
The Analytics of Human Sociality

The whole earth had one language. Men said, “Come, let us build ourselves a city, and a tower with its top in the heavens.” The Lord said, “Behold, they are one people, and they have all one language. Nothing will now be impossible for them. Let us go down and confuse their language.” The Lord scattered them over the face of the earth, and they ceased building the city.

Genesis 11:1

An economic transaction is a solved political problem. Economics has gained the title of Queen of the Social Sciences by choosing solved political problems as its domain.

Abba Lerner

10.1 Explaining Cooperation: An Overview

It is often said that sociology deals with cooperation and economics deals with competition. Game theory, however shows that cooperation and competition are neither distinct nor antithetical. Cooperation involves aligning the beliefs and incentives of agents with distinct interests, competition among groups requires cooperation within these groups, and competition among individuals may be mutually beneficial.

A major goal of economic theory is to show the plausibility of wide-scale cooperation among self-regarding individuals. In an earlier period, this endeavor centered on the Walrasian model of general market equilibrium, culminating in the celebrated Fundamental Theorem of Welfare Economics (Arrow and Debreu 1954, Debreu 1959, Arrow and Hahn 1971). However, the theorem’s key assumption that market exchange can be enforced at zero cost to the exchanging parties is often violated (Arrow 1971, Bowles and Gintis 1993, Gintis 2002, Bowles 2004).

The game theory revolution replaced reliance on exogenous enforcement with repeated game models in which punishment of defectors by cooperators secures cooperation among self-regarding individuals. Indeed, when a game \mathcal{G} is repeated an indefinite number of times by the same players, many of the anomalies associated with finitely repeated games (§5.1, 5.7, 4.11) dis-

appear. Moreover, Nash equilibria of the repeated game arise that are not Nash equilibria of \mathcal{G} . The exact nature of these equilibria is the subject of the Folk Theorem (§10.3), which shows that when individuals are Bayesian rational, self-regarding, have sufficiently long time-horizons, and there is adequate public information concerning who obeyed the rules and who did not, efficient social cooperation can be achieved in a wide variety of cases.

The Folk Theorem requires that a defection on the part of a player carry a signal that is conveyed to other players. We say a signal is *public* if all players receive the same signal. We say the signal is *perfect* if it accurately reports whether or not the player in question defected. The first general Folk Theorem that does not rely on incredible threats was proved by Fudenberg and Maskin (1986) for the case of perfect public signals (§10.3).

We say a signal is *imperfect* if it sometimes mis-reports whether or not the player in question defected. An imperfect public signal reports the same information to all players, but it is at times inaccurate. The Folk Theorem was extended to imperfect public signals by Fudenberg, Levine and Maskin (1994), as will be analyzed in §10.4.

If different players receive different signals, or some receive no signal at all, we say the signal is *private*. The case of private signals has proved much more daunting than that of public signals, but Folk Theorems for private but near-public signals (i.e., where there is an arbitrarily small deviation ϵ from public signals) have been developed by several game theorists, including Sekiguchi (1997), Piccione (2002), Ely and Välimäki (2002), Bhaskar and Obara (2002), Hörner and Olszewski (2006), and Mailath and Morris (2006). It is difficult to assess how critical the informational requirements of these Folk Theorems are, because generally the theorem is proved for “sufficiently small ϵ ,” with no discussion of the actual order of magnitude involved.

The question of the signal quality required for efficient cooperation to obtain is especially critical when the size of the game is considered. Generally, the Folk Theorem does not even mention the number of players, but in most situations, in real life, the larger the number of players participating in a cooperative endeavor, the lower the average quality of the cooperation vs. defection signal, because generally a player only observes a small number of others with a high degree of accuracy, however large the group involved. We explore this issue in §10.4, which illustrates the problem by applying the Fudenberg et al. (1994) framework to the Public Goods Game (§3.9) which in many respects is representative of contexts for cooperation in humans.

10.2 Bob and Alice Redux

Suppose Bob and Alice play the Prisoner's Dilemma shown on the right. In the one-shot game there is only one Nash equilibrium, in which both parties defect. However, suppose the same players play the game at times $t = 0, 1, 2, \dots$. This is then a new game, called a *repeated game*, in which the payoff to each is the sum of the payoffs over all periods, weighted by a *discount factor* δ , with $0 < \delta < 1$. We call the game played in each period the *stage game* of a *repeated game* in which at each period the players can condition their moves on the complete history of the previous stages. A strategy that dictates cooperating until a certain event occurs and then following a different strategy, involving defecting and perhaps otherwise harming one's partner, for the rest of the game is called a *trigger strategy*.

	<i>C</i>	<i>D</i>
<i>C</i>	5,5	-3,8
<i>D</i>	8,-3	0,0

Note that we have exactly the same analysis if we assume that players do not discount the future, but in each period the probability that the game continues at least one more period is δ . In general, we can think of δ as some combination of discount factor and probability of game continuation.

We show that the cooperative solution (5,5) can be achieved as a subgame perfect Nash equilibrium of the repeated game if δ is sufficiently close to unity and each player uses the trigger strategy of cooperating as long as the other player cooperates, and defecting forever if the other player defects on one round. To see this, consider a repeated game that pays 1 now and in each future period to a certain player, and the discount factor is δ . Let x be the value of the game to the player. The player receives 1 now and then gets to play exactly the same game in the next period. Because the value of the game in the next period is x , its present value is δx . Thus $x = 1 + \delta x$, so $x = 1/(1 - \delta)$.

