

# Social Norms as Choreography

Herbert Gintis\*

July 20, 2009

## 1 Introduction

This paper extends the seminal contributions of David Lewis (1969), Michael Taylor (1976, 1982, 1987), Robert Sugden (1986, 1989), Cristina Bicchieri (1992, 1993, 1997, 1999, 2006), and Ken Binmore (1993, 1998, 2005) in treating social norms as Nash equilibria of non-cooperative games played by rational agents. The insight underlying all these contributions is that if agents play a game  $\mathcal{G}$  with several Nash equilibria, a social norm can serve to choose among these equilibria. While this insight applies to several important social situations, it does not apply to most. In this paper I will suggest a more general principle, according to which a social norm is a *choreographer* of a supergame  $\mathcal{G}^+$  of  $\mathcal{G}$ . By the term ‘choreographer’ I mean a correlating device that implements a correlated equilibrium of  $\mathcal{G}$  in which all agents play strictly pure strategies (these terms are defined below).

The social norms as choreographer has two attractive properties lacking in the social norms as Nash equilibria. First, the conditions under which rational agents play Nash equilibria are generally complex and implausible, whereas rational agents in a very natural sense play correlated equilibria (see Section 7). Second, the social norms as Nash equilibria approach cannot explain why compliance with social norms is often based on other-regarding and moral preferences in which agents are willing to sacrifice on behalf of compliance with social norms. We can explain this association between norms and morality in terms of the incomplete information possessed by the choreographer. Morality, in this view, is doing the right thing even if no one is looking.

---

\*Santa Fe Institute and Central European University. Prepared for a special issue of *Politics, Philosophy and Economics* edited by Cristina Bicchieri. This material is taken from my book *The Bounds of Reason: Game Theory and the Unification of the Behavioral Sciences* (Princeton, 2009).

## 2 From Nash to Correlated Equilibrium

		Alice	
		$m$	$f$
Bob	$m$	2,1	0,0
	$f$	0,0	1,2

**Figure 1:** The Battle of the Sexes Game

Consider a society in which men prefer the company of women and vice-versa, but when they consort, their two forms of entertainment,  $m$  and  $f$ , are favored by men and women, respectively. Their payoffs are described in Figure 1. There are two pure-strategy equilibria and one mixed-strategy equilibrium for this game. Clearly both sexes would be better off if they stuck to either of their pure-strategy equilibria than by choosing the mixed strategy equilibrium, in which each plays his or her favorite choice with probability  $1/3$ , resulting in the payoff  $2/3$  to each.

No principle of rational choice objectively favors any of these three equilibria, but a good case could be made for the mixed strategy equilibrium, as it conforms to *a priori* symmetry principles that are likely to hold in the absence of other information. Suppose there is a social norm that says “the man always gets to decide where to go.” Then, if both men and women believe that this social norm is in effect, each knows the other will choose  $m$ , and hence each will choose  $m$ , thus validating the social norm.

For a more complicated but more realistic example, consider a town with a rectangular North-South-East-West array of streets. In the absence of a social norm, whenever two cars find themselves in a condition of possible collision, both stop and each waits for the other go first. Obviously not a lot of driving will get done. So, consider a social norm in which (a) all cars drive on the right, (b) at an intersection both cars stop and the car that arrived first proceeds forward, and (c) if both cars arrive at an intersection at the same time, the car that sees the other car on its left proceeds forward. This is one of several social norms that will lead to an efficient use of the system of streets, provided there is not too much traffic.

Suppose, however, that there is so much traffic that cars spend much of their time stopping a crossings. We might then prefer the social norm in which we amend the above social norm to say that cars traveling North-South always have the right of way and need not stop at intersections. However, if there is really heavy

traffic, East-West drivers may never get a chance to move forward at all. Moreover, if some intersections violate the strict North-South-East-West orientation, it may be unclear who does have the right of way in some cases.

Suppose, then, we erect a set of signals at each intersection that indicate “Go” or “Stop” to drivers moving in one direction and another set of “Go” or “Stop” signals for drivers moving in the crossing direction. We can then correlate the signals so that when one set of drivers see “Go”, the other set of drivers see “Stop.” The social norm then says that “if you see Go, do not stop at the intersection, but if you see Stop, then stop and wait for the signal to change to Go.” We add to the social norm that the system of signals alternates sufficiently rapidly and there is a sufficiently effective surveillance system that no driver has an incentive to disobey the social norm.

This would appear to be a perfect example of a social norm, indeed a convention. However, the original game does not have a system of signals, and the proposed social norm does not single out a Nash equilibrium of the original game. Indeed, it is easy to see that there is a wide array of payoffs in the original game in which the only Nash equilibrium is for both cars to stop when an encounter occurs.

The system of signals in fact represents what is called a *correlated equilibrium* of  $\mathcal{G}$  (Aumann 1974, 1987). Basically, a correlated equilibrium adds a new player, whom I shall call the *choreographer* (Gintis 2009), who has signal set  $S = \{\text{Go}, \text{Stop}\}$  who views a set  $\Omega$  of “states of nature,” and for each intersection  $i$  in the town, and each state of nature  $\omega \in \Omega$ , chooses a signal  $s_{in}(\omega) \in S$  for the North-South drivers at  $i$  and a signal  $s_{ie}(\omega) \in S$  with  $s_{ie}(\omega) \neq s_{in}(\omega)$  for the East-West drivers. For simplicity, we may think of  $\omega \in \Omega$  as a time of day, or as a time elapsed since last signal flip, so the choreographer flips the signals to the two groups of drivers according to some time schedule. The social norm is then the strict Nash equilibrium of the expanded game  $\mathcal{G}^+$  in which all rational agents obey the traffic laws.

