  is the long-run equilibrium (i.e. whether Ss comply partially or completely). However, once in the  $MEE_I$ , since compliance reaches its maximum possible rate and it remains equal to one regardless of the value of  $p^S$ , it does not determine the sign or the size of

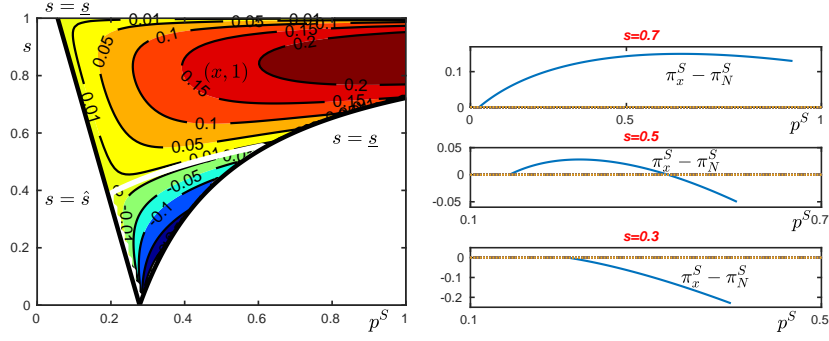


Figure 10: Level curves of  $\pi_x^S - \pi_{NE}^S$  f(left) and  $\pi_x^S - \pi_{NE}^S$  for  $s = .3, .5, .7$  (right).

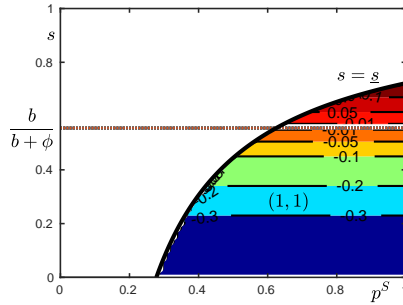


Figure 11: Level curves of  $\pi_1^S - \pi_{NE}^S$ .

the gap in social welfare, which only depends on whether  $s$  is above or below  $b/(b + \phi)$ .

## 5 Conclusions

In this paper we have analyzed a collective action problem involving heterogeneous agents, distinguishing two populations with different preferences. In a two-population evolutionary game we assume inter and intra-population interaction, implying that the welfare in one population depends on the compliance decisions made by the individuals within and in the other population. Our most interesting contribution is to remove the constraint of agents who only imitate their own kind. Our social life is not only restricted to other individuals sharing our exact same preferences. In fact, this is seldom the case, and more commonly we mix with people whose values and motives differ from ours.

Thus, we focus on the situations in which a “revising” individual is paired with someone with different preferences and who is acting differently (i.e. chooses a different strategy). In standard evolutionary game theory the revising agent would recognize his partner as belonging to a different “species” and would ignore him. However, when we do not refer to different natural species, but to two individuals with different preferences, then, for the revising agent it is more difficult to discern whether this alternative action/strategy will provide the expected

payoff before this encounter, or the observed payoff in his partner. Under these circumstances, the agent conjectures the payoff to the new strategy as the combination between what he observes his partner is getting and what he expected given the current populations states. If this conjecture exceeds what he actually gets with his current strategy, he will be willing to imitate the observed strategy.

This is applied to a collective action problem involving two type of agents. Standard *homo oeconomicus*, selfish or perfectly rational Sanchos; and the norm-using or more socially-oriented Quixotes. Depending on how different their preferences are and depending on the weights that the conjecture process gives to expectations and to observations, the unique stable equilibrium of the evolutionary dynamics can be of different types. Norm-using Quixotes might imitate (partially or completely) the non-compliant behavior of selfish Sanchos; conversely, selfish Sanchos might find it attractive to imitate (partially or completely) compliant Quixotes; finally, there exists an equilibrium at which neither defecting Sanchos imitate compliant Quixotes and vice versa. The equilibria in which agents from one population imitate agents from the alternative population differ from the NE (which coincides with the ESS equilibrium without inter-population imitation). When comparing the different equilibria against the NE, we can highlight the following results.

Whenever individuals in one population comply above (resp. below) the NE, the agents in the other population benefit correspondingly (resp. suffer). In the equilibria in which Quixotes comply less because they imitate non-compliant Sanchos, they are worse off. However, in the equilibria in which Sanchos comply more because they imitate compliant Quixotes, their social welfare can increase or decrease.