Now suppose both agents play the trigger strategy. Then, the payoff to each is $5/(1 - \delta)$. Suppose a player uses another strategy. This must involve cooperating for a number (possibly zero) of periods, then defecting forever; for once the player defects, his opponent will defect forever, the best response to which is to defect forever. Consider the game from the time t at which the first player defects. We can call this $t = 0$ without loss of generality. A player who defects receives 8 immediately and nothing thereafter. Thus the cooperate strategy is Nash if and only if $5/(1 - \delta) \geq 8$, or $\delta \geq 3/8$. When δ satisfies this inequality, the pair of trigger strategies

is also subgame perfect, because the situation in which both parties defect forever is Nash subgame perfect.

This gives us an elegant solution to the problem, but in fact there are lots of other subgame perfect Nash equilibria to this game. For instance, Bob and Alice can trade off defecting on each other as follows. Consider the following trigger strategy for Alice: alternate C, D, C, \dots as long as Bob alternates D, C, D, \dots . If Bob deviates from this pattern, defect forever. Suppose Bob plays the complementary strategy: alternate D, C, D, \dots as long as Alice alternates C, D, C, \dots . If Alice deviates from this pattern, defect forever. These two strategies form a subgame perfect Nash equilibrium for δ sufficiently close to unity.

To see this, note that the payoffs are now $-3, 8, -3, 8, \dots$ for Alice and $8, -3, 8, -3, \dots$ for Bob. Let x be the payoffs to Alice. Alice gets -3 today, 8 in the next period and then gets to play the game all over again starting two periods from today. Thus, $x = -3 + 8\delta + \delta^2 x$. Solving this, we get $x = (8\delta - 3)/(1 - \delta^2)$. The alternative is for Alice to defect at some point, the most advantageous time being when it is her turn to get -3 . She then gets zero in that and all future periods. Thus, cooperating is Nash if and only if $x \geq 0$, which is equivalent to $8\delta - 3 \geq 0$, or $\delta \geq 3/8$.

For an example of a very unequal equilibrium, suppose Bob and Alice agree that Bob will play C, D, D, C, D, D, \dots and Alice will defect whenever Bob supposed to cooperate, and vice-versa. Let v_B be the value of the game to Bob when it is his turn to cooperate, provided he follows his strategy and Alice follows hers. Then, we have

$$v_B = -3 + 8\delta + 8\delta^2 + v_B\delta^3,$$

which we can solve, getting $v_B = (8\delta^2 + 8\delta - 3)/(1 - \delta^3)$. The value to Bob of defecting is 8 now and zero forever after. Hence, the minimum discount factor such that Bob will cooperate is the solution to the equation $v_B = 8$, which gives $\delta \approx 0.66$. Now let v_A be the value of the game to Alice when it is her first turn to cooperate, assuming both she and Bob follows their agreed strategies. Then we have

$$v_A = -3 - 3\delta + 8\delta^2 + v_A\delta^3,$$

which give $v_A = (8\delta^2 - 3\delta - 3)/(1 - \delta^3)$. The value to Alice of defecting rather than cooperating when it is her first turn to do so is then given by $v_A = 8$, which we can solve for δ , getting $\delta \approx 0.94$. With this discount factor, the value of the game to Alice is 8 , but $v_B \approx 72.47$, so Bob gains more than nine times as much as Alice.

10.3 The Folk Theorem

The *Folk Theorem* is so called because no one knows who first thought of it—it is just part of the “folklore” of game theory. We shall first present a stripped-down analysis of the Folk Theorem with an example and provide a more complete discussion in the next section.

Consider the stage game in §10.2. There is a subgame perfect Nash equilibrium in which each player gets zero. Moreover, neither player can be forced to receive a negative payoff in the repeated game based on this stage game, because zero can be assured simply by playing *D*. Also, any point in the region OEABCF in Fig. 10.1 could be attained in the stage game, assuming the players could agree on a mixed strategy for each. To see this, note that if Bob uses *C* with probability α and Alice uses *C* with probability β , then the expected payoff to the pair is $(8\beta - 3\alpha, 8\alpha - 3\beta)$, which traces out every point in the quadrilateral OEABCF for $\alpha, \beta \in [0, 1]$. Only the points in OABC are superior to the universal defect equilibrium $(0,0)$, however.

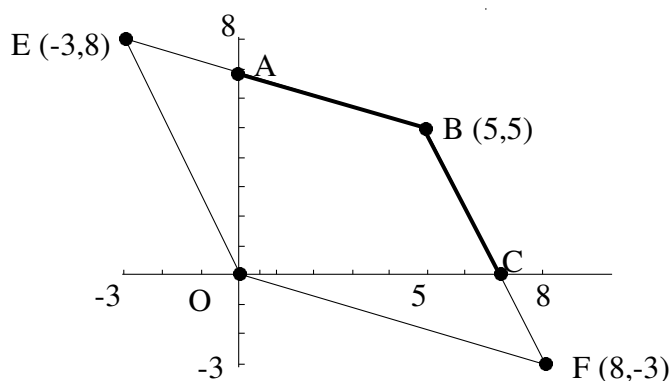


Figure 10.1. The Folk Theorem: any point in the region OABC can be sustained as the average per-period payoff the subgame perfect Nash equilibrium of the repeated game based on the stage game in §10.2.