In no sense is this a Nash equilibrium of the original game. Nor is  $\mathcal{G}^+$  unique; we can propose many alternative correlating devices, based on different state spaces  $\Omega$ , that produce substantially different patterns of traffic. For instance, we can include in each  $\omega \in \Omega$  a measure of the volume of traffic in the two directions as the intersection, and the choreography can increase the Go time for the drivers that are currently in the more congested direction.

For another example, consider a society in which agents contest for possession of valuable territory. The game  $\mathcal{G}$  has two possible strategies. The hawk ( $H$ ) strategy is to escalate battle until injured or your opponent retreats. The dove ( $D$ ) strategy is to display hostility but retreat before sustaining injury if your opponent escalates. The payoff matrix is given in the figure, where  $v > 0$  is the value of territory,  $w > v$  is the cost of injury, and  $(v - w)/2$  is the payoff when two hawks

	$H$	$D$
$H$	$\frac{v-w}{w}, \frac{v-w}{w}$	$v, 0$
$D$	$0, v$	$\frac{v}{2}, \frac{v}{2}$

**Figure 2:** The Hawk-Dove Game

meet. The agents can play mixed strategies, but they cannot condition their play on whether they are player 1 or player 2, and hence players cannot condition their behavior on being player 1 or player 2. The payoffs are shown in Figure 2.

The Hawk-Dove game has a unique symmetric equilibrium, determined as follows. Let  $\alpha$  be the probability of playing hawk. The payoff to playing hawk is then  $\pi_h = \alpha(v - w)/2 + (1 - \alpha)v$ , and the payoff to playing dove is  $\pi_d = \alpha(0) + (1 - \alpha)v/2$ . These two are equal when  $\alpha^* = v/w$ , so the unique symmetric Nash equilibrium occurs when  $\alpha = \alpha^*$ . The payoff to each player is thus

$$\pi_d = (1 - \alpha)\frac{v}{2} = \frac{v}{2} \left( \frac{w - v}{w} \right).$$

Note that when  $w$  is close to  $v$ , almost all the value of the territory is dissipated in fighting.

Clearly, because there is only one symmetric Nash equilibrium, the only possible social norm associated with a Nash equilibrium is extremely inefficient.

Suppose, however, that when two players contest, each knows which of the two happened upon the territory first. We may call the former the “incumbent” and the latter the “contester.” Consider the social norm that signals to the incumbent to play hawk and to the contester to play dove. Following the social norm, which we may call the *property rights* strategy, is not even a strategy of  $\mathcal{G}$ , but it is a third strategy to the augmented game  $\mathcal{G}^+$ . Note that if all individuals obey the property rights social norm, then there can be no efficiency losses associated with the allocation of property.

To see that we indeed have a correlated equilibrium, it is sufficient to show that if we add the property rights strategy  $P$  to the Hawk-Dove-Game, then  $P$  is a strict best response to itself. With this addition, we get the game depicted in figure 3. Note that the payoff to property against property,  $v/2$ , is greater than  $3v/4 - w/4$ , which is the payoff to hawk against property, and is also greater than  $v/4$ , which is the payoff to dove against property. Therefore, property is a strict Nash equilibrium. It is also efficient, because there is never a hawk-hawk confrontation in the property correlated equilibrium, so there is never any injury.

	H	D	P
H	$\frac{v-w}{2}, \frac{v-w}{2}$	$v, 0$	$\frac{3v-w}{4}, \frac{v-w}{4}$
D	$0, v$	$\frac{v}{2}, \frac{v}{2}$	$\frac{v}{4}, \frac{3v}{4}$
P	$\frac{v-w}{4}, \frac{3v-w}{4}$	$\frac{3v}{4}, \frac{v}{4}$	$\frac{v}{2}, \frac{v}{2}$

**Figure 3:** The Hawk-Dove-Property Game

The property equilibrium is a highly efficient correlated equilibrium  $\mathcal{G}^+$  of the Hawk-Dove game  $\mathcal{G}$ , and corresponds to the classical political economy defense of property rights, but it applies as well to non-human territorial animals and explains *status quo* bias and loss aversion in humans (Gintis 2007).

### 3 Nash Equilibrium and Correlated Equilibrium

There are important implications of the fact that a social norm is the choreographer of a correlated equilibrium rather than a Nash equilibrium selection device. A simple game  $\mathcal{G}$  may have many qualitatively distinct correlated extensions  $\mathcal{G}^+$ , which implies that life based on social norms can be significantly qualitatively richer than the simple underlying games that they choreograph. The correlated equilibrium concept thus indicates that social theory goes beyond game theory to the extent that it supplies dynamical and equilibrium mechanisms for the constitution and transformation of social norms. At the same time, the power of the correlated equilibrium interpretation of social norms indicates that social theory that rejects game theory is likely to be significantly handicapped.