A wider gap in preferences between selfish Sanchos and socially-oriented Quixotes widens the region (in the parameter space) where Sanchos imitate the compliant behavior of Quixotes. However, the more socially-oriented Quixotes are, the greater is the region where (partial or complete) compliance by Sanchos reduces their welfare in relative terms to the region where they are better off. In consequence, if Quixotes are very socially-oriented, there are good chances that some Sanchos will imitate compliance. However, it is also likely that they experience an average loss rather than a gain in welfare.

In the region (in the parameter space) where the stable equilibrium is compatible with compliant Sanchos, the more they trust in what they see, the more they will imitate compliance. Three situations can occur. If the ratio of selfish individuals in the overall population is large, they are better off if they imitate compliance than if they do not (under the NE). Conversely, they become worse off if the ratio of Sanchos is small. Therefore, a greater confidence in what they see benefits Sanchos when they are many, but harms them when they are few. Finally, if the ratios of Sanchos and Quixotes are similar, a higher confidence in what they see could lead Sanchos from a better off to a worse off situation with compliance, hence, increasing initially,



and later reducing, the social welfare of Sanchos, as a larger number of them comply.

## References

- [1] Alger, I., Weibull, J.W., 2013. Homo Moralis-Preference Evolution under Incomplete Information and Assortative Matching. *Econometrica* 81, 2269-2302. doi:10.3982/ECTA10637
- [2] Andreoni, J., 1988. Why Free Ride? *Journal of Public Economics* 37, 291-304. doi:10.1016/0047-2727(88)90043-6.
- [3] Andreoni, J., 1990. Impure altruism and donations to public goods: a theory of warm-glow giving. *The Economic Journal* 100, 464-477. doi:10.2307/2234133.
- [4] Andreoni, J., 1995. Warm-Glow versus Cold-Prickle: The effects of positive and negative framing on Cooperation in experiments. *The Quarterly Journal of Economics* 110, 1-21.
- [5] Antoci, A., Borghesi, S., Russu, P., 2012. Environmental protection mechanisms and technological dynamics. *Economic Modelling* 29, 840-847. doi:10.1016/j.econmod.2011.10.004.
- [6] Antoci, A., Galeotti, M., Radi, D., 2011. Financial tools for the abatement of traffic congestion: a dynamical analysis. *Computational Economics* 38, 389-405. doi:10.1007/s10614-011-9294-7.
- [7] Antoci, A., Dei, R., Galeotti, M., 2009. Financing the adoption of environment preserving technologies via innovative financial instruments: an evolutionary game approach. *Nonlinear Analysis: Theory, Methods & Applications* 71, 952-959. doi: 10.1016/j.na.2009.01.077.
- [8] Apesteguia, J., Huck, S., Oechssler, J., Weidenholzer, S., 2010. Imitation and the evolution of Walrasian behavior: Theoretically fragile but behaviorally robust. *Journal of Economic Theory* 145, 1603-1617. doi: 10.1016/j.jet.2010.02.014.
- [9] Bontems, P., Rotillon, G., 2000. Honesty in environmental compliance games. *European Journal of Law and Economics* 10, 31-41. doi:10.1023/A:1018786721348.
- [10] Cabo, F., García-González, A., 2019. Interaction and imitation in a world of Quixotes and Sanchos. *Journal of Evolutionary Economics* 29, 1037-1057. doi: 10.1007/s00191-019-00620-3.
- [11] Cabo, F., García-González A., Molpeceres-Abella, M., 2020. Compliance with social norms as an evolutionary stable equilibrium. In: Pineau, P.O., Sigué, S., Taboubi, S. (Eds.). *Games in Management Science. Essays in Honor of Georges Zaccour*. Switzerland: Springer Nature, 283-313.

