Thesis Submitted to the University of Nottingham for the

Degree of Doctor of Philosophy

# The Socially Embedded Individual

Xueheng Li

School of Economics

April 18, 2018

# Abstract

This thesis contains three studies. They are connected by the idea that "no man is an island": each individual contributes to shaping, and is constrained by, the social and economic structures of the organization or the society that the individual is embedded in.

The *first* study, Chapter 2, examines optimal networks with weighted and directed links under complementarities. A group of agents take actions that are endogenously determined by which network the planner implements. Complementarities mean that the best-response action of each agent is increasing in the actions of those who have a link with positive weight pointing to the agent (representing the direction and intensity of influence). Optimal networks are those maximizing the planner's objective function which is an increasing function in the effort of each agent, subject to the constraint that the total weight of the links of the network does not exceed a certain level. The agents' best-response function and the planner's objective function can be convex or concave. We show that every optimal network exhibits dramatic concentration of influence so that a very small number of agents impose significant impact on the productivity of the whole organization.

The *second* study, Chapter 3, investigates how cooperative norms emerge and evolve over time. I construct a stochastic dynamic model based on the idea that cooperation in one-shot interactions is sustained by endogenous social norms. The model shows how cooperation and punishment of defectors co-evolve. It reveals the conditions under which cooperation emerges and persists in the long run. In

particular, recent empirical studies find that cooperation in one-shot interactions is positively correlated with law enforcement across societies, and that cooperation is higher in large, modern societies with higher degrees of market integration compared to small-scale societies. I extend the model to explain these regularities. I show that the ability to "vote with feet" is the key to understanding the difference in cooperation between small-scale societies and large, modern societies.

The *third* study, Chapter 4, is an experimental project, a joint work with Lucas Molleman and Dennie van Dolder. Previous studies suggest that whether individuals perceive a behavior as fair depends on its frequency in the population. Using a prisoner's dilemma game, we test experimentally whether informing individuals of a higher proportion of cooperators in the population affects the fairness perception about free riding and changes individuals' punishment of free riders. Different from previous studies, we use the strategy method to obtain each participant's complete punishment strategy. We find a remarkable heterogeneity among participants: some participants increase punishment of free riders as the proportion of cooperators increases, suggesting that they consider free riding to be more unfair when more cooperators are around; yet, many others punish independently of the proportion of cooperators. We show that the heterogeneity cannot be captured by any single existing theory.

# Acknowledgments

So many people have helped me during the past four PhD years and contributed to the completion of this thesis in one way or another. First of all, I thank my supervisors Abigail Barr and Silvia Sonderegger. I have benefited greatly from them throughout the PhD program. Whenever I needed their advices and insights, they have been always available.

I thank my officemates Despoina Alempaki and Hanna Fromell, for their daily companion, encouragements, academic and nonacademic discussions, and many indispensable pushes throughout the past two years. They have also given countlessly many detailed comments that have led to significant improvement of this thesis.

Almost every faculty member in the Center for Decision Research and Experimental Economics (CeDEx) has given useful advices to me at several stages of my research. In particular, I would like to thank Alex Possajennikov and Daniele Nosenzo, who have provided important comments on the studies contained in this thesis. I also thank my coauthors, Lucas Molleman, Dennie van Dolder, and Miachel Muthukrishna, for having many insightful discussions with them.

This thesis would not exist without funding from the Network for Integrated Behavioral Science (NIBS), which I gratefully acknowledge.

The four PhD years is not a short period, and it is not always easy. Let me finally thank Till Weber, Tom Lane, Arno Hantzsche, Cindy Fu, Vasileios Kotsidis, and Jesal Sheth. The experience of being friends and colleagues with them demonstrates that complementarities exist in social networks.

# Contents

# List of Figures

# Chapter 1

# Introduction

There used to be a tension between the economic approach and the sociological approach of human behavior: "Economics is all about how people make choices; sociology is all about how they don't have any choices to make" (Duesenberry, 1960, p. 233). In *Economic Action and Social Structure: the Problem of Embeddedness*, Mark Granovetter (1985) criticizes the under-socialized account of human actions used by many economists, as well as the over-socialized account proposed by many sociologists. This thesis contains three studies: two theoretical and one experimental. While research methods vary across the three studies, there is a common underlying assumption. That is, a social or an economic organization is not a simple collection of actions or preferences of individuals. On the one hand, the structure of an organization constrains the choices of individuals, and it might also change individuals' preferences over their actions. On the other hand, an organization's structure is itself endogenous: it emerges from the interactions of individuals who are endowed with more primitive social psychologies, the "underlying preferences" (Becker, 1978, p. 5), that are independent of specific structures of the organization.

In *Concentration of Influence under Complementarities*, I use a network approach to study the optimal structure of organizations. A motivation for this study is Mas and Moretti (2009)'s empirical observation that: (a) complementarities,

such as positive productivity spill-overs, exist in the work place, (b) by carefully designing the direction of flow of information relevant for production, a planner can control the direction of the complementarity effects, and (c) the intensity of those complementarity effects depends on the frequency of interactions among the workers. Hence, to capture the possibility of controlling the direction and intensity of complementarities, we investigate optimal networks with unweighted and directed links. This departs from most of the existing works of networks in economics which mostly focus on unweighted networks (e.g., Baetz, 2015; Belhaj et al., 2016; Galeotti and Goyal, 2010; Hiller, 2017).

Specifically, we investigate a setting in which there are a set of agents, each of which needs to take a production action. Before the agents take the actions, the planner chooses a network to implement. The network consists of the agents as nodes and a set of unweighted and directed links that connect the agents. We can think of the links as representing any one-way flow of information relevant for productivity. We examine a setting with complementarities such that each agent's action is increasing in those who have a link with positive weighted pointing towards that agent. The strength of the complementarities depends on the weight of the links. The planner's problem we aim to solve is what kind of networks that the planner would like to implement, given that (a) the planner aims to maximize an objective function that is increasing in the action of each agent, and that (b) the planner facts the constraint that the total weight of the links of the implemented network cannot exceed a certain amount. The planner objective function may be convex or strictly concave. We call those solution networks *optimal networks*.

What we find is that every optimal network exhibits a *dramatic concentration of influence* such that a very small number of agents impose significant impact on the productivity of the whole group. More precisely, in every optimal network, regardless of the size of whole group, there are two and only two agents who would have outward-link with positive weight pointing towards others, reflecting their influence on the rest.

Our finding is related to Galeotti and Goyal (2010)'s observation about information sharing networks that "a large majority of individuals get most of their information from a very small subset of the group". However, the scenario their model captures is distinct from ours. Besides focusing on unweighted networks, they consider a situation in which the actions of neighbor agents exhibit *substitutability* instead of complementarity. Also, they investigate *equilibrium* networks (formed by agents who choose their own links), whereas we characterize *optimal* networks from a planner's perspective. Yet, despite the differences, surprisingly, we observe a very similar outcome to the law of the few. Hence, together with Galeotti and Goyal (2010), we contribute to the literature by providing different rationales to the pervasiveness of concentration of influence observed in the real world, be it in decentralized information sharing networks or centralized production organizations.

In *Norm-based Resentment and the Evolution of Cooperative Norms*, I investigate how cooperative norms emerge and evolve over time. I analyze a stochastic dynamic model constructed on *norm-based resentment* (Sugden, 1984, 2000, 2004; Bicchieri, 2006; Cooper and Dutcher, 2011; Falk and Ichino, 2006; Herz and Taubinsky, 2017; Kahneman et al., 1986; Peysakhovich and Rand, 2016). The idea is that individuals have empirical expectations about others' behavior. They feel frustrated and may punish a defector if, and only if, they expect that most others cooperate. Norm-based resentment leads to two locally stable equilibria: the *defection equilibrium* in which each individual in the population defects and does not punish defectors, and the *cooperation equilibrium* in which each individual cooperates and punishes defectors. A social norm is defined as a locally stable equilibrium. Following Kandori et al. (1993), Ellison (1993, 2000), and Young (1993, 2001), I set up a myopic best-response dynamic in which a population of individuals is randomly matched to play a cooperation game recurrently and infinitely over time. Individuals may also make mistakes such that they deviate from their best-responses with positive probability.

In the first part of the analysis I characterize the most likely equilibrium in an infinite span of time when the probability of making a mistake goes to zero, the *stochastically stable equilibrium* (Ellison, 1993, 2000; Kandori et al., 1993; Young, 1993, 2001). The basic result is that, given norm-based resentment, the cooperation equilibrium can be stochastically stable. Whether and when this is the case depends on the trade-off between *the intolerance of defection* generated by norm-based resentment and *the temptation to defect*.

The second part of the analysis investigates two particular empirical regularities observed in recent cross-cultural experiments: *i*) cooperation is higher in societies with better law enforcement (Gächter and Schulz, 2016; Herrmann et al., 2008; Tabellini, 2008), and *ii*) cooperation is also higher in large, modern societies with higher degrees of market integration compared to small-scale societies (Henrich et al., 2010). I show that norm-based resentment combined with the ability to vote with feet generates a migration effect and a fitting-in effect that, together, explain the higher cooperation level in large, modern societies.

Chapter 4, *Fairness Perceptions and Punishment Types*, is an experimental study, a joint work with Lucas Molleman and Dennie van Dolder. Previous theoretical and empirical literature (Bicchieri, 2006; Cooper and Dutcher, 2011; Herz and Taubinsky, 2017; Kahneman et al., 1986; Sugden, 2000, 2004) suggests that individuals' fairness perceptions and punishment of selfish behavior may depend on the prevalence of the behavior in society. We design an experiment to further investigate this issue. Using a prisoner's dilemma game experiment, we examine whether informing individuals of a higher level of cooperation in the population affects their punishment of free riders. The novelty of our experiment is that we use the strategy method to elicit each participant's complete punishment strategy, i.e., how each individual punishes defectors in response to each possible prevalence level of cooperation in our sample. The advantage of the strategy method is that it allows us to reveal any potential heterogeneity in punishment behavior among participants. In contrast, previous studies only show the average response to dif-

ferent levels of cooperation at the aggregate level. The individual-level data also provide a new opportunity to test some existing theories relevant to understanding public goods contribution that are otherwise difficult to disentangle.

We find that participants are heterogeneous in their punishment strategies. 21% of our participants punish free riders independently of the proportion of cooperators. These participants' behavior is in line with standard models of inequality aversion (Fehr and Schmidt, 1999) and intention-based reciprocity (Dufwenberg and Kirchsteiger, 2004; Rabin, 1993). In contrast, 13% increase their punishment as the percentage of cooperation goes up, suggesting that these participants consider free riding more unfair when more participants cooperate. Interestingly, we also find 10% of participants who punish free riders less as cooperation becomes more common. These participants' behavior can only be explained by a model assuming diminishing marginal return to punishment.

# Chapter 2

# Concentration of Influence under Complementarities

## 2.1 Introduction

Various forms of complementarities, such as positive productivity spillovers, are often present in the workplace. For example, Mas and Moretti (2009) examine empirically the productivity of supermarket workers and find "strong evidence of positive productivity spillovers from the introduction of highly productive personnel into a shift". Moreover,

> "Worker effort is positively related to the productivity of workers who see him, but not workers who do not see him. Additionally, workers respond more to the presence of coworkers with whom they frequently interact." (Mas and Moretti, 2009)

That is, not only complementarities exist in the workplace, but we can also control the *direction* and *intensity* of complementarities by carefully designing the direction of the flow of information relevant for productivity and the frequency of interactions between employees. Then a key question is—what is the optimal arrangement of an organization given the complementarities, and especially given that we can control the direction and intensity of the complementarities? This

paper studies this question.

Alongside the complementarities, many real-world networks and organizations exhibit dramatic *concentration of influence*—i.e., a very small number of individuals can impose significant impact on the productivity of the whole organization. For example, there is only one CEO in every firm, and one president in every country. Empirical studies show that a remarkable payment gap exists between CEOs and ordinary employees (Connelly et al., 2013; Edmans and Gabaix, 2016; Faleye et al., 2013; Mishel and Sabadish, 2013). This large payment gap partly reflects the great difference in the extent of influence between a CEO and ordinary employees on a firm's productivity.

More specifically, this paper uses a network approach to study the optimal arrangement of organization given complementarities. We show that the coappearance of complementarities and concentration of influence might not be a coincidence. Instead, optimization given complementarities leads to dramatic concentration of influence to a very small number of individuals, and this holds under broad conditions.

The setting we investigate is the following. There is a group of homogenous agents; we can think of them as employees in a firm. Each agent exerts an effort to perform a production task. We can potentially increase an agent's effort by increasing the effort of some others, due to the complementarity we assume (e.g, the peer effects documented by Mas and Moretti (2009)).[1] Whether or not this occurs, and the strength of the increase, depend on the *weighted* and *directed links* the planner implements to connect the agents. These links could represent any one-way flow of knowledge or information relevant for production. The agents and the links then form a network with weighted and directed links. The planner's problem is to find a weighted and directed network to maximize an objective function that is increasing in the effort of each agent—we call these networks

[1]The best-response function we consider is a general one; an agent's effort can be convex, or strictly concave, in the effort of those who influence them. The basic network game we analyze thus nests many existing models in the literature, including Baetz (2015); Ballester et al. (2006); Belhaj et al. (2016); Hiller (2017).

*optimal networks.* The planner's objective function can be convex, as previously studied (e.g., Belhaj et al. 2016; Hiller 2017), or strictly concave, which, to our knowledge, has not been examined. The constraint the planner faces is that the total weight of the links of the network does not exceed a certain amount.

We show that every optimal network is a what we call *ABC-form network.* Every ABC-form network exhibits dramatic concentration of influence to a very small number of agents. Figure 2.1 below displays a typical ABC-form network. Regardless of the total number of agents, there is only one agent, A, who influences the rest (i.e., having outward links with positive weight pointing towards others). There is also an agent, B, who influences A and is influenced by A. Intuitively, we can think of B as A's work assistant. The rest of agents, Cs, never influence any others, i.e, they never have an outward link with positive weight pointing towards any other agent. The numbers, $a, b, c$, next to the links in Figure 2.1 are the weight of those links. The surprising observation is that, even if all agents are ex ante identical, it is optimal to concentrating all resources to enhance a single agent's influences on the rest and another agent's influence on the single center agent.[2]



**Figure 2.1:** A regular ABC-form network: $a > b > c \geq 0$.

The two studies that are closest to ours are Belhaj et al. (2016) and Hiller (2017), who also investigate optimal networks under complementarities. However, they restrict the search of optimal networks from those with unweighted and undirected links. This restriction imposes two implicit assumptions. One is that,

---

[2]The class of ABC-form networks also includes those with only two agents influencing each other while all the rest are isolated.

since the weight of each link is either 0 or 1, this means that an upper bound is imposed exogenously on the weight of each link. Thus, concentration of influence is suppressed in their models. Second, when restricting to undirected links with symmetric two-way flow of complementarities, they rule out the possibility of asymmetric influences such that some agents influence the rest while the rest do not influence the former. In contrast, we allow for symmetric influences as well as asymmetric ones. We show that, indeed, highly asymmetric influences could be optimal.

Also, Belhaj et al. (2016) and Hiller (2017) only examine objective functions that are convex in each agent's effort. A convex objective function covers the case where the planner is an utilitarian one who wants to maximize the sum of the utilities of the agents, given that the agents have a linear-quadratic utility function as in Ballester et al. (2006). However, it does not capture the cases where the planner has a fairness concern so that she wants to trade-off equality against aggregate efficiency. In our model, the objective function can be convex or strictly concave, covering both cases.

More broadly, this paper, together with a recent paper by Galeotti and Goyal (2010), provides rationales to the dramatic concentration of influence we often observe in networks and organizations in reality. Galeotti and Goyal (2010) make a similar observation regarding information sharing networks, namely, *the law of the few* such that "a large majority of individuals get most of their information from a very small subset of the group". However, the scenario their model captures is very different from ours. Besides focusing on unweighted networks, they consider a situation in which the actions of neighbor agents exhibit *substitutability* instead of complementarity. Also, they investigate *equilibrium* networks (formed by agents who choose their own links), whereas we characterize *optimal* networks from a planner's perspective. Yet, despite the differences, surprisingly, we observe a very similar outcome to the law of the few. Hence, together with Galeotti and Goyal (2010), we contribute to the literature by providing different rationales to

the pervasiveness of concentration of influence observed in the real world, be it in decentralized information sharing networks or centralized production organizations.

Section 2.2 presents the model. Section 2.3 characterizes the optimal networks in the general model. Section 2.4 applies our main result to the special case where the best-response function of agents is linear. Section 2.5 concludes. All proofs are provided in Appendix.

## 2.2 Model

Consider a set of agents $N = \{1, 2, \ldots, n\}$ with $n \geq 3$. For each $i, j \in N$, $i \neq j$, let $g_{ij} \geq 0$ denote the **weighted link** pointing from $j$ to $i$. The links are directed, i.e., $g_{ij}$ need not equal to $g_{ji}$. Let $G = (g_{ij})$ be the $n$-by-$n$ matrix whose $ij$th element is $g_{ij}$ and has zeros at its main diagonal, i.e., $g_{ii} = 0$ for each $i \in N$.

Let $x_i(G)$ be the **effort** of $i$ given network $G$. We assume that $\mathbf{x}(G) = (x_1, \ldots, x_n)$ is a solution to the following best-response system:

$$x_i = \phi \left( \sum_{j \in N} g_{ij} x_j \right).$$

In general, a fixed point of $\phi$ may not exist, or it exists but is not unique. Throughout this paper, we restrict our attention to the cases where a solution $\mathbf{x}(G)$ exists and is unique (we will present the assumption that guarantees this later). The problem we consider is that which network $G$ we would like to implement to maximize a planner's objective function $f : \mathbb{R}^n \to \mathbb{R}$ that is increasing in each agent's effort. There is also a constraint on the networks that the planner can implement.

In what follows, we properly define the planner's problem and present the assumptions regarding the objective function $f$ and the constraint the planner faces.

**Assumption 2.1** (*Complementarity*). $\phi \geq 0$ and $\phi' > 0$.

**Example 2.1** (Ballester et al., 2006; Belhaj et al., 2016). Suppose each agent takes effort $x_i$ to maximize the utility function

$$u_i(\mathbf{x}) = \left(1 + \sum_{j \in N} g_{ij} x_j\right) x_i - \frac{1}{2} x_i^2.$$

Then we have the *linear* best-response function

$$x_i = 1 + \sum_{j \in N} g_{ij} x_j$$

and the closed-form solution (if a solution exists)

$$\mathbf{x}(G) = [I - G]^{-1} \cdot \mathbf{1},$$

where $I$ is the $n$-by-$n$ identity matrix, and $\mathbf{1} = (1, 1, \ldots, 1)$ is the column vector of $n$-folds of 1s. Note also that, given each $i$ exerts effort $x_i(G)$, we have

$$u_i(\mathbf{x}(G)) = \frac{1}{2} x_i(G)^2.$$

(End of Example 2.2)

**Example 2.2** (The baseline model of Baetz (2015)). Consider

$$u_i(\mathbf{x}) = 2 \left(\sum_{j \in N} g_{ij} x_j\right)^{\frac{q}{2}} x_i^{\frac{1}{2}} - x_i.$$

with $0 < q < 1$. Then we have a *strictly concave* best-response function

$$x_i = \left(\sum_{j \in N} g_{ij} x_j\right)^q,$$

and, given each $i$ exerts $x_i(G)$, we have

$$u_i(\mathbf{x}(G)) = x_i(G).$$

(End of Example 2.2)

Let $\eta(G)$ be the cost of implementing network $G$, with $\eta \geq 0$ and $\partial\eta(G)/\partial g_{ij} > 0$ for each $i, j \in N$. Let $\bar{\eta} > 0$ represent the total resources that the planner can spend to implement the networks. That is, a network $G$ is feasible if and only if $\eta(G) \leq \bar{\eta}$. Let $\mathcal{G}$ denote the collection of networks with $\eta(G) \leq \bar{\eta}$. The planner would like to choose a network from $\mathcal{G}$ to maximize an objective function $f : \mathbb{R}^n \to \mathbb{R}$ which is increasing in the effort of each agent. More precisely, our task is to characterize the following:

**Definition.** An **optimal network** $G$ is one that maximizes $f(\mathbf{x}(G))$ subject to $G \in \mathcal{G}$.

**Assumption 2.2.** $\eta(G) = \sum_{i,j \in N} g_{ij}$.

Our model is quite general in various dimensions except for the above assumption that the link-cost function is linear. Indeed, as we will discuss further after we show our main result, the structure of optimal networks would change if we consider alternative cost functions such as $\eta(G) = \sum_{i,j \in N} h(g_{ij})$ or $\eta(G) = \sum_{i \in N} h(\sum_{j \in N} g_{ij})$ with $h' > 0, h'' < 0$. However, assuming linear link-cost is a reasonable benchmark. Moreover, it is by far the most widely used link-cost function in the literature (e.g., Baetz, 2015; Belhaj et al., 2016; Galeotti and Goyal, 2010; Goyal and Bala, 2000; Hiller, 2017). Hence, it is worth investigating how far we can push our result by maintaining this assumption.

We also assume the following to ensure that the problem is well-defined. For each $G \in \mathcal{G}$, define $\Phi_G : \mathbb{R}^n \to \mathbb{R}^n$ as

$$\Phi_G(\mathbf{x}) = (\phi(\sum_{j \in N} g_{1j} x_j), \phi(\sum_{j \in N} g_{2j} x_j), \ldots, \phi(\sum_{j \in N} g_{nj} x_j)).$$

**Assumption 2.3.** $\phi$ and $\mathcal{G}$ are such that $\Phi_G$ is a contraction mapping for each $G \in \mathcal{G}$.

By the Banach fixed-point theorem, the above assumption implies that a unique $\mathbf{x}(G)$ exists for each $G \in \mathcal{G}$, and that for each $\mathbf{x} \geq 0$ we have $\lim_{t \to \infty} \Phi_G^t(\mathbf{x}) =$

$\mathbf{x}(G)$, where $\Phi_G^t(\mathbf{x})$ is the $t$-th functional power of $\Phi_G(\mathbf{x})$. Roughly speaking, this assumption is satisfied if $\phi$ is not too convex and $\bar{\eta}$ is not too large. For example, consider the linear best-response function in Example 2.1. The above assumption is satisfied if $\bar{\eta} < 2$. In Example 2.2, $\phi$ is strictly concave and $\phi'(y) \to 0$ as $y \to \infty$. Thus, it is easy to check that, for each $\bar{\eta} > 0$ and $G \in \mathcal{G}$, a unique (strictly positive) $\mathbf{x}(G)$ exists and $\lim_{t \to \infty} \Phi_G^t(\mathbf{x}) = \mathbf{x}(G)$.

The assumptions regarding the objective function $f$ are as follows. Let $f_i(\mathbf{x})$ denote the partial derivative of $f$ with respect to its $i$-th argument.

**Assumption 2.4.**   *1. $f_i(\mathbf{x}) > 0$ for each $i \in N$; and*

*2. $f$ is symmetric, i.e., $f_i(\mathbf{x}) = f_j(\mathbf{x})$ if $x_i = x_j$.*

In other words, $f$ is strictly increasing in the effort of each agent. The symmetry assumption is not necessary but improves exposition. Note that $f$ can be linear, strictly convex, or strictly concave in the effort of each agent. Moreover, no assumption is imposed on the separability of the arguments of $f$.

**Example 2.3.** Consider a *utilitarian* planner with $f(\mathbf{x}(G)) \equiv \sum_{i \in N} u_i(\mathbf{x}(G))$. Then, in the case where the agents have the utility function in Example 2.1, we have

$$f(\mathbf{x}(G)) = \frac{1}{2} \sum_{i \in N} x_i(G)^2.$$

That is, the planner wants to maximize the sum of squares of the agents' effort, and $f$ is strictly convex in the effort of each agent.

In the case where agents have the utility function in Example 2.2, we have

$$f(\mathbf{x}(G)) = \sum_{i \in N} x_i(G).$$

In this case, the planner wants to maximize the sum of the agents' effort, and $f$ is linear. (End of Example 2.3)

**Example 2.4.** Agents have the utility function in Example 2.2, and the planner has a taste for *fairness*:

$$f(\mathbf{x}(G)) \equiv \sum_{i \in N} \ln u_i(\mathbf{x}(G)) = \sum_{i \in N} \ln x_i(G).$$

In this case, $f$ is *strictly concave* in each $x_i$. Also, maximizing $\sum \ln x_i(G)$ is equivalent to maximizing $\Pi x_i(G)$. Hence, this example also shows that linear separability of $f$ is not necessary. (End of Example 2.4)

Finally, we assume the following to rule out the possibility that networks with only one link $g_{ij} = \bar{\eta}$ could be optimal. Let $\underline{x} = \phi(0)$ and $\bar{x} = \phi(\bar{\eta}\underline{x})$.

**Assumption 2.5.** $\phi'(\bar{\eta}\underline{x})\bar{\eta} + \frac{f_2(\bar{x},\underline{x},...,\underline{x})}{f_1(\bar{x},\underline{x},...,\underline{x})} > \frac{\underline{x}}{\bar{x}} \frac{\phi'(\bar{\eta}\underline{x})}{\phi'(0)}$.

All our previous examples satisfy this assumption. For example, consider the utility function of agents in Example 2.1, and the planner is an utilitarian one, so that we have $x_i = 1 + \sum_j g_{ij}x_j$ and $f(\mathbf{x}) = \frac{1}{2}\sum x_i^2$. Then: $\underline{x} = 1$, $\bar{x} = 1 + \bar{\eta}$, $\phi' = 1$, $f_2(\bar{x},\underline{x},...,\underline{x}) = 1$, and $f_1(\bar{x},\underline{x},...,\underline{x}) = 1 + \bar{\eta}$. Thus, the left-hand side of the inequality in the assumption equals to $\bar{\eta} + \frac{1}{1+\bar{\eta}}$, while the right-hand side $\frac{1}{1+\bar{\eta}}$.

Before we proceed to the analysis, let us comment that it is equally valid to call $x_i(G)$ the **centrality of $i$ under $G$**. This interpretation connects our analysis to the large sociology literature on networks. For example, if $\phi(y) = 1 + y$ so that $x_i = 1 + \sum_j g_{ij}x_j$, then $x_i(G)$ is simply the famous *Bonacich centrality* (Bonacich, 1987; Bonacich and Lloyd, 2001). Then the question we investigate is—which network maximizes an objective function that is increasing in the centrality of each agent, subject to the constraint that the total weight of the links of the network does not exceed $\bar{\eta}$.

## 2.3  Analysis

For expositional purpose, we index the agents so that

$$x_1(G) \geq x_2(G) \geq \ldots \geq x_n(G).$$

Given homogenous agents, this is without loss of generality. We aim to show that, for every $f$ and $\phi$ that satisfy Assumptions 2.1 to 2.5, all optimal networks take the following form.

**Definition.** $G$ is an **ABC-form network** if

   a) $g_{12}, g_{21} > 0$;

   b) $g_{i2} = 0$ for each $i > 2$;

   c) $g_{ji} = 0$ for each $i > 2$, $j \in N$; and

   d) $x_1(G) \geq x_2(G) > x_i(G)$ for each $i > 2$.

The first condition says that there are two agents, 1 and 2, influencing each other. The second condition says that: $a)$ one of the two connecting agents might have outward-links pointing towards other agents; $b)$ and $c)$ say that no other links with strictly positive weight exist. The fourth condition says that all ABC-form networks are *asymmetric* in the sense that there are agents who exert *strictly* higher effort than the rest. We call the networks satisfying the above conditions ABC-form networks because we can think of agent 1 as the director at the top of the organization (the A), agent 2 the director's assistant or consultant (the B), and the rest lower-level subordinates (the Cs).

For expositional purpose, we also define the following weaker notion.

**Definition.** $G$ is a **weak ABC-form network** if

   a) $g_{i2} = 0$ for each $i > 2$,

   b) $g_{ji} = 0$ for each $i > 2$, $j \in N$.

That is, a weak ABC-form network admits the possibility that $g_{12} = 0$ or $g_{21} = 0$, and that some agent $i > 2$ performs the same as agent 2. To show that all

optimal networks are ABC-form networks, we start with the following observation, which we call *one-link-switch principle.*

**Lemma 2.1** (One-link-switch principle). *Suppose that Assumptions 2.1, 2.2 and 2.3 hold. Consider three distinct agents $i, j, k \in N$. Let $G'$ be such that $g'_{kj} = 0$ and $g'_{ki} = g_{ki} + g_{kj}$, while $g'_{pq} = g_{pq}$ for all other elements in $G'$.*

*a) If $x_i(G) \geq x_j(G)$, then $x_\ell(G') \geq x_\ell(G)$ for each $\ell \in N$.*

*b) If $x_i(G) > x_j(G)$ and $g_{kj} > 0$, then $x_k(G') > x_k(G)$, and $x_\ell(G') \geq x_\ell(G)$ for each $\ell \in N$.*

All proofs are provided in Appendix. The one-link-switch principle says that, for any network $G$, if we can find an agent $j$ with an outward-link $g_{kj} > 0$, but there is another agent, $i$, who performs better than $j$, then we can construct a better network $G'$ in which everyone performs at least as good as before and some does strictly better. Such $G'$ is obtained by re-allocating the weight of the link $g_{kj}$ to the link $g_{ki}$.

In what follows, we consider a network $G = (g_{ij})$ which is not an ABC-form network. Building on the one-link-switch principle, we show that $G$ is not optimal. Our argument involves three steps. *First,* Lemma 2.2 below shows that we can usually obtain a network $G'$ using the one-link-switch principle such that $G'$ performs strictly better than $G$. If we indeed obtain a better $G'$ by switching a link, that means $G$ is not optimal. *Next,* if we cannot obtain a better network $G'$ by switching a link of $G$ using the one-link-switch principle, then $G$ must satisfy Condition (*). Lemma 2.3 shows that, given Condition (*), we can construct a weak ABC-form $\hat{G}$ on the basis of $G$ so that every agent performs exactly the same in $\hat{G}$ as they do in $G$. Moreover, certain properties must hold for $\hat{G}$. *Finally,* in Lemma 2.4, we examine closer $\hat{G}$ and show that $\hat{G}$ is not optimal, because we can further find an ABC-form network $\hat{G}^\epsilon$ that does strictly better than $\hat{G}$. But $\hat{G}$ generates the same outcome as the network $G$ we initially consider. This establishes that the optimal networks are among the set of ABC-form networks.

First, observe that, if $G = (g_{ij})$ is not a weak ABC-form network, then one of
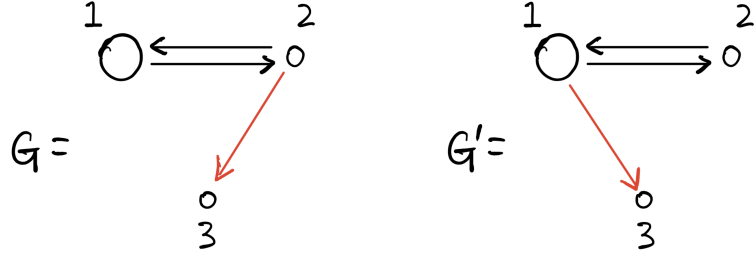
23

**Figure 2.2:** If $x_1(G) > x_2(G)$, then $x_3(G') > x_3(G)$, and $x_\ell(G') \geq x_\ell(G)$ for each $\ell \in N$.
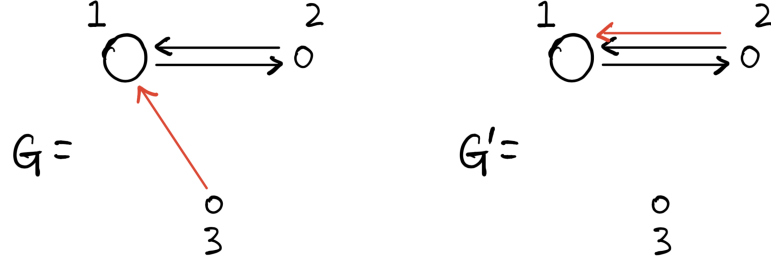


**Figure 2.3:** If $x_2(G) > x_3(G)$, then $x_1(G') > x_1(G)$, and $x_\ell(G') \geq x_\ell(G)$ for each $\ell \in N$.

the following must hold: (a) $g_{i2} > 0$ for some $i > 2$, (b) $g_{1i} > 0$ for some $i > 2$, or, (c) $g_{ji} > 0$ for some $i > 2, j \neq 1$. However, as Figures 2.2 and 2.3 illustrate, these conditions imply that, by switching a link of $G$ using the one-link-switch principle, we usually obtain a strictly better network $G' = (g'_{ij})$ with $\sum_{i,j \in N} g'_{ij} = \sum_{i,j \in N} g_{ij}$.

**Lemma 2.2.** *Suppose that Assumptions 2.1, 2.2 and 2.3 hold. Consider a network* $G = (g_{ij})$.

1. *If $g_{i2} > 0$ for some $i > 2$, and $x_1(G) > x_2(G)$, then there is a network $G' = (g'_{ij})$ with $\sum_{i,j \in N} g'_{ij} = \sum_{i,j \in N} g_{ij}$ such that $x_i(G') > x_i(G)$ and $x_\ell(G') \geq x_\ell(G)$ for each $\ell \in N$.*

2. *If $g_{1i} > 0$ for some $i > 2$, and $x_2(G) > x_i(G)$, then there is a network $G' = (g'_{ij})$ with $\sum_{i,j \in N} g'_{ij} = \sum_{i,j \in N} g_{ij}$ such that $x_1(G') > x_1(G)$ and $x_\ell(G') \geq x_\ell(G)$ for each $\ell \in N$.*

3. *If $g_{ji} > 0$ for some $i > 2$, $j \neq 1$, and $x_1(G) > x_i(G)$, then there is a network $G' = (g'_{ij})$ with $\sum_{i,j \in N} g'_{ij} = \sum_{i,j \in N} g_{ij}$ such that $x_j(G') > x_j(G)$ and $x_\ell(G') \geq x_\ell(G)$ for each $\ell \in N$.*

Lemma 2.2 implies that, if $G$ is considered to be optimal so that we cannot obtain a strictly better network $G'$ by simply switching a link, then the following properties must hold for $G$:

**Condition (\*).**  *a) For each $i > 2$ with $g_{1i} > 0$, we have $x_i(G) = x_2(G)$;*

*b) for each $i > 2$ with $g_{ji} > 0$ for some $j \neq 1$, we have $x_i(G) = x_1(G)$; and*

*c) if $g_{i2} > 0$ for some $i > 2$, then $x_1(G) = x_2(G)$.*

Given the above condition, we can then construct a weak ABC-form network $\hat{G} = (\hat{g}_{ij})$ with $\sum_{i,j \in N} \hat{g}_{ij} = \sum_{i,j \in N} g_{ij}$ such that $x_i(\hat{G}) = x_i(G)$ for each $i \in N$. That is, $\hat{G}$ performs exactly the same as $G$. We obtain such $\hat{G}$ by:

a) for each $i > 2$ with $g_{1i} > 0$, let $\hat{g}_{1i} = 0$ and $\hat{g}_{12} = g_{12} + g_{1i}$;

b) for each $i > 2$ with $g_{ji} > 0$ for some $j \neq 1$, let $\hat{g}_{ji} = 0$ and $\hat{g}_{j1} = g_{j1} + g_{ji}$;

c) for each $i > 2$ with $g_{i2} > 0$, let $\hat{g}_{i2} = 0$ and $\hat{g}_{i1} = g_{i1} + g_{i2}$; and

d) let $\hat{g}_{pq} = g_{pq}$ for all other elements in $\hat{G}$.

**Lemma 2.3.** *Suppose that Assumptions 2.1, 2.2 and 2.3 hold. Consider a network $G = (g_{ij})$ and the weak ABC-form network $\hat{G}$ as constructed above. Suppose $G$ is not a weak ABC-form network and Condition (\*) holds for $G$. Then, first, we have $x_i(\hat{G}) = x_i(G)$ for each $i \in N$. Moreover, one of the following holds for $\hat{G}$:*

*a) $x_i(\hat{G}) = x_2(\hat{G})$ for some $i > 2$, and $\hat{g}_{12} > 0$;*

*b) $x_1(\hat{G}) = x_2(\hat{G})$, and $\hat{g}_{i1} > 0$ for some $i > 2$.*

Finally, Lemma 2.4 below shows that, if condition (\*) indeed holds for $G$ so that the conclusions of Lemma 2.3 hold for $\hat{G}$, then there must exist room for improvement. More precisely, we can find an ABC-form network $\hat{G}^\epsilon$ that performs strictly better than $\hat{G}$ and thus also better than $G$. Hence, $G$ is not optimal.

**Lemma 2.4.** *Suppose that Assumptions 2.1 to 2.5 hold. Consider a weak ABC-form network $\hat{G} = (\hat{g}_{ij})$, and suppose $\sum_{i,j \in N} \hat{g}_{ij} = \bar{\eta}$. Then $\hat{g}_{12} > 0$.*

*Moreover, suppose one of the following holds:*

*a) $\hat{g}_{21} = 0$;*

*b) $x_i(\hat{G}) = x_2(\hat{G})$ for some $i > 2$; or*

*c) $x_1(\hat{G}) = x_2(\hat{G})$, and $\hat{g}_{i1} > 0$ for some $i > 2$.*

*Then there is an ABC-form network $\hat{G}^\epsilon = (\hat{g}_{ij}^\epsilon)$ with $\sum_{i,j\in N} \hat{g}_{ij}^\epsilon = \bar{\eta}$ such that $f(\mathbf{x}(\hat{G}^\epsilon)) > f(\mathbf{x}(\hat{G}))$.*

Lemma 2.4 also implies that any network that is *only* a weak ABC-form network (but not an ABC-form network) is not optimal. This is because, for any network that is only a weak ABC-form network, we must have $\hat{g}_{12} = 0$ or $\hat{g}_{21} = 0$, or $x_i(\hat{G}) = x_2(\hat{G})$ for some $i > 2$.