Indeed, in a fundamental sense the correlated equilibrium is more basic than the Nash equilibrium. The epistemic conditions under which rational agents will play a Nash equilibrium are extremely confining and cannot be expected to hold in any but a small subset of even the simplest games, such as games with very few strategies per player that are solvable by the iterated elimination of strongly dominated strategies (Aumann and Brandenburger 1995, Basu 1994, Gintis 2009). By contrast, Aumann (1987) has shown that Bayesian rationality in a game-theoretic setting is effectively isomorphic with correlated equilibrium. I shall here sketch his argument, which is extremely straightforward, once the proper machinery is set up.

## 4 Epistemic Games

An *epistemic game*  $\mathcal{G}$  consists of a normal form game with players  $i = 1, \dots, n$  and a finite pure-strategy set  $S_i$  for each player  $i$ , so  $S = \prod_{i=1}^n S_i$  is the set of pure-strategy profiles for  $\mathcal{G}$ , with payoffs  $\pi_i : S \rightarrow \mathbf{R}$ . In addition,  $\mathcal{G}$  includes a set of possible states  $\Omega$  of the game, a *knowledge partition*  $\mathcal{P}_i$  of  $\Omega$  for each player  $i$ , and a *subjective prior*  $p_i(\cdot; \omega)$  over  $\Omega$  that is a function of the current state  $\omega$ . A state  $\omega$  specifies, possibly among other aspects of the game, the strategy profile  $s$  used in the game. We write this  $s = \mathbf{s}(\omega)$ . Similarly, we write  $s_i = \mathbf{s}_i(\omega)$  and  $s_{-i} = \mathbf{s}_{-i}(\omega)$ .

The subjective prior  $p_i(\cdot; \omega)$  represents  $i$ 's beliefs concerning the state of the game, including the choices of the other players, when the actual state is  $\omega$ . Thus,  $p_i(\omega'; \omega)$  is the probability  $i$  places on the current state being  $\omega'$  when the actual state is  $\omega$ . A *partition* of a set  $X$  is a set of mutually disjoint subsets of  $X$  whose union is  $X$ . We write the cell of the knowledge partition  $\mathcal{P}_i$  containing state  $\omega$  as  $\mathbf{P}_i\omega$ , and we interpret  $\mathbf{P}_i\omega \in \mathcal{P}_i$  as the set of states that  $i$  considers possible (i.e., among which  $i$  cannot distinguish) when the actual state is  $\omega$ . Therefore, we require that  $\mathbf{P}_i\omega = \{\omega' \in \Omega \mid p_i(\omega'|\omega) > 0\}$ . Because  $i$  cannot distinguish among states in the cell  $\mathbf{P}_i\omega$  of his knowledge partition  $\mathcal{P}_i$ , his subjective prior must satisfy  $p_i(\omega''; \omega) = p_i(\omega''; \omega')$  for all  $\omega'' \in \Omega$  and all  $\omega' \in \mathbf{P}_i\omega$ . Moreover, we assume a player believes the actual state is possible, so  $p_i(\omega|\omega) > 0$  for all  $\omega \in \Omega$ .

The possibility operator  $\mathbf{P}_i$  has the following two properties: for all  $\omega, \omega' \in \Omega$ ,

$$\begin{aligned} \text{(P1)} \quad & \omega \in \mathbf{P}_i\omega \\ \text{(P2)} \quad & \omega' \in \mathbf{P}_i\omega \Rightarrow \mathbf{P}_i\omega' = \mathbf{P}_i\omega \end{aligned}$$

P1 says that the current state is always possible (i.e.,  $p_i(\omega|\omega) > 0$ ), and P2 follows from the fact that  $\mathcal{P}_i$  is a partition: if  $\omega' \in \mathbf{P}_i\omega$ , then  $\mathbf{P}_i\omega'$  and  $\mathbf{P}_i\omega$  have nonempty intersection, and hence must be identical.

We call a set  $E \subseteq \Omega$  an *event*, and we say that player  $i$  *knows* the event  $E$  at state  $\omega$  if  $\mathbf{P}_i\omega \subseteq E$ ; i.e.,  $\omega' \in E$  for all states  $\omega'$  that  $i$  considers possible at  $\omega$ . We write  $\mathbf{K}_i E$  for the event that  $i$  knows  $E$ .

Given a possibility operator  $\mathbf{P}_i$ , we define the *knowledge operator*  $\mathbf{K}_i$  by

$$\mathbf{K}_i E = \{\omega \mid \mathbf{P}_i\omega \subseteq E\}.$$

The most important property of the knowledge operator is  $\mathbf{K}_i E \subseteq E$ ; i.e., if an agent knows an event  $E$  in state  $\omega$  (i.e.,  $\omega \in \mathbf{K}_i E$ ), then  $E$  is true in state  $\omega$  (i.e.,  $\omega \in E$ ). This follows directly from P1.

We can recover the possibility operator  $\mathbf{P}_i\omega$  for an individual from his knowl-

edge operator  $\mathbf{K}_i$ , because

$$\mathbf{P}_i \omega = \bigcap \{E \mid \omega \in \mathbf{K}_i E\}. \quad (1)$$

To verify this equation, note that if  $\omega \in \mathbf{K}_i E$ , then  $\mathbf{P}_i \omega \subseteq E$ , so the left hand side of (1) is contained in the right hand side. Moreover, if  $\omega'$  is not in the right hand side, then  $\omega' \notin E$  for some  $E$  with  $\omega \in \mathbf{K}_i E$ , so  $\mathbf{P}_i \omega \subseteq E$ , so  $\omega' \notin \mathbf{P}_i \omega$ . Thus the right hand side of (1) is contained in the left.