- [12] Cui, Z., Wang, R., 2016. Collaboration in networks with randomly chosen agents. *Journal of Economic Behavior & Organization* 129, 129-141. doi:10.1016/j.jebo.2016.06.015.
- [13] de Young, R., 1996. Some psychological aspects of reduced consumption behavior: the role of intrinsic satisfaction and competence motivation. *Environment and Behavior* 28, 358-409. doi:10.1177/0013916596283005.
- [14] Doebeli, M., Hauert, C., Killingback, T., 2004. The evolutionary origin of cooperators and defectors. *Science* 5697, 859-862. doi: 10.1126/science.1101456.
- [15] Eshel, I., Samuelson, L., Shaked, A., 1998. Altruists, egoists, and hooligans in a local interaction model. *The American Economic Review* 88, 157-79. doi: www.jstor.org/stable/116823.
- [16] Di Giovinazzo, V., Naimzada, A., 2015. A model of fashion: endogenous preferences in social interaction. *Economic Modelling* 47, 12-17.
- [17] Gokhale, C.S, Freat, M., Rainey, P.B., 2019. Eco-evolutionary dynamics of mutualisms. Mimeo.
- [18] Lulu, G., Gao, J., Cao, M., 2018. Evolutionary game dynamics for two interacting populations under environmental feedback. arXiv.org:1806.03194.
- [19] Grafton, R.Q., Kompas, T. VanLong, Ngo. 2017. A brave new World? Kantian–Nashian interaction and the dynamics of global climate change mitigation. *European Economic Review* 99, 31-42. doi: https://doi.org/10.1016/j.euroecorev.2017.04.002.
- [20] Güth, W., 1995. An evolutionary approach to explaining cooperative behavior by reciprocal incentives. *International Journal of Game Theory* 24, 323-344. doi:10.1007/BF01243036
- [21] Güth, W., Yaari, M.E., 1992. Explaining reciprocal behavior in simple strategic games: an evolutionary approach. In: Witt, U. (Eds.). *Explaining Process and Change: Approaches to Evolutionary Economics*. University of Michigan Press, Ann Arbor.
- [22] Hofbauer, J., Sigmund, K., 1998. *Evolutionary Games and Population Dynamics*. Cambridge: Cambridge University Press.
- [23] Kandori, M., Mailath, G.J., Rob, R., 1993. Learning, mutation, and long run equilibria in games. *Econometrica* 61, 29-56. doi:10.2307/2951777.
- [24] Khan, A., 2014. Coordination under global random interaction and local imitation. *International Journal of Game Theory* 43, 721-745. doi: 10.1007/s00182-013-0399-1.

- [25] Matsuyama, K., 1991. Custom versus fashion: path-dependence and limit cycles in a random matching game. Discussion paper 83, Institute for Empirical Macroeconomics. Federal Reserve Bank of Minneapolis.
- [26] Miller, J.H., Andreoni, J., 1991. Can evolutionary dynamics explain free riding in experiments? *Economics Letters* 36, 9-15. doi:10.1016/0165-1765(91)90047-O.
- [27] Naimzada, A.K., Pireddu, M., 2018. Fashion cycle dynamics in a model with endogenous discrete evolution of heterogeneous preferences. *Chaos* 28, 055907. <https://doi.org/10.1063/1.5024931>.
- [28] Olson, M., 1965. *The Logic of Collective Action: Public Goods and the Theory of Groups*. Harvard University Press.
- [29] Radi, D., Gardini, L. 2018. A piecewise smooth model of evolutionary game for residential mobility and segregation. *Chaos* 28, 055912. doi: 10.1063/1.5023604.
- [30] Robson, A.J., Vega-Redondo, F., 1996. Efficient equilibrium selection in Evolutionary Games with random matching. *Journal of Economic Theory* 70, 65-92. doi: 10.1006/jeth.1996.0076.
- [31] Ostrom, E., 2000. Collective action and the evolution of social norms. *The Journal of Economic Perspectives* 14, 137-158. doi:10.1257/jep.14.3.137
- [32] Samuelson, L., Zhang, J., 1992. Evolutionary stability in asymmetric games. *Journal of Economic Theory* 57, 363-391. [https://doi.org/10.1016/0022-0531\(92\)90041-F](https://doi.org/10.1016/0022-0531(92)90041-F).
- [33] Sandholm, W.H., 2010. *Population Games and Evolutionary Dynamics*. Cambridge: MIT Press.
- [34] Selten, R., 1980. A note on evolutionarily stable strategies in asymmetric animal conflicts. *Journal of Theoretical Biology* 84, 93-101. [https://doi.org/10.1016/S0022-5193\(80\)81038-1](https://doi.org/10.1016/S0022-5193(80)81038-1).
- [35] Schlag, K.H., 1998. Why imitate, and if so, how?: A boundedly rational approach to multi-armed bandits. *Journal of Economic Theory* 78, 130-156. doi: 10.1006/jeth.1997.2347.