We can now state the main result of this paper.

**Theorem 2.1.** *Suppose that Assumptions 2.1 to 2.5 hold. Then every optimal network is an ABC-form network.*

Three remarks follow. *First*, the set of ABC-form networks constitutes a very tiny subset of networks in $\mathcal{G}$. Hence, the result is sharp. *Second*, the proposition says that all optimal networks are *asymmetric* in the sense that some agents exert strictly more effort than others. This holds for every increasing $f$ and $\phi$ that satisfy our assumptions, regardless of their concavity, and it holds given that all agents are assumed ex ante identical. Nevertheless, the assumption that the link-cost function $\eta(G)$ is linear is critical to this result. We provide further discussion below. *Third*, as we previously mentioned, an ABC-form network can be one in which there are only two agents connecting with each other, while the rest of agents are all isolated. Whether or not this occurs does depend on the concavity of $\phi$ and $f$. We return to this issue in Section 2.4.

Now we comment on the importance of the linear link-cost assumption, $\eta(G) = \sum_{i,j\in N} g_{ij}$, in supporting the above argument. The linear link-cost assumption implies that, starting from a network within the constraint $\eta(G) \leq \bar{\eta}$, we can do *any* re-allocation of the weight of links without violating the constraint. Hence, from the cost side, there is no bound on the degree of concentration of the weight of links to just a few agents. However, suppose $\eta(G) = \sum_{i\in N} h(\sum_{j\in N} g_{ji})$ with $h' > 0, h'' > 0$. This convex link-cost function would then limit the concentration

26

of resources to enhance the influence of just a single agent. If $h$ is sufficiently convex, then ABC-form networks are not optimal.[3]

## 2.4 Linear best-response

In this section, we apply our general result to the following special case:

**Definition.** We call the problem $\max_G f(\mathbf{x}(G))$ s.t. $\sum_{i,j \in N} g_{ij} \leq \bar{\eta}$ **the linear best-response model** if the following assumptions hold:

a) $\phi(y) = 1 + y$, so that $x_i = 1 + \sum_{j \in N} g_{ij} x_j$;

b) $f(\mathbf{x}) = \sum_{i \in N} v(x_i)$ with $v' > 0$; and

c) $\bar{\eta} < 2$.

That is, we assume a linear best-response function in the linear model. A linear best-response function is widely assumed in the literature about network activities (e.g., Ballester et al., 2006; Belhaj et al., 2016; Bramoullé and Kranton, 2007; Bramoullé et al., 2014; Galeotti and Goyal, 2010). Hence, it is a good benchmark and provides a concrete demonstration how our result sheds new light on the literature. For analytical convenience, we also assume that the effort of agents exhibits linear separability in the objective function $f$. The assumption $\bar{\eta} < 2$ guarantees that a unique solution exists for every $G$ with $\sum_{i,j \in N} g_{ij} \leq \bar{\eta}$.

In what follows, we examine the conditions under which the following most stylized ABC-form network (see Figure 2.4 for illustration) is optimal.

---

[3]Despite being widely used in the literature (e.g., Baetz, 2015; Belhaj et al., 2016; Galeotti and Goyal, 2010; Goyal and Bala, 2000; Hiller, 2017), the linear link-cost assumption is not always plausible. For example, suppose that outward links represent "talking to" actions. Then it seems more reasonable to assume $\eta(G) = \sum_{i \in N} h(\sum_{j \in N} g_{ji})$ with $h' > 0, h'' > 0$. That is, the longer an agent talks to others, the more and more tired she feels and the more compensation she requires at the margin.
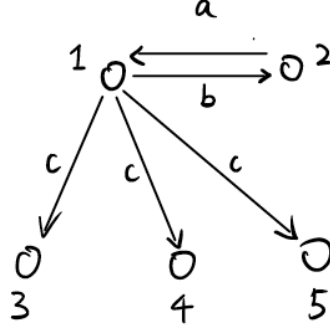
**Figure 2.4:** A regular ABC-form network: $a > b > c > 0$.

**Definition.** $G$ is a **regular ABC-form network** if

a) $G$ is an ABC-form network and

b) there is $c > 0$ such that $g_{12} > g_{21} > c$ and $g_{i1} = c$ for each $i > 2$.

In other words, a regular ABC-form network is an ABC-form network in which agent 1 exerts *strictly* higher effort than agent 2, and agent 2 exerts *strictly* higher effort than the rest. It also requires that the top agent has links with exactly the same positive weight pointing to all $i > 2$.

The following statement characterizes the optimal networks of the linear best-response model.

**Proposition 2.1.** *Consider the linear best-response model.*

*a) There is $\lambda(\bar{\eta}) < 0$ with $\lambda'(\bar{\eta}) < 0$ such that, if $v'' < \lambda(\bar{\eta})$, then all optimal networks are regular ABC-form networks.*

*b) If $v'' \geq 0$, then in every optimal network we have $g_{12} + g_{21} = \bar{\eta}$.*

The proof is provided in Appendix. The proposition says thsat, if the best-response function is linear, and the objective function is sufficiently concave, then all optimal networks are regular ABC-form networks. In contrast, if the objective function is linear or strictly convex, then in every optimal network there are two and only two agents connecting with each other. The bound on the concavity of $v$ for a regular ABC-form network to be optimal, $\lambda(\bar{\eta})$, is decreasing in $\bar{\eta}$. This reflects that, the stronger the overall complementarities, the stronger the force of

28

concentration of weight of links to just two agents.[4]

The following presents a numerical example.

**Example 2.5.** Suppose $N = \{1, 2, 3, 4\}$ and $\bar{\eta} = 1$. Consider the networks in Figure 2.5. We have the following results:

a) if $v(x_i) = x_i$ then $G_1$ is an optimal network;

b) if $v(x_i) = x_i^8$ then $G_2$ is an optimal network;

c) if $v(x_i) = -\frac{1}{x_i}$ then $G_3$ is an optimal network.



**Figure 2.5:** Optimal networks under different concavity assumptions of the objective function.

## 2.5 Conclusion

We study optimal networks under complementarities in this paper. The distinct feature of our analysis is that we allow for weighted and directed links to capture the possibilities that the planner can control the direction and intensity of complementarities. We characterize the set of optimal networks that maximize an objective function that is increasing in the effort of each agent in the network,

---

[4]We conjecture that there is a unique $\lambda(\bar{\eta}) < 0$ such that if $v'' < \lambda(\bar{\eta})$ then all optimal networks are regular ABC-form networks, and if $v'' > \lambda(\bar{\eta})$ all optimal networks involve only two agents connecting with each other. But it turns out that this is not a straightforward result.

subject to the constraint that the total weight of the links of the network does not exceed a certain amount. We find that, under rather weak conditions on the planner's objective function and the agents' best-response function, every optimal network exhibits dramatic concentration of influence such that only two agents would have outward links with positive weight which represent their influence on others.

## 2.6 Appendix

### 2.6.1 Proof of Theorem 2.1

**Proof of Lemma 2.1.** Suppose $x_i(G) \geq x_j(G)$. Consider the sequence $\mathbf{x}^t = (x_1^t, \ldots, x_n^t)$, for $t = 0, 1, 2, \ldots$, such that, for each $\ell \in N$,

$$x_\ell^0 = x_\ell(G),$$

$$x_\ell^t = \phi \left( \sum_{p \in N} g'_{\ell p} x_p^{t-1} \right) \quad \text{for } t = 1, 2, \ldots.$$

Given Assumption 2.1, we have

$$x_k^1 = \phi \left( (g_{ki} + g_{kj}) x_i(G) + \sum_{\ell \neq i, j} g_{k\ell} x_\ell(G) \right)$$

$$\geq \phi \left( g_{ki} x_i(G) + g_{kj} x_j(G) + \sum_{\ell \neq i, j} g_{k\ell} x_\ell(G) \right)$$

$$\geq x_k(G),$$

and $x_\ell^1 = x_\ell(G)$ for each $\ell \in N, \ell \neq k$. Now, for $t \geq 1$, suppose $\mathbf{x}^t \geq \mathbf{x}^{t-1}$, i.e., $x_\ell^t \geq x_\ell^{t-1}$ for each $\ell \in N$. Then, for each $\ell \in N$,

$$x_\ell^{t+1} = \phi \left( \sum_{p \in N} g'_{\ell p} x_p^t \right) \geq \phi \left( \sum_{p \in N} g'_{\ell p} x_p^{t-1} \right) = x_\ell^t.$$

Hence, $\mathbf{x}^t$ is a (weakly) increasing sequence.

Also, notice $\mathbf{x}^t = \Phi_{G'}^t(\mathbf{x}^0)$ for $t \geq 1$. Hence, by Assumption 2.3, $\lim_{t \to \infty} \mathbf{x}^t = \lim_{t \to \infty} \Phi_{G'}^t(\mathbf{x}^0) = \mathbf{x}(G')$.

Therefore, $x_\ell(G') \geq x_\ell^1 \geq x_\ell(G)$ for each $\ell \in N$.

In the case of $x_i(G) > x_j(G)$, we have

$$x_k^1 = \phi\left((g_{ki} + g_{kj})x_i(G) + \sum_{\ell \neq i,j} g_{k\ell}x_\ell(G)\right)$$

$$> \phi\left(g_{ki}x_i(G) + g_{kj}x_j(G) + \sum_{\ell \neq i,j} g_{k\ell}x_\ell(G)\right)$$

$$\geq x_k(G),$$

Therefore, $x_k(G') > x_k(G)$. *Q.E.D.*

**Proof of Lemma 2.2.** In case (a), let $G' = (g'_{ij})$ be such that $g'_{i1} = g_{i1} + g_{i2}$ and $g'_{i2} = 0$, while $g'_{pq} = g_{pq}$ for all other element in $G'$; in the case (b), let $G' = (g'_{ij})$ be such that $g'_{12} = g_{1i} + g_{12}$ and $g'_{1i} = 0$, while $g'_{pq} = g_{pq}$ for all other element in $G'$; in case (c), let $G' = (g'_{ij})$ be such that $g'_{j1} = g_{j1} + g_{ji}$ and $g'_{ji} = 0$, while $g'_{pq} = g_{pq}$ for all other element in $G'$. Then, by applying Lemma 2.1 we obtain the desired conclusions. *Q.E.D.*

**Proof of Lemma 2.3.** Consider the contrapositive of Lemma 2.2; the conclusion that $x_i(\hat{G}) = x_i(G)$ for each $i \in N$ then follows.

Now, observe that, if $G$ is not a weak ABC-form network, then one of the following must hold: (a) $g_{i2} > 0$ for some $i > 2$, or (b) $g_{ji} > 0$ for some $i > 2, j \in N$.

First, suppose that (a) $g_{i2} > 0$ for some $i > 2$, and fix that $i$. Then, by Condition (*)-c, we have $x_1(G) = x_2(G)$, and thus: $x_1(\hat{G}) = x_2(\hat{G})$, and $\hat{g}_{i1} = g_{i1} + g_{i2} > 0$.

Next, suppose that (b) $g_{ji} > 0$ for some $i > 2, j \in N$. Fix that $i$ and $j$. Then, if $j = 1$, by Condition (*)-a, we obtain $x_i(G) = x_2(G)$. Thus, $x_i(\hat{G}) = x_2(\hat{G})$ and $\hat{g}_{12} = g_{12} + g_{1i} > 0$. In the case of $j \neq 1$, we apply Condition (*)-b and obtain $x_i(G) = x_1(G)$. It follows that $x_i(\hat{G}) = x_1(\hat{G})$. Thus, $x_2(\hat{G}) = x_1(\hat{G})$, because $x_1(\hat{G}) \geq x_2(\hat{G}) \geq x_i(\hat{G})$. Hence, that $x_1(\hat{G}) = x_2(\hat{G})$ and $\hat{g}_{k1} = g_{k1} + g_{ki} > 0$ holds.

In either case, one of the claimed property for $\hat{G}$ holds. *Q.E.D.*

**Proof of Lemma 2.4.** Consider a weak ABC-form network $\hat{G} = (\hat{g}_{ij})$ with $\sum_{i,j \in N} \hat{g}_{ij} = \bar{\eta}$. First, we show that $\hat{g}_{12} > 0$. To see this, observe that we have $x_1 = \phi(\hat{g}_{12}x_2)$ in a weak ABC-form network. Hence, if $\hat{g}_{12} = 0$, then $x_1 = \phi(0) = \underline{x}$. But $x_i(\hat{G}) \leq x_1(\hat{G})$ for each $i > 2$. Thus, $x_i = \underline{x}$ for each $i \in N$ and $\hat{g}_{ij} = 0$ for each $i, j \in N$, violating $\sum_{i,j \in N} \hat{g}_{ij} = \bar{\eta}$.

Now, we examine the three cases claimed in the lemma one by one.

**Case a)** $\hat{g}_{21} = 0$.

Suppose $\hat{g}_{21} = 0$. In a weak ABC-form network, we have $x_2 = \phi(\hat{g}_{21}x_1)$. Hence, $x_2(\hat{G}) = \phi(0) = \underline{x}$, and thus $x_i(\hat{G}) \leq x_2(\hat{G}) = \underline{x}$. It follows that $\hat{g}_{i1} = 0$ for each $i \geq 2$. Thus, there is only one link with positive weight in the network: $\hat{g}_{12} = \bar{\eta}$.

Now we show that, given Assumption 2.5, a network $\hat{G}$ with $g_{12} = \bar{\eta}$ and $g_{ij} = 0$ for all the other links is not optimal. First, notice that we have $x_i(\hat{G}) = \phi(0) = \underline{x}$ for each $i > 1$, while $x_1(\hat{G}) = \phi(\bar{\eta}\phi(0)) = \bar{x}$. Now, let $\epsilon \geq 0$ and consider the network $\hat{G}^{\epsilon}$ with $\sum_{i,j \in N} \hat{g}_{ij}^{\epsilon} = \bar{\eta}$ such that $\hat{g}_{21}^{\epsilon} = \epsilon$ and $\hat{g}_{12}^{\epsilon} = \bar{\eta} - \epsilon$. Then $x_2 = \phi(\epsilon x_1)$ and $x_1 = \phi((\bar{\eta} - \epsilon)x_2)$. For brevity, denote $x_i' \equiv \frac{\partial x_i(\hat{G}^{\epsilon})}{\partial \epsilon}$. Then

$$x_2' = \phi'(\epsilon x_1)(x_1 + \epsilon x_1'), \quad x_1' = \phi'((\bar{\eta} - \epsilon)x_2)\left[(\bar{\eta} - \epsilon)x_2' - x_2\right].$$

Hence,

$$x_2'|_{\epsilon=0} = \phi'(0)\bar{x}, \quad x_1'|_{\epsilon=0} = \phi'(\bar{\eta}\underline{x})\left[\bar{\eta}\phi'(0)\bar{x} - \underline{x}\right].$$

Therefore,

$$\frac{\partial f(\mathbf{x}(\hat{G}^{\epsilon}))}{\partial \epsilon}\bigg|_{\epsilon=0} = f_1(\bar{x}, \underline{x}, \ldots, \underline{x})\phi'(\bar{\eta}\underline{x})\left[\bar{\eta}\phi'(0)\bar{x} - \underline{x}\right] + f_2(\bar{x}, \underline{x}, \ldots, \underline{x})\phi'(0)\bar{x}.$$

The above expression is strictly positive if and only if

$$\phi'(\bar{\eta}\underline{x})\bar{\eta} + \frac{f_2(\bar{x}, \underline{x}, \ldots, \underline{x})}{f_1(\bar{x}, \underline{x}, \ldots, \underline{x})} > \frac{\underline{x}}{\bar{x}}\frac{\phi'(\bar{\eta}\underline{x})}{\phi'(0)},$$

which is Assumption 2.5. It follows that, given Assumption 2.5, there is $\epsilon > 0$ such that $f(\mathbf{x}(\hat{G}^{\epsilon})) > f(\mathbf{x}(\hat{G}))$.

**Case b)** $x_i(\hat{G}) = x_2(\hat{G})$ *for some* $i > 2$.

Fix $i > 2$ with $x_i(\hat{G}) = x_2(\hat{G})$. Since $x_i(\hat{G}) = \phi(1 + \hat{g}_{i1}x_1(\hat{G}))$, $x_2(\hat{G}) = \phi(1 + \hat{g}_{21}x_1(\hat{G}))$ and $x_i(G) = x_2(G)$, we have $\hat{g}_{i1} = \hat{g}_{21} > 0$. Also, since $x_1(\hat{G}) \geq x_2(\hat{G})$, we have $\hat{g}_{12} \geq \hat{g}_{21} > 0$. Now, consider an ABC-form network $\hat{G}^\epsilon$ with $\sum_{i,j \in N} \hat{g}_{ij}^\epsilon = \bar{\eta}$ constructed as follows. Let $\epsilon \in [0, \hat{g}_{i1}]$. $\hat{G}^\epsilon$ is such that $\hat{g}_{i1}^\epsilon = \hat{g}_{i1} - \epsilon$ and $\hat{g}_{21}^\epsilon = \hat{g}_{21} + \epsilon$, while $\hat{g}_{pq}^\epsilon = \hat{g}_{pq}$ for all other elements in $\hat{G}^\epsilon$. In what follows, we show that there is $\epsilon > 0$ such that $f(\mathbf{x}(\hat{G}^\epsilon)) > f(\mathbf{x}(\hat{G}))$.

For brevity, denote $x_k' \equiv \frac{\partial x_k(\hat{G}^\epsilon)}{\partial \epsilon}$ for each $k \in N$. First, observe that, for any $\epsilon > 0$, we have $x_1(\hat{G}^\epsilon) > x_1(\hat{G})$. Hence, we also have $x_2(\hat{G}^\epsilon) > x_2(\hat{G})$, and $x_k(\hat{G}^\epsilon) \geq x_k(\hat{G})$ for each $k > 2, k \neq i$. Thus we have $x_1', x_2' > 0$, and $x_k' \geq 0$ for each $k > 2, k \neq i$. It follows that

$$
\begin{aligned}
w(\epsilon) &\equiv f_1(\mathbf{x}(\hat{G}^\epsilon))x_1' + \sum_{k>2, k\neq i} f_k(\mathbf{x}(\hat{G}^\epsilon))x_k' \\
&\geq f_1(\mathbf{x}(\hat{G}^\epsilon))x_1' \\
&> 0.
\end{aligned}
$$

Second, given $x_2 = \phi((\hat{g}_{21} + \epsilon)x_1)$ and $x_i = \phi((\hat{g}_{i1} - \epsilon)x_1)$, we have

$$
\begin{aligned}
x_2' &= \phi((g_{21} + \epsilon)x_1)[x_1 + (g_{21} + \epsilon)x_1'] \\
x_i' &= \phi((g_{21} - \epsilon)x_1)[-x_1 + (g_{21} + \epsilon)x_1'].
\end{aligned}
$$

Let $r(\epsilon) \equiv f_2(\mathbf{x}(\hat{G}^\epsilon))x_2' + f_i(\mathbf{x}(\hat{G}^\epsilon))x_i'$. Then, given $x_2(\hat{G}) = x_i(\hat{G})$, $\hat{g}_{i1} = \hat{g}_{21}$ and the symmetry assumption of $f$ (Assumption 2.4), we obtain

$$
\begin{aligned}
r(0) &= f_2(\mathbf{x}(\hat{G}))\phi(g_{21}x_1)(x_1 + g_{21}x_1') + f_i(\mathbf{x}(\hat{G}))\phi(g_{21}x_1)(-x_1 + g_{21}x_1') \\
&= 2f_2(\mathbf{x}(\hat{G}))g_{21}x_1' \\
&> 0.
\end{aligned}
$$

Therefore, $\left.\frac{\partial f(\mathbf{x}(\hat{G}^\epsilon))}{\partial \epsilon}\right|_{\epsilon=0} = w(0) + r(0) > 0$; thus, there is $\epsilon > 0$ such that

34

$f(\mathbf{x}(\hat{G}^\epsilon)) > f(\mathbf{x}(\hat{G}))$.

**Case c)** $x_1(\hat{G}) = x_2(\hat{G})$, *and* $\hat{g}_{i1} > 0$ *for some* $i > 2$.

Notice that $x_1(\hat{G}) = x_2(\hat{G})$ implies $\hat{g}_{12} = \hat{g}_{21} > 0$. Now consider an ABC-form network $\hat{G}^\epsilon \in \mathcal{G}$ with $\sum_{i,j\in N}\hat{g}_{ij}^\epsilon = \bar{\eta}$ constructed as follows. Let $\epsilon \in [0, \hat{g}_{21}]$. Let $\hat{G}^\epsilon$ be such that $\hat{g}_{12}^\epsilon = \hat{g}_{12} + \epsilon$ and $\hat{g}_{21}^\epsilon = \hat{g}_{21} - \epsilon$, while $\hat{g}_{pq}^\epsilon = \hat{g}_{pq}$ for all other elements in $\hat{G}^\epsilon$. We shall show that there is $\epsilon > 0$ such that $f(\mathbf{x}(\hat{G}^\epsilon)) > f(\mathbf{x}(\hat{G}))$.

First, denote $x_k' \equiv \frac{\partial x_k(\hat{G}^\epsilon)}{\partial \epsilon}$ for each $k \in N$. Also, let

$$r(\epsilon) \equiv f_1(\mathbf{x}(\hat{G}^\epsilon))x_1' + f_2(\mathbf{x}(\hat{G}^\epsilon))x_2' + f_i(\mathbf{x}(\hat{G}^\epsilon))x_i'.$$

Given $x_1(\hat{G}^\epsilon) = \phi((\hat{g}_{12}+\epsilon)x_2(\hat{G}^\epsilon))$, $x_2(\hat{G}^\epsilon) = \phi((\hat{g}_{21}-\epsilon)x_1(\hat{G}^\epsilon))$, $x_k(\hat{G}^\epsilon) = \phi(\hat{g}_{k1}x_1(\hat{G}^\epsilon))$, we have

$$x_1'|_{\epsilon=0} = \phi'(\hat{g}_{12}x_2(\hat{G}))(g_{12}x_2' + x_2), \quad x_2'|_{\epsilon=0} = \phi'(\hat{g}_{21}x_1(\hat{G}))(g_{21}x_1' - x_1),$$

and $x_k'|_{\epsilon=0} = \phi'(\hat{g}_{k1}x_1(\hat{G}))\hat{g}_{k1}x_1'|_{\epsilon=0}$ for each $k > 2$. Thus, $(x_1' + x_2')|_{\epsilon=0} = 0$ and $x_1'|_{\epsilon=0} > 0$. Then it follows from $x_1(\hat{G}) = x_2(\hat{G})$ and Assumption 2.4 that

$$
\begin{aligned}
r(0) =& f_1(\mathbf{x}(\hat{G}))x_1'|_{\epsilon=0} + f_2(\mathbf{x}(\hat{G}))x_1'|_{\epsilon=0} + f_i(\mathbf{x}(\hat{G}))x_i'|_{\epsilon=0} \\
=& f_1(\mathbf{x}(\hat{G}))\left(x_1' + x_2'\right)|_{\epsilon=0} + f_i(\mathbf{x}(\hat{G}))\phi'(\hat{g}_{i1}x_1(\hat{G}))\hat{g}_{i1}x_1'|_{\epsilon=0} \\
=& f_i(\mathbf{x}(\hat{G}))\phi'(\hat{g}_{i1}x_1(\hat{G}))\hat{g}_{i1}x_1'|_{\epsilon=0} \\
>& 0.
\end{aligned}
$$

Next, denote $w(\epsilon) \equiv \sum_{k>2, k\neq i} f_k(\mathbf{x}(\hat{G}^\epsilon))x_k'$. We have

$$w(0) = \sum_{k>2, k\neq i} f_k(\mathbf{x}(\hat{G}))\phi'(\hat{g}_{k1}x_1(\hat{G}))\hat{g}_{k1}x_1'|_{\epsilon=0} \geq 0.$$

Therefore, $\left.\frac{\partial f(\mathbf{x}(\hat{G}^\epsilon))}{\partial \epsilon}\right|_{\epsilon=0} = w(0) + r(0) > 0$; thus, there is $\epsilon > 0$ such that $f(\mathbf{x}(\hat{G}^\epsilon)) > f(\mathbf{x}(\hat{G}))$. *Q.E.D.*

**Proof of Theorem 2.1.**   Consider a network $G = (g_{ij}) \in \mathcal{G}$ with $\sum_{i,j \in N} g_{ij} = \bar{\eta}$.

**Case 1:** *Suppose $G$ is not a weak ABC-form network.*

If we cannot obtain a network $G'$ with $f(\mathbf{x}(G')) > f(\mathbf{x}(G))$ by switching a link using the one-link-switch principle, then $G$ is not optimal. By Lemma 2.2, if we cannot obtain a network $G'$ with $f(\mathbf{x}(G')) > f(\mathbf{x}(G))$ using the one-link-switch principle, then condition (*) holds for $G$. But then, by Lemma 2.3, given condition (*), we can find a weak ABC-form network $\hat{G}$ with $\sum_{i,j \in N} \hat{g}_{ij} = \bar{\eta}$ such that $f(\mathbf{x}(\hat{G})) = f(\mathbf{x}(G))$. Moreover, the network $\hat{G}$ is such that either (a) $x_i(\hat{G}) = x_2(\hat{G})$ for some $i > 2$ and $\hat{g}_{12} > 0$, or (b) $x_1(\hat{G}) = x_2(\hat{G})$ and $\hat{g}_{i1} > 0$ for some $i > 2$. By Lemma 2.4, this implies that there is an ABC-form network $\hat{G}^\epsilon = (\hat{g}_{ij}^\epsilon)$ with $\sum_{i,j \in N} \hat{g}_{ij}^\epsilon = \bar{\eta}$ such that $f(\mathbf{x}(\hat{G}^\epsilon)) > f(\mathbf{x}(\hat{G}))$. Given $f(\mathbf{x}(\hat{G})) = f(\mathbf{x}(G))$, we thus obtain $f(\mathbf{x}(\hat{G}^\epsilon)) > f(\mathbf{x}(G))$. Hence, $G$ is not optimal.

**Case 2:** *Suppose $G$ is already a weak ABC-form network.*

Then let $\hat{G} = G$. Since $G$ is not an ABC-form network but a weak one, we have $\hat{g}_{12} = 0$, or $\hat{g}_{21} = 0$, or $x_i(\hat{G}) = x_2(\hat{G})$ for some $i > 2$. However, by Lemma 2.4, for a weak ABC-form network with $\sum_{i,j \in N} \hat{g}_{ij} = \bar{\eta}$, we have $\hat{g}_{12} > 0$. And, if $\hat{g}_{21} = 0$, or $x_i(\hat{G}) = x_2(\hat{G})$ for some $i > 2$, then there is an ABC-form network $\hat{G}^\epsilon = (\hat{g}_{ij}^\epsilon)$ with $\sum_{i,j \in N} \hat{g}_{ij}^\epsilon = \bar{\eta}$ such that $f(\mathbf{x}(\hat{G}^\epsilon)) > f(\mathbf{x}(\hat{G}))$. Hence, $G$ is not optimal. *Q.E.D.*

## 2.6.2   Proof of Proposition 2.1

Proposition 2.1 follows from Lemma 2.5 and Lemma 2.6 below.

**Lemma 2.5.** *Consider the linear best-response model. Then there is a threshold $\lambda(\bar{\eta}) < 0$, with $\lambda'(\bar{\eta}) < 0$, such that if $v'' < \lambda(\bar{\eta})$, then all optimal networks are regular ABC-form networks.*

**Proof of Lemma 2.5.**   Consider the linear best-response model, and $v'' < 0$. Suppose that $G \in \mathcal{G}$ is an optimal network. By Theorem 2.1, all optimal networks are the ABC-form networks. Hence $G$ is an ABC-form network. We establish our

conclusion through the following steps.

**Step 1:** *$G$ is such that $g_{i1} = g_{j1}$ for each $i, j > 2$ and $i \neq j$.*

Consider $i, j > 2$, $i \neq j$. Given $G$ is an ABC-form network, and $\phi(y) = 1 + y$, we have $x_i = 1 + g_{i1}x_1$ and $x_j = 1 + g_{j1}x_1$. Consider a network $G'$ with $g'_{i1} + g'_{j1} = g_{i1} + g_{j1}$ while $g'_{pk} = g_{pk}$ for all other elements in $G'$. Since $v'' < 0$, $G'$ maximizes $v(x_i(G')) + v(x_j(G'))$ within the constraint $g'_{i1} + g'_{j1} = g_{i1} + g_{j1}$ if and only if $g'_{i1} = g'_{j1} = (g_{i1} + g_{j1})/2$. Hence, if $g_{i1} \neq g_{j1}$, then there is a network $G' \in \mathcal{G}$ such that $\sum_{k \in N} v(x_k(G')) > \sum_{k \in N} v(x_k(G))$.

**Step 2:** *Suppose $G$ is such that $g_{i1} = 0$ for each $i > 2$. Then, given $v'' < 0$, we have $g_{12} = g_{21} = \frac{\bar{\eta}}{2}$.*

Consider a network $G \in \mathcal{G}$ such that $g_{12} = a$ and $g_{21} = \bar{\eta} - a$. Then $x_1 = 1 + ax_2$ and $x_2 = 1 + (\bar{\eta} - a)x_1$. Thus

$$x_1(a) = \frac{1 + a}{1 - a(\bar{\eta} - a)}, \quad x_2(a) = \frac{1 - a + \bar{\eta}}{1 - a(\bar{\eta} - a)},$$

and

$$x_1(a) + x_2(a) = \frac{2 + \bar{\eta}}{1 - a(\bar{\eta} - a)}.$$

Notice $a = \frac{\bar{\eta}}{2}$ is the maximum of $a(\bar{\eta} - a)$, and thus the minimum of the denominator of the expression above. Hence, $a = \frac{\bar{\eta}}{2}$ is the maximum of $x_1(a) + x_2(a)$. Given $v'' < 0$, we thus have, for each $a \neq \frac{\bar{\eta}}{2}$:

$$\frac{1}{2}v(x_1(a)) + \frac{1}{2}v(x_2(a)) \leq v(\frac{1}{2}x_1(a) + \frac{1}{2}x_2(a))$$
$$< v(x_1(\frac{\bar{\eta}}{2}))$$
$$= \frac{1}{2}v(x_1(\frac{\bar{\eta}}{2})) + \frac{1}{2}v(x_2(\frac{\bar{\eta}}{2})).$$

Hence, $a = \frac{\bar{\eta}}{2}$ is the maximum of $v(x_1(a)) + v(x_2(a))$.

**Step 3:** *Suppose $G$ is such that $g_{i1} = c$ for each $i > 2$. Then there is $\lambda(\bar{\eta}) < 0$, with $\lambda'(\bar{\eta}) < 0$, such that, if $v'' < \lambda(\bar{\eta})$, then $c > 0$.*

Denote $x^* = x_1(\frac{\bar{\eta}}{2})$, which is the effort of agent 1 as well as agent 2 when

$g_{12} = g_{21} = \frac{\bar{\eta}}{2}$ and $g_{ij} = 0$ for all other links. Consider an ABC-form network $G \in \mathcal{G}$ such that $g_{i1} = c$ for each $i > 2$. In what follows, we show that, if $v'(1) > \frac{2}{2-\bar{\eta}} v'(x^*)$, then we must have $c > 0$. Notice $x^* > 1$. Hence, we can always find a sufficiently negative $\lambda$ so that, given $v'' < \lambda$, we have $v'(1) > \frac{2}{2-\bar{\eta}} v'(x^*)$. Moreover, note that the right hand side of the inequality, $\frac{2}{2-\bar{\eta}} v'(x^*)$, is increasing in $\bar{\eta}$. Hence, for greater $\bar{\eta}$, we need a tighter (more negative) $\lambda$. Our claim then follows.

We now show that, if $v'(1) > \frac{2}{2-\bar{\eta}} v'(x^*)$, then we have $c > 0$. Consider the network $G$ such that $g_{12} = \frac{\bar{\eta}}{2}$, $g_{21} = \frac{\bar{\eta}}{2} - (n-2)c$, and $g_{i1} = c$ for each $i > 2$. For brevity, denote $a = \frac{\bar{\eta}}{2}$. Then $x_1 = 1 + ax_2$, $x_2 = 1 + [a - (n-2)c] x_1$, and $x_i = 1 + cx_1$ for each $i > 2$. Again for brevity, let $x_i' \equiv \frac{\partial x_i(G)}{\partial c}$ for each $i \in N$. Then we have

$$x_1' = ax_2'$$
$$x_2'|_{c=0} = \frac{(n-2)}{a^2 - 1} x^*$$
$$x_3' = cax_2' + x_1.$$

Thus,

$$
\begin{aligned}
\frac{\partial \sum_{i \in N} v(x_i)}{\partial c}\bigg|_{c=0} &= v'(x^*)ax_2'|_{c=0} + v'(x^*)x_2'|_{c=0} + (n-2)v'(1)x^* \\
&= v'(x^*)(a+1)\frac{(n-2)}{a^2-1} x^* + (n-2)v'(1)x^* \\
&= (n-2)x^* \left[ v'(x^*)\frac{a+1}{a^2-1} + v'(1) \right] \\
&= (n-2)x^* \left[ v'(x^*)\frac{1}{a-1} + v'(1) \right].
\end{aligned}
$$

Substituting with $a = \frac{\bar{\eta}}{2}$, we obtain

$$\frac{\partial \sum_{i \in N} v(x_i)}{\partial c}\bigg|_{c=0} > 0 \text{ if and only if } v'(1) > \frac{2}{2-\bar{\eta}} v'(x^*).$$

**Step 4**: *Suppose $G$ is such that $g_{i1} = c > 0$ for each $i > 2$. Then $g_{12} > g_{21} >$*

*c.*

By Theorem 2.1, we have $g_{21} > c$ (otherwise $x_2 \leq x_3$, but then $G$ would not be an ABC-form network). Now suppose $g_{12} = g_{21} = a > c$, and consider an alternative network $G^\epsilon$ with $g_{12}^\epsilon = a + \epsilon$, $g_{21}^\epsilon = a - \epsilon$. For brevity, denote $x_i' \equiv \frac{\partial x_1(G^\epsilon)}{\partial \epsilon}$ for each $i \in N$, and let $x^*$ be the effort of 1 given $g_{12} = g_{21} = a$. Then, given $x_1 = 1 + (a+\epsilon)x_2, x_2 = 1 + (a-\epsilon)x_1$, and $x_i = 1 + cx_1$ for $i > 2$, we obtain

$$(x_1' + x_2')|_{\epsilon=0} = 0$$

and

$$x_1'|_{\epsilon=0} = \frac{x_2 - (a+\epsilon)x_1}{1 - (a+\epsilon)(a+\epsilon)}\Big|_{\epsilon=0} = \frac{(1-a)x^*}{1-a^2} = \frac{x^*}{1+a} > 0.$$

Therefore,

$$
\begin{aligned}
\frac{\partial \sum_{i \in N} v(x_i)}{\partial \epsilon}\bigg|_{\epsilon=0} &= v'(x^*)x_1'|_{\epsilon=0} + v'(x^*)x_2'|_{\epsilon=0} + (n-2)v'(x_3)cx_1'|_{\epsilon=0} \\
&= v'(x^*)(x_1' + x_2')|_{\epsilon=0} + (n-2)v'(x_3)cx_1'|_{\epsilon=0} \\
&= (n-2)v'(x_3)cx_1'|_{\epsilon=0} \\
&> 0
\end{aligned}
$$

implying $g_{12} > g_{21}$ at optimum, which completes the proof. *Q.E.D.*

**Lemma 2.6.** *Consider the linear best-response model. If $v'' \geq 0$, then in every optimal network we have $g_{12} + g_{21} = \bar{\eta}$.*

**Proof of Lemma 2.6.** Suppose $v'' \geq 0$, and that $G$ is an optimal network. By Theorem 2.1, all optimal networks are ABC-form networks. Hence $G$ is an ABC-form network, such that $g_{12} \geq g_{21} > g_{i1}$ for each $i > 2$. We show that $g_{i1} = 0$ for each $i > 2$, and therefore we have $g_{12} + g_{21} = \bar{\eta}$. Suppose that, without loss of generality, $g_{31} > 0$. Then consider an alternative network $G'$ with $g_{31}' = 0$ and $g_{21}' = g_{21} + g_{31}$, while $g_{kp}' = g_{kp}$ for all other links. Given $v'' \geq 0$, we certainly have $\sum_{k \in N} v(x_k(G')) > \sum_{k \in N} v(x_k(G))$. To see this, observe $x_2 = 1 + (g_{21} + g_{31})x_1$,

$x_3 = 1$, and $x_1 = 1 + g_{12}x_2$. Hence, $x_2$ increases while $x_3$ might decrease, but the increase of $x_2$ is at least as great as the reduction of $x_i$. Meanwhile, $x_1$ also increases due to the increase of $x_2$. But then all $x_i$ with $i > 3$ increase as well. Altogether, and given $v'' \geq 0$, we have $\sum_{k \in N} v(x_k(G')) > \sum_{k \in N} v(x_k(G))$. *Q.E.D.*

# Chapter 3

# Norm-based Resentment and the Evolution of Cooperative Norms

## 3.1   Introduction

It is well known that the British have strong expectation for everyone to take queues in every possible instance (and instances not possible to form a queue, e.g., waiting at the bus stop alone). Queuing is an example of many social rules governing the *cooperation in one-shot interactions* between individuals in Britain: if an individual jumps a queue, it benefits the individual but harms all the ones behind him; thus, sticking to a queue is a cooperative action that benefits others at a personal cost. However, that the British always resort to queuing is not due to any formal regulation imposed by a central authority. Instead, it is because queuing is a *social norm*, one that is "shared by other people and partly sustained by their approval and disapproval... sustained by the feelings of embarrassment, anxiety, guilt and shame that a person suffers at the prospect of violating them" (Elster, 1989, p. 99-100).