If  $\mathbf{P}_i$  is a possibility operator for  $i$ , the sets  $\{\mathbf{P}_i \omega \mid \omega \in \Omega\}$  form a partition  $\mathcal{P}$  of  $\Omega$ . Conversely, any partition  $\mathcal{P}$  of  $\Omega$  gives rise to a possibility operator  $\mathbf{P}_i$ , two states  $\omega$  and  $\omega'$  being in the same cell iff  $\omega' \in \mathbf{P}_i \omega$ . Thus, a knowledge structure can be characterized by its knowledge operator  $\mathbf{K}_i$ , by its possibility operator  $\mathbf{P}_i$ , by its partition structure  $\mathcal{P}$ , or even by the subjective priors  $p_i(\cdot \mid \omega)$ .

Since each state  $\omega$  in epistemic game  $\mathcal{G}$  specifies the players' pure strategy choices  $\mathbf{s}(\omega) = (\mathbf{s}_1(\omega), \dots, \mathbf{s}_n(\omega)) \in S$ , the players' subjective priors must specify their beliefs  $\phi_1^\omega, \dots, \phi_n^\omega$  concerning the choices of the other players. We have  $\phi_i^\omega \in \Delta S_{-i}$ , which allows  $i$  to assume other players' choices are correlated. This is because, while the other players choose independently, they may have communalities in beliefs that lead them independently to choose correlated strategies.

We call  $\phi_i^\omega$  player  $i$ 's *conjecture* concerning the behavior of the other players at  $\omega$ . Player  $i$ 's conjecture is derived from  $i$ 's *subjective prior* by defining  $\phi_i^\omega(s_{-i}) = p_i([s_{-i}]; \omega)$ , where  $[s_{-i}] \subset \Omega$  is the event that the other players choose strategy profile  $s_{-i}$ . Thus, at state  $\omega$ , each player  $i$  takes the action  $\mathbf{s}_i(\omega) \in S_i$  and has the subjective prior probability distribution  $\phi_i^\omega$  over  $S_{-i}$ . A player  $i$  is deemed *Bayesian rational* at  $\omega$  if  $\mathbf{s}_i(\omega)$  maximizes  $\pi_i(s_i, \phi_i^\omega)$ , where

$$\pi_i(s_i, \phi_i^\omega) =_{\text{def}} \sum_{s_{-i} \in S_{-i}} \phi_i^\omega(s_{-i}) \pi_i(s_i, s_{-i}). \quad (2)$$

In other words, player  $i$  is Bayesian rational in epistemic game  $\mathcal{G}$  if his pure-strategy choice  $\mathbf{s}_i(\omega) \in S_i$  for every state  $\omega \in \Omega$  satisfies

$$\pi_i(\mathbf{s}_i(\omega), \phi_i^\omega) \geq \pi_i(s_i, \phi_i^\omega) \quad \text{for } s_i \in S_i. \quad (3)$$

## 5 Example: A Simple Epistemic Game

Suppose Alice and Bob each choose heads (h) or tails (t), neither observing the other's choice. We can write the universe as  $\Omega = \{\text{hh}, \text{ht}, \text{th}, \text{tt}\}$ , where  $xy$  means Alice chooses  $x$  and Bob chooses  $y$ . Alice's knowledge partition is then  $\mathcal{P}_A = \{\{\text{hh}, \text{ht}\}, \{\text{th}, \text{tt}\}\}$ , and Bob's knowledge partition is  $\mathcal{P}_B = \{\{\text{hh}, \text{th}\}, \{\text{ht}, \text{tt}\}\}$ . Alice's possibility operator  $\mathbf{P}_A$  satisfies  $\mathbf{P}_A \text{hh} = \mathbf{P}_A \text{ht} = \{\text{hh}, \text{ht}\}$  and  $\mathbf{P}_A \text{th} = \mathbf{P}_A \text{tt} =$

{th, tt}, whereas Bob’s possibility operator  $\mathbf{P}_B$  satisfies  $\mathbf{P}_B hh = \mathbf{P}_B th = \{hh, th\}$  and  $\mathbf{P}_B ht = \mathbf{P}_B tt = \{ht, tt\}$ .

In this case, the event “Alice chooses h” is  $E_A^h = \{hh, ht\}$ , and because  $\mathbf{P}_A hh, \mathbf{P}_A ht \subset E$ , Alice knows  $E_A^h$  whenever  $E_A^h$  occurs (i.e.,  $E_A^h = \mathbf{K}_i E_A^h$ ). The event  $E_B^h$  expressing “Bob chooses h” is  $E_B^h = \{hh, th\}$ , and Alice does not know  $E_B^h$  because at th Alice believes tt is possible, but  $tt \notin E_B^h$ .

## 6 Correlated Strategies and Correlated Equilibria

We want to show that if players are Bayesian rational in an epistemic game  $\mathcal{G}$  and have a common prior over  $\Omega$ , the strategy profiles  $\mathbf{s} : \Omega \rightarrow S$  that they play form a correlated equilibrium (Aumann 1987). The converse also holds: for every correlated equilibrium of a game, there is an extension to an epistemic game  $\mathcal{G}$  with a common prior  $p \in \Omega$  such that in every state  $\omega$  it is rational for all players to carry out the move indicated by the correlated equilibrium.