Consider the repeated game \mathcal{R} based on the stage game \mathcal{G} of §10.2. The Folk Theorem says that under the appropriate conditions concerning the cooperate/defect signal available to players, any point in the region OABC can be sustained as the average per-period payoff of a subgame perfect

Nash equilibrium of \mathcal{R} , provided the discount factors of the players are sufficiently near unity.

More formally, consider an n -player game with finite strategy sets S_i for $i = 1, \dots, n$, so the set of strategy profiles for the game is $S = \prod_{i=1}^n S_i$. The payoff for player i is $\pi_i(s)$, where $s \in S$. For any $s \in S$ we write s_{-i} for the vector obtained by dropping the i th component of s , and for any $i = 1, \dots, n$ we write $(s_i, s_{-i}) = s$. For a given player j , suppose the other players choose strategies m_{-j}^j such that j 's best response m_j^j gives j the lowest possible payoff in the game. We call the resulting strategy profile m^j the *maximum punishment payoff* for j . Then, $\pi_j^* = \pi_j(m^j)$ is j 's payoff when everyone else "gangs up on him." We call

$$\pi^* = (\pi_1^*, \dots, \pi_n^*), \quad (10.1)$$

the *minimax point* of the game. Now define

$$\Pi = \{(\pi_1(s), \dots, \pi_n(s)) \mid s \in S, \pi_i(s) \geq \pi_i^*, i = 1, \dots, n\},$$

so Π is the set of strategy profiles in the stage game with payoffs at least as good as the maximum punishment payoff for each player.

This construction describes a stage game \mathcal{G} for a repeated game \mathcal{R} with discount factor δ , common to all the agents. If \mathcal{G} is played in periods $t = 0, 1, 2, \dots$, and if the sequence of strategy profiles used by the players is $s(1), s(2), \dots$, then the payoff to player j is

$$\tilde{\pi}_j = \sum_{t=0}^{\infty} \delta^t \pi_j(s(t)).$$

Let us assume that information is *public* and *perfect*, so that when a player deviates from some agreed-upon action in some period, a signal to this effect is transmitted with probability one to the other players. If players can use mixed strategies, then any point in Π can be attained as payoffs to \mathcal{R} by each player using the same mixed strategy in each period. However, it is not clear how a signal indicating deviation from a strictly mixed strategy should be interpreted. The simplest assumption guaranteeing the existence of such a signal is that there is a *public randomizing device* that can be seen by all players and that players use to decide which pure strategy to use, given that they have agreed to use a particular mixed strategy. Suppose, for instance, the randomizing device is a circular disc with a pointer that can

be spun by a flick of the finger. Then, a player could mark off a number of regions around the perimeter of the disc, the area of each being proportional to the probability of using each pure strategy in a given mixed strategy to be used by that player. In each period, each player flicks his pointer and chooses the appropriate pure strategy, this behavior is recorded accurately by the signaling device, and the result is transmitted to all players.

With these definitions, we have the following, where for $\pi \in \Pi$, $\sigma_i(\pi) \in \Delta S_i$ is a mixed strategy for player i such that $\pi_i(\sigma_1, \dots, \sigma_n) = \pi_i$:

THEOREM 10.1 Folk Theorem. *Suppose players have an available public randomizing device and the signal indicating cooperation or defection of each player is public and perfect. Then, for any $\pi = (\pi_1, \dots, \pi_n) \in \Pi$, if δ is sufficiently close to unity, there is a Nash equilibrium of the repeated game such that π_j is j 's payoff for $j = 1, \dots, n$ in each period. The equilibrium is effected by each player i using $\sigma_i(\pi)$ as long as no player has been signaled as having defected, and playing the minimax strategy m_i^j in all future periods after player j is first detected defecting.*

The idea behind this theorem is straightforward. For any such $\pi \in \Pi$, each player j uses the strategy $\sigma_j(\pi)$ that gives payoffs π in each period, provided the other players do likewise. If one player deviates, however, all other players play the strategies that impose the maximum punishment payoff on j forever. Because $\pi_j \geq \pi_j^*$, player j cannot gain from deviating from $\sigma_j(\pi)$, so the profile of strategies is a Nash equilibrium.

Of course, unless the strategy profile (m_1^j, \dots, m_n^j) is a Nash equilibrium for each $j = 1, \dots, n$, the threat to minimax even once, let alone forever, is not a credible threat. However, we do have the following:

THEOREM 10.2 The Folk Theorem with Subgame Perfection. *Suppose $y = (y_1, \dots, y_n)$ is the vector of payoffs in a Nash equilibrium of the underlying one-shot game, and $\pi \in \Pi$ with $\pi_i \geq y_i$ for $i = 1, \dots, n$. Then, if δ is sufficiently close to unity, there is a subgame perfect Nash equilibrium of the repeated game such that π_j is j 's payoff for $j = 1, \dots, n$ in each period.*

To see this, note that for any such $\pi \in \Pi$, each player j uses the strategy s_j that gives payoffs π in each period, provided the other players do likewise. If one player deviates, however, all players play the strategies that give payoff vector y forever.

10.4 The Folk Theorem with Imperfect Public Information

An important model due to Fudenberg et al. (1994) extends the Folk Theorem to many situations in which there is public imperfect signaling. Although their model does not discuss the n -player Public Goods Game, we shall here show that this game does satisfy the conditions for applying their theorem.