## Appendix

### Proof of Proposition 1.

The conditional imitation rates when an  $h$ -type agent is paired with a  $-h$ -type agent, immediately follow from (7) and (6):

$$R_C^S(q) = [E_C^S(q) - \pi_D^S(q)]_+ = \sigma [U^S - q]_+, \quad (25)$$

$$R_D^S(q) = [E_D^S(q) - \pi_C^S(q)]_+ = \sigma [q - L^S]_+, \quad (26)$$

$$R_C^Q(q) = [E_C^Q(q) - \pi_D^Q(q)]_+ = \sigma [U^Q - q]_+, \quad (27)$$

$$R_D^Q(q) = [E_D^Q(q) - \pi_C^Q(q)]_+ = \sigma [q - L^Q]_+. \quad (28)$$

with

$$U^S = \frac{\varepsilon \alpha p^S - (d - \phi)}{\sigma}, \quad L^S = \frac{(1 - \alpha) \varepsilon p^S - (d - \phi)}{\sigma}, \quad (29)$$

$$U^Q = \frac{\varepsilon(1 - \alpha p^Q) - (d - \phi)}{\sigma}, \quad L^Q = \frac{\varepsilon[1 - (1 - \alpha)p^Q] - (d - \phi)}{\sigma}. \quad (30)$$

The first and last items in the Proposition trivially hold from (25)-(28). To prove the second item, note that the bounds in (29)-(30) can be compared attending to the values of  $\alpha$ ,  $p^S$  and  $p^Q$ . In particular:

- $\alpha > \frac{1}{2}$

Revising S:  $L^S < U^S$  and  $R_C^S(q), R_D^S(q) > 0, \forall q \in (L^S, U^S)$ .

Revising Q:  $U^Q < L^Q$  and  $R_C^Q(q) = R_D^Q(q) = 0, \forall q \in (U^Q, L^Q)$ .

- $\alpha < \frac{1}{2}$

Revising S:  $U^S < L^S$  and  $R_C^S(q) = R_D^S(q) = 0, \forall q \in (U^S, L^S)$ .

Revising Q:  $L^Q < U^Q$  and  $R_C^Q(q), R_D^Q(q) > 0, \forall q \in (L^Q, U^Q)$ .

In the particular case  $\alpha = 1/2$ ,  $L^h = U^h$  and therefore,  $R_{-i}^h > 0 \Rightarrow R_i^h = 0$ , for all  $h \in \{S, Q\}$  and all  $i \in \{C, D\}$ . ■

### Proof of Proposition 2.

The ordering of the upper and lower bounds  $U^h, L^h$  determines different regions for  $q \in [0, 1]$  within which the system evolves differently. Depending on parameters values, we distinguish four different scenarios:

- $\alpha > \frac{1}{2}$  and  $p^S + p^Q > \frac{1}{\alpha}$ , with  $L^S < U^Q < U^S < L^Q$ , as shown in Figure 12 up.
- $\alpha > \frac{1}{2}$  and  $p^S + p^Q < \frac{1}{\alpha}$ , with  $L^S < U^S < U^Q < L^Q$ , as shown in Figure 12 down.
- $\alpha < \frac{1}{2}$  and  $p^S + p^Q > \frac{1}{1-\alpha}$ , with  $U^S < L^Q < L^S < U^Q$ , as shown in Figure 13 up.

iv)  $\alpha < \frac{1}{2}$  and  $p^S + p^Q < \frac{1}{1-\alpha}$ , with  $U^S < L^S < L^Q < U^Q$ , as shown in Figure 13 down.

Figures 12 and 13 help us summarize the different behavior of the evolutionary dynamics depending on  $q$  and parameters values. First recall that  $\Delta \geq \max_h\{U^h, L^h\}$ . Thus, if  $\Delta < q < 1$  compliant Ss and Qs will switch to defection whenever paired with defecting agents from within or from the other population. Because the number of compliant Ss and Qs decreases, no equilibrium with positive compliance is feasible in this region.

It always holds that  $r_C^S = 0$  ( $r_D^S > 0$ ) (because defection is a dominant strategy in this population) and  $R_C^S = 0$  if  $q > U^S$  (encircled in Figures 12 and 13). Above this bound,  $x$  undoubtedly decreases. Therefore, an equilibrium above this bound is only possible with zero compliance within Ss. This is compatible with Qs showing zero compliance in  $MEE_O = (0, 0)$ , partial compliance in  $MEE_y = (0, y^*)$ , or complete compliance in  $MEE_{01} = (0, 1)$ .

Below  $\Delta$ ,  $r_D^Q = 0$  ( $r_C^Q > 0$ ) moreover,  $R_D^Q = 0$  if  $q < L^Q$  (squared in Figures 12 and 13). In consequence, below this bound  $y$  increases, and an equilibrium is only possible with full compliance within this population. This is compatible with a population of Ss with zero compliance, in  $MEE_{01} = (0, 1)$ , partial compliance, in  $MEE_x = (x^*, 1)$ , or full compliance in  $MEE_I = (1, 1)$ .