Informally, a correlated equilibrium of an epistemic game  $\mathcal{G}$  is a Nash equilibrium of a game  $\mathcal{G}^+$ , which is  $\mathcal{G}$  augmented by an initial move by Nature, who observes a random variable  $\gamma$  on a probability space  $(\Gamma, p)$  and issues a directive  $f_i(\gamma) \in S_i$  to each player  $i$  as to which pure strategy to choose. Following Nature’s directive is a best response, if other players also follow Nature’s directives, provided players have the *common prior*  $p$ .

Formally, a *correlated strategy* of epistemic game  $\mathcal{G}$  consists of a finite probability space  $(\Gamma, p)$ , where  $p \in \Delta\Gamma$ , and a function  $f : \Gamma \rightarrow S$ . If we think of a choreographer who observes  $\gamma \in \Gamma$  and directs players to choose strategy profile  $f(\gamma)$ , then we can identify a correlated strategy with a probability distribution  $\tilde{p} \in \Delta S$ , where, for  $s \in S$ ,  $\tilde{p}(s) = p([f(\gamma) = s])$  is the probability that the choreographer chooses  $s$ . We call  $\tilde{p}$  the *distribution* of the correlated strategy. Any probability distribution on  $S$  that is the distribution of some correlated strategy  $f$  is called a *correlated distribution*.

Suppose  $f^1, \dots, f^k$  are correlated strategies and let  $\alpha = (\alpha_1, \dots, \alpha_k)$  be a lottery (i.e.,  $\alpha_i \geq 0$  and  $\sum_i \alpha_i = 1$ ). Then,  $f = \sum_i \alpha_i f^i$  is also a correlated strategy defined on  $\{1, \dots, k\} \times \Gamma$ . We call such an  $f$  a *convex sum* of  $f^1, \dots, f^k$ . Any convex sum of correlated strategies is clearly a correlated strategy. It follows that any convex sum of correlated distributions is itself a correlated distribution.

Suppose  $\sigma = (\sigma_1, \dots, \sigma_n)$  is a Nash equilibrium of a game  $\mathcal{G}$ , where for each  $i = 1, \dots, n$ ,

$$\sigma_i = \sum_{k=1}^{n_i} \alpha_{ki} s_{ki}$$

where  $n_i$  is the number of pure strategies in  $S_i$  and  $\alpha_{ki}$  is the weight given by  $\sigma_i$  on the  $k^{\text{th}}$  pure strategy  $s_{ki} \in S_i$ . Note that  $\sigma$  thus defines a probability distribution  $\tilde{p}$  on  $S$  such that  $\tilde{p}(s)$  is the probability that pure strategy profile  $s \in S$  will be chosen when mixed strategy profile  $\sigma$  is played. Then,  $\tilde{p}$  is a correlated distribution of an epistemic game associated with  $\mathcal{G}$ , which we will call  $\mathcal{G}$  as well. To see this, define  $\Gamma_i$  as a set with  $n_i$  elements  $\{\gamma_{1i}, \dots, \gamma_{n_i i}\}$  and define  $p_i \in \Delta S_i$  that places probability  $\alpha_{ki}$  on  $\gamma_{ki}$ . Then, for  $s = (s_1, \dots, s_n) \in S$ , define  $p(s) = \prod_{i=1}^n p_i(s_i)$ . Now, define  $\Gamma = \prod_{i=1}^n \Gamma_i$  and let  $f: \Gamma \rightarrow S$  be given by  $f(\gamma_{k_1 1}, \dots, \gamma_{k_n n}) = (s_{k_1 1}, \dots, s_{k_n n})$ . It is easy to check that  $f$  is a correlated strategy with correlated distribution  $\tilde{p}$ . In short, every Nash equilibrium is a correlated strategy, and hence any convex combination of Nash equilibria is a correlated strategy.

If  $f$  is a correlated strategy, then  $\pi_i \circ f$  is a real-valued random variable on  $(\Gamma, p)$  with an expected value  $\mathbf{E}_i[\pi_i \circ f]$ , the expectation taken with respect to  $p$ . We say a function  $g_i: \Gamma \rightarrow S_i$  is *measurable with respect to  $f_i$*  if  $f_i(\gamma) = f_i(\gamma')$ , then  $g_i(\gamma) = g_i(\gamma')$ . Clearly, player  $i$  can choose to follow  $g_i(\gamma)$  when he knows  $f_i(\gamma)$  iff  $g_i$  is measurable with respect to  $f_i$ . We say that a correlated strategy  $f$  is a *correlated equilibrium* if for each player  $i$  and any  $g_i: \Gamma \rightarrow S_i$  that is measurable with respect to  $f_i$ , we have

$$\mathbf{E}_i[\pi_i \circ f] \geq \mathbf{E}_i[\pi_i \circ (f_{-i}, g_i)].$$

A correlated equilibrium induces a *correlated equilibrium probability distribution* on  $S$ , whose weight for any strategy profile  $s \in S$  is the probability that  $s$  will be chosen by the choreographer. Note that a correlated equilibrium of  $\mathcal{G}$  is a Nash equilibrium of the game generated from  $\mathcal{G}$  by adding Nature, whose move at the beginning of the game is to observe the state of the world  $\gamma \in \Gamma$ , and to indicate a move  $f_i(\gamma)$  for each player  $i$  such that no player has an incentive to do other than comply with Nature's recommendation, provided that the other players comply as well.