This paper presents a model of the social evolution of cooperation based on the idea that cooperation in one-shot interactions is supported by the endogenous social norms of cooperation. That is, individuals might cooperate because they

might perceive that a norm of cooperation exists in the society. When a norm of cooperation does exists, people feel the emotion of resentment against, and thus punish, those who defect. The explicit role of emotions in our model departs from most of the existing evolutionary models of cooperation. Most existing models do not make an explicit distinction between social evolution from biological evolution. Instead, they view social evolution as a complete analogue to biological evolution (e.g., Abramson and Kuperman, 2001; Eshel et al., 1998; Nowak and May, 1992; Nowak et al., 2010; Ohtsuki, 2006; Santos and Pacheco, 2005; Szabó and Fáth, 2007).[1] This paper shows that analyzing the evolution of cooperation from a social norms' perspective matters. More specifically, I develop a dynamic model of the norm of cooperation that explains the regularities. The model shows how cooperation and punishment of defectors co-evolve. It reveals the conditions under which cooperation emerges and persists in the long run. The model explains naturally why there is a positive correlation between cooperation and the quality of law enforcement, as documented by Herrmann et al. (2008), Gächter and Schulz (2016) and Tabellini (2008). The model also predicts that the level of cooperation should be higher in societies with higher mobility, which is in sharp contrast with the prediction of a typical model that assume the behavioral rule of *imitating the best*, i.e., individuals have a tendency to imitate the behavior of those earning higher material payoffs (e.g., Abramson and Kuperman, 2001; Eshel et al., 1998; Nowak and May, 1992; Nowak et al., 2010; Ohtsuki, 2006; Santos and Pacheco, 2005; Szabó and Fáth, 2007).

The model builds on what I call **norm-based resentment** (Sugden, 1984, 2000, 2004; Bicchieri, 2006; Cooper and Dutcher, 2011; Falk et al., 2006; Herz and Taubinsky, 2017; Kahneman et al., 1986; Peysakhovich and Rand, 2016).[2] The idea is that individuals have empirical expectations on the behavior of others in a society. If an individual expects that most others cooperate, but he meets a

---

[1] Exceptions are Boyd and Richerson (1985) and Henrich and Boyd (1998, 2001), who investigate cultural evolution models that assume a behavioral rule of conformity, i.e., imitating the most frequent behavior in the population.

[2] I use **bold** font for new definitions and *italic* for highlights in this paper.

defector, then the individual is frustrated: the defector is acting unkindly to him. To release his frustration, the individual may punish the defector. However, if the individual holds the expectation that defection is commonplace, then defection is acceptable—that is what everyone else does anyway. In this case, there is no impulse to punish the defector. Now suppose that a population of individuals are randomly matched into pairs to play the following **stage game**. One individual moves first and decides whether to cooperate or defect. The other moves second and decides whether to punish the first-mover if the first-mover defects. Given norm-based resentment, there are two (sequential) equilibria: **the defection equilibrium** in which each individual in the population defects *and* does not punish defectors, and **the cooperation equilibrium** in which each individual cooperates *and* punishes defectors.[3] I define **social norm** as a profile of expectations of individuals consistent with one of the equilibria.

My goal is to analyze how behavior and expectations of individuals evolve over time. The dynamics are adapted from Young (1993, 2001). Individuals are randomly matched to play the stage game recurrently and infinitely over discrete time periods. In each period, each individual forms expectations about others' behavior based on strategies used in the last period. Given expectations, each individual plays a best-response with a high probability; best-responses are constructed on norm-based resentment. With a small probability, individuals make mistakes, in which case they randomly pick a strategy. I call the dynamics **adaptive dynamics with norm-based resentment**. A **population state** of the dynamics specifies the strategies *and* the expectations of individuals in the population. The defection equilibrium and the cooperation equilibrium are two different population states. I examine which population state is more likely to emerge and persist in the long run. Formally, I characterize the **stochastically stable equilibrium**: it is the *most likely* population state in an infinite span of time when the probability of making a mistake goes to zero (Ellison, 1993, 2000; Kandori et al., 1993; Young,

---

[3]The model keeps track of expectations (beliefs) explicitly; thus, I apply the solution concept of sequential equilibrium.

1993, 2001).

The basic result is that, given norm-based resentment, the cooperation equilibrium can be stochastically stable. Whether and when this is the case depends on two statistics: the **intolerance of defection**, defined by the maximum proportion of defectors that is consistent with punishing defectors, and the **temptation to defect**, defined by the minimum proportion of punishers (i.e., those who punish defectors) that is required to induce cooperation. The intolerance of defection is increasing in the psychological parameter of resentment and the harm that punishment generates, and decreasing in the cost of conducting punishment. The temptation to defect is increasing in the individual cost of cooperation, and decreasing in the suffering from being punished. If the intolerance of defection is greater than the temptation to defect, then the cooperation equilibrium is the unique stochastically stable equilibrium, otherwise the defection equilibrium is the unique stochastically stable equilibrium.

The second part of the paper investigates two key differences between large, modern societies and small-scale societies. In small-scale societies, an individual mostly interacts with his relatives, which is a fixed, small subset of individuals in the population. In contrast, in large societies with higher market integration, $i$) individuals interact with larger groups of people, and $ii$) they can choose where to live. First, I compare **global interactions** with **local interactions**: in global interactions, every individual interacts with everyone else in the population; in local interactions, by contrast, each individual only interacts with a small subset of others in the population.[4] I obtain a neutrality result: whether interactions are local or global does not affect whether the cooperation equilibrium is stochastically stable. The reason is that whether interactions are local or global affects neither side of the trade-off between the intolerance of defection and the temptation to defect.

However, *mobility*—the ability to migrate—leads to to two effects, both sug-

---

[4]More precisely, for local interactions, individuals are located on a two-dimensional lattice, and each individual only interacts with the four direct neighbors around them.

gesting that there are positive correlations between the norm of cooperation and the size of a society. One is the migration effect in the intermediate run. Since cooperation leads to higher aggregate efficiency, individuals will move from societies with the defection norm to the societies with the cooperation norm. Hence, societies with the cooperation norm *become* larger. The other effect is about the selection of norms in the long run. In the long run the norms in a society might change. Then what matters is the relative difficulty of transiting from the defection norm to the cooperation norm compared to transitions in the opposite direction. Again, mobility matters: since societies with the cooperation norm are larger, it is more difficult—i.e., it requires a larger number of mistakes—to overthrow a society with the cooperation norm than to overthrow a society with the defection norm. The key to the argument is the **fitting-in effect** generated by norm-based resentment: when migrating to a society, individuals adjust their behavior and expectations to make them compatible with the prevailing norms in the society.

The result on the positive effect of mobility stands in contrast with previous analysis on the evolution of cooperation (Abramson and Kuperman, 2001; Eshel et al., 1998; Nowak and May, 1992; Nowak et al., 2010; Ohtsuki, 2006; Santos and Pacheco, 2005; Szabó and Fáth, 2007). Previous analysis considers the social evolution of cooperation as an analogue to biological evolution. It assumes the behavioral rule of **imitating the best**, i.e., individuals have a tendency to imitate the behavior of those earning higher material payoffs. Imitating the best implies that interactions in small neighborhoods and without mobility—the small-scale societies—are the ideal setting for cooperation to emerge. The reason is that, in local interactions with fixed matching, cooperators can form clusters and separate themselves from defectors. As a result, cooperators earn higher material payoffs; their behavior is therefore imitated. However, if mobility is possible, then defectors can move to the center of the clusters of cooperators and exploit the cooperators. Hence, defectors always earn higher material payoffs. Eventually, if mobility is

possible, a dynamic of imitating the best leads to universal defection.[5] Hence, a dynamic of imitating the best does not explain why cooperation and punishment of defectors are higher in large, modern societies. The difference between a dynamic of imitating of best and our model is that the fitting-in effect present in our model is missing in a dynamic of imitating of best.

Finally, I examine the relationship between cooperation and law enforcement. I extend our model to analyze the co-evolution of cooperative norms and the quality of law enforcement. The analysis shows that there could only be two population states to be stochastically state: one in which everyone cooperates and punishes defectors, and there is high quality of law enforcement, and the other in which everyone defects, no one punishes defectors, and the quality of law enforcement is low. This result explains the observed correlation between cooperation and the quality of law enforcement across societies (Gächter and Schulz, 2016; Herrmann et al., 2008; Tabellini, 2008).

The paper is organized as follows. Section 3.2 introduces norm-based resentment and defines the stage game that individuals play recurrently in each period. Section 3.3 presents and analyzes an adaptive dynamic without mistakes, the "unperturbed" adaptive dynamic. The results for the unperturbed dynamic form the basis for later analysis. Section 3.4 presents the adaptive dynamics where mistakes are possible, the "perturbed" adaptive dynamics, and characterizes the stochastically stable equilibrium. All analyses so far assume global interactions. Section 3.5 examines local interactions and the effect of mobility. Section 3.6 considers the existence of a central monitor who enforces laws.

---

[5]Introducing the opportunity to punish defectors would not change the conclusion. The reason is that, even the cooperation norm is initially established in a society and the initial residents punish defectors, the new comers would cooperate but they would not punish. Eventually, the new comers "dilute" the norm to an extent such that everyone would rather defect. Axelrod's (1986) simulation shows this "drift".

## 3.2 The stage game and norm-based resentment

### 3.2.1 The stage game

Let $N = \{1, 2, \ldots, n\}$ be a population of individuals. $n$ is an even number. Think of individuals in $N$ as generic members of a society who do not know each other in person. The individuals are randomly matched into pairs to play the following game. Within each matched pair, one individual moves first and the other moves second; the identify of the first-mover is randomly determined. They then play the game in Figure 3.1. The first-mover decides between *cooperate* and *defect*. If the first-mover cooperates, the game ends. If the first-mover defects, the second-mover can *punish* the first-mover, or *not punish*. Figure 1 gives the *material payoffs* for the first-mover and the second-mover, respectively.

**Assumption 3.1.** $\bar{a} > a > \underline{a}$ *and* $\bar{b} > b > \underline{b}$.
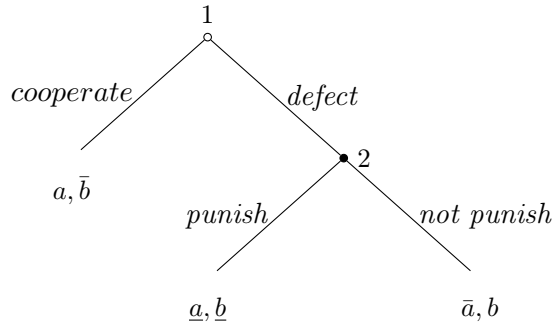


**Figure 3.1:** The stage game. The numbers at the terminal nodes represent *material payoffs*.

Since the identity of the first-mover is randomly determined, a strategy for $i \in N$ is a pair $s_i = (x_i, y_i) \in \{1, 0\} \times \{1, 0\}$: if $i$ cooperates in the role of the first-mover, $x_i = 1$, otherwise $x_i = 0$; if $i$ punishes in the role of the second-mover, $y_i = 1$, otherwise $y_i = 0$. Let $s = (s_1, s_2, \ldots, s_n)$ be the strategy profile of all individuals, and $S = \{1, 0\}^{2n}$ denote the set of all strategy profiles. I call the game **stage game**. In the next section, I introduce a discrete dynamic, $t = 0, 1, 2, \ldots$, such that individuals play the stage game recurrently over time.

For now, let us focus on the stage game and clarify the meaning of norm-based resentment in the specific context of the base game. If the second-mover purely seeks to maximize material payoffs, she does not punish. Anticipating this, and if the first-mover also seeks to maximize material payoffs, she defects. However, a large body of public goods game experiments show that many participants punish defectors at the expense of their own earnings (Fehr and Schmidt, 2000). Why punish? And why does the tendency to punish varies across societies as observed by Herrmann et al. (2008) and Henrich et al. (2010; 2006)? According to Sugden (2000; 2004) and Bicchieri (2006), individuals punish defectors because defection is *socially unacceptable*, i.e., considered as a violation of a social norm. Critically, Sugden (2000; 2004) and Bicchieri (2006) suggest that people consider defection as *more* unacceptable when cooperation is common, compared to when many others defect—what I call norm-based resentment. However, Sugden and Bicchieri have not provided an formal account of the punishment mechanism.[6] I provide a formulation and investigate its implications here.

**The second-mover's decision.** I call those individuals who choose to cooperate as first-movers **cooperators**. Consider individual $i \in N$. Let $q_i \in [0, 1]$ be $i$'s expectation of the proportion of cooperators in the society. Suppose that $i$ is matched with $j$, and $i$ is in the role of the second-mover. In $i$'s eyes, $j$ is not different from any others in the society. Hence, $i$ expects her payoffs to be bounded *below* by

$$\theta(q_i) = q_i \bar{b} + (1 - q_i)\underline{b}.$$

The expected payoff, $\theta(q_i)$, is $i$'s **reference payoffs**. Since $\bar{b} > \underline{b}$, $\theta(q_i)$ is strictly increasing in $q_i$.

Suppose now that $j$ defects. In this case the payoff for $i$ is bounded *above* by

---

[6]The models by the two authors that are closest to ours are Sugden (2000) and Bicchieri (2006, Ch. 6).

b. Comparing the payoff $b$ with $\theta(q_i)$, $i$ is certainly worse-off if

$$q_i > \frac{b - \underline{b}}{\overline{b} - \underline{b}}.$$

The fact that $i$ obtains less than her reference payoffs makes $i$ frustrated. To release her frustration, $i$ exhibits the preferences to *act unkind to unkind*—a tendency analogous to Rabin's (1993) and Dufwenberg and Kirchsteiger's (2004) models of reciprocity. In particular, I focus on the negative domain of reciprocity. A justification for this is that, empirically, "hurting hurts more than helping helps" (Offerman, 2002), i.e., it appears that people reciprocate others' hurtful choices much more often than reciprocating others' helpful choices.[7]

More precisely, by defecting, $j$ produces $\theta(q_i) - b$ unkindness to $i$. In response to $j$'s defection, if $i$ punishes $j$, $i$ is unkind to $j$; the unkindness is $\frac{\bar{a}-\underline{a}}{2}$. If $i$ does not punish $j$, $i$ is kind to $j$; the kindness is $\frac{\bar{a}-\underline{a}}{2}$. Let $\lambda > 0$ be the **resentment parameter**; it measures $i$'s frustration given $j$'s unkindness. Conditional on $j$ defecting, and $b < \theta(q_i)$, if $i$ chooses to punish, her *utility* is

$$\underline{b} + \lambda \left[ b - \theta(q_i) \right] \left( \frac{\underline{a} - \bar{a}}{2} \right); \tag{3.1}$$

if $i$ instead chooses not to punish, her utility is

$$b + \lambda \left[ b - \theta(q_i) \right] \left( \frac{\bar{a} - \underline{a}}{2} \right). \tag{3.2}$$

Hence, $i$ prefers to punish if $i$'s reference payoffs are sufficiently high:

$$\theta(q_i) > b + \left( \frac{1}{\lambda} \right) \frac{b - \underline{b}}{\bar{a} - \underline{a}}.$$

---

[7]On the other hand, it will be clear that no any qualitative result derived in this chapter relies on the reciprocity-type specification of punishment. All results would hold as long as there is a threshold such that individuals punish defectors if and only if the proportion of cooperators is higher than the threshold.

The above condition can be expressed by

$$q_i > \pi(\bar{a}, \underline{a}, b, \bar{b}, \underline{b}, \lambda) \quad \text{where} \quad \pi(\bar{a}, \underline{a}, b, \bar{b}, \underline{b}, \lambda) = \left[ 1 + \frac{1}{\lambda(\bar{a} - \underline{a})} \right] \frac{b - \underline{b}}{\bar{b} - \underline{b}}. \quad (3.3)$$

$\pi(\bar{a}, \underline{a}, b, \bar{b}, \underline{b}, \lambda)$ is the threshold proportion of cooperators needed to activate punishment towards defectors. Note that $q_i > \pi(\bar{a}, \underline{a}, b, \bar{b}, \underline{b}, \lambda)$ implies $\theta(q_i) > b$. Hence, $q_i > \pi(\bar{a}, \underline{a}, b, \bar{b}, \underline{b}, \lambda)$ is a sufficient and necessary condition for $i$ to strictly preferring to punish a defector.[8] I call the preferences specified above **norm-based resentment**.

The threshold $\pi(\bar{a}, \underline{a}, b, \bar{b}, \underline{b}, \lambda)$ for punishment has the following properties. First, it is increasing in punishment cost $(b - \underline{b})$. Second, it is decreasing in the resentment parameter $(\lambda)$. It is also decreasing in the effectiveness of punishment $(\bar{a} - \underline{a})$ and the payoff difference between meeting a defector and meeting a cooperator $(\bar{b} - b)$. Note that, if $\lambda$ is sufficiently small, or $\bar{a} - \underline{a}$ is small, then we have $\pi(\bar{a}, \underline{a}, b, \bar{b}, \underline{b}, \lambda) > 1$. In that case, the condition for $i$ to punish a defector is never satisfied regardless of her expectation $q_i$. The sufficient and necessary condition for $\pi(\bar{a}, \underline{a}, b, \bar{b}, \underline{b}, \lambda) < 1$ is

$$\lambda > \left( \frac{1}{\bar{a} - \underline{a}} \right) \frac{b - \underline{b}}{\bar{b} - b}.$$

In Rabin's and Dufwenberg and Kirchsteiger's original models, the reference point to determine the first-mover's kindness and unkindess is *independent* of $i$'s expectation on the behavior of other individuals. This is the distinction between their models and ours.

**The first-mover's decision.**    Now, consider the decision of the first-mover. Each individual chooses her first-mover decision to maximize expected material

---

[8]I can also introduce an individual-specific resentment parameter $\lambda_i > 0$ such that the utility for choosing to punish is $\underline{b} + \lambda_i \left[ b - \theta(q_i) \right] \left( \frac{\underline{a} - \bar{a}}{2} \right)$, and the utility for choosing not to punish is $b + \lambda_i \left[ b - \theta(q_i) \right] \left( \frac{\bar{a} - \underline{a}}{2} \right)$. Since I do not analyze the effects of the heterogeneity in $\lambda_i$ across individuals in this paper, I normalize $\lambda_i = \lambda$ to simplify the exposition.

payoffs. This means that, if enough individuals punish defectors, individuals prefer to cooperate as first-movers, otherwise they defect. I call the individuals who choose to punish as second-movers **punishers**. Let $p_i \in [0,1]$ be $i$'s expectation of the proportion of punishers in the society. Suppose $i$ is the first-mover. If $i$ cooperates, she gets material payoffs $a$. If $i$ defects, her expected material payoffs are

$$p_i \underline{a} + (1 - p_i)\bar{a}.$$

Hence, $i$ prefers to cooperate if

$$p_i > \varphi(a, \bar{a}, \underline{a}) \quad \text{where} \quad \varphi(a, \bar{a}, \underline{a}) = \frac{\bar{a} - a}{\bar{a} - \underline{a}}. \tag{3.4}$$

If $p_i < \varphi(a, \bar{a}, \underline{a})$, $i$ prefers to defect. If $p_i = \varphi(a, \bar{a}, \underline{a})$, $i$ is indifferent between the two options. $\varphi(a, \bar{a}, \underline{a})$ is increasing in temptation to defect, $\bar{a} - a$, while decreasing in damage suffered from punishment, $\bar{a} - \underline{a}$. Given $\underline{a} < a < \bar{a}$, the range of $\varphi$ is $(0, 1)$.

The thresholds $\pi(\bar{a}, \underline{a}, b, \bar{b}, \underline{b}, \lambda)$ and $\varphi(a, \bar{a}, \underline{a})$ are important determinants of the long-run dynamics of the population. To simply notations, I drop the arguments and use $\pi$ and $\varphi$ to denote their values in subsequent analyses.

### 3.2.2 Equilibria of the stage game

I characterize the set of *strict* equilibria of the stage game in this subsection. I always have an equilibrium in which each $i \in N$ defects and does not punish defectors and, for some parameter values, another equilibrium in which each $i \in N$ cooperates and punishes defectors.

Let $q = (q_1, \ldots, q_n) \in [0, 1]^n$ be the individuals' expectations on the proportion of cooperators in the society. Let $p = (p_1, \ldots, p_n) \in [0, 1]^n$ be the individuals' expectations on the proportion of punishers. I call the pair $(p, q)$ **expectation profile**. A **population state** is a tuple $(s, q, p) \in S \times [0, 1]^n \times [0, 1]^n$; it consists of the strategy profile $s$ and the expectation profile $(q, p)$. Let $Z = S \times [0, 1]^n \times$

$[0,1]^n$ be the collection of all population states. We apply the solution concept of sequential equilibrium (Kreps and Wilson, 1982) to our specific context.[9]

**Definition.** A (sequential) **equilibrium** of the stage game is a population state $(s, q, p) \in Z$ such that the following two conditions hold:

1. the expectation profile $(q, p)$ is consistent with $s$, namely, for each $i \in N$,

$$q_i = \frac{1}{n-1} \sum_{j \in N, j \neq i} x_i \quad \text{and} \quad p_i = \frac{1}{n-1} \sum_{j \in N, j \neq i} y_j;$$

2. their strategies are **sequentially rational** given their expectations, namely, for each $i \in N$,

$$x_i = \begin{cases} 1 & \text{if } p_i > \varphi \\ 0 & \text{if } p_i < \varphi \end{cases}, \qquad y_i = \begin{cases} 1 & \text{if } q_i > \pi \\ 0 & \text{if } q_i < \pi \end{cases}. \tag{3.5}$$

An equilibrium of the stage game, $(s, q, p)$, is **strict** if no individual is indifferent between any two choices at any decision node of the stage game. In our model, a strict equilibrium is one in which for each $i \in N$, we have $p_i \neq \varphi$ and $q_i \neq \pi$. I focus on strict equilibria because, as I shall show, they are the only steady states of the society in the long run.

In particular, let $s^C \in S$ denote the strategy profile in which each $i \in N$ cooperates and punishes defectors, and $s^D \in S$ denote the strategy profile in which each $i \in N$ defects and does not punish defectors. With abuse of notation, I use $(s^C, 1, 1) \in Z$ to denote the population state in which the strategy profile is $s^C$ and each $i \in N$ holds the expectations $p_i = 1$ and $q_i = 1$. Likewise, let $(s^D, 0, 0) \in Z$ denote the population state in which the strategy profile is $s^D$ and each $i \in N$ holds the expectations $p_i = 0$ and $q_i = 0$. **Generic cases of the stage game** are such that, for the thresholds $\varphi$ and $\pi$, we have either $\varphi + \pi \geq \frac{n}{n-1}$

---

[9]Fudenberg and Tirole (1991) provide a discussion about the relations between sequential equilibrium and perfect Bayesian equilibrium.

or $\varphi + \pi \leq \frac{n-2}{n-1}$. The following proposition identifies $(s^C, 1, 1)$ and $(s^D, 0, 0)$ as two strict equilibria in the generic cases of the stage game.[10]

**Proposition 3.1.** *Consider the generic cases of the stage game.*

1. $(s^D, 0, 0)$ *is a strict equilibrium.*

2. *If $\pi < 1$, then $(s^C, 1, 1)$ is also a strict equilibrium.*

3. *No other strict equilibrium exists.*

*Proof.* All proofs are provided in the appendix. $\qquad\square$

I define a **social norm** as an expectation profile $(q, p) \in [0, 1]^n \times [0, 1]^n$ in an strict equilibrium $(s, q, p) \in Z$ of the game that takes norm-based resentment into account. Hence, social norms are *self-fulfilling expectations* such that, once the expectations are established among individuals, individuals take strategies consistent with them, which validates their initial expectations. In our context, there exist two social norms that can be established. One supports cooperation and motivates people to punish defectors: $p_i = 1$ and $q_i = 1$ for each $i \in N$. The other is associated with defection, leaving defectors unpunished: $p_i = 0$ and $q_i = 0$ for each $i \in N$. Both are self-fulfilling.

When $\pi \geq 1$, the only strict equilibrium is the defection one. I assume the following throughout our subsequent analyses, so that multiple strict equilibria exist.

**Assumption 3.2.** $\pi(\bar{a}, \underline{a}, b, \bar{b}, \underline{b}, \lambda) < 1$.

In the next section, I analyze which equilibrium is more likely to emerge and persist in a society.

---

[10]Note that $\lim_{n \to \infty} \frac{n}{n-1} = \lim_{n \to \infty} \frac{n-2}{n-1} = 1$. Hence, the generic cases essentially require that $\varphi + \pi$ is not exactly equal to 1.

## 3.3 The unperturbed adaptive dynamics

This section introduces the unperturbed adaptive dynamics to study the evolution of norms. The dynamics I consider are adapted from Young (1993). The departure from Young (1993) is that the dynamics I consider explicitly keep track of the expectation profile over time alongside the strategy profile. Investigating the dynamic enhances our understanding regarding the determinants of the emergence and persistence of social norms.

### 3.3.1 The adaptive dynamic

Let time unfold in discrete units and be indexed by $t = 0, 1, 2, \ldots$. I use superscript to denote variables in period $t$. The strategy of $i \in N$ in period $t$ is denoted by $s_i^t = (x_i^t, y_i^t)$. The strategy profile of the population in period $t$ is

$$s^t = \left( s_1^t, s_2^t, \ldots, s_n^t \right) = \left( (x_1^t, y_1^t), (x_2^t, y_2^t), \ldots, (x_n^t, y_n^t) \right).$$

I call $z_i^t = (s_i^t, p_i^t, q_i^t)$ $i$'s **individual state** at $t$, and $z^t = (s^t, p^t, q^t)$ is the **population state** or simply the **state** at $t$. $Z = S \times [0,1]^n \times [0,1]^n$ is the collection of all population states. Note that the state variable $z^t$ not only keeps track of the strategy profile $s^t$ of the population, but also explicitly keeps track of the expectation profile $(p^t, q^t)$. In what follows, I describe how $z^t$ evolves over time. Roughly speaking, in each period $t$, individuals form expectations regarding the proportions of cooperators and punishers based on the strategy profile of the previous period, $s^{t-1}$. Given their expectations, individuals play *myopic best-responses* to maximize their utilities at each decision node of the stage game in the current period. The best-responses are myopic in the sense that individuals do not consider using their current actions to influence payoffs from future periods. In addition, individuals may draw small and possibly biased samples from the previous population state when forming their expectations. I also allow for *inertia behavior* (discussed in greater detail below). These elements improve the model's generality and allow
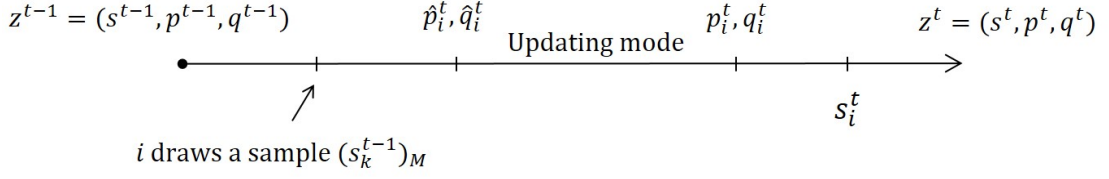
us to obtain more results.



$$z^{t-1} = (s^{t-1}, p^{t-1}, q^{t-1}) \qquad \hat{p}_i^t, \hat{q}_i^t \qquad\qquad p_i^t, q_i^t \qquad\qquad z^t = (s^t, p^t, q^t)$$

Updating mode

$s_i^t$

$i$ draws a sample $(s_k^{t-1})_M$

**Figure 3.2:** Timeline of the updating procedure.

More precisely, let $z^0 \in Z$ be the **initial population state**. Consider $t \geq 1$, and let the population state at $t-1$ be $(s^{t-1}, p^{t-1}, q^{t-1})$. Let $M$ be a subset of individuals, i.e., $M \subset N$. A **sample of the strategy profile** $s^{t-1} = (s_1^{t-1}, s_2^{t-1}, \dots, s_n^{t-1})$ is a sub-sequence $(s_j^{t-1})_{j \in M}$ of it. The **size of the sample**, denoted by $m$, is defined by the number of individuals in $M$. Let $\mathcal{M}$ be the collection of all subsets of $N$ whose size is $m$, i.e., $\mathcal{M}$ contains all those $M \subset N$ with $|M| = m$. To form expectations $p_i$ and $q_i$ in the current period, each $i \in N$ draws a sample $(s_j^{t-1})_{j \in M}$ of size $m$. The samples are drawn randomly and independently across individuals. By counting how many cooperators and punishers there are in the sample, each $i$ obtains estimates about the proportions of cooperators and punishers in the population:

$$\hat{q}_i^t = \frac{1}{m} \sum_{j \in M} x_j^{t-1}, \qquad \hat{p}_i^t = \frac{1}{m} \sum_{j \in M} y_j^{t-1}.$$

These estimates form the basics of $i$'s expectations in period $t$. The implicit assumption here is that individuals can observe the complete strategies of the individuals in their sample—i.e., the sample individuals' cooperation decisions as well as their punishment decisions—in the previous period.

Individuals' behavior may exhibit **inertia**, i.e., there is a positive probability that some individuals may behave in the same way as the past. I model inertia behavior by allowing some individuals to *not* update their expectations. There are four **updating modes**:

1. $q_i^t = \hat{q}_i^t$ and $p_i^t = \hat{p}_i^t$;

2. $q_i^t = \hat{q}_i^t$ and $p_i^t = p_i^{t-1}$;

3. $q_i^t = q_i^{t-1}$ and $p_i^t = \hat{p}_i^t$;

4. $q_i^t = q_i^{t-1}$ and $p_i^t = p_i^{t-1}$.

The updating mode for each individual is determined randomly and independently across individuals and time. Each updating mode occurs with positive probability for each individual at each period.

More precisely, let $\sigma_i(p_i, q_i | s^{t-1}, p^{t-1}, q^{t-1})$ denote the probability that $i$ has expectations $(p_i, q_i)$ conditional on previous state $(s^{t-1}, p^{t-1}, q^{t-1})$. I assume that

$$\sigma_i(p_i, q_i | s^{t-1}, p^{t-1}, q^{t-1}) > 0$$

if, and only if, the following conditions hold:

1. if $p_i \neq p_i^{t-1}$ then there is $M \in \mathcal{M}$ such that $p_i = \frac{1}{m} \sum_{k \in M} y_k^{t-1}$;

2. if $q_i \neq q_i^{t-1}$ then there is $M \in \mathcal{M}$ such that $q_i = \frac{1}{m} \sum_{k \in M} x_k^{t-1}$.

If the above conditions do not hold, then $\sigma_i(p_i, q_i | s^{t-1}, p^{t-1}, q^{t-1}) = 0$. Note that, if $p_i = p_i^{t-1}$, or if $q_i = q_i^{t-1}$, then the above conditions place no restriction on the conditional probability $\sigma_i$; in this case, we again have $\sigma_i(p_i, q_i | s^{t-1}, p^{t-1}, q^{t-1}) > 0$. Intuitively, in this case $i$ is in the mode of *not* updating at least one side of her expectations, which occurs with positive probability.[11]

Now I describe the probability of transiting from one population state to another. Conditional on the previous population state $z^{t-1}$, let $P_i(z_i^t | z^{t-1})$ denote the probability that $i \in N$ is in state $z_i^t = (s_i^t, p_i^t, q_i^t)$ in period $t$. According to our previous definition, the individual state $z_i^t = (s_i^t, p_i^t, q_i^t)$ is sequentially rational if and only if it satisfies condition (3.5). Let $P_i(z_i^t | z^{t-1})$ assign positive probability

---

[11]It is not necessary for $\sigma_i(.|.)$ to be the same across individuals. The analysis only requires them to be stationary (invariant) across time.

only to sequentially rational states. That is,

$$P_i(z_i^t|z^{t-1}) = \begin{cases} \sigma_i(p_i^t, q_i^t|s^{t-1}, p^{t-1}, q^{t-1}) & \text{if } s_i^t \text{ is sequentially rational w.r.t. } (p_i^t, q_i^t), \\ 0 & \text{otherwise.} \end{cases}$$

(3.6)

For each two population states $z^t, z^{t-1} \in Z$, let $P^0(z^t, z^{t-1})$ denote the probability of transiting from state $z^{t-1}$ to state $z^t$. We have

$$P^0(z^t, z^{t-1}) = \prod_{i \in N} P_i(z_i^t|z^{t-1}).$$

$P^0$ defines a Markov chain on finite state space $Z$. Following Young (1993), I call $P^0$ **unperturbed adaptive dynamic with sample size** $m$. "Unperturbed" refers to the assumption that individuals always play sequentially rational strategies given their expectations. Nevertheless, the path $\{z^t\}_{t=0}^{\infty}$ is *not* deterministic due to the uncertainty involved in the process of forming expectations.

## 3.3.2 Emergence of a social norm

I now examine whether the unperturbed adaptive dynamic necessarily converges to one of the equilibria described by Proposition 3.1. In other words, will a social norm—cooperation or defection—necessarily become established in the society? The answer is yes.

The theorem below establishes that the adaptive dynamic $P^0$ converges almost surely either to a population state in which every $i \in N$ cooperates and punishes defectors and this is commonly expected, *or*, to a state in which every $i \in N$ defects and no one punishes defectors. The convergence occurs no matter how disordered the initial population state might be and regardless of the sample size $m$. Note that when $m$ is small relative to $n$, different individuals may draw completely different samples from the past and form different expectations, leading to different strategies. Hence, it is not obvious how a social norm can always

become established spontaneously.

As before, I use $(s^C, 1, 1) \in Z$ to denote the population state in which each $i \in N$ cooperates and punishes defectors and each $i \in N$ holds the expectations $p_i = 1$ and $q_i = 1$. And $(s^D, 0, 0) \in Z$ is the opposite population state in which each $i \in N$ defects and does not punish defectors and each $i \in N$ expects $p_i = 0$ and $q_i = 0$. Let $\{z^t\}_{t=0}^{\infty} \subset Z$ be a sequence of *random* variables such that: *i*) $z^0 \in Z$, and *ii*) each $z^t$ with $t \geq 1$ is generated according to the probability system $P^0$, i.e., for each $t \geq 1$ and $z^t, z^{t-1} \in Z$, we have $\text{Prob}\{z^t|z^{t-1}\} = P^0(z^t, z^{t-1})$. Let $L \subset Z$ be a subset of population states. We say that $P^0$ **converges almost surely** to $L$ if, for *each* $z^0 \in Z$, the event

$$\left( \lim_{t \to \infty} z^t \right) \in L$$

has probability one.

**Theorem 3.1.** *For every sample size $m$ with $1 \leq m \leq n$, $P^0$ converges almost surely to $\{(s^C, 1, 1), (s^D, 0, 0)\}$.*

The proof of the theorem builds on a general observation regarding Markov chains on finite state spaces. For a Markov chain, a population state $z \in Z$ is an **absorbing state** of the dynamic if, once reaching $z$, it stays at the state for all future periods with probability one. *First*, observe that the absorbing states of the adaptive dynamic $P^0$ are exactly the strict equilibria of the stage game: $(s^C, 1, 1)$ and $(s^D, 0, 0)$. The reason for this is that, starting from any state *not* a strict equilibrium of the stage game, some individual $i$ would change her strategy with positive probability in the next period. *Next*, observe that, since the state space $Z$ is finite, there is a positive probability $\delta > 0$ that the dynamic $P^0$ transits from any initial state to $(s^C, 1, 1)$ or $(s^D, 0, 0)$ within a finite number of periods. Let that finite number of periods be $T$. It follows that the probability of *not* transiting to $\{(s^C, 1, 1), (s^D, 0, 0)\}$ within $kT$ periods is at most $(1 - \delta)^{kT}$, which shrinks to zero as $k$ gets large.

The theorem says that the adaptive dynamic converges almost surely to a strict equilibrium of the stage game. However, there are two distinct strict equilibria—one in which everyone cooperates and the other in which everyone defects. I investigate the determinants of which equilibrium the dynamic converges to in the next subsection.

### 3.3.3 Selection of norms in the intermediate run

I identify the factors that affect which equilibrium the unperturbed dynamic $P^0$ converges to in this subsection. We discuss two factors. First, where the dynamic starts from, *the initial population state*, matters. Intuitively, if the dynamic starts from a state "close" enough to one of the two strict equilibria, the dynamic will converge to that equilibrium almost surely. This leads to the notion of basin of attraction. Second, when forming expectations, what samples individuals draw from the past matter. I discuss some of the implications.

For a population state $z \in Z$, let $B(z) \subset Z$ denote the **basin of attraction of $z$ under the unperturbed dynamic** $P^0$, which is defined as follows. $B(z)$ is a subset of population states such that, if the dynamic $P^0$ starts from a population state in $B(z)$, it converges almost surely to $z$. The (relative) sizes of the basins of attraction of the two equilibria—$(s^C, 1, 1)$ and $(s^D, 0, 0)$—are critical in determining which equilibrium the dynamic tends to converge to. To simplify notations, I use $z^C = (s^C, 1, 1)$ and $z^D = (s^D, 0, 0)$ to denote the two equilibria of the stage game. Let $\mathbb{1}\{.\}$ be the indicator function, which takes the value of 1 if the statement within the parenthesis holds, and zero otherwise. For example, $\mathbb{1}\{p_i \geq \varphi\} = 1$ if and only if $i$ holds the expectation $p_i \geq \varphi$. Then

$$\sum_{i \in N} \mathbb{1}\{p_i \geq \varphi\}$$

counts the number of individuals in the population who expect the proportion of individuals, $p_i$, to be at least $\varphi$.