The  $MEE_O$  can occur in Figures 12 and 13, in every subinterval where  $(0, y^*)$  is possible, provided that  $L^Q < 0$ . Similarly,  $MEE_I$  can occur in every subinterval where  $(x^*, 1)$  is possible, provided that  $U^S$  is sufficiently large above 1 (see Condition for equilibrium  $(1, 1)$  later in this section).

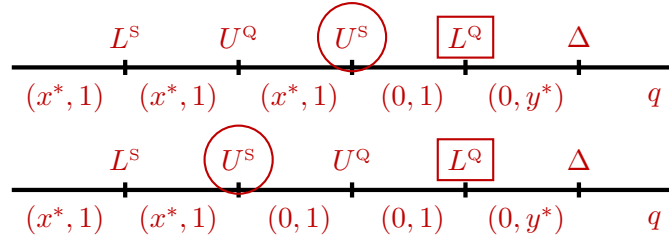


Figure 12: Regions for  $\alpha > \frac{1}{2}$ :  $p^S + p^Q > \frac{1}{1-\alpha}$  (up);  $p^S + p^Q < \frac{1}{1-\alpha}$  (down)

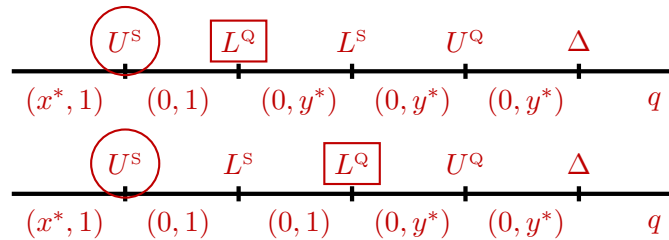


Figure 13: Regions for  $\alpha < \frac{1}{2}$ :  $p^S + p^Q > \frac{1}{1-\alpha}$  (up);  $p^S + p^Q < \frac{1}{1-\alpha}$  (down)

For each of the four scenarios *i*), *ii*), *iii*), *iv*), one can analytically compute the five equilibria  $MEE_O$ ,  $MEE_y$ ,  $MEE_{01}$ ,  $MEE_x$  and  $MEE_I$ . Within each interval it is easy although tedious to compute the system of differential equations defined in (8) and (9), together with (5), (7) and (10)-(11). Given this system, if  $\dot{y} > 0$  within a particular interval, then the only possible stable equilibrium requires  $y^* = 1$ . Taking this into account, one can solve the  $\dot{x} = 0$  equation, getting  $x^*$  in (13) or  $x^* = 1$ . Following similar reasoning, if conversely  $\dot{x} < 0$ , then the only possible stable equilibrium should satisfy  $x^* = 0$ , and plugging this in  $\dot{y} = 0$  the value of  $y^*$  in (13) or  $y^* = 0$  follows. Finally, in the intervals where  $\dot{y} > 0$  and  $\dot{x} < 0$  the only feasible equilibrium is  $(0, 1)$ .

The conditions which characterize each equilibrium type can be derived as follows:

- Condition for equilibrium  $(0, 0)$ .

At state  $(0, y)$ , the system dynamics reads:

$$\begin{aligned}\dot{x} &= y(1-s)R_C^S, \\ \dot{y} &= (1-y)y(1-s)r_C^Q - y\{(1-y)(1-s)r_D^Q + sR_D^Q\}.\end{aligned}$$

Then,  $\dot{x} = 0$  for any  $y > 0$  if and only if  $R_C^S = 0$ , i.e.  $U^S < 0$  or equivalently  $\alpha\varepsilon p^S < d - \phi$ .

Furthermore, we need to prove that  $\dot{y}/y < 0$  for any  $y > 0$ . This is straightforward if  $r_C^Q = 0$  (i.e.  $r_D^Q > 0$ ). If, conversely  $r_C^Q > 0$  (i.e.  $r_D^Q = 0$ ), from (25)-(30), after some simplification one gets:

$$\frac{\dot{y}}{y} = [\varepsilon - (d - \phi) - y(1-s)\sigma][1 - y(1-s)] - s(1-\alpha)p^Q\varepsilon.$$

And

$$\lim_{y \rightarrow 0} \frac{\dot{y}}{y} = \varepsilon - (d - \phi) - s(1-\alpha)p^Q\varepsilon.$$

Thus,  $\lim_{y \rightarrow 0} \dot{y}/y < 0$  if and only if  $\varepsilon(1 - (1-\alpha)p^Qs) \leq b - \sigma$  as stated in condition (14).

- Condition for equilibrium  $(0, y^*)$ .