## 7 Correlated Equilibrium and Bayesian Rationality

We now show that if the players in epistemic game  $\mathcal{G}$  are Bayesian rational at  $\omega$ , have a common prior  $p(\cdot; \omega)$  in state  $\omega$ , and each player  $i$  chooses  $\mathbf{s}_i(\omega) \in S_i$  in state  $\omega$ , then the distribution of  $\mathbf{s} = (\mathbf{s}_1, \dots, \mathbf{s}_n)$  is a correlated equilibrium distribution given by correlating device  $f$  on probability space  $(\Omega, p)$ , where  $f(\omega) = \mathbf{s}(\omega)$  for all  $\omega \in \Omega$ .

To prove this theorem, we identify the state space for the correlated strategy with the state space  $\Omega$  of  $\mathcal{G}$ , and the probability distribution on the state space with

the common prior  $p$ . We then define the correlated strategy  $f : \Omega \rightarrow S$  by setting  $f(\omega) = (\mathbf{s}_1(\omega), \dots, \mathbf{s}_n(\omega))$ , where  $\mathbf{s}_i(\omega)$  is  $i$ 's choice in state  $\omega$ . Then, for any player  $i$  and any function  $g_i : \Omega \rightarrow S_i$  that is  $\mathcal{P}_i$ -measurable (i.e., that is constant on cells of the partition  $\mathcal{P}_i$ ), because  $i$  is Bayesian rational, we have

$$\mathbf{E}[\pi_i(\mathbf{s}(\omega))|\omega] \geq \mathbf{E}[\pi_i(\mathbf{s}_{-i}(\omega), g_i(\omega))|\omega].$$

Now, multiply both sides of this inequality by  $p(P)$  and add over the disjoint cells  $P \in \mathcal{P}_i$ , which gives, for any such  $g_i$ ,

$$\mathbf{E}[\pi_i(\mathbf{s}(\omega))] \geq \mathbf{E}[\pi_i(\mathbf{s}_{-i}(\omega), g_i(\omega))].$$

This proves that  $(\Omega, f(\omega))$  is a correlated equilibrium. Note that the converse clearly holds as well.

## 8 Common Priors and Social Norm Equilibria

The isomorphism between correlated equilibrium distributions and Bayesian rationality demonstrated in Section 7 highlights an assumption that lies at the heart of a game-theoretic concept of social norms. This is the requirement that the players have a common prior over the state space  $\Omega$ . If the correlated equilibrium assigns a strict best response to each player, it is clear that some amount of preference heterogeneity will not destroy the equilibrium (the reader is invited to verify this). Moreover, if there are known “types” of players (e.g., Optimists and Pessimists) whose priors are distinct but commonly known, and the population composition is commonly known, it may be possible to redefine the state space so that there are common priors over the new state space, to which the correlated equilibrium theory then applies. The reader is invited to develop this theme.

However, when common priors are lacking and the actual composition and frequency distribution of priors are not held in common for some suitably enlarged state space, the social norm analysis will fail to apply. Rational agents with fundamental disagreements as to the actual structure of their social life do not dance to a choreographer’s instructions.

## 9 The Omniscient Choreographer and Social Preferences

The isomorphism between correlated equilibrium distributions and Bayesian rationality also requires that the choreographer be *omniscient* in the sense of having a knowledge partition is at least as fine as each of the player’s knowledge partition. This latter requirement was not explicitly mentioned in the proof, but is implicit in the requirement that  $f(\omega) = \mathbf{s}(\omega)$  for all  $\omega \in \Omega$ .

When this assumption fails, a correlated equilibrium may still obtain, provided the players have sufficiently strong prosocial preferences. Despite the fact that we have placed no restrictions on preferences other than Bayesian rationality, many social norms modelers, including Bicchieri (2006), predicate their analysis on the fact that rational individuals may have other-regarding preferences and/or may value certain moral virtues so that they voluntarily conform to a social norm in a situation where as perfectly self-regarding and amoral agent would not. In such cases, the choreographer may be obeyed even at a cost to the players, provided that the cost of doing so is not excessive.

For instance, each agent's payoff might consist of a *public component* that is known to the choreographer and a *private component* that reflects the idiosyncrasies of the agent and is unknown to the choreographer. Suppose the maximum size of the private component in any state for an agent is  $\alpha$ , but the agent's inclination to follow the choreographer has strength greater than  $\alpha$ . Then, the agent continues to follow the choreographer's directions whatever the state of his private information. Formally, we say an individual has an  $\alpha$ -*normative predisposition* towards conforming to the social norm if he strictly prefers to play his assigned strategy so long as all his pure strategies have payoffs no more than  $\alpha$  greater than when following the choreographer. We call an  $\alpha$ -normative predisposition a *social preference* because it facilitates social coordination but violates self-regarding preferences for  $\alpha > 0$ . There are evolutionary reasons for believing that humans have evolved such social preferences for fairly high levels of  $\alpha$  in a large fraction of the population through gene-culture coevolution (Gintis 2003).