Proposition 3.2 below characterizes the basins of attraction of the defection equilibrium $z^D$ and the cooperation equilibrium $z^C$, respectively. Notice that we have requirements on individuals' *actions*, $x_i$ and $y_i$, as well as their *expectations*, $p_i$ and $q_i$. The basin of attraction of the defection equilibrium $z^D$ contains those states in which *a)* the number of individuals who cooperate or punish defectors is small, and *b)* the number of individuals having the expectations $p_i \geq \varphi$ or $q_i \geq \pi$ is small. In contrast, the basin of attraction of the cooperation equilibrium $z^C$ contains those states in which *a)* the number of individuals who cooperate and punish defectors is large, and *b)* the number of individuals having the expectations $p_i > \varphi(a, \bar{a}, \underline{a})$ and $q_i > \pi$ is large.

**Proposition 3.2.** *Consider the unperturbed adaptive dynamic $P^0$ with sample size $m$.*

1. *We have $(s, p, q) \in B(z^D)$ if and only if the strategy profile $s = (x, y)$ and the expectation profile $(p, q)$ are such that*

$$\sum_{i \in N} x_i < \pi m, \quad \sum_{i \in N} y_i < \varphi m, \tag{3.7}$$

   *and*

$$\sum_{i \in N} \mathbb{1}\{p_i \geq \varphi\} < \pi m, \quad \sum_{i \in N} \mathbb{1}\{q_i \geq \pi\} < \varphi m.$$

2. *We have $(s, p, q) \in B(z^C)$ if and only if the strategy profile $s = (x, y)$ and the expectation profile $(p, q)$ are such that*

$$\sum_{i \in N} x_i > n - (1 - \pi)m, \quad \sum_{i \in N} y_i > n - (1 - \varphi)m, \tag{3.8}$$

   *and*

$$\sum_{i \in N} \mathbb{1}\{p_i > \varphi\} > n - (1 - \pi)m, \quad \sum_{i \in N} \mathbb{1}\{q_i > \pi\} > n - (1 - \varphi)m.$$

Figure 3.3 illustrates the conditions on the strategy profile $(x, y)$ of the basins of
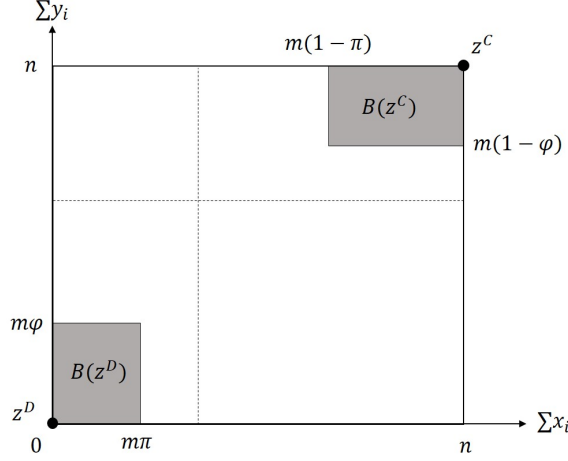
**Figure 3.3:** The bottom-left corner is the defection equilibrium $z^D$ and the upper-right corner is the cooperation equilibrium $z^C$. The gray areas display the basins of attraction of the two equilibria under unperturbed adaptive dynamic $P^0$ with sample size $m$.

attraction of the two equilibria. The $x$-axis represents the number of cooperators in the population, and the $y$-axis represents the number of punishers. When the population is in the defection equilibrium $z^D$, we have $\sum x_i = \sum y_i = 0$. In contrast, in the cooperation equilibrium $z^C$, we have $\sum x_i = \sum y_i = n$. The gray areas show the basins of attraction of the two equilibrium in terms of strategies. We have the following comparative statics. First, the area of the basin of attraction of the defection equilibrium, $B(z^D)$, is increasing in $\pi$ and $\varphi$. By contrast, the area of the basin of attraction of the cooperation equilibrium, $B(z^C)$, is decreasing in $\pi$ and $\varphi$. Hence, lower values of $\pi$ and $\varphi$ favor the emergence of the cooperation equilibrium.

Second, the smaller the sample size $m$ when forming expectations, the smaller the areas of both $B(z^D)$ and $B(z^C)$. In fact, if the sample size $m$ is much smaller than the population size $n$, then the basins of attraction of the two equilibria only consist of a small number of population states among all population states. The reason for this is that, if sample size is small, there is positive probability that the individuals draw samples that are not representative of the proportion of cooperators and punishers in the whole population. That is, they may exhibit **small sample bias**. In that case, the equilibrium that the adaptive dynamic converges to depends on the exact samples that the individuals draw along the

61

way. The small sample bias may explain why governments around the world often invest considerable efforts in advertising altruistic behavior of national heroes and other role models. By exploiting the small sample bias, they manipulate people's expectations on the population state and thereby affect social norms. The small sample bias also has other implications. For instance, government officials are often seen as more "visible" individuals in the society, either because their behaviors are revealed more often to the public than an average citizen, or people choose to pay more attention to these government officials. It follows that, if government officials engage in corruption activities, it not only affects the functioning of formal institutions of the country, but also reduce individuals' willingness to punish defectors by negatively affecting their perceptions on the population state. This latter effect on informal institutions further increases the incentives to corrupt. While causality is not clear, evidence from the cross-cultural experiments on dishonesty is consistent with our observeation. For example, using data from a lying game experiment, Gächter and Schulz (2016) show that, in countries where misconducts of government officials such as corruption and fraudulent politics are more pervasive, participants exhibit higher dishonesty.

## 3.4 The perturbed adaptive dynamics and stochastic stability

In this section, I investigate a perturbed version of the adaptive dynamic. In a perturbed dynamic, individuals may make mistakes and deviate from rational strategies. The motivations of investigating a perturbed dynamic are the following. First, it is not realistic to assume that individuals always play rationally in the long run. Some individuals may "tremble", or deliberately experiment new strategies. As previously shown by Young (1993; 2001), Kandori et al. (1993), and Ellison (1993; 2000), these mistakes play a key role in determining which equilibrium is more likely in the long run. Second, without knowing the initial population state

and the exact samples that the individuals draw along the way, it is difficult to predict which equilibrium the unperturbed dynamic converges to. However, if mistakes are possible, the influence of the initial population state eventually dies out. The theory then generates sharper results regarding the long-run trend of the system.

### 3.4.1    The perturbed dynamics

We introduce the perturbed dynamic in this subsection. How individuals form expectations is the same as in the unperturbed dynamic. However, under the perturbed dynamic, each individual makes a mistake with probability $\varepsilon > 0$ such that each of the four pure strategies in the stage game (cooperate or not as a first-mover and punish or not as a second-mover) is taken with positive probability.

To simplify exposition, assume that, conditional on the event of $i$ making a mistake, $i$ takes each of the four strategies with probability $\frac{1}{4}$.[12] Let $J \subset N$ be a subset of individuals, and let $|J|$ denote the number of individuals in $J$. Making a mistake or not is a random event independent across individuals and time. Hence, within a certain period $t$, the probability that the individuals in $J$ make mistakes while others play rationally is $\varepsilon^{|J|}(1-\varepsilon)^{n-|J|}$. Let $Q(z^t, z^{t-1}, J)$ be the probability of transiting from state $z^{t-1} \in Z$ to a certain state $z^t \in Z$ in the next period, *conditional* on that exactly the individuals in $J$ make mistakes. We have

$$Q(z^t, z^{t-1}, J) = \left[ \left( \frac{1}{4} \right)^{|J|} \prod_{i \in J} \sigma_i(p_i^t, q_i^t | s^{t-1}, p^{t-1}, q^{t-1}) \right] \left[ \prod_{i \notin J} P_i(z_i^t | z^{t-1}) \right].$$

The conditional probabilities $\sigma_i(.,.|.,.,.)$ and $P_i(.|.)$ are the ones defined in Section 3.3.1. Within the first square bracket of the above expression is the probability that individuals in $J$ have the expectations and strategies that $z^t$ specifies. Notice that individuals only make mistakes in the sense that their strategies are

---

[12]This assumption is not necessary. We can even allow for each individual $i$ to draw strategies based on different distributions. All our subsequent theorems hold. The only requirement is that the distributions to randomly choose strategies are invariant across time.

not rational with respect to their expectations. However, they do not invent new expectations without a foundation. The function $\sigma_i(.,.|.,.,.)$ requires every individual to either update expectations based on the strategy profiles $s^{t-1}$, or have the same expectations as the ones in the previous period (depending on the updating mode).

Let $P^\varepsilon(z^t, z^{t-1})$ denote the probability of transiting from $z^{t-1}$ in period $t-1$ to $z^t$ in period $t$. Summing over different subsets of individuals who make mistakes, we obtain:

$$P^\varepsilon(z^t, z^{t-1}) = (1-\varepsilon)^n P^0(z^t, z^{t-1}) + \sum_{J \subset N, J \neq \emptyset} (\varepsilon)^{|J|} (1-\varepsilon)^{n-|J|} Q(z^t, z^{t-1}, J).$$

$P^0(z^t, z^{t-1})$ is transition probability under the unperturbed adaptive dynamic. We call $P^\varepsilon$ with $\varepsilon > 0$ **the perturbed dynamic**. $P^\varepsilon$ is also a finite Markov process with $Z$ as state space.

Observe from the expression of $P^\varepsilon(z^t, z^{t-1})$ that, as the probability of making a mistake $\varepsilon$ goes to zero, the second term on the right-hand side goes to zero. Hence, $P^\varepsilon$ converges to $P^0$ as $\varepsilon$ goes to zero. However, there is a qualitative difference between the unperturbed dynamic $P^0$ and the perturbed one $P^\varepsilon$ with $\varepsilon > 0$. That is, under the unperturbed dynamic, once the society falls into a strict equilibrium of the stage game, it locks in that state. In contrast, under the perturbed dynamic, there is always a positive probability of the society transiting from any state to any another one within finite periods.[13] Hence, even if a norm—i.e., a strict equilibrium of the stage game—is established in the society, there is a positive probability of escaping from it and tipping to another population state.

We shall characterize the *most frequently visited* population state by the dynamic $P^\varepsilon$ in the long run. When the probability of making a mistake $\varepsilon$ is small, the most frequent population state is called stochastically stable equilibrium (Young, 1993). Formally, let $\mu^\varepsilon : Z \to [0,1]$ denote the stationary distribution of $P^\varepsilon$ with

---

[13]That is, with $\varepsilon > 0$, $P^\varepsilon$ is a finite, irreducible Markov process.

$\varepsilon > 0$.[14] A unique stationary distribution exists for all finite irreducible Markov processes. Hence $\mu^\varepsilon$ is well-defined.

**Definition.** The population state $z \in Z$ is a **stochastically stable equilibrium (SSE)** if

$$\lim_{\varepsilon \to 0} \mu^\varepsilon(z) = 0.$$

### 3.4.2 Selection of norms in the long run

I characterize SSE in this subsection. Building on Young (1993), I show that, for small $\varepsilon$, which equilibrium—$z^C$ or $z^D$—is the most frequent state visited by $P^\varepsilon$ in the long run is determined by the numbers of mistakes required to leave $B(z^C)$ and $B(z^D)$, the basins of attraction of $P^0$. The reason is that, once reaching $z^C$ or $z^D$, the only way that the population transits to another equilibrium is to have individuals making enough number of mistakes to escape from the basin of attraction of the established equilibrium. When $\varepsilon$ is small, mistakes are rare events. Hence, once reaching $z^C$ or $z^D$, the population will stay within the basin of attraction of the established equilibrium for a very long span of time. The number of mistakes required to escape from $B(z^C)$ compared with the number of mistakes required to escape $B(z^D)$ then determines which equilibrium is the most frequent one in an infinite span of time.

Now, let $R(z^C)$ denote the minimum number of mistakes required to escape from $B(z^C)$, given that the dynamic $P^\varepsilon$ starts from $z^C$. Analogously, let $R(z^D)$ be the minimum number of mistakes required to escape from $B(z^D)$, given that the dynamic $P^\varepsilon$ starts from $z^D$. For a real number $v$, $\lceil v \rceil$ denotes the smallest integer equal to or greater than $v$. We have the following lemma about $R(z^C)$ and $R(z^D)$.

**Lemma 3.1.** *Consider the unperturbed adaptive dynamic $P^0$ with sample size $m$.*

---

[14]$P^\varepsilon$ is actually a transition probability matrix with dimension $|Z| \times |Z|$. Its element $P^\varepsilon(z',z)$ specifies the probability of transiting from state $z$ to state $z'$, with $z, z' \in Z$. Let $[P^\varepsilon]^T$ denote the $T$th power of the matrix $P^\varepsilon$. Let $v \in \mathbb{R}_+^{|Z|}$ be $|Z|$-dimentional probability distribution vector such that its elements sum to 1. Then the stationary distribution of $P^\varepsilon$ can be defined by $\mu^\varepsilon \equiv \lim_{T \to \infty} [P^\varepsilon]^T v$.

1. $R(z^D) = \lceil m \min\{\pi, \varphi\} \rceil$.

2. $R(z^C) = \lceil m \min\{1 - \pi, 1 - \varphi\} \rceil$.

If $R(z^C) > R(z^D)$, then it is more difficult to leave $B(z^C)$ than to leave $B(z^D)$, and thus $z^C$ is the unique SSE. Conversely, if $R(z^C) < R(z^D)$, then $z^D$ is the unique SSE. This leads to the following theorem, which characterizes SSE for all generic cases.

**Theorem 3.2.** *Consider the adaptive dynamic $P^\varepsilon$ with sample size $m$. Consider the generic cases with $\lceil \varphi m \rceil \neq \lceil (1 - \pi)m \rceil$ and $\lceil \pi m \rceil \neq \lceil (1 - \varphi)m \rceil$.*

1. *If $1 - \pi < \varphi$, then $z^D$ is the unique SSE.*

2. *If $1 - \pi > \varphi$, then $z^C$ is the unique SSE.*

*Remark.* If $\lceil \varphi m \rceil = \lceil (1 - \pi)m \rceil$ and $\lceil \pi m \rceil = \lceil (1 - \varphi)m \rceil$, then both $z^D$ and $z^C$ are SSEs, and no other SSE exists.

The theorem says that, if the sum of the two thresholds, $\pi + \varphi$, is greater than 1, then the unique stochastically stable equilibrium is the defection equilibrium $z^D$. In this equilibrium, each $i \in N$ defects and does not punish defectors. Moreover, each $i \in N$ holds the expectations $q_i = 0$ and $p_i = 0$, i.e., everyone expects everyone else to defect and not punish defectors. By contrast, if the sum of $\pi$ and $\varphi$ is less than 1, then the unique stochastically stable equilibrium is the cooperation equilibrium $z^C$. In this equilibrium, everyone cooperates and punishes defectors, as well as expecting all others to cooperate and punish defectors. This comparative static result is not surprising. Recall that $\pi$ is the threshold such that, if $i$ expects the proportion of cooperators is greater than $\pi$, $i$ punishes defectors. $\varphi$ is the threshold such that, if $i$ expects the proportion of punishers in greater than $p_i$, $i$ cooperates. The formulas for the two thresholds are given by equations (3.3) and (3.4). It is both intuitive and following from our previous characterization of basins of attraction (Proposition 3.2) that lower values of the two thresholds favor the emergence and persistence of the cooperation equilibrium.

Nevertheless, the above theorem makes a new observation: the technologies that shape $\pi$ and the technologies that shape $\varphi$ are perfect substitutes. It does not require that the values of $\pi$ and $\varphi$ are both low enough; it is sufficient to have the sum of them being low enough.

Expressing the condition in terms of underlying parameters, we obtain

$$\left[1 + \frac{1}{\lambda(\bar{a} - \underline{a})}\right] \frac{b - \underline{b}}{\bar{b} - \underline{b}} + \frac{\bar{a} - a}{\bar{a} - \underline{a}} < 1.$$

This leads to the following comparative statics. The social norm of cooperation tends to emerge and persist in the long run if the resentment parameter $\lambda$ is large, the loss of meeting a defector $\bar{b} - \underline{b}$ is large, punishment cost $b - \underline{b}$ is small, the damage of punishment $\bar{a} - \underline{a}$ is large, or, the temptation to defect $\bar{a} - a$ is small.

The observation that the cooperation equilibrium can ever be stochastically stable is not trivial. Consider instead the popular cultural evolution model based on *simple conformity* (Boyd and Richerson, 1985; Henrich and Boyd, 1998, 2001). Simple conformity means that individuals have a tendency to adopt the strategy that is most frequently used in the population. If material incentives favor a different strategy from the most frequently used strategy, then individuals face trade-offs between material incentives and conformity (Henrich and Boyd, 2001). This model is often seen as a legitimate reduced-form model of social norms transmission: it abstracts away specific psychology and emotions that are considered important in sustaining social norms. Instead, it models individuals as machines programmed to put decision weights on popular cultural traits. However, simple conformity implies that the the defection equilibrium is always the unique SSE, regardless of the strength of conformity. The reason is that conformity loses its power when half of population cooperate while the other half defect. In this case, material incentives dictate choices of individuals in favor of defection and no punishment. As a result, the defection equilibrium always has a larger basin of attraction than the cooperation equilibrium.[15]

---

[15]More precisely, the defection equilibrium is always the unique $\frac{1}{2} - dominant\ equilibrium$ (El-

Hence, the conformity model leads to a qualitatively different conclusion from the one reached by our model that takes explicitly the psychology of norm-based resentment into account. The difference suggests that, if a researcher does believe that "Social norms have a grip on the mind that is due to the strong emotions they can trigger" (Elster, 1989, p. 99-100), then he should model these emotions explicitly rather than relying on the conformity model, for they generally make different predictions. In *The Grammar of Society*, Bicchieri (2006) argues that an essential component of social norms is individuals' *normative expectations*, possibly with *sanctions*, such that they expect other people to expect them to conform to certain behavior. In other words, to show the existence of a social norm, it is not sufficient to observe that individuals have the preferences to conform to what others do. We must also show that individuals expect others to expect them to behave in a certain way. If they fail to fulfill such expectations of others, they are sanctioned, or suffer from negative social emotions such as shame and guilt as Elster (1989) emphasizes. According to Wrong (1961), the idea that sanctions and social emotions are essential in sustaining social orders at least dates back to Durkheim. However, these authors emphasize the importance of sanctions and social emotions as a fact. They have not yet shown *why* sanctions and social emotions are important. Bicchieri shows that the presence of sanctions and social emotions transform a prisoner's dilemma game into a coordination game and, thereby, make cooperation among unrelated individuals possible. However, simple conformity would do the job. Why guilt, shame, and sanctions? Here I show that sanctions and social emotions matter because they lead to qualitatively different long-run dynamics from that of simple conformity.

---

lison (2000)): suppose that half of the population defect and do not punish defectors; then the best-response of everyone is to defect and leave the defectors unpunished. By Ellison (2000), a $\frac{1}{2}$−dominant equilibrium is a SSE.

## 3.5 Local interactions and mobility

This section examines $i$) local interactions, i.e., each individual only interacts with a small subset of others in the population, and $ii$) the effects of mobility. By conducting a series of experiments in 15 diverse societies around the world, Henrich et al. (2010) show that individuals exhibit stronger tendency for cooperation and are more willing to punish selfish behaviors in large societies with higher degree of market integration. This section aims at providing a stylized analysis to Henrich et al. (2010)'s findings by examining the long-run evolution of societies. I show that the ability to vote with feet is critical in understanding the relationships between cooperation and community size. An interesting observation is that it may not be the case that the larger size of a society leads to a higher level of cooperation. Instead, that a society is larger may be due to the emergence of cooperative norms in the society.

### 3.5.1 Local interactions



**Figure 3.4:** Local interactions



**Figure 3.5:** The updating procedure under local interactions.

For tractability, I follow Ellison (2000) and consider the following local interaction structure. Let individuals in $N$ locate at the vertexes of a two-dimentional lattice on the surface of a torus (see Figure 3.4). Let $n_1$ and $n_2$ be two integers no less than

three. Let $\tilde{N} = \{1, 2, \ldots, n_1\} \times \{1, 2, \ldots, n_2\}$ be a discrete coordinate system to indicate the vertexes of the lattice. Each coordinate pair $i = (l, k) \in \tilde{N}$ represents an individual in the population. Two individuals $i, j \in \tilde{N}$ are **neighbors** if the distance between them is exactly one:

1. $|i - j| = 1$, or

2. $i = (k, \ell)$ and $j = (k', \ell')$ are such that $k = 1$, $k' = n_1$ and $\ell = \ell'$, or

3. $i = (k, \ell)$ and $j = (k', \ell')$ are such that $\ell = 1$, $\ell' = n_2$ and $k = k'$.

Hence, each individual has exactly four neighbors. Every individual in $\tilde{N}$ only observes and react to the strategies of his four neighbors. Let $N_i \subset \tilde{N}$ denote the set of $i$'s four neighbors. The **adaptive dynamics $P^0$ and $P^\epsilon$ in local interactions** are such that, when forming estimates about the proportions of cooperators and punishers at the current period, each $i \in \tilde{N}$ pays attention precisely to her neighbors. That is, for each $t \geq 1$ and $i \in \tilde{N}$, we have

$$\hat{q}_i^t = \frac{1}{4} \sum_{j \in N_i} x_j^{t-1}, \qquad \hat{p}_i^t = \frac{1}{4} \sum_{j \in N_i} y_j^{t-1}.$$

Now, I characterize SSE under local interactions. As shown by the graph on the right of Figure 3.6, when $\pi + \varphi$ is sufficiently small, the unique SSE is the cooperation equilibrium $z^C$ in which each $i \in \tilde{N}$ cooperates and punishes defects and expect everyone else to do so. By contrast, when $\pi + \varphi$ is sufficiently large, the unique SSE is the defection equilibrium $z^D$ in which each $i \in \tilde{N}$ defects and does not punish defectors and expect everyone else to defect and not punish defectors. The white squares on the line $\pi + \varphi = 1$ in the graph are non-generic cases. The non-generic cases occupy a non-negligible area; this is due to the integer issue arising from that $\hat{q}_i^t$ and $\hat{p}_i^t$ can only take values from $\{0, \frac{1}{4}, \frac{2}{4}, \frac{3}{4}, 1\}$.

**Theorem 3.3.** *Consider the adaptive dynamics in local interactions. Consider the cases with $\lceil 4\varphi \rceil \neq \lceil 4(1 - \pi) \rceil$ and $\lceil 4\pi \rceil \neq \lceil 4(1 - \varphi) \rceil$.*

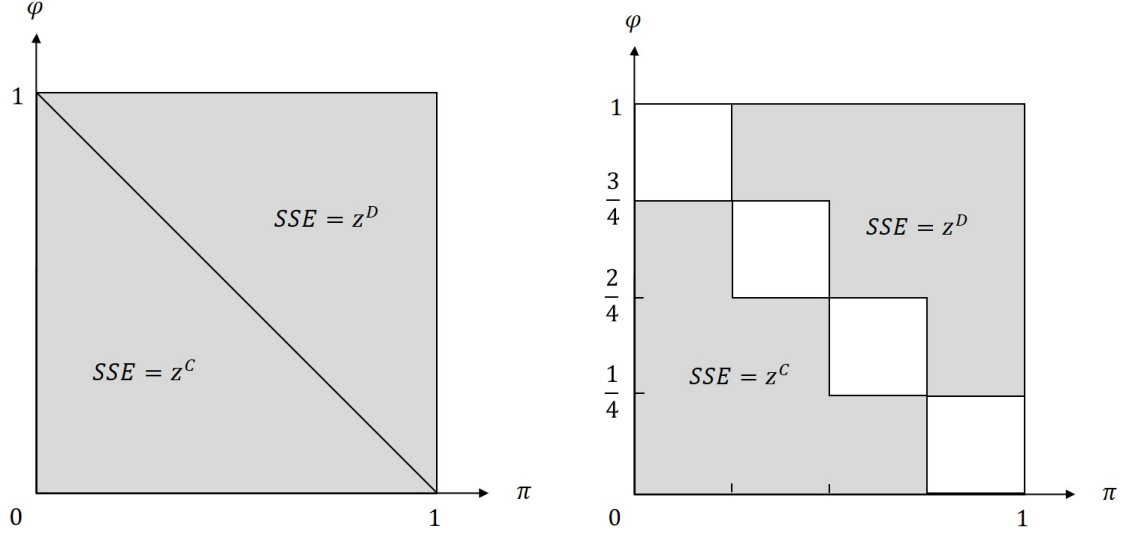1. *If $1 - \pi < \varphi$, then $z^D$ is the unique SSE.*

70

**Figure 3.6:** Global interactions (on the left) versus local interactions (on the right).

2. *If $1 - \pi > \varphi$, then $z^C$ is the unique SSE.*

*Remark.* The white spare areas (excluding broundaries) in the right-graph of Figure 3.6 are the cases with $\lceil 4\varphi \rceil = \lceil 4(1 - \pi) \rceil$ and $\lceil 4\pi \rceil = \lceil 4(1 - \varphi) \rceil$. If the combination of $\pi$ and $\varphi$ falls into these areas, then either *i*) both $s^C$ and $s^D$ are SSEs, or *ii*) neither of them is a SSE.

In previous sections, I investigate the dynamics in which, at each period, each individual has a positive chance of drawing any another individual's strategy into her sample to form estimates $\hat{q}_i^t$ and $\hat{p}_i^t$. To compare with local interactions, I call the dynamics examined in previous sections **adaptive dynamics in global interactions.** Figure 3.6 illustrates the SSE under global interactions with the graph on the left, alongside the case of local interactions on the right. The comparison reveals that whether the dynamic with norm-based resentment selects the cooperation equilibrium $z^C$ or the defection equilibrium $z^D$ is *independent* of whether interactions are local or global.

In what follows, I argue that the dynamics under norm-based resentment helps explain recent cross-cultural studies (Henrich et al., 2010) and network experiments (Cassar, 2007; Gracia-Lázaro et al., 2012; Grujic et al., 2010; Kirchkamp and Nagel, 2007; Rand et al., 2011; Suri and Watts, 2011; Traulsen et al., 2010) that

are difficult to explain by previous evolutionary models of cooperation (Abramson and Kuperman, 2001; Eshel et al., 1998; Nowak and May, 1992; Nowak et al., 2010; Ohtsuki, 2006; Santos and Pacheco, 2005; Szabó and Fáth, 2007). Henrich et al. (2010) show that egalitarian behavior, corresponding to choosing cooperation in our model, is positively correlated with the degree of market integration of a society (measured as the percentage of food obtained from market transactions). Moreover, punishment of defectors covaries positively with community size across societies. Henrich et al. (2010)'s study includes 15 diverse societies around the world, covering small-scale societies as well as large, modern ones. When comparing small-scale societies with large, modern ones, two things change. First, in small-scale societies, interactions are restricted to relatives and small neighborhoods, corresponding to the local interactions structure I examine. In contrast, anonymous, long distant interactions feature many market transactions in large societies, corresponding to the global interactions structure I examine. Second, in small-scale societies, people interact mostly with their relatives: neighborhoods are fixed there. In contrast, in large modern societies people *choose* where to live and who to be their friends: neighborhoods are formed endogenously.

Previous evolutionary theories of cooperation often assume the behavioral rule of *imitating the best*, i.e., people imitate the behavior of those earning higher material payoffs. However, imitating the best implies that cooperation should thrive and only thrive when interactions are local and fixed (Abramson and Kuperman, 2001; Eshel et al., 1998; Nowak and May, 1992; Nowak et al., 2010; Ohtsuki, 2006; Santos and Pacheco, 2005; Szabó and Fáth, 2007). The reason is that local interactions allow cooperative behavioral types to form clusters and earn higher material payoffs. Their behavior are then imitated by others, leading to high level of cooperation. However, in global interactions, defectors can always exploit cooperators. Also, for one-shot anonymous interactions, it is difficult to identify who are the punishers and who are not. Hence, it is difficult, if not impossible, to justify punishment of defectors under global interactions in term of material

payoffs.

In contrast, the dynamics under norm-based resentment provide a better explanation to Henrich et al. (2010)'s findings. First, our theorems show that, if individuals are motivated by norm-based resentment, then whether interactions are local or global does not affect the selection of cooperation or defection. Recent controlled experiments confirm this prediction (Cassar, 2007; Gracia-Lázaro et al., 2012; Grujic et al., 2010; Kirchkamp and Nagel, 2007; Rand et al., 2011; Suri and Watts, 2011; Traulsen et al., 2010). For example, in Grujic et al. (2010)'s experiment, participants are located on the same lattice structure as the one I examine. In one treatment, participants play a prisoner's dilemma game with neighbors in the lattice. In another, neighbors are randomly reallocated for each round. The experiment shows that the cooperation level is not distinguishable between the two treatments, suggesting that the lattice structure does not influence cooperation. Rand et al. (2011) consider more arbitrary networks and examine three treatment conditions: random-link condition (the network is randomly regenerated for each round), fixed-link condition (the network is fixed during the experiment), and endogenous-link condition (participants can rewire links). They find that cooperation decays over time in both the random-link condition and in the fixed-link condition, and average cooperation level is not statistically different between the two conditions. However, if participants can frequently rewire links, then cooperation persists at a high level. The pattern is clear: whether interactions are local or global does not matter, but whether neighborhoods are fixed or formed endogenously matters—the ability to vote with feet matters.

### 3.5.2 Mobility

We now extend the model to explore the ability to migrate to different societies. Let the total population $N$ consist of $2n$ individuals. Index time periods by $t = 0, 1, 2, \ldots$. For simplicity, consider two societies, 1 and 2, that individuals can

choose to live in.[16] Let $p_{i\ell}^t$ denote individual $i$'s expectation of the proportion of punishers in society $\ell = 1, 2$ in period $t$. Likewise, let $q_{i\ell}^t$ denote individual $i$'s expectation of the proportion of cooperators in society $\ell = 1, 2$ in period $t$. Let $x_i^t \in \{0, 1\}$ denote whether $i$ cooperates and $y_i^t \in \{0, 1\}$ denote whether $i$ punishes defectors in period $t$.

At the beginning of each period, some individuals are randomly selected to have the opportunity to choose which society to live for the current period. The timeline of events within each period $t$ is as follows:

1. $2m$ individuals are randomly selected to have the opportunity to choose which society to live in for the current period, with $m < n$. The rest $2(n-m)$ individuals are randomly and equally divided into the two societies.[17] Hence, at each period, each society has at least $n - m$ individuals (assume $n - m$ is an even integer). The ratio $\frac{m}{n}$ reflects the degree of freedom of mobility. Let $N_\ell^t \subset N$ be the set of individuals in society $\ell$ in period $t$.

2. Based on each individual $i$'s expectations on the proportion of cooperators and punishers in her society, $p_{i\ell}^t$ and $q_{i\ell}^t$, $i$ chooses actions $x_i$ and $y_i$ according to behavioral rule (3.5) supported by norm-based resentment.

3. For each $i \in N$ and $\ell = 1, 2$, each of the following expectation-updating events occurs with positive probability: 1) $q_{i\ell}^{t+1} = q_{i\ell}^t$ and $p_{i\ell}^{t+1} = p_{i\ell}^t$, 2) $q_{i\ell}^{t+1} = q_{i\ell}^t$ and $p_{i\ell}^{t+1} = \frac{1}{N_\ell} \sum_{j \in N_\ell} y_j^t$, 3) $q_{i\ell}^{t+1} = \frac{1}{N_\ell} \sum_{j \in N_\ell} x_j^t$ and $p_{i\ell}^{t+1} = p_{i\ell}^t$, and 4) $q_{i\ell}^{t+1} = \frac{1}{N_\ell} \sum_{j \in N_\ell} x_j^t$ and $p_{i\ell}^{t+1} = \frac{1}{N_\ell} \sum_{j \in N_\ell} y_j^t$.

How each individual chooses where to live is determined as follows. At the beginning of each period, each $i$ calculates her expected material payoffs for each society based on the expectations $p_{i\ell}^t$ and $q_{i\ell}^t$ about the society and the actions that she anticipates to take. Upon having the opportunity to move, each $i$ chooses

---

[16]The argument extends to the case of more societies without extra efforts.

[17]That we do not allow all individuals to choose the societies to migrate is because, in that case, a society may end up with zero population and thus a norm of that society is not well defined. An alternative assumption to prevent this is to have a upper bound $\bar{n}$ on the size of each society, with $\bar{n} < n$.

to live in the society with higher expected material payoffs. If $i$ has the same expected material payoffs for both societies, she randomly chooses a society, with each society chosen with positive probability. All random events are independent across individuals and time.

We call the dynamic specified above **unperturbed adaptive dynamic** $P^0$ **with mobility of degree** $m$. We define the associated **perturbed adaptive dynamic** $P^\varepsilon$ **with mobility of degree** $m$ as follows: at each period, with probability $\varepsilon > 0$, an individual takes cooperation and punishment actions that are not consistent with their expectations.

The state variable of the dynamics with mobility is

$$ z = \Big( (x_i), (y_i), (p_{i\ell}), (q_{i\ell}), (N_\ell) \Big). $$

Let $Z$ denote the collection of all states with $|N_\ell| \geq n - m$, $\ell = 1, 2$. Let $Z^d \subset Z$ denote the subset of states such that, for each $i \in N$, we have $x_i = y_i = 0$ and $p_{i\ell} = q_{i\ell} = 0$. Let $Z^c \subset Z$ denote the subset of states such that, for each $i \in N$, we have $x_i = y_i = 1$ and $p_{i\ell} = q_{i\ell} = 1$.

Moreover, let $\tilde{Z}^{c1} \subset Z$ denote the subset of states such that: society 1 is larger, consisting of $n + m$ individuals, whereas society 2 is smaller, consisting of $n - m$ individuals; in society 1, each $i$ has $x_i = y_i = 1$ and $p_{i\ell} = q_{i\ell} = 1$; in society 2, each $i$ has $x_i = y_i = 0$ and $p_{i\ell} = q_{i\ell} = 0$. Let $\tilde{Z}^{c2} \subset Z$ denote the subset of states with the same pattern but society 2 is larger: society 1 has $n - m$ individuals, whereas society 2 has $n + m$ individuals; in society 1, each $i$ has $x_i = y_i = 0$ and $p_{i\ell} = q_{i\ell} = 0$; in society 2, each $i$ has $x_i = y_i = 1$ and $p_{i\ell} = q_{i\ell} = 1$. In other words, for a state in $\tilde{Z}^{c1}$ or in $\tilde{Z}^{c2}$, one society is larger than the other. The cooperation norm prevails in the larger society, while the defection norm prevails in the smaller one. The greater the degree of mobility $m$, the greater the size of the society where the cooperation norm prevails.

I assume that the outcome of cooperation generates higher aggregate efficiency than the outcome of defection and no punishment:

**Assumption 3.3.** $a + \bar{b} > \bar{a} + b$.

The above assumption is used and only used in proving the following theorem.

**Theorem 3.4.** *Consider the adaptive dynamics with mobility of degree m. Consider the generic cases with $\lceil \varphi(n-m) \rceil \neq \lceil (1-\pi)(n-m) \rceil$ and $\lceil \pi(n-m) \rceil \neq \lceil (1-\varphi)(n-m) \rceil$.*

*1. If $1 - \pi > \varphi$, then SSEs are states in $Z^c$.*

*2. If $1 - \pi < \varphi < 1 - \pi + \frac{m}{n}(1 - |\pi - \varphi|)$, then SSEs are states in $\tilde{Z}^{c1} \cup \tilde{Z}^{c2}$.*

*3. If $\varphi > 1 - \pi + \frac{m}{n}(1 - |\pi - \varphi|)$, then SSEs are states in $Z^d$.*

The above theorem shows that mobility has two effects: the migration effect in the intermediate run, leading to larger cooperation societies in states in $\tilde{Z}^{c1} \cup \tilde{Z}^{c2}$; and the norm selection effect in the long run, increasing the chance of states in $\tilde{Z}^{c1} \cup \tilde{Z}^{c2}$ being stochastically stable. First, consider the migration effect in the intermediate run. Suppose that we have the cooperation norm in one society and the defection norm in the other. Since the cooperation norm is more efficient, everyone who has a chance to move migrates to the society where the cooperation norm prevails. Hence, the society with the cooperation norm *becomes* larger. Moreover, the greater the degree of mobility $m$, the greater the size of the society with the cooperation norm. What is critical in the above argument is the **fitting-in effect** generated by norm-based resentment, namely, when migrating to a society, individuals adjust their behavior to make sure that it is compatible with the prevailing norms in the society.

Now consider the long-run effect. Which equilibrium—cooperation or defection—is more likely to emerge and persist in the long run depends on the difficulty of transitions from states in $Z^d$ to states in $\tilde{Z}^{c1} \cup \tilde{Z}^{c2}$ and $Z^c$, relative to transitions in the opposite direction. The above theorem shows that the greater the degree of mobility $m$, the more difficult are transitions from $\tilde{Z}^{c1} \cup \tilde{Z}^{c2}$ to $Z^d$, and the easier are transitions from $Z^d$ to $\tilde{Z}^{c1} \cup \tilde{Z}^{c2}$. Hence, mobility favors the emergence and the persistence of the cooperation norm. The intuition is as follows.