At state  $(0, 1)$  the system moves to equilibrium  $(0, y^*)$  if some compliant Qs switch to defection. From (28), at state  $(0, 1)$ ,  $\dot{y} = -sR_D^Q(0s + (1-s)) < 0$  if and only if  $1 - s - L^Q > 0$ . And this is equivalent to condition (15) left. Moreover, to guarantee that the rate of compliance among Qs does not decrease leading the system to  $(0, 0)$ , the opposite to condition (14) is required in condition (15) right.

- Condition for equilibrium  $(0, 1)$ .

The system will remain at state  $(0, 1)$  if neither a revising compliant Q switches to defection, nor a revising defecting S moves to compliance. This is equivalent to impose  $R_D^Q = R_C^S = 0$ . At this point this is equivalent to  $1 - s - L^Q \leq 0$  and  $U^S - (1-s) \leq 0$ . These are the two conditions in (16).

- Condition for equilibrium  $(x^*, 1)$ .

At state  $(0, 1)$  the system will evolve towards  $(x^*, 1)$ , with  $x^* > 0$ , if defecting Ss are willing to switch to compliance when paired with compliant Qs. This is equivalent to  $R_C^S > 0$ , or condition  $U^S - (1 - s) > 0$ , which can be rewritten as in (17).

- Condition for equilibrium  $(1, 1)$ .

At state  $(x, 1)$ , the dynamics for  $x$  reads:

$$\dot{x} = (1 - x)(1 - s)R_C^S - x(1 - x)sr_D^S.$$

And  $(1 - s)R_C^S - xsr_D^S > 0$  for  $x$  tending to one if:

$$(1 - s)\sigma(U^S - 1) - sb > 0,$$

which leads to condition (18).

Note that condition (17) leads to  $x^* > 0$ , while condition (18) implies to  $x^* \geq 1$ , and since  $x^*$  cannot be greater, it must be equal to 1.

Under condition (2) it is easy to verify that  $y^* < 1$  and  $y^* > 0$  are equivalent to the two conditions in (15). Conversely,  $y^* \leq 0$  in (12) gives condition (14), which characterized null compliance among Qs. ■

### Proof of Proposition 3.

Because  $p^S \in [0, 1]$ ,  $\alpha\varepsilon \leq b - \sigma s$  immediately implies  $\varepsilon \leq \varepsilon_{01}$ , hence, only  $MEE_O$ ,  $MEE_y$ , and  $MEE_{01}$  are feasible. The conditions in Proposition 3.1 immediately follows from (14)-(16). Likewise, since  $p^Q \in [0, 1]$ , then  $\varepsilon(1 - (1 - \alpha)p^Q) \geq \alpha\varepsilon$  and condition  $\alpha\varepsilon > b - \sigma s$  implies  $\varepsilon > \varepsilon_y$  and hence, only  $MEE_{01}$ ,  $MEE_x$ , and  $MEE_I$  are feasible. The conditions in Proposition 3.2 immediately follow from (16)-(18). ■

### Proof of Proposition 4.

Proposition 2 holds for any  $p^S = p^Q \in [0, 1]$  and, in particular, for  $p^S = p^Q = 0$ . Under this assumption, conditions (14), (17) and (18) can never be satisfied, equilibria  $MEE_O$ ,  $MEE_x$  and  $MEE_I$  are not feasible and hence,  $x^* = 0$ . Because  $b - \sigma s > 0$ , if  $\varepsilon < b - \sigma s$  condition (15) holds and  $MEE_y$  is the stable equilibrium. But when  $p^Q = 0$ , then  $y^*$  in (12) equals  $\Delta/(1 - s)$ . Conversely, if  $b - \sigma s \leq \varepsilon$ , then condition (16) holds and  $MEE_{01}$  is the stable equilibrium. The combination of these two values of  $y^*$  with  $x^* = 0$  is precisely the NE given in (21). ■

### Proof of Proposition 6.

Plugging  $(x^{NE}, y^{NE})$  into (23), after some simplifications, the social welfare in the population

of Ss and Qs reads:

$$\begin{aligned}\pi_{NE}^S &= \varepsilon - \frac{d(b-\varepsilon)}{\sigma}, & \pi_{NE}^Q &= \alpha\varepsilon - \frac{d(b-\varepsilon)}{\sigma}, & \text{if } y^{NE} < 1, \\ \pi_{NE}^S &= b - (b+\phi)s, & \pi_{NE}^Q &= \alpha\varepsilon - ds, & \text{if } y^{NE} = 1.\end{aligned}\tag{31}$$

*MEE<sub>O</sub>*: Under condition (14), the social welfare in each population is:

$$\pi_O^S = -\phi, \quad \pi_O^Q = -\phi - (1-\alpha)\varepsilon.$$

Comparing these expressions against (31), the gap between them is the same in both populations:

$$\pi_O^h - \pi_{NE}^h = -\frac{b+\phi}{\sigma}[\varepsilon - (d-\phi)], \quad h \in \{S, Q\}.$$

And this expression is negative under condition (2).