We shall see that the apparent power of the Folk Theorem in this case comes from letting the discount factor δ go to one *last*, in the sense that for any desired level of cooperation (by which we mean the level of *intended*, rather than *realized* cooperation), for any group size n and for any error rate ϵ , there is a δ sufficiently near unity that this level of cooperation can be realized. However, given δ , the level of cooperation may be quite low when n and ϵ are relatively small. Throughout this section, we shall assume that the signal imperfection takes the form of players defecting by accident with probability ϵ and hence failing to provide the benefit b to the group, although they expend the cost c .

The Fudenberg, Levine, and Maskin stage game consists of players $i = 1, \dots, n$, each with a finite set of pure actions $a_1, \dots, a_{m_i} \in A_i$. A vector $a \in A \equiv \prod_{j=1}^n A_j$ is called a pure action *profile*. For every profile $a \in A$ there is a probability distribution $y|a$ over the m possible public signals Y . Player i 's payoff, $r_i(a_i, y)$, depends only on his own action and the resulting public signal. If $\pi(y|a)$ is the probability of $y \in Y$ given profile $a \in A$, i 's expected payoff from a is given by

$$g_i(a) = \sum_{y \in Y} \pi(y|a) r_i(a_i, y).$$

Mixed actions and profiles, as well as their payoffs are defined in the usual way, and denoted by greek letters, so α is a mixed action profile, and $\pi(y|\alpha)$ is the probability distribution generated by mixed action α .

Note that in the case of a simple Public Goods Game, in which each player can cooperate by producing b for the other players at a personal cost c , each action set consists of the two elements $\{C,D\}$. We will assume that players choose only pure strategies. It is then convenient to represent the choice of C by 1 and D by 0. Let A be the set of strings of n zeros and ones, representing the possible pure strategy profiles of the n players, the k th entry representing the choice of the k th player. Let $\tau(a)$ be the number of ones in $a \in A$, and write a_i for the i th entry in $a \in A$. For any

$a \in A$, the random variable $y \in Y$ represents the imperfect public information concerning $a \in A$. We assume defections are signaled correctly, but intended cooperation fails and appears as defection with probability $\epsilon > 0$. Let $\pi(y|a)$ be the probability that signal $y \in A$ is received by players when the actual strategy profile is $a \in A$. Clearly, if $y_i > a_i$ for some i , then $\pi(y|a) = 0$. Otherwise

$$\pi(y|a) = \epsilon^{\tau(a)-\tau(y)}(1-\epsilon)^{\tau(y)} \quad \text{for } \tau(y) \leq \tau(a). \quad (10.2)$$

The payoff to player i who chooses a_i and receives signal y is given by $r_i(a_i, y|a) = b\tau(y)(1-\epsilon) - a_i c$. The expected payoff to player i is just

$$g_i(a) = \sum_{y \in Y} \pi(y|a)r_i(a_i, y) = b\tau(a)(1-\epsilon) - a_i c. \quad (10.3)$$

Moving to the repeated game, we assume in each period $t = 0, 1, \dots$, the stage game is played with public outcome $y^t \in Y$. The sequence $\{y^0, \dots, y^t\}$ is thus the *public history* of the game through time t , and we assume that the strategy profile $\{\sigma^t\}$ played at time t depends only on this public history (Fudenberg, Levine, and Maskin show that allowing agents to condition their play on their previous private profiles does not add any additional equilibrium payoffs). We call a profile $\{\sigma^t\}$ of public strategies a *perfect public equilibrium* if, for any period t and any public history up to period t , the strategy profile specified for the rest of the game is a Nash equilibrium from that point on. Thus, a public perfect equilibrium is subgame perfect Nash equilibrium implemented by public strategy profiles. The payoff to player i is then the discounted sum of the payoffs from each of the stage games.

The *minimax* payoff for player i is largest payoff i can attain if all the other players collude to choose strategy profiles that minimize i 's maximum payoff—see equation 10.1. In the Public Goods Game, the minimax payoff is zero for each player, because the worst the other players can do is universally defect, in which case i 's best action is to defect himself, giving payoff zero. Let V^* be the convex hull of stage game payoffs that dominate the minimax payoff for each player. A player who intends to cooperate and pays the cost c (which is not seen by the other players) can fail to produce the benefit b (which is seen by the other players) with probability $\epsilon > 0$. In the two-player case, V^* is the quadrilateral ABCD in Figure 10.2, where $b^* = b(1-\epsilon) - c$ is the expected payoff to a player if everyone cooperates.

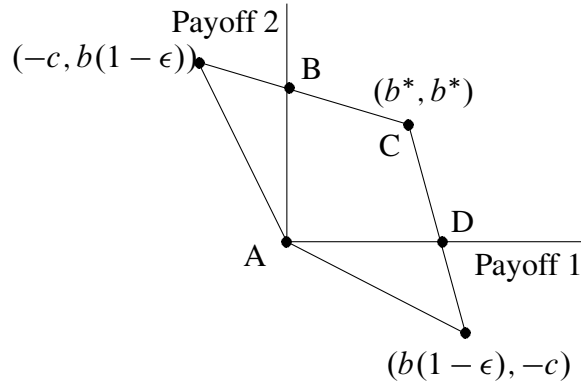


Figure 10.2. Two-player Public Goods Game

The Folk Theorem (Theorem 6.4, p. 1025 in Fudenberg, Levine, and Maskin) is then as follows.¹ We say $W \subset V^*$ is *smooth* if W is closed and convex, has a nonempty interior, and such that each boundary point $v \in W$ has a unique tangent hyperplane P_v that varies continuously with v (e.g., a closed ball with center interior to V^*). Then if $W \subset V^*$ is smooth, there is a $\underline{\delta} < 1$ such that for all δ satisfying $\underline{\delta} \leq \delta < 1$, each point in W corresponds to a strict perfect public equilibrium with discount factor δ , in which a pure action profile is played in each period. In particular, we can choose W to have a boundary as close as we might desire to $\mathbf{v}^* \equiv (b^*, \dots, b^*)$, in which case the full cooperation payoff can be approximated as closely as desired.