Suppose, for example, that police in a certain town are supposed to apprehend criminals, where it costs police officer  $i$  a variable amount  $f_i(\omega)$  to file a criminal report. For instance, if the identified perpetrator is in the same ethnic group as  $i$ , or if the perpetrator offers a bribe to be released,  $f_i(\omega)$  might be very high, whereas an offender from a different ethnic group, or one who does not offer a bribe, might entail a low value of  $f_i(\omega)$ . How can this society erect incentives to induce the police to act in a non-corrupt manner?

Assuming police officer  $i$  is self-regarding and amoral,  $i$  will report a crime only if  $f_i(\omega) \leq w$ , where  $w$  is the reward for filing an accurate criminal report (accuracy can be guaranteed by fact-checking). A social norm equilibrium that requires that all apprehended criminals be prosecuted cannot then be sustained because all officers for whom  $f_i(\omega) < w$  with positive probability will at least at times behave corruptly. Suppose however officers have a *normative predisposition* to behave honestly, in the form of a police culture favoring honesty that is internalized by all officers. If  $f_i(\omega) < w + \alpha$  with probability one for all officers  $i$ , where  $\alpha$  is the strength of police culture, the social norm equilibrium can be sustained, despite the fact that the choreographer has incomplete information concerning events

in which criminal behavior is detected.

The following is a more complex example of a social norm equilibrium that requires a normative predisposition of honesty.

## 10 A Reputational Model of Honesty and Corruption

Consider a society in which sometimes people are Needy, and sometimes others help the Needy. In the first period, a pair of members is selected randomly, one of the pair being designated Needy and the other Giver. Giver and Needy then play a game  $\mathcal{G}$  in which if Giver helps, a benefit  $b$  is conferred on Needy at a cost  $c$  to Giver, where  $0 < c < b$ ; or, if Giver defects, both players receive 0. In each succeeding period, Needy from the previous period becomes Giver in the current period. Giver is then paired with a new, random Needy, and the game  $\mathcal{G}$  is played by the new pair. If we assume that helping behavior is common knowledge, there is a Nash equilibrium of the following form, provided the discount factor  $\delta$  is sufficiently close to unity. At the start of the game, each player is labeled “in good standing.” In every period Giver helps if and only if his partner Needy is in good standing. Failure to do so puts a player “in bad standing,” where he remains for the rest of the game.

To see that this is a Nash equilibrium in which every Giver helps in every period for  $\delta$  sufficiently close to 1, let  $v_c$  be the present value of the game to a Giver, and let  $v_b$  be the present value of the game for an individual who is not currently Giver or Needy. Then we have  $v_c = -c + \delta v_b$  and  $v_b = p(b + \delta v_c) + (1 - p)\delta v_b$ , where  $p$  is the probability of begin chosen as Needy. The first equation reflects the fact that a Giver must pay  $c$  now and becomes a candidate for Needy in the next period. The second equation expresses the fact that a candidate for Needy is chosen with probability  $p$  and then gets  $b$ , plus is Giver in the next period, and with probability  $1 - p$  remains a candidate for Needy in the next period. If we solve these two equations simultaneously, we find that  $v_c > 0$  precisely when  $\delta > c / (c + p(b - c))$ . Because the right hand side of this expression is strictly less than 1, there is a range of discount factors for which it is a best response for a Giver to help, and thus remain in good standing.

Suppose, however, the informational assumption is that each new Giver knows only whether his partner Needy did or did not help his own partner in the previous period. If Alice is Giver and her partner Needy is Bob, and Bob did not help when he was Giver, it could be because when he was Giver, Carole, his Needy partner, had defected when she was Giver, or because Bob failed to help Carole even though she had helped Donald, her previous Needy partner when she was Giver. Because Alice cannot condition her action on Bob’s previous action, Bob’s best response

is to defect on Carole, no matter what she did. Therefore, Carole will defect on Donald, no matter what he did. Thus, there can be no Nash equilibrium with the pure strategy of helping.

This argument extends to the richer informational structure where a Giver knows the previous  $k$  actions for any finite  $k$ . Here is the argument for  $k=2$ , which the reader is encouraged to generalize. Suppose the last five players are Alice, Bob, Carole, Donald, and Eloise, in that order. Alice can condition her choice on the actions taken by Bob, Carole, and Donald, but not on Eloise's action. Therefore, Bob's best response to Carole will not be conditioned on Eloise's action, and hence Carole's response to Donald will not be conditioned on Eloise's action. So, finally, Donald's response to Eloise will not be conditioned on her action, so her best response is to defect when she is Giver. Thus, there is no helping Nash equilibrium.

Suppose, however, back in the  $k=1$  case, that instead of defecting unconditionally when facing a Needy who has defected improperly, a Giver helps with probability  $p = 1 - c/b$  and defects with probability  $1 - p$ . The gain from helping unconditionally is then  $b - c$ , while the gain from following this new strategy is  $p(b - c) + (1 - p)pb$ , where the first term is the probability  $p$  of helping times the reward  $b$  in the next period if one helps minus the cost  $c$  of helping in the current period, and the second term is the probability  $1 - p$  of defecting times the probability  $p$  that you will be helped anyway when you are Needy, times the benefit  $b$ . Equating this expression with  $b - c$ , the cost of helping unconditionally, we get  $p = 1 - c/b$ , which is a number strictly between zero and one and hence a valid probability.