*MEE<sub>y</sub>*: Under condition (15) the equilibrium *MEE<sub>y</sub>* is asymptotically approached. From (12) and (21),  $y^* < y^{NE} \leq 1$ , for all  $p^Q \in (0, 1]$ . Moreover, from (23), the social welfare in each population with a zero compliance among Ss and  $y$  compliance among Qs reads:

$$\pi^S(y) = (b+\phi)y(1-s) - \phi, \quad \pi^Q(y) = y[\varepsilon - (d-\phi) + (b+\phi)(1-s) - \sigma(1-s)y] - \phi - (1-\alpha)\varepsilon.$$

After some computation, the average payoff in the population of Ss and in the population of Qs for the *MEE<sub>y</sub>* and for the *NE* compares as:

$$\pi_y^S - \pi_{NE}^S = (b+\phi)(1-s)(y^* - y^{NE}),\tag{32}$$

$$\pi_y^Q - \pi_{NE}^Q = \sigma(1-s) \left( \frac{\sigma+d}{\sigma} - y^* \right) (y^* - y^{NE}).\tag{33}$$

For  $p^Q \in (0, 1]$  it is known that  $y^* - y^{NE} < 0$  and hence, both  $\pi_y^S - \pi_{NE}^S < 0$  and  $\pi_y^Q - \pi_{NE}^Q < 0$ .

Moreover, from these expressions immediately follows:

$$\frac{d(\pi_y^S - \pi_{NE}^S)}{dy^*} = (b+\phi)(1-s) > 0, \quad \frac{d(\pi_y^Q - \pi_{NE}^Q)}{dy^*} = \sigma(1-s) \left[ \left( \frac{\sigma+d}{\sigma} - y^* \right) - (y^* - y^{NE}) \right].$$

Because  $0 < y^* \leq y^{NE}$ ,  $\sigma > 0$  and  $s \in (0, 1)$ , then  $d(\pi_y^Q - \pi_{NE}^Q)/dy^* > 0$ . Moreover, from (12)  $dy^*/dp^Q < 0$ , which ends the proof of the second statement in Proposition 6.

■

### Proof of Proposition 7.

The social welfare under the NE was computed in (31).

*MEE<sub>x</sub>*: Condition (17) characterizes *MEE<sub>x</sub>* =  $(x^*, 1)$ , while the NE is  $(0, 1)$ .



(a) The comparison for Qs is straightforward:

$$\pi_x^Q = \alpha\varepsilon - (1 - x^*)ds > \alpha\varepsilon - ds = \pi_{NE}^Q. \quad (34)$$

The effect of  $p^S$  on  $\pi_x^Q - \pi_{NE}^Q$  follows straightforwardly. From (13),  $dx^*/dp^S > 0$  and from (34),  $\pi_x^Q - \pi_{NE}^Q = x^*ds$ . Hence,  $d(\pi_x^Q - \pi_{NE}^Q)/dp^S > 0$ .

(b) For Ss, the average payoff within this population in equilibrium  $MEE_x$  reads:

$$\pi_x^S = (1 - x^*)[b + \sigma s x^* - (b + \phi)s],$$

and the comparison against  $\pi_{NE}^S$  in (31) with  $y^{NE} = 1$ , can be written as a second order polynomial  $p(x^*)$ :

$$\pi_x^S - \pi_{NE}^S = p(x^*) = -x^* \left[ x^* - \frac{\sigma s + (b + \phi)s - b}{\sigma s} \right] \sigma s. \quad (35)$$

i) If  $b > \sigma s + (b + \phi)s$ , the polynomial  $p(x^*)$  has a null and a negative root.

Consequently, for any  $x^* > 0$ , it holds that  $\pi_x^S - \pi_{NE}^S < 0$  and decreasing in  $x^*$ .