The only condition of the theorem that must be verified in the case of the Public Goods Game is that the full cooperation payoff $\mathbf{v}^* = \{b^*, \dots, b^*\}$ is on the boundary of an open set of payoffs in \mathbf{R}^n , assuming players can use mixed strategies. Suppose player i cooperates with probability x_i , so the payoff to player i is $v_i = \pi_i - cx_i$, where

$$\pi_i = b \sum_{j=1}^n x_j - x_i.$$

If J is the Jacobian of the transformation $x \rightarrow v$, it is straightforward to show that

$$\det[J] = (-1)^{n+1}(b - c) \left(\frac{b}{n - 1} + c \right)^{n-1},$$

¹I am suppressing two conditions on the signal y that are either satisfied trivially or irrelevant in the case of a Public Goods Game.

which is non-zero, proving the transformation is not singular.

The method of recursive dynamic programming used to prove this theorem in fact offers an equilibrium construction algorithm, or rather, a collection of such algorithms. Given a set $W \subset V^*$, a discount factor δ , and a strategy profile α , we say α is *enforceable* with respect to W and δ if there is a payoff vector $v \in \mathbf{R}^n$ and a *continuation function* $w: Y \rightarrow W$ such that for all i ,

$$v_i = (1 - \delta)g_i(a_i, \alpha_{-i}) + \delta \sum_{y \in Y} \pi(y|a_i, \alpha_{-i})w_i(y)$$

for all a_i with $\alpha_i(a_i) > 0$, (10.4)

$$v_i \geq (1 - \delta)g_i(a_i, \alpha_{-i}) + \delta \sum_{y \in Y} \pi(y|a_i, \alpha_{-i})w_i(y)$$

for all a_i with $\alpha_i(a_i) = 0$. (10.5)

We interpret the continuation function as follows. If signal $y \in Y$ is observed (the same signal will be observed by all, by assumption), each player switches to a strategy profile in the repeated game that gives player i the long-run average payoff $w_i(y)$. We thus say that $\{w(y)_{y \in Y}\}$ *enforces* α with respect to v and δ , and that the payoff v is *decomposable* with respect to α , W , and δ . To render this interpretation valid, we must show that $W \subseteq E(\delta)$, where $E(\delta)$ is the set of average payoff vectors that correspond to equilibria when the discount factor is δ .

Equations (10.4) and (10.5) can be used to construct an equilibrium. First, we can assume that the equations in (10.4) and (10.5) are satisfied as equalities. There are then two equations for $|Y| = 2^n$ unknowns $\{w_i(y)\}$ for each player i . To reduce the underdetermination of the equations, we shall seek only pure strategies that are symmetrical in the players, so no player can condition his behavior on having a particular index i . In this case, that $w_i(y)$ depends only on whether or not i signaled cooperate, and the number of other players who signaled cooperate. This reduces the number of strategies for a player at this point from 2^n to $2(n - 1)$. In the interests maximizing efficiency, we assume that in the first period all players cooperate, and as long as y indicates universal cooperation, players continue to play all cooperate.

To minimize the amount of punishment meted out in the case of observed defections while satisfying (10.4) and (10.5), we first assume that if more than one agent signals defect, all continue to cooperate. If there is a single defection, this is punished by all players defecting an amount that just

satisfies the incentive compatibility equations (10.4) and (10.5). There is of course no assurance that this will be possible, but if so, there will be a unique punishment γ that is just sufficient to deter a self-regarding player from intentionally defecting. This level of γ is determined by using (10.2) and (10.3) to solve (10.4) and (10.5) using the parameters for the Public Goods Game. The calculations are quite tedious, but the solution for γ in terms of the model parameters is

$$\gamma = \frac{c(1 - \delta)}{\delta(1 - \epsilon)^{n-1}(1 - n\epsilon)}. \quad (10.6)$$

Note that γ does not depend on b . This is because the amount of punishment must only induce a player to expend the cost c . The fraction of a period of production entailed by a given γ will, of course, depend on b . Note also that (10.6) only holds for $n\epsilon < 1$. We deal with the more general case below.

We can now calculate $v = v_i(\forall i)$, the expected one-period payoff. Again, the calculations are tedious, but we find

$$v = b(1 - \epsilon) - c - \frac{n\epsilon c(1 - \delta)}{1 - n\epsilon}. \quad (10.7)$$

This shows that when $n\epsilon$ approaches unity, the efficiency of cooperation plummets.

The above solution is meaningful only when $n\epsilon < 1$. Suppose k is a positive integer such that $k - 1 \leq n\epsilon < k$. An extension of the above argument shows that if no punishment is meted out unless exactly k defections are signaled, then (10.7) becomes

$$v = b(1 - \epsilon) - c - \frac{n\epsilon c(1 - \delta)}{k - n\epsilon}. \quad (10.8)$$

Again, for δ sufficiently close to unity, we can approximate Pareto-efficiency as closely as desired.