Consider transitions from $Z^d$ to $\tilde{Z}^{c1} \cup \tilde{Z}^{c2}$, corresponding to the emergence of the cooperation norm. For the norm of cooperation to emerge in a society, we need enough individuals to cooperate or punish defectors. When the norm of defection prevails, these events only occur in the form of mistakes. Hence, the easiest way to escape from $Z^d$ is to wait until a society becomes smaller, so that we need less mistakes for the norm of cooperation to emerge in the society. Once this is done, the migration effect in the intermediate run automatically takes place, leading to a state in $\tilde{Z}^{c1} \cup \tilde{Z}^{c2}$. The size of the smallest possible society is $n - m$. Hence, the greater the degree of mobility $m$, the easier transitions from $Z^d$ to $\tilde{Z}^{c1} \cup \tilde{Z}^{c2}$. Next, consider transitions from $\tilde{Z}^{c1} \cup \tilde{Z}^{c2}$ to $Z^d$, whose difficulty determines the persistence of the cooperation norm. For a transition from $\tilde{Z}^{c1} \cup \tilde{Z}^{c2}$ to $Z^d$ to occur, we need enough proportion of individuals in the society where the norm of cooperation initially prevails to make mistakes. The greater the degree of mobility $m$, the larger the society where the cooperation norm prevails, and thus the more mistakes are required for the transition to take place.

Finally, let me compare the adaptive dynamics with norm-based resentment with a dynamic of imitating the best. Observe that the fitting-in effect is missing in a dynamic of imitating the best. As a result, a dynamic of imitating the best cannot explain why mobility matters for the emergence and persistence of cooperation. To see this, suppose initially that the cooperation norm prevails in one society and the defection norm prevails in the other. Since the cooperation norm is more efficient, individuals migrate to the society where everyone cooperates. However, individuals soon discover that choosing not to punish defectors leads to a material payoff at least as good as punishing defectors, and it would be strictly better if there indeed are defectors. Hence, as shown by Home et al. (2006)'s simulation, the population drifts to a state where no individual punishes defectors. As a result, defectors obtain higher material payoffs than cooperators. Defection is then imitated and spreads. Eventually, everyone defects and no one punishes defectors in the whole population. Hence, defection prevails under a dynamic of

77

imitating the best regardless of whether mobility is possible.

## 3.6 Law enforcement and cooperative norms

This section extends the basic model in another dimension in order to analyze the interplay between law enforcement and cooperative norms. Using public goods game experiments, Herrmann et al. (2008) show that individuals behave more cooperative in societies with better law enforcement. Using survey data, Tabellini (2008) also find positive relationships between quality of formal institutions and individuals' cooperative attitudes.

First, I model law enforcement as an exogenous central monitor who conducts punishment towards defectors independently of the decentralized punishment supported by norm-based resentment. Such centralized punishment reduces individuals' temptation to defect and thus help cooperative norms emerge and persist in the long run. In the second step, I endogenize the quality of law enforcement by exploiting the fact that shirking or corruption of the central monitor is an instance of defection behaviors against the common interests of the public.

### 3.6.1 Exogenous law enforcement

Consider the following law enforcement of cooperation implemented by a central monitor. If an individual defects, the central monitor detects her defection with probability $\delta > 0$. If detected, the defector pays a fine $v > 0$. Law enforcement is independent of the decentralized punishment conducted by matched individuals. As Figure 3.7 shows, if an individual defects and the matched second-mover does not punish the defector, then the defector's expected payoffs are $\bar{a} - \delta v$. If the matched second-mover punishes the defector, the defector's expected payoffs are $\underline{a} - \delta v$.

**Figure 3.7:** The stage game with law enforcement.

Applying our previous analysis of the stage game in Section 3.2.1, we obtain the threshold proportion of cooperators to activate decentralized punishment of defectors:

$$\pi = \left[ 1 + \frac{1}{\lambda(\bar{a} - \underline{a})} \right] \frac{b - \underline{b}}{\bar{\bar{b}} - \underline{b}},$$

which is the same threshold (3.3) without law enforcement. Hence, law enforcement does not affect the threshold proportion of cooperators to activate punishment of defectors. However, the threshold proportion of punishers to deter defectors becomes

$$\frac{\bar{a} - a - \delta v}{\bar{a} - \underline{a}} = \varphi - \frac{\delta v}{\bar{a} - \underline{a}},$$

where $\varphi = \frac{\bar{a} - a}{\bar{a} - \underline{a}}$ is the threshold (3.4) without law enforcement. Therefore, law enforcement lowers the threshold proportion of punishers to deter defectors. We thus have the following result, which is a corollary of Theorem 3.2.

**Theorem 3.5.** *Consider the adaptive dynamics with sample size m in global interactions with law enforcement. Let $\pi$ and $\varphi$ be the thresholds defined by (3.3) and (3.4), corresponding to the case without law enforcement. Consider the generic cases with $\lceil (\varphi - \frac{\delta v}{\bar{a} - \underline{a}})m \rceil \neq \lceil (1 - \pi)m \rceil$ and $\lceil \pi m \rceil \neq \lceil (1 - \varphi + \frac{\delta v}{\bar{a} - \underline{a}})m \rceil$.*

1. *If $(1 - \pi) + \frac{\delta v}{\bar{a} - \underline{a}} < \varphi$, then $z^D$ is the unique SSE.*

2. *If $(1 - \pi) + \frac{\delta v}{\bar{a} - \underline{a}} > \varphi$, then $z^C$ is the unique SSE.*

In the inequality conditions of the above theorem, $\varphi$ is the threshold proportion of punishers to deter defectors when there is no law enforcement; it measures

the temptation to defect. $1 - \pi$ is the maximum proportion of defectors in the population that individuals still think the norm of cooperation exists and are willing to punish defectors; it measures the strength of norm-based resentment to sustain cooperation. $\frac{\delta v}{\bar{a} - \underline{a}}$ measures the strength of law enforcement of cooperation. The theorem says that, if the combining force of norm-based resentment and law enforcement is strong enough, then the cooperation equilibrium is the unique SSE. Hence, law enforcement improves the chance of the norm of cooperation to emerge and persist in the long run.

### 3.6.2 Endogenous law enforcement

The central monitoring institute that enforces laws itself consists of humans. They might shirk or be involved in corruption. Although difficult to establish causality in empirical studies (Gächter and Schulz, 2016; Herrmann et al., 2008; Tabellini, 2008), it is hard to believe that the influences between centralized law enforcement and decentralized social norms only go in one direction. I examine how they might co-evolve in this subsection.

Suppose that the central monitor can choose whether to shirk. Shirking reduce $\delta$, the probability of the central monitor detecting defectors. For an average individual $i$ in the population, $p_i$ is $i$'s expectation of the proportion of punishers in the population. $i$ also has an expectation on $\delta$. I assume that $i$'s expectation on the probability of the central monitor detecting defectors, $\delta(p_i)$, is increasing in $p_i$, with $\delta(0) = \underline{\delta}$ and $\delta(1) = \bar{\delta} > \underline{\delta}$. The idea is that, when individual $i$ assesses the probability of detecting defectors, $i$ reasons as follows: when the monitor chooses whether or not to shirk, others consider the monitor as facing the same choice between cooperation and defection as in the stage game: the monitor could have chosen an action (not shirking) that contributes to the public good and benefits others. Hence, $i$ expects that, the more individuals punish defectors in the population, the more individuals punish the central monitor if the monitor shirks. Therefore, if $i$ expects a higher proportion of individuals who punish defectors,

then $i$ would also expect a higher probability of the central monitor enforcing the law and detecting defectors.

Let $\tilde{\varphi} \in [0,1]$ denote the expected proportion of punishers that makes $i$ indifferent between cooperation and defection. We have

$$a = \tilde{\varphi}\underline{a} + (1 - \tilde{\varphi})\bar{a} - \delta(\tilde{\varphi})v. \tag{3.9}$$

If $p_i > \tilde{\varphi}$, $i$ cooperates. If $p_i < \tilde{\varphi}$, $i$ defects. In particular, $\tilde{\varphi} = \tilde{\varphi}(v)$ is decreasing in $v$, the fine that a detected defector needs to pay. Think of $v$ as punishment terms written in laws, and $\delta$ as the quality of law enforcement. The observation is that, when the quality of law enforcement is endogenous, laws might or might not be executed. However, laws still matter. The reason is that, controlling for decentralized punishment towards defectors $(p_i)$, a higher $v$ does increase the threat to defectors. Thus, it effectively lowers $\tilde{\varphi}(v)$, the barrier of the population tipping to the cooperation equilibrium. Applying Theorem 3.1, we obtain the following result.

**Theorem 3.6.** *Consider the adaptive dynamics with sample size $m$ in global interactions with law enforcement. Let $\pi$ be defined by (3.3) and $\tilde{\varphi}(v)$ by (3.9). Consider the generic cases with $\lceil \tilde{\varphi}(v)m \rceil \neq \lceil (1-\pi)m \rceil$ and $\lceil \pi m \rceil \neq \lceil (1-\tilde{\varphi}(v)m \rceil$.*

1. *If $1 - \pi < \tilde{\varphi}(v)$, then $z^D$ is the unique SSE in which the probability of detecting defectors is $\underline{\delta}$.*

2. *If $1 - \pi > \tilde{\varphi}(v)$, then $z^C$ is the unique SSE in which the probability of detecting defectors is $\bar{\delta}$.*

The above theorem says that the terms in laws, $v$, are fundamentals that affect the selection of social norms in the long run. The social norm then feeds back to influence $\delta$, the quality of law enforcement.

## 3.7 Conclusion

In this chapter, I have examined the conditions under which cooperative norms emerge and persist in the long run given norm-based resentment. I conclude by providing a remark on the evolutionary foundation of norm-based resentment. As Henrich and Boyd (1998) argue, due to the ability to generate multiple equilibria and for group beneficial norms to emerge and spread, conformist bias is flexible and adaptive. As a result, the conformist bias would maximize fitness in a changing environment in the gene-culture co-evolution history of humans. However, in terms of the ability to generate multiple equilibria, norm-based resentment is at least as good as the conformist bias. Moreover, I show that norm-based resentment leads to a population dynamic that is *more* adaptive to environmental changes than the conformist bias. Altogether, norm-based resentment should perform strictly better than the conformist bias in the evolution of humans. Exploring this idea formally would be a task for future research.

## 3.8 Appendix: proofs

### Proof of Proposition 3.1

A strict equilibrium must involve only pure strategies. Hence, for any candidate strategy profile to be a strict equilibrium, it includes at most four types of individuals. Define $C, C^-, D, D^- \subset N$ as follows:

*i)* $C = \{i \in N | x_i = 1, y_i = 1\}$,

*ii)* $C^- = \{i \in N | x_i = 1, y_i = 0\}$,

*iii)* $D = \{i \in N | x_i = 0, y_i = 0\}$, and

*iv)* $D^- = \{i \in N | x_i = 0, y_i = 1\}$.

We use $|A|$ to denote the number of elements in the finite set $A$. The proof involves several steps.

*Step 1: $0 < \varphi < 1$ and $\pi > 0$.*

To see $0 < \varphi < 1$, recall $\underline{a} < a < \bar{a}$. Hence, $0 < \bar{a} - a < \bar{a} - \underline{a}$. Therefore, $0 < \frac{\bar{a}-a}{\bar{a}-\underline{a}} < 1$. Since $\varphi = \frac{\bar{a}-a}{\bar{a}-\underline{a}}$, the desired result follows.

To see that $\pi > 0$, observe

$$\pi = \left[1 + \frac{1}{\lambda(\bar{a} - \underline{a})}\right]\frac{b - \underline{b}}{\bar{b} - \underline{b}} > \frac{b - \underline{b}}{\bar{b} - \underline{b}} > 0,$$

since $\underline{b} < b < \bar{b}$.

*Step 2: There is a strict equilibrium in which $|D| = n$.*

In the strategy profile with $|D| = n$, we have $p_i = 0$ and $q_i = 0$ for each $i \in N$. We have shown that $\varphi > 0$ and $\pi > 0$. Hence, we have $p_i < \varphi$ and $q_i < \pi$ for each $i \in N$, and thus choosing $x_i = 0$ and $y_i = 0$ is sequentially rational, as desired.

*Step 3: If $\pi < 1$, then there is a strict equilibrium in which $|C| = n$.*

In the strategy profile with $|C| = n$, we have $p_i = 1$ and $q_i = 1$ for each $i \in N$. Given the assumption $\pi < 1$, we have $q_i > \pi$ for each $i \in N$. We also have $p_i = 1 > \varphi$ for each $i \in N$. Hence, choosing $x_i = 1$ and $y_i = 1$ is sequentially rational, as desired.

*Step 4: If there is a strict equilibrium in which $|D| \geq 1$, then we have $|D| = n$ in that equilibrium.*

Consider a strict equilibrium in which $|D| \geq 1$. Then, for $i \in D$, we have $p_i < \varphi$ and $q_i < \pi$. Moreover, $p_i$ and $q_i$ are consistent with the strategies of others:

$$p_i = \frac{|C| + |D^-|}{n-1}, \qquad q_i = \frac{|C| + |C^-|}{n-1}.$$

Now, suppose $|C| \geq 1$. Then for $j \in C$ we have

$$q_j = \frac{|C| + |C^-| - 1}{n-1}.$$

But then $q_j < q_i < \pi$. Hence, $j$ would rather deviate to $x_j = 0$. Therefore, the supposition $|C| \geq 1$ does not hold; we must have $|C| = 0$. Likewise, if $|D^-| \geq 1$, then for $j \in D^-$ we have $q_j = \frac{|C| + |C^-|}{n-1} > \pi$, which contradicts $q_i < \pi$. Hence, we must also have $|D^-| = 0$. Finally, suppose $|C^-| \geq 1$. Then for $j \in |C^-|$ we have $p_j = \frac{|C| + |D^-|}{n-1} > \varphi$, which contradicts $p_i < \varphi$. Thus $|C^-| = 0$, as desired.

*Step 5: If there is a strict equilibrium in which $|C| \geq 1$, then we have $|C| = n$ in that equilibrium.*

This can be analogously shown as Step 4.

*Step 6: If $\varphi + \pi \geq \frac{n}{n-1}$ or $\varphi + \pi \leq \frac{n-2}{n-1}$, then there is no strict equilibrium in which $|C^-| \geq 1$ or $|D^-| \geq 1$.*

Suppose $|C^-| \geq 1$ or $|D^-| \geq 1$. By Steps 4 and 5, this implies $|C| = 0$ and $|D| = 0$. Now pick $i \in C^-$ and $j \in D^-$. For their expectations to be consistent, it requires

$$p_i = \frac{|D^-|}{n-1}, \quad q_i = \frac{|C^-| - 1}{n-1}, \quad p_j = \frac{|D^-| - 1}{n-1}, \quad q_j = \frac{|C^-|}{n-1}.$$

For their strategies to be sequential rational, it requires $p_j < \varphi < p_i$ and $q_i < \pi <$

$q_j$. Moreover, given $|C| = 0$ and $|D| = 0$, we have $|C^-| + |D^-| = n$. It follows that

$$\frac{n-2}{n-1} < \varphi + \pi < \frac{n}{n-1},$$

which is ruled out by the assumption that we either have $\varphi + \pi \geq \frac{n}{n-1}$ or $\varphi + \pi \leq \frac{n-2}{n-1}$.

## Proof of Theorem 3.1

We shall show that there is $\delta \in (0, 1]$ such that, with at least probability $\delta$, the dynamic $P^0$ transits from *any* population state $z^t \in Z$ to a population state in $\{(s^C, 1, 1), (s^D, 0, 0)\}$ within *three* periods. It follows that the probability of the dynamic *not* converging to $\{(s^C, 1, 1), (s^D, 0, 0)\}$ within three periods is at most $1 - \delta < 1$. Hence, the probability of *not* converging to $\{(s^C, 1, 1), (s^D, 0, 0)\}$ within $3k$ periods is at most $(1 - \delta)^k$, which goes to zero as $k$ goes to infinity.

Now I show that there exists such $\delta > 0$ by explicitly showing a transition path starting from an arbitrary state $z^t \in Z$ to either $(s^C, 1, 1)$ or $(s^D, 0, 0)$ within three periods, and that this path occurs with positive probability. Let the population state in period $t$ be $z^t = (s^t, p^t, q^t) \in Z$. At $t+1$, let all individuals update their expectations on the proportions of cooperators and punishers—i.e., $q_i$ and $p_i$, and they draw exactly the same samples $(s_i^t)_{i \in M}$ from $s^t$. Then at $t+1$, each $i \in N$ has the same $q_i^{t+1}$ and $p_i^{t+1}$. This results in the same sequentially rational strategy $s_i^{t+1}$ for each $i \in N$. Since $s_i^{t+1} = (x_i^{t+1}, y_i^{t+1})$ is the same for each $i \in N$, there are only four cases: 1) $x_i^{t+1} = y_i^{t+1} = 1$ for each $i \in N$, corresponding to $s^C$, 2) $x_i^{t+1} = y_i^{t+1} = 0$ for each $i \in N$, corresponding to $s^D$, 3) $x_i^{t+1} = 1$ and $y_i^{t+1} = 0$ for each $i \in N$, and finally 4) $x_i^{t+1} = 0$ and $y_i^{t+1} = 1$ for each $i \in N$. If it falls into one of the first two cases, let all individuals update their expectations at $t+2$. In the case of $s_i^{t+1} = s^C$, each $i \in N$ will then have $p_i^{t+2} = q_i^{t+2} = 1$, and thus we obtain $z^{t+2} = (s^C, 1, 1)$. In the case of $s_i^{t+1} = s^D$, each $i \in N$ will have $p_i^{t+2} = q_i^{t+2} = 0$. Thus we obtain $z^{t+2} = (s^D, 0, 0)$, as desired.

Next consider the case $x_i^{t+1} = 1$ and $y_i^{t+1} = 0$ for each $i \in N$. At $t + 2$, let all individuals update but *only* update their expectation of the proportion of punishers, so that we obtain $p_i^{t+2} = 0$ and $q_i^{t+2} = q_i^{t+1}$ for each $i \in N$; this occurs with positive probability. Conditional on this event, each $i \in N$ at $t + 2$ defects, i.e., $x_i^{t+2} = 0$. Moreover, since $y_i^{t+1} = 0$ is a best-response to $q_i^{t+1}$, $y_i^{t+2} = 0$ must also be a best-response to $q_i^{t+2}$ for each $i \in N$. Thus, it occurs with positive probability that we have $x_i^{t+2} = 0$ and $y_i^{t+2} = 0$ for each $i \in N$. At $t + 3$, let all individuals update their expectations on the proportion of cooperators, so that we obtain $q_i^{t+3} = 0$. As a result, the dynamic reaches $z^{t+3} = (s^D, 0, 0)$.

Finally, consider the case $x_i^{t+1} = 0$ and $y_i^{t+1} = 1$ for each $i \in N$. Similar to the last case, at $t + 2$, let all individuals update but *only* update their expectation of the proportion of punishers, so that we obtain $p_i^{t+2} = 1$ and $q_i^{t+2} = q_i^{t+1}$ for each $i \in N$. Conditional on this event, each $i \in N$ at $t + 2$ would cooperate, i.e., we have $x_i^{t+2} = 1$ for each $i \in N$. Meanwhile, since $y_i^{t+1} = 1$ is a best-response to $q_i^{t+1}$, $y_i^{t+2} = 1$ must also be a best-response to $q_i^{t+2}$ for each $i \in N$. Thus, it occurs with positive probability that we have $x_i^{t+2} = 1$ and $y_i^{t+2} = 1$ for each $i \in N$. At $t + 3$, let all individuals update their expectations on the proportion of cooperators. This results in $q_i^{t+3} = 1$ and thus we obtain $z^{t+3} = (s^C, 1, 1)$ with positive probability.

Observe that each step involved in the transition described above occurs with positive probability. Hence there is $\delta > 0$ such that, with at least probability $\delta$, dynamic transits from any population state to a population in $\{(s^C, 1, 1), (s^D, 0, 0)\}$ within three periods, as I claim.

## Proof of Proposition 3.2

*Step 1: Let $\{z^t\}_{t=0}^{\infty} \subset Z$ be a sequence of random states by such that $z^0 \in Z$ satisfies: a) $\sum x_i^0 < \pi m$, $\sum y_i^0 < \varphi m$, $\sum \mathbb{1}\{p_i^0 \geq \varphi\} < \pi m$ and $\sum \mathbb{1}\{q_i^0 \geq \pi\} < \varphi m$, and b) for each $t \geq 1$ we have $Prob\{z^t | z^{t-1}\} = P^0(z^t, z^{t-1})$. Then $Prob\{\lim_{t \to \infty} z^t = z^D\} = 1$.*

First, observe that, for each $z^0 \in Z$ that satisfies the specified conditions, there is at least probability $\delta > 0$ of having $z^1 = z^D$, namely, the transition from $z^0$ to $z^D$ is completed within one period with positive probability. This is achieved by requiring all individuals update their expectations—both $p_i$ and $q_i$—at $t = 1$, so that we have $p_i^1 < \varphi$ and $q_i^1 < \pi$ for each $i \in N$. Then, even if an individual includes all cooperators and punishers in her sample, she would take actions $x_i^1 = 0$ and $y_i^1 = 0$ at $t = 1$. The event that all individuals update expectations occurs with positive probability. Hence there exists $\delta > 0$ such that $z^1 = z^D$ occurs with at least probability $\delta$.

Next, I show that, if $z^0$ satisfies the specified conditions, then $z^1$ must also satisfy the specified conditions. By induction, $z^t$ satisfies the specified conditions for each $t \geq 1$. First, let $J = \{i \in N | q_i^0 \geq \pi\}$. Given $\sum x_i^0 < \pi m$, if an individual $j \in N$ updates her expectation of the proportion of cooperators at $t = 1$, then she has $q_j^1 < \pi$. Hence, if and only if $j \in J$ and $j$ does not update her expectation of the proportion of cooperators, then we have $q_j^1 \geq \pi$. Let $\lceil \varphi m \rceil$ denote the smallest integer equal to or greater than $\varphi m$. Since $|J| < \varphi m$, the number of individuals with $q_j^1 \geq \pi$ at $t = 1$ is at most $\lceil \varphi m \rceil - 1$, i..e, $\sum \mathbb{1}\{q_i^1 \geq \pi\} \leq \lceil \varphi m \rceil - 1$. Hence, we at most have $\lceil \varphi m \rceil - 1$ individuals who take $y_i^1 = 1$ at $t = 1$. Thus, $\sum y_i^1 < \varphi m$.

Similarly, given $\sum y_i^0 < \varphi m$, if an individual $j \in N$ updates her expectation of the proportion of punishers at $t = 1$, then she has $p_j^1 < \varphi$. Hence, we have $p_j^1 \geq \varphi$ if and only if $j$ is such that $p_j^0 \geq \varphi$ and $j$ does not update her expectation of the proportion of cooperators. Therefore, given $\sum \mathbb{1}\{p_i^0 \geq \varphi\} < \pi m$, we must also have $\sum \mathbb{1}\{p_i^1 \geq \varphi\} < \pi m$ at $t = 1$. Hence, $\sum x_i^1 < \pi m$, as desired.

To conclude, starting any state $z^0$ satisfying the specified conditions, we have $z^t$ satisfying the specified conditions for all future periods. Meanwhile, there is at least probability $\delta > 0$ of transiting from a state $z^t$ satisfying the specified conditions to $z^{t+1} = z^D$ within one period. Hence, the probability of not transiting from $z^0$ to $z^D$ within $k$ periods is at most $(1 - \delta)^k$, which goes to zero as $k$ goes to infinity.

*Step 2: There is positive probability that the dynamic $P^0$ transits to $z^C$ within four periods if one of the following conditions holds for the initial state $z^0$:* a) $\sum x_i^0 \geq \pi m$, b) $\sum y_i^0 \geq \varphi m$, c) $\sum \mathbb{1}\{p_i^0 \geq \varphi\} \geq \pi m$ and d) $\sum \mathbb{1}\{q_i^0 \geq \pi\} \geq \varphi m$.

Case 1: Let $\sum x_i^0 \geq \pi m$. Then at $t = 1$ let all individuals draw the exact sample $(s_i^0)_{i \in M}$ with $M = \{i \in N | x_i^0 = 1\}$ and update their expectations. As a result, we have $q_i^1 = \frac{\sum x_i^0}{m} \geq \pi$ for each $i \in N$. Thus, each $i \in N$ takes $y_i^1 = 1$ with positive probability. At the next period, $t = 2$, let all individuals update their expectation of the proportion of punishers but *not* about cooperators. We then have $p_i^2 = 1$ and $q_i^2 = q_i^1 \geq \pi$ for each $i \in N$. Each individual then takes $x_i^2 = 1$ and $y_i^2 = 1$ with positive probability at $t = 2$. At $t = 3$, let all individuals update their expectations on the proportions of cooperators as well as the punishers; we obtain $z^3 = z^C$.

Case 2: Let $\sum y_i^0 \geq \varphi m$. Similarly, at $t = 1$ let all individuals draw the exact sample $(s_i^0)_{i \in M}$ with $M = \{i \in N | y_i^0 = 1\}$ and update their expectations. As a result, we have $p_i^1 = \frac{\sum y_i^0}{m} \geq \varphi$ for each $i \in N$. Thus, each $i \in N$ takes $x_i^1 = 1$ with positive probability. At the next period, $t = 2$, let all individuals update their expectation of the proportion of cooperators but *not* about punishers. We then have $q_i^2 = 1$ and $p_i^2 = p_i^1 \geq \varphi$ for each $i \in N$. Each individual then takes $x_i^2 = 1$ and $y_i^2 = 1$ with positive probability.

Case 3: Let $\sum \mathbb{1}\{p_i^0 \geq \varphi\} \geq \pi m$. At $t = 1$, let all individuals do not update their expectation of the proportion of punishers, i.e., $p_i^1 = p_i^0$ for each $i \in N$. Then at $t = 1$ we also have $\sum \mathbb{1}\{p_i^1 \geq \varphi\} \geq \pi m$. Since the individuals with $p_i^1 \geq \varphi$ take $x_i^1 = 1$ with positive probability, it occurs with positive probability that $\sum x_i^1 \geq \pi m$. But then, the argument for Case 1 applies, leading to $z^4 = z^C$ with positive probability.

Case 4: Let $\sum \mathbb{1}\{q_i^0 \geq \pi\} \geq \varphi m$. Similarly to Case 3, let all individuals do not update their expectation of the proportion of cooperators at $t = 1$, i.e., $q_i^1 = q_i^0$ for each $i \in N$. Then at $t = 1$ we have $\sum \mathbb{1}\{q_i^1 \geq \pi\} \geq \varphi m$. The individuals with $q_i^1 \geq \pi$ take $y_i^1 = 1$ with positive probability. Hence, it occurs with positive

88

probability that $\sum y_i^1 \geq \varphi m$. But then, the argument for Case 2 applies, leading to $z^4 = z^C$ with positive probability.

*Step 3: Let $\{z^t\}_{t=0}^{\infty} \subset Z$ be a sequence of random states by such that $z^0 \in Z$ satisfies: a) $n - \sum x_i < (1-\pi)m$, $n - \sum y_i < (1-\varphi)m$, $\sum \mathbb{1}\{p_i \leq \varphi\} < (1-\pi)m$ and $\sum \mathbb{1}\{q_i \leq \pi\} < (1-\varphi)m$, and b) for each $t \geq 1$ we have $Prob\{z^t|z^{t-1}\} = P^0(z^t, z^{t-1})$. Then $Prob\{\lim_{t\to\infty} z^t = z^C\} = 1$.*

The proof is completely analogous to the proof of Step 1, so omitted.

*Step 4: There is positive probability that the dynamic $P^0$ transits to $z^D$ within two periods if one of the following conditions holds for the initial state $z^0$: a) $n - \sum x_i \geq (1-\pi)m$, b) $n - \sum y_i \geq (1-\varphi)m$, c) $\sum \mathbb{1}\{p_i \leq \varphi\} \geq (1-\pi)m$ and d) $\sum \mathbb{1}\{q_i \leq \pi\} \geq (1-\varphi)m$.*

The proof is completely analogous to the proof of Step 2, so omitted.

## Proof of Lemma 3.1

*Step 1: $R(z^D) = \lceil m \min\{\pi, \varphi\} \rceil$.*

Let the dynamic start from $z^0 = z^D$. By Proposition 3.2, to leave $B(z^D)$ we need to either have a) $\sum x_i \geq \pi m$ or $\sum y_i \geq \varphi m$, or b) $\sum \mathbb{1}\{p_i \geq \varphi\} \geq \pi m$ or $\sum \mathbb{1}\{q_i \geq \pi\} \geq \varphi m$. However, notice that it is not possible to have the latter condition regarding the expectations before having the former condition regarding the actions. To see this, suppose that we have the state in period $t$, for some $t \geq 1$, is such that $\sum \mathbb{1}\{p_i^t \geq \varphi\} \geq \pi m$, and that, for each $t' < t$, we have $\sum y_i^{t'} < \varphi m$. Then, in period $t-1$, we must have $\sum \mathbb{1}\{p_i^{t-1} \geq \varphi\} \geq \pi m$ or $\sum y_i^{t-1} \geq \varphi m$. Since $\sum y_i^{t'} < \varphi m$ for each $t' < t$, we have $\sum \mathbb{1}\{p_i^{t-1} \geq \varphi\} \geq \pi m$. But then, by induction, we have $\sum \mathbb{1}\{p_i^{t'} \geq \varphi\} \geq \pi m$ for each $t' < t$. However, in the starting period, we have $z^0 = z^D$ so that $\sum \mathbb{1}\{p_i^0 \geq \varphi\} = 0$, which is a contradiction. Hence, our initial supposition does not hold. If we have $\sum \mathbb{1}\{p_i^t \geq \varphi\} \geq \pi m$ for some $t \geq 1$, then we must have $\sum y_i^{t'} \geq \varphi m$ for some $t' < t$. Likewise, if we have $\sum \mathbb{1}\{q_i^t \geq \pi\} \geq \varphi m$ for some $t \geq 1$, then we must have $\sum x_i^{t'} \geq \pi m$ for some $t' < t$. Hence, for a dynamic that starts from $z^D$ to leave $B(z^D)$, it must first

require enough mistakes to have $\sum x_i \geq \pi m$ or $\sum y_i \geq \varphi m$. Before this is done, we always have $q_i < \pi$ and $p_i < \varphi$ for each $i \in N$; under these expectations, if some $i$ takes $x_i = 1$ or $y_i = 1$, she must be making a mistake. Hence, the shortest path—the path involving the minimum number of mistakes—to leave $B(z^D)$ is to make sufficient mistakes to have either $\sum x_i \geq \pi m$ or $\sum y_i \geq \varphi m$ all at once, i.e., all mistakes occur in the period. It follows that, if $\pi \leq \varphi$, then the shortest path to leave $B(z^D)$ is to have $\lceil \pi m \rceil$ individuals to make mistakes and cooperate ($x_i = 1$). If $\pi > \varphi$, then the shortest path to leave $B(z^D)$ is to have $\lceil \varphi m \rceil$ individuals to make mistakes and punish ($y_i = 1$). Hence, $R(z^D) = \lceil m \min\{\pi, \varphi\} \rceil$, as desired.

*Step 2:* $R(z^c) = \lceil m \min\{1 - \pi, 1 - \varphi\} \rceil$.

Analogous to the proof of Step 1, so omitted.

## Proof of Theorem 3.2

**Preliminaries.** To prove the statement, I apply a general result about stochastic dynamics of Young (1993). We start by introducing new concepts. A **recurrent class** of a Markov process is a subset of states $L \subset Z$ such that there is a positive probability of transiting between any two states in $L$ within finite periods, but the probability of transiting from a state in $L$ to a state outside of $L$ is zero within any finite periods. Theorem 3.1 implies that there are two and only two absorbing states of $P^0$: the two strict equilibria of the stage game, and no other recurrent class exists (otherwise it cannot be probability one of transiting from any other state to the two strict equilibria within finite periods).

For a Markov process, a **stochastic tree** is a directed tree defined on the collection of all recurrent classes of the process, which specifies how to transit *from* all other recurrent classes *to* the recurrent class at the root of the tree. Formally, let $\Omega$ be the collection of all recurrent classes of $P^0$. For a recurrent class $L$ of $P^0$, a $L$**-tree** $t : \Omega \to \Omega$ is such that $t(L) = L$, and that, for each $L' \in \Omega, L' \neq L$, there is a natural number $k$ such that $t^k(L') = L$ (where $t^k$ is the $k$th functional power of $t$). For two recurrent classes $L$ and $L'$, let $c(L, L')$ denote

the minimum number of mistakes required for $P^\varepsilon$ to transit from $L$ to $L'$. The **cost of a stochastic L-tree** $t$ is defined by

$$c(t) = \sum_{L' \in \Omega, L' \neq L} c(L', t(L')).$$

**Theorem** (Young, 1993)**.** *A state $z \in Z$ is stochastically stable if and only if it is contained in a recurrent class $L \in \Omega$ such that there is a stochastic L-tree with minimum cost among all stochastic trees defined on $\Omega$.*

In our context, there are two and only two recurrent classes of $P^0$: one contains $z^C$, and the other contains $z^D$. Hence a stochastic tree of $P^0$ contains only two vertexes, and there are only two stochastic trees. One stochastic tree is rooted at $z^C$, which consists of *i*) the state $z^C$, *ii*) the state $z^D$, and *iii*) a directed edge weighted with the minimum number of mistakes required to transit from $z^D$ to $z^C$. The other stochastic tree is rooted at $z^D$, which is obtained by reversing the direction of the edge and weighting the edge with the minimum number of mistakes to transit from $z^C$ to $z^D$. Therefore, to identify the stochastic tree with minimum resistance over all stochastic trees is simple in our setting. It is done by comparing the minimum numbers of mistakes required to transit from one strict equilibrium to the other. Observe that the minimum number of mistakes required to transit from a strict equilibrium to the other is simply the minimum number of mistakes required to leave the basin of attraction of the incumbent equilibrium. Hence, we have the following lemma.

**Lemma 3.2.** *If $R(z^C) > R(z^D)$, then $s^C$ is the unique SSE. Conversely, if $R(z^C) < R(z^D)$, then $z^D$ is the unique SSE.*

*Proof.* The statement follows from Young's theorem and Theorem 3.1. $\qquad\square$

To establish theorem 3.2, what remains to show is the following:

**Lemma 3.3.** *Consider the generic cases with $\lceil \varphi m \rceil \neq \lceil (1-\pi)m \rceil$ and $\lceil \pi m \rceil \neq \lceil (1-\varphi)m \rceil$.*

*1. If $\pi + \varphi > 1$, then $\lceil \min\{\pi, \varphi\}m \rceil > \lceil \min\{1 - \pi, 1 - \varphi\}m \rceil$.*

*2. If $\pi + \varphi < 1$, then $\lceil \min\{\pi, \varphi\}m \rceil < \lceil \min\{1 - \pi, 1 - \varphi\}m \rceil$.*

*Proof.* Case 1: Let $\pi \geq \varphi$. Then $\min\{\pi, \varphi\} = \varphi$ and $\min\{1 - \pi, 1 - \varphi\} = 1 - \pi$. In this case, suppose $\pi + \varphi > 1$. Then we have $\varphi m > (1 - \pi)m$ and thus

$$\lceil \min\{\pi, \varphi\}m \rceil = \lceil \varphi m \rceil \geq \lceil (1 - \pi)m \rceil = \min\{1 - \pi, 1 - \varphi\}.$$

Since I consider the generic cases with $\lceil \varphi m \rceil \neq \lceil (1 - \pi)m \rceil$, we obtain $\lceil \min\{\pi, \varphi\}m \rceil > \min\{1 - \pi, 1 - \varphi\}$. Now suppose $\pi + \varphi < 1$. Then $\varphi m < (1 - \pi)m$. Given $\lceil \varphi m \rceil \neq \lceil (1 - \pi)m \rceil$,

$$\lceil \min\{\pi, \varphi\}m \rceil = \lceil \varphi m \rceil < \lceil (1 - \pi)m \rceil = \min\{1 - \pi, 1 - \varphi\}.$$

Case 2: Let $\pi < \varphi$. Then $\min\{\pi, \varphi\} = \pi$ and $\min\{1 - \pi, 1 - \varphi\} = 1 - \varphi$. In this case, suppose $\pi + \varphi > 1$. Then we have $\pi m > (1 - \varphi)m$. Combined with $\lceil \pi m \rceil \neq \lceil (1 - \varphi)m \rceil$, we obtain

$$\lceil \min\{\pi, \varphi\}m \rceil = \lceil \pi m \rceil > \lceil (1 - \varphi)m \rceil = \min\{1 - \pi, 1 - \varphi\}.$$

Conversely, if $\pi + \varphi < 1$, then $\pi m < (1 - \varphi)m$, and thus

$$\lceil \min\{\pi, \varphi\}m \rceil = \lceil \pi m \rceil < \lceil (1 - \varphi)m \rceil = \min\{1 - \pi, 1 - \varphi\}.$$

$\square$

## Proof of Theorem 3.3

**Preliminaries.** To prove the statement, I apply Ellison's theorem (2000, p. 30). We introduce the notions of radius and coradius of a state.

Let $(z^t)_{t=0}^{\infty}$, with $z^t \in Z$ for each $t \geq 0$, be a sample sequence of $P^\varepsilon$ with initial state $z^0$. Let $\gamma = (z^0, z^1, \ldots, z^t)$ denote a path from state $z^0$ to state $z^t$. Let

$\Gamma(X, Y)$ denote the set of all paths from $X \subset Z$ to $Y \subset Z$, i.e.,

$$\Gamma(X, Y) = \{(z^0, z^1, \dots, z^t) | z^0 \in X, z^t \in Y \text{ for some integer } t \geq 0\}.$$

Let $c(\gamma)$ denote the minimum number of mistakes that is required to go through the path $\gamma$. We call $c(\gamma)$ the **cost of path** $\gamma$.