Consider the following strategy. In each round, Giver helps if his partner helped in the previous period, and otherwise helps with probability  $p$  and defects with probability  $1 - p$ . With this strategy each Giver  $i$  is indifferent to helping or defecting, because helping costs  $i$  the amount  $c$  when he is Giver but  $i$  gains  $b$  when he is Needy, for a net gain of  $b - c$ . However, defecting costs zero when Giver, but gives  $bp = b - c$  when he is Needy. Because the two actions have the same payoff, it is incentive-compatible for each Giver to help when his partner Needy helped, and to defect with probability  $p$  otherwise. This strategy thus gives rise to a Nash equilibrium with helping in every period (Bhaskar 1998),

However, there is no reason for rational self-regarding players to implement this Nash equilibrium. To see this, note that each player chooses a totally mixed strategy as Giver, and hence is indifferent between helping and defecting. Therefore, Givers have no incentive to play the equilibrium strategy. Erecting a social norm does not improve the situation because a Giver who is instructed by the choreography to help or to defect has no incentive to follow the instruction. However, if we add an  $\epsilon > 0$  of normative predisposition, an omniscient choreographer could

implement this Nash equilibrium by acting as the appropriate randomizing device. Moreover, suppose Givers have private preferences that, for instance, favor some players (e.g., friends or coreligionists) over others (e.g., enemies or infidels). In this case the choreographer's instructions will be followed only if players have a commitment to norm-following that is greater than their personal preferences to give or withhold aid to particular individuals.

## 11 Conclusion

In the first pages of *The Grammar of Society*, Cristina Bicchieri (2006) asserts that “social norms... *transform* mixed-motive games into coordination ones.” As we have seen in Section 2, this transformation is not always the case, but Bicchieri's affirmation is generally on the mark, and indeed, as shown in this paper, is the key to understanding the relationship between rational choice theory and social norms. Section 7 developed the central principle (Aumann 1987) that every state of an epistemic game  $\mathcal{G}$  at which players are rational can be implemented as a correlated equilibrium distribution, provided the appropriate epistemic conditions hold (common priors and choreographer omniscience). The associated correlated equilibrium is indeed a Nash equilibrium of an augmented game  $\mathcal{G}^+$ , in which an additional player, the choreographer (aka social norm) who implements the equilibrium, is added.

I believe that the epistemic game theoretic analysis of social norms presented in this paper can serve as the theoretical core for a general social theory of human strategic interaction. This analysis shows precisely where classical game theory goes wrong: it focuses on Nash as opposed to correlated equilibria, and hence ignores the rich social fabric of potential conditioning devices  $(\Omega, f)$ , each corresponding to a distinct social structure of interaction. Moreover, the theory renders salient the epistemic conditions for the existence of a social norm, conditions that are fulfilled only in an idealized, fully equilibrated, social system. In general, social norms will be contested and only partially implemented, and the passage from one choreographed equilibrium to another will be mediated by forms of collective action and individual heroism that cannot be currently explicated in game theoretic terms.

## REFERENCES

Aumann, Robert J., “Subjectivity and Correlation in Randomizing Strategies,” *Journal of Mathematical Economics* 1 (1974):67–96.

- , “Correlated Equilibrium and an Expression of Bayesian Rationality,” *Econometrica* 55 (1987):1–18.
- and Adam Brandenburger, “Epistemic Conditions for Nash Equilibrium,” *Econometrica* 65,5 (September 1995):1161–1180.
- Basu, Kaushik, “The Traveler’s Dilemma: Paradoxes of Rationality in Game Theory,” *American Economic Review* 84,2 (May 1994):391–395.
- Bhaskar, V., “Noisy Communication and the Evolution of Cooperation,” *Journal of Economic Theory* 82,1 (September 1998):110–131.
- Bicchieri, Cristina, *The Grammar of Society: The Nature and Dynamics of Social Norms* (Cambridge: Cambridge University Press, 2006).
- Binmore, Kenneth G., *Game Theory and the Social Contract: Playing Fair* (Cambridge, MA: MIT Press, 1993).
- , *Game Theory and the Social Contract: Just Playing* (Cambridge, MA: MIT Press, 1998).
- , *Natural Justice* (Oxford: Oxford University Press, 2005).
- Gintis, Herbert, “The Hitchhiker’s Guide to Altruism: Genes, Culture, and the Internalization of Norms,” *Journal of Theoretical Biology* 220,4 (2003):407–418.
- , “The Evolution of Private Property,” *Journal of Economic Behavior and Organization* 64,1 (September 2007):1–16.
- , *The Bounds of Reason: Game Theory and the Unification of the Behavioral Sciences* (Princeton, NJ: Princeton University Press, 2009).
- Lewis, David, *Conventions: A Philosophical Study* (Cambridge, MA: Harvard University Press, 1969).
- Sugden, Robert, *The Economics of Rights, Co-operation and Welfare* (Oxford: Basil Blackwell, 1986).
- , “Spontaneous Order,” *Journal of Economic Perspectives* 3,4 (Fall 1989):85–97.
- Taylor, Michael, *Anarchy and Cooperation* (London: John Wiley and Sons, 1976).
- , *Community, Anarchy, and Liberty* (Cambridge, UK: Cambridge University Press, 1982).
- , *The Possibility of Cooperation* (Cambridge, UK: Cambridge University Press, 1987).