By inspecting (10.4) and (10.5), we can gain some insight into what the Folk Theorem is really saying in this case. When $n\epsilon$ is large, punishment is kept down by requiring several defections to trigger the punishment, in which case the punishment continues over several periods, during which payoffs are zero. However, with positive probability when cooperation resumes there will be few defections, and if δ is near unity, the sum of payoffs over these rare periods of cooperation will be high, so the second term in

(10.4) will be large. Moreover, for δ near unity, the first term, representing the expected current period payoff, is near zero, so the present value of cooperation will be determined by the second term as $\delta \rightarrow 1$. There is clearly no sense in which this can be considered a solution to the problem of cooperation in large groups.

10.5 Cooperation with Private Signaling

Repeated game models with private signals, including Bhaskar and Obara (2002), Ely and Välimäki (2002), and Piccione (2002), are subject to the critique of the previous section, but private signaling models are complicated by the fact that no sequential equilibrium can support full cooperation, so strictly mixed strategies are necessary in equilibrium. To see this, consider the first period. If each player uses the full cooperation strategy, then if a player receives a defection signal from another player, with probability one this represents a bad signal rather than an intentional defection. Thus, with very high probability, no other member received a defection signal. Therefore no player will react to a defect signal by defecting, and hence the always defect strategy will have a higher payoff than the always cooperate strategy. To deal with this problem, *all players defect with positive probability in the first period*.

Now, in any Nash equilibrium, the payoff to any two pure strategies that are used with positive probability by a player must have equal payoffs against the equilibrium strategies of the other players. Therefore, the probability of defecting must be chosen so each player is indifferent between cooperating and defecting at least on the first round. Sekiguchi (1997) and Bhaskar and Obara (2002) accomplish this by assuming players randomize in the first round and play the grim trigger strategy in each succeeding round—cooperate as long as you receive a signal that your partner cooperated in the previous round, and after receiving a defect signal, defect yourself in each succeeding round. After the first round, it is possible that a defect signal really means a player defected, because that player, who is also playing a trigger strategy, could have received a defect signal in the previous round.

This model is plausible when the number n of players is small, especially for $n=2$. However, when the error rate approaches $1/n$, the model becomes inefficient, because the probability of at least one agent receiving a defect signal approaches unity, so the expected number of rounds of the game is

close to one, where the benefits of game repetition disappear. Moreover, it will be true in most cases that the quality of the private signal deteriorates with increasing group size (e.g., if each individual receives a signal from a fixed maximum number of other players). As the authors show, we can improve the performance of the model by “restarting” cooperation with a positive probability in each round after cooperation has ceased, but this only marginally improves efficiency, because this process does not increase the incentive of players to cooperate on any given round.

Ely xxand Välimäki (2002) have developed a different approach to the problem, following the lead of Piccione (2002), who showed how to achieve coordination in a repeated game with private information without the need for the sort of belief updating and grim triggers used by Sekiguchi, Bhaskar, and Obara. They construct an equilibrium in which at every stage, each player is *indifferent* between cooperating and defecting no matter what his fellow members do. Such an individual is thus willing to follow an arbitrary mixed strategy in each period, and the authors show that there exists such a strategy for each player that ensures close to perfect cooperation, provided individuals are sufficiently patient and the errors are small.

One problem with this approach is that it uses mixed strategies in every period, and unless the game can be purified (§6.2), there is no reason for players to play such strategies or to believe that their partners will do so either. Bhaskar (2000) has shown that most repeated game models that use mixed strategies cannot be purified, and Bhaskar, Mailath and Morris (2004) have shown that purification is generally impossible in the Ely-Välimäki approach to the Prisoner’s Dilemma when the signal is public. The case of private signals is much more difficult, and there is no known example of purification in this case.

Without a choreographer, there is no mechanism that coordinates activities of large numbers of people so as to implement a repeated game equilibrium with private information. It follows that the issue of whether or not such games, the Nash equilibria of which invariably require strictly mixed strategies, can be purified is not of fundamental importance. Nevertheless, it is useful to note that there are no examples of purification of such games, and at least two examples of the impossibility of purification. These examples, due to Bhaskar (1998,2000), make clear why purification is not likely: the Nash equilibria in repeated games with private information are “engineered” so that players are indifferent between acting on information concerning defection and ignoring such information. A slight change in

payoffs, however, destroys this indifference, so players behave the same way whatever signal they receive. Such behavior is not compatible with a cooperative Nash equilibrium.

It might be thought that the lack of purification is not a fatal weakness, however, because we have already shown that the social instantiation of a repeated game model requires a choreographer, and there is no reason, at least in principle, that a social norm could not implement a mixed strategy σ by suggesting each of the pure strategies in the support of σ with a probability equal to its weight in σ . This idea is, however, incorrect, as we exhibited in §6.3. Unless players have a sufficiently strong normative predisposition, small random changes in payoffs will induce players to deviate from the choreographer's suggestion. The lack of purification for these models is virtually fatal—they cannot be socially instantiated.

10.6 One Cheer for the Folk Theorem

The Folk Theorem is the most promising analytically rigorous theory of human cooperation in the behavioral sciences. Its strength lies in its transformation of Adam Smith's *invisible hand* into analytical model of elegance and clarity. The Folk Theorem's central weakness is that it is only an existence theorem with no consideration for how the Nash equilibria whose existence it demonstrates can actually be instantiated as a social process. Certainly these equilibria cannot be implemented spontaneously or through a process of player learning. Rather, as we have stressed throughout this volume, strategic interaction must be socially structured by a choreographer—a social norm with the status of common knowledge, as outlined in chapter 7.