Let $L \subset Z$ be a subset of states. Let $B(L) \subset Z$ be the basin of attraction of $L$ such that if $P^0$ starts from a state within $B(L)$ then it converges almost surely to states within $L$. Let integer $R(L)$ denote the minimum number of mistakes required to leave $B(L)$. Following Ellison (2000), I call $R(L)$ the **radius of** $L$. More precisely, we have

$$R(L) = \min_{\gamma \in \Gamma(L, Z \setminus B(L))} c(\gamma).$$

In words, $R(L)$ is the minimum cost that is required to leave $B(L)$. With abuse of notation, if $L$ is singleton containing only the state $z \in Z$, then I write $R(z)$ instead of $R(\{z\})$.

For a path $\gamma$, let $(L_1, L_2, \dots, L_K)$ denote the sequence of recurrent classes of $P^0$ through which the path passes consecutively. The **modified cost of path** $\gamma$ is defined by

$$c^*(\gamma) = c(\gamma) - \sum_{k=2}^{K-1} R(L_k).$$

Let $c^*(z, B(L))$ denote the modified cost of the path $\gamma \in \Gamma(\{z\}, B(L))$ that has minimum modified cost among all paths in $\Gamma(\{z\}, \{z'\})$. The **coradius of** $L$ is defined by

$$CR(L) = \max_{z \in Z} c^*(z, B(L)).$$

If $L$ is singleton containing only the state $z \in S$, I write $CR(z)$.

The radius $R(L)$ measures the difficulty of leaving $B(L)$, whereas the coradius $CR(L)$ measures the difficulty of reaching $B(L)$ from a state most distant from $B(L)$.

**Theorem** (Ellison, 2000). *If $z \in Z$ is such that $R(z) > CR(z)$, then $z$ is the unique SSE.*

Now I apply the above result to prove our statement. We divide the proof into a series of lemmas. We use the following notations in our proofs. We frequently use the discrete coordinate system

$$\tilde{N} = \{1, 2, \ldots, n_1\} \times \{1, 2, \ldots, n_2\}$$

to refer to individuals in the population. Let $\Delta^k(z) \subset Z$ denote the set of states that can be transited from $z$ in exactly $k \geq 1$ period(s) with positive probability *under $P^0$*, i.e., the transition involves no mistake. Also, define $\Delta^0(z) = \{z\}$.

For state $z$, let $\tau_1^d(z)$ denote the minimum number of *adjacent rows* in the lattice $\tilde{N}$ such that each individual $i$ locating within these rows has $q_i < \pi$, $p_i < \varphi$ and $x_i = y_i = 0$. Let $\tau_2^d(z)$ denote the minimum number of *adjacent columns* such that each $i$ within the columns has $q_i < \pi$, $p_i < \varphi$ and $x_i = y_i = 0$. Define $\tau^d(z) = \min\{\tau_1^d(z), \tau_2^d(z)\}$. Also, let $\Lambda^d(z) \subset \tilde{N}$ denote the subset of individuals who locate at the adjacent rows or columns as specified. More precisely, if $\tau^d(z) = 0$, let $\Lambda^d(z) = \emptyset$. If $\tau^d(z) \geq 1$, we can always re-arrange the origin of the coordinate system so that

$$\Lambda^d(z) = \{1, \ldots, \tau_1^d(z)\} \times \{1, 2, \ldots, n_2\}$$
$$\cup \{1, 2, \ldots, n_1\} \times \{1, \ldots, \tau_2^d(z)\}.$$

For each $i \in \Lambda^d(z)$ we have $q_i < \pi$, $p_i < \varphi$ and $x_i = y_i = 0$. Re-arrangement of the origin is without loss of generality because vertexes are symmetric everywhere in the local interactions structure I consider.

Analogously, let $\tau_1^c(z)$ denote the minimum number of adjacent rows such that each $i$ within these rows has $q_i > \pi$, $p_i > \varphi$ and $x_i = y_i = 1$. And let $\tau_2^c(z)$ denote the minimum number of adjacent columns such that each $i$ within the columns has $q_i > \pi$ and $p_i > \varphi$ and $x_i = y_i = 1$. Define $\tau^c(z) = \min\{\tau_1^c(z), \tau_2^c(z)\}$. Let

$\Lambda^c(z) \subset \tilde{N}$ denote the subset of individuals who locate at the adjacent rows or columns and have $q_i > \pi$, $p_i > \varphi$ and $x_i = y_i = 1$. If $\tau^c(z) = 0$, let $\Lambda^c(z) = \emptyset$. If $\tau^c(z) \geq 1$, we re-arrange the origin of the coordinate system so that

$$\Lambda^c(z) = \{1, \ldots, \tau_1^c(z)\} \times \{1, 2, \ldots, n_2\}$$

$$\cup \{1, 2, \ldots, n_1\} \times \{1, \ldots, \tau_2^c(z)\}.$$

**Lemma 3.4.** *If one of the following holds then $\{z \in Z | \tau^d(z) \geq 2\} \subset B(z^D)$:*

1. $\pi > \frac{1}{2}$ and $\varphi > \frac{1}{2}$;

2. $\pi > \frac{3}{4}$ and $\varphi > \frac{1}{4}$;

3. $\pi > \frac{1}{4}$ and $\varphi > \frac{3}{4}$.

*Proof. Step 1: If $z \in Z$ is such that $\tau^d(z) \geq 2$, then $\tau^d(z') \geq \tau^d(z)$ for each $z' \in \Delta^t(z)$, $t \geq 1$.*

Let the dynamic $P^0$ start from $z^0 = (x^0, y^0, p^0, q^0) \in Z$ with $\tau^d(z^0) \geq 2$. Take $i \in \Lambda^d(z^0)$. Then we have $q_i^0 < \frac{1}{4}$, $p_i^0 < \frac{1}{4}$, $\sum_{j \in N_i} x_j^0 \leq 1$, and $\sum_{j \in N_i} y_j^0 \leq 1$. Let $z^1 = (x^1, y^1, p^1, q^1) \in \Delta^1(z^0)$, i.e., $z^1$ is a state in the next period that occurs with positive probability under $P^0$. In the case where $i$ does not update her expectations, we have $q_i^1 = q_i^0 < \frac{1}{4}$ and $p_i^1 = p_i^0 < \frac{1}{4}$. In the case where $i$ update her expectations, we have $q_i^1 = \frac{1}{4} \sum_{j \in N_i} x_j^0 \leq \frac{1}{4}$ and $p_i^1 = \frac{1}{4} \sum_{j \in N_i} y_j^0 \leq \frac{1}{4}$. In either case, given conditions $\pi > \frac{1}{4}$ and $\varphi > \frac{1}{4}$, we have $q_i^1 < \pi$ and $p_i^1 < \varphi$. Hence, $x_i^1 = y_i^1 = 0$. $i$ is an arbitrary individual taken from $\Lambda^d(z^0)$. We thus prove that $\Lambda^d(z^0) \subset \Lambda^d(z^1)$ and therefore $\tau^d(z^1) \geq \tau^d(z^0)$, for each $z^1 \in \Delta^1(z)$. By induction, we have $\tau^d(z') \geq \tau^d(z^0)$ for each $z' \in \Delta^t(z^0)$, $t \geq 1$.

*Step 2: Starting from $z$ with $\tau^d(z) \geq 2$, $\tau_1^d(z) < n_1$ and $\tau_1^d(z) < n_2$, there is a positive integer $t$ such that $\tau^d(z') > \tau^d(z)$ for some $z' \in \Delta^t(z)$.*

There are three cases.

Case 1: $\pi > \frac{1}{2}$ and $\varphi > \frac{1}{2}$. Let the initial state be $z^0 = (x^0, y^0, p^0, q^0) \in Z$ with $\tau^d(z^0) \geq 2$, $\tau_1^d(z) < n_1$ and $\tau_2^d(z) < n_2$. Consider individual $i = (\tau_1^d(z^0) +$

$1, \tau_2^d(z^0) + 1)$. We have $\sum_{j \in N_i} x_j^0 \leq 2$, and $\sum_{j \in N_i} y_j^0 \leq 2$. Let $i$ update expectations in the next period, so that her expectations in the next period are $p_i^1 = \frac{1}{4} \sum_{j \in N_i} y_j^0 \leq \frac{1}{2}$ and $q_i^1 = \frac{1}{4} \sum_{j \in N_i} x_j^0 \leq \frac{1}{2}$. We then have $p_i^1 < \varphi$ and $q_i^1 < \pi$, and thus $x_i^1 = y_i^1 = 0$. Hence, there is $z^1 \in \Delta^1(z^0)$ such that we have $p_i^1 < \varphi$, $q_i^1 < \pi$ and $x_i^1 = y_i^1 = 0$ for $i = (\tau_1^d(z^0) + 1, \tau_2^d(z^0) + 1)$. Now, consider integer $k$ with $0 \leq k \leq n_1 - \tau_1^d(z^0) - 1$. Assume $z^k \in \Delta^k(z^0)$ is such that $p_i^k < \varphi$, $q_i^k < \pi$ and $x_i^k = y_i^k = 0$ for $i = (\tau_1^d(z^0) + k + 1, \tau_2^d(z^0) + 1)$. By the argument above, there is $z^{k+1} \in \Delta^1(z^k)$ such that $p_i^{k+1} < \varphi$, $q_i^{k+1} < \pi$ and $x_i^{k+1} = y_i^{k+1} = 0$ for $i = (\tau_1^d(z^0) + k + 1, \tau_2^d(z^0) + 1)$. This inductive argument implies that there is $z^t = (x^t, y^t, p^t, q^t) \in \Delta^t(z^0)$ for some integer $t \geq 0$ such that $p_i^t < \varphi$, $q_i^t < \pi$ and $x_i^t = y_i^t = 0$ for each $i \in \{1, 2, \ldots, n_1\} \times \{1, 2, \ldots, \tau_2^d(z^0) + 1\}$. By the same argument, starting from such state $z^t$, there is $z^{t+k} = (x^{t+k}, y^{t+k}, p^{t+k}, q^{t+k}) \in \Delta^k(z^t)$ for some integer $k \geq 0$ such that $p_i^{t+k} < \varphi$, $q_i^{t+k} < \pi$ and $x_i^{t+k} = y_i^{t+k} = 0$ for each $i \in \{1, 2, \ldots, n_1\} \times \{1, 2, \ldots, \tau_2^d(z^0) + 2\}$. By induction, there is integer $t \geq 0$ such that for some $z^t = (x^t, y^t, p^t, q^t) \in \Delta^t(z^0)$ we have $p_i^t < \varphi$, $q_i^t < \pi$ and $x_i^t = y_i^t = 0$ for each $i \in \{1, 2, \ldots, n_1\} \times \{1, 2, \ldots, n_2\}$; such state $z^t$ is the defection equilibrium $z^D$, as desired.

Case 2: $\pi > \frac{3}{4}$ and $\varphi > \frac{1}{4}$. Let the initial state be $z^0 = (x^0, y^0, p^0, q^0) \in Z$ with $\tau^d(z^0) \geq 2$, $\tau_1^d(z) < n_1$ *and* $\tau_2^d(z) < n_2$. Take $i \in \{1, 2, \ldots, n_1\} \times \{\tau_2^d(z^0) + 1\}$. We have $\sum_{j \in N_i} x_j^0 \leq 3$. Let $i$ update expectations in the next period, so that her expectation of the proportion of cooperators in the next period is $q_i^1 = \frac{1}{4} \sum_{j \in N_i} x_j^0 \leq \frac{3}{4}$. We then have $q_i^1 < \pi$ and thus $y_i^1 = 0$. Individual $i$ is taken arbitrarily from $\{1, 2, \ldots, n_1\} \times \{\tau_2^d(z^0) + 1\}$. Hence, there is $z^1 = (x^1, y^1, p^1, q^1) \in \Delta^1(z^0)$ such that $q_i^1 < \pi$ and $y_i^1 = 0$ for each $i \in \{1, 2, \ldots, n_1\} \times \{\tau_2^d(z^0) + 1\}$. But then, for such $z^1$ we also have $\sum_{j \in N_i} y_j^1 \leq 1$ for each $i \in \{1, 2, \ldots, n_1\} \times \{\tau_2^d(z^0) + 1\}$. Let all individuals at the $(\tau_2^d(z^0) + 1)$th column update expectations in the next period. We obtain $p_i^2 \leq \frac{1}{4} < \varphi$ and $x_i^2 = 0$, alongside $q_i^2 < \pi$ and $y_i^2 = 0$, for each $i \in \{1, 2, \ldots, n_1\} \times \{\tau_2^d(z^0) + 1\}$. We have just shown that there is $z^2 \in \Delta^2(z^0)$ such that for each $i \in \{1, 2, \ldots, n_1\} \times \{1, 2, \ldots, \tau_2^d(z^0) + 1\}$ we have $p_i^2 < \varphi$,

$q_i^2 < \pi$ and $x_i^2 = y_i^2 = 0$. By induction, there is integer $t \geq 0$ such that there is $z^{2t} = (x^{2t}, y^{2t}, p^{2t}, q^{2t}) \in \Delta^{2t}(z^0)$ with $p_i^{2t} < \varphi$, $q_i^{2t} < \pi$ and $x_i^{2t} = y_i^{2t} = 0$ for each $i \in \{1, 2, \ldots, n_1\} \times \{1, 2, \ldots, n_2\}$, reaching the defection equilibrium $z^D$, as desired.

Case 3: $\pi > \frac{1}{4}$ and $\varphi > \frac{3}{4}$. The proof for this case is obtained from the proof for Case 2 by replacing: $\pi$ with $\varphi$, $\varphi$ with $\pi$, $x$ with $y$, $y$ with $x$, $p$ with $q$, and $q$ with $p$. $\qquad \square$

**Lemma 3.5.** *Consider local interactions. If one of the following holds then $R(z^D) \geq \lceil \min\{n_1, n_2\}/2 \rceil$:*

1. *$\pi > \frac{1}{2}$ and $\varphi > \frac{1}{2}$;*

2. *$\pi > \frac{3}{4}$ and $\varphi > \frac{1}{4}$;*

3. *$\pi > \frac{1}{4}$ and $\varphi > \frac{3}{4}$.*

*Proof.* By lemma 3.4, $\{z \in Z | \tau^d(z) \geq 2\} \subset B(z^D)$. It requires at least $\lceil \min\{n_1, n_2\}/2 \rceil$ mistakes to transit from $z^D$ to a state $z$ with $\tau^d(z) < 2$. Hence, $R(z^D) \geq \lceil \min\{n_1, n_2\}/2 \rceil$. $\qquad \square$

**Lemma 3.6.** *Consider local interactions. If $\pi > \frac{1}{2}$ and $\varphi > \frac{1}{2}$, then $CR(z^D) \leq 2$.*

*Proof.* Let $z^0 \in Z$ be an arbitrary initial state. We shall construct a path $\gamma = (z^0, z^1, \ldots, z^t)$ from $z^0$ to $z^t$ with $\tau^d(z^t) \geq 2$ such that $c^*(\gamma) \leq 2$, i.e., the modified cost of the path is at most two. Hence, $CR(z^D) \leq 2$.

*Step 1: Emergence of a local defection norm.*

Let $z^0$ be the initial state in period $t = 0$. Let each individual $i \in \{(1,1), (2,2)\}$ take $x_i^1 = y_i^1 = 0$ at $t = 1$. This event requires at most two mistakes. At $t = 2$, let each individual $i \in \{(1,2), (2,1)\}$ update expectations so that she has $p_i^2 = \frac{1}{4} \sum_{j \in N_i} y_j^1 \leq \frac{1}{2}$ and $q_i^2 = \frac{1}{4} \sum_{j \in N_i} x_j^1 \leq \frac{1}{2}$, and thus takes $x_i^2 = y_i^2 = 0$. This event does not require mistakes. At $t = 3$, let each $i \in \{(1,2), (2,1)\}$ do not update expectations so that we have $p_i^3 = p_i^2 < \varphi$, $q_i^3 = q_i^2 < \pi$, $x_i^3 = y_i^3 = 0$. Meanwhile, let each $i \in \{(1,1), (2,2)\}$ update expectations so that she has $p_i^3 =$

$\frac{1}{4} \sum_{j \in N_i} y_j^2 < \varphi$ and $q_i^3 = \frac{1}{4} \sum_{j \in N_i} x_j^2 < \pi$, and thus takes $x_i^3 = y_i^3 = 0$. This event does not require mistakes. Hence, starting from any state $z^0$, the minimum number of mistakes to reach a state $z^3$ in which each individual $i$ in the 2-by-2 block $\{(1, 2), (2, 1), (1, 2), (2, 1)\}$ holds $p_i^3 < \varphi$ and $q_i^3 < \pi$ and takes $x_i^3 = y_i^3 = 0$ is at most two.

*Step 2: Expansion of the defection norm.*

Let the initial state $z^0$ be such that $p_i^0 < \varphi$, $q_i^0 < \pi$ and $x_i^0 = y_i^0 = 0$ for each individual $i$ in the 2-by-2 block $\{(1, 2), (2, 1), (1, 2), (2, 1)\}$. Note that, for each integer $t \geq 0$ and $z = (x, y, p, q) \in \Delta^t(z^0)$, we have $p_i < \varphi$, $q_i < \pi$ and $x_i = y_i = 0$ for each $i \in \{(1, 2), (2, 1), (1, 2), (2, 1)\}$. Now let the dynamic follow a path under $P^0$, i.e., no mistake occurs, and transit to a recurrent class $L_1 \subset Z$ of $P^0$. Suppose the dynamic reaches state $z^t = (x^t, y^t, p^t, q^t) \in L_1$ in period $t \geq 3$. We have $p_i^t < \varphi$, $q_i^t < \pi$ and $x_i^t = y_i^t = 0$ for each $i \in \{(1, 2), (2, 1), (1, 2), (2, 1)\}$. Let $z^{t+1}$ be the state in the next period such that individual $i = (3, 1)$ takes $x_i^{t+1} = y_i^{t+1} = 0$, while all other individuals play rationally based on their expectations. This transition requires at most one mistake. Then, for $j = (3, 2)$, we have $\frac{1}{4} \sum_{i \in N_j} x_i^{t+1} \leq \frac{1}{2}$ and $\frac{1}{4} \sum_{i \in N_j} y_i^{t+1} \leq \frac{1}{2}$. In the next period, let $j = (3, 2)$ update expectations, leading to $q_j^{t+2} < \pi$ and $p_j^{t+2} < \varphi$, and thus $x_j^{t+2} = y_j^{t+2} = 0$. In the next period, $t + 3$, let individual $(3, 1)$ update expectations while individual $(3, 2)$ do not. Then for each individual $i \in \{1, 2, 3\} \times \{1, 2\}$ we have $p_i^{t+3} < \varphi$, $q_i^{t+3} < \pi$ and $x_i^{t+3} = y_i^{t+3} = 0$. Note that $R(L) \geq 1$ for any recurrent class $L$ of $P^0$. Hence, the modified cost of transiting from $z^0$ to the specified state $z^{t+3}$ is zero! This is because, to calculate the modified cost of the transition, we need to extract $R(L_1)$ (which is at least 1) from the minimum number of mistakes to reach $z^{t+3}$ (which is at most 1).

By induction, we can construct a path $(z^0, \ldots, z^t, \ldots, z^{t+k})$ to a state $z^{t+k}$ with zero modified cost such that $p_i^{t+k} < \varphi$, $q_i^{t+k} < \pi$ and $x_i^{t+k} = y_i^{t+k} = 0$ for each $i \in \{1, 2, \ldots, n_1\} \times \{1, 2\}$.

By the same inductive argument, we can construct a path

$$\left(z^0, \ldots, z^t, \ldots, z^{t+k}, \ldots, z^{t+k+\ell}\right)$$

to a state $z^{t+k+\ell}$ with zero modified cost such that $p_i^{t+k+\ell} < \varphi$, $q_i^{t+k+\ell} < \pi$ and $x_i^{t+k+\ell} = y_i^{t+k+\ell} = 0$ for each $i$ in the set

$$\Lambda^d(z^{t+k+\ell}) = \{1,2\} \times \{1,2,\ldots,n_2\}$$

$$\cup \{1,2,\ldots,n_1\} \times \{1,2\}.$$

We have $\tau^d(z^{t+k+\ell}) \geq 2$. By lemma 3.4, $z^{t+k+\ell} \in B(z^D)$.

Hence, within $3 + t + k + \ell$ periods, the dynamic transits from an arbitrary initial state $z$ to $z^{t+k+\ell} \in B(z^D)$. The modified cost of the constructed transition path is at most 2, which is the minimum number of mistakes required for the 2-by-2 block of defection and no punishment to emerge in Step 1. □

**Lemma 3.7.** *Consider local interactions. If $\pi > \frac{3}{4}$ and $\varphi > \frac{1}{4}$, or if $\pi > \frac{1}{4}$ and $\varphi > \frac{3}{4}$, then $CR(z^D) \leq 1$.*

*Proof.* We show the case of $\pi > \frac{3}{4}$ and $\varphi > \frac{1}{4}$. The case of $\pi > \frac{3}{4}$ and $\varphi > \frac{1}{4}$ can be analogously shown. Let $z^0 \in Z$ be an arbitrary initial state. Similar to the proof of lemma, I shall construct a path $\gamma = (z^0, z^1, \ldots, z^t)$ from $z^0$ to $z^t$ with $\tau^d(z^t) \geq 2$ such that $c^*(\gamma) \leq 1$. Hence, $CR(z^D) \leq 1$.

*Step 1: Emergence of a local defection norm.*

Let $z^0$ be the initial state in period $t = 0$. In the next period, $t = 1$, let individual $i = (2,2)$ defect, i.e., $x_i^1 = 0$. This event requires at most one mistake. Then, for $i \in \{(2,1),(1,2),(3,2),(2,3)\}$, we have $\sum_{j \in N_i} x_j^1 \leq 3$. At period $t = 2$, let each $i \in \{(2,1),(1,2),(3,2),(2,3)\}$ update expectations. We obtain $q_i^2 \leq \frac{3}{4} < \pi$ and $y_i^2 = 0$ for each $i \in \{(2,1),(1,2),(3,2),(2,3)\}$. Transition from period $t = 1$ to $t = 2$ does not require a mistake. In the next period, $t = 3$, let individual $(2,2)$ update expectations; then we have $p_{(2,2)}^3 = 0 < \varphi$ and $x_{(2,2)}^3 = 0$. Meanwhile, let

99

the four neighbors of $(2,2)$ do not update expectations. Then we have $q_i^3 < \pi$ and $y_i^3 = 0$ for each $i \in \{(2,1),(1,2),(3,2),(2,3)\}$. Notice that since the state $z^3$, individual $(2,2)$ defects and the four neighbors of $(2,2)$ do not punish defectors in all future periods, unless some of them make a mistake. More precisely, for each $k \geq 1$, $z = (x,y,p,q) \in \Delta^k(z^3)$, we have $p_{(2,2)} < \varphi$, $x_{(2,2)} = 0$, as well as $q_i < \pi$ and $y_i = 0$ for each $i \in \{(2,1),(1,2),(3,2),(2,3)\}$.

*Step 2: Expansion of the defection norm.*

Let the initial state $z^0$ be such that $p_{(2,2)}^0 < \varphi$, $x_{(2,2)}^0 = 0$, $q_i^0 < \pi$ and $y_i^0 = 0$ for each $i \in \{(2,1),(1,2),(3,2),(2,3)\}$. Now let the dynamic follow a path under $P^0$, i.e., no mistake occurs, and transit to a recurrent class $L_1 \subset Z$ of $P^0$. Suppose the dynamic reaches state $z^t = (x^t, y^t, p^t, q^t) \in L_1$ in period $t \geq 3$. For $z^t$, we have $p_{(2,2)}^t < \varphi$, $x_{(2,2)}^t = 0$, $q_i^t < \pi$ and $y_i^t = 0$ for each $i \in \{(2,1),(1,2),(3,2),(2,3)\}$. In the next period, $t+1$, let individual $(3,2)$ take $x_{(3,2)}^{t+1} = 0$, leading to state $z^{t+1}$. This event requires at most one mistake. Following the argument of Step 1, there is some $z^{t+3} \in \Delta^2(z^{t+1})$ such for each $i \in \{(2,2),(3,2)\}$ and $j \in N_i$: 1) $p_i^{t+3} < \varphi$, $x_i^{t+3} = 0$ , and 2) $q_j^{t+3} < \pi$ and $y_j^{t+3} = 0$. Notice that individuals $(2,2)$ and $(3,2)$ are defecting as well as not punishing now, and all other neighbors of them are not punishing. The modified cost of transiting from $z^0$ to the specified state $z^{t+3}$ is zero. This is because the modified cost of the transition equals to the minimum number of mistakes to reach $z^{t+3}$ (which is at most 1) *less* $R(L_1)$ (which is at least 1).

By induction, we can construct a path $(z^0, \ldots, z^t, \ldots, z^{t+k})$ to a state $z^{t+k}$ with zero modified cost such that for each $i \in \{1,2,\ldots,n_1\} \times \{2\}$ and $j \in N_i$: 1) $p_i^{t+k} < \varphi$, $q_i^{t+k} < \pi$ and $x_i^{t+k} = y_i^{t+k} = 0$, and 2) $q_j^{t+k} < \pi$ and $y_j^{t+k} = 0$. This means that, at state $z^{t+k}$, all individuals at the 2nd column defect and do not punish, and all individuals at the 1st and the 3rd columns do not punish. Then, each individual at the 1st and the 3rd columns has at most one neighbor punishing. That is, for each $i \in \{1,2,\ldots,n_1\} \times \{1,3\}$ we have $\sum_{j \in N_i} y_j^{t+k} \leq 1$. Let each individual at the 1st and the 3rd columns update expectations in the next period. We obtain

$p_i^{t+k+1} \leq \frac{1}{4} < \varphi$ and $x_i^{t+k+1} = 0$ for each $i \in \{1, 2, \ldots, n_1\} \times \{1, 3\}$. To sum up, in period $t + k + 1$, the dynamic reaches a state $z^{t+k+1}$ such that $p_i^{t+k+1} < \varphi$, $q_i^{t+k+1} < \pi$ and $x_i^{t+k+1} = y_i^{t+k+1} = 0$ for each $i \in \{1, 2, \ldots, n_1\} \times \{1, 2, 3\}$.

By the same inductive argument, we can construct a path

$$\left( z^0, \ldots, z^t, \ldots, z^{t+k}, \ldots, z^{t+k+\ell} \right)$$

to a state $z^{t+k+\ell}$ with zero modified cost such that $p_i^{t+k+\ell} < \varphi$, $q_i^{t+k+\ell} < \pi$ and $x_i^{t+k+\ell} = y_i^{t+k+\ell} = 0$ for each $i$ in the set

$$\Lambda^d(z^{t+k+\ell}) = \{1, 2, 3\} \times \{1, 2, \ldots, n_2\}$$
$$\cup \{1, 2, \ldots, n_1\} \times \{1, 2, 3\}.$$

Since $\tau^d(z^{t+k+\ell}) \geq 3$, by lemma 3.4 we have $z^{t+k+\ell} \in B(z^D)$.

Hence, within $3 + t + k + \ell$ periods, the dynamic transits from an arbitrary initial state $z$ to $z^{t+k+\ell} \in B(z^D)$. The modified cost of the constructed transition path is at most 1, which is the minimum number of mistakes required for the small block of defection and no punishment in Step 1 to emerge. $\qquad \square$

**Lemma 3.8.** *If one of the following holds then $\{z \in Z | \tau^c(z) \geq 2\} \subset B(z^C)$:*

1. *$\pi < \frac{1}{2}$ and $\varphi < \frac{1}{2}$;*

2. *$\pi < \frac{3}{4}$ and $\varphi < \frac{1}{4}$;*

3. *$\pi < \frac{1}{4}$ and $\varphi < \frac{3}{4}$.*

*Proof.* Completely analogous to the proof of Lemma 3.4, so omitted. $\qquad \square$

**Lemma 3.9.** *Consider local interactions. We have $R(z^C) \geq \lceil \min\{n_1, n_2\}/2 \rceil$ if one of the following holds:*

1. *$\pi < \frac{1}{2}$ and $\varphi < \frac{1}{2}$;*

2. *$\pi < \frac{3}{4}$ and $\varphi < \frac{1}{4}$;*

*3.* $\pi < \frac{1}{4}$ *and* $\varphi < \frac{3}{4}$.

*Proof.* Completely analogous to the proof of Lemma 3.5, so omitted. □

**Lemma 3.10.** *Consider local interactions. If* $\pi < \frac{1}{2}$ *and* $\varphi < \frac{1}{2}$, *then* $CR(z^C) \leq 2$.

*Proof.* Completely analogous to the proof of Lemma 3.6, so omitted. □

**Lemma 3.11.** *Consider local interactions. If* $\pi < \frac{1}{4}$ *and* $\varphi < \frac{3}{4}$, *or if* $\pi < \frac{3}{4}$ *and* $\varphi < \frac{1}{4}$, *then* $CR(z^C) \leq 1$.

*Proof.* Completely analogous to the proof of Lemma 3.7, so omitted. □

## Proof of Theorem 3.4

**Lemma 3.12.** *Consider the unperturbed adaptive dynamic* $P^0$ *with mobility of degree* $m$. *Suppose* $\pi < 1$. *Consider the generic cases with* $\lceil \varphi(n-m) \rceil \neq \lceil (1-\pi)(n-m) \rceil$ *and* $\lceil \pi(n-m) \rceil \neq \lceil (1-\varphi)(n-m) \rceil$. *There are exactly four recurrent classes:* $Z^c$, $Z^d$, $\tilde{Z}^{c1}$, *and* $\tilde{Z}^{c2}$.

*Proof.* For a state $z \in Z$ and $t \geq 1$, let $\Delta^t(z)$ denote the set of states that the dynamic $P^0$ with mobility can transit to from $z$ in exactly $t$ periods with positive probability.

*Step 1: $Z^d$ is a recurrent class of $P^0$ with mobility.*

First, I show that if $P^0$ starts from a state in $Z^d$, it stays in $Z^d$ thereafter. Let $z^0$ be an initial state in which for each $i \in N$, $\ell = 1, 2$, we have $x_i^0 = y_i^0 = 0$ and $p_{i\ell}^0 = q_{i\ell}^0 = 0$. Consider the next period and an arbitrary individual $i \in N$. For each $\ell = 1, 2$, we have either $p_{i\ell}^1 = p_{i\ell}^0$ or $p_{i\ell}^1 = \frac{1}{N_\ell^0} \sum_{j \in N_\ell^0} y_j^0$; in either case, we have $p_{i\ell}^1 = 0 < \varphi$. Likewise, we have $q_{i\ell}^1 = 0 < \pi$. Hence, we have $x_i^1 = y_i^1 = 0$. By induction, for each $t \geq 0$ and state $z \in \Delta^t(z^0)$, we have $x_i = y_i = 0$ and $p_{i\ell} = q_{i\ell} = 0$ for each $i \in N$, $\ell = 1, 2$.

Second, there is positive probability of transiting from any state in $Z^d$ to any another one within finite periods. The difference between the states in $Z^d$ is that they have different partitions of the population $N$ into the two societies. However,

102

when an individual $i$ has $p_{i\ell} = q_{i\ell} = 0$ for both $\ell = 1$ and $\ell = 2$, $i$ has expected material payoffs $\frac{1}{2}(\bar{a} + b)$ in either of the societies. Hence, for each period, each $i$ randomly chooses a society to live. Therefore, starting from any $z \in Z^d$, there is a positive probability of transiting to any $z' \in Z^d$ in the next period.

*Step 2: If $\pi < 1$, then $Z^c$ is a recurrent class of $P^0$ with mobility.*

Analogous to the proof of Step 1a, so omitted. The condition $\pi < 1$ ensures that, if $q_{i\ell} = 1$, then $q_{i\ell} > \pi$.

*Step 3: If $\pi < 1$, then $\tilde{Z}^{c1}$ is a recurrent class of $P^0$ with mobility.*

First, I show that, if $P^0$ starts from a state in $\tilde{Z}^c$, it stays in $\tilde{Z}^c$ thereafter. Let $z^0$ be an initial state such that: there are $n + m$ individuals in society 1 and $n - m$ individuals in society 2; in society 1, each $i$ has $x_i = y_i = 1$ and $p_{i1} = q_{i1} = 1$; in society 2, each $i$ has $x_i = y_i = 0$ and $p_{i2} = q_{i2} = 0$. Consider the next period and an arbitrary individual $i \in N$. Either because $p_{i\ell}^1 = p_{i\ell}^0$ or because $p_{i\ell}^1 = \frac{1}{N_\ell^0} \sum_{j \in N_\ell^0} y_j^0$, we have $p_{i1}^1 = 1$ for society 1 and $p_{i2}^1 = 0$ for society 2. Likewise, $q_{i1}^1 = 1$ for society 1 and $q_{i2}^1 = 0$ for society 2. Hence, $i$ takes $x_i^1 = y_i^1 = 1$ if $i$ lives in society 1, while taking $x_i^1 = y_i^1 = 0$ if $i$ lives in society 2. The expected payoffs of living in society 1 is $\frac{1}{2}(a + \bar{b})$, while the expected payoffs of living in society 2 is $\frac{1}{2}(\bar{a} + b)$. By Assumption 3.3, cooperation yields higher aggregate efficiency, i..e, $a + \bar{b} > \bar{a} + b$. Hence, $i$ chooses to live in society 1 upon having the opportunity to choose. Since every individual chooses society 1 upon having the opportunity to choose, society 1 has $n + m$ individuals. By contrast, society 2 only has $n - m$ individuals, those who do not have an opportunity to choose where to live. By induction, for each $t \geq 1$, we have $\Delta^t(z^0) \subset \tilde{Z}^{c1}$.

Second, each individual has positive probability of *not* having the opportunity to choose where to live and being assigned to society 2. Thus, $\Delta^t(z^0) = \tilde{Z}^{c1}$ for each $t \geq 1$, $z^0 \in \tilde{Z}^{c1}$.

*Step 4: If $\pi < 1$, then $\tilde{Z}^{c2}$ is a recurrent class of $P^0$ with mobility.*

The same as the proof of Step 3, so omitted.

*Step 5: Starting from any initial state $z \in Z$, there is a positive probability*

*that $P^0$ transits to $Z^c \cup Z^d \cup \tilde{Z}^{c1} \cup \tilde{Z}^{c2}$ within finite periods.*

Let $z^0 \in Z$ be an arbitrary initial state. From period 0 to period 1, let all individuals update expectations. Then we have $q_{i\ell}^1 = q_{j\ell}^1$ and $p_{i\ell}^1 = p_{j\ell}^1$ for each $i, j \in N$ in period 1. Let all individuals take the same actions in response to $q_{i\ell}^1$ and $p_{i\ell}^1$. Let $\ell = 1, 2$. There are four cases about society $\ell$: 1) $q_{i\ell}^1 \geq \pi$, $p_{i\ell}^1 \geq \varphi$ and $x_i^1 = y_i^1 = 1$ for each $i \in N_\ell^1$; 2) $q_{i\ell}^1 \leq \pi$, $p_{i\ell}^1 \leq \varphi$ and $x_i^1 = y_i^1 = 0$ for each $i \in N_\ell^1$; 3) $q_{i\ell}^1 \leq \pi$, $p_{i\ell}^1 \geq \varphi$, $x_i^1 = 1$ and $y_i^1 = 0$ for each $i \in N_\ell^1$; and 4) $q_{i\ell}^1 \geq \pi$, $p_{i\ell}^1 \leq \varphi$, $x_i^1 = 0$ and $y_i^1 = 1$ for each $i \in N_\ell^1$. In the first two cases listed above, let all individuals update expectations in the next period. Then in period 2, the dynamic reaches a state $z^2 \in Z^c \cup Z^d \cup \tilde{Z}^{c1} \cup \tilde{Z}^{c2}$.

Now consider the third case; fourth case can be analogously shown. At period 2, let all individuals update *only* their expectation of the proportion of *punishers*, so that $p_{i\ell}^2 = \frac{1}{N_\ell^1} \sum_{j \in N_\ell^1} y_i^1 = 0$ and $q_{i\ell}^2 = q_{i\ell}^1 \leq \pi$ for each $i \in N$. Then it occurs with positive probability that $x_i^2 = y_i^2 = 0$ for each $i \in N_\ell^2$. Let this event occur. At period 3, let all individuals update expectations, so that $q_{i\ell}^3 = p_{i\ell}^3 = 0$ and $x_i^3 = y_i^3 = 0$ for each $i \in N_\ell^3$. By the same construction, for the other society $\ell' = 2 - \ell$, we can have either $q_{i\ell'}^3 = p_{i\ell'}^3 = 0$ and $x_i^3 = y_i^3 = 0$ for each $i \in N_{\ell'}^3$, or $q_{i\ell'}^3 = p_{i\ell'}^3 = 1$ and $x_i^3 = y_i^3 = 1$ for each $i \in N_{\ell'}^3$ in period 3. Then in period 3, the dynamic reaches a state $z^3 \in Z^c \cup Z^d \cup \tilde{Z}^{c1} \cup \tilde{Z}^{c2}$. $\qquad \square$

In what follows, I divide the proof of Theorem 3.4 into three lemmas, each for a different case. Let $\Omega = \{Z^c, Z^d, \tilde{Z}^{c1}, \tilde{Z}^{c2}\}$ be the collection of all recurrent classes of $P^0$ with mobility. For each case, I identify the stochastic tree that has minimum cost among all stochastic trees defined on $\Omega$. The desired results then follow from Young's (1993) theorem. Let $\alp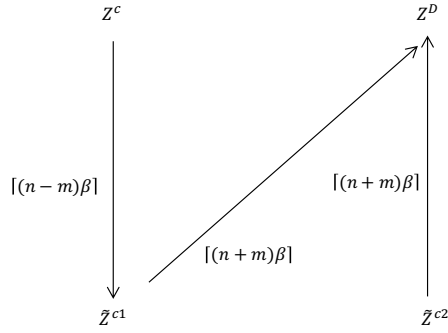ha = \min\{\pi, \varphi\}$ and $\beta = \min\{1-\pi, 1-\varphi\}$. Figure 3.8 shows the minimum numbers of mistakes required to transit from one recurrent class to another.