Therefore, since  $dx^*/dp^S > 0$ , then  $d(\pi_x^S - \pi_{NE}^S)/dp^S < 0$ .

ii) If  $b < \sigma s + (b + \phi)s$ , the polynomial  $p(x^*)$  has a null and a positive root at:

$$x^+ = \frac{\sigma s + (b + \phi)s - b}{\sigma s}.$$

Under condition (16),  $x^* > 0$  and hence,

$$\pi_x^S - \pi_{NE}^S \geq 0 \Leftrightarrow x^* \leq x^+.$$

And from the definition of  $x^*$  in (13), this is equivalent to:

$$\sqrt{(d - \phi)^2 + 4\varepsilon\sigma p^S \alpha(1 - s)} \leq 2(b + \phi)s - (d - \phi).$$

Because we are under assumption  $b < \sigma s + (b + \phi)s$ , the RHS of this expression is positive. Hence, it can be rewritten as:

$$A(s) = A_2 s^2 - s[A_1 + 2A_0] + A_0 \geq 0,$$

with  $A_2 = (b + \phi)^2 > 0$ ,  $A_1 = (d - \phi)(b + \phi) > 0$ , and  $A_0 = -\alpha\varepsilon p^S \sigma < 0$ .

This second order polynomial in  $s$  presents two real roots of opposite sign. The positive root is given by:

$$\hat{s} = \frac{A_0 + A_1 + \sqrt{(A_0 + A_1)^2 - 4A_0A_2}}{2A_2} \in (0, 1). \quad (36)$$

It is easy to prove that  $\hat{s} < 1$ . Thus, any positive  $s < \hat{s}$  implies  $A(s) < 0$  while for  $s \in (\hat{s}, 1)$  it holds that  $A(s) > 0$ . Note finally that  $MEE_x$  only occurs under condition (17), which requires  $s > \underline{s}$ , with

$$\underline{s} = \begin{cases} \frac{b - \alpha\varepsilon p^s}{\sigma} & \text{if } b > \alpha\varepsilon p^s, \\ \frac{b - \alpha\varepsilon p^s}{\alpha\varepsilon p^s} & \text{if } b \leq \alpha\varepsilon p^s. \end{cases}$$

If  $\underline{s} < \hat{s}$  then any  $s \in (\underline{s}, \hat{s})$  (for which  $A(s)$  is negative) implies  $\pi_x^s < \pi_{NE}^s$ , while any  $s \in (\hat{s}, 1)$  (for which  $A(s)$  is positive) implies  $\pi_x^s > \pi_{NE}^s$ . If  $\underline{s} > \hat{s}$  the first subinterval disappears and this solution is only compatible with  $s \in (\underline{s}, 1)$ , leading to  $\pi_x^s > \pi_{NE}^s$ .

In this scenario,  $p(x^*)$  increases from 0 to its maximum value as  $x^*$  runs from 0 to  $x^+/2$ , it decreases from this point to reach 0 at  $x^+$ , and it continues decreasing below 0 as  $x^*$  moves to the right of  $x^+$ . The value of  $p^s$  where  $p(x^*)$  reaches its maximum, i.e. where  $x^* = x^+/2$  is denoted by:

$$\hat{p}^s = \frac{b^2 - 2b(d - \phi) + (ds)^2 + 2ds\sigma}{4\alpha\varepsilon\sigma(1 - s)}.$$

Since  $x^*$  is monotonously increasing in  $p^s$ , immediately follows that  $\pi_x^s - \pi_{NE}^s$  increases with  $p^s$  as it moves from the minimum value at which  $MEE_x$  is the long-run equilibrium,  $(b - \sigma s)/(\alpha\varepsilon)$ , to  $\hat{p}^s$ . And it decreases from this value on (becoming negative at the value of  $p^s$  where  $x^* = x^+$ ).

Thus, as  $p^s$  increases  $d(\pi_x^s - \pi_{NE}^s)/dp^s > 0$  up until  $\hat{p}^s$ . And to the right of this value  $d(\pi_x^s - \pi_{NE}^s)/dp^s < 0$ .

Finally, note that if  $b > \sigma s + ds$ , then  $\hat{p}^s < 1$ , while if  $b < \sigma s + ds$ , then  $\hat{p}^s > 1$  and the interval with  $\pi_x^s < \pi_{NE}^s$  does not exist.

$MEE_I$  Condition (18) characterizes  $MEE_I = (1, 1)$ , while the NE is  $(0, 1)$ .

- (a) The comparison for Qs is straightforward:  $\pi_1^Q = \alpha\varepsilon > \alpha\varepsilon - ds = \pi_{NE}^Q$ .
- (b) For Ss, reasoning as in the previous item,  $\pi_1^S - \pi_{NE}^S = (b + \phi)s - b$ . And this is clearly positive (resp. negative) if  $s$  is above (resp. below)  $b/(b + \phi)$ .

■