This weakness is analytically trivial, but scientifically monumental. Correcting it both strengthens repeated game models and suggests how they may be empirically tested—namely, by looking for the choreographer, and where it cannot be found, determining what premise of the repeated game model is violated, and proposing an alternative model. Recognizing the normative dimension of social cooperation has the added benefit of explaining why repeated game models have virtually no relevance beyond *Homo sapiens* (Clements and Stephens 1995, Stephens et al. 2002, Hammerstein 2003), the reason being that normative behavior is extremely primitive, at best, for nonhuman species.

A second weakness of repeated game theory is its preoccupation with situations in which players are almost perfectly future-oriented (i.e., use a dis-

count factor close to unity) and noise in the system (e.g., signaling stochasticity or player error) is arbitrarily small. The reason for this preoccupation is simple: the Folk Theorem with self-regarding agents fails when either agents are present-oriented, signals are imperfect, or players are likely to err.

The correct response to this weakness is to (a) observe how cooperation really occurs in society and (b) alter the characteristics of the repeated game model to incorporate what one has discovered. We learn from biology that there are huge gains to cooperation for an organism, but the challenges of coordinating behavior and keeping defection to manageable levels are extreme, and solve only by rare genetic innovation (Maynard Smith and Szathmáry 1997). The notion that human cooperation has a strong biological element, as we stressed in chapter 3, is in line with this general biological point. We present in this chapter, for illustrative purposes, a model of cooperation based on observed characteristics of humans that are not captured either by Bayesian rationality or the social epistemology developed in earlier chapters (§10.7).

Both common observation and behavioral experiments suggest that humans are disposed to behave in prosocial ways when raised in appropriate cultural environments (Gintis 2003a). This disposition includes having other-regarding preferences, such as empathy for others, and the predisposition to embrace cooperative norms and to punish the violators of these norms, even at personal cost. It also includes upholding such character virtues as honesty, promise-keeping, trustworthiness, bravery, group-loyalty, and considerateness. Finally, it includes valuing self-esteem and recognizing that self-esteem depends on how one is evaluated by those with whom we strategically interact. Without these prosocial, biologically rooted traits, human language could not develop, because there would then be no means of maintaining veridical information transmission. Without high quality information, efficient cooperation based on repeated game Nash equilibria would be impossible. Indeed, it is probably rare that information is of sufficient quality to sustain cooperation of self-regarding actors.

10.7 Altruistic Punishing in the Public Goods Game

This section develops a model of cooperation in the Public Goods Game, in which each agent is motivated by self-interest, unconditional altruism, and strong reciprocity, based on Carpentier et al. (2009). We investigate the

conditions for a cooperative equilibrium, as well as how the efficiency of cooperation depends on the level of altruism and reciprocity. We show that if there is a stable interior equilibrium (i.e., including both cooperation and shirking), an increase in either altruism or reciprocity motives will generate higher efficiency.

Consider a group of size $n > 2$, where member i supplies an amount of effort $1 - \sigma_i \in [0, 1]$. We call σ_i the *level of shirking* of member i , and write $\bar{\sigma} = \sum_{j=1}^n \sigma_j / n$ for the average level of shirking. We assume shirking at level σ_i adds $q(1 - \sigma_i)$ dollars to group output, where $q > 1$, while the cost of working is a quadratic function $s(1 - \sigma) = (1 - \sigma)^2/2$. We call q the *productivity of cooperation*. We assume the members of the group share their output equally, so member i 's payoff is given by

$$\pi_i = q(1 - \bar{\sigma}) - (1 - \sigma_i)^2/2. \quad (10.9)$$

The payoff loss to each member of the group from one member shirking is $\beta = q/n$. We assume $1/n < \beta < 1$.

We assume member i can impose a cost on $j \neq i$ with monetary equivalent s_{ij} at cost $c_i(s_{ij})$ to himself. The cost s_{ij} results from public criticism, shunning, ostracism, physical violence, exclusion from desirable side-deals, or another form of harm. We assume $c_i(0) = c_i'(0) = c_i''(0) = 0$ and $c_i(s_{ij})$ is increasing and strictly convex for all i, j when $s_{ij} > 0$.

Member j 's cooperative behavior b_j depends on j 's level of shirking and the harm that j inflicts on the group, which we assume is public knowledge. Specifically, we assume

$$b_j = \beta(1 - 2\sigma_j). \quad (10.10)$$

Thus, $\sigma_j = 1/2$ is the point at which i evaluates j 's cooperative behavior as neither good nor bad.

To model cooperative behavior with social preferences, we say that individual i 's utility depends on his own material payoff π_i and the payoff π_j to other individuals $j \neq i$ according to:

$$u_i = \pi_i + \sum_{j \neq i} [(a_i + \lambda_i b_j)(\pi_j - s_{ij}) - c_i(s_{ij})] - s_i(\sigma_i) \quad (10.11)$$

where $s_i(\sigma_i) = \sum_{j \neq i} s_{ji}(\sigma_i)$ is the punishment inflicted upon i by other group members, and $\lambda_i \geq 0$. The parameter a_i is i 's level unconditional altruism if $a_i > 0$ and unconditional spite if $a_i < 0$, and λ_i is i 's strength

of reciprocity motive, valuing j 's payoffs more highly if j conforms to i 's concept of good behavior, and conversely (Rabin 1993, Levine 1998). If λ_i and a_i are both positive, the individual is termed a strong reciprocator, motivated to reduce the payoffs of an individual who shirks even at a cost to himself.

Players maximize (10.11), and because b_j can be negative, this may lead i to increase his shirking σ_i and/or to punish j by increasing s_{ij} in response to a higher level of shirking by j . This motivation for punishing a shirker values the punishment *per se* rather than the benefits likely to accrue to the punisher if the shirker responds positively to the punishment. Moreover, members derive utility from punishing the shirker, not simply from observing that the shirker was punished. This means that punishing is *warm glow* rather than instrumental towards affecting j 's behavior (Andreoni 1995, Casari and Luini 2007).

This model requires only that a certain fraction of group members be reciprocators. This is in line with the evidence from behavioral game theory evidence presented in chapter 3, which indicates that in virtually every experimental setting a certain fraction of the subjects do not act reciprocally, either because they are self-regarding or they are purely altruistic. Note also that the punishment system could elicit a high level of cooperation, yet a low level of net material payoff. This is because punishment is not strategic in this model. In real societies, the amount of punishment of shirkers is generally socially regulated, and punishment beyond the level needed to secure compliance is sanctioned (Wiessner 2005).

In this model, i will choose $s_{ij}^*(\sigma_j)$ to maximize utility in (10.11), giving rise to the first order condition (assuming an interior solution)

$$c'_i(s_{ij}^*) = \lambda_i \beta (2\sigma_j - 1) - a_i. \tag{10.12}$$

If $\lambda_i > 0$ and

$$\sigma_j \leq \sigma_i^0 = \frac{1}{2} \left[\frac{a_i}{\lambda_i \beta} + 1 \right], \tag{10.13}$$

the maximization problem has a corner solution in which i does not punish. For $\lambda_i > 0$ and $\sigma_j > \sigma_i^0$, denoting the right hand side of (10.12) by ϕ and differentiating (10.12) totally with respect to any parameter x , we get

$$\frac{ds_{ij}^*}{dx} = \frac{\partial \phi}{\partial x} \frac{1}{c''_i(s_{ij}^*)}. \tag{10.14}$$

In particular, setting $x = a_i$, $x = \lambda_i$, $x = \sigma_j$, $x = \beta$ and $x = n$ in turn in (10.14), we see that

THEOREM 10.3 For $\lambda_i > 0$ and $\sigma_j > \sigma_i^0$, the level of punishment by i imposed on j , s_{ij}^* , is (a) decreasing in i 's unconditional altruism a_i ; (b) increasing in i 's reciprocity motive, λ_i ; (c) increasing in the level σ_j of j 's shirking; (d) increasing in the harm β that j inflicts i by shirking; and (e) decreasing in group size.

The punishment $s_j(\sigma_j)$ inflicted upon j by the group is given by

$$s_j(\sigma_j) = \sum_{i \neq j} s_{ij}^*(\sigma_j), \quad (10.15)$$

which is then differentiable and strictly increasing in σ_j over some range, provided there is at least one reciprocator i ($\lambda_i > 0$).

The first order condition on σ_i from (10.11) is given by

$$1 - \sigma_i - \beta = \beta \sum_{j \neq i} (a_i + \lambda_i b_j) + s'_i(\sigma_i), \quad (10.16)$$

so i shirks up to such point as the net benefits of shirking (the left hand side) equal i 's valuation of the cost imposed on others by his shirking (the first term on the right hand side) plus the marginal cost of shirking entailed by the increased level of punishment that i may expect. This defines i 's optimal shirking level σ_i^* for all i , and hence closes the model, assuming the second order conditions $s''_i(\sigma_i) > -1$. Whether there is an interior solution depends on the array of parameters of the problem. For instance, if reciprocity is very weak, there could be complete shirking by every player, or if very strong, zero shirking by every player. We assume an interior solution to investigate the comparative statics of the problem.

The average shirking rate of i 's partners is given by

$$\bar{\sigma}_{-i} = \frac{1}{n-1} \sum_{j \neq i} \sigma_j.$$

We say that i 's partners *shirk on balance* if $\bar{\sigma}_{-i} > 1/2$ and they *work on balance* if the opposite inequality holds. We then have the following theorem, which is proved in Carpentier et al. (2009):

THEOREM 10.4 *Suppose there is an stable interior equilibrium under a best response dynamic. Then (a) an increase in i 's unconditional altruism a_i leads i to shirk less; and (b) an increase in i 's reciprocity motive λ_i leads i to shirk more when i 's partners shirk on balance, and to shirk less if i 's partners work on balance.*

While this is a simple one-shot model, it could easily be developed into a repeated game model in which some of the parameters evolve endogenously, and where reputation effects strengthen the other-regarding motives on which the above model depends.

10.8 The Failure of Models of Self-regarding Cooperation

Providing a plausible game-theoretic model of cooperation among self-regarding agents would vindicate methodological individualism (§8.8), and render economic theory virtually independent from, and foundational to, the other behavioral disciplines. In fact, this project is not a success. A fully successful approach is likely to require a psychological model of social preferences and a social epistemology, as well as an analysis of social norms as correlating devices that choose among a plethora of Nash equilibria and choreograph the actions of heterogeneous agents into a harmonious operational system.