**Figure 3.8:** The costs (the minimum numbers of mistakes) of transiting from one recurrent class to another under the dynamics with mobility; $\alpha = \min\{\pi, \varphi\}$; $\beta = \min\{1 - \pi, 1 - \varphi\}$.

(a) The $Z^c$-tree with minimum cost among all $Z^c$-trees.



(b) The $Z^d$-tree with minimum cost among all $Z^d$-trees.



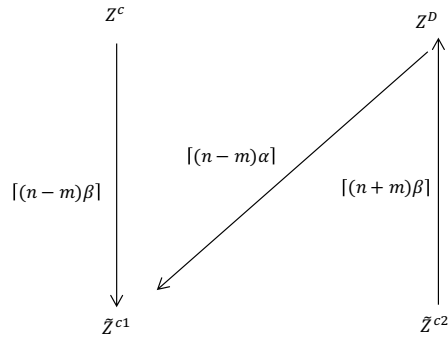(c) The $\tilde{Z}^{c1}$-tree with minimum cost among all $\tilde{Z}^{c1}$-trees.



**Figure 3.9:** Candidates for the stochastic tree with minimum cost if $\pi + \varphi < 1 + \frac{m}{n}(1 - |\pi - \varphi|)$.

**Lemma 3.13.** *Consider the unperturbed adaptive dynamic $P^0$ with mobility of degree $m$. Consider the generic cases with $\lceil \varphi(n-m) \rceil \neq \lceil (1-\pi)(n-m) \rceil$ and $\lceil \pi(n-m) \rceil \neq \lceil (1-\varphi)(n-m) \rceil$. If $\pi + \varphi < 1$, then the SSEs are states in $Z^c$.*

*Proof.* Given $\pi + \varphi < 1$, and for the generic cases, we have $\lceil (n+m)\beta \rceil > \lceil (n-m)\beta \rceil > \lceil (n-m)\alpha \rceil$. Figure 3.9 shows three candidate stochastic trees. Given $\lceil (n+m)\beta \rceil > \lceil (n-m)\beta \rceil > \lceil (n-m)\alpha \rceil$, the stochastic tree with minimum cost among all stochastic trees defined on $\Omega$ must come from one of the three candidates. The cost of the $Z^c$-tree with minimum cost among all $Z^c$-trees is $3\lceil (n-m)\alpha \rceil$. The cost of the $Z^d$-tree with minimum cost among all $Z^d$-trees is $\lceil (n-m)\alpha \rceil + \lceil (n-m)\beta \rceil + \lceil (n+m)\beta \rceil$. The cost of the $\tilde{Z}^{c1}$-tree with minimum cost among all $\tilde{Z}^{c1}$-trees is $2\lceil (n-m)\alpha \rceil + \lceil (n-m)\beta \rceil$. Hence, the $Z^c$-tree shown in Figure 3.9(a) has minimum cost among all stochastic trees defined on $\Omega$. The desired result then follows from Young's (1993) theorem. $\qquad\square$

**Lemma 3.14.** *Consider the unperturbed adaptive dynamic $P^0$ with mobility of degree $m$. Consider the generic cases with $\lceil \varphi(n-m) \rceil \neq \lceil (1-\pi)(n-m) \rceil$ and $\lceil \pi(n-m) \rceil \neq \lceil (1-\varphi)(n-m) \rceil$. If $1 < \pi + \varphi < 1 + \frac{m}{n}(1 - |\pi - \varphi|)$, then the SSEs are states in $\tilde{Z}^{c1} \cup \tilde{Z}^{c2}$.*

*Proof.* The condition $1 < \pi + \varphi < 1 + \frac{m}{n}(1 - |\pi - \varphi|)$ implies $\lceil (n-m)\beta \rceil < \lceil (n-m)\alpha \rceil < \lceil (n+m)\beta \rceil$ for the generic cases. The candidate stochastic trees with minimum costs are the ones shown in Figure 3.9; they are the same ones as the case of $\pi + \varphi < 1$. The cost of the $Z^c$-tree in Figure 3.9 is $3\lceil (n-m)\alpha \rceil$. The cost of the $Z^d$-tree is $\lceil (n-m)\alpha \rceil + \lceil (n-m)\beta \rceil + \lceil (n+m)\beta \rceil$. The cost of the $\tilde{Z}^{c1}$-tree is $2\lceil (n-m)\alpha \rceil + \lceil (n-m)\beta \rceil$. Given $\lceil (n-m)\beta \rceil < \lceil (n-m)\alpha \rceil < \lceil (n+m)\beta \rceil$, the $\tilde{Z}^{c1}$-tree has the minimum cost among the three. But since $\tilde{Z}^{c1}$ and $\tilde{Z}^{c2}$ are completely symmetric, there must also be a $\tilde{Z}^{c2}$-tree having the same cost as the $\tilde{Z}^{c1}$-tree in Figure 3.9. By Young's (1993) theorem, the SSE states are contained in $\tilde{Z}^{c1} \cup \tilde{Z}^{c2}$. $\qquad\square$

(a) The $Z^c$-tree with minimum cost among all $Z^c$-trees.

$$Z^c \qquad\qquad Z^D$$

$[(n-m)\alpha]$

$[(n-m)\alpha]$

$[(n+m)\beta]$

$$\tilde{Z}^{c1} \qquad\qquad \tilde{Z}^{c2}$$

(b) The $Z^d$-tree with minimum cost among all $Z^d$-trees.

$$Z^c \qquad\qquad Z^D$$

$[(n-m)\beta]$

$[(n+m)\beta]$

$[(n+m)\beta]$

$$\tilde{Z}^{c1} \qquad\qquad \tilde{Z}^{c2}$$

(c) The $\tilde{Z}^{c1}$-tree with minimum cost among all $\tilde{Z}^{c1}$-trees.

$$Z^c \qquad\qquad Z^D$$

$[(n-m)\beta]$

$[(n-m)\alpha]$

$[(n+m)\beta]$

$$\tilde{Z}^{c1} \qquad\qquad \tilde{Z}^{c2}$$

**Figure 3.10:** Candidates for the stochastic tree with minimum cost if $\pi + \varphi > 1 + \frac{m}{n}(1 - |\pi - \varphi|)$.

**Lemma 3.15.** *Consider the unperturbed adaptive dynamic $P^0$ with mobility of degree $m$. Consider the generic cases with $\lceil \varphi(n - m) \rceil \neq \lceil (1 - \pi)(n - m) \rceil$ and $\lceil \pi(n - m) \rceil \neq \lceil (1 - \varphi)(n - m) \rceil$. If $\pi + \varphi > 1 + \frac{m}{n}(1 - |\pi - \varphi|)$, then the SSEs are states in $Z^d$.*

*Proof.* $\pi + \varphi > 1 + \frac{m}{n}(1 - |\pi - \varphi|)$ implies $\lceil (n-m)\beta \rceil < \lceil (n+m)\beta \rceil < \lceil (n-m)\alpha \rceil$ for the generic cases. Figure 3.10 shows the candidate stochastic trees of the one with minimum cost among all stochastic trees defined on $\Omega$. The cost of the $Z^c$-tree in Figure 3.10 is $2\lceil (n - m)\alpha \rceil + \lceil (n + m)\beta \rceil$. The cost of the $Z^d$-tree is $\lceil (n - m)\beta \rceil + 2\lceil (n + m)\beta \rceil$. The cost of the $\tilde{Z}^{c1}$-tree is $\lceil (n - m)\alpha \rceil + \lceil (n - m)\beta \rceil + \lceil (n + m)\beta \rceil$. The $Z^d$-tree has the minimum cost among the three. Thus, we obtain the desired result from Young's (1993) theorem. $\square$

# Chapter 4

# Fairness Perceptions and Punishment Types[1]

## 4.1 Introduction

Now there has been a large body of experimental evidence showing that individuals' behavior is influenced by fairness and reciprocity concerns (for reviews of the evidence see Fehr and Schmidt (2006) and O'Connell and Siafarikas (2010)). Studies have also shown that fairness and reciprocity may play a key role in market transactions (Kahneman et al., 1986) and optimal contracts in workplaces (e.g., Barr and Serneels, 2009; Fehr et al., 1997; Fehr and Schmidt, 2000, 2004). However, it is not clear what determines which behavior people perceive as fair and others not, their *fairness perceptions*. Cross-cultural studies show that cooperation and punishment of selfish behavior vary dramatically across societies (Herrmann et al., 2008; Henrich et al., 2006, 2010). These studies indicate that fairness perceptions are not fixed; they vary across societies and could be changed. However, to date, still relatively little is known about how fairness perceptions are formed and what policies could change them. Building on previous theoretical and empirical literature (Bicchieri, 2006; Cooper and Dutcher, 2011; Herz and Taubinsky, 2017; Kahneman et al., 1986; Sugden, 2000, 2004), we hypothesize that providing

---

[1]This chapter is a joint work with Lucas Molleman and Dennie van Dolder.

different information regarding a behavior's prevalence in society would change individuals' fairness perception about the behavior. We examine this hypothesis experimentally in a prisoner's dilemma game involving choices between cooperation and free riding. Specifically, we test: if we inform people that the proportion of cooperators in the population is higher, it would reduce the legitimacy of free riding and, thereby, it is more likely that individuals punish free riders.

It contributes to the literature in the following aspects. *First*, if the more cooperators in a society the more people punish free riders, then we can explain the dramatic variations in cooperation and punishment observed across societies (Herrmann et al., 2008; Henrich et al., 2006, 2010) through belief channels: people from different societies may not have different preferences (Stigler and Becker, 1977); it could be that they merely have different beliefs about the behavior of others. *Second*, the tested hypothesis has direct policy implications. If the mere information about the proportion of cooperators can affect punishment of free riders, then it suggests that a pure informational policy might be effective in changing fairness perceptions and promoting cooperation. In contrast, if we fail to find the effect, it suggests that a pure informational policy might have limited power. *Third*, since Olson (1965), economists have developed a variety of theories that are relevant to understanding private contributions of public goods (Andreoni, 1990; Bergstrom et al., 1986; Dufwenberg and Kirchsteiger, 2004; Fehr and Schmidt, 1999; Rabin, 1993; Sugden, 1982, 1984). Examining whether the proportion of cooperators affects punishment of free riders allows us to disentangle these theories that are otherwise difficult to disentangle.

Our experiment is as follows. We match participants into pairs and ask each pair to play the following two-stage game. In the first stage, the two participants play a prisoner's dilemma game, i.e., they choose simultaneously whether to cooperate or to free ride. In the second stage, participants are given the opportunity to punish their free riding partners by assigning up to 10 deduction points. For each deduction point that a participant assigns to her partner, one point is deducted

from the participant and three points are deducted from her partner. Hence, punishment is costly, but the damage to the other party is greater than the cost of conducting it. We first run some sessions to obtain an estimate of the proportion of cooperators in the subject pool. In later sessions, we inform each participant about the proportion of cooperators before they make punishment decisions. More specifically, we randomly draw a sample of participants from previous sessions, and inform participants in the current sessions *how many of the randomly selected previous participants are cooperators*. We provide this information to change participants' belief about the proportion of cooperators in the population. We use the *strategy method* (Selten, 1967) to elicit each participant's *complete punishment strategy*. We classify all possible cases into five situations regarding the proportion of cooperators in the selected sample. We ask participants to indicate, for *each* of the five situations, how many deduction points they assign. A participant's complete punishment strategy includes her deduction points to assign in all five situations.

We only allow participants to punish their partners if their partners free ride. Thus, if a cooperator chooses to punish, it is not because she expects her partner to punish her. Cooperators would punish only for the following reasons: *i*) the cooperator wants to sustain a community-level threat to deter free riders, and such threat has diminishing marginal return—*the diminishing marginal return hypothesis* (Becker, 1974; Bergstrom et al., 1986; Kandori, 1992), *ii*) the cooperator considers her partner's behavior as morally wrong—*norm-based reciprocity* (Bicchieri, 2006; Sugden, 1984), or, *iii*) the cooperator considers her partner as having bad intentions—*intention–based reciprocity* (Dufwenberg and Kirchsteiger, 2004; Rabin, 1993) or she wants to reduce payoff difference between her partner and herself—*inequality aversion* (Fehr and Schmidt, 1999). We observe that these theories make different predictions regarding how punishment of free riders changes according to the provided information about the proportion of cooperators. The diminishing marginal return hypothesis predicts that the higher the proportion of

cooperators, the *less* people would punish free riders. To the contrary, norm-based reciprocity predicts that the higher the proportion of cooperators, the *more* people would punish free riders. Different from both, inequality aversion and intention-based reciprocity predict that punishment of free riders is *independent* of the proportion of cooperators. Hence, our experiment provides an opportunity to disentangle these theories. Indeed, the only key difference between intention-based reciprocity and norm-based reciprocity is whether the perceived fair behavior depends on the behavior of the majority of the population. Hence, examining how punishment of free riders changes in response to a change in the proportion of cooperators is the key to know which theory provides a better account for the formation of fairness perceptions.

Our results are as follows. On average, the proportion of cooperators does *not* significantly affect punishment of free riders. However, this is not because no participant responds to the provided information about the proportion of cooperators. Instead, it is due to the remarkable heterogeneity among participants. 42% of our participants do not punish free-riding partners regardless of the proportion of cooperators, as consistent with the null hypothesis that they seek to maximize their material payoffs. 21% of our participants behave in line with inequality aversion and intention-based reciprocity: they assign the same positive number of deduction points to free riders, *independent* of the proportion of cooperators. By contrast, 23% of participants do condition their punishment on the proportion of cooperators. However, on the one hand, 13% behave in line with norm-based reciprocity: they *increase* deduction points as the proportion of cooperators increases; on the other, 10% *decrease* deduction points as the proportion of cooperators increases, as consistent with the diminishing marginal return hypothesis. Hence, we observe co-existence of four distinct punishment types. Each punishment type is associated with a category of theory we consider, plus there is a type of participants consistent with material payoffs maximization. The four punishment types add up to 86% of the whole sample, while the rest 14% are not consistent with

any specified theory.[2]

There have been several experimental studies investigating how fairness perceptions are formed and might be changed according to the frequency of behaviors in the population. The experiment of Falk et al. (2006) shows that introducing the law of minimum wage—forcing firms to make high wage offers—increases the reservation wages of workers even after the restriction is removed at a later stage. The result suggests that the minimum wage law increases workers' reference point regarding what offers are considered as fair. Cooper and Dutcher (2011) conduct a meta-analysis on ultimatum game experiments with multiple rounds. They find that, on average, responders who have experienced high offers in early rounds reject low offers more often in later rounds. Hence, which offer is regarded as fair seems to depend on past experience. Peysakhovich and Rand (2016) conduct an experiment to see whether they can induce different "norms" of cooperation and defection in the lab. Their experiment involves two parts. In part one, participants play repeated prisoner's dilemma games; this part is used to induce more cooperation in some sessions and more defection in the others. In part two, they measure the impacts of the experience in part one on participants' pro-social attitudes and willingness to punish selfish behavior in a series of one-shot games. They find that, on average, participants having experienced cooperation in part one behave more pro-socially and are more willing to punish selfish behavior in part two. Finally, Herz and Taubinsky (2017) examine transaction behavior in market games using the same two-part design as Peysakhovich and Rand's. They find that transaction experience in part one influences behavior in part two through two channels: $i$) personal payoff experience, i.e., individuals who are used to receiving certain payoffs feel entitled to obtain the same payoffs in similar transactions; and $ii$) the information effect, i.e., individuals make inference of social norms from observing the transactions of others.

---

[2]We verify that no any of the four punishment types can be explained by pure randomness: if each participant had rolled a dice to make a decision, the four punishment types would have altogether been about only 4% of the sample, instead of 86%.

Our experiment differs from the experiments mentioned above and thus complements their studies in two aspects. *First*, while previous experiments have the merit of showing path-dependent behavior directly, their results admit several explanations. For example, Peysakhovich and Rand's result that participants who experienced mutual cooperation in part one behave more pro-socially and punish selfish behavior more later can be due to any of the following reasons: *i*) people regard their own behavior as fair and expect others to do the same; *ii*) people consider a distribution of earnings similar to the past as fair; and *iii*) people infer social norms from observing the interactions of others. Our experiment does not have an explicitly dynamic structure and, thus, rules out the first two possibilities. Hence, our experiment allows us to test the pure information effect of knowing the frequency of a behavior in the population. *Second*, we use the strategy method to elicit the complete punishment strategy for each individual. In contrast, the aforementioned studies all focus on examining the *aggregate effect* across individuals. As a result, our experiment reveals the co-existence of four systematic punishment types of individuals that previous experiments cannot reveal. Our results suggest that unpacking the aggregate effect and revealing the behavior for each individual indeed matters. If we had focused on the aggregate effect only, it would lead us to conclude that policies of manipulating the information regarding the proportion of cooperators results in no changes in people's fairness perceptions and punishment of free riders. However, thinking in this way is a mistake. We show that having more cooperators crowds out the incentives to punish free riders for some individuals, while it increases the willingness to punish for some others. Studies from economics (Currarini et al., 2009) and sociology (McPherson et al., 1992, 2001) both show that *homophily* is a basic pattern of social networks in the real world: individuals with similar preferences and attitudes tend to be friends with each other or live in the same neighborhoods. Co-existence of different punishment types then suggests that manipulating the information regarding the proportion of cooperators would have significant and dramatically different effects on different

social groups—a new and striking conclusion that we would not have obtained if we had focused on the aggregate effect only.

The paper is organized as follows. Section 4.2 reviews theories and derives hypotheses. Section 4.3 presents experimental design and procedures. Section 4.4 presents the results. We conclude by discussing the implications of our finding in Section 4.5.

## 4.2 Theory

We review existing theories and derive hypotheses in this section. We consider the following stylized game. Two players, Ann and Bob, who do not know each other but encounter on a street. They face decisions as the ones in Figure 4.1. Ann wants to get rid of a banana skin but fails to find a trash bin on the street. Ann can take the banana skin all the way back home (*cooperate*), or throw the banana skin on the street (*free ride*). Bob is an observer who happens to pass by. If Ann throws the banana skin on the street, Bob can either walk away as if nothing special occurs (*not punish*), or, Bob stops Ann, gives Ann a warming, and asks Ann to pick up the banana skin (*punish*). The numbers at the terminal nodes of the game tree in Figure 4.1 are the *material payoffs*.



**Figure 4.1:** The stylized game. The numbers at the terminal nodes represent material payoffs.

If Bob is a material-payoff maximizer, Bob will never punish. This leads to the following null hypothesis.

116

**Null Hypothesis** (Material payoffs maximization). *When conducting punishment is costly, people never punish free riders, regardless of the proportion of cooperators in a society.*

However, a large body of public goods game experiments show that many participants punish free riders at the expense of their own earnings (e.g., Fehr and Schmidt, 2000; Herrmann et al., 2008). Why punish? And why the tendency to punish varies across societies as observed by Herrmann et al. (2008)? In what follows, we review existing theories that economists have developed in the past fifty years related to punishment of free riders. We compare three categories of theories: *i*) the conventional view among economists up to the 1980s, *ii*) the theories of inequality aversion (Fehr and Schmidt, 1999) and intention-based reciprocity (Dufwenberg and Kirchsteiger, 2004; Rabin, 1993), and *iii*) theories of norm-based motives (Bicchieri, 2006; Sugden, 1984, 2000). We show that each category generates a different prediction regarding how punishment of free riders may change in response to a change in the public goods contribution level in the society.

### 4.2.1   Diminishing marginal return to punishment

We introduce the conventional view among economists up the 1980s regarding private provision of public goods in this subsection and derive the corresponding prediction. A fully specified model would be a combination of Bergstrom et al. (1986) and Kandori (1992). For the purpose of this paper, we only sketch the idea.

Economists have long recognized that the assumption of private consumption maximization cannot explain the observed level of public goods contributions in many sectors (Becker, 1974; Olson, 1965). According to Bergstrom et al. (1986), the dominant view among economists on public goods contribution up to the 1980s is:

> "...the case where people are concerned only about their private consumptions and the *total* supply of public goods... is the model which

has received the most attention so far in the literature and is, we sus-
pect, the one on which many economists base their intuitions. " (p.
26; italic original)

This conventional view of public good contributions assumes that individuals have
intrinsic concerns for the total supply of public goods. It also assumes that
the marginal utility from public goods is diminishing. In addition, we assume
that, following the spirit of Kandori (1992), punishment of free riders generates
a community-level punishment threat to deter free riding for future interactions
and, thereby, benefits every member of society.

More precisely, let $q$, with $0 \le q \le 1$, be the proportion of cooperators in a
society. It determines the total supply of public goods in the society. Let $y \in \{0, 1\}$
denote Bob's choice of punishment: $y = 1$ indicates that Bob punishes, and $y = 0$
for not punishing. Let $Y$ denote the number of individuals in the society who
punish free riders other than Bob. How many individuals cooperate in the society
depends on how many people punish free riders. That is, $q$ is increasing in $Y + y$.
Bob cares about the proportion of cooperators in a society, $q$, not only because it
affects his own earnings, but also because it is a public good that benefits others.
Let $V(q)$ denote the part of Bob's utility that is derived from $q$, with $V' > 0$ and
$V'' < 0$. This term $V(q)$ summarizes Bob's individual benefits from $q$ as well as
his intrinsic concern for the total supply of public goods. When deciding whether
to punish Ann, the immediate material consequence for Bob is $4y + 9(1 - y)$. Bob's
total utility includes the immediate consequence and the term $V(q)$:

$$u = 4y + 9(1 - y) + V(q).$$

The proportion of cooperators, $q$, is an increasing function of $Y + y$. Hence we
have

$$u_{yq} = V''q' < 0.$$

That is, a higher proportion of cooperators *crowds out* Bob's incentive to punish

free riders.

**Prediction 1** (Diminishing marginal return to punishment). *The higher the proportion of cooperators in a society, the* less *people punish free riders.*

### 4.2.2 Inequality aversion and intention-based reciprocity

The theories of norm-based motives (Sugden, 1986) are introduced before the introduction of inequality aversion (Fehr and Schmidt, 1999) and intention-based reciprocity (Dufwenberg and Kirchsteiger, 2004; Rabin, 1993) into economics. Nevertheless, for pedagogical reasons, we first discuss the theories of inequality aversion and intention-based reciprocity.

**Inequality aversion** (Fehr and Schmidt, 1999). Suppose Ann free rides. Bob has *inequality aversion* such that each unit of difference in material earnings between Ann and Bob incurs $\alpha \geq 0$ units of loss of *utility* to Bob. Hence, if Bob does not punish, he obtains $9 - 16\alpha$ units of utility. In contrast, if Bob punishes, he obtains $4 - 6\alpha$. Bob finds it better off to punish to reduce earning difference between Ann and him, if

$$\alpha > {}^{1}\!/\!{}_{2}.$$



**Figure 4.2:** The game with reciprocity preferences. $\theta \geq 0$ is the reciprocity parameter of Bob. The payoffs at the terminal nodes of the game tree for Bob show the *utility* of Bob, given Bob's belief that Ann knows that Bob chooses *punish*. The numbers for Ann show the material payoffs for Ann.

**Intention-based reciprocity** (Dufwenberg and Kirchsteiger, 2004; Rabin,

1993). The following argument is an application of Dufwenberg and Kirchsteiger (2004).[3] Suppose both Ann and Bob have the preferences to respond to unkind with unkind. Let us verify that, if Bob's preference for retaliate is strong enough, punishment is the best-response of Bob.[4] Bob expects that Ann correctly expects that Bob will punish. Upon observing Ann free riding, Bob thinks that Ann could have cooperated and given 18 points to Bob, but instead Ann free rides, resulting in 4 points to Bob. Clearly, Ann is unkind to Bob. If Bob has sufficiently strong preferences for retaliate, he would punish. Figure 4.2 shows the utility for Bob for taking each action, given the specified beliefs of Bob, and $\theta \geq 0$ is the reciprocity parameter of Bob.[5] Bob prefers punishing if

$$\theta > {}^1\!/{}_{21}.$$

There is a common feature of inequality aversion and reciprocity: Bob's punishment of Ann only depends on the behavior of Ann, not on the behavior of any payoff-irrelevant third-parities. More explicitly, suppose that the game is played out in our experiment, and Ann is the randomly assigned partner of Bob. Based on inequality aversion and intention-based reciprocity, Bob's punishment decision should not be affected by the proportion of cooperators among other participants who are *not* matched with Bob.

**Prediction 2** (Intention-based reciprocity and inequality aversion): *Individuals punish free riders, and their punishment is* independent *of the proportion of co-operators in a society.*

---

[3]Rabin's (1993) model fails to eliminate some unintuitive equilibria in sequential games.

[4]More precisely, if Ann is a material-payoff maximizer, Bob has sufficiently strong preferences to retaliate, and their preferences are commonly known, then *(cooperate, punish)* is a *sequential reciprocity equilibrium* (SRE) defined by Dufwenberg and Kirchsteiger (2004). If Ann also has the preferences to retaliate, then (*free ride, punish*) is a SRE.

[5]The numbers are calculated based on the model of Dufwenberg and Kirchsteiger (2004). The perceived kindness of Ann to Bob is $-(18-4)/2 = -7$. The kindness of Bob to Ann by taking *punish* is $-(25-10)/2 = -{}^{15}\!/{}_2$, whereas the kindness of Bob to Ann is ${}^{15}\!/{}_2$.

### 4.2.3 Norm-based reciprocity

According to Sugden (1984; 1986; 2000) and Bicchieri (2006), Bob may punish because Ann deviates from Bob's expectation on the *norm (majority) behavior* of the population. However, neither Sugden nor Bicchieri provides an explicit formulation of the punishment mechanism. We provide one specification as follows.

Bob has his expectation, $q_i \in [0, 1]$, regarding the proportion of individuals who would cooperate if they are in the position of Ann. Bob does not know Ann in person and, in Bob's eyes, Ann is not so different from any other unrelated individuals in the society. Hence, Bob expects his earnings to be *at least*

$$18q_i + 4(1 - q_i) = 4 + 14q_i. \tag{4.1}$$

However, Ann free rides, in which case Bob can *at most* earn 9. If $q_i > 5/14$, Bob is certainly worse off than what he expected. Bob is frustrated about this. In order to release his frustration, Bob may punish Ann. Bob has the same impulse to act unkind to unkind as in the reciprocity model of Dufwenberg and Kirchsteiger (2004). The numbers at the terminal nodes of the game tree in Figure 4.3 represent the preferences of Bob. $\lambda \geq 0$ measures Bob's frustration. The *norm-based unkindness* of Ann towards Bob is $14q_i - 5$, which is obtained by subtracting the most earnings that Bob can get if Ann free rides, 9, from the least earnings expected by Bob given belief $q_i$, i.e., equation (4.1). It follows that Bob prefers punishing if

$$\lambda(14q_i - 5) > 1/3. \tag{4.2}$$

**Figure 4.3:** The game with norm-enforcement preferences. $\lambda \geq 0$ is the norm-enforcement parameter of Bob. The payoffs at the terminal nodes of the game tree for Bob show the *utility* of Bob, conditional on $q_i$ (Bob's belief on the fraction of individuals in the society who would *cooperate*). The numbers for Ann show the material payoffs for Ann.

The distinct feature of the above norm-based reasoning is that Bob's punishment decision is conditional on $q_i$, his expectation on the proportion of cooperators in the society. If Bob believes that no one is going to cooperate in the society, then condition (4.2) is never satisfied, and thus Bob does not punish Ann. The intuition is that, when $q_i = 0$, Ann's free riding is completely expected by Bob. Thus, Bob is not frustrated. In contrast, the greater the expected proportion of cooperators, the greater is the frustration of Bob upon observing Ann free riding. Hence, norm-based reciprocity predicts the following.

**Prediction 3** (Norm-based reciprocity): *The higher the proportion of cooperators in a society, the more people punish free riders.*

## 4.3 Experimental design and procedures

### 4.3.1 Design

**The experimental game.** We randomly match participants in pairs to play the following two-stage game. In the first stage, the two players simultaneously choose whether to *cooperate* or *free ride*. Throughout the experiment, we use the word "share" to indicate cooperation and "keep" to indicate free riding in order to improve participants' understanding and involvement. Table 4.1 below shows

the earnings from the first stage. If both players cooperate, each player gets 18 points. If both free ride, each player gets 16 points. Hence, the outcome of both players cooperating Pareto-dominates the outcome of both free riding.[6] If one player cooperates and the other free rides, the free rider obtains more points than the cooperator.

**Table 4.1:** The first stage of the game

|  | *cooperate* | *free ride* |
|---|---|---|
| *cooperate* | 18, 18 | 9, 25 |
| *free ride* | 25, 9 | 16, 16 |

If both players cooperate in the first stage, then the game ends and each player obtains 18 points. If at least one player free rides, the game proceeds to the second stage. Specifically, if a player, $i$, free rides in the first stage, then in the second stage her partner, $j$, can assign up to 10 *deduction points* to $i$. Player $j$ may cooperate, or free ride. The *deduction ratio* is 1-to-3, i.e., each deduction point that $j$ assigns to $i$ costs one point to $j$ but reduces $i$'s payoff by three points. Notice that a player can assign deduction points to the other player only when the other player free rides.

More precisely, let $x_i \in \{1, 0\}$ denote the action of player $i$ taken in the first stage: $x_i = 1$ if $i$ cooperates, and $x_i = 0$ if $i$ free rides. Let $y_i \in \{0, 1, \dots, 10\}$ denote the deduction points that $i$ assigns to the matched partner $j$, conditional on $j$ choosing to free ride. Likewise, let $x_j$ and $y_j$ be the corresponding actions taken by $j$. Moreover, let $\pi_i^1$ be the material earnings for $i$ from the first stage as shown in table 4.1. Let $\mathbb{1}\{x_j = 0\}$ be the indicator function taking the value of one if $x_j = 0$, i.e., $j$ free rides, and zero otherwise. We analogously define $\mathbb{1}\{x_i = 0\}$. Then the total material earnings for $i$ from both stages are

$$\pi_i = \pi_i^1 - \mathbb{1}\{x_j = 0\} \times y_i - \mathbb{1}\{x_i = 0\} \times 3 \times y_j.$$

---

[6]The payoffs when both players cooperate are not so much different from the payoffs when both players free ride. The reason is that, in order to analyze punishment behavior towards defectors, we need a non-negligible sample of free rider. After running some pilot sessions, we found that if the payoffs from both players cooperating are much better than both free riding, then most participants would choose to cooperate.

**The strategy method.** Before each participant decides how many deduction points to assign, we inform them about the proportion of cooperators among the participants they are *not* matched with. This information should change participants' belief regarding the proportion of cooperators in the population. Our experiment uses the strategy method to elicit, for *each* participant, whether her deduction points' decision depends on her belief regarding the proportion of cooperators in the population. The implementation is as follows. We randomly draw $n$ participants from *previous* sessions, and tell participants at the current session how many among the $n$ randomly selected previous participants chose to cooperate. Instead of providing a single piece of such information, we present to each participant multiple possible situations that might occur. They need to indicate the deduction points to assign in *each* of the situations presented to them. After they provide their answer, we reveal to them about the situation that actually occurred. The actually occurred situation then determines the earnings of them and their partners'. Thus, the participants are incentivized to provide their decision in each possible situation presented to them.

To be more concrete, in three experimental sessions, we have $n = 4$, the *4-peer condition*. We ask each participant to indicate their deduction points' decision in response to each of the following five situations:

**Table 4.2:** The five situations in the 4-peer condition.

| | |
|---|---|
| Situation 1 | None of the four previous participants chose to cooperate. |
| Situation 2 | One of the four previous participants chose to cooperate. |
| Situation 3 | Two of the four previous participants chose to cooperate. |
| Situation 4 | Three of the four previous participants chose to cooperate. |
| Situation 5 | All of the four previous participants chose to cooperate. |

After participants provide their decisions in all five situations, we inform the participants about how many among the four randomly selected previous participants have actually cooperated. The deduction points' decision of each participant in the actually occurred situation is implemented to calculate the earnings of the participant and his partner. All the details of the elicitation procedure are thor-

oughly explained to participants. We also include several control questions to make sure that participants understand the rules.

In another two sessions, we have $n = 50$, the *50-peer condition*. We present to each participant the following five situations:

**Table 4.3:** The five situations in the 50-peer condition.

| | |
|---|---|
| Situation 1 | Between 0 to 10 of the fifty previous participants chose to cooperate. |
| Situation 2 | Between 10 to 20 of the fifty previous participants chose to cooperate. |
| Situation 3 | Between 20 to 30 of the fifty previous participants chose to cooperate. |
| Situation 4 | Between 30 to 40 of the fifty previous participants chose to cooperate. |
| Situation 5 | Between 40 to 50 of the fifty previous participants chose to cooperate. |

We now discuss several concerns about our design. First, we agree that, strictly speaking, the provided information in the 4-peer condition is not very statistically informative. However, the provided information may generate additional priming effects. Thus, we expect participants' belief to effectively change in the 50-peer condition as well as in the 4-peer condition. Second, regarding the validity of the strategy method, Brandts and Charness (2011) conduct a meta-analysis to compare the strategy method and the direct-response method. In all twenty-nine experiments documented in their study, if a treatment effect is found when using the strategy method, it is also observed when using the direct-response method. Third, one may be concerned with potential experimenter demand effects in our experiment, namely, participants take particular actions only because they think that is what the experimenter expects them to do (Zizzo, 2010). However, it is not obvious from reading our instructions what the experimenter's expectation is. On the one hand, we present to participants with multiple situations. Participants might think that this suggests that the experimenter expects them to condition decision on the proportion of cooperators. On the other hand, however, we highlight in various places in the instruction that the choices of participants from previous sessions are completely payoff-irrelevant to them. We also add control questions to test and reinforce participants' understanding of this fact. Hence, it is not obvious whether the experimenter expects the participants to condition

decision on the proportion of cooperators or not. Moreover, there is no way to guess how the experiment expects a participant to exactly condition decision on the proportion of cooperators (increasingly or decreasingly?). The bottom line is that, unless a participant has a strong prior of her own about what constitutes appropriate behavior in similar contexts in real life, it is hard to imagine that she would have a strong expectation about what the experimenter expects her to do.

### 4.3.2  Procedures

We used the online platform Amazon Mechanical Turk to recruit and pay participants. We conducted our experiment on LIONESS Lab (Arechar et al., 2018) during February 2017. LIONESS Lab is a web-based platform for online interactive experiments.[7] Each participant only played the game once in our experiment. To reduce drop-out rate—which might particularly be an issue for online experiments, we conducted the matching once every two participants submitted all decisions. In total we had 246 individuals registered our experiment, while 203 went through to the end (there was a participant who we could not find another one to match with). Participants' age varied between twenty to seventy years old. Average age was 34.7. 51% of our participants are male. Among the 203 participants, 91 participated in the three sessions of the 4-peer condition, and 112 participated in the two sessions of the 50-peer condition. Our experiment lasted for about ten to twenty minutes. Earnings were paid out in US dollar. Each US dollar corresponded to 20 experimental points in the experiment. Earnings ranged from \$1 to \$2.5. Average earning was \$1.95.

---

[7]Similar to our experiment, Arechar et al. (2018) conduct a public goods game experiment on LIONESS Lab with participants recruited from Amazon Mechanical Turk. They show that all main patterns of cooperation and punishment observed in previous laboratory experiments are replicable in the online experiment.

## 4.4 Results

We present the results of our experiment in this section. The focus is on participants' punishment behavior.

### 4.4.1 Main results

**Result 1—No significant effect at the aggregate level.** *On average, the proportion of cooperators among previous participants does* not *significantly affect punishment of free riders.*

Using the strategy method, we obtain each participant's deduction point decision for each of the five situations regarding the proportion of cooperators (Table 4.2 for participants in the 4-peer condition and Table 4.3 for participants in the 50-peer condition). Figure 4.4 presents the overall effects. The figure shows that, on average, a higher proportion of cooperators among previous participants does not significantly increase deduction points assigned to free riders no matter whether in the 4-peer condition or in the 50-peer condition. This result is confirmed by the regressions in Table 4.5. In the regressions in Table 4.5, the dependent variable is the deduction points that participants assign to their free riding partners, ranging from 0 to 10. The independent variable *Proportion of cooperators* indicates the percentage of cooperators among previous participants, taking values of 0, 0.25, 0.5, 0.75, and 1. The coefficient for *Proportion of cooperators* is, although positive, not statistically different from zero in all five regressions at any conventional level for all the six specifications presented. Moreover, to test whether the participants who cooperate differ from those who free ride in terms of their sensitivity to the *Proportion of cooperators*, regression (iii) includes an interaction term between *Proportion of cooperators* and *Cooperate*, where the dummy *Cooperate* indicates whether the participant cooperates or not. We do not see a significant difference.

**Result 2—Coexistence of different punishment types.** *Participants are heterogeneous in their punishment behavior. Most of them can be classified as one*

127

**Figure 4.4:** Average deduction points assigned to free riding partners. Provided information about the proportion of cooperators does not significantly affect deduction points. The spikes indicate the standard errors of means (SEM).

**Figure 4.5:** OLS regressions. Dependent variable—Deduction points assigned to free riding partner (0 to 10).

|  | (i) | (ii) | (iii) | (iv) | (v) |
|---|---|---|---|---|---|
| Proportion of cooperators | 0.219 | 0.219 | 0.229 | 0.219 | 0.132 |
|  | (0.174) | (0.175) | (0.226) | (0.175) | (0.229) |
| Cooperate |  | 2.981** | 2.988** | 2.983** | 2.983** |
|  |  | (0.319) | (0.319) | (0.323) | (0.323) |
| Proportion of cooperators × Cooperate |  |  | -0.014 |  |  |
|  |  |  | (0.324) |  |  |
| 50-peer |  |  |  | -0.029 | -0.107 |
|  |  |  |  | (0.421) | (0.433) |
| Proportion of cooperators × 50-peer |  |  |  |  | 0.157 |
|  |  |  |  |  | (0.344) |
| Constant | 2.455** | 0.897** | 0.892 | 0.922 | 0.966 |
|  | (0.847) | (0.795) | (0.788) | (0.886) | (0.893) |
| $N$ | 1015 | 1015 | 1015 | 1015 | 1015 |

Notes: Within parentheses are robust standard errors clustered by participants. $^{+}p < 0.1$, $^{*}p < 0.05$ and $^{**}p < 0.005$ for two-tailed tests. *Proportion of cooperators* takes values of $0, 0.25, 0.5, 0.75$ and $1$, corresponding to the situations in Table 4.2 and 4.3. *Cooperate* is a dummy indicating whether the participant cooperates him- or herself. *50-peer* is a dummy indicating the 50-peer condition. All regressions include age and gender as controls.

*of the four punishment types shown in Figure 4.6.*



**Figure 4.6:** Deduction points assigned to free riding partners given the provided information about the proportion of cooperators. The spikes indicate standard errors of means (SEM). Each line displays a different punishment type. The four types add up to 86% of the sample. The remaining 14% (not shown) do not fall into any of the four types displayed.

That we do not find a significant effect at the aggregate level is not because no participant responds to the provided information about the proportion of cooperators. On the contrary, about a fourth of our participants (23%) do condition their punishment decision on how many cooperators are around. However, there is remarkable heterogeneity among participants. Figure 4.6 presents the pooled data for the two conditions and shows that most of the participants are classified as one of the following four distinct *punishment types*:

1. *Never punish.* 85 among all 203 participants assign zero deduction point to their free riding partners for all five situations, constituting 42% of the whole sample. These participants behave in line with the Null Hypothesis of maximizing material payoffs of themselves. On the other hand, 58% of our participants do punish free riders by assigning at least 1 deduction point in at least one situation.

2. *Punish independently* of the proportion of cooperators. 21% of participants assign one or more deduction points to their free riding partners and they

make identical decisions in all five situations. Hence, these participants' punishment decisions are independent of the proportion of cooperators among previous participants. Their behavior is in line with intention-based reciprocity and inequality aversion.

3. *Punish increasingly* in the proportion of cooperators. 13% of participants do not decrease deduction points to assign as the proportion of cooperators among previous participants goes up, and they strictly *increase* deduction points to assign for at least one instance. The behavior of these participants is consistent with norm-based reciprocity.

4. *Punish decreasingly* in the proportion of cooperators. 10% of participants do not increase their deduction points assigned as the proportion of cooperators among previous participants goes up, and they strictly *decrease* deduction points to assign for at least one instance. The behavior of these participants can be explained by the reasoning underpinning the conventional view of public goods contributions. That is, their punishment is motivated by a desire to generate a community-level threat to deter free riding and they are intrinsically concerned with the total supply of public goods.

The four punishment types add up to 86% of the whole sample, while the behavior of the rest of participants is not explained by previously specified theory. To test whether the four punishment types of participants arise for systematic reasons rather than from pure randomness, we perform a simulation exercise as follows. We generate a simulated data of one-million individuals. Each of them randomly chooses deduction points between 0 to 10 to assign in each of the five popularity levels of cooperation. That is, they each time roll an eleven-sided die to give a decision. The simulation generates the following result: among the one-million individuals, 0.00% never punish, 0.01% punish independently of the proportion of cooperators, 1.83% punish increasingly in the proportion of cooperators, and 1.88% punish decreasingly in the proportion of cooperators. The remaining 96.28%

do not fall into any of the four types, in contrast to the 14% we observe in the experiment. To conclude, the participants of *each* of the four punishment types in our experiment are significantly more than they would arise from pure randomness.

### 4.4.2 Further results

**Result 3—Punishers are cooperators.**

69% of our participants (140/203) are cooperators, i.e., they cooperate in the first stage of the game, while 31% free ride.[8] Figure 4.7 shows the distributions of punishment types among cooperators and among free riders, respectively. Among free riders, 73% never assign a positive deduction point to their free riding partners. In contrast, 72% of cooperators assign 1 or more deduction points to their free riding partners in at least one situation. Statistically we reject the hypothesis that cooperators and free riders have the same probability of punishing free riding partners (Fisher's exact $p < 0.001$). Regressions in Table 4.5 provide further estimates of the difference between cooperators and free riders. The dummy *Cooperate* indicates whether the participant cooperates or not. The results show that cooperators assign significantly more deduction points to their free riding partners than free riders would do, and the difference is sizable—about 3 deduction points on average.

---

[8]The 203 participants are those who registered and completed the experiment. In addition, there were 10 participants who only completed the first stage of the game. Among all 213 participants who we have data on their stage-one decision, 147 (69%) choose to cooperate.

**Figure 4.7:** Percentage of each punishment type among cooperators and among free riders.

**Result 4.** *Assigning deduction points to equalize earnings is the mode behavior of the cooperators who punish independently.*

Figure 4.8 presents the percentage distributions of deduction points among *cooperators*; we separate the distributions according to their different punishment behavior. In our experiment, if a cooperator meets a free rider, and the cooperator aims at equalizing the final earnings between her and the free rider, then the cooperator would assign 8 deduction points. Hence, assigning exactly 8 deduction points indicates that a cooperator is motivated by inequality aversion (Fehr and Schmidt, 1999). Interestingly, we find that assigning 8 deduction points is the most popular choice among those, and only among those, who punish free riders independently of how many other cooperators are around.

**Figure 4.8:** Distribution of deductions points among cooperators, separated by punishment types. Assigning 8 deduction points indicates that a cooperator aims at equalizing the earnings between her and her free riding partner.

**Result 5.** *More participants condition their punishment of free riders on the proportion of cooperators in the 50-peer condition than in the 4-peer condition; however, the difference is not statistically significant.*

Figure 4.9 displays the percentage of each punishment type for each of the two conditions. The participants who punish free riders increasingly in the proportion of cooperators increases from 11% in the 4-peer condition to 15% in the 50-peer condition. Meanwhile, the participants who punish decreasingly in the proportion of cooperators increases from 8% in the 4-peer condition to 12% in the 50-peer condition. Nevertheless, we cannot reject the hypothesis that the sample distributions for the 4-peer condition and for the 50-peer condition are drawn from the same population distribution (Fisher's exact $p = 0.397$). Regression (v) in Table 4.5 also shows that the coefficient for the interaction term between Proportion of cooperator and the dummy for 50-peer condition is not significantly different from zero.

**Figure 4.9:** Percentage of each punishment type for the 4-peer condition and the 50-peer condition, respectively.

## 4.5 Conclusion

We ask whether manipulating the information regarding the frequency of a behavior affects people's fairness perception about the behavior. To answer the question, we design an experiment to test whether informing people of a higher proportion of cooperators affects the legitimacy of free riding and the punishment of free riders. Our results indicate that there is no simple answer to the question. On the one hand, we don't observe that the information about the proportion of cooperators significantly affects punishment of free riders *on average*. This result speaks to previous studies (Cooper and Dutcher, 2011; Falk et al., 2006; Herz and Taubinsky, 2017; Peysakhovich and Rand, 2016) about how fairness perceptions are formed and could be changed. Previous studies show that, on average, people punish selfish behavior of others more after experiencing more cooperative interactions. However, previous studies admit several interpretations. In contrast, our experiment focuses on the pure information effect of a change in the cooperation level in a society. Comparing our result with previous ones, we conclude that a

large part of the average change in fairness perceptions previously observed is due to other channels than the pure information effect. The other channels include that people regard their own behavior as fair and expect others to do the same, and that people consider a distribution of earnings similar to the past as fair.

On the other hand, however, it would be a mistake to think that our result implies that a policy of changing the information about the cooperation level of a society would not affect punishment of free riders and the cooperation level in the field. Using the strategy method, our experiment reveals the remarkable heterogeneity among individuals that previous studies cannot reveal. We discover co-existence of four distinct punishment types: $i$) individuals who never punish, $ii$) those who punish but independently of the proportion of cooperators, $iii$) those who punish increasingly in the proportion of cooperators, and $iv$) those who punish decreasingly in the proportion of cooperators. What does the heterogeneity imply in the field? Studies from economics (Currarini et al., 2009) and sociology (McPherson et al., 1992, 2001) both indicate that homophily features many social networks: individuals with similar preferences and attitudes tend to be friends with each other or live in the same neighborhoods. Co-existence of different punishment types then implies that changing the information regarding the proportion of cooperators would lead to significant and dramatically different impacts on different social groups. Hence, our finding suggests that implementing an informational policy could be fruitful if, and only if, it is combined with policies to identify the punishment type of individuals in the targeted groups.

## 4.6   Appendix: experimental instructions

This appendix provides the complete screen shots of our experiment.

**The 4-peers treatment**

### Welcome

The setup of this HIT is different from HITs that you might be used to completing via MTurk. You will be participating in an interactive task with another MTurker.

As you are completing this task at the same time, it is important that you complete this HIT *without interruptions*.

Including the time for reading these instructions, the HIT will take about 8 minutes to complete. During the HIT, please do not close this window or leave the HIT's web pages in any other way.
If you do close your browser or leave the HIT, you will NOT be able to re-enter the HIT and we will NOT be able to pay you!

During this HIT you can earn **Points**. You will receive **25 Points to start with**. Depending on your choices and the choices of the other MTurker in the task, you may earn additional Points.

At the end of the task, your Points will be converted into real money according to the exchange rate:

**20 Points = 1 Dollar.**

You will receive a code to collect your payment via MTurk upon completion.

| Continue |

## Instructions

You and **another MTurker** will form a pair and participate in this task at the same time.

We will refer to the MTurker in your pair as **your partner**.

You and your partner have received these same instructions. You two will play a game.

The game has two stages: **Stage 1** and **Stage 2**.

Next

## Stage 1

In Stage 1, you and your partner at the same time choose between *Keep* and *Share*.

| If you choose: | and your partner chooses: | then you get: | and your partner gets: |
|---|---|---|---|
| *Keep* | *Keep* | 16 | 16 |
| *Share* | *Share* | 18 | 18 |
| *Keep* | *Share* | 25 | 9 |
| *Share* | *Keep* | 9 | 25 |

After you and your partner have each made your choice, you will proceed to Stage 2. ♂

Your earnings from Stage 1 will carry over to Stage 2.

Next

**Stage 2**

Once you and your partner have each submitted your choice for Stage 1, you proceed to Stage 2.

If your partner chose *Keep* in Stage 1, you can choose to assign **up to 10 Deduction Points** to your partner.



For each Deduction Point you assign, **1 Point** will be deducted from **your earnings**, and **3 Points** will be deducted from **your partner's earnings**.



Similarly, if you chose *Keep* in Stage 1, your partner can choose to assign up to 10 Deduction Points to you. For each Deduction Point your partner chooses to assign, 1 Point will be deducted from your partner's earnings and 3 Points will be deducted from your earnings.



| Previous | Next |
|----------|------|

---

**Stage 2**

If your partner chose *Share* in Stage 1, you **cannot** assign any Deduction Points to your partner.



Similarly, if you chose *Share* in Stage 1, your partner **cannot** assign Deduction Points to you.



| Previous | Next |
|----------|------|

**Here is an example:**

In Stage 1, you choose *Share* and your partner chooses *Keep*. At the end of Stage 1, **you** have **9 Points** and **your partner** has **25 Points**. In Stage 2, you assign **8 Deduction Points** to your partner.

Then at the end of the game...

Your earnings are **1 Point** (9 points from Stage 1 *minus* 8 points deducted from Stage 2).
Your partner's earnings are **1 Point** (25 points from Stage 1 *minus* 3 x 8 points deducted from Stage 2).



9 − 8 = 1 Point          25 − 3*8 = 1 Point

| Previous | Next |
|---|---|

---

In summary, in Stage 1, you and your partner choose between *Share* and *Keep*.

| If you choose: | and your partner chooses: | then you get: | and your partner gets: |
|---|---|---|---|
| *Keep* | *Keep* | 16 | 16 |
| *Share* | *Share* | 18 | 18 |
| *Keep* | *Share* | 25 | 9 |
| *Share* | *Keep* | 9 | 25 |

If your partner chooses *Keep*, you can assign up to 10 **Deduction Points** to your partner in Stage 2.

For each Deduction Point you assign, 1 Point will be deducted from your earnings, and 3 Points will be deducted from your partner's earnings.

If your partner chooses *Share* in Stage 1, then you CANNOT assign Deduction Point to your partner in Stage 2.

Your partner faces the same choices as you do.

**Control questions**

Before you make your choices, please answer the following questions to make sure you completely understand the game.

**1. Suppose that in Stage 1, you choose *Share* and your partner chooses *Keep*.**

(a) How many Points would you earn from Stage 1?

(b) How many Points would your partner earn from Stage 1?

**2. Subsequently in Stage 2, you assign 5 Deduction Points to your partner.**

(a) How many Points would you have at the end of the game?

(b) How many Points would your partner have at the end of the game?

Next

---

**Your choice for Stage 1**

Remember:

| If you choose: | and your partner chooses: | then you get: | and your partner gets: |
|---|---|---|---|
| *Keep* | *Keep* | 16 | 16 |
| *Share* | *Share* | 18 | 18 |
| *Keep* | *Share* | 25 | 9 |
| *Share* | *Keep* | 9 | 25 |

Now please make your choice:

Share

Keep

After you made your choice, press "Submit" to proceed to Stage 2.

Once you press "Submit", you **cannot** go back and change your choice for Stage 1.

Submit

## Stage 2

Your choice for Stage 1 has been recorded. Stage 2 begins now.

Before we inform you about the actual choice made by your partner, imagine that

**your partner has chosen *Keep*.**

Hence:

You can assign up to 10 Deduction Points to your partner.

For each Deduction Point you assign, 1 Point will be deducted from your earnings, and 3 Points will be deducted from your partner's earnings.

If your partner indeed chose *Keep*, your choice regarding the Deduction Points will be used to calculate the the earnings of you and your partner.

Next

---

## Stage 2

In Stage 2 you will decide to assign Deduction Points if in Stage 1 your partner has chosen Keep.

To give you an idea of what people usually choose in Stage 1, we will inform you of **how many of four randomly selected MTurkers (NOT your partner)** chose *Share* in the HITs prior to this one.

Depending on the the four previous MTurkers' choices, there are five possible situations:

**Situation 1: None** of them chose *Share.*

**Situation 2: One** of them chose *Share.*

**Situation 3: Two** of them chose *Share.*

**Situation 4: Three** of them chose *Share.*

**Situation 5: All** of them chose *Share.*

For each of these situations you will indicate how many Deduction Points to assign to your partner.

After you complete your choices for all five situations, we will let you know which situation actually occurred.

Your choice in that actual situation will be used to calculate the earnings of you and your partner.

**Note:**

*Your partner is NOT one of these previous MTurkers.*
*You CANNOT assign Deduction Points to any of these previous MTurkers.*
*Your choice on the Deduction Points only affects the earnings of you and your partner.*

Previous

Next

**Control Questions**

Before you make your choices, please answer the following questions to make sure you understand your task.

Suppose that in Stage 1, your partner chose *Keep*. The situation is that in Stage 1, all four previous MTurkers (NOT your partner) all chose *Keep*. In Stage 2 you choose to assign 3 Deduction Points to your partner.

1. How many Points will be deducted from your earnings?

2. How many Points will be deducted from the previous MTurkers' earnings?

3. How many Points will be deducted from your partner's earnings?

| Previous | Next |

# Your choice for Stage 2

Situation 1

**None** of the four previous MTurkers (NOT your partner) chose *Share*.



Enter the Deduction Points you assign to your partner if your partner chose *Keep* in this situation:



Next

# Your choice for Stage 2

Situation 2

**One** of the four previous MTurkers (NOT your partner) chose *Share*.

Enter the Deduction Points you assign to your partner if your partner chose *Keep* in this situation:

Previous          Next

---

# Your choice for Stage 2

Situation 3

**Two** of the four previous MTurkers (NOT your partner) chose *Share*.

Enter the Deduction Points you assign to your partner if your partner chose *Keep* in this situation:

Previous          Next

# Your choice for Stage 2

Situation 4

**Three** of the four previous MTurkers (NOT your partner) chose *Share*.

Enter the Deduction Points you assign to your partner if your partner chose *Keep* in this situation:

| Previous | Next |
|----------|------|

---

# Your choice for Stage 2

Situation 5

**All** of the four previous MTurkers (NOT your partner) chose *Share*.

Enter the Deduction Points you assign to your partner if your partner chose *Keep* in this situation:

To go back and change your decisions for previous situations, press "Previous".

To submit your final decisions, press "Submit".

| Previous | Submit |
|----------|--------|

**Questionnaire**

You have now completed the decision making part of this HIT. Please fill out this brief questionnaire.
Once your partner is ready, we will calculate you and your partner's earnings and display them on your screen.

You will then receive a code to collect your earnings on MTurk.

What is your age? [          ]

What is your gender?

| Male |
|------|
| Female |

Could you briefly describe your reasoning for choosing to either *Share* or *Keep*?

[                                                                    ]

remaining characters [ 200 ]

Could you briefly describe the reasoning you used when allocating **Deduction Points**? In particular, why did or didn't you make your choices dependent on **how many previous MTurkers chose *Share*?**

[                                                                    ]

remaining characters [ 400 ]

| Continue |
|----------|

---

# Earnings

You and your partner have now finished making your choices. Find the results below.

You started with **25 Points**.

In **Stage 1**, you chose **Keep**.
Your partner chose **Share**.

This means that your earnings from Stage 1 are: **25 Points**.
Your partner's earnings from Stage 1 are 9 Points.

Among the four MTurkers in the previous session, **1** chose *Share*.

As your partner chose *Share*, any Deduction Points you assigned in Stage 2 **will not** be implemented.

As you chose *Keep*, the Deduction Points your partner assigned to you in Stage 2 **will** be implemented.

Your partner has assigned **3** Deduction Points to you.

This reduces your final earnings with **9 Points** and reduces your partner's earnings with 3 Points.

Your total earnings in this HIT are

25 Points you started with
*plus* 25 Points from Stage 1
*minus* 9 Points from Stage 2
= **41 Points.**

Your partner's total earnings in this HIT are
25 Points they started with
*plus* 9 Points from Stage 1
*minus* 3 Points from Stage 2
= 31 Points.

Your Points are worth **$ 2.05**.

Your guaranteed participation fee is **$ 1**.

To collect your earnings, please copy the below code and paste it into MTurk. We will try to pay you as soon as possible.

**2000306**

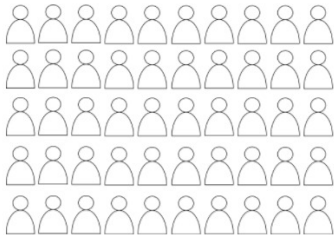Once you have done that, you can close this window. Thank you for your participation!

## The 50-peers treatment

The welcome page and the instructions for the first stage of the game are identical to the 4-peers treatment. Hence they are omitted. The instructions for the second stage of the game—the punishment stage—are provided below.



**Stage 2**

In Stage 2 you will decide to assign Deduction Points if in Stage 1 your partner has chosen *Keep*.

To give you an idea of what people usually choose in Stage 1, we will inform you of **how many of fifty randomly selected MTurkers (NOT your partner) chose** *Share* in the HITs prior to this one.

Depending on the the fifty previous MTurkers' choices, there are five possible situations:

**Situation 1**: Between **0 to 10** of them chose *Share*.

**Situation 2**: Between **10 to 20** of them chose *Share*.

**Situation 3**: Between **20 to 30** of them chose *Share*.

**Situation 4**: Between **30 to 40** of them chose *Share*.

**Situation 5**: Between **40 to 50** of them chose *Share*.

For each of these situations you will indicate how many Deduction Points to assign to your partner.

After you complete your choices for all five situations, we will let you know which situation actually occurred.

Your choice in that actual situation will be used to calculate the earnings of you and your partner.

**Note:**

*Your partner is NOT one of these previous MTurkers.*
*You CANNOT assign Deduction Points to any of these previous MTurkers.*
*Your choice on the Deduction Points only affects the earnings of you and your partner.*

| Previous | Next |
|---|---|

**Control Questions**

Before you make your choices, please answer the following questions to make sure you understand your task.

Suppose that in Stage 1, your partner chose *Keep*. The situation is that in Stage 1, between 40 to 50 of the fifty selected previous MTurkers chose *Keep*. In Stage 2 you choose to assign 3 Deduction Points to your partner.

1. How many Points will be deducted from your earnings?

2. How many Points will be deducted from the previous MTurkers' earnings?

3. How many Points will be deducted from your partner's earnings?

Previous

Next

---

# Your Choice for Stage 2

Situation 1

Between **0 to 10** of the fifty previous MTurkers (NOT your partner) chose *Share.*



Enter the Deduction Points you assign to your partner if your partner chose *Keep* in this situation:



Next

## Your Choice for Stage 2

Situation 2

Between **10 to 20** of the fifty previous MTurkers (NOT your partner) chose *Share*.



Enter the Deduction Points you assign to your partner if your partner chose *Keep* in this situation:
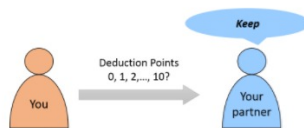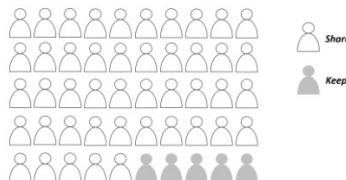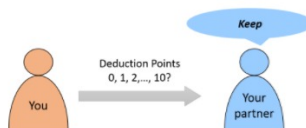


Previous

Next

---

## Your Choice for Stage 2

Situation 3

Between **20 to 30** of the fifty previous MTurkers (NOT your partner) chose *Share*.



Enter the Deduction Points you assign to your partner if your partner chose *Keep* in this situation:
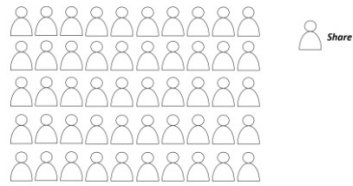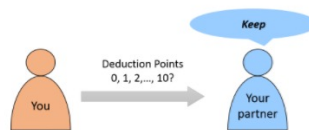


Previous

Next

# Your Choice for Stage 2

... In addition, we ask you to consider the following situation:

*All* of the fifty previous MTurkers (NOT your partner) chose *Share*.



Enter the Deduction Points you assign to your partner if your partner chose *Keep* in this situation:



Submit

# Chapter 5

# Conclusion

In this concluding chapter, I briefly mention some ideas for future research. In *Concentration of Influence under Complementarities*, I analyze optimal networks by setting aside incentive problems that might arise due to asymmetric information. However, asymmetric information appears in many natural settings and new complexities may arise. For example, when a manager designs the production network between a set of workers, the manager may not know each worker's ability ex ante, and the workers' efforts may be unobservable to the manager. When a financial regulator wants to intervene in the inter-bank network to improve the stability of the financial system, the regulator may not have precise information about each bank's financial status. There are also situations in which the agents may build their own collusion networks within the interaction and incentive constraints set by the planner, and the collusion networks are not observable to the planner. To analyze these settings, we need to combine contract theory and network analysis, which constitutes an important agenda for future research.

# Bibliography

Abramson, G. and Kuperman, M. (2001). Social Games in a Social Network. *Physical Review E*, 63(3):030901.

Andreoni, J. (1990). Impure Altruism and Donations to Public Goods: A Theory of Warm-Glow Giving. *Economic Journal*, 100(401):464.

Arechar, A. A., Gächter, S., and Molleman, L. (2018). Conducting Interactive Experiments Online. *Experimental Economics*, 21(1):99–131.

Baetz, O. (2015). Social activity and network formation. *Theoretical Economics*, 10(2):315–340.

Ballester, C., Calvo-Armengol, A., Zenou, Y., Calvó-Armengol, A., and Zenou, Y. (2006). Who's Who in the Network. Wanted: the Key Player. *Econometrica*, 74(5):1403–1417.

Barr, A. and Serneels, P. (2009). Reciprocity in the workplace. *Experimental Economics*, 12(1):99–112.

Becker, G. (1974). A Theory of Social Interactions. *Journal of Political Economy*, 82(6):1063–1093.

Becker, G. (1978). *The Economic Approach to Human Behavior*. University of Chicago press.

Belhaj, M., Bervoets, S., and Deroïan, F. (2016). Efficient networks in games with local complementarities. *Theoretical Economics*, 11(1):357–380.

Bergstrom, T., Blume, L., and Varian, H. (1986). On the private provision of public goods. *Journal of Public Economics*, 29(1):25–49.

Bicchieri, C. (2006). *The grammar of society: the nature and origins of social*

*norms*. Cambridge University Press.

Bonacich, P. (1987). Power and Centrality: A Family of Measures. *American Journal of Sociology*, 92(5):1170–1182.

Bonacich, P. and Lloyd, P. (2001). Eigenvector-like Measures of Centrality for Asymmetric Relations. *Social Networks*, 23(3):191–201.

Boyd, R. and Richerson, P. J. (1985). *Culture and the Evolutionary Process*, volume 175. University of Chicago Press.

Bramoullé, Y. and Kranton, R. (2007). Public Goods in Networks. *Journal of Economic Theory*, 135(1):478–494.

Bramoullé, Y., Kranton, R., and D'Amours, M. (2014). Strategic Interaction and Networks. *American Economic Review*, 104(3):898–930.

Brandts, J. and Charness, G. (2011). The Strategy versus the Direct-response Method: A First Survey of Experimental Comparisons. *Experimental Economics*, 14(3):375–398.

Cassar, A. (2007). Coordination and Cooperation in Local, Random and Small World Networks: Experimental Evidence. *Games and Economic Behavior*, 58(2):209–230.

Connelly, B. L., Haynes, K. T., Tihanyi, L., Gamache, D. L., and Devers, C. E. (2013). Minding the Gap: Antecedents and Consequences of Top Management-To-Worker Pay Dispersion. *Journal of Management*, 42(4):862–885.

Cooper, D. J. and Dutcher, E. G. (2011). The Dynamics of Responder Behavior in Ultimatum Games: A Meta-Study. *Experimental Economics*, 14(4):519–546.

Currarini, B. S., Jackson, M. O., Pin, P., Currarini, S., Jackson, M. O., and Pin, P. (2009). An Economic Model of Friendship: Homophily, Minorities, and Segregation. *Econometrica*, 77(4):1003–1045.

Duesenberry, J. S. (1960). Comment on "An Economic Analysis of Fertility" by Gary S. Becker. *Demographic and Economic Change in Developed Countries*, pages 231–34.

Dufwenberg, M. and Kirchsteiger, G. (2004). A Theory of Sequential Reciprocity. *Games and Economic Behavior*, 47(2):268–298.

Edmans, A. and Gabaix, X. (2016). Executive Compensation: A Modern Primer. *Journal of Economic Literature*, 54(4):1232–1287.

Ellison, G. (1993). Learning, Local Interaction, and Coordination. *Econometrica*, 61(5):1047–1071.

Ellison, G. (2000). Basins of Attraction, Long-Run Stochastic Stability, and the Speed of Step-by-Step Evolution. *Review of Economic Studies*, 67(1):17–45.

Elster, J. (1989). Social Norms and Economic Theory. *Journal of Economic Perspectives*, 3(4):99–117.

Eshel, I., Samuelson, L., and Shaked, A. (1998). Altruists, Egoists, and Hooligans in a Local Interaction Model. *American Economic Review*, 88(1):157–179.

Faleye, O., Reis, E., and Venkateswaran, A. (2013). The Determinants and Effects of CEO-employee Pay Ratios. *Journal of Banking and Finance*, 37(8):3258–3272.

Falk, A., Fehr, E., and Zehnder, C. (2006). Fairness Perceptions and Reservation Wages–the Behavioral Effects of Minimum Wage Laws. *Quarterly Journal of Economics*, 121(4):1347–1381.

Falk, A. and Ichino, A. (2006). Clean Evidence on Peer Effects. *Journal of Labor Economics*, 24(1):39–57.

Fehr, E., Gachter, S., and Kirchsteiger, G. (1997). Reciprocity as a Contract Enforcement Device: Experimental Evidence. *Econometrica*, 65(4):833–860.

Fehr, E. and Schmidt, K. M. (1999). A Theory of Fairness, Competition, and Cooperation. *Quarterly Journal of Economics*, 114(3):817–868.

Fehr, E. and Schmidt, K. M. (2000). Fairness, Incentives, and Contractual Choices. *European Economic Review*, 44(4):1057–1068.

Fehr, E. and Schmidt, K. M. (2004). Fairness and incentives in a multi-task principal-agent model. *Scandinavian Journal of Economics*, 106(3):453–474.

Fehr, E. and Schmidt, K. M. (2006). Chapter 8 The Economics of Fairness, Reci-

procity and Altruism - Experimental Evidence and New Theories. *Handbook of the Economics of Giving, Altruism and Reciprocity*, 1:615–691.

Fudenberg, D. and Tirole, J. (1991). Perfect Bayesian Equilibrium and Sequential Equilibrium. *Journal of Economic Theory*, 53(2):236–260.

Gächter, S. and Schulz, J. F. (2016). Intrinsic honesy and the prevalencce of rule violations across societies. *Nature*, 531(7595):1–11.

Galeotti, A. and Goyal, S. (2010). The Law of the Few. *American Economic Review*, 100(4):1468–1492.

Goyal, S. and Bala, V. (2000). a Noncooperative Model of Network Formation. *Econometrica*, 68(5):1181–1229.

Gracia-Lázaro, C., Cuesta, J. A., Sánchez, A., and Moreno, Y. (2012). Human Behavior in Prisoner's Dilemma Experiments Suppresses Network Reciprocity. *Scientific Reports*, 2:2–5.

Granovetter, M. (1985). Economic Action and Social Structure: The Problem of Embeddedness. *American Journal of Sociology*, 91(3):481–510.

Grujic, J., Fosco, C., Araujo, L., Cuesta, J. A., and Sanchez, A. (2010). Social Experiments in the Mesoscale: Humans Playing a Spatial Prisoner's Dilemma. *PLoS ONE*, 5(11):e13749.

Henrich, J. and Boyd, R. (1998). The Evolution of Conformist Transmission and the Emergence of Between-Group Differences. *Evolution and Human Behavior*, 19(4):215–241.

Henrich, J. and Boyd, R. (2001). Why People Punish Defectors: Weak Conformist Transmission Can Stabilize Costly Enforcement of Norms in Cooperative Dilemmas. *Journal of Theoretical Biology*, 208(1):79–89.

Henrich, J., Ensminger, J., McElreath, R., Barr, A., Barrett, C., Bolyanatz, A., Cardenas, J. C., Gurven, M., Gwako, E., Henrich, N., Lesorogol, C., Marlowe, F., Tracer, D., and Ziker, J. (2010). Markets, Religion, Community Size, and the Evolution of Fairness and Punishment. *Science*, 327(5972):1480–1484.

Henrich, J., McElreath, R., Barr, A., Ensminger, J., Barrett, C., Bolyanatz, A., Cardaroas, J. C., Gurven, M., Gwako, E., Henrich, N., Lesoronol, C., Marlowe, F., Tracer, D., and Ziker, J. (2006). Costly Punishment across Human Societies. *Science*, 312(5781):1767–1770.

Herrmann, B., Thöni, C., and Gächter, S. (2008). Antisocial Punishment across Societies. *Science*, 319(5868):1362–1367.

Herz, H. and Taubinsky, D. (2017). What Makes a Price Fair? An Experimental Study of Transaction Experience and Endogenous Fairness Views. *Journal of the European Economic Association*, 16(2):316–352.

Hiller, T. (2017). Peer Effects in Endogenous Networks. *Games and Economic Behavior*, 105:349–367.

Home, I., Forum, I., Resources, R., and Hanft, A. (2006). An Evolutionary Approach to Innovation. *Innovation*, 80(4):1–3.

Kahneman, D., Knetsch, J. L., and Thaler, R. H. (1986). Fairness as a Constraint on Profit Seeking: Entitlements in the Market. *American Economic Review*, 76(4):728–741.

Kandori, M. (1992). Social Norms and Community Enforcement. *Review of Economic Studies*, 59(1):63.

Kandori, M., Mailath, G. J., and Rob, R. (1993). Learning, Mutation, and Long Run Equilibria in Games. *Econometrica*, 61(1):29–56.

Kirchkamp, O. and Nagel, R. (2007). Naive Learning and Cooperation in Network Experiments. *Games and Economic Behavior*, 58(2):269–292.

Kreps, D. M. and Wilson, R. (1982). Sequential Equilibria. *Econometrica*, 50(4):863–894.

Mas, A. and Moretti, E. (2009). Peers at Work. *American Economic Review*, 99(1):112–145.

McPherson, M., Popielarz, P. a., and Drobnic, S. (1992). Social Networks and Organizational Dynamics. *American Sociological Review*, 57(2):153–170.

McPherson, M., Smith-Lovin, L., and Cook, J. M. (2001). Birds of a Feather:

Homophily in Social Networks. *Annual Review of Sociology*, 27(1):415–444.

Mishel, L. and Sabadish, N. (2013). Ceo Pay in 2012 Was Extraordinarily High Relative To Typical Workers and Other High Earners. *Economic Policy Institute. Issue Brief*, (367).

Nowak, M. A. and May, R. M. (1992). Evolutionary Games and Spatial Chaos. *Nature*, 359(6398):826–829.

Nowak, M. A., Tarnita, C. E., and Antal, T. (2010). Evolutionary Dynamics in Structured Populations. *Philosophical Transactions of the Royal Society B: Biological Sciences*, 365(1537):19–30.

O'Connell, S. and Siafarikas, A. (2010). *Addison Disease: Diagnosis and Initial Management*, volume 39. Princeton University Press.

Offerman, T. (2002). Hurting Hurts More than Helping Helps. *European Economic Review*, 46(8):1423–1437.

Ohtsuki, H. (2006). A Simple Rule for the Evolution of Cooperation on Graphes and Social Networks. *Nature*, 441(1):1–11.

Olson, M. (1965). *Logic of Collective Action: Public Goods and the Theory of Groups*. Harvard University Press.

Peysakhovich, A. and Rand, D. G. (2016). Habits of Virtue: Creating Norms of Cooperation and Defection in the Laboratory. *Management Science*, 62(3):631–647.

Rabin, M. (1993). Incorporating Fairness into Game Theory and Economics. *American Economic Review*, 83(5):1281–1302.

Rand, D. G., Arbesman, S., and Christakis, N. A. (2011). Dynamic Social Networks Promote Cooperation in Experiments with Humans. *Proceedings of the National Academy of Sciences*, 108(48):19193–19198.

Santos, F. C. and Pacheco, J. M. (2005). Scale-free Networks Provide a Unifying Framework for the Emergence of Cooperation. *Physical Review Letters*, 95(9):098104.

Selten, R. (1967). Die Strategiemethode zur Erforschung des eingeschränkt ra-

tionalen Verhaltens im Rahmen eines Oligopolexperiments. In *Beiträge zur experimentellen Wirtschaftsforschung*, pages 136–168. Seminar für Mathemat. Wirtschaftsforschung u. Ökonometrie.

Stigler, G. J. and Becker, G. S. (1977). De Gustibus Non Est Disputandum. *American Economic Review*, 67(5):76–90.

Sugden, R. (1982). On the Economics of Philanthropy. *Economic Journal*, 92(366):341–350.

Sugden, R. (1984). Reciprocity: The Supply of Public Goods Through Voluntary Contributions. *Economic Journal*, 94(376):772–787.

Sugden, R. (1986). *The Economics of Rights, Co-operation, and Welfare*. B. Blackwell.

Sugden, R. (2000). The Motivating Power of Expectations. *Rationality, Rules and Structure*, pages 103–129.

Sugden, R. (2004). *The Economics of Rights, Cooprration and Welfare*. Palgrave Macmillan, New York, 2005 edition.

Suri, S. and Watts, D. J. (2011). Cooperation and Contagion in Web-based, Networked Public Goods Experiments. *PLoS ONE*, 6(3):3–8.

Szabó, G. and Fáth, G. (2007). Evolutionary Games on Graphs. *Physics Reports*, 446(4-6):97–216.

Tabellini, G. (2008). Presidential Address: Institutions and Culture. *Journal of the European Economic Association*, 6(2-3):255–294.

Traulsen, A., Semmann, D., Sommerfeld, R. D., Krambeck, H.-J., and Milinski, M. (2010). Human strategy updating in evolutionary games. *Proceedings of the National Academy of Sciences*, 107(7):2962–2966.

Wrong, D. H. (1961). The Oversocialized Conception of Man in Modern Sociology. *American Sociological Review*, 26(2):183–193.

Young, H. P. (1993). The Evolution of Conventions. *Econometrica*, 61(1):57–84.

Young, P. (2001). *Individual Strategy and Social Structure: An Evolutionary Theory of Institutions*. Princeton University Press.

Zizzo, D. J. (2010). Experimenter Demand Effects in Economic Experiments. *Experimental Economics*, 13(1):75